

BRNO UNIVERSITY OF TECHNOLOGY FACULTY OF ELECTRICAL ENGINEERING AND COMMUNICATION

PROCEEDINGS II OF THE 30TH STUDENT EEICT 2024

ELECTRICAL ENGINEERING, INFORMATION SCIENCE, AND COMMUNICATION TECHNOLOGIES SELECTED PAPERS

April 23, 2024 Brno, Czech Republic



Title	Proceedings II of the 30 th Conference STUDENT EEICT 2024
Editor	Assoc. Prof. Vítězslav Novák
Publisher	Brno University of Technology, Faculty of Electrical Engineering and Communication
Year released	2024
	1 st Edition

Content and language issues are the responsibility of the authors. ISBN 978-80-214-6230-4 ISSN 2788-1334



SPONSORS Diamond tier



Golden tier



onsemi

ThermoFisher SCIENTIFIC



SPONSORS

Silver tier











Honorary tier



TECHNICAL SPONSORS & MEDIA PARTNERS







Contents

Foreword

Miriam Nagyová, Michal Nohel	
Deep learning model for segmentation of trabecular tissue on CT data of the lumbar spine	8
Adéla Fialová, Karel Sedlář	
BacSeger: Read simulator for bacterial RNA-Seg	12
David Podrazký	
Laboratory database for bacterial sample data cataloguing in the hospital environment	16
Václav Kubeš, Tomáš Urbanec	
Four Channel Active Antenna Switch for UHF Band Satellite Reception	20
Martin Ťavoda	
Development Module for Radar Safety Sensor in Single-Track Vehicles	25
Simon Buchta, Ladislav Polak	
Measurement of the DVB-T2-based MISO Signal influenced by I/Q-errors in the OFDM Modulator	
S. Orgoň	
Device for galvanic plating of 3D printed parts	
Miloslav Kužela, Tomáš Frýza	
Detection of parking space availability based on video	
Roman Vaněk, Jan Mikulka	
Modular system for electrical impedance tomography	40
Matej Grega	
Design of the Electronic Target for Shooting Sports and Sensor Suitability Analysis	
Sebastian Simanský	
Enhancement of Vehicle Dynamics Through Adaptive Torque Vectoring Control with PMSM Powertrain	
Stěrba Radek, Stanislav Pikula	
Shock severity comparison using Shock Response Spectrum and Pseudo-Velocity Shock Spectrum in LabVIEW	53
Valentina Hrtonova, Marina Filipenska, Petr Klimes	
Graph Neural Networks in Epilepsy Surgery	57
Katerina Ingrova, Larisa Chmelikova, Inna Zumberg	<i>c</i> 1
Quantitative analysis of Scratch assay and Colony formation assay data using MAILAB	61
Adriana Spakova, Vaishali Pankaj, Inderjeet Bhogal, Sudeep Koy	65
Probing natural molecules with PPAR-γ to reveal potent agonist against Cancer	
Jan Drska Madamination of Tauma Diala and Diaga Mashing	70
Modernization of Tenryu Pick and Place Machine	
Automated testing device for turboist ECU	74
Automated testing device for furbojet ECO	/4
Semiautomatic crimping machine	77
Matěi Podaný Jaroslav Láčík	••••••
Design of Miniaturized Logarithmic-Periodic Antenna	81
Michal Baranek, Ladislav Polak, Jan Kufa	
Mapping and analyzing of signal coverage of 4G/5G mobile networks	
Tomáš Kříčka. Jan Král	
Platform for Digital Predistortion Based on RFSoC	
Patrik Pis, Willi Lazarov	
Utilizing Dynamic Analysis for Web Application Penetration Testing	
Marco Pintér, Petr Marcoň	
Navigation of UAV in GNSS denied area	
Jiří Veverka, Aleš Povalač, Kamil Jaššo	
Educational PocketQube Satellite Demonstrator	
Patrik Staroň	
Electrical characterization of graphene sensors	

Daniel Havránek, Michael Pešek, Daniel Paulitsch, Magdalena Kowalska, Mark Lloyd Bissel, Jan Král	
Development of Beta-NMR Detection Electronics for VITO Beamline at ISOLDE Facility at CERN	109
Jaromír Jarušek, Jan Brodský, Imrich Gablech	
Transparent materials for planar microelectrode arrays	113
Jakub Vašíček, Petr Vyroubal	
Estimating the Equivalent Circuit of Lithium-Ion Batteries During Operation	117
Silvia Batorova, Miroslav Zatloukal	101
Corrosion potential analysis of iron-magnesium alloys	121
Vaciav Kutnar, Kamil Jasso	10.4
LCA of different types of cars in Czech Republic	124
Michai Nonei, Jiri Unmelik	120
Calculation of Bone Mineral Density from Dual-energy C1 and its Application on Patient with Multiple Myeloma	129
David Kallipo Assessing Diversity in Dradictive Equations for Pody Comportment Estimation	124
Assessing Diversity in Fredictive Equations for Body Compartment Estimation	134
Ontimizing of ma processing englysis for Illuming DNA. Sog data in Anghidensis thaligns	140
Opuninzing of pre-processing analysis for munima KIVA-seq data in <i>Arabiaopsis manana</i>	142
Quantitative Analysis of Vacal Tract Decomances in Patients with Parkinson's Disease	146
Qualitative Analysis of vocal fract Resonances in fatients with fatkinson's Disease	140
I ong Term Effect of Repetitive Transcranial Magnetic Stimulation on Parkinson's Disease Patients with Different Severit	w
of Hypokinetic Dysortheia	.y 151
Š Slavaril	151
Workplace buildup for measurement of low frequency absorption of acoustic structures	156
Pavel Palurik Radim Dyorak	150
Hardware Design of Mohile Probe for Validating 5G-IoT Technologies	161
Aneta Kolackova Jan Jerahek	101
Strategic Capacity Planning and Ontimization in Communication Networks: A Case Study	166
Patrik Dobiáš Lukáš Malina	100
Ontimizing components for Dilithium and Kyber unified hardware implementation	171
Samer I. Daradkeh. Oscar Recalde. Marwan S. Mousa, Dinara Sobola	
Overdoping effect with Zr and Hf on the oxidation behaviour of FeCrAl-Hf by means of Atom Probe Tomography	176
Mohammad M. Allaham. Zuzana Košelová. Dinara Sobola. Zdenka Fohlerová. Alexandr Knánek	
An enhanced theoretical approach for accurate measurements of the optical and energy characteristics of semiconducting	
materials	182
David Trochta. Ondřei Klvač. Tomáš Kazda	
A Review of the Li-ion Battery in-situ Experiments in Scanning Electron Microscope	187
Helena Picmausová, Jan Farlík, Marc Eichhorn, Christelle Kieleck	
Lidar systems testing considerations for field use	192
Tadeáš Zbožínek, Michal Jelínek, Břetislav Mikel	
Evaluation of the gamma-ray spectrum transmitted upon alteration in scintillator shading	197
Zuzana Košelová, Zdenka Fohlerová	
Enzyme-Based Impedimetric Biosensor dotted with gold nanoparticles	202
Stepan Jezek, Radim Burget	
Deployment of deep learning-based anomaly detection systems: challenges and solutions	207
Dominik Ricanek	
Amplitude measurement of small displacement using video magnification	212
Jan Klouda, Petr Marcoň	
The Human-machine interface for UAV ground control station	218
Thao Dinh Le, Pavel Masek	
Machine Learning-based Fingerprinting Localization in 5G Cellular Networks	222
Vladimír Bílek	
Comparative Analysis of Gaussian Process Regression Modeling of an Induction Machine: Continuous vs. Mixed-Input	
Approaches	227
Marek Šťastný, Tomáš Dytrych, Vladimír Dániel, Kryštof Mrózek, Adam Obrusník	
Feastibility study and comparative study of air breathing electric propulsion systems operating in very low Earth orbit	
conditions	232

OPENING WORD OF THE DEAN

These Proceedings contain papers presented during the **30th annual STUDENT EEICT conference**, held at the Faculty of Electrical Engineering and Communication, Brno University of Technology, on April 23, 2024. The fruitful tradition of joining together creative students and seasoned science or research specialists and industry-based experts was not discontinued, providing again a valuable opportunity to exchange information and experience.

The EEICT involves multiple corporate partners, collaborators, and evaluators, whose intensive support is highly appreciated. Importantly, the competitive, motivating features of the conference are associated with a practical impact: In addition to encouraging students to further develop their knowledge, interests, and employability potential, the forum directly offers career opportunities through the affiliated PerFEKT JobFair, a yearly job-related workshop and exhibition complementing the actual EEICT sessions. In this context, the organizers acknowledge the long-term assistance from the Ministry of Education, Youth and Sports of the Czech Republic, which has proved essential for refining the scope and impact of the symposium.

In total, 137 peer-reviewed full papers distributed between 18 sessions were submitted, before examining boards with industry and academic specialists. The presenting authors exhibited a very high standard of knowledge and communication skills, and the best competitors received prize money and/or valuable gifts. These Proceedings comprise 54 award winning full papers, all selected by the conference's evaluation boards.

Our sincere thanks go to the sponsors, experts, students, and collaborators who participated in, contributed to, and made the conference a continued success.

Considering all the efforts and work invested, I hope that the 30th STUDENT EEICT (2024) has been beneficial for all the participants.

I believe that the inspiration gathered during the event will contribute towards a further rise of open science and research, giving all the attendees a chance to freely discuss their achievements and views.

Prof. Vladimír Aubrecht Dean of the Faculty of Electrical Engineering and Communication

Deep learning model for segmentation of trabecular tissue on CT data of the lumbar spine

Miriam Nagyová Department of Biomedical Engineering Brno University of Technology, FEEC Brno, Czech Republic 231091@vut.cz

Abstract—This paper focuses on training a deep learning model for vertebral body segmentation of the lumbar spine. The nnU-Net model was trained and tested on a publicly available dataset LumVBCanSeg consisting of 185 lumbar CT scans. Dice coefficient was used to evaluate the accuracy of the trained model. The mean Dice coefficient of the testing dataset was 0.949 with a standard deviation of 0.103. The model was also tested on clinical data containing various abnormalities, such as lytic lesions in multiple myeloma patients and metallic implants. Results were evaluated visually. While the model showed high accuracy on the testing dataset, the results on scans with anomalies showed a decline in accuracy.

Index Terms—multiple myeloma, osteolytic lesions, nnU-Net, segmentation

I. INTRODUCTION

The spine, a crucial and intricate structure, provides support for the body's mechanical functions, safeguards the delicate spinal cord, and facilitates movement. Maintaining a healthy spine is crucial for overall well-being and mobility. The human spine consists of small individual bones called vertebrae. The human spine's vertebrae are named based on their location along the spinal column. Vertebrae exhibit variations in shape and size depending on their placement within the spine. Despite these differences, they share a remarkably similar structure, consisting of three parts: a body, a vertebral arch, and processes. [1]

Multiple myeloma (MM) represents a form of bone marrow cancer and stands as the most prevalent cancer within the skeletal system. MM causes the development of osteolytic bone lesions, characterized by areas of compromised bone tissue that render the bone more fragile and susceptible to fractures. Common symptoms of MM include excruciating bone pain, typically concentrated around the chest and back, hypercalcemia, spinal cord compression, and pathological fractures. [2]

While MM can impact any bone, it tends to manifest more prominently in bones with a higher concentration of red marrow. Consequently, bones like vertebral bodies, the skull, pelvis, and ribs are particularly susceptible to MM's effects. [2]

This paper focuses on training our own nnU-Net [3] model on the publicly available dataset LumVBCanSeg [4]. This dataset was constructed by the authors of [5]. They proposed Michal Nohel Department of Biomedical Engineering Brno University of Technology, FEEC Brno, Czech Republic xnohel04@vutbr.cz

a deep learning-based, two-stage, coarse-to-fine solution designed for automating the segmentation of cancellous bone in the lumbar vertebral bodies from CT scans of the lumbar spine.

The accuracy of our trained model was evaluated using the Dice coefficient and Hausdorff distance. The model was later tested on clinical data of patients with MM.

II. DATASETS DESCRIPTION

A. Publicly available dataset LumVBCanSeg

For training and testing purposes, publicly available dataset LumVBCanSeg (A Lumbar Vertebral Body Cancellous Bone Segmentation Dataset) was used [4]. The comprehensive dataset, consisting of 185 lumbar CT scans, was acquired using multiple CT scanners, including those manufactured by Philips and Siemens. All scans were sourced from ShengJing Hospital of China Medical University. Cases involving vertebral fractures, metal implants, bone tumors, and foreign materials were purposefully omitted from the dataset. Additionally, to ensure uniformity, all data underwent resampling to achieve an isotropic resolution of $1 \times 1 \times 1$ mm. Additionally to the CT scans, the dataset contains corresponding segmentation masks (see Fig. 1). Annotations cover the bodies of five lumbar vertebrae, L1 to L5, and are labeled sequentially from 1 to 5.

As we can see in Fig. 1, the mask only includes the trabecular part of the vertebral body, omitting the periostenum and processes. The dataset was divided in an 80:20 ratio, with 80% of the dataset being used for training, leaving the 20% for testing purposes. This means that we got 143 scans in the training set and 37 were used for testing purposes.

B. Clinical data

In addition to the publicly available dataset, a dataset containing clinical data is also available to test the accuracy of the model. The dataset consists of 20 scans and their corresponding segmentation masks. However, in contrast to the publicly available dataset, most of the images are wholebody CT scans. The dataset contains scans of 10 individuals without any pathology in spine and 10 patients who have been diagnosed with multiple myeloma. Unlike the previous



Fig. 1: An example of lumbar CT scan with the corresponding segmentation mask from publicly available dataset

dataset, this contains all kinds of inconsistencies, including metal implants and bone lesions (see Fig. 2 and Fig. 3).

Data were acquired following approval from the Ethics Committee, under the application registration number NU23J-08-00027. All patients gave their consent after being informed. The data were acquired utilizing the Philips Healthcare IQon spectral CT system in collaboration with the University Hospital Brno, Department of Radiology and Nuclear Medicine. The scanning parameters included a peak tube voltage of 100 kV, tube current of 10 mA, matrix size of 512×512 , and a slice thickness of 0.9 mm using a sharp reconstruction kernel and hybrid iterative reconstruction technique (iDose4, set to level 4). Scans of the myeloma patients covered the region from the head to the knee. The scans were examined using a dedicated workstation (Intellispace Portal version 12.1; Philips Healthcare) by two independent readers, with at least one being board certified. Patients were diagnosed with MM based on elevated monoclonal immunoglobulin in the blood and an increased plasma cell count in the bone marrow.



Fig. 2: Close up look at multiple myeloma lesions



Fig. 3: An example of scans in the clinical dataset, (a) - patient with lesions, (b) - patient with a metallic implant

III. MODEL IMPLEMENTATION

In this work, a deep learning model was trained based on the nnU-Net framework [3], which is an automatic semantic segmentation technique designed to adjust to diverse datasets. This means that even without specialized knowledge, the model can be effortlessly trained and applied to a range of applications. It excels particularly in semantic segmentation, demonstrating proficiency in processing both 2D and 3D images across diverse input modalities and channels. Its adaptive nature extends to accommodating differences in voxel spacings and anisotropies, demonstrating resilient performance even in cases with substantial class imbalances. nnU-Net utilizes supervised learning, which means that it is necessary to submit a set of training cases specific to the application.

As mentioned above, the training dataset only contained CT scans of the lumbar region. This meant that to achieve good results, the data from the clinical dataset had to be cropped to the desired size. This was done by creating a bounding box with intended dimensions (see Fig. 4).



Fig. 4: Scans after being cropped, (a) - no spinal pathologies, (b) - MM patient

nnU-Net automatically devised the optimal architecture and learning parameters. Configuration was set to 3d fullres, the learning rate was initialized at 0.01 and progressively decreased throughout the training process, the batch size was configured to 2, patch size was $128 \times 128 \times 128$ px, featuring a kernel size of $3 \times 3 \times 3$. A total of 1000 learning epochs were executed. Training of the model took place on a dedicated computer at the Metacenter.

IV. USED METRICS

Dice coefficient [6] was used to evaluate the lumbar spine segmentation. Dice coefficient is a metric used to quantify the similarity between two sets. It can either be calculated for binary or multilabel classification. In our case, we used the multilabel approach, because of the nature of the available data. In this case, the result is calculated using the formula:

$$Dice(P,T) = \frac{1}{N} \sum_{i=1}^{N} \frac{2|P_i \cap T_i|}{|P_i| + |T_i|}$$
(1)

where P stands for predicted values, T represents the ground truth and i is an index of the N^{th} vertebrae. The final dice score ranges between 0 and 1 and has no unit.

The Hausdorff distance [7] is a metric used to assess the similarity between two sets of points within a given metric space. This distance is calculated by determining the maximum distance from any point in one set to the closest point in the other set. It serves as a measure of dissimilarity or mismatch between the two sets, finding applications in fields like computer vision, image analysis, pattern recognition, and shape matching. Hausdorff distance is given by the formula:

$$HD(A,B) = \max(\max_{a \in A} \min_{b \in B} d(a,b), \max_{b \in B} \min_{a \in A} d(a,b)[px], (2)$$

where HD(A, B) represents the Hausdorff distance of sets A and B, the symbol max represents the maximum, the symbol min represents the minimum, d(a, b) represents the distance of point a and point b, where it can be for example the L2 norm, the symbol $a \in A$ indicates that the variable a takes values from the set A and the symbol $b \in B$ indicates that the variable b takes values from the set B. In our case, we calculated the HDs for each vertebral mask separately and then calculated their average value to evaluate the success of the segmentation model.

V. RESULTS AND DISCUSSION

The trained model was evaluated using the testing database, which consisted of 37 images. After a brief inspection of the results, we can say that the model performed well. Tab. I summarizes the calculated metrics. After inspecting the results, few outliers were observed, where both the metrics differ greatly from the rest of the testing dataset.

TABLE I: Calculated metrics

Metric	Dice coefficient [-]	Hausdorff distance [px]
Minimum value	0.526	1.166
Maximum value	0.986	173.816
Mean	0.949	26.189
Median	0.980	2.159
Standard deviation	0.103	48.527

Figures 5 and 6 visualize the results using box-and-whisker plots. It is visible that except for a few outliers, the model performed well. For better visualisation, 6 biggest outliers were eliminated.



Fig. 5: Multilabel Dice coefficient



Fig. 6: Multilabel Hausdorff distance



Fig. 7: Results of the segmentation by nn-Unet on the testing dataset, (a) - successful segmentation, (b) - failed segmentation

Fig. 7 shows us two results of the segmentation on the testing dataset. The failed segmentation shows that instead of 5 lumbar vertebrae, 6 vertebrae were segmented, including one thoracic vertebra, and that two of the vertebrae were assigned the same label.

The model was then applied to the clinical data. The results were evaluated visually. In Fig. 8 we can see the results of the model being applied to the clinical data. It is clear that the vertebrae with any anomalies were segmented incorrectly, whereas the healthy vertebrae were segmented accurately.

VI. CONCLUSION

This paper aimed to train a deep learning model for the segmentation of trabecular tissue on CT data of the lumbar spine. nnU-Net model was trained for this purpose. The trained model was tested on a test database and the outcomes were evaluated. The model was later tested on clinical data containing different anomalies. Overall, the model achieved high levels of accuracy on testing database. Vertebrae with anomalies were segmented incorrectly.



Fig. 8: Results of the segmentation on clinical data, problematic vertebrae

ACKNOWLEDGMENT

The paper and the research were supported by Philips Healthcare company and the University Hospital Brno, Department of Radiology and Nuclear Medicine. Computational resources were provided by the Ministry of Education, Youth, and Sports of the Czech Republic under the Projects CESNET (Project No. LM2015042) and CERIT-Scientific Cloud (Project No. LM2015085) provided within the program Projects of Large Research, Development, and Innovations Infrastructures. The Titan Xp GPU used for the implementation of the algorithms was donated by NVIDIA Corporation to support academic research.

REFERENCES

- J. Gordon Betts, Kelly A. Young, James A. Wise, Eddie Johnson, Brandon Poe, Dean H. Kruse, Oksana Korol, Jody E. Johnson, Mark Womble, Peter DeSaix. Houston, Texas: OpenStax, 2013.
- [2] R. Silbermann and G. D. Roodman, "Myeloma bone disease: Pathophysiology and management", *Journal of Bone Oncology*, vol. 2, no. 2, pp. 59-69, 2013.
- [3] Isensee, F., Jaeger, P. F., Kohl, S. A., Petersen, J., Maier-Hein, K. H. (2021). "nnU-Net: a self-configuring method for deep learning-based biomedical image segmentation." Nature methods, 18(2), 203-211.
- [4] Z. Yingdi, S. Zelin, W. Haun, C. Shaoqian, Z. Lei, L. Jiachen et. al. "LumVBCanSeg: A Lumbar Vertebral Body Cancellous Bone Segmentation Dataset"
- [5] Z. Yingdi, S. Zelin, W. Haun, C. Shaoqian, Z. Lei, and L. Jiachen, "LumVertCancNet: A novel 3D lumbar vertebral body cancellous bone location and segmentation method based on hybrid Swin-transformer", *Computers in Biology and Medicine*. vol. 171, 2024
- [6] L. R. Dice, "Measures of the Amount of Ecologic Association Between Species", *Ecology*, vol. 26, no. 3, pp. 297-302, 1945.
- [7] D. P. Huttenlocher, G. A. Klanderman, and W. J. Rucklidge, "Comparing images using the Hausdorff distance", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 15, no. 9, pp. 850-863.

BacSeqer: Read simulator for bacterial RNA-Seq

A. Fialová, K. Sedlář

Department of Biomedical Engineering Faculty of Electrical Engineering and Communication Brno University of Technology Czech Republic 240499@vut.cz, sedlar@vut.cz

Abstract-The RNA-Seq method has become a key tool in transcriptomics, providing deep insight into the genetic functioning of organisms. It is essential for understanding how genes are expressed, for discovering new gene variants and for understanding pathologies at the genetic level. RNA-Seq simulators are great for testing and validating bioinformatics tools, they offer a way to compare different computational algorithms and allow examining different sequencing protocols, parameter settings and sample sizes to ensure the most efficient use of resources. This study introduces BacSeger, an RNA-Seg simulator specifically designed for bacterial transcriptomes with the aim of producing data that faithfully resembles those obtained in real experiments. It allows to take into account the specificities of the bacterial genome, especially an existence of only exons and their organization in operons that are not normally present in eukarvotes.

Index Terms—RNA-Seq, simulation, sequencing, bacteria, gene expression

I. INTRODUCTION

RNA-Seq [1] has become a key technique in the field of genomics and functional genomics and transriptomics, providing a detailed view of an organism's transcriptome on a genome wide scale. This method allows high resolution quantification and characterization of gene expression, making it invaluable tool for studying regulatory mechanisms of cellular development, identifying new gene variants and understanding pathological conditions at the level of gene expression.

Simulations of RNA-Seq data have therefore become an important element in bioinformatics research and development. There are several reasons to simulate RNA-Seq data. The first reason is to validate bioinformatics tools and analyses. Simulated data allow testing and validating algorithms for assembly, quantification of expression and differential analysis of gene expression before applying them to real experimental data. This contributes to ensure the reliability of results and interpretations obtained from experimental data.

The second reason for simulation is the ability to generate data for experiments, that would otherwise be difficult or expensive. Simulations allow us to test different experimental conditions, changes in experimental design and compare the effect of technical parameters on the results of analyses. This increases the flexibility and efficiency of planning and interpretation of actual RNA-Seq experiments.

Several RNA-Seq simulators have been developed in the past, each offering unique features and capabilities. Notable examples are tools such as Polyester [2], Flux Simulator [3],

RNASeqRead Simulator [4], and BEERS2 [5]. Polyester is widely recognized for its flexibility and ability to model differential expression between conditions. It can generate data mimicking the distribution of transcript lengths and expression levels found in real biological samples. Flux Simulator is a comprehensive RNA-Seq data simulator that has gained popularity for its ability to model the complete workflow of RNA-Seq experiments. It is especially noted for its detailed emulation of the sequencing process, including RNA transcription, fragmentation, adapter ligation, sequencing, and read mapping. RNASeqSim offers a more targeted approach and specializes in simulating differential expression data. It is particularly useful for researchers who wish to evaluate the performance of different methods of differential expression analysis under controlled conditions. BEERS2 is an update to its predecessor BEERS, and improves on the original tool by providing more realistic simulation of read counts and splice variation. It can simulate different types of sequencing errors and biases.

In spite of all the differences, the simulators mentioned above have one common drawback - none of them offers consideration of the specifics of the bacterial genome, particularly the occurrence of operons in the bacterial genome or different length of untranslated regions (UTRs). An operon is a cluster of functionally related genes that are transcribed together from a single promoter into a single mRNA strand in prokaryotic organisms. This arrangement allows for coordinated regulation of gene expression, typically in response to environmental changes. Polycistronic RNA, which results from such transcription, contains the information for multiple proteins, with each distinct protein-coding region known as a cistron, allowing for simultaneous expression of several genes from the same mRNA molecule.

II. MATERIALS AND METHODS

A. Data sets

The data set from an experimental RNA-Seq analysis, which corresponded to the transcriptional response of *Clostridium beijerinckii* strain NRRL B-598 to butanol shock [6], was obtained from the NCBI Sequence Read Archive database under the accession number SRP033480.

As input data for the simulation with our tool we used the lastest version of *Clostridium beijerinckii* NRRL B-598 genome available from the NCBI GenBank database under the accession number CP011966.3.

B. Input

The BacSeqer simulator accepts any annotated bacterial genome sequences as input. These can be provided as sequences in FASTA format representing the genome, supplemented by annotation provided in a GTF or GFF file.

C. Functions

In addition to functions for opening input files and creating output files, the simulator includes several functions for sequential simulation of RNA-Seq data. This allows the user to go through the process step by step and set custom parameters.

The function *extract_operons* is designed to parse a GTF or GFF file and extract information about operons. It reads through the file and searches for attributes indicating the presence of operons. When it finds an operon, the function records its name and location (including the start and end points) in a dictionary. This dictionary, with operon names as keys and their locations as values, is then returned. The output of this function serves as the input parameter of the main simulation function.

The *parse_for_strand* function processes a GTF or GFF file to create a dictionary containing gene names or IDs as keys and corresponding strand information ('+' or '-') as values. This function is useful for mapping gene names to their respective strand orientations. This function plays a critical role in the *simulate_reads* function, particularly in providing strand orientation information for each sequence. When simulating reads, the *simulate_reads* function takes this strand information into account, especially when deciding whether to reverse complement the reads based on the strand orientation. This adds realism to the simulated reads, ensuring they reflect the actual transcription process.

The most important function is *simulate_reads*, designed to simulate reads from given sequences. It accepts multiple parameters, including a dictionary of sequences, desired read length in bp, number of reads, GC content (in decimals), operon locations (provided by *extract_operons*), and strand information (provided by *parse_for_strand*). The function iterates over each sequence and generates the specified number of reads. Firstly, the GC content is calculated using the *calculate_gc_content* function if not provided. If operon locations are given, the function considers them in generating reads. If strand information is provided, the reverse complement of the read is generated for sequences on both or reversed strands. The function returns a list of simulated reads.

The *calculate_rrna_percentage* function allows the user to calculate a rRNA percentage in the input sequences. It collects essential data like the sequence ID, start and end positions of the rRNA. Following this extraction, the function determines the proportion of rRNA in the total genomic sequence. It takes paths to both a FASTA file and a GFF file as an input, calculates the total length of the genomic sequences and the lengths of all rRNA segments and comparing this to the total

genomic sequence length, it calculates the percentage of the genome that is composed of rRNA.

The write_output function records simulated reads into a file in either FASTA or FASTQ format. For FASTQ format, it also allows setting maximum and minimum quality scores, generating these scores using the generate_quality_scores function if needed; this function synthesizes a descending sequence of Phred quality scores represented in ASCII characters. If the maximum and minimum values are not provided, default ASCII values based on typical Illumina quality scores are used. An auxiliary function to generate_quality_scores is a phred_to_ascii function, which translates a Phred quality score that quantitatively represents the accuracy of a nucleotide call in DNA sequencing into its corresponding ASCII character.

An overview of the functions can be seen in the block diagram in figure 1.



Fig. 1. Simulator flow chart.

D. Output

The simulator output is optionally a FASTA or FASTQ file with simulated RNA-Seq experiment data.

E. Evaluation tools

The FastQC and MultiQC tools were used to evaluate the simulated data and compare them with the results of real experiments. FASTQC [7] is a bioinformatics tool for quality control of raw sequencing data, FastQC generates reports on metrics like sequence quality and duplication, helping to quickly identify issues in high-throughput sequencing projects. MultiQC [8] aggregates data from various bioinformatics tools into a single report, providing a comprehensive overview of

multiple analyses in an accessible format, essential for largescale genomic studies.

III. RESULTS AND DISCUSSION

The simulated RNA-Seq data were analysed using FastQC, followed by a comparison of these results with the FastQC outputs from actual experimental data, utilizing MultiQC. The experimental data parameters were as follows: total sequences 43414028, sequence length 75, GC content 44% and per base sequence quality 30. The parameters of the simulated data were therefore set to the same values of sequence length and GC content, the sequence quality was set to decrease linearly from 32 to 28. The total number of sequences was chosen to be lower to reduce due to computational memory, 10,000 to be precise.

Figure 2 shows plots that exhibit two distributions of GC content across sequencing reads. The green line represents the GC content distribution from the simulated dataset, while the red line corresponds to the GC content of the actual experimental data. The shape and width of the peaks provide additional insights. A sharp, narrow peak could imply a uniform GC content across the reads, while a broader peak may indicate greater variability. In this instance, both lines show a significant degree of variability, with the experimental data (red line) exhibiting a wider spread, which is typical for biological variability and sequencing bias caused by rRNA contamination. The simulated data (green line) appears to have a slightly tighter distribution, hinting at a more controlled, but less varied, GC composition within the simulated sequences.

Figure 3 shows a MultiQC Mean Quality Scores plot. The two lines are each representing quality scores across the sequence read positions: the straight line represents the experimental data when the reading quality was steady at 30. The simulator parameters were set so that the quality score gradually decreased with a mean value of precisely 30.

Figure 4 (located on the next page) shows two Per Base Sequence Content charts. For the experimental data, labeled "SRR10556741", there are noticeable fluctuations in the percentage of each nucleotide (T, C, A, G) at the start of the reads. Such variations are typical in real sequencing data due to random sequencing errors or biases introduced during library preparation, such as ligation or PCR amplification artifacts. However, these fluctuations are not taken into account in the real evaluation and are therefore not included in the simulation. The chart simulated data, labeled "output," presents a more uniform distribution of nucleotides across all positions, which is characteristic for simulated data or a well-controlled sequencing process. Although the simulated data lack the initial variability observed in the experimental data, the consistent representation of nucleotide percentages throughout the read length suggests that the simulation algorithm quite effectively models the base composition of the transcriptome. However, replicating the balance of nucleotide frequencies that occurs in real biological samples is challenging.



Fig. 2. Distributions of GC content across sequencing reads for experimental (red line) and simulated (green line) data.



Fig. 3. Quality scores plot for both experimental (straight line) and simulated data (descending line).

IV. CONCLUSION

The main theme of this paper was to develop a tool that allows the simulation of RNA-Seq data taking into account the specificities of the bacterial genome, in particular the presence of operons. The BacSeqer includes several functions allowing the simulation of just such data, while the input formats are widely available and common. The FASTA file as output of the simulator was compared with the results of a real experiment using the FastQC tool, both datasets contained sequences of *Clostridium beijerinckii*. In summary, we designed a simulator that is a simple and straightforward tool for simulating bacterial RNA-Seq reads, allowing the user to set custom simulation parameters. In the future, it would be desirable to make the simulator a package that allows calls from the command line.

BacSeqer is freely available from GitHub: https://github.com/adelafi/BacSeqer.git.



Fig. 4. Per Base Sequence Content for both experimental ("SRR10556741") and simulated data ("output").

REFERENCES

- [1] WANG, Zhong; GERSTEIN, Mark a SNYDER, Michael. RNA-Seq: a revolutionary tool for transcriptomics. Online. Nature Reviews Genetics. Available: https://doi.org/10.1038/nrg2484.
- [2] FRAZEE, Alyssa C.; JAFFE, Andrew E.; LANGMEAD, Ben a LEEK, Jeffrey T. Polyester: simulating RNA-seq datasets with differential transcript expression. Online. Bioinformatics. Available: https://doi.org/10.1093/bioinformatics/btv272. Version: 1.38.0.
- [3] GRIEBEL, Thasso; ZACHER, Benedikt; RIBECA, Paolo; RAINERI, Emanuele; LACROIX, Vincent et al. Modelling and simulating generic RNA-Seq experiments with the flux simulator. Online. Nucleic Acids Research. Available: https://doi.org/10.1093/nar/gks666.
- [4] RNASeqReadSimulator [software]. Available from: http://alumni.cs.ucr.edu/ liw/rnaseqreadsimulator.html#history. Version: 2012.04.30.
- [5] BROOKS, Thomas G.; LAHENS, Nicholas F.; MRČELA, Antonijo; SARANTOPOULOU, Dimitra; NAYAK, Soumyashant; NAIK, Amruta; SENGUPTA, Shaon; CHOI, Peter S.; GRANT, Gregory R. BEERS2: RNA-Seq simulation through high fidelity in silico modeling. Online. PubMed Central. Available: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC10168222.
- [6] SEDLAR, Karel; KOLEK, Jan; GRUBER, Markus; JURECKOVA, Katerina; BRANSKA, Barbora et al. A transcriptional response of Clostridium beijerinckii NRRL B-598 to a butanol shock. Online. Biotechnology for Biofuels. Available: https://doi.org/10.1186/s13068-019-1584-7.
- [7] FastQC: A quality control tool for high throughput sequence data. Available: https://www.bioinformatics.babraham.ac.uk/projects/fastqc/.
- [8] Philip Ewels, Måns Magnusson, Sverker Lundin, Max Käller, MultiQC: summarize analysis results for multiple tools and samples in a single report, Bioinformatics. Available: https://doi.org/10.1093/bioinformatics/btw354.

Laboratory database for bacterial sample data cataloguing in the hospital environment

David Podrazký

Department of Biomedical Engineering, Faculty of Electrical Engineering and Communication Brno University of Technology Brno, Czech Republic 240549@vutbr.cz

Abstract—The proper management of the results of typing methods for bacterial classification is crucial for monitoring the spread of infections in the hospital environment. The solution for their management is a relational database. The problem is the storage of large amounts of multidimensional data due to the variability of typing methods and the need to adapt the working of the database to the requirements of a specific healthcare facility. This paper introduces the design and implementation of a relational database for the Centre of Molecular Biology and Genetics at the University Hospital Brno. It discusses the key steps in database development, including data structuring, table and relation creation, and user interface design. The use of the database has significantly improved the management and accessibility of typing data within the department.

Index Terms—Database, bacterial typing, healthcareassociated infections, MS Access

I. INTRODUCTION

Healthcare-associated infections (HAIs) are diseases that are associated with patients staying in medical facilities. They are defined as infections that are not present or incubated at the time of patient admission and do not manifest until 48 hours or more after admission. These diseases can be caused, for example, by invasive surgery, the use of medical devices (e.g. catheters), transmission between patients and medical staff, etc. [1] The European Centre for Disease Prevention and Control reports that up to 3.5 million patients in the European are infected with HAIs annually, resulting in more than 90,000 deaths. [2]

Identifying the origin of the infection is important to prevent the spread of HAIs and to establish appropriate treatment, thereby reducing the risk of prolonged hospitalization and increasing the cost of treatment. [1] Typing tests, which are used to distinguish bacteria at the species level, are an essential tool for identifying HAIs. [3]

Typing methods and clinical testing present a significant challenge because they generate large amounts of variable and often incompatible data that are specific to different bacterial species. Commonly used hospital databases, which are primarily patient-oriented, offer a limited view of the spread of infections. Therefore, creating a bacteria-specific database and typing results is crucial for detailed mapping and a deeper understanding of the dynamics of infection spread in the hospital environment.

II. BACTERIAL TYPING

Typing methods can be divided into phenotyping and genotyping. Genotyping methods examine the genetic material while phenotyping methods examine the products of gene expression. Genotyping can be further divided into sequencing and non-sequencing. [3]

A. Phenotyping methods

Phenotyping methods are older than genotyping methods. Their main disadvantages are time-consuming and low discriminatory power within bacterial species. Individual phenotyping methods differ from each other according to the substance they analyze. [3]

These include five basic methods. Biotyping uses biochemical and physiological tests. Phage typing monitors the response of bacteria to bacteriophages. Proteomic typing analyzes bacterial proteins. Serotyping uses fluorescently labeled antibodyantigen binding. [3] Antibiogram examines the resistance of bacteria to antibiotics. [4]

B. Non-sequencing genotyping methods

Most non-sequencing methods are based on PCR and electrophoresis. The simplest method is Restriction Fragment Length Polymorphism (RFLP). [3] The DNA is amplified, and split by restriction enzymes, and the fragments evaluated by electrophoresis. Pulsed-field gel electrophoresis (PFGE), which uses multiple electrodes and allows the movement of fragments in multiple directions. [5]

Other methods are also evaluated by electrophoresis, but differ in which fragments are selected by primers for amplification. [3] Random Amplified Polymorphic DNA (RAPD) uses random primers of 10 bp in length. [6] Amplified fragment length polymorphism (AFLP) uses custom adaptors that are complementary to the primers used and attach after cleavage by restriction enzymes. [7] Repetitive PCR uses primers that are identical to repetitive genome sequences. [8] Variable number tandem repeats (VNTR), on the other hand, is a method that has repetitive sequences called tandems as the target of amplification, and primers are chosen to fit in front of the tandems. [9]

C. Sequencing genotyping methods

Sequencing methods typify bacteria based on knowledge of the order of nucleotides in the DNA. [3]

Multi-Locus Sequence Typing (MLST) uses housekeeping genes for analysis. Typically, 7 genes are sequenced in one run and assigned a number based on the nucleotide order, and the combination of these numbers forms the resulting sequence type for the bacterium. [3] [10]

Mini-MLST is a method that builds on the previous method. Again, we examine housekeeping genes, this time applying high-resolution melting curve analysis (HRM) to them, which determines at what temperature the bonds between the strands of the DNA double helix are broken. [11] Again, each gene is assigned a number and the entire bacterium is assigned a Melt type. [3]

Whole-genome sequencing enables typing, thanks to nextgeneration sequencing methods that can read the entire genome of a bacterium due to the large number of parallel reads. [12] Using bioinformatic analysis, we then look for characteristic differences in them. [3]

III. DEMANDS ON THE DATABASE

In bacterial typing, the processing and storage of large amounts of multidimensional data is a major problem. These arise due to the large variability of typing methods. This can be specifically illustrated in Table I, which gives an overview of the basic tests for *Klebsiella pneumoniae* and *Staphylococcus aureus*. The point is that not all tests are performed for every bacterium, e.g. the test for enterotoxins. And the tests that are done on multiple bacteria have different attributes that they look at. For example, in MLST, we typify different bacterial species according to different genes.

A suitable solution for managing this data is a bacterialevel relational database. The main advantages of a database over spreadsheet software are lower memory requirements, and elimination of data redundancy while maintaining a tabular view of the data. It must be taken into account that it is not only enough to access the data but also to manage it and be able to manage the database development interface. This must be dealt with at the user level to reflect the specific needs of the healthcare facility. [13]

Therefore, MS Access was chosen for the database development. This offers simple integrated tools for designing tables, forms, queries, or print reports. And, in combination with the Visual Basic for Applications (VBA) programming language and SQL, it can also implement more complex database operations. So, once a database is complete, its report can be provided by a user without in-depth knowledge of database systems. [14]

IV. DATA STRUCTURE

The data received from the Centre for Molecular Biology and Genetics (CMBG) reflect the functioning of the department and can be divided into four groups according to the subject they describe within the laboratory process.

 TABLE I

 EXAMPLE OF TYPING METHODS FOR Klebsiella pneumoniae AND

 Staphylococcus aureus

	Klebsiella pneumoniae	Staphylococcus aureus			
Antibiogram	21 types of antibiotics	12 types of antibiotics			
	gapA	arcC			
	infB	aroE			
	mdh	glpF			
MI ST	pgi	gmk			
WILST	phoE	pta			
	rpoB	tpi			
	tonB	yqiL			
	ST	ST			
	infB729	arcC78-210			
	mdh1197	aroE88-155			
	phoE2013	gmk286			
Mini-MLST	rpoB2227	pta294			
	tonB2693	tpi36			
	tonB2886	tpi241-243			
	MelT	MelT			
	ybtS				
	mrkD				
	entB				
	rmpA				
Multiplex PCR	К2	NA			
	kfu				
	allS				
	iutA				
	magA				
WGS	Illumina	Illumina			
	Nanopore	Nanopore			
Sna-HRM	NA	Alel's lenght			
opa-mon	1474	Spa-HRM type			
PVL/STS	NA	PVL			
1 11/010	1474	STS			
		sea			
		seb			
		sec			
		sed			
Enterotoxins	NA	see			
		femA			
		mecA			
		eta			
		etb			
		tst			
Spa typing	NA	Spa type			

The first group, patient informations, contains basic biological and identification data of patients. The second group focuses on collected samples, capturing details of the biological material and its collection. The third group, informations about identification of cultured bacteria. The last group of informations focuses on typing test results, where the database holds data on 7 bacterial species and 10 typing tests in 27 variants, depending on the bacterial type. Patient and sample informations is included in the request forms from external medical facilities that request testing from CMBG. Informations on bacteria and test results is generated directly within the department.

Part of the data is covered by dials. This is data that does not change in the long term and helps to maintain data integrity. The specific dials used in the database store informations about external requesting devices, biological materials, and known sequencing types for MLST and mini-MLST typing methods.



Fig. 1. Relational scheme

V. RELATIONAL SCHEME

The design of the relational schema was based on the data structure, specific CMBG requirements, and the expected user access. The relational schema is shown in Fig. 1

The tables in the database are divided into three categories: main tables, dials, and supplementary tables that solve the conditional relation problem.

The main tables are for storing the informations described in the previous sections, such as data about patients, collected samples, cultured bacteria and typing test results.

When accessing the database, the user first enters informations about the biological material into the **Samples** table. Integrity constraints check for the existence of a patient record. If the patient is not in the database, the patient informations must be added to the **Patients** table. There is a 1:N relation between the **Patients** and **Samples** tables, allowing the typing of multiple types of biological material for a single patient. When a sample is added, the informations about the cultured bacteria is entered into the **Typization** table with a 1:N relation because multiple bacteria can be isolated from a single sample. Information on the tests required is entered into the **Request list** table with a 1:N relation, as multiple tests can be performed on a single bacterium.

The relation between the **Typization** table and the **Test's** tables is essential to solve the problem of data variability. This relation is defined as a conditional 1:0 or 1:1 relation, reflecting the fact that not all possible tests may be performed on one bacterial type. The establishment of a typing-testing relation occurs only when a request for a specific test for a given sample is recorded in the **Request list** table.

		MA		ME	NU			
DATA ADDITION	DATA PREVI	EW		IMPORT		EXPC	ORT DATA	
Add patient	Patient det	ail List of pat	tients (Sequenc	e type dial	Sp	readsheets	
Add sample	Sample det	ail List of sar	nples	Melt t	ype dial	Pr	int reports	
TEST RESULTS Klebsiella pneumoniae	Request list Staphylococcus aureus	Escherichia coli	Enteroco faeci	occus ium	Clostridioi difficile	des 9	Pseudomonas aeruginosa	Serratia marcescens
Phenotype KLPN	Phenotype STAU	Phenotype ESCO	Phenotype	e ENFA	Phenotype (CLDI	Phenotype PSAE	Phenotype SEMA
MLST KLPN	Spa-HRM	MLST ESCO	MLSTE		Ribotypin	ıg	MLST PSAE	WGS SEMA
Minim KLPN	MLST STAU	Minim ESCO	WGS EI	NFA	WGS CLE	וכ	Minim PSAE)
Multiplex PCR	Minim STAU						WGS PSAE	5
WGS KLPN	PVL/STS							
	Enterotoxins							
	WGS STAU							
	Spa typing							

Fig. 2. Main menu GUI

VI. DATABASE IMPLEMENTATION

Based on the designed relational schema, the main tables were created and the relations between them were defined. A user interface was created to simplify access to the data.

In MS Access, the user interface consists of forms. These forms contain text fields based on the tables and control elements. You can select preset functions for individual elements or program your own. [14] An example of the user interface is shown in Fig. 2

Preview forms have been created for each of the main tables. These forms are divided into two basic types: list forms and detail preview forms. List forms provide an overview of all records with only basic attributes. Detail view forms show only one record but with all attributes. Forms for adding and editing records are similar in structure to the detailed view forms but differ in that they allow data entry and editing, while the preview forms have user interface elements locked.

The data export forms allow the generation of tabular datatypes. The user can specify parameters such as date, bacteria type or test type to determine exactly which data they want to export. In addition to the ability to export data, report printing offers additional features such as: basic statistical analysis, a list of bacteria waiting to be tested for each examination, or for MLST and mini-MLST examinations, a list of unknown sequence types not yet in the dials.

Key functions in the database had to be programmed in VBA, beyond the basic MS Access environment. The development of a system of conditional relations between the **Typization** table and the **Test's** tables and the implementation of a system for exporting data according to user-defined parameters were essential. Other advanced features included the generation of primary keys, filtering in combined fields at the form/subform level, dynamic display of examination requests and import of dials.

VII. CONCLUSION

This paper presents the design and implementation of a comprehensive bacterial-oriented clinical database for multidimensional typing of bacterial pathogens. In the development process, a database structure was created to meet the specific needs and functioning of the department, involving the design and creation of tables and their interrelationships. In addition, a user interface was developed, which was subsequently enhanced with additional functionality beyond the basic MS Access offering through VBA programming.

Currently, the database supports a structure for storing informations on the seven most important bacterial pathogens occurring in the hospital environment, namely *Klebsiella pneumoniae*, *Staphylococcus aureus*, *Escherichia coli*, *Enterococcus faecium*, *Clostridioides difficile*, *Pseudomonas aeruginosa* and *Serratia marcescens*. For future development, it is possible to expand the database to include more bacterial species, which would allow it to cover a broader spectrum of pathogens and thus improve its usability in practice. Another possible extension is the introduction of direct import of results from online tools. For example, to find resistance and virulence genes.

The proposed database represents a key tool for more accurate mapping of the occurrence and transmission of bacterial infections in the hospital environment, providing valuable informations to healthcare professionals for better diagnosis, prevention, and treatment of infectious diseases.

ACKNOWLEDGMENT

I would like to express my sincere gratitude to my supervisor, Ing. Helena Vítková, Ph.D., for her invaluable help, advice, and guidance throughout this project.

REFERENCES

- [1] VOIDAZAN, Septimiu; ALBU, Sorin; TOTH, Réka; GRIGORESCU, Bianca; RACHITA, Anca et al., 2020. Healthcare Associated Infections—A New Pathology in Medical Practice? Online. International Journal of Environmental Research and Public Health. Roč. 17, č. 3. Available at: https://doi.org/10.3390/ijerph17030760.
- [2] EUROPEAN CENTRE FOR DISEASE PREVENTION AND CONTROL, 2021. Healthcare-associated infections. Online. European Centre for Disease Prevention and Control. Available at: https://www.ecdc.europa.eu/en/healthcare-associated-infections.
- [3] RAMADAN, Asmaa A. Bacterial typing methods from past to present: A comprehensive overview. Gene Reports [online]. 2022, 29, 1-13 [cit. 2023-11-28]. ISSN 24520144. Available at: https://www.sciencedirect.com/science/article/pii/S2452014422001832
- [4] TRUONG, William R; HIDAYAT, Levita and BOLARIS, Michael A, 2021. The antibiogram: key considerations for its development and utilization. Online. JAC-Antimicrobial Resistance. 2021-06-01, roč. 3, č. 2, s. 1-9. Available at: https://doi.org/10.1093/jacamr/dlab060.
- [5] NEOH, Hui-min; TAN, Xin-Ee; SAPRI, Hassriana Fazilla a TAN, Toh Leong, 2019. Pulsed-field gel electrophoresis (PFGE): A review of the "gold standard" for bacteria typing and current alternatives. Online. Infection, Genetics and Evolution. Roč. 74. ISSN 15671348. Available at: https://doi.org/10.1016/j.meegid.2019.103935. [cit. 2023-12-15].
- [6] BABU, Kantipudi Nirmal; SHEEJA, Thotten Elampilay; MINOO, Divakaran; RAJESH, Muliyar Krishna; SAMSUDEEN, Kukkamgai et al., 2021. Random Amplified Polymorphic DNA (RAPD) and Derived Techniques. Online. Molecular Plant Taxonomy. Methods in Molecular Biology. S. 219-247. ISBN 978-1-0716-0996-5. Available at: https://doi.org/10.1007/978-1-0716-0997-2-13. [cit. 2024-01-03].
- [7] PAUN, Ovidiu a SCHÖNSWETTER, Peter, 2012. Amplified Fragment Length Polymorphism: An Invaluable Fingerprinting Technique for Genomic, Transcriptomic, and Epigenetic Studies. Online. Plant DNA Fingerprinting and Barcoding. Methods in Molecular Biology. S. 75-87. ISBN 978-1-61779-608-1. Available at: https://doi.org/10.1007/978-1-61779-609-8-7. [cit. 2024-01-03].
- [8] DOMBEK, Priscilla E.; JOHNSON, LeeAnn K.; ZIMMERLEY, Sara T. a SADOWSKY, Michael J., 2000. Use of Repetitive DNA Sequences and the PCR To Differentiate Escherichia coli Isolates from Human and Animal Sources. Online. Applied and Environmental Microbiology. Roč. 66, č. 6, s. 2572-2577. ISSN 0099-2240. Available at: https://doi.org/10.1128/AEM.66.6.2572-2577.2000. [cit. 2024-01-03].
- [9] MARSHALL, Jack NG; LOPEZ, Ana Illera; PFAFF, Abigail L; KOKS, Sulev; QUINN, John P et al., 2021. Variable number tandem repeats – Their emerging role in sickness and health. Online. Experimental Biology and Medicine. Roč. 246, č. 12, s. 1368-1376. ISSN 1535-3702. Available at: https://doi.org/10.1177/15353702211003511. [cit. 2024-01-03].
- [10] IBARZ PAVÓN, Ana Belén a MAIDEN, Martin C.J., 2009. Multilocus Sequence Typing. Online. Molecular Epidemiology of Microorganisms. Methods in Molecular Biology. S. 129-140. ISBN 978-1-60327-998-7. Available at: https://doi.org/10.1007/978-1-60327-999-4-11. [cit. 2023-12-15].
- [11] ANDERSSON, Patiyan; TONG, Steven Y. C.; BELL, Jan M.; TURNIDGE, John D.; GIFFARD, Philip M. et al., 2012. Minim Typing – A Rapid and Low Cost MLST Based Typing Tool for Klebsiella pneumoniae. Online. PLoS ONE. 2012-3-12, roč. 7, č. 3. ISSN 1932-6203. Available at: https://doi.org/10.1371/journal.pone.0033530. [cit. 2023-12-15].
- [12] JOENSEN, Katrine Grimstrup; SCHEUTZ, Flemming; LUND, Ole; HASMAN, Henrik; KAAS, Rolf S. et al., 2014. Real-Time Whole-Genome Sequencing for Routine Typing, Surveillance, and Outbreak Detection of Verotoxigenic Escherichia coli. Online. Journal of Clinical Microbiology. Roč. 52, č. 5, s. 1501-1510. ISSN 0095-1137. Available at: https://doi.org/10.1128/JCM.03617-13. [cit. 2023-12-15].
- [13] PASTORINO, Roberta; DE VITO, Corrado; MIGLIARA, Giuseppe; GLOCKER, Katrin; BINENBAUM, Ilona et al., 2019. Benefits and challenges of Big Data in healthcare: an overview of the European initiatives. Online. European Journal of Public Health. 2019-10-01, roč. 29, č. Supplement-3, s. 23-27. Available at: https://doi.org/10.1093/eurpub/ckz168.
- [14] GEORGE, Nathan, 2022. Mastering Access 365: An Easy Guide to Building Efficient Databases for Managing Your Data. GTech Publishing. ISBN 978-1-915476-00-5.

Four Channel Active Antenna Switch for UHF Band Satellite Reception

Václav Kubeš Department of Radio Electronics FEEC, Brno University of Technology Czech Republic 240644@vutbr.cz

Abstract—The aim of this work is to design and implement a complex electronic device, known as an antenna switch, that allows to receive a high-frequency satellite signals in the UHF band via one of five connected antennas at a given time. This device is capable of automatic antenna switching according to non-geostationary satellite position, from which signal is received, thus with suitable antenna selection, is alternative to motorized directional antennas. The device consists of outdoor and indoor unit and Human-Machine interface (HMI). The outdoor unit is driven by a microcontroller which collect diagnostic data about proper function and controls the RF switch. The outdoor unit amplifies the UHF signal too. indoor unit with power inserter, diplexer and USB to serial converter, allows diagnostic data and commands transfer between outdoor unit and users computer with Human-Machine interface.

Index Terms—Antenna switch, RF switch, RF circuit, satellite, UHF, outdoor unit, indoor unit, signal reception

I. INTRODUCTION

In recent decades, there has been and is still a rapid development of wireless communication. Its possibilities of use are enormous. A cascade of electronic devices that enable wireless communication is very long, but among its most basic building blocks, we can include an antenna, without which a signal would not be possible to transmit and to receive and antenna switch. Antenna switch enables flexibility of choice of connected antennas or even choice whether the antenna will be used for transmitting or receiving. This functionality is for example used in smartphones with Wi-Fi and Bluetooth connectivity for which is only one common antenna used, which greatly simplifies miniaturization of the device.

In case of this work, an antenna switch is designed, that it will enable signal receiving from small satellites that are on non-geostationary orbit. This means, that the satellite is moving relative to the reception point. In order to maintain good reception of the weak signal from small satellite, whose orbit is usually in 350 - 700 km altitude [1] and maximum transmission power is 30 dBm [2], directional antennas are necessary.

Motorization of the directional antenna enables to follow the movement of satellite so that main lobe of the directional antenna is heading to the signal source, thus the biggest gain is achieved. Another way of satellite position following is gradual switching of antennas. Antennas have to be carefully Tomáš Urbanec Department of Radio Electronics FEEC, Brno University of Technology Czech Republic urbanec@vut.cz

oriented, so that in certain time, signal is received by antenna, whose main lobe orientation corresponds to current satellite position so that the biggest gain is ensured. The main focus of this work is the design and implementation of such automatic antenna switch, which enables preamplification of the received signal and enables diagnostic data of the device collection.

II. DEVICE REQUIREMENTS

The frequency of the received signal was chosen to be 435 MHz, as it corresponds to the amateur frequencies used to communicate with CubeSats [3]. Another requirement for the device is diagnostic data about working state collection, so that user can be warned if some error, such as RF amplifier failure or low power voltage, occurs. The next requirement is to design the device in such way, that by assembling some components, second variant of the device arises and by connection of these two variants, switching of both polarizations, horizontal and vertical, can be performed. This requirement should lead to lower expenses of manufacturing, because only one printed circuit board (PCB) is needed and the number of components is minimized.

III. CONCEPTUAL DESIGN

The device itself have to be divided to two parts: external unit and internal unit. The external unit will be placed near to the antennas and its purpose is as follows:

- Preamplification of the received signal
- Signal filtration
- Signal from antennas switching
- Signal diplexing
- Diagnostic data collection and transmission
- Ensuring a stable power supply
- Command reception

The internal units purpose is:

- Diagnostic data reception and command transmission
- Power insertion
- Signal diplexing

The simplified diagram of the resulting device is depicted in Fig. 1. As diagnostic data, current to RF amplifiers, which carries information about their proper function, phantom power voltage, vital for powering the whole outdoor unit, temperature and azimuth of the outdoor unit, should be measured. For power savings, the power supply for each RF amplifier should be switchable, as their current draw is usually high.



Fig. 1. Simplified diagram of designed device.

IV. CIRCUIT DESIGN AND COMPONENT SELECTION

A. RF switch

Because the heart of this device is the switching element itself, it is necessary to pay great attention on its selection. It can affect the device behavior and the quality of RF signal. There are three main options when selecting the RF switch: mechanical RF relay, MEMS RF switch and semiconductor RF switch. Each of them has its own advantages and disadvantages.

The integrated semiconductor RF switches are very easy to use and require small number of components to work. They offer small insertion losses, fast switching times and high range of switchable frequencies. The ease of use, good qualities of the integrated semiconductor RF switches and great selection, led to choosing of this type of RF switch to be used in designed device.

From the range of integrated RF switch on the market, AS195-306LF [4] from Skyworks was chosen. The main criteria for this selection were the capability of switching up to five channels (ST5P configuration), 5 V control voltage and possibility of manual soldering and that all while maintaining the lowest insertion loss. The main specifications of AS195-306LF as well as other possible integrated RF switches, which meet the requirements are in Table I.

 TABLE I

 RF SWITCHES WHICH MEET THE REQUIREMENTS^a

RF Switch	Max. Inser- tion loss [dB]	Max. Swit- ched power [dBm]	Freq. Range [GHz]	SMD Package
AS195-306LF	0.5	27.0	0.1 - 2.0	QFN-16
ADRF5250	1.3	33.0	0.1 - 6.0	QFN-24
HMC252AQS24E	0.8	29.8	DC - 3.0	SOP-24
SKY13415-485LF	0.4	37.5	0.1 - 3.8	QFN-14
F2915	0.93	37.0	0.05 - 5.0	QFN-24

^aValues taken from [4] [5] [6] [7] [8], at frequencies close to 435 MHz.

B. Other components

As the microcontroller for the outdoor unit, ATmega328PB [9] was selected because of the simplicity of use and vast user base which leads to big amount of available information. ATmega328PB has I2C and UART peripherals needed for communication between sensors and indoor unit as well as 10 bit A/D converter which is needed for diagnostic data collection.

MCP9808 [10] digital temperature sensor is used for measuring temperature of the surroundings of outdoor unit and outdoor unit itself. This sensor communicates with microcontroller via I2C. Module of digital compass [11] with HMC5883L integrated circuit is used for antenna orientation evaluating, and it communicates with microcontroller via I2C too. The main purpose of this sensor is the antenna's azimuth detection after it is installed so that the right order of switching is achieved. It is also used to detect displacement of the antenna.

The communication between outdoor unit and indoor unit needs to be at least half duplex and must be feasible on one line, thus can be diplexed with RF signal to one coaxial cable so that the installation of the outdoor unit is as convenient and easy as possible. LIN (Local Interconnect Network) bus is well suited for this. In this application only physical layer is used and classical serial protocol with 9600 baud speed is used. The transmitting and receiving on the physical layer are ensured by TJA1020 LIN transceiver integrated circuit [12].

Integrated circuit INA1802A was chosen for measurement of current to RF amplifiers. It is integrated differential amplifier and its amplification is 50 [13]. For phantom voltage stabilization was selected LM2940 low-dropout voltage regulator because it is capable to supply current up to 1 A. Its output voltage is 5 V and it can regulate input voltages in range from about 6 V to 28 V [14]. With that in mind the range of supply voltage was selected from 7 V to 12 V. The higher limit was set to maintain the power losses on the voltage regulator low so that it would not overheat.

C. Circuit design

The input part of the RF circuit and the RF amplifier with stabilisation circuit was taken from [4] and was adjusted for this circuit design (see Fig. 2). Antennas can be connected to the outdoor unit via SMA connectors. The input part which is resonant circuit used as input band pass filter, is used for filtering other close signals in this 435 - 438 MHz frequency band [15]. PGA103+ was chosen as the low noise RF amplifier (LNA). Its typical gain is 22.1 dB (@400 MHz) and typical noise figure is 0.5 dB (@400 MHz). Its supply voltage is 5 V and typical current consumption is 97 mA [4]. All the high frequency traces on final PCB needed to be done by microstrips, which ensures the impedance of 50 Ω .

Power switching for each LNA is provided by a SSM3J355R unipolar P-MOSFET transistor. The power is switched in high-side configuration, so the unexpected behaviour of ungrounded LNA is prevented. That is why the P type of MOSFET is used. It also means, that the powering of



Fig. 2. Input part with resonant circuit and the low noise amplifier.

the LNA has to be driven from the microcontroller by negative logic. The voltage from the drain of the unipolar transistor is also used as driving voltage of the RF switch AS195-306LF, so the LNA and corresponding input of the RF switch are switched at the same time by one control signal from the microcontroller.

After the RF switch there is surface acoustic wave filter (SAW), which ensures filtering of only narrow band of frequencies. To ensure the lowest noise figure (NF) of antenna switch circuit, there is another LNA after SAW filter. By adding this second amplifier the noise figure at the output of internal unit drops from 1.72 dB to 1.04 dB.

To combine received RF signal and LIN communication, there is a diplexer at the output so that one coaxial cable can be used to route the signals to the indoor unit. There is also possibility to diplex GPS signal from active GPS antenna and to separate the phantom power.

To measure the current to LNAs, the INA180A2 is used as mentioned above. The amplified differential voltage is created on sensing resistors through which the current flows. There are used two 1 Ω resistors in parallel to ensure small voltage drop so that the power voltage at the LNAs is not highly affected by it, even when high current is drawn.

The phantom voltage is measured using a resistive divider, so the voltage is in safe range for microcontrollers A/D converter, which should use an internal 1.1 V voltage reference [9]. This ensures that the measured voltage is not dependent on the potentially volatile power voltage of the microcontroller so that the measurements are accurate.

D. Possible assembling variants

As was written above, the circuit was devised so that two variants (variant A and variant B) of assembly of designed PCB are possible. By connecting these two variants via pin header or ribbon cable, device which is capable of switching signals from five antennas in both polarization (horizontal and vertical polarization of one antenna simultaneously) arises.

While variant A includes microcontroller wiring and its purpose is the switching control and communication with internal unit, the B variant ensures regulating the incoming phantom powering and it adds feature of connecting active GPS antenna and diplex its signal into one coaxial cable together with RF signal from antennas. The variant A can be independent unit and for its working only external powering is needed. The resulting PCB design can be seen in Fig. 3.



Fig. 3. Designed custom PCB for outdoor unit.

E. Indoor unit

The main purpose of the indoor unit is to split the diplexed signals from the outdoor unit. There are two RF inputs, for A and B variant of the outdoor unit. On one input there is diplexer which splits LIN communication and the RF signal coming from the outdoor unit assembled as variant A. On the second input there is diplexer splitting incoming RF signal and GPS signal from the variant B. There is also connected the phantom power inserter to this part of the indoor unit. Separated RF signals are fed to the SMA connectors so that software defined radio (SDR) receiver can be connected.

Command from the HMI running on the user's computer and diagnostic data from the outdoor unit are exchanged throughout USB/Serial converter, which is present on the indoor unit. Integrated circuit MCP2200 [16] was used for this and to connect to physical layer of LIN bus, TJA1020 is used in the same way as on the outdoor unit.

V. HUMAN-MACHINE INTERFACE: CONTROL APPLICATION FOR USER'S COMPUTER

To enable easy interfacing with basic settings of the antenna switch and viewing the diagnostic data, application for user's computer was written in Python programming language. The screenshot of antenna switch control application is depicted in Fig. 4.

🍄 Antenna switch controller	- 🗆 X				
Antenna switch controller					
Satellite position	Status info				
Azimuth: 267.8 °	Phantom voltage [V]: 11.0				
Elevation: 7.7 °	Current to LNAs [mA]: 251.0 326.0				
Antenna in use:	Temperature [°C]: 17.0 15.0				
1	Ant. orientation measured/set: 87.0 ° 276.0 °				
	Warning!				
Antenna switch	Manual ant. Set current orientation set. ant. orientation				
O Manual:	Serial port:				
Ant. 1	COM6 ~				
Ant. 3	App settings				
Ant. 5	App info				
Orbitron: Connected	Switch: Connected				

Fig. 4. Screenshot of created antenna switch control app.

The application communicates with the outdoor unit via the USB connected to indoor unit with USB/Serial converter and LIN transceiver. User can choose between manual antenna switching, when the user selects the switched-on antennas by clicking on buttons and auto switching, when the antennas are switched according to satellite position.

To ensure switching of the antennas according to selected satellite, data about its current position must be available. For this reason, Orbitron, satellite position prediction software [17], is used. By connecting to its DDE server, data about selected satellite are passed to the application and corresponding antennas are selected automatically.

The created application also displays received diagnostic data and notifies the user, when they are not in normal working range. Diagnostic data are also saved to file when an error occurs.

VI. CONCLUSION

This article presents the design and implementation of a device that enables the switching of high-frequency signals from satellites received by five antennas and provides their pre-amplification using LNAs and it provides filtration of the signal too. The device also allows for simple diagnostics by measuring the current drawn by the LNAs, the phantom power supply voltage, and the temperature of the device. Control and data transfer are ensured via the physical layer of the LIN communication bus.

Furthermore, the indoor unit, which consists of a diplexer with a USB/UART converter, was designed, which split the incoming high-frequency signal from the switch. It also serves to supply the phantom power to the switch and allows communication between the switch and the user's computer.

Suitable components (microcontroller ATmega328PB, temperature sensor MCP9808, differential amplifier INA180A2, LDO regulator LM2940...) were selected for the device with regard to the requirements. Great attention was paid to the selection of the switching element itself. After a thorough study of the high-frequency switch issue, the integrated circuit AS195-306LF from Skyworks was chosen.

The circuit diagrams of the outdoor unit and indoor unit were designed in Autodesk Eagle, where printed circuit boards (PCBs) were designed too. The outdoor unit consists of two combined circuits for switching vertical and horizontal polarization of the received signal. This combination allowed for the design of a single PCB for both variants. The variant can be selected when assembling components. By connecting these two variants together, a fully functional device arises. This, among other things, reduces production costs.

By the time of submission of this paper, indoor and outdoor unit PCBs has been manufactured and partly assembled. Now it is worked on the tuning of input resonant circuit for highest gain and the lowest noise figure and on firmware for the microcontroller of outdoor unit. Both of the activities are expected to be finished very soon.

For user control, a simple application for the computer was created. It enables automatic switching according to satellite position as well as manual selection of switched-on antennas and it displays the incoming diagnostic data from outdoor unit.

References

- H. Polat, J. Virgili-Llop, and M. Romano, "Survey, statistical analysis and classifica- tion of launched cubesat missions with emphasis on the attitude control method," *JoSS*, vol. 5, pp. 513–530, 2016. [Online]. Available: https://www.jossonline.com/wp-content/uploads/2016/10/Fin al-Survey-Statistical-Analysis-and-Classification-of-Launched-CubeS at-Missions-with-Emphasis-on-the-Attitude-Control-Method3.pdf
- [2] N. Saeed, A. Elzanaty, H. Almorad, H. Dahrouj, T. Y. Al-Naffouri, and M.-S. Alouini, "Cubesat communications: Recent advances and future challenges," *IEEE Communications Surveys Tutorials*, vol. 22, no. 3, pp. 1839–1862, 04 2020.
- [3] Český telekomunikační úřad, "Pásmo 432 438 mhz detail kmitočtového pásma," spektrum.ctu.gov.cz, 01 2024. [Online]. Available: https://spektrum.ctu.gov.cz/kmitocty/432-438-mhz?filter%5B frequencyFrom%5D=435&filter%5BfrequencyFromUnit%5D=MHz&f ilter%5BfrequencyTo%5D=435&filter%5BfrequencyToUnit%5D=MHz
- [4] Skyworks Solutions, Inc., AS195-306LF: PHEMT GaAs IC High-Power SP5T Switch 0.1 to 2 GHz, 06 2016. [Online]. Available: https://www.skyworksinc.com/-/media/SkyWorks/Documents/Products/ 1-100/AS195_306LF_200187D.pdf
- [5] Analog Devices, Inc., ADRF5250: 0.1 GHz to 6 GHz Silicon SP5T Switch, revision 0 ed., 2017. [Online]. Available: https://www.analog.c om/media/en/technical-documentation/data-sheets/ADRF5250.pdf
- [6] —, HMC252AQS24E: GaAs MMIC SP6T NON-REFLECTIVE SWITCH, DC - 3 GHz, v01.0316 ed. [Online]. Available: https://www. analog.com/media/en/technical-documentation/data-sheets/hmc252a.pdf
- [7] Skyworks Solutions, Inc., SKY13415-485LF: 0.1 to 3.8 GHz SP5T Antenna Switch, 11 2016. [Online]. Available: https://www.skyworksin c.com/-/media/SkyWorks/Documents/Products/701-800/SKY13415_4 85LF_201704I.pdf
- [8] Renesas Electronics Corporation, F2915 Datasheet, 08 2021. [Online]. Available: https://www.renesas.com/us/en/document/dst/f2915-datasheet
- [9] Microchip Technology Inc., ATmega328PB: AVR® Microcontroller with Core Independent Peripherals and PicoPower® Technology, revision c ed., 02 2018. [Online]. Available: https://ww1.microchip.com/downlo ads/aemDocuments/documents/MCU08/ProductDocuments/DataSheets /40001906C.pdf
- [10] —, MCP9808: ±0.5°C Max. Accuracy Digital Temp. Sensor, revision b ed., 05 2018. [Online]. Available: https://ww1.microchip.com/down loads/aemDocuments/documents/OTH/ProductDocuments/DataSheets

/MCP9808-0.5C-Maximum-Accuracy-Digital-Temperature-Sensor-Dat a-Sheet-DS20005095B.pdf

- [11] LaskaKit, "3-osý magnetometr a kompas hmc58831 laskakit," laskakit.cz, 09 2021. [Online]. Available: https://www.laskakit.cz/3-osy -magnetometr-a-kompas-hmc58831/?gclid=Cj0KCQjwmvSoBhDOARI sAK6aV7hGcETBH1c4WRcwR_PqNbA_309J4XdzoGnZQoKA_iZiQq 2aMNTGcn0aAoWmEALw_wcB#relatedFiles
- [12] NXP Semiconductors N.V., *TJA1020 LIN transceiver*, revision 5 ed., 06 2004. [Online]. Available: https://www.nxp.com/docs/en/data-sheet /TJA1020.pdf
- [13] Texas Instruments, Inc., INAx180 Low- and High-Side Voltage Output, Current-Sense Amplifiers, 07 2022. [Online]. Available: https://www.ti.com/lit/ds/symlink/ina180.pdf?ts=1698465700154&ref_u rl=https%253A%252F%252Fwww.ti.com%252Fproduct%252FINA180
- [14] —, LM2940x 1-A Low Dropout Regulator, 12 2014. [Online]. Available: https://www.ti.com/lit/ds/symlink/lm2940-n.pdf?ts=170039 9384615
- [15] Cadence System Analysis, "Rf mems switches provide superior performance over solid-state switches," resources.systemanalysis.cadence.com, 2021. [Online]. Available: https://resources.syst em-analysis.cadence.com/blog/msa2021-rf-mems-switches-provide-sup erior-performance-over-solid-state-switches
- [16] Microchip Technology Inc., MCP2200: USB 2.0 to UART Protocol Converter with GPIO, revision e ed., 06 2021. [Online]. Available: https://ww1.microchip.com/downloads/aemDocuments/documents/API D/ProductDocuments/DataSheets/MCP2200-USB-2.0-to-UART-Proto col-Converter-with-GPIO-DS20002228E.pdf
- [17] S. Stoff, "Orbitron software." [Online]. Available: http://www.stoff.pl/

YTVARET?

Objev kariérní příležitosti v největším evropském R&D centru společnosti Honeywell v České republice. Nabízíme široké možnosti uplatnění v oblasti výzkumu, vývoje a IT se zaměřením na letectví a vesmírné technologie, řešení pro zvýšení bezpečnosti práce a produktivity.



www.honeywell.com



careers.honeywell.com/Brno careers.honeywell.com/EarlyCareersCZ

THE FUTURE IS WHAT WE MAKE IT | HONEYWE

Development Module for Radar Safety Sensor in Single-Track Vehicles

Martin Ťavoda FEEC Department of Radio Electronics Brno University of Technology Brno, Czech Republic martin.tavoda@vutbr.cz

Abstract—This article describes the hardware and software design of a development module for the Radar Safety Sensor (RSS). RSS uses a Frequency-Modulated Continuous Wave (FMCW) radar to track moving objects behind single-track vehicles. The current radar configuration and antenna have significant object tracking deficiency when the vehicle makes a turn and the radar Field of view (FOV) tilts. The assumed solution is to add Inertial Measurement Unit (IMU) data to the tracking algorithm. The IMU is implemented in the designed development module, which also provides an efficient development and debugging environment.

Index Terms—embedded device, cyclo-safety equipment, ROS 2, micro-ROS, FMCW radar

I. INTRODUCTION

The designed module will serve the company ALPS ALPINE Co., Ltd. in the development of Radar Safety Sensor (RSS). Together with the ALPS Generic Radar 5, it should result in a similar device as the already existing Ride Safety System RS 1000 or GARMIN Radar Tail Light. These devices are safety equipment for cyclists that offer information about vehicles behind them. If the vehicle is approaching too quickly or it is driving too close the cyclist and also the driver are warned [1].

The module purpose is to process data from an IMU and send them to a Single Board Computer (SBC) through a Robot Operating system 2 (ROS2) domain. As an output, it collects result data from a ROS2 domain and triggers an alert. Additionally, it provides robust connectors, power management, ANT+ connection, and physical mount to singletrack vehicles, specifically a bicycle. The whole system will be supplied from an external 18 V Li-Ion accumulator.

II. PROBLEM FORMULATION

ALPS ALPINE Co., Ltd. is developing a new cyclo product using a proprietary radar system with a new tracking algorithm. However, they ran into a problem. In the default scenario, the radar is mounted under the bicycle seat and the bicycle is moving in a straight line (no tilt). The radar sees multiple points, where some of them are reflections from a static background (the blue dots in Figure 1), and the rest are reflections from actual moving targets behind the bicycle (the green dot in Figure 1).

In this scenario, these two conditions are easily distinguishable in processed radar data output and the ego-motion vector



Fig. 1. Correct tracking - straight ride

can be calculated from these reflections of static background. Ego-motion vector describes the motion of an object relative to the rigid scene [2]. In this context, it consists of two elements, direction (forward or backward), and speed (approaching or receding). In the current radar configuration, the radar FOV is only a horizontal line, there is no elevation data, and the egomotion vector is only one-dimensional. The reflections from moving targets are clustered and classified, and a tracking algorithm is used. From these data, it is possible to estimate the moving target trajectory, and speed (approaching or receding), classify if it's a car or truck, and calculate the possibility of a potential collision.



Fig. 2. Tracking degradation - ride in turn

However, if the bicycle tilts, for example, when turning, the radar FOV also tilts with the same angle and this condition cannot be detected from radar data alone. Since the egomotion vector is only a single dimensional and the height of the radar from the ground changes, the previously tracked object could be tracked incorrectly (see in Figure 2) or can disappear altogether.

A similar problem is described in [3], but the company requested a simpler approach. One possible solution to this problem seemed to be adjusting the tracking algorithm parameters to be more lenient, but this resulted in insufficient accuracy. The development team decided to try to add to the tracker algorithm some external information about radar tilt and position in space to enhance the tracking estimation. This information can be obtained from an IMU located in a development module.

After the successful implementation of IMU data into the tracking algorithm by the ALPS team they will need to easily show the performance of this product. For this purpose, the development module offers output peripherals such as highly luminous Light-Emitting Diodes (LEDs), a loud buzzer, and also in the future the ability to connect to external devices such as wireless ANT+ fitness accessories.

III. SYSTEM DESIGN

The system design is shown in Figure 3. The development module for the RSS consists of multiple blocks. The main part is the microcontroller ESP32 which collects the raw data from the IMU using the I²C bus. These data are then processed and sent to the ROS2 domain via Ethernet. Ethernet connection is provided by Media access control sublayer (MAC) and Physical layer (PHY) interfaces. MAC is included in ESP32 and PHY is an external component. They communicate with each other by Reduced Media-independent interface (RMII). Ethernet was preferred over Wi-Fi due to its lower latency. As IMU is used the MPU5060 chip and the PHY interface is provided by IP101GR Integrated circuit (IC).

The SBC connected to the ROS2 domain records data from the development module and the radar Printed Circuit Board (PCB). These data are stored in the format of ROS2 database.



Fig. 3. Development module block diagram



Fig. 4. Power management block diagram

The Alps team can afterward playback these data and use them to train the neural network that currently presents the tracking algorithm. The output (a threat level) of the finished algorithm should be sent back to the ESP32, whose task is to provide an alert by visual or sound peripherals.

As visual output the programmable RGB LEDs WS2812B mini are used. They are placed on the perimeter of the PCB and the brightness of these LEDs is adjusted by the software according to the amount of ambient light. This level is read from the external photoresistor by the Analog-to-digital converter (ADC) included in the ESP32.

The output threat level from the tracking algorithm will be also wirelessly sent to an external smart device such as a phone or commercially available fitness accessories. This communication is provided by the Wi-Fi peripheral inside ESP32 or the external ANT+ module D52 by GARMIN.

The power management block is shown in Figure 4 and it consists of four power lines. The external accumulator has a nominal voltage of U = 18 V but the radar needs to be supplied with a voltage of U = 12 V. Therefore the RT2589 DC/DC step-down converter is used. For supplying LEDs the AP63205 DC/DC converter with fixed output voltage U = 5 V was chosen. The same converter but with a different fixed output voltage of U = 3.3 V is used to supply all chips on the PCB.

Besides the fact that the USB connector serves as a programming interface for the ESP32, it also serves as a second power source. The seamless transition between 5 V converted from the battery voltage and the voltage from the USB is provided by the TPS2121 power multiplexer. When both sources are available the USB source has priority, but when it's not present the module is powered by the battery.

The input power line and also the 12 V power line are monitored by the INA3221 chip. It senses the current by measuring voltage on shunt resistors and it also measures voltages on these lines. From these two values, the power consumption can be calculated. This chip communicates with ESP32 by the I²C bus. The maximum power consumption is calculated to be P = 13.6 W. The whole schematic diagram can be found in the GitHub repository.

IV. SOFTWARE

The software for this project includes firmware in ESP32 and ROS2 middleware running on the SBC. The whole ROS2 graph can be seen in Figure 5. The SBC runs two ROS2



Fig. 5. RSS ROS2 graph (nodes are in ellipses, topics are in rectangles)

nodes. One of them, the *tracking algorithm*, subscribes to topics containing the Euler angles and raw data from the IMU. As an output, it publishes the threat level. The other node, the *control* node, subscribes to additional information, such as temperature, power consumption, and amount of ambient light, and publishes instructions for power management IC. All the topics use standard ROS2 messages. Development of this system was done on Microsoft Windows Subsystem for Linux (WSL) with running the Ubuntu distribution. All ROS2 principles were sourced from the ROS2 wiki [4].

The firmware for ESP32 was written in C language using the ESP-IDF framework. While ROS2 is a powerful tool for non-resource-constrained devices, for embedded devices it cannot be used. Fortunately, there exists a lightweight version of ROS2 for microcontrollers, called the *micro-ROS* [5]. It can be compiled for ESP32 but to communicate with the ROS2 graph it needs an *agent*, which runs also on SBC. The micro-ROS libraries were imported to the ESP-IDF environment as a component and were used to create the only node on ESP32. The whole firmware is supported by the FreeRTOS system.

The ESP32 takes the raw data from the IMU as linear acceleration and angular velocity and transforms them to Euler angles using the complementary filter. In further development, the Kalman filter will be used [6]. The node publishes both the raw data and the Euler angles. On the other hand, the subscribed threat level is used to light up the LEDs or beep the buzzer. The additional information on power consumption can be calculated from voltage and current values obtained from the power monitor IC. This information together with others previously mentioned is published to the \addl_data topic. The voltage converter for the radar can be disabled according to the message from the \control topic.

V. HARDWARE DESIGN

The PCB was designed in KiCad freeware software. Since the design includes 50 MHz clock lines and differential pairs, the impedance control and delay matching principles had to be followed when designing the PCB. The PCB has four layers with the following stack-up: high-speed traces, ground plane,



Fig. 6. 3D design of the enclosure with the module and radar board

power domains, and low-speed signals. The visuals of the PCB can be found in the GitHub repository.

The ALPS Generic Radar 5 is connected to the development module by spring-loaded pins and Flexible flat cable (FFC) used for future debugging. Both boards are enclosed in 3D printed covering mounted on the seat pole of a bicycle. The design of the complete device is shown in Figure 6.

VI. CONCLUSION

The development module was designed, manufactured, and assembled according to the needs of the company. The software concept was presented and is currently being programmed. Finally, the ALPS development team will use this module for future development of the ALPS Generic Radar 5, but mainly to prove the concept of enhancing the tracking algorithm by the external IMU data.

REFERENCES

- P. Norman, "What are rearview radar bike lights and should you use one?", *Bike Radar*, 2023.
- [2] N. H. Khan and A. Adnan, "Ego-motion estimation concepts, algorithms and challenges: an overview", *Multimedia Tools and Applications*, vol. 76, no. 15, pp. 16581-16603, 2017.
- [3] L. Seongwook, S. Lee, S. Lim, and S. -C. Kim, "Machine Learning-Based Estimation for Tilted Mounting Angle of Automotive Radar Sensor", *EEE Sensors Journal*, vol. 20, no. 6, pp. 2928-2937, 2020.
- [4] S. Macenski, T. Foote, B. Gerkey, C. Lalancette, and W. Woodall, "Robot Operating System 2: Design, architecture, and uses in the wild", *Science Robotics*, vol. 7, no. 66, May 2022.
- [5] P. Nguyen, "MICRO-ROS FOR MOBILE ROBOTICS SYSTEMS", Master thesis, Västerås, 2022.
- [6] H. Ferdinando, H. Khoswanto, and D. Purwanto, "Embedded Kalman Filter for Inertial Measurement Unit (IMU) on the ATMega8535", in 2012 International Symposium on Innovations in Intelligent Systems and Applications, 2012, pp. 1-5.

Measurement of the DVB-T2-based MISO Signal influenced by I/Q-errors in the OFDM Modulator

Simon Buchta

Department of Radioelectronics Faculty of Electrical Engineering and Communication Brno University of Technology Brno, Czech Republic xbucht30@vutbr.cz

Abstract—The Second Generation Digital Video Broadcasting Terrestrial (DVB-T2) system stands out as the most advanced system among the second generation DVB standards. Among others, it supports the multiple-input single-output (MISO) transmission technique to enhance terrestrial TV broadcasting. This paper focuses on a performance study of the DVB-T2 system utilizing the MISO technique for various fixed transmission scenarios, which are modeled by different transmission channel models. Additionally, it considers imperfections of the transmitter, due to errors in the Orthogonal Frequency-Division Multiplexing (OFDM) modulator. To facilitate this study, a laboratory measurement workstation with an appropriate measurement methodology is established. The proposed concept enables the interchangeability of measurement equipment and set-top-boxes (STBs).

Index Terms—DVB-T2, MISO, I/Q-errors, power imbalance, BER, MER, QEF.

I. INTRODUCTION

The Second Generation Digital Video Broadcasting Terrestrial (DVB-T2) system signifies a remarkable leap forward in DVB technology [1]–[3]. Engineered to address the evolving needs of contemporary television broadcasting in Europe, DVB-T2 provides enhanced flexibility and efficiency compared to its predecessor, the DVB-T system [4]. Notable advancements in the DVB-T2 system include the incorporation of multiple-input single-output (MISO) transmission mode and the utilization of rotated constellations, both of them, at specific transmission scenarios, can improve the overall performance of the DVB-T2 signal transmission [5].

In the past decade, numerous studies have delved into examining the performance of DVB-T2 transmission employing the MISO technique. Tormos et al. [6], [7] presented comprehensive experimental findings on the performance of DVB-T2 in classical mobile channels, particularly the typical urban (TU6) profile, across various Doppler frequencies within a conventional Single Frequency Network (SFN) compared to distributed MISO setups. Their results, obtained with two, three, or four transmission antennas, indicate that SFNs employing distributed Alamouti MISO outperform classical SFNs when utilizing two or three antennas.

This work was supported by the Internal Grant Agency of the Brno University of Technology under project no. FEKT-S-23-8191.

Ladislav Polak

Department of Radioelectronics Faculty of Electrical Engineering and Communication Brno University of Technology Brno, Czech Republic polakl@vut.cz

Particularly, with two antennas, it is possible to achieve higher the performance gain.

Morgade et al. [5] and Polak et al. [8]–[10] investigated the performance of DVB-T2 in both SISO and MISO configurations across a range of fixed TV channels and specialized transmission scenarios. Their findings reveal a significant improvement in performance when employing the MISO technique compared to SISO setups. Interestingly, the utilization of a rotated constellation did not lead to a significant enhancement in the performance of the DVB-T2 signal.

In [11], researchers conducted a simulation-based study on the performance of DVB-T2 MISO systems for three channel models, namely Additive White Gaussian Noise (AWGN) Ricean (RC20) and Rayleigh (RL20), while considering six DVB-T2 code rates and two constellation types. Their findings indicate that employing the MISO mode in DVB-T2 can lead to improvement of signal quality, particularly in SFN areas. The modified Alamouti MISO scheme effectively reduces spectral interference at the receiver, resulting in improved reception. Additionally, optimal MISO gain is achieved with a low constellation rate and high code rate. Another study by Gbadamassi et al. [12] investigated the performance of two distinct MISO schemes: co-located and distributed. The results of the study demonstrate that the MISO system with a distributed topology outperforms the co-located topology in both Rayleigh fading scenarios and SFN environments.

Contribution: Compared to previous studies, this paper offers two main contributions. Firstly, it explores the influence of various imperfections in the Orthogonal Frequency-Division Multiplexing (OFDM) modulator on the DVB-T2 MISO signal transmitted under different transmission scenarios and conditions. Secondly, it provides an analysis of the measurement results obtained by three different equipment and set-top-boxes (STBs). To the best of our knowledge, a similar study has not been conducted thus far.

Remaining parts of this paper are organized as follows. Section II contains a brief description of the DVB-T2 MISO transmission and I/Q-errors in the OFDM modulator. The laboratory measurement setup is described in Section III. The measurement results are evaluated in Section IV. Section V concludes this paper.

II. DVB-T2 MISO AND I/Q-ERRORS

III. MEASUREMENT WORKPLACE

A. DVB-T2 MISO transmission

The DVB-T2 MISO network stands apart from the standard SFN [1] by concurrently transmitting two slightly different versions of the intended signals through multiple (minimally two) transmitting antennas. Usually, these transmitters are geographically dispersed, providing an ideal configuration for maximizing MISO network efficiency. By utilizing multiple transmitters, the MISO network can leverage transmit diversity, resulting in improved Carrier-to-Noise Ratio (C/N), data rate, and network coverage [2].

In the DVB-T2 system, MISO transmission relies on a modified version of the Alamouti scheme [1]. Transmitters are paired together, concurrently transmitting payload data in pairs as well. This setup enhances signal diversity, thereby improving reception characteristics within the DVB-T2 MISO network. A significant advantage of using the modified Alamouti scheme is its relatively simple implementation at both the transmitting and receiving ends of the network [2].

B. I/Q - errors in the DVB-T2 OFDM modulator

The OFDM modulation process involves several essential steps. Initially, data is transformed from the frequency domain to the time domain through an Inverse Fast Fourier Transform (IFFT) operation. Additionally, a Guard Interval (GI) is inserted between consecutive OFDM symbols to mitigate potential inter-symbol interference (ISI) and ensure reliable transmission. Before transmission, the signal undergoes RF modulation utilizing an In-phase and Quadrature (I/Q) modulator. During this process, the *I* component is multiplied by a sine signal. Such a scheme facilitates efficient signal transmission over the channel while preserving its integrity [1].

Ensuring precise amplitude ratios and accurate phase configuration between the I and Q signal branches is crucial to avoid I/Q-errors. Two primary types of IQ-errors exist: Amplitude Imbalance (AI) and Phase Imbalance (PI). AI occurs when the amplitude levels in the I and Q branches are unequal, whereas PI arises when the phase difference deviates from the ideal 90° alignment (marking by red lines in Fig. 1).



Fig. 1: DVB-T2 modulator with I/Q-errors (based on [1])

The realized laboratory workplace to measure the influence of different I/Q-errors on DVB-T2 MISO signal is captured in Fig. 2. This setup is also suitable for testing and measuring different STBs. The basic concept of this measurement setup is originated in [10].

The SFU and SFE DVB-T2 signal generators from Rohde & Schwarz (R&S), representing transmitters TX1 and TX2, are utilized to generate a DVB-T2 MISO radio frequency (RF) signal. To broadcast a correct DVB-T2 MISO signal, both signal generators must be synchronized at the modulator interface level [10]. The video transport stream (TS) used for the DVB-T2 MISO signal broadcasting is generated in the SFU generator. The DVB-T2 system parameters employed for broadcasting the DVB-T2 MISO RF signal are summarized in Table I. Such a configuration is suitable for so-called fixed transmission scenarios [5]. In the SFU and SFE generators, the power level of the signals was set to -35 dBm and -30 dBm, respectively. Both signals are RF modulated and combined using a TEROZ T 226 K RF signal combiner. During the measurement, the SFU generator allows for setting different I/Q-errors and use different fading channel models, while the SFE generator permits only varying the value of C/N.

The DVB-T2 RF signal, affected by *I/Q*-errors and emulated transmission channel, is divided into two paths. The first path leads to equipment that measures various parameters such as bit and modulation error ratio (BER and MER). In this paper, three pieces of equipment were utilized: the R&S ETL TV analyzer, the Sefram TV analyzer, and the DVMS1 DTV monitoring system. The second RF signal path leads to STB connected to a TV, monitoring the condition of Quasi-Error Free (QEF) reception [1]. In this study, three STBs were employed: the Thomson THT712, the Sencor SDB 5002T, and the STC6000HD PVR.

In these measurement, we are considering the following scenario: there is a power imbalance of 5 dB between TX1 $(P_{TX} = -35 \, dBm)$ and TX2 $(P_{TX} = -30 \, dBm)$. During the measurement, the values of C/N simultaneously adjusted on both SFU and SFE signal generators. Additionally, two types of I/Q-errors are considered: AI = 10%, AI = 10% and $PI = 10^{\circ}$.



Fig. 2: Measurement workplace $(1 - SFU \text{ signal generator}, 2 - SFE \text{ signal generator}, 3 - attenuator}, 4 - R&S ETL - TV analyzer, 5 - DVMS1 - DTV monitoring system + PC monitor, 6 - STB, 7 - TV)$

Frequency	570 MHz
Bandwidth	8 MHz
FFT mode	32K extended
Guard Interval	1/16
Code Rate	2/3
Pilot Pattern	PP2
Transmission technique	MISO
Modulation	256-QAM
Rotation of constellation	On
Channel models	AWGN, RC20, RL20 [10]

TADLE I. System parameters of the DVD-12 signal with	TABLE I: System	parameters	of the	DVB-T2	signal	MISC
--	-----------------	------------	--------	--------	--------	------

IV. MEASUREMENT RESULTS

The extensive laboratory-based measurement results are graphically presented in Figs. 3-5. These measurements cover three different transmission scenarios. In the first (reference) scenario, a transmission with AWGN conditions was considered. The results obtained (see Fig. 3) highlight that the simultaneous presence of high I/O-errors, namely AI = 10%and $PI = 10^{\circ}$, exerts the most significant influence on the DVB-T2 MISO performance. Additionally, a noticeable difference is observed among the measured BER and MER values obtained using different equipment. The DVMS1 DTV monitoring system provides the least accurate results, possibly because it is primarily designed for monitoring the video TS rather than conducting long-term measurements of objective parameters on the physical layer of the DVB-T2 system. The Sefram TV analyzer did not measure the number of repeated channel decoding (iterations) per FEC Frame (FECFRAME) due to its lack of support for this parameter.

The performance of the DVB-T2 MISO signal affected by I/Q-errors in the Ricean channel (RC20) [10], in terms of the measured objective parameters, is depicted in Fig. 4. Once again, the impact on the DVB-T2 MISO system is minimal when only AI = 10% is considered. However, overall, the measured BER and MER values are slightly better than in the AWGN channel. This phenomenon is likely due to the higher power level of the TX2 signal (generated in SFE) compared to TX1 (generated in SFU) by 5 dB. Consequently, the transmission path for the TX2 signal exhibits only AWGN features, allowing for similar MER values to those observed in the first measurements.

Finally, the same measurements were repeated for the case where Rayleigh fading channel (RL20) conditions [10], emulated by 20 echoes, were applied to the TX1 path. Interestingly, the obtained results are very similar to the previous case when the RC20 fading channel model was considered. This similarity is further confirmed by the measurement of the number of FEC decoding iterations (see Fig. 5 (c)). It is evident that the DVB-T2 MISO system exhibits solid resistance against I/Q-errors and noise when one of the signal paths features an AWGN channel.

The comparison of the C/N values required for QEF reception (BER after FEC decoding $\leq 10^{-7}$) of the DVB-T2 MISO signal is listed in Table II.

V. CONCLUSION

In this paper, the performance of the DVB-T2 system utilizing MISO transmission mode under different channel conditions was investigated, focusing on *I/Q*-errors. To facilitate this study, a comprehensive laboratory setup was established, suitable also for measuring and testing various STBs.

The results revealed that the overall performance of the DVB-T2 configuration is significantly affected by the occurrence of I/Q-errors in the OFDM modulator. Particularly noteworthy was the detrimental impact observed under simultaneous AI and PI conditions, representing the worst-case scenario. However, the adverse effects of fading channels, notably RL20, can be "mitigated" when favorable channel conditions are present in the second signal path. Furthermore, the measurements highlighted variations in BER and MER values across different measurement equipment. Conversely, when measuring the C/N values required for QEF reception (as detailed in Table II), minimal differences were observed among the tested STBs.

This work will extend to include additional measurements encompassing various I/Q-errors, power imbalances, and transmission scenarios (mobile and portable) [10].

REFERENCES

- FISCHER, Walter. Digital video and audio broadcasting technology: a practical engineering guide. 3rd ed. Berlin: Springer, 2010. ISBN 978-3-642-11611-7.
- [2] SERIES, B. T. "Frequency and network planning aspects of DVB-T2. 2012."
- [3] EIZMENDI, Inaki, et al. "DVB-T2: The second generation of terrestrial digital video broadcasting system," *IEEE Trans. on Broad.* 2014. 60.2: 258-271.
- [4] BARUFFA, Giuseppe, et al. "Real-Time Generation of Standard-Compliant DVB-T Signals," *Radioengineering* June 2018, vol. 27., no. 2, pp. 476–484.
- [5] MORGADE, Javier, et al., "SFN-SISO and SFN-MISO Gain Performance Analysis for DVB-T2 Network Planning," *IEEE Trans. on Broad.* June 2014, vol. 60., no. 2, pp. 272–286.
- [6] TORMOS, Mokhtar, et al. "Experimental performance of mobile DVB-T2 in SFN and distributed MISO network," In: 2012 19th International Conference on Telecommunications (ICT). IEEE, 2012. p. 1-5.
- [7] TORMOS, Mokhtar, et al. "Modeling and performance evaluations of Alamouti technique in a single frequency network for DVB-T2." EURASIP J. on Wireless Commun. and Net., 2013, 2013: 1-12.
- [8] POLAK, Ladislav; KALLER, Ondrej; KRATOCHVIL, Tomas. "SISO/MISO performances in DVB-T2 and fixed TV channels." In: 2015 38th International Conference on Telecommunications and Signal Processing (TSP), IEEE, 2015. p. 768-771.
- [9] POLAK, Ladislav, et al. "Influence of the LTE system using cognitive radio technology on the DVB-T2 system using diversity technique," *Automatika* 2016, 57.2: 496-505.
- [10] POLAK, Ladislav, et al. "On the Performance of DVB-T2 MISO System: Special Fixed Transmission Scenarios," In: 2021 31st International Conference Radioelektronika. IEEE, 2021. p. 1-4.
- [11] GHAYYIB, Hamzah Sabr; MOHAMMED, Samir Jasim. "Performance improvement of DVB-T2 SFNS by using MISO transmission scheme." In: AIP Conference Proceedings. AIP Publishing, 2022.
- [12] GBADAMASSI, Abdoul-Warris, et al. "Co-located and distributed MISO techniques in DVB-T2 Single Frequency Networks." In: 2021 4th International Conference on Advanced Communication Technologies and Networking (CommNet), IEEE, 2021. p. 1-9.
- [13] PROAKIS, John G. Digital communications. 5th ed. Boston: McGraw-Hill, 2008. ISBN 978–0–07–295716–7.
- [14] POLAK, Ladislav, et al. "Single Frequency Networks for DVB-T2: Analysis of Real Case Scenarios in Czech Republic," In: 2023 33rd International Conference Radioelektronika. IEEE, 2023. p. 1-6.



Fig. 3: AWGN channel (solid line: no I/Q-errors, dashed line: AI = 10%, dotted line: AI = 10% and $PI = 10^{\circ}$)



Fig. 4: RC20 channel (solid line: no I/Q-errors, dashed line: AI = 10%, dotted line: AI = 10% and $PI = 10^{\circ}$)



Fig. 5: RL20 channel (solid line: no I/Q-errors, dashed line: AI = 10%, dotted line: AI = 10% and $PI = 10^{\circ}$)

Channal madal		STB			Measurement Equipment		
	ng-enois	Thomson	Sencor	STC	ETL	Sefram	DVMS1
	no I/Q-errors	19	19	18	19	18	19
AWGN	AI = 10%	19	19	19	19	19	19
	AI = 10% and PI = 10°	21	20	20	20	20	20
	no I/Q-errors	19	19	18	18	18	18
RC20	AI = 10%	19	19	19	19	19	18
	AI = 10% and PI = 10°	19	20	19	19	19	19
	no I/Q-errors	19	19	18	18	18	18
RL20	AI = 10%	19	19	19	19	18	18
	AI = 10% and PI = 10°	19	19	19	19	19	19

TABLE II: Required C/N in unit of dB for QEF reception

Device for galvanic plating of 3D printed parts

S. Orgoň

Faculty of Electrical Engineering and Communication, Brno University of Technology, Brno, Czech Republic E-mail: 240869@vutbr.cz

Abstract— This work aims to bring closer the process of galvanic plating of 3D prints, it is dedicated to individual parts of 3D printing where we go through different types of printing and materials. Also points out the importance of pre-treatments before galvanic plating and discusses the individual options for these pretreatments. It also deals with galvanic plating itself, where it analyzes the principle of this technology of thin layer deposition and differentiates between the individual types of possible coatings used. At the end, it deals with the concept of the device for such galvanic plating.

Keywords— Galvanic plating, 3D print, Surface treatment

I. INTRODUCTION

Nowadays, more and more emphasis is placed on the used materials and their price. We want to achieve the best possible properties for the final product, while we want to make it as cheap as possible. Properties such as low weight, easy manufacturability, sufficiently durable surface, all at a low price. In many cases, we can ensure improvement of these properties by electroplating 3D prints.

Surface treatment with galvanic plating is a relatively simple technology that is widely used in industry as well as 3D printing at home. When we combine these two technologies, we can get a light and cheap final product with improved properties. This technology can be used in the home environment as well as in industry, as the number of users of 3D printers in the home environment is currently growing. When combined with electroplating equipment, greater possibilities can be achieved.

In addition to surface treatment, electroplating also gives individual products a better visual appearance. Currently, there is no complex equipment on the market that is intended for galvanic plating. In practice, laboratory supplies are usually used, a bath is placed next to them and plated element is immersed in the electrolyte with the material that is applied to it. Therefore, this work is devoted to such a device and to the theories of pre-treatments and galvanic plating itself.

II. PROCEDURE FOR ELECTROPLATING 3D PRINTED PARTS

In the case of galvanic plating of plastics or in general galvanic plating of non-conductive components, the process is

more lengthy compared to the plating of an already conductive component, as it is necessary to make the given components conductive. The procedure is shown in the block diagram figure no. 1.



Fig. 1. Block diagram of procedure for electroplating 3D printed parts

The most important requirement for electroplating is the cleanliness of the surface on which the layers are applied. Except for exceptions, parts cannot be plated without being placed in a solution without preliminary treatment. Contamination often prevents deposition of the material. There are two types of surface treatment methods: Physical cleaning and chemical cleaning. Chemical cleaning consists of using solvents that are either surface-active chemicals or chemicals that react with the metal/surface. Physical cleaning uses mechanical energy to remove dirt. Physical cleaning includes, for example, sanding or grinding.

Next step is making the surface conductive, that is possible in several ways, the most used method is the use of a spray with conductive paint such as graphite, copper and silver paint, which are most often used given their good specifications and price. The biggest difference between the individual paints is in their price and conductivity.

Another method is airbrush technology, or classic application with a brush. In these cases, the same types of paints are used as with spray.

The least common way is to mix conductive material into plastic, for example in 3D prints that use material with carbon admixture. This method of plating is not as efficient and requires more effort and quality of other factors in the plating process.

III. GALVANIC PLATING

Electroplating itself involves passing an electric current through a solution called an electrolyte as shown in Figure 2. We achieve this by placing two electrodes in the solution and connect them to the circuit so that copper becomes the positive electrode, i.e. the anode, and brass becomes the negative electrode, i.e. the cathode.



Fig. 2. The principial scheme of electrodeposition copper

The electrode on which we apply a thin layer is generally made of a cheaper metal or non-metal coated with a conductive material such as graphite, in any case it must be conductive, otherwise the electric current will not flow and the plating will not occur. When we turn on the voltage source, the copper sulfate solution splits into ions. If we talk about copper which is most often used we usually use an electrolyte based on copper sulfate, which is dissolved in diluted sulfuric acid, which serves not only to improve conductivity, but also helps the layers to be uniform. Copper sulfate dissociates in solution into Cu2+ and (SO4)2- ions. Cu2+ ions are reduced to copper at the cathode and deposited there as a coating. The sulfate ions move to the copper anode where new copper sulfate is formed by dissolving the anode in the solution. With copper, we can also use a solution based on cyanide, which has very good adhesive properties, but is not used often in practice due to its toxicity.

The time during which galvanized atoms accumulate on the surface of the negative electrode depends on the magnitude of the electric current and the concentration of the electrolyte, the temperature have the strong share in the process therefore affects the final product. By increasing one of them, we achieve increase of the movement speed of electrons and ions through the circuit and also an increase of the plating speed.

The goal of this technology is usually to apply a thin layer of metal on the material for the purpose of protection, improvement of mechanical properties or possibly aesthetic properties. Properties depend on the thickness of the layer, the typical thickness of electroplating ranges from 0.5 microns up to 20 microns.[2][4][5][7]

A. Local Galvanic plating

The technology of local galvanic plating or so-called tamponing is a special type of technology in which only a certain part of the object's surface is plated. Most often, this technology is used for hard-to-reach places on the surface, or for repairing the surface. The main advantages include lower consumption of electricity, electrolyte, and a faster plating process for small areas. The principle is similar to the technology of classical plating. The plated object is connected as a cathode, while the anode is placed in a holder and covered with absorbent material, this material is impregnated with electrolyte. After applying the anode to the plated object, the circuit is closed and the metal ions begin to be captured on the surface of the object. We can supply the electrolyte manually or automatically. The setup of local galvanic plating is shown in Figure 3.

Parts of the equipment for local electroplating:

- 1) Plated component (cathode)
- 2) Tampon
- 3) Electrolyte
- 4) Electrolyte supply to tampons (only with automatic)
- 5) Anode

6) Movement of the anode to ensure uniform distribution of the coating



Fig. 3. The principial scheme of local galvanic plating

IV. CONCEPT

The design of the device as it is shown in Figure 3 includes a 12V 20A direct current source, which serves to power other components and at the same time is used as a source for the galvanic plating itself. For proper galvanic plating, it is necessary to regulate the current value so that we achieve an optimal result, this is provided by the control circuit, which is marked as a constant current source in the block diagram. The second constant current source is designed for the galvanic pen using the principle of local plating, which is to serve for the final modifications and repair of plated parts. Furthermore, the design includes a switched power source powering other parts of the device such as the user interface and the microcontroller.

The microcontroller will control the individual blocks in the circuit, at the same time it will be used for the necessary calculations for plating, such as the calculation of time and current values based on the properties of the electrolyte and the plated object. We will use the ESP32-PICO-D4 as a microcontroller, which should have enough input-output ports and the necessary functions such as PWM, DAC and others.

As a user interface, an LCD display and an encoder with a button are used to set the current value for plating and other functions such as turning the motor, turning on the light and others. At the same time, the LCD display provides information about the plating process, such as bath temperature, duration of deposition, current and voltage.

The design also includes a sander that will be used for surface treatment so that the best possible surface is guaranteed before being placed in the bath and the metal can settle on the surface. The plating bath features a motor for rotating the plated component to achieve even plating and electrolyte itself is mixed by air pump to further enhance the quality of the process. We will measure the temperature of the bath because of the theoretical change in the behavior of the plating process with temperature change. Lighting of the bath is ensured by led strip so that the product under the electrolyte is clearly visible and can be checked.

The device itself will have to be cooled due to the current load, we will try to achieve this by passive cooling using coolers and active cooling using fans.



Fig. 3. The block diagram of final product

V. DESIGN OF THE DEVICE

The 3D parts of the design were made in Solidworks 2020, final design is combination of 3D printed parts and the main metal part which is mainly used because of better EMC and thermal properties. The full setup of galvanic plastic device is shown in Figure 4. On the right is the device itself, at the front is display with encoder for easy user interface. On the right side of device is galvanic pen for local plating and on the top is the sand grinder. On the left is bath with mount on top to attach the 3D printed product, the motor for rotating the object is inside the mount.



Fig. 4. Full design of galvanic plating setup

In the Figure 5 we can see the device itself with the layout shown inside.



Fig. 5. Device for galvanic plating with shown inside

In the Figure 6 which is shown below, we can see configuration of back panel with all connectors and master power off button.



Fig. 6. Device for galvanic plating - back panel configuration

The electronic part was designed in Altium designer, the final PCB is shown in Figure 7. The PCB has 2 layers with mostly SMD components mounted from only one side to be more easily placed into final product.



Fig. 7. The design of PCB

VI. CONCLUSION

First part of this paper deals with overall procedure of electroplating the 3D printed parts. At the end of this part we talk about the quick theory of galvanic plating itself. Second part introduces the device for galvanic plating which was developed as the first part of bachelor thesis. We talk about individual blocks that are needed for such a device to have good final results of plated object. In the future there is room for improvement of the device, because we use ESP32-PICO-D4 and can send the notification to user or the user could manage the device through phone or computer.

ACKNOWLEDGMENT

I would like to thank my supervisor doc. Ing. Petr Vyroubal Ph.D. for his valuable comments and guidance during work on this project.

REFERENCES

 URBANEC, Filip. Online, Návrh zkušebního zařízení pro automatické galvanické tamponování. Diplomová práce. Liberec: Technická univerzita v Liberci, Fakulta strojní, 2011. Dostupné z: https://dspace.tul.cz/server/api/core/bitstreams/7e44b98a-b894-46c9-8e4b-0df846de0135/content [cit. 2023-11-13].

- Microchemicals. Electro-plating of certain metals. Online.. Dostupné z: https://www.microchemicals.com/products/electroplating/electroplating_metals.html. [cit. 2023-11-18].
- [3] LAMAČ, Martin. Lokální galvanické pokovení v extrémním prostředí Online, bakalárska práce. Praha: České vysoké učení technické v Praze, Fakulta strojní, 2020. Dostupné z: https://dspace.cvut.cz/bitstream/handle/10467/86255/F2-BP-2020-Lamac-Martin-BP%20-%20Martin%20Lamac1.pdf?sequence=-1&isAllowed=y [cit. 2023-11-13].
- [4] Chem.libretexts. Electroplating. [online]. Dostupné z: https://chem.libretexts.org/Bookshelves/Analytical_Chemistry/Supplem enta_Modules_(Analytical_Chemistry)/Electrochemistry/Electrolytic_C ells/Electro plating [cit. 2023-11-13].
- [5] WOODFORD, Chris. Electroplating. Online. [2021]. Dostupné z: https://www.explainthatstuff.com/electroplating.html [cit. 2023-11-13].
- [6] ŽÁK, V., KUDLÁČEK J.: Technologie lokálního galvanického pokovování (tampónování). Povrcháři online. Praha, 2008 vol.3, s.1-4 Dostupné z: http://povrchari.cz/kestazeni/200803_povrchari.pdf [cit. 2023-11-13].
- [7] Elektrolab. Galvanizácia a galvanické pokovenie. Online. [2022].
 Dostupné z: https://www.elektrolab.eu/blog/galvanizacia-a-galvanicke-pokovenie [cit. 2023-11-13].
Detection of parking space availability based on video

1st Miloslav Kužela Department of Radio Electronics FEEC, Brno University of Technology Brno, Czech Republic 240648@vut.cz

Abstract—This paper deals with the use of Machine vision and ML (Machine Learning) for a parking lot occupation detection. It presents and compares an already existing technology that solves such a problem with an AI (Artificial Intelligence) usecase. It introduces tools used to train and create such models and their subsequent results as well as a dataset that was used to verify the trained networks and discusses the future of how such a technology could be used to effectively and more affordably detect occupied parking spaces on parking lots.

Index Terms—Machine Vision, Machine learning, Parking occupancy, Python

I. INTRODUCTION

The use of parking sensors to monitor larger, enclosed parking lots is the norm at this day and age. But with the use of ceiling or wall mounted sensors, a lot of parking lots are left out due to their locations. The cost of such devices must be taken into account. With the recent rise of AI, especially Machine learning and Machine vision, tackling such a problem with these tools could be of benefit. A CCTV (Closed-circuit television) camera and a computer is all that is needed for the solution of this problem. And due to a lot of already existing parking lots having a camera in use, it is more cost efficient. The aim of this article is to show and introduce the reader to a concept of such a system and its tools to create a simple parking detection algorithm using Machine learning and Machine vision. It present tools that helps a user train, test or fine-tune an existing model for such a use case. It also consults the possible ways of utilizing such a system in real world application.

A. Existing commercial solutions

The current state of systems that are used for parking occupancy monitoring can be divided into two groups. One that counts only incoming and outgoing cars, compares it with the total parking spaces available and shows the result on an information board of some sorts. The other type is using active sensors placed near all parking places which detects the presence of a car. Such solution is accurate although expensive and is hard to realize in an open space parking lot. Ground sensors do exist, but due to battery requirements they are not as effective. An example of such sensors include an Ultrasonic sensor, Magnetic or MMW (Millimeter wave) radar. For example an ultrasonic sensor, mostly used in enclosed 2nd Tomáš Frýza Department of Radio Electronics FEEC, Brno University of Technology Brno, Czech Republic fryza@vut.cz

parking lots, due to the fact that they often require to be mounted on the roof so they can measure the distance between it and the floor. When a car parks under such sensor it detect the change of distance, this information is then sent to a central unit which keeps track of all the parking spaces and their occupancy. When compared to a solution proposed later in this paper, such sensors are more accurate, although that advantage is outweighed by the latter's price, and the possibility of utilizing an already existing infrastructure such as cameras.

II. METHOD

A. Use of Cameras

When relying on data acquired by the means of a digital camera. It is necessary to be aware of possible inaccuracies that arise when using such technology. If the weather conditions change there is a high chance that the camera's automatic functions will make the resulting photos look much different than what it captured last time. A simple machine vision algorithm could have a problem with such a fact. So an image preprocessing stage is in order.

B. Initial idea

The use of machine learning in such a problem is simple. Having a model decide between two classes: Car and background. Then apply such a model to a parking lot images and tests its accuracy. The main problem was deciding what existing machine vision network/model to use. The main network candidates were either faster R-CNN (Regional Based Convolution Neural Network) [?] or Retina-net [?]. Considered models were Mobilenet V3 [?] and ResNet50 [?].

C. R-CNN

As a regional proposal network, meaning that the first stage of prediction is locating regions of interest which are then used in prediction, which could be a benefit as the camera is not moving and the regions are staying mostly at a constant location in the image. MobileNetV3 was the newest model available that uses the Convolution network as its base and is designed to be ran on a low power mobile devices. There are two versions, small and large, which correlate to its number of layers. [?]

D. Retinanet

This network uses a one-stage approach instead of the two stage of the previous network mentioned. It uses a FPN(Feature Pyramid Network) and Focal Loss to predict the correct class. Such network is suitable for detecting multiple small object in an image. It was chosen due to its faster one stage approach. Creators of this network recommend the use of a Resnet50 FPN model [?].

E. Dataset

A custom dataset was created to train pre-trained models and models from scratch called T10LOT. It consists of photos of a parking lot located next to a FEEC's (faculty of electronic engineering and communication) building. Photos were captured on an iPhone X camera with an application that took a photo every 30 minutes. The phone was held up by a 3D printed arm that was fixed to a window frame. The dataset contains over 100 images including different weather conditions such as fog, night and rainy weather. This dataset was used to train and verify the models. All of the images have a resolution of 1980x1080px.

F. Custom tools

Custom tools and scripts were created to assist with the creation of datasets, training and testing. Coded in Python with the use of Jupyter notebook. These tools allow the user to create their own datasets that can be latter used with the training CLI (Command Line Interface) applications. Dataset creation tools run in a Jupyter notebook and feature a GUI (Graphical User Interface). The training function includes a configuration stage that allows user to specify how to train the network, this includes the number of datasets to use, network and model combination, number of epochs, learning



Fig. 1. Parking space occupancy detection with water droplets

rate and more. There is also a testing script that test the trained network on testing images of a selected dataset that were not used during training, this eliminates any possible bias. The script outputs the accuracy and speed of inferencing a single image, this data is later used to decide the fastest network. The annotating application was reworked from a another paper that featured comparing trained models on a multiple datasets [?].

G. Labeling

As mentioned in chapter **??** a tool was created that assist with the labelling/annotating of images. User puts all the images into a specific folder and runs the application, then by clicking around individual parking slots crates polygons that surround it. If all of the images were captured in the same angle and position, there is a button for labelling all the images same way as the initial one. The second stage of labelling is marking each parking space as either occupied or not. This is done with another included tool, that allows the user to annotate individual spaces by just a click. Once the user is finished, program creates a dataset in a format that can be later used with the training and testing scripts.

H. Training

To be able to accurately use a neural network, it firstly needs to be trained on a dataset according to its specific use-case. In this case, images with parking lots. There is a difference between training a model with randomized weights and fine tuning a model with pre-retrained weight. All of the pretrained models featured in this paper were trained on a COCO(Common Objects in Context) dataset, which features a lot of different images featuring different objects. This trains the network to properly do region proposals and is easier to retrain on a dataset that fits the use-case. Models that started with randomized weights values were trained on two datasets that are really similar to the wanted use case, meaning that the camera was positioned above the parking lots and behind a window, the two datasets used are called CNRPark [?] and PKLot [?]. Lastly both pretrained and randomized weights models were trained on a T10LOT dataset for 20 epochs. No augmentation was performed on the training set of images, only the conversion to tensors were applied, although the use of color augmentation would benefit if training with images taken by an IR (Infra-red) cameras.

I. Testing

The testing script, loads the wanted model and runs it in evaluation mode through a set of testing images, that were not used in training. These are randomly selected when creating a dataset with the creation script included in this work. The testing is evaluated based on the F1 score, which is calculated from precision and recall. Lastly there is an option to enable exporting of the results into images, where each individual detection is visible. The current model is currently designed to differentiate between two classes: Car parked in a parking slot and background. When trying to understand the testing results an example of a badly trained model is presented in Fig.??.



Fig. 2. A simple diagram describing the workings of the system

Red dot is a false positive detection. Squares that highlight the individual parking spaces turn green if a valid detection is present. Red squares are occupied parking spaces that were not classified as and hence are categorized as false negatives.

$$precision = \frac{TP}{TP + FP} \tag{1}$$

$$recall = \frac{TP}{TP + FN} \tag{2}$$

$$F1 = \frac{2 \cdot precission \cdot recall}{precission + recall}$$
(3)

III. RESULTS

A handful of model combinations were tested. The testing was focused on their speed and accuracy. Images were provided from the T10LOT datasets, which contains 25 of pictures dedicated for testing. After comparing the results of individual models it is notable that the fastest and most effective ones are from the MobileNetV3 family, where the small version is faster but less accurate and vice-versa **??**.



Fig. 3. An example of correct and false detection

A. Difficult scenarios

It is worth noting the results of inferencing the networks on images that include a difficult scenario such as fog, rain or dark environment. It was surprising to see, that even with water droplets Fig. ??, that were in front of the camera and worked as a simulation of moisture accumulating on the lens, both MobileNetV3 networks were efficient at detecting parked cars, same applies to fog and even darkness Fig. ??, where it is hard even for a human to recognize if a parking slot is occupied or not.

IV. HARDWARE SELECTION AND SOLUTION PROPOSAL

The final Parking system will use a raspberry Pi V4 with a camera that overlooks the selected parking lot. The single board computer would need to be connected to a network by either a Wi-Fi connection or an Ethernet cable. Multiple cameras can be used to cover more of the parking place. A map of the parking lot would need to be created by labeling the individual position of parking places in the image. This would ensure that the final algorithm would label these places as occupied or not. The single board computer would run the object detection model and would mark the detected places a occupied, this information would be sent to a database server that could be accessed at any time by an IoT device that would display the current occupancy status of the parking lot.



Fig. 4. Detection in a dark scenario

A. Possible shortcomings

As the system is relying solely on cameras it could be easily disturbed if just one camera were to be slightly shifted. The

Model	Pretrained	Detector	Training loss	Validation loss	F1 score	Inference time
MobileNetV3 small	NO	FasterRCNN	0.0236	0.0256	0.972	11.7 ms
MobileNetV3 large	NO	FasterRCNN	0.0259	0.0284	0.970	13.8 ms
ResNet50	NO	FasterRCNN	0.481	0.431	0.218	16.4 ms
MobileNetV3 small	NO	RetinaNet	0.0136	0.0237	0.881	10.2 ms
MobileNetV3 large	NO	RetinaNet	0.0134	0.0232	0.873	15.1 ms
ResNet50	NO	RetinaNet	0.0949	0.0984	0.206	35.9 ms
MobileNetV3 large	YES	FasterRCNN	0.078	0.120	0.990	22.5 ms
ResNet50V2	YES	FasterRCNN	0.056	0.0637	0.962	97.9 ms
ResNet50V2	YES	RetinaNet	0.047	0.0758	0.961	73 ms

TABLE I TABLE OF RESULTS FOR INDIVIDUAL MODELS

neural network part of the algorithm wouldn't mind this fact, but the detection one would loose the correct positions of the individual parking slots. A system that would combat this problem needs to be proposed. The current system is designed to work only during the day, as it doesn't include IR light emitters or cameras that have an option to disable their IR filters. If such a night vision camera were to be used the network would need to be retrained on a new set of images taken in dark scenario.

V. CONCLUSION

As the number of cars and parking places increases in the world a more accessible and affordable space monitoring system will be required. With the increase use of Machine learning utilizing it in such a use case is possible as was shown in this paper. The use of a single board computer to monitor parking spaces can be cheap and efficient. The tools used for this research are available on a GitHub repository: https://github.com/slavajda02/parking-research-argon as well as a download link containing the testing result images and the created dataset previously mentioned. A further research and test will be conducted on a prototype of a system previously mentioned to see how quickly and effectively can such a computer process individual images.

ACKNOWLEDGMENT

The Author would like to acknowledge the contribution and valuable time provided by the supervisor, doc. Ing. Tomáš Frýza Ph.D., as well as to thanks the authors of other researches focusing on parking detection algorithm.

REFERENCES

- REN, Shaoqing, HE, Kaiming, GIRSHICK, Ross and SUN, Jian. "Faster r-cnn: towards real-time object detection with region proposal networks" Online. 6 January 2016.
- [2] LIN, Tsung-Yi, GOYAL, Priya, GIRSHICK, Ross, HE, Kaiming and DOLLÁR, Piotr, "Focal Loss for Dense Object Detection." Online. 7 February 2018.
- [3] HOWARD, Andrew, SANDLER, Mark, CHU, Grace, CHEN, Liang-Chieh, CHEN, Bo, TAN, Mingxing, WANG, Weijun, ZHU, Yukun, PANG, Ruoming, VASUDEVAN, Vijay, LE, Quoc V. and ADAM, Hartwig. "Searching for mobilenetv3" Online. 20 November 2019.
- [4] HE, Kaiming, ZHANG, Xiangyu, REN, Shaoqing and SUN, Jian, "Deep Residual Learning for Image Recognition". Online. 10 December 2015.

- [5] Anastasia Martynova, Mikhail Kuznetsov, Vadim Porvatov, Vladislav Tishin. "Revising deep learning methods in parking lot occupancy detection" Online. 2023
- [6] LIN, Tsung-Yi, MAIRE, Michael, BELONGIE, Serge, BOURDEV, Lubomir, GIRSHICK, Ross, HAYS, James, PERONA, Pietro, RA-MANAN, Deva, ZITNICK, C. Lawrence and DOLLÁR, Piotr. "Microsoft COCO: Common Objects in Context" Online. 20 February 2015.
- [7] DE ALMEIDA, Paulo R. L., OLIVEIRA, Luiz S., BRITTO, Alceu S., SILVA, Eunelson J. and KOERICH, Alessandro L. "PKLot – A robust dataset for parking lot classification", "Expert Systems with Applications" 1 July 2015. Vol. 42, no. 11, p. 4937–4949. DOI 10.1016/j.eswa.2015.02.009.

We make what matters work.



At Eaton, we believe that power is a fundamental part of just about everything people do. That's why we're dedicated to helping our customers find new ways to manage electrical, hydraulic and mechanical power more efficiently, safely and sustainably. To improve people's lives, the communities where we live and work, and the planet our future generations depend upon. Because this is what really matters. And we're here to make sure it works.

To learn more go to: Eaton.com/WhatMatters



Dáváme smysl věcem, na kterých záleží.*





✤ Naše vize

Pomocí technologií a služeb pro správu napájení zlepšit kvalitu života a prostředí.

Modular system for electrical impedance tomography

1st Roman Vaněk Department of Radio Electronics FEEC, Brno University of Technology Brno, Czech Republic roman.vanek3@vut.cz 2st Jan Mikulka

Department of Theoretical and Experimental Electrical Engineering FEEC, Brno University of Technology Brno, Czech Republic mikulka@vut.cz

Abstract—This paper presents a novel approach to an electrical impedance tomograph's system and hardware architecture design. It delves into the strategic decisions that shaped the hardware layer of a measurement system comprising multiple units and enabling intercommunication between them. The primary objective was to develop a cost-effective method for switching a current source among numerous electrodes and measuring voltage across each with a wide dynamic range while maintaining a minimal phase shift between the channel and input of an analog-to-digital converter. The measurement system is designed and documented, serving as a valuable reference for developing new additional units.

Index Terms—electrical impedance tomography, active electrode, modular system

I. INTRODUCTION

Electrical impedance tomography (EIT) is a method of imaging the distribution of a material's specific conductivity (or impedivity). These images are constructed from voltage between electrodes that make direct contact with the material. The voltage is formed by injecting an alternating current into the material by the tomograph. The method is actively used in the medical field, utilizing lung function monitoring. There is a big potential for extending its applications to other fields. Assuming the material is not biological, the option of supplying higher currents to electrodes replaces the issue of securing non-invasive measurement, enabling the measurement of materials with a wider range of impedance.

The central objective of this project is to develop easily adaptable measurement system. The primary strategy was to create a block design of multiple units and establish interconnections between them. The whole measurement chain consists of signal synthesis, a voltage-controlled current source, current multiplexing, voltage measurement by an active electrode obtaining a high dynamic range, voltage multiplexing, and main calculation of impedance using measured voltage, current values and phase between them. Each part/unit must meet specific requirements to ensure the system's modifiability and compatibility with possible extensions.

II. SYSTEM REQUIREMENTS

In order to expand the field of use to the materials with higher impedance or measurements taken over greater distances with materials of relatively low specific impedivity, it is necessary to increase the current injected into the material. This will allow electrodes placed further from the injecting electrodes to detect a voltage drop greater than any noise present. To further enhance the signal-to-noise ratio (SNR) and low voltage measurement capabilities, we chose to use the active electrodes that will also compensate for the unwanted noise potentially coupled to the cables carried from the central system to the contact of the material.

With increasing current and the total impedance of the material, we will need to provide greater voltages for the constant-current source, multiplexer (MUX), and other active devices to avoid saturation. A study of the current electronics market was carried out to determine the maximal voltage we can achieve without using mechanical devices, as was required in the project assignment. The most sensitive part of the measurement chain will be the analog MUXs, which will handle the switching of current source output to the electrodes; they must withstand both the current passed through the material and the voltage. The TMUX8108, with a switchable voltage of 100 V peak-to-peak, was chosen for its easy-tosolder package, wide frequency band, and relatively low price compared to other solutions. The maximum current value injected into the material by the device is also determined by the current capability of the MUX, as stated above, and constant-current source circuitry.

III. VOLTAGE MEASUREMENT

There are two types of voltage measurement: single-ended and differential, as shown in Fig. 1. The single-ended measures voltage with respect to ground potential, and the differential measures the difference between two voltages utilizing an instrumentation amplifier as a converter to single-ended output in the case of using analog-to-digital converter (ADC) with single-ended input. Differential measurement is typically used in passive EIT systems, where the measured voltage between adjacent electrodes is smaller than in the case of the singleended type. This allows a reduction of the ADC's dynamic range requirements [1]. The advantage of differential measurement is suppression of noise detected at both inputs of the instrumental amplifier; the problem with using the differential measurement in this application is that we cannot suspect that the gain of two measuring active electrodes will be the



Fig. 1. Two types of voltage measurement using active electrodes: differential (a) and single-ended (b).

same, in case of different gain settings on each active electrode there will be uncertainty inserted into the system because an instrumental amplifier located on the central system could cause phase inversion. However, an even more significant problem would be the potential phase delay introduced by one of the amplifiers in the active probe. This difference in phases at the input of the instrumental amplifier would cause frequency-dependant common-mode error at the output of the amplifier; these imbalances between electrodes were considered the principal sources of error in the differential voltage reading [2].

The single-ended type of voltage measurement lacks the advantages of common mode voltage rejection and typically low offset error of the differential voltage measurement. However, the significant advantage of this method is that we do not have to deal with the need for phase error correction. Furthermore, since we only need to set one gain, we can avoid the potential problem of phase inversion caused by the instrumental amplifier.

In order to meet the EIT method requirements of measuring the voltage drop across material without creating a low impedance path to the ground of the current supply affecting the measurement, the negative input of the operational amplifier (OpAmp) in the single-ended configuration needs to be connected to the second ground galvanically isolated from the ground/negative pole of the current supply. With galvanic isolation done in the power supply section, we can use a simple high-voltage OpAmp (not isolated one) in the active electrodes to address both the high input impedance needed when using any attenuator and the high SNR addressed in Chapter II.

Having two ADCs available is essential; we can use dual regular simultaneous mode and sample two channels at the same time [3]. One channel is for current measurement, and one is for voltage measurement. The 12-bit ADC in the microcontroller (MCU) converts these values and stores them in a preallocated buffer utilizing the direct memory access (DMA) channel [4]. With a use of the internal ADC of the MCU we get 24 Samples per period, when samling the highest acceptable frequency 100 kHz, which should be enough, but there is always option to add an aditions ADC unit to the system. Essentially, a calculation of the impedance \vec{Z} from determined magnitudes and phase shifts of voltage (U, ϕ_U) and current (I, ϕ_I) flowing through the measured phantom (1) [5].

$$\vec{Z} = \frac{U}{I}e^{j}(\phi_{U} - \phi_{I}).$$
⁽¹⁾

IV. SYSTEM DESIGN

The system was split into individual units, each focusing on a specific task and completing the measurement chain together. This makes minor redesigns and contributions to projects much easier and more cost-efficient than redesigning or improving the system from scratch. These are the units that the whole system consists of:

- Power supply module
- Main controller module
- Signal generator module
- Constant current supply module
- Multiplexer module
- Active electrodes

A block diagram was created to facilitate an overview of the system and its functioning between units, see Fig. 2. The block diagram divides all the units into larger blocks with names in them. They all have the same connector, carrying the bus interfaces through the system. If the unit uses some interface or power rails from the bus, then the interface is in the bolted.



Fig. 2. System block diagram

The arrows indicate the logical direction of communication. Some pins/interfaces are intended as possible extensions and upgrades of the system; therefore, they are included on the bus but not used by any units.

Galvanically isolated sections are identified in the block diagram accordingly: "GND" for the measurement circuitry and "GND2" as a ground reference of the signal generator unit and constant current source.

A. Interconnection

The bus is formed as a pack of multiple communication interfaces, pins for analog signals and power supply rails. The main controller unit drives the communication between units. Reading digital signalization or writing low-speed data, like address inputs of the MUXs, is done using the generalpurpose input/output (GPIO) expanders MCP23017 are used, connected to the I2C bus with a manually selectable address by a DIP switch on the unit. This expander can be populated on units that need some low-speed configuration process in the form of bit settings.

As described in Chapter III, there must be galvanic separation circuitry between the main controller module and any module referencing the GND2 (signal generator or constant current source unit). It is more cost-efficient to galvanically isolate two wires of I2C communication than all the individual GPIOs, which could be routed directly to the bus from an MCU. The bus also provides standard serial interfaces such as SPI, I2C, and UART for direct communication with components mounted on some units, e.g., the signal generator unit's direct-digital-synthesis (DDS) integrated circuit (IC).

B. Offset adder for analog signals

The analog pins must carry a signal within the range of ± 1.65 V. This range is crucial as it relates to the internal ADC of the STM32 MCU, which will be used. To ensure the input signal is correctly converted to a proper range, an offset adder circuit must be added to shift the signal to a 1.65 V DC offset - which is the exact half of the power supply voltage of the ADC. This offset helps to convert symmetrical voltage values to an acceptable range by STM's ADC, which is 0 V-3.3 V.

To demonstrate the disparities between a classical DC voltage adder with capacitive coupling Fig. 3(b) and a voltage adder that uses a summing amplifier Fig. 3(a) simulations were conducted on both of those.

Both types employ identical wiring to obtain the same offset voltage. This is achieved through a simple resistive divider made up of identical resistors R_3 , R_4 , R_5 , and R_6 , which divide VDD or V_{ref} (the MCU pin) voltages in half. To ensure the accuracy of this assumption, the resistor values must have a small tolerance. The active voltage adder has the benefit of working down to DC instead of the passive offset adder, which has a non-zero lower limit cutoff frequency introduced



Fig. 3. Offset adder simulation: active (a) and passive (b) voltage adder and low impedance simulation of the active electrode (c).

by capacitive coupling. However, a larger value of C_1 will result in a very low lower limit cutoff frequency (2).

$$f_{1} = \frac{1}{2\pi \cdot \left(\frac{R_{5} \cdot R_{6}}{R_{5} + R_{6}}\right) \cdot C_{1}}$$

$$= \frac{1}{2\pi \cdot \left(\frac{100 \cdot 100}{100 + 100}\right) \cdot 10^{3} \cdot 100 \cdot 10^{-6}} = 0.032 \,\mathrm{Hz}.$$
(2)

A second-order low-pass filter, consisting of R_{AIN} and C_{AIN} , was added after the offset adder circuitry to limit the noise potentially coupled to the signal and an aliasing effect of ADC. According to the application note, this second-order filter was chosen as sufficient enough [3]. Its cutoff frequency is calculated in (3) and selected to be lower than the Nyquist limit $(\frac{f_{samp}}{2})$, which for a sample rate of 2.4 MSps will be 1.2 MHz.

$$f_{2} = \frac{1}{2\pi \cdot R_{\text{AIN}} \cdot C_{\text{AIN}}}$$

$$= \frac{1}{2\pi \cdot 10 \cdot 10^{3} \cdot 16 \cdot 10^{-12}} = 994.718 \text{ kHz.}$$
(3)

The simulation resulted in the upper limit cutoff frequency being the $f_2 = 960.557 \text{ kHz}$ for the passive offset adder. This topology was preferred because it is more straightforward than the active offset adder, and our application does not need to measure DC voltage.

V. ACTIVE ELECTRODES

In this section, a design process for an active electrode is described, resulting in full hardware implementation on a board, which is shown in Fig. 5. The upper limit of the input voltage on the active electrode is proportional to the power supply voltage of the constant current source, which is 48 V peak. The lower limit is not exactly stated. A suitable programmable gain amplifier (PGA) supporting one of the interfaces provided by the system bus was required. After assessing the options, it was determined that the PGA280 was the most fitting choice for the intended application.



Fig. 4. Schematic diagram of the active electrode

It has GPIOs that are used to control the attenuator and provide additional features, such as LED indication of active electrode modes directly on the probe. That feature was proven helpful in the debugging stage of the development. Commonly, these PGAs do not have the high-voltage power supply capabilities needed for our purpose. To solve the issue, a high voltage switch TMUX8108 (the same on the multiplexer module) is added before the PGA's input, selecting between R_1 to R_7 with R_8 forming resistor dividers to ensure attenuation from approximately 1 V/V to 0.28 V/V respectively, where the attenuation of $0.28 \,\mathrm{V/V}$ will decrease the maximum input voltage ± 48 V to a further processable value of 13.44 V by PGA, refer to the schematic diagram in Fig. 4. The resistor dividers are paced after high voltage OpAmp the OPA454 to ensure high input impedance of the active electrode, securing correct voltage measurement.

To further extend the active electrode's overall gain, another OpAmp was paced at the PGA's second input. Theoretically,



Fig. 5. Design render of the active electrode

with this gain network, the lowest selectable range that the active electrode could measure is $350 \,\mu\text{V}$.

Several voltage rails need to be provided for the ICs used in the active electrode, and several voltages are available on the system bus. It was decided to provide the high voltage rail ± 48 V as power for the electrodes as it was the highest needed voltage, and it is cheaper to use linear regulators instead of switching regulators needed to create larger voltage levels from lover voltage. The central system cannot provide other voltages because the cable connecting the active probe to the multiplexer unit lacks available wires. Due to the low power consumption of the circuitry, we can use the linear regulators without worrying about the efficiency and thermal stress of the components.

VI. CONCLUSION

After conducting market research, we selected the necessary components for our system, including the multiplexer module and active electrode. This system enables broader exploration of EIT methods, such as performing single-ended measurements with active electrodes. Ongoing research is being conducted to develop new applications for the method, and the system's modularity and wide range of potential use cases make it a promising option for future research initiatives.

REFERENCES

- [1] A. Adler and D. Holder, Electrical impedance Tomography: Methods, history and applications, Second edition. CRC Press, 2021.
- [2] A. McEwan, G. Cusick, and D. S. Holder, "A review of errors in multifrequency EIT instrumentation", Physiological Measurement, vol. 28, no. 7, pp. S197-S215, Jul. 2007.
- [3] "STM32TM's ADC modes and their applications", Application note, vol. 2010, no. AN3116, p. 9, 2010.
- [4] T. Piasecki, K. Chabowski, and K. Nitsch, 'Design, calibration and tests of versatile low frequency impedance analyser based on ARM microcontroller', Measurement, vol. 91, pp. 155–161, Sep. 2016, doi: 10.1016/j.measurement.2016.05.057.
- [5] P. Horowitz and W. Hill, The Art of Electronics Third Edition, Third edition. New York: Cambridge University Press, 2015.

Design of the Electronic Target for Shooting Sports and Sensor Suitability Analysis

Matej Grega

Department of Microelectronics, Faculty of Electrical Engineering and Communication Brno University of Technology Brno, Czech republic 240839@vut.cz

Abstract-Electronic scoring targets (ESTs) are designed to overcome the drawbacks of classic paper targets, particularly the inability to score individual hits in groups if they overlap and the time-consuming manual scoring process. This paper presents the design of a prototype of an acoustically based EST for 10m air pistol discipline and examines the suitability of microelectromechanical system (MEMS) microphones and sealed flexural ultrasonic transducers (FUTs) as hit point localization sensors. The proposed prototype of the EST is mobile and battery-powered, with built-in illumination and radiofrequency communication. The position of the hit point is calculated using a closed-form, combined weighted method based on time difference of arrival (TDOA) measurements. FUTs were used as sensors due to their filtering properties of shot and ambient noise and overall higher signal-to-noise ratio than MEMS microphones, without saturation of the output signal. The sensor positions for TDOA localization were accurately obtained using an iterative calibration method. The proposed EST prototype achieved a mean position error of 0.29 mm and a standard deviation of 0.19 mm for hit point localization.

Index Terms—hit point localization, electronic scoring target (EST), airgun, shooting, time difference of arrival (TDOA), multilateration, acoustic sensors, automatic shot scoring system

I. INTRODUCTION

Paper targets have been used for a long time, but they have many disadvantages. Firstly, each shot creates a hole in the paper target making it almost impossible to score individual shots if they are in groups of more than three and they overlap. Secondly, the shooter cannot see hits on the target with the naked eye, so it is necessary to use an additional optical device. This requires a change in shooting position which could negatively impact shooter's performance. Furthermore, scoring hits on conventional paper targets is time-consuming process as it must be done manually, one hit at a time, or by using a specialized scanner, one target at a time.

Electronic scoring target (EST) have been developed to mitigate the disadvantages of paper targets. They evaluate hits individually and display their position instantly on a monitor near the shooter. This allows for shooting at long distances and eliminates the problem of multiple hits at a single location. Additionally, precision shooting has become more attractive to spectators, and finals of competitions can now be organized with successive elimination of shooters. The primary drawback of commercially available ESTs is their cost ranging from approximately $800 \notin$ to $3300 \notin$. Therefore,

there was a motivation to develop, design and construct an affordable EST specifically for the International Shooting Sports Federation (ISSF) discipline of 10m air pistol for home use.

II. HIT POINT LOCALIZATION

The primary function of the EST is to determine the location of a hit on a target. This involves utilizing a localization principle, sensing necessary physical parameters, and calculating the coordinates of the hit point. The accuracy of the EST must be at least one half of a decimal scoring ring according to the ISSF General Technical Rules [1], which is ± 0.4 mm in the case of the ISSF 10m air pistol discipline.

A. Principle of Hit Point Localization

There are various principles to determine the location of the hit point, including the use of acoustic and mechanical waves, as well as optical localization. This paper describes the EST that uses the acoustic principle due to its sufficient accuracy and less critical mechanical construction of measuring parts compared to the optical principle.

Acoustic hit point localization is based on the detection of a sound wave by multiple sensors at known locations. From the time differences of arrival (TDOA) of the wave at each sensor, the position of the hit on the target is calculated. A membrane is located on the face of the EST to acoustically separate the external environment from the internal measurement space. The perforation of the membrane generates a sound wave from that point which propagates spherically into the space behind the membrane and is detected by sensors located in the frame of EST. For the principle of sound wave generation to work, the projectile must always hit the membrane. This is achieved using a black tape that moves behind the paper target or white paper mask after each shot, covering the holes.

B. Mathematical Method Based on TDOA

The arrival times (TOAs) of the sound wave at each sensor are measured. By subtracting the TOA of the reference sensor from the others, time differences of arrival (TDOAs) are obtained. Multiplying TDOAs by the speed of sound yields range differences of arrival (RDOAs). Each RDOA defines a possible location of the hit point on a hyperbola with the two sensors as foci [2]. This results in a set of nonlinear positioning equations:

$$vt_{m1} = \|\mathbf{u} - \mathbf{w}_m\| - \|\mathbf{u} - \mathbf{w}_1\|,$$
 (1)

where v denotes the speed of sound, t_{m1} is the TDOA between the *m*-th sensor and the reference sensor, \mathbf{w}_m and \mathbf{w}_1 are the position vectors of the *m*-th and reference sensors respectively, and \mathbf{u} is the position vector of the hit point which value is to be computed.

The optimization problem of computing coordinates of a hit point from TDOAs with measurement errors is wellknown and widely used in various fields such as GPS, passive radars, and sonars. There are two main groups of mathematical methods for this purpose: iterative algorithms and closedform algorithms. Paper [3] reviews numerous mathematical approaches and their properties. It is important to note that iterative positioning methods generally require a good initial position estimate, and convergence is not always guaranteed. Additionally, these methods can be time-consuming and require higher computing power. In contrast, closed-form algorithms are fast and lightweight, making them suitable for the proposed EST.

The Combined Weighted (COM-W) positioning method described in [4] is applied to calculate the coordinates of the hit point on the target. This method computes multiple preliminary results based on different minimal configuration combinations (three sensors) that are always consistent. A final result is then obtained as a weighted average based on the accuracy of estimation of preliminary results defined by Cramér-Rao lower bound (CRLB). The CRLB defines the maximum achievable accuracy of hit point position estimation, taking into account measurement uncertainties and the position of the hit point relative to the sensors.

III. DESIGN OF THE ELECTRONIC SCORING TARGET

Commercial ESTs are designed for fixed indoor installations, primarily in shooting clubs. They often require special one-purpose monitors to display the locations of hits. In contrast, the proposed EST is intended for more universal use among hobbyists. This means it is a mobile EST with no cables, a shooting range of up to 100 m, and the ability to display hits on a handheld device such as a tablet or mobile phone. The concept of proposed system is shown in Fig. 1.

The additional requirements for EST include an accuracy of hit point localization of at least 0.4 mm, detection of shots to the chassis, adjustable intensity and temperature of



Fig. 1. Concept of scoring system with proposed EST.

illumination, and a battery life of at least 3 hours on a single charge. The electronics of EST are divided into two printed circuit boards (PCBs): the sensor PCB and the main PCB. This division allows further upgrade of the sensor PCB and use EST also for other shooting disciplines. The block diagram of the proposed EST is shown in Fig. 2.

A. Main Printed Circuit Board

The primary function of the main PCB is to measure TOAs by a timer with input capture channels, compute hit point position and send result over universal serial bus (USB) and radiofrequency (RF) transceiver. The TDOA measurement resolution must be precise enough to not degrade assumed hit point localization accuracy. To achieve an equivalent RDOA resolution of 0.01 mm, the counter's clock frequency must be at least 35 MHz, assuming a speed of sound of 350 ms⁻¹. The designed EST uses two 16-bit counters in series, operating at 72 MHz, which are integrated into the STM32F103VC microcontroller. The result is transmitted by the RF transceiver RFM69HW and sent via USB.

A shot to chassis of the EST results in intense vibrations with a sharp peak at the beginning, which is detected, and the shot is scored 0 points. The vibrations are detected by a piezoelectric sensor mechanically coupled to the front panel of the EST.

Another important function of the main PCB is to move the black cover tape by a motor. The rated voltage of the motor should be as close to the battery voltage as possible since the EST is battery-powered. Therefore, a 6 V DC motor JGA25 370 with a built-in gearbox was used. Motor is PWM controlled at 15 kHz by a half H bridge with resettable fuse and overcurrent protection.

The LUXEON 2835E series of light emitting diodes (LEDs) with light temperatures of 6500 K and 2700 K were chosen for adjustable illumination. The illumination temperature is determined by the ratio of the relative brightness of the two



Fig. 2. Block schematics of proposed EST.

types of LEDs. This ratio is controlled by applying two voltages from the digital-to-analog converter to a LED driver TPS61150A. This driver integrates an electronically controlled current source with a step-up voltage converter.

To ensure versatility, the proposed EST should be independent of mains power supply. Therefore, a rechargeable Li-Ion 18650 battery with capacity of 3 Ah was selected. Battery protection include overvoltage protection (OVP), undervoltage protection (UVP), and overcurrent protection for both charging and discharging currents.

Designed EST is compatible with USB 2.0 data transfer and USB 3.0 and USB-C power delivery by implementing following functions:

- a default charging current of 0.5 A,
- an increase in the charging current by 0.9 A when the EST is powered from a USB-C type power supply,
- a decrease in the charging current by 0.5 A when the EST is powered on.

The USB-C standard supports a voltage of up to 20 V while the EST only supports 5 V. Therefore, additional OVP and OCP were designed.

The power section of the EST provides 6 V for motor by step-up converter and two separate 3.3 V branches for digital and analog parts by linear low dropout regulators (LDOs) with high power supply rejection ratio (PSRR).

B. Sensors for Hit Point Localization

The accuracy and reliability of hit point localization depends on the properties of the signal produced by sensors. Their nature also determines the requirements for the following electronics for signal processing. Therefore, the ideal sensor type is one that can most accurately and easily determine the exact moment of arrival of the sound wave, which is generated by a membrane perforation. It is desirable for the signal level to change steeply when the sound wave is detected, without a gradual rise. Additionally, the sensor should be insensitive to other signals, particularly ambient and shot noise.

Two types of acoustic sensors were examined: microelectromechanical system (MEMS) microphones and ultrasonic transducers. The advantage of MEMS microphones is a small diameter of acoustic port (0.3 mm to 0.8 mm), which allows precise determination of the sound wave sensing point coordinates required for mathematical methods. Four types of MEMS microphones were selected, and their output signals measured while shooting to paper membrane of EST:

- SPU0410LR5H-QB average performance (for reference purposes),
- ICS-40800 two acoustic ports for differential sensing,
- IMP23ABSUTR bandwidth of 60 kHz, designed for ultrasonic applications,
- IM73A135V01XTSA1 bandwidth of more than 80 kHz, differential output, IP57.

Fig. 3 shows the typical waveforms for individual microphones. Each waveform is a combination of three signal types. Firstly, the shot noise is captured while the projectiles are subsonic. Secondly, mechanical vibrations caused by EST's membrane perforation are captured, as they travel faster in solid materials than sound waves in the air. Finally, the sound wave of membrane perforation is captured, which is the only useful signal. In order to accurately demonstrate the response of individual sensors to this sound wave, waveforms with negligible vibration effects have been selected in Fig. 3. However, if there is direct mechanical coupling between the target and sensor, the vibration signal can reach an amplitude as high as 1.5 V. Additionally, the noise of the pellet trap is present after the detection of the sound wave.

In general, MEMS microphones have high sensitivity and an integrated preamplifier. Therefore, the amplitude of the output signal was sufficient, and the moment of sound wave reception was clearly visible on an oscilloscope. However, the slew rate of the signal varied among different types of microphones, which most probably depends on the properties of their preamplifiers. The signal-to-noise ratio (SNR) decreased due to high sensitivity and output voltage swing limitation. This is because ambient and shot noise generate signals with similar amplitudes and frequency contents as the useful sound wave, making it almost impossible to differentiate them in some cases.

The mean position error for hit point localization was 0.7 mm and 1.8 mm for SPU0410LR5H-QB and IM73A135V01XTSA1 microphones, respectively. Despite the drawbacks of MEMS microphones, they can still be used for hit point localization in a cost-effective manner. However, it is important to note that their use should be limited to indicative purposes only.

The second type of sensors examined were sealed flexural ultrasonic transducers (FUTs). FUTs are robust and tuned to a particular resonant frequency, making them insensitive to ambient noises and thus increasing SNR. Additionally, FUTs have the advantage of being able to fully swing the output signal without saturation, and they produce almost identical output signals across FUTs from different manufacturers with the same resonant frequency. However, the amplitude of the output signal is significantly lower than that of microphones, requiring the use of an amplifier.

FUTs with operating resonant frequencies of 40 kHz, 200 kHz, and 1 MHz were tested. The FUTs operating at 40 kHz produced the highest amplitude signal and were capable of detecting the hit point from the longest distance out of all the tested FUTs. As a result, FUTs of this type were used in the proposed EST. Fig. 3 shows the 100x amplified waveform of the FUT in comparison to those of MEMS microphones.

C. Sensor Printed Circuit Board

The Sensor PCB comprises a measuring channel for each sensor, settable voltage references, and a temperature and humidity sensor SHT41. Each measuring channel amplifies the input signal by a factor of 100 in two stages, 10x in each stage. The amplified signal is then compared with the reference voltage. The output of the comparator is a digital



Fig. 3. Waveforms of membrane perforation detected by different types of MEMS microphones and 40 kHz flexural ulrasonic transucer (FUT).

signal used for measuring TOA. In addition, a peak voltage of the 10x amplified signal is measured and used for software normalization of amplitudes that affect TOA measurements.

IV. CALIBRATION AND ACCURACY ANALYSIS

Due to the nature of TDOA, localization is highly sensitive to perturbations in sensor position [5]. Therefore, it is mandatory to have the exact positions of the sensors. As the sensor diameter is 16mm and the mounting angle is high, simple length measurement is not sufficient to obtain exact positions. Instead, a mathematical calibration method proposed in [5] was applied. The method is based on TDOA measurements of hits at known positions.

The proposed EST's accuracy was analyzed using 39 shots on millimeter paper. Ten evenly distributed shots were used for calibration, and Fig. 4 displays the computed and actual hit point positions of the fully calibrated EST. The mean position error (MPE) and standard deviation σ of MPE for the uncalibrated and calibrated EST can be found in Table I. The MPE achieved is likely to depend on the impact of mechanical vibration on the measured signal. Therefore, flexible or foam sensor mounting could potentially increase accuracy and should be investigated in the future.

V. CONCLUSION

This paper presents the design of the electronic scoring target (EST) and an examination of sensors for hit point localization. The proposed EST localizes the hit point by measuring the time difference of arrival (TDOA) of the acoustic wave generated by the projectile perforating the paper target. In conclusion, flexural ultrasonic transducers (FUTs) were used

 TABLE I

 MEAN POSITION ERROR AND ITS STANDARD DEVIATION

	raw	offset compensated	calibrated
MPE [mm]	1.05	0.45	0.29
σ [mm]	0.43	0.22	0.19



Fig. 4. Calculated and actual hit point positions after calibration.

instead of MEMS microphones because they are less sensitive to noise and do not saturate the output signal, resulting in a higher signal-to-noise ratio. Furthermore, performance of FUTs is consistent across different manufacturers. A prototype of the proposed EST has been built and it localizes the hit point with a mean error of 0.29 mm and a standard deviation of 0.19 mm after calibration. This level of accuracy is sufficient to practice the 10 m ISSF air pistol discipline, although the MPE is close to the limit set by the ISSF international rules. With a material cost of less than $300 \in$, which is a quarter of the mean cost of similar products, the proposed EST is an affordable alternative to commercially available ESTs.

ACKNOWLEDGMENT

I would like to express my sincere gratitude to my supervisor Ing. Pavel Tomíček for his professional advice and guidance throughout this project.

REFERENCES

- International Shooting Sport Federation, "6 General Technical rules", [Onine], Jan. 2023. Available: https://www.issfsports.org/getfile.aspx?mod=docf&pane=1&inst=458&file=ISSF_-Technical-Rules-Rule_Book_2023_Approved_Version.pdf
- [2] Y. -Tong Chan, H. Yau Chin Hang, and P. -chung Ching, "Exact and approximate maximum likelihood localization algorithms", *IEEE transactions on vehicular technology*, vol. 55, no. 1, pp. 10-16, 2006, DOI: 10.1109/TVT.2005.861162.
- [3] X. Li, Z. D. Deng, L. T. Rauchenstein, and T. J. Carlson, "Contributed Review: Source-localization algorithms and applications using time of arrival and time difference of arrival measurements", *Review of Scientific Instruments*, vol. 87, no. 4, pp. 041502-041502, 2016, DOI: 10.1063/1.4947001.
- [4] S. Cao, X. Chen, X. Zhang, and X. Chen, "Combined Weighted Method for TDOA-Based Localization", *IEEE transactions on instrumentation and measurement*, vol. 69, no. 5, pp. 1962-1971, 2020, DOI: 10.1109/TIM.2019.2921439.
- [5] D. Wang, J. Yin, X. Chen, C. Jia, and F. Wei, "On the use of calibration emitters for TDOA source localization in the presence of synchronization clock bias and sensor location errors", *EURASIP journal* on advances in signal processing, vol. 2019, no. 1, pp. 1-34, 2019, DOI: 10.1186/s13634-019-0629-1.

Enhancement of Vehicle Dynamics Through Adaptive Torque Vectoring Control with PMSM Powertrain

Sebastian Šimanský Department of Control and Instrumentation Brno University of Technology Brno, Czechia Sebastian.Simansky@vut.cz

Abstract—This paper presents an adaptive control torque vectoring algorithm designed for enhancing vehicle dynamics and stability through precise manipulation of individual wheel torques. The algorithm utilizes the vehicle's yaw rate as the primary input parameter to dynamically adjust the torque distribution among the wheels, thereby optimizing the usage of traction, overall manoeuvrability, and handling performance. By integrating adaptive control techniques, the proposed algorithm continuously adapts to varying driving conditions and vehicle dynamics, ensuring robust and effective torque vectoring across a wide range of scenarios. Lastly, the functionality of the developed algorithm is demonstrated through simulations with various driving scenarios.

Index Terms—Torque Vectoring, Formula Student, Yaw Rate control, Adaptive control, Simulink

I. INTRODUCTION

Formula Student is a worldwide racing competition where each year, teams from universities converge to design, build, and race formula-style cars in a thrilling showcase of automotive ingenuity. A typical characteristic of Formula Student tracks is that they are narrow and highly technical.

In this context, the core challenge addressed in this paper revolves around the absence of torque vectoring in the team TU BRNO Racing electric monopost called Dragon e4. This absence leads to a uniform distribution of torque to the wheels, regardless of their varying capabilities during cornering. When navigating corners, load transfer affects the capacity of each wheel to transmit power to the ground, as indicated by the concept of the friction circle [2]. In essence, without torque vectoring or a different solution such as Differential Braking, Steer-by-Wire [4] or Rear-wheel steering [6], drivers are unable to fully exploit the tyre's limit, thus compromising their ability to drive at the edge of traction. Hence this algorithm is utilized not only by Formula Student teams [3] but also by high-performance electric cars such as Rimac Nevera [5].

To address this challenge, a novel approach is employed, leveraging adaptive control alongside a feed-forward disturbance compensator.

A. Acronyms

• *CoG* ... Centre of gravity

- DoF ... Degrees of freedom
- TV ... Torque Vectoring
- 4WD ... Four wheel drive
- 2WD ... Two wheel drive
- *DiL* ... Driver in loop

II. TORQUE VECTORING PRINCIPLE AND THE RESULTING MATHEMATICAL MODEL

As mentioned in previous chapter, since each wheel experiences different load forces F_z , their cornering capabilities differ accordingly. Adapting on this tyre behaviour, the main goal of the presented algorithm is to increase the cornering speed. This can be achieved by increasing the yaw moment of inertia of the vehicle, thus making the vehicle rotate faster around its Z axis.

The yaw moment equilibrium equation is given by:

$$J_{zz}\psi = l_f F_{yf} - l_r F_{yr} \tag{1}$$

Where:

- J_{zz} ... Yaw moment of inertia $[kg \cdot m^2]$
- l_f ... Distance from CoG to the front axle [m]
- l_r ... Distance from CoG to the rear axle [m]
- F_{yf} ... Front axle forces [N]
- F_{yr} ... Rear axle forces [N]

Additional yaw moment contribution by TV is represented by:

$$M_{ZTV} = \Delta T \frac{t_r}{2} \tag{2}$$

- M_{ZTV} ... Additional yaw moment by TV [Nm]
- ΔT ... Left and right wheel torque difference [Nm]
- t_r ... Distance between left and right wheel [m]

Putting these equations together: $L = \frac{1}{2} M = \frac{1}{2} E = \frac{1}{2} E$

$$J_{zz}\psi - M_{ZTV} = l_f F_{yf} - l_r F_{yr} \tag{3}$$

This equation illustrates the relationship between torque applied to the left and right wheel and the resultant increase in yaw moment.

Based on Newton's second law and moment equilibrium (3). Following state-space model is presented [4], [3]:

$$\dot{x} = Ax + Bu \tag{4}$$

$$y = Cx + Du \tag{5}$$

$$A = \begin{bmatrix} -\frac{C_{\alpha f} + C_{\alpha r}}{mV_X} & -V_X - \frac{C_{\alpha f} l_f - C_{\alpha r} l_r}{mV_X} \\ \frac{-C_{\alpha f} l_f + C_{\alpha r} l_r}{J_{zz} V_X} & -\frac{C_{\alpha f} l_f^2 + C_{\alpha r} l_r^2}{J_{zz} V_X} \end{bmatrix}$$
$$B = \begin{bmatrix} \frac{C_{\alpha f}}{m} & 0 \\ \frac{C_{\alpha f} l_f}{J_{zz}} & \frac{t_r}{2 \cdot J_{zz}} \end{bmatrix}$$
$$C = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$D=0.$$

TABLE I System parameters

C_f	Cornering stiffness of front tire	600 [N/deg]
C_r	Cornering stiffness of rear tire	600 [N/deg]
m	Vehicle weight with driver	250 [kg]
J_{zz}	Yaw moment of inertia	109.1 [Nm]
V_x	X - axis vehicle speed	- [m/s]
V_y	Y - axis vehicle speed	- [m/s]
$\dot{\psi}$	Yaw rate	- [rad/s]
δ	Steering angle	- [deg]
t_r	Distance between rear wheels	1.17 [m]

The input and state vectors are: $u = \begin{bmatrix} \delta \\ \Delta T \end{bmatrix} x = \begin{bmatrix} \delta \\ \Delta T \end{bmatrix}$

A small caveat that this representation has is the non-linear behaviour of cornering stiffness, which indicates what forces the tyre can input to the ground in a lateral direction for a certain angle of slip. Hence the model is linearized around the chosen point.

Extraction of the transfer functions is done by following equation (6)

$$G(s) = C(sI - A)^{-1}B + D$$
 (6)

which results in a $2x^2$ transfer function matrix. The intention is to observe the behaviour of the yaw rate based on the steering angle and the torque difference by the TV.

•
$$G(s)_{2,2} = \frac{Y(s)_{\dot{\psi}}}{U(s) \Delta T}$$

• $G(s)_{2,1} = \frac{Y(s)_{\dot{\psi}}}{U(s)_{\dot{\psi}}}$

Derived transfer function based on the torque difference at $V_x = 1m/s$

$$G(s)_{2,2} = \frac{Y(s)_{\dot{\psi}}}{U(s)_{\Delta T}} = \frac{0.005362(s+1440)}{(s+4600)(s+\frac{6.42}{V_{\pi}})} \tag{7}$$

Assuming the seemingly equivalent pole and zero cancel out and neglecting the lost DoF, the resulting transfer function (8) is a first-order system where the pole is determined by the longitudinal velocity V_x

$$G(s)_{2,2} = \frac{0.005362}{(s + \frac{6.42}{V_r})}.$$
(8)

It was proved that this cancelation is possible for any given longitudinal velocity.

The transfer function based on the type steering angle δ at $V_x = 1m/s$ is given by:

$$G(s)_{2,1} = \frac{Y(s)_{\dot{\psi}}}{U(s)_{\delta}} = \frac{4.21(s+1440)^2(s+\frac{6.42}{V_x})}{(s+1440)^2(s+\frac{6.42}{V_x})^2}$$
(9)

III. ALGORITHN STRUCTURE AND COMPENSATOR DESIGN

Building on the equations (1),(2),(3) Simulink schematic shown in Figure 1 was developed.



Fig. 1. Model of the Controller

The input to the controller is the difference between the measured and desired yaw rate value. The output is the torque difference of the left and right rear wheels. Saturation blocks serve as a limitation for the actuator, in this case, it is the torque limit of the motor being 348 Nm.

To regulate the torque difference output of the algorithm, a PI controller is employed. The parameters of the PI controller, including the time constant Ti and gain K, are adjusted as functions of V_x to compensate for the specified pole of transfer function (8).

$$F_r(s) = \frac{K(Tis+1)}{Tis} \tag{10}$$

$$Ti = \frac{V_x}{6.45} \tag{11}$$

Static gain of controller is chosen through optimizing the value of quadratic integral criterion [1]

$$J_K = \int_0^\infty e(t)^2 dt \tag{12}$$

With the intention to keep the value of integral criterion problem-related. DiL simulation is used on an FSUK Nottingham track shown in Figure 7. After each lap the criterion value is calculated and the static gain is adjusted with the goal of minimizing this value, resulting in K = 318.75.

To further address the issue of the driver's capability to get to the tyre limit a feed-forward component of the algorithm is proposed. This component is based around the transfer function (9) since it is possible to take the driver's steering as a disturbance to the system. This feed-forward disturbance component is given by:

$$F_R(s)_U = -\frac{G(s)_{2,1}}{G(s)_{2,2}} = -785.154$$
(13)

Poles and zeros of both systems cancel out, and the resulting compensator is a static gain with a given value in (13). The gain of the disturbance compensator is then retuned, due to the particular overcompensation discovered from simulation, further improvement could potentially include retuning this gain with the help of the quadratic integral criterion (12).



Fig. 2. PI Controller with Anti-Windup

The windup problem caused by the integral part of the controller is solved by the Back-calculation method shown in Figure 2. The dynamic saturation point is set by maximal wheel torque and torque applied by the accelerator of the driver. Lastly, the tracking time constant is also dynamically changing with the adaptation of the controller.

The calculation of the desired yaw rate is based on the kinematics of the vehicle

$$\dot{\psi_{des}} = \frac{V_X}{(l_f + l_r) + K_u V_X^2} \delta, \tag{14}$$

where K_u is oversteer gradient [4].

$$K_{u} = \frac{ml_{r}}{(l_{f} + l_{r})C_{\alpha f}} - \frac{ml_{f}}{(l_{f} + l_{r})C_{\alpha r}}.$$
 (15)

- $K_u < 0$ Understeer behaviour
- $K_u > 0$ Oversteer behaviour
- $K_u = 0$ Neutral

Setting this gradient to a certain value will result in TV trying to compensate for the specific behaviour based on the sign of the gradient. Lastly, the final value of the desired yaw rate is dynamically saturated by

$$\dot{\psi_{lim}} = \frac{\mu g}{V_X}.$$
(16)

This equation introduces the limit to the desired value based on the friction coefficient of the tyre μ .



Fig. 3. IPG Carmaker linked with Simulink

IV. SIMULATIONS

The car is simulated via Simulink with the help of IPG Carmaker software. The software can communicate with the Simulink algorithm enabling DiL simulation.

The car model used for the simulation is the default Formula Student type car developed by IPG shown in Figure 3, with a remodeled powertrain to match the one of team TU BRNO Racing Dragon e4 monopost. This model error is accounted for and less impactful results of the TV algorithm are expected.

A. Slalom simulation



Fig. 4. Zoomed comparison of yaw rate - Slalom simulation

Figure 4 illustrates the comparison of yaw rate with and without TV. The angular velocity contributed by TV is relatively small, which aligns with expectations given the low static gain of the transfer function (8). In other words, the torque difference's contribution to angular velocity is anticipated to be minor compared to the driver's steering input.

Figure 5 illustrates the difference in wheel torques with and without TV engaged. Across each corner of the slalom, this difference in wheel torques amounts to $2\Delta T$. The controller's output, depicted by the green line, reaches its peak at approximately 20 Nm. Notably, the controller exhibits discontinuities at these peaks, attributed to the implementation of a dead zone aimed at filtering out steering inputs near 0 degrees



Fig. 5. Wheel torques comparison - Slalom simulation

of steering angle. However, this design feature introduces potential disturbances in control due to sudden changes. Future consultation with the drivers might result in removing this feature from the physical vehicle.



Fig. 6. G-G Diagram - Slalom simulation

Diagram 6 depicts measured points of achieved acceleration along the X and Y axes. The acceleration's contribution is notably observed along the Y axis, demonstrating the algorithm's ability to harness the added grip of the outer wheel for faster cornering acceleration.

B. Nottingham - FSUK 2018 simulation

While the short slalom track is representative of the capabilities of the controller, the functionality of the mentioned Anti-Windup method might not be visible, due to the track being symmetrical along the X axis. Thus the algorithm is simulated on uneven track FSUK Nottingham.

Figure 8 shows the comparison of torques and the controller output at the mentioned track. The output does not reach and most importantly does not stay saturated, causing no disturbance in vehicle dynamics.



Fig. 7. Nottingham FSUK 2018



Fig. 8. Nottingham FSUK 2018 - right and left wheel torques

V. CONCLUSION

The proposed algorithm adeptly utilizes the grip, enhancing the angular rate of the Z axis, along with the acceleration of both the X and Y axes. While the individual contributions of these factors may appear modest, it is essential to recognize that in motorsport, incremental enhancements often translate into more significant gains in faster lap times. It is noteworthy, however, that the observed impact of these enhancements may be somewhat subdued, due to the error between IPG car model and the bicycle model of the monopost Dragon e4.

Further improvements of this algorithm might rest in the incorporation of traction control to limit the desired torque of the wheels and the comparison between the contribution of TV with 2WD and 4WD. Overall, through various simulations, the capability of vehicle dynamics improvement was proven and the next step is the implementation of this algorithm to the physical car, which is one of the goals for the season 2024 of team TU BRNO Racing.

REFERENCES

- BLAHA, Petr a VAVŘÍN, Petr. Řízení a regulace I. Online, Skripta. Brno, Česká Republika: Vysoké Učení Technické v Brně. Dostupné z: https://moodle.vut.cz/pluginfile.php/585752/course/section/68732/bpcrr1.pdf. [cit. 2023-12-30].
- [2] PACEJKA, Hans B. *Tyre and vehicle dynamics*. 2nd ed. Oxford: Butterworth-Heinemann, 2006. ISBN 07-506-6918-7.
- [3] ANTUNES, João. *Torque Vectoring for a Formula Student Prototype*. Master thesis. Lisboa, Portugalsko: Instituto Superior Técnico, 2017.
- [4] RAJAMANI, Rajesh. Vehicle Dynamics and Control. Online. 2nd edition. Springer, 2012. ISBN 978-1-4614-1432-2. Dostupné z:
- [5] RIMAC ALL WHEEL TORQUE VECTORING. Online. *Rimac Newsroom.* Roč. 2016, s. 1. Dostupné z: https://www.rimac-newsroom.com/press-releases/rimac-automobili/rimac-all-wheel-torque-vectoring. [cit. 2024-04-01].
- [6] DOSTÁL, Marek. SYSTÉM ADAPTIVNÍHO ŘÍZENÍ ZADNÍ NÁPRAVY. Online, Diplomová práce, vedoucí Ing. Jiří Míša. Brno, Česká Republika: Vysoké Učení Technické v Brně, 2022. Dostupné z: https://www.vut.cz/www_base/zav_prace_soubor_verejne.php?file_id=241667. [cit. 2024-04-01]. https://ftp.idu.ac.id/wp-content/uploads/ebook/tdg/ TERRAMECHAN-ICONTECTOR DEPUTITION of the second sec

ICS%20AND%20MOBILITY/ epdf.pub_vehicle-dynamics-and-control-2nd-edition.pdf. [cit. 2023-12-26].

Shock severity comparison using Shock Response Spectrum and Pseudo-Velocity Shock Spectrum in LabVIEW

1st Štěrba Radek Department of Control and Instrumentation Brno University of Technology Brno, Czech Republic 230191@vut.cz

Abstract—This article deals with the comparison of Shock Response Spectrum (SRS) and Pseudo Velocity Shock Spectrum (PVSS). Publicly available data as well as values from an implemented experiment are used. The experiment consisted of acquisition of accelerations when a model car hit various braking pads. The data was recorded and evaluated in LabVIEW and for this purpose the PVSS needed to be implemented in same programing environment. The results show that the impact severity rating can be seen from SRS as well from PVSS, but on the 4CP paper the results are much easier to spot.

Index Terms—shock, shock response spectrum, pseudo velocity shock spectrum, LabVIEW, data acquisition

I. INTRODUCTION

A mechanical shock can be defined as a sudden and abrupt change in the state of motion of the components or particles of a body or environment. due to the sudden application of a relatively large external influence on the particles such as a blow or impact [1].

Shocks are all around us and it is important to know their severity. The study of shocks and their severity has its origins mainly in military and seismic engineering. One of the most famous cases of poorly designed impact systems is the crash of the space shuttle Columbia in 2003 [2].

There are several methods for assessing the severity of impacts. The **Shock Response Spectrum** (**SRS**) has become the standard in impact assessment. Shock response spectrum analysis is the maximum response of a series of single-unit systems with different natural frequencies and with the same damping to a given excitation shock. It is used, for example, by NASA [3] in the analysis of the survival of payloads in spaceships when shocks caused by pyrotechnic bolts occur. Another major standard that uses SRS for the evaluation of shaking was developed by the US Department of Defense for military purposes [4].

The second metric used especially in recent years is the **Pseudo-Velocity Shock Spectrum** (PVSS). The PVSS can be derived from the SRS:

$$v_{max} = \frac{a_{max}}{\omega_n} [\frac{m}{s}] \tag{1}$$

2nd Stanislav Pikula Department of Control and Instrumentation Brno University of Technology Brno, Czech Republic pikula@vut.cz

where v_{max} is maximal pseudo velocity, a_{max} is maximal acceleration in SRS and ω_n is the $2 \cdot \pi \cdot natural frequency$ of the system [7]. Along with traditional SRS using absolute acceleration, PVSS is one of the most commonly used methods. It is primarily used in the United States Navy and is part of standard [8]. It is suggested that PVSS is much more advantageous than traditional SRS for the analysis of shock severity.

The mechanical stress is proportional to the modal velocity, as describes Gaberson [10] by the equation:

$$\sigma_{max} = \rho c v_{max} [Pa] \tag{2}$$

Where: $\sigma_{max} =$ is the maximum stress anywhere in the bar, c[m/s] is wave speed and $\rho[\frac{kg}{m^{-3}}]$ is mass density.

Since the tabulated values of the maximum stress are known, it is possible to determine the maximum velocity v_{max} and check the PVSS figure for exceeding this value.

PVSS is often displayed on 4CP paper. In addition to the pseudo-velocity at the SDOF natural frequency, the maximum deflection per excitation shock as well as the maximum acceleration per excitation shock can be read from the figure. The maximum deflection is determined by the asymptote at low frequencies and the maximum acceleration by the asymptote at high frequencies.

This paper will compare SRS and PVSS methods. Methodology of comparison as well as results and discussion follows in sections bellow.

II. METHODS

The main objective of this paper is to compare the severity of impacts by means of SRS and PVSS in LabVIEW which is often used for data acquisition. Following section describes implementation of SRS and PVSS in LabVIEW. Folowing sections decribe used data (publicly available data and data obtained from experiment). The publicly available data were also used to verify the correctness of the PVSS implementation, i.e. to compare our results with those available in the literature.

A. SRS and PVSS implementatiton in LabVIEW

The implementation of SRS in programming environments is done by approximating the SDOF response transfer function to the excitation shock by a digital filter with given coefficients, where the transformation is the system natural frequency and damping. A possible implementation is shown for example in [5]. In LabVIEW, SRS is already implemented in the Sound and Vibration package [6].

Unlike SRS, PVSS is not implemented in LabVIEW. The implementation of PVSS can be done in two ways:

- 1) PVSS by filtration: the sampled acceleration data enters the filter with coefficients given by the standard [11]
- 2) SRS values are divided by $2 \cdot \pi \cdot f_n$ frequency domain integration [7].

Since both methods give identical results except for the very low frequency region, and the second method is straightforward for offline data analysis, we used the second method.

We then implemented method to plot the 4CP paper in LabVIEW. We realised three possibilities: plot data via Matlab, via Python or directly into LabVIEW figure. Last option can be seen in figure 1.



Fig. 1. Example of 4CP figure plotted directly in LabVIEW

B. Publicly available data

We used time histories of four acceleration shocks (Motorcycle crash, Gun-stock, Punching bag, Elcentro earthquake), which are publicly available [1]. The time histories of data are shown in figure 2. One can see, that higher acceleration shocks are shorter, and longer shocks have less acceleration. Determining impact severity solely from acceleration time histories is impossible; SRS and PVSS will reveal shock severity later.

C. Experiment

To get real life data for evaluation of severity of impacts, we used lineup of a laboratory assignment of the course Sensors at Department of Control and Instrumentation [9]. The experiment uses a toy car hitting a brake pad. The design of the experiment is described in Figure 3. A model car is pulled by a weight towards the brake pads. It is critical to maintain



Fig. 2. Acceleration time history of publicly available data.

the same impact velocity for all pads. This is achieved by maintaining the same launch distance and the same weight. The weight has a mass of 70 g and the distance chosen is 20 cm in order not to exceed the range of the acceleration sensor.



Fig. 3. Experiment layout based on laboratory assignemnt of the course Sensors at Department of Control and Instrumentation [9].

1) Used braking pads: The exact properties of materials used are not known, but the materials were selected to have different stiffness and therefore probably different impact braking influence. Used materials:

- 1) (carpet) 5x carpet pad 5 mm (25 mm total)
- 2) (PS) Polystyrene 20 mm
- 3) (PE1) Polyethylene foam 1 20 mm
- 4) (PE2) Polyethylene foam 2 25 mm
- 5) (Lino) 6x linoleum pad 3 mm (18 mm total)

Photos of the materials are in figure 4.

2) Data acquisition: The measuring system consists of two main parts: the acceleration sensor ADXL150 (single-axis MEMS accelerometer with a range of \pm 50 g and a sensitivity of 38 mV/g) and measurement card NI-9234 in NI-cDAQ 9181 slot with USB connection to PC. Used sampling rate was 51.2 kS/s.

Acquisition software was realised on the basis of the LabVIEW continuous measurement and logging tamplate. The acceleration data are stored in TDMS file. Because measurement system doesn't contain hardware trigger, we implemented software trigger by means of measured acceleration treshold. In the postprocessing, a .csv file is created from the TDMS data with a column of timestamps and corresponding



Fig. 4. Photos of used braking pads: 1) carpet 2) Polystyrene 3) Polyethylene 1 4) Polyethylene 2 5) Linoleum.

accelerations which is the same format as the one used in the publicly available data.

III. RESULTS AND DISCUSSION

Acceleration data acquired from experiment are in figure 5.



Fig. 5. Time history of acceleration of the experimental impact into the brake pads.

Maximum acceleration values are given in the table I. The highest time history acceleration is achieved by the impact with pad *carpet*, while the lowest acceleration is achieved by the impact with pad *PE1*.

Braking pad	Peak Accel [g]	Peak accel $\left[\frac{m}{s^2}\right]$	max PVSS $\left[\frac{m}{s}\right]$
Carpet	22.5	220.6	0.90
PS	18.4	180.4	0.97
PE1	7.7	75.5	0.80
PE2	9.9	97.0	0.85
Lino	25.9	253.9	0.75

TABLE I MAX. ACCEL VALUE

A. SRS impact severity evaluation

Figure 6 shows the SRS for the publicly available acceleration data shown previously in figure 2. The SRS for individual shocks shows that for low natural frequencies the SDOF dampens shocks at low frequencies. However, the low frequencies are different for each shock, for the Elcentro earthquake the low frequencies are up to 0.02 Hz. Consequently, it can be observed that the maximum acceleration reaches values up to 3.3 g at 0.06 Hz. From the time course, the maximum acceleration is only 0.43 g, at these frequencies the shock signal is amplified. This is the resonance frequency region. For the other shocks the same similar principle can be observed, but at different frequencies.



Fig. 6. SRS of publicly available data. Note: highest SRS calculated frequency is ten times lower than sampling frequency, therefore the SRS of Elcentro earthquake is only displayed up to 70 Hz.

SRS of experimentaly acquired data are shown in figure 7. It is possible to observe the same principles as for the publicly available data. At low frequencies the shock is damped, then there is a region around the resonance frequency where the shock is amplified and reaches acceleration values higher than the maximum amplitude of the acceleration time history. It then sets up a simple transmission where, for example, for a carpet, a maximum acceleration of 22 g can be read out. This corresponds to the maximum value of the acceleration time course. The same can be said for the other pads.



Fig. 7. SRS of impact into different pads

B. PVSS impact severity evaluation

The figure 8 shows one crucial thing. Despite the fact that the time course of the earthquake reaches the smallest values of maximum acceleration, this shock is the most severe because the PVSS plateau reaches the highest values of pseudo-velocity.

Figure 9 with PVSS for the pads shows the same situation as for the previous shocks. In the introduction we insist that the maximal pseudo-velocity is a measure of the severity of the shock (because the pseudo-velocity is proportional to the mechanical stress). The shock where the damping pad was PS appears to be the most severe. A look at the shock time histories does not show the highest acceleration for PS but for Lino.

It is evident that this observation is also reflected in the SRS. One needs to watch the left asymptote, where SRS of pad *PS* have biggest SRS value among tested pads around 30 Hz. In comparison this result can be much more easily seen in PVSS figure.



Fig. 8. PVSS of publicly available data



Fig. 9. PVSS of impact into different pads. Note: The deflection asymptote for PVSS is only for shocks with a limit approaching zero mean value. Since these shocks do not have a zero mean, the asymptote for the deflection is not plotted

IV. CONCLUSION

The article showed the implementation of a 4CP plot in LabVIEW for display of PVSS in section II-A. It was checked on publicly available data that the implementation of the PVSS calculation in LabVIEW was done correctly. Furthermore, an experiment was described to verify the hardness of the impact on different brake pads.

It was shown that a most severe impact may not have a maximum acceleration. This is typical of the earthquake shock, which reaches much smaller deflections than other shocks, but according to PVSS is the most severe. For the experiment, the most serious shock occurred when using a polystyrene brake pad, despite the fact that the time history did not reach the highest acceleration.

It has been shown that PVSS is more suitable for determining the severity of an impact because the velocity is proportional to the mechanical stress. This severity can also be determined from the traditional SRS, but using an asymptote. From this point of view, the analysis using PVSS seems easier.

REFERENCES

- Alexander, J. Edward. Shock Response Spectrum A Primer. Sound and Vibration 43 (2009): 6-15.
- [2] Sisemore, C., & Babuška, V. (2020). The science and engineering of mechanical shock. In Springer eBooks. https://doi.org/10.1007/978-3-030-12103-7
- [3] National Aeronautics and Space Administration. NASA-STD-7003a, Pyroshock test criteria.
- [4] Department of Defense test method standard. MIL-STD-810H, Environmental Engineering Considerations and Laboratory Tests.
- [5] Tuma, Jiri a Koci, Petr. Calculation of a shock response spectrum. Online. 2011 12th International Carpathian Control Conference (ICCC). 2011, s. 404-409. ISBN 978-1-61284-360-5.
- [6] Emerson. SVT Shock Response Spectrum VI. Online. https://www.ni.com/docs/en-US/bundle/labview-sound-and-vibrationtoolkit/page/sndvibtk/shock_resp_spec.html
- [7] Gaberson HA. *Shock severity estimation*. Sound & vibration. 2012 Jan 1:46(1):12-20.
- [8] ANSI/ASA S2.62-2009 Test Requirements for Equipment in a Rugged Shock Environment.
- [9] Instructions for the laboratory task *Measurement of Vibrations* of course Sensors 2023, Departement of Control and Instrumentation, Faculty of Electrical Engineering and Communication, Brno University of Technology
- [10] Gaberson, H. A. and Chalmers, R. H., Modal Velocity as a Criterion of Shock Severity. Shock and Vibration Bulletin, Vol. 40, Part 2, pp 31-49, Dec. 1969.
- [11] ČSN ISO 18431-4. Vibrace a rázy Zpracování signálů Část 4: Analýza spektra rázové odezvy. 12/2007.



VÁHÁŠ, JAKÝM SMĚREM SE VYDAT?

VRHNI SE DO SVĚTA POLOVODIČŮ A UTVÁŘEJ NOVÉ TECHNOLOGICKÉ TRENDY!



onsemi je světový technologický lídr v oblasti polovodičů, a to od jejich vývoje až po finální výrobu součástek. Jsme mezinárodní společnost s bohatou historií v České republice, kde působíme v Rožnově pod Radhoštěm a Brně. Zabýváme se návrhem integrovaných obvodů, výrobou desek z křemíku i karbidu křemíku a také výrobou polovodičových čipů, ale také samotným vývojem technologií výroby polovodičových materiálů i polovodičových součástek.

~60 stážistů za rok
 zaměstnanců v Rožnově a Brně
 2200

STUDUJEŠ OBOR V OBLASTI?

elektro

chemie

PŘIJĎ K NÁM NA STÁŽ

nateriály

- Můžeš k nám na letní a celoroční stáž Jako stážistu tě ohodnotíme 200–280
 Kč/hod. (podle ročníku a zkušeností)
- Můžeš získat stipendium až **5 000 Kč/měs**.

fvzika

- Napiš si u nás svou bakalářskou nebo diplomovou práci
- U dlouhodobých stáží si můžeš vyzkoušet práci v různých týmech
- Zajistíme ti bezplatné ubytování (Rožnov p. R.) nebo příspěvek na bydlení (Brno)

PŘIJĎ K NÁM PRACOVAT

Nabízíme ti velmi atraktivní startovní plat a zrychlený platový růst (navýšení platu až o 20 % po 18 měsících)

IT

technologie

- 25 dní dovolené + 3 dny sick days + bridge days
- Přispějeme ti na dopravu, bydlení, stravování, penzijko do výše 6 % z tvého platu a dostaneš od nás výhodný mobilní tarif
- Podpoříme tvé další vzdělávání, jazykové dovednosti, možnost získat zkušenosti v zahraničí a mnoho dalších výhod

Zažij, jaké je to být součástí moderní firmy.



Výroba polovodičových materiálů (SiC a Si)



Špičkově vybavené laboratoře



Výroba čipů



Středisko návrhu součástek







Vlastní výzkum a vývoj





Graph Neural Networks in Epilepsy Surgery

Valentina Hrtonova Department of Biomedical Engineering Department of Biomedical Engineering FEEC, Brno University of Technology Brno, Czech Republic

xhrton02@vutbr.cz

Marina Filipenska FEEC, Brno University of Technology Institute of Scientific Instruments of CAS Brno, Czech Republic ronzhina@vut.cz

Petr Klimes Department of Medical Signals Brno, Czech Republic klimes@isibrno.cz

Abstract—Epilepsy surgery presents a viable treatment option for patients with drug-resistant epilepsy, necessitating precise localization of the epileptogenic zone (EZ) for optimal outcomes. As the limitations of currently used localization methods lead to a seizure-free postsurgical outcome only in about 60% of cases, this study introduces a novel approach to EZ localization by leveraging Graph Neural Networks (GNNs) for the analysis of interictal stereoelectroencephalography (SEEG) data. A GraphSAGE-based model for identifying resected seizure-onset zone (SOZ) electrode contacts was applied to a clinical dataset comprising 17 patients from two institutions. This study uniquely focuses on the use of interictal SEEG recordings, aiming to streamline the presurgical monitoring process and minimize risks and costs associated with prolonged SEEG monitoring. Through this innovative approach, the GNN model demonstrated promising results, achieving an Area Under the Receiver Operating Characteristic (AUROC) score of 0.830 and an Area Under the Precision-Recall Curve (AUPRC) of 0.432. These outcomes along with the potential of GNNs in leveraging the patient-specific electrode placement highlight their potential in enhancing the accuracy of EZ localization in drug-resistant epilepsy patients.

Index Terms—graph neural networks, deep learning, epilepsy, intracranial EEG, epileptogenic zone, seizure-onset zone, interictal biomarkers

I. INTRODUCTION

Epilepsy stands as a formidable challenge to global health, affecting an estimated 70 million people worldwide with recurrent seizures that severely impact brain function. Despite the wide availability of antiseizure medications, about 30-40% of patients fail to achieve adequate seizure control, underscoring the need for more definitive interventions such as epilepsy surgery. In the presurgical phase, diagnostic modalities such as stereo-electroencephalography (SEEG) are indispensable for the precise delineation of the epileptogenic zone (EZ). Patients undergoing SEEG must tolerate the invasive nature of the procedure, face associated risks, and endure the process of reducing anti-seizure medication to minimize seizures during monitoring lasting up to 4 weeks. Moreover, only around 60% of well-selected patients are seizure-free after surgery, partially due to the inability to correctly identify the EZ [1] highlighting a critical need for the development of more sophisticated diagnostic tools capable of localizing the EZ with great speed and precision from SEEG data.

Supported by project nr. 22-28784S: Financed by the Czech Science Foundation.

The Graph Neural Networks (GNNs) in particular have emerged as promising tools for more efficient and effective analysis of intracranial EEG in epilepsy research [2], [3]. The representation of signals as vertices within a graph, which enables GNNs to capture the implicit topological and functional relationships between signals, is especially relevant in the study of complex epileptic networks.

This study proposes an innovative approach by utilizing well-established interictal biomarkers - namely interictal epileptiform discharges (IEDs) and relative entropy [4], [5] - in conjunction with GNNs to localize the EZ, therefore, addressing the challenges associated with prolonged SEEG monitoring. By employing the GraphSAGE [6] inductive framework to analyze interictal SEEG recordings, we aim to provide a more efficient and effective approach to EZ localization which would be able to reflect the patient-specific electrode placement. This approach could potentially streamline the presurgical workup with the ultimate goal of enhancing treatment options for patients facing drug-resistant epilepsy.

II. METHODS

A. Patients

All consecutive adult patients with drug-resistant focal epilepsy who underwent SEEG and subsequent resective surgery at St. Anne's University Hospital in Brno (SAUH) between 3/2012 and 3/2022 and the Montreal Neurological Institute & Hospital (MNI) between 1/2010 and 12/2015 were analyzed for the study. The inclusion criteria for the patient cohort were: (i) availability of SOZ and resected area information; (ii) availability of ≥ 24 hours of continuous SEEG recording for sleep staging; (iii) availability of scalp EEG, electro-oculography, and electromyography or video for sleep staging; and (iv) good postsurgical outcome (defined as Engel IA-ID) after a minimum follow-up period of 1 year. The final patient cohort consisted of 17 patients (6 from SAUH, and 11 from the MNI). The study was approved by the SAUH Research Ethics Committee and the MNI Ethics Review Board. All patients granted written informed consent.

B. Data

1) SEEG Recordings: Patients in the study received stereotactic depth electrode implantation for evaluating drugresistant focal epilepsy. SAUH used platinum electrodes from DIXI or ALCIS, recording SEEGs with a 192-channel system,



Fig. 1. Schematics of the Graph Neural Network model pipeline. The stereo-electroencephalography (SEEG) recordings are pre-processed (1) and features (relative entropy and interictal epileptiform discharge rate) are extracted (2). Next, the data of each patient is represented in a separate graph structure (3). The entire dataset is used for training and testing of the model (4) with leave-one-patient-out cross-validation. During training, graphs of training subjects are passed through the GNN, and the optimizer updates model parameters based on calculated loss. During testing, the graph of the left-out patient is fed into the GNN and a probability of each node being in the resected seizure-onset zone (SOZ) is predicted. The model is evaluated (5) using the Receiver Operating Characteristic (ROC) and Precision-Recall (PR) curves.

low-pass filtered and downsampled to 5 kHz. MNI recorded SEEGs using MRI-compatible electrodes from DIXI or inhouse, with the Harmonie system, filtered at 500 Hz and sampled at 2 kHz.

2) Electrode Contact Annotations: The seizure-onset electrode contacts were determined based on the earliest changes at seizure onset irrespective of the fast activity content by boardcertified epileptologists. Pre- and post-resection MRI was used to identify areas of the brain and marking of resected electrode contacts. The localization of electrode contacts was done using the SEEGAtlas plugin of the openly available software IBIS [7].

C. EZ Localization Pipeline

Fig. 1 shows a visualization of the proposed EZ localization pipeline.

1) Signal Pre-Processing: Visual sleep staging was conducted in 30-second epochs and segments corresponding to 5 minutes of N3 sleep [8], [9] were selected from the complete recording. The segments were required to be at least 1 hour away from epileptic seizures. The selection of channels for analysis was performed with the exclusion of channels from contacts with locations outside of the brain, and channels of poor recording quality based on visual inspection and automated detection [10]. An average from all SEEG contacts with a confirmed location inside the brain was used as an SEEG reference and subtracted from each SEEG signal to suppress far-field potentials caused predominantly by volume conduction.

2) Feature Extraction: The ElectroPhYsiology COmputational Module (EPYCOM) [11] open-source Python library was used for the extraction of features from SEEG signals. IEDs were detected using the Janca detector [12], and the number of IEDs per minute was calculated and normalized to fall between 0 and 1 for each patient. REN was calculated in the low gamma band (20-45 Hz) between all channel pairs X and Y for each patient as $REN = \sum_{i=1}^{n} p(X_i) \cdot \log\left(\frac{p(X_i)}{p(Y_i)}\right)$, where p(X) and p(Y) are the probability distribution functions of X and Y. REN values were normalized to fall between 0 and 1 and a functional connectivity matrix was constructed for each patient. To prevent over-smoothing which occurs in densely connected graphs, REN values below the empirically set threshold of 0.2 were set to zero.

3) Graph Representation: A data object was constructed for each patient with IED rate as the node feature, REN as the edge feature, and patient ID as a graph attribute. The resected SOZ in patients with a good postsurgical outcome was used to approximate the EZ and binary labels were assigned to each electrode contact, corresponding to 1 for SOZ contacts resected during epilepsy surgery ("Resected&SOZ") and 0 for the remaining contacts.

4) GNN Classifier: The architecture of the model, based on a GraphSAGE neural network, was designed for node classification tasks. The model was implemented using the PyTorch Geometric (PyG) library and it consisted of a GraphSAGE layer followed by a linear classifier.

The **GraphSAGE layer** transformed node features, edge indices, and edge weights into new node embeddings with the use of a SAGEConv operator for message passing within the graph. The message passing was performed with sampling, aggregation, and concatenation steps as

$$h_i^k \leftarrow \sigma \left(W_1 h_i^{k-1} + W_2 \cdot \text{MEAN} \left(\{ h_j^{k-1} \} \cup \{ h_j^{k-1}, \forall j \in N_i^{k-1} \} \right), \forall j \in N_i^{k-1} \quad (1)$$

where h_i^k are new embeddings for node i, h_i^{k-1} embeddings for node i from a previous layer, $\{h_j^{k-1}, \forall j \in N_j^{k-1}\}$ are node embeddings for neighbourhood N_j^{k-1} of node i, W_1 and W_2 are trained weights, and σ is a nonlinear activation function. The number of hidden channels in the GraphSAGE layer, which controls the dimensionality of the node embeddings during the message-passing process, was set to 4, and the Rectified Linear Unit (ReLU) was used as the activation function. The model used a single GraphSAGE messagepassing layer.

The node embeddings generated by the GraphSAGE layer were input into a **linear classifier** with a **sigmoid activation function** for binary classification, which can be expressed as

$$\hat{y} = \sigma(W_{\text{linear}} \times h + b_{\text{linear}}), \qquad (2)$$

where \hat{y} are the predicted probability scores, σ is the sigmoid activation function, W_{linear} represents the weight matrix for the linear layer, h denotes the node embeddings generated by the GraphSAGE layer, and b_{linear} is the bias for the linear layer. The output scores indicate the likelihood of each node belonging to the resected SOZ.

5) Model Training and Evaluation: The model was trained and tested with leave-one-patient-out cross-validation on the complete dataset of 17 patients, meaning that the model was trained on data from all patients except one, and tested on withheld patient data.

A new model was trained in 400 training epochs for each cross-validation fold. The loss was calculated using the Binary Entropy loss function with class weights calculated as inversely proportional to class frequencies to correct for class imbalance in the dataset (resected SOZ contacts are underrepresented to the rest of electrode contacts). The Adam optimizer was employed with $\alpha = 0.001$, $\beta_1 = 0.9$, and $\beta_2 = 0.999$ to update the model's parameters.

Model performance was comprehensively analyzed with the Receiver Operating Characteristic (ROC) curve and the Precision-Recall (PR) curve, focusing on their areas under the curve (AUROC and AUPRC, respectively) to evaluate the trade-offs between sensitivity (also called "true positive rate" or "recall"), specificity (1 - false positive rate), and precision at various thresholds. The AUROC provides a comprehensive measure of the model's discriminative ability, while the AUPRC, more sensitive to class imbalance, reflects the precision-recall balance. AUROC and AUPRC were calculated using standard methods, with the chance level at 0.5 for AUROC and the proportion of positive samples for AUPRC.

III. RESULTS

The GNN model for EZ localization was evaluated separately for each patient and the results were averaged. In Fig. 2, the ROC and PR curves for each patient, as well as the average curves are plotted. The model achieved a mean AUROC of 0.840 ± 0.22 standard deviation (std) and AUPRC of 0.431 ± 0.29 std. The chance level for AUROC is 0.5, whereas for AUPRC, the chance level performance for this dataset is 0.044.

IV. DISCUSSION

This study's findings underscore the potential of GNNs in enhancing the localization of the EZ from interictal SEEG data. The proposed model's performance, surpassing chance levels in both the AUROC and the AUPRC, demonstrates its capacity to differentiate between resected SOZ contacts and the remaining contacts (contacts which were not resected SOZ). Notably, the ROC and PR curve analyses (Figure 2) reveal exemplary performance for some patients, while indicating a need for model refinement to enhance generalizability across diverse patient datasets, as evidenced by performances near or below random classification in some cases.

Contrasting with previous efforts, such as the multi-feature Support Vector Machine (SVM) approach for localizing epileptogenic tissue from interictal SEEG data [4], which reported an AUROC of 0.838, our method integrates GNNs to leverage the complex topological structure of SEEG signals. This approach offers a nuanced understanding of the dynamics of epileptic networks, potentially uncovering patterns that conventional models, including SVMs, may overlook. Although our results do not surpass all existing models, with another study achieving an AUPRC of 0.480 using SVM on 5-minute N3 sleep segments [9], the novel use of GNNs presents a significant advantage by efficiently analyzing the networked nature of epileptic signals.

The novel application of GNNs for EZ localization, particularly using only interictal data, marks a significant distinction from prior research, such as Grattarola et al.'s work [3], which focused on ictal intracranial EEG recordings using an attention-based GNN network. Our strategy avoids the complexities and limitations associated with seizure recording and detection by employing only interictal segments, specifically a single 5-minute segment from N3 sleep for each patient, suggesting a pathway to significantly reduce the duration of presurgical monitoring. To our knowledge, this is the first study to explore EZ localization using GNNs exclusively with



Fig. 2. Receiver Operating Characteristic (ROC) (A) and Precision-Recall (PR) curves (B) for the GNN model. Curves for each patient are visualized in gray, average curve in blue.

interictal data, potentially setting the stage for a more efficient and less invasive presurgical evaluation process for epilepsy patients.

However, the observed variability in model performance across patients underscores the imperative need for further research. Enhancing the robustness and adaptability of GNN models, through methods such as architecture optimization and the integration of novel features, remains crucial. Future investigations should also consider larger, more heterogeneous datasets to refine the understanding of epilepsy's underlying mechanisms and improve EZ localization accuracy.

V. CONCLUSIONS

This research underscores the promising capabilities of GNNs in the precise localization of the EZ, critical for achieving good outcomes in epilepsy surgery that often offers the last hope of seizure freedom to patients with drug-resistant epilepsy. By potentially reducing SEEG monitoring durations and enhancing the precision of EZ localization, this research aspires to provide tools for a more timely and effective intervention for these patients, enhancing both efficiency and patient outcomes in epilepsy surgery meanwhile reducing costs and risks associated with prolonged SEEG monitoring.

REFERENCES

- FRAUSCHER, Birgit, 2020. Localizing the epileptogenic zone. Current Opinion in Neurology. Vol. 33, no. 2, pp. 198–206. DOI 10.1097/WCO.00000000000790.
- [2] LIAN, Qi et al., 2020. Learning graph in graph convolutional neural networks for robust seizure prediction. Journal of Neural Engineering. Vol. 17, no. 3, p. 035004. DOI 10.1088/1741-2552/ab909d.

- [3] GRATTAROLA, Daniele et al., 2022. Seizure localisation with attentionbased graph neural networks. Expert Systems with Applications. Vol. 203, p. 117330. DOI 10.1016/j.eswa.2022.117330.
- [4] CIMBALNIK, J. et al., 2019. Multi-feature localization of epileptic foci from interictal, intracranial EEG. Clin Neurophysiol. Vol. 130, no. 10, pp. 1945–1953. DOI 10.1016/j.clinph.2019.07.024.
- [5] TRAVNICEK, Vojtech et al., 2023. Relative entropy is an easy-to-use invasive electroencephalographic biomarker of the epileptogenic zone. Epilepsia. Vol. 64, no. 4, pp. 962–972. DOI 10.1111/epi.17539.
- [6] HAMILTON, William L., YING, Rex and LESKOVEC, Jure, 2018. Inductive Representation Learning on Large Graphs [online]. Retrieved from : http://arxiv.org/abs/1706.02216 [accessed 23 December 2023]. arXiv:1706.02216
- [7] ZELMANN, Rina et al., 2023. SEEGAtlas: A framework for the identification and classification of depth electrodes using clinical images. Journal of Neural Engineering. Vol. 20, no. 3, p. 036021. DOI 10.1088/1741-2552/acd6bd.
- [8] KLIMES, P. et al., 2019. NREM sleep is the state of vigilance that best identifies the epileptogenic zone in the interictal electroencephalogram. Epilepsia. Vol. 60, no. 12, pp. 2404–2415. DOI 10.1111/epi.16377.
- [9] CHYBOWSKI, Bartlomiej et al., 2024. Timing matters for accurate identification of the epileptogenic zone. Clinical Neurophysiology. Vol. 161, pp. 1–9. DOI 10.1016/j.clinph.2024.01.007.
- [10] NEJEDLY, Petr et al., 2019. Intracerebral EEG Artifact Identification Using Convolutional Neural Networks. Neuroinformatics. Vol. 17, no. 2, pp. 225–234. DOI 10.1007/s12021-018-9397-6.
- [11] CIMBALNIK, J. et al., 2019. Epycom: ElectroPhYsiology COmputational Module [software]. Version 1.3.4. 2019. Retrieved from : https://github.com/ICRC-BME/epycom
- [12] JANCA, Radek et al., 2015. Detection of interictal epileptiform discharges using signal envelope distribution modelling: application to epileptic and non-epileptic intracranial recordings. Brain Topography. Vol. 28, no. 1, pp. 172–183. DOI 10.1007/s10548-014-0379-1.

Quantitative analysis of Scratch assay and Colony formation assay data using MATLAB

1st Katerina Ingrova Department of Biomedical Engineering Faculty of Electrical Engineering and Communication Brno University of Technology Brno, Czech Republic 219244@vut.cz 2nd Larisa Chmelikova Department of Biomedical Engineering Faculty of Electrical Engineering and Communication Brno University of Technology Brno, Czech Republic chmelikoval@vut.cz 3rd Inna Zumberg Department of Biomedical Engineering Faculty of Electrical Engineering and Communication Brno University of Technology Brno, Czech Republic zumberg@vutbr.cz

Abstract—Utilizing Scratch and Colony formation assays is crucial for understanding cell migration and proliferation dynamics. In this study, we present novel approaches utilizing MATLAB for the processing of images obtained from these assays. We demonstrate the effectiveness of our approaches through validation experiments on human osteosarcoma SaOS-2 cells, showcasing their applicability across different experimental conditions. Moreover, the inhibitory effect of two concentrations of gentamicin (0.1 mg/ml, 0.3 mg/ml) on the migration and clonogenicity of SaOS-2 cells was declared.

Keywords—Gentamicin, SaOS-2, Scratch assay, Colony formation assay, MATLAB

I. INTRODUCTION

The Scratch assay, also known as the Wound healing assay or migration assay, is a commonly used *in vitro* technique in cell biology to study cell migration dynamics. It involves creating a "scratch" or artificial wound on a cell monolayer, typically using a sterile pipette tip or specialized tool, and then monitoring the movement of cells to close the gap over time. The Scratch assay is utilized for various purposes, including a study of cell migration and investigation of wound healing processes. In a Scratch assay, images of a cell-free area are typically captured using a microscope equipped with a camera at regular time intervals, starting from the initial time point, when the "scratch" is created until it is completely or partially closed.

There are various approaches for analyzing images from Scratch assays. Usually, the wound area is measured or wound closure percentage is calculated. Manual processing of the Scratch assay images involves visually inspecting the images to identify the scratch area, and then tracing the edges of the scratch by manually drawing lines along the edges of the scratch to outline its borders [1, 2]. After tracing the gap borders with lines, the manually processed images can be quantitatively analyzed (e.g. using ImageJ [3, 4]). While manual processing of Scratch assay images using lines provides precise control over the analysis, it can be time-consuming, labor-intensive, and subjective. To address this problem, semi-automated approaches can be utilized. Wound healing size tool [5] is an ImageJ plugin that allows the analysis of brightfield, phase contrast, and fluorescence images, but the major drawback of this approach is a necessity to adjust the input parameters to fit accurately the wound area. In our previous work [6], a MATLAB approach for fluorescence image processing was described. The main shortcoming of this method is a presence of fluorescent cell labeling that can alter cell viability and behavior and increase preparation time. There are also some commercial software for automated processing of Scratch assay images [7, 8], but a significant limitation is that they are not free.

A Colony formation assay is an easy and well-established laboratory technique commonly used to assess the ability of single cells to proliferate and form colonies under specific conditions. It is particularly useful for studying the clonogenic potential of cancer cells, which refers to their capacity to grow and produce colonies of daughter cells. In this assay, single cells are seeded and allowed to grow into colonies over time. After fixation and staining, colonies are counted to measure the cells' clonogenic potential and growth characteristics. The results of the Colony formation assay are typically expressed as the number of colonies per well. Moreover, some additional parameters can be calculated, such as mean colony size and the area covered by cell colonies. Manual approach is acceptable for counting the colonies, although it may be tedious. On the other hand, utilizing image analysis systems and algorithms can dramatically simplify and accelerate data processing. The most commonly used approaches for Colony formation assay image analysis are based on ImageJ software. Examples include ColonyArea plugin [9] and ColonyCountJ [10] add-on program, which measure both the number of colonies and the area covered, along with average colony intensity. OpenCFU [11] is an open-source software programmed in C++ that enables the enumeration of circular objects, including bacterial and cell colonies. Additionally, there are commercially available solutions [12, 13] that provide automated, high-throughput analysis, albeit at a higher cost. While numerous free solutions are available, each has its drawbacks. The ImageJ plugins usually require human intervention, such as converting to grayscale and thresholding. Here, the parameter selection process demands significant user experience. Furthermore, many programs lack the capability to exclude colonies that are too small (less than 50 cells or 1 mm size), potentially resulting in misinterpretation of surface dust as cell colonies.

Therefore, the aim of this study was to develop robust algorithms in MATLAB for analyzing Scratch assay and Colony

formation assay images to address these limitations. MATLAB is a powerful platform for data processing with significant potential for application in biological research. Despite being a commercial software, MATLAB is widely popular among academic institutions. Here, the effect of gentamicin on migration and clonogenicity of human osteosarcoma SaOS-2 cells was evaluated, and the Scratch assay and Colony formation assay images were acquired. The accuracy of image processing using created algorithms was compared with selected ImageJ plugins.

II. MATERIALS AND METHODS

A. Description of the experiment

Human osteosarcoma SaOS-2 cells were purchased from CLS (Cell Lines Service GmbH) and were maintained in high glucose Dulbecco's Modified Eagle Medium (Sigma-Aldrich) and Ham's F-12 medium with L-Glutamine (Serana) (DMEM:F12, 1:1) supplemented with 10% fetal bovine serum (Biosera) and 1% of penicillin-streptomycin (Sigma-Aldrich) at 37 °C and 5% CO₂. The medium was changed every 2-3 days, and cells were subcultured after reaching 90% confluence. Gentamicin used in this study was purchased from Sigma-Aldrich (10 mg/ml in deionized water).

For Scratch assay, 23.10⁴ cells were seeded into each of three tissue culture dishes (Deltalab) with a growth area of 8.5 cm^2 . Cells were maintained in culture medium for 3-4 days until a confluent cell monolayer was formed. Then a linear "scratch" was performed in each dish using a p200 pipette tip. After scratching, the culture medium was removed and cells were washed twice with phosphate buffer saline (Sigma-Aldrich). Subsequently, the dishes were filled with 1.5 ml of culture medium. First dish was used as control (no gentamicin present), while the cells in remaining two dishes were treated with different concentrations of gentamicin (0.1 mg/ml and 0.3 mg/ml, respectively). The microscopic images (1920×1440 pixels) of cell-free area were acquired using Leica DMi8 phase contrast microscope equipped with a CCD camera (LEICA DFC7000 T) immediately after scratching (0 hours), then in 24 and 48 hours using a 10X objective magnification. The image acquisition parameters were the same in all groups.

For Colony formation assay, SaOS-2 cells were seeded in 6-well plates (1,100 cells/well), allowed to adhere for 24 hours, and treated with mentioned concentrations of gentamicin from day 2. During each experiment, three wells were used, one as a control and the other for exposure to 0.1 mg/ml and 0.3 mg/ml of gentamicin. After 14 days of incubation, colonies were fixed with ice-cold 70% methanol for 20 min and stained with 0.5% crystal violet for 15 min. Then, the cells were washed three times with distilled water. Finally, the wells were photographed using a Samsung Galaxy S20 FE smartphone. The images have a spatial resolution of 3024×4032 pixels and were taken with an aperture of f/1.8, ISO 40 and a focal length of 26 mm.

B. Data analysis

The functions for the evaluation of Scratch and Colony formation assay (available at: <u>https://github.com/KaterinaIngrova/ScratchColonyQuantificati</u> <u>on</u>) are implemented in the MATLAB programming environment (version R2020a).



Fig. 1. Individual steps of the algorithm: (A) original image, (B) image edges, (C) brightness profile of the image, (D) possible dead cells, (E) image with removed dead cells, (F) image after morphological operations, (G) segmented scratch, (H) scratch boundaries on the original image

The Scratch assay algorithm analyzes the phase microscopy images (with vertical scratch) and calculates the average scratch width in micrometers and the scratch area in mm². The user has to provide two input variables to the function, namely the path to the microscopy images and the pixel size in micrometers. First, the function retrieves a list of images and initializes a field to store the results including the image name. Then, the individual images are scanned to load the current image and determine its size. Several operations are then performed, which could be categorized as image processing, brightness profile analysis, dead cell removal, scratch segmentation, scratch width and area calculation, and saving the results (see Fig. 1). During image processing, the contrast is enhanced, edges are detected using a Canny detector and circles with a given radius corresponding to dead cells are found. In the next stage, a brightness profile of the image is created to remove dead cells from the scratch region only. To remove the dead cells from the edge representation, they must also have an average brightness value greater than the set threshold. The following segmentation uses morphological operations including hole filling. In case the scratch area is already significantly overgrown, a local map of standard deviations and entropy is used for segmentation, also with further use of morphological operations and filling holes formed by dead cells. The scratch width is further converted to micrometers using the given pixel size. The result of the scratch area is expressed in mm². In the last step, the edge of the scratch is added to the original image to visualize the result. The obtained results are stored in the results field mentioned before. While the algorithm is running, the information about the image name and the results obtained are printed to the Command

window. The results of each step are also displayed in a Figure window as the algorithm progresses.

The algorithm for Colony formation assay analyzes the smartphone images and calculates the area of the well covered by colonies in cm² and the number of colonies. The user has to enter three input variables into the function, namely the path to the images, the well size in cm² and the valid colony area threshold in mm². First, a list of images is loaded and a field for storing the results including the image name is initialized. Subsequently, the current image is loaded and a new window is displayed in which the user interactively selects a circle representing the culture well (see Fig. 2). According to the selected circle, the image is cut and the size of the cropped region is determined. To segment the colonies, k-means clustering is applied to the cropped image. Since the colonies may be connected, they are separated using watershed algorithm. The colony counting is then performed using a subordinate function that analyzes the connected components, filters the cell colonies by minimum size, produces a binary image containing only valid colonies and displays information about the number of colonies. For a colony to be considered valid, it must be bigger than the valid colony area threshold. Here, a value 0.03 mm² was used, as it adequately covers a growth area of 50 cells. Then, the number of pixels to create a 1 mm colony is entered into the subordinate function as the minimum size. Finally, the valid (red) and non-valid colonies (cyan) are color distinguished and the area occupied by the colony in cm² is calculated. In the last step, the obtained results are stored in the results field mentioned above. During the algorithm execution, information about the image name and the obtained results are displayed to the user in the Command window. The results of each step are also displayed in a Figure window during the algorithm execution.



Fig. 2. Individual steps of the algorithm: (A) original image with labeled well, (B) well cutout, (C) segmented and separated colonies, (D) colonies counting, (E) colonies color separation

III. RESULTS AND DISCUSSION

Here, the effect of gentamicin on migration and clonogenic activity of SaOS-2 cells was evaluated. First presented algorithm was used to quantify the cell migration after treatment with gentamicin. The Fig. 3 and 4 show the process of wound closure within the groups. It can be seen that SaOS-2 cells exhibit slower migration in both tested groups compared to the control group. This conclusion was also confirmed using statistical analysis,



Fig. 3. Scratch assay images with illustrated scratch boundary: control (left), 0.1 mg/ml of gentamicin (middle), 0.3 mg/ml of gentamicin (right): (A) t = 0 hours, (B) t = 24 hours, (C) t = 48 hours. Scale bar: 300 μ m



Fig. 4. The effect of various concentrations of gentamicin on the migration of SaOS-2 cells

where the parametric one-way ANOVA was used to analyze the cell-free areas. In order to be able to use parametric one-way ANOVA, certain conditions have to be met, which were tested using a multi-sample test for equality of variances. The results indicated a significant difference between the control group and both tested groups at 24 hours ($p = 1.05 \cdot 10^{-4}$) and at 48 hours ($p = 2.110 \cdot 10^{-7}$).



Fig. 5. Color distinguished colonies (red for valid and cyan for nonvalid): control (left), 0.1 mg/ml of gentamicin (middle) and 0.3 mg/ml of gentamicin (right)

TABLE I. A QUANTITATIVE ANALYSIS OF COLONY FORMATION ASSAY IMAGES

Group	Control	Gentamicin 0.1 mg/ml	Gentamicin 0.3 mg/ml
Number of valid colonies	257	208	180
Colony area [cm ²]	1.8015	1.3045	1.0929

The effect of gentamicin on clonogenic activity of SaOS-2 cells was analyzed by utilizing the second algorithm. It was found that the cell clonogenicity was suppressed in all

gentamicin groups compared to the control group (see Fig. 5 and Tab. 1). Furthermore, the effect of gentamicin was dose-dependent, with the lowest number of colonies and total colony area observed at the highest gentamicin concentration.

The accuracy of the designed algorithms was compared with freely available ImageJ plugins. For Scratch assay images, Wound healing size tool [5] was used to process the images acquired in 0 hours and 24 hours after scratching. Additionally, the consistency of cell-free area detection was analyzed using images (2048×1536 pixels) acquired with a ZEISS Axio Observer microscope equipped with a CCD camera (PROMICAM 3-3CP). Fig. 6 shows that the scratch area detection accuracy is comparable in 0 hours images for both approaches, while it is significantly worse for 24 hours images analyzed with Wound healing size tool. Better segmentation can be achieved, however, by tedious selection of appropriate segmentation parameters. In contrast, our algorithm is ready-touse and only requires input of a pixel size in micrometers, while the other parameters including the threshold will be calculated automatically.



Fig. 6. Comparison of scratch area detection using our algorithm (green) and Wound healing size tool (red): (A1-2) 0 hours images, (B1-2) 24 hours images

The ColonyArea plugin [9] was used to examine the cell colony detection accuracy. An advantage of our algorithm is its ability to process standard images of individual wells, captured with any type of camera. These images may or may not include information about other wells on the plate. In contrast, ColonyArea is primarily created to analyze images of entire well plates acquired using flatbed scanners. To make the comparative analysis possible, such images were captured using a Samsung Galaxy S20 FE smartphone. Fig. 7 illustrates the results of cell colony segmentation using both approaches in well №2. Our algorithm detected 253 colonies (187 colonies were considered as valid) with 14.87 % area covered, while the ColonyArea declared an area coverage of 11.96 %. Unfortunately, this plugin lacks the capability to count the number of cell colonies. Unlike ColonyArea, our algorithm can not only count the total number of colonies but also identify and eliminate the smallest ones, similar to manual counting. Another advantage is that our approach does not require manual parameter tuning for image segmentation and can work automatically. However, it's crucial to remember that the quality of the original image is a significant



Fig. 7. Comparison of cell colony detection in well №2: (A) original image of entire well plate (B) our algorithm (C) ColonyArea plugin

factor in successful segmentation. For optimal results, it's recommended to use a suitable device and capture images against a white background.

CONCLUSION

In summary, this study has demonstrated the development and application of novel algorithms for the analysis of Scratch and Colony formation assays. By utilizing them on real data, we have shown that our algorithms can accurately quantify Scratch assay and Colony formation assay images with improved efficiency and reliability compared to existing methods. In addition, the evaluation of the impact of gentamicin on migration and clonogenicity of human osteosarcoma SaOS-2 cells was provided. The results declare that gentamicin inhibits cell migration and suppresses clonogenic activity of the cells in a dose-dependent manner.

REFERENCES

- H. Tan, M. Zhang, L. Xu, X. Zhang, and Y. Zhao, "Gypensapogenin H suppresses tumor growth and cell migration in triple-negative breast cancer by regulating PI3K/AKT/NF-κB/MMP-9 signaling pathway", Bioorganic Chemistry, vol. 126, pp. 1–9, September 2022.
- [2] C. -S. Cheng et al., "Paeonol Inhibits Pancreatic Cancer Cell Migration and Invasion Through the Inhibition of TGF-β1/Smad Signaling and Epithelial-Mesenchymal-Transition", Cancer Management and Research, vol. 12, pp. 641–651, January 2020.
- [3] N. Wang, T. Feng, X. Liu, and Q. Liu, "Curcumin inhibits migration and invasion of non-small cell lung cancer cells through up-regulation of miR-206 and suppression of PI3K/AKT/mTOR signaling pathway", Acta Pharmaceutica, vol. 70, no. 3, pp. 399–409, September 2020.
- [4] C. T. Rueden et al., "ImageJ2: ImageJ for the next generation of scientific image data", BMC Bioinformatics, vol. 18, no. 1, pp. 3–26, November 2017.
- [5] A. Suarez-Arnedo et al., "An image J plugin for the high throughput image analysis of in vitro scratch wound healing assays", PLOS ONE, vol. 15, no. 7, July 2020.
- [6] L. Baiazitova et al., "The Effect of Rhodamine-Derived Superparamagnetic Maghemite Nanoparticles on the Motility of Human Mesenchymal Stem Cells and Mouse Embryonic Fibroblast Cells", Molecules, vol. 24, no. 7, March 2019.
- [7] WIMASIS IMAGE ANALYSIS. WimScratch. Online. Available at: https://www.wimasis.com/wound-healing-assay.
- [8] IBIDI GMBH, 2019. Wound Healing and Migration Assays. Online. Available at: <u>https://ibidi.com/img/cms/resources/AG/FL_AG_033_Wound_Healing_150dpi.pdf</u>.
- [9] C. Guzman, M. Bagga, A. Kaur, J. Westermarck, D. Abankwa, and R. Rota, "ColonyArea: An ImageJ Plugin to Automatically Quantify Colony Formation in Clonogenic Assays", PLoS ONE, vol. 9, no. 3, pp. 1–9, March 2014.
- [10] D. Kumar Maurya, "ColonyCountJ: A User-Friendly Image J Add-on Program for Quantification of Different Colony Parameters in Clonogenic Assay: A User-Friendly Image J Add-on Program for Quantification of Different Colony Parameters in Clonogenic Assay", Journal of Clinical Toxicology, vol. 07, no. 04, pp. 1–4, July 2017.
- [11] Q. Geissmann and R. M. H. Merks, "OpenCFU, a New Free and Open-Source Software to Count Cell Colonies and Other Circular Objects", PLoS ONE, vol. 8, no. 2, pp. 1-10, February 2013.
- [12] CYTOSMART TECHNOLOGIES AND AXION BIOSYSTEMS, INC. Omni. Online. Available at: https://cytosmart.com/products/systems/omni.
- [13] SCINTICA INSTRUMENTATION, INC. GelCount. Online. Available at: <u>https://scintica.com/product/cell-and-isolated-tissue/automatedcolony-counter/.</u>

Probing natural molecules with PPAR-γ to reveal potent agonist against Cancer

1st Adriána Špaková Department of Biomedical Engineering Faculty of Electrical Engineering and Communication, Brno University of Technology Brno, Czech republic xspako00@vutbr.cz 2nd Vaishali Pankaj Department of Biomedical Engineering Faculty of Electrical Engineering and Communication, Brno University of Technology Brno, Czech republic 234084@vut.cz 3rd Inderjeet Bhogal Department of Biomedical Engineering Faculty of Electrical Engineering and Communication, Brno University of Technology Brno, Czech republic 234190@vut.cz 4th Sudeep Roy Department of Biomedical Engineering Faculty of Electrical Engineering and Communication, Brno University of Technology Brno, Czech republic roy@vut.cz

II. EASE OF USE

A. $PPAR-\gamma$

Abstract - The work focuses on searching for a molecule with potential agonistic properties against cancer. Molecules from six databases were screened and docked to peroxisome proliferatoractivated receptor gamma (PPAR- γ) by using computer aided drug design approach. Hits underwent further exploration, including dynamic simulation and safety verification. Piperlongumine - naturally occurring small molecule, derived from long pepper (Piper longum) showed after comparison with standards Troglitazone and Rosiglitazone promising results.

Keywords - Cancer, PPAR- γ , target, natural molecules, ligand, virtual screening, molecular docking, dynamic simulation, ADMET

I. Introduction

According to World health organization (WHO), cancer is one of a leading cause of death worldwide, accounting for nearly 10 million deaths in 2020, or nearly one in six deaths. [1] The numbers vary between countries. In the US, cancer is the second most common cause of death, exceeded only by heart disease. [2] Just like in the US, cancer is the second most common cause of death after cardiovascular diseases also in the Czech Republic, where around 100,000 people fall ill with cancer every year and 30 000 people die due to cancer. [3]

Cancer, as a complex and heterogeneous group of diseases, is characterized by uncontrolled cell growth and division, often leading to the formation of malignant tumours. Treatment options include surgery, chemotherapy, radiotherapy, personalized targeted therapy, and immunotherapy. [4]

Despite considerable advancements in cancer research and treatment modalities, the need for innovative approaches remains paramount. The exploration of natural molecules as potential agonists against cancer has emerged as a promising avenue, with particular emphasis on their interaction with the PPAR- γ . [2][4] The reason is that many drugs have their origin in natural sources, such as plants, animals, fungi, and microorganisms. E.g. aspirin which has analgesic and antipyretic properties came from tree bark, and morphine, known as pain reliever, has its origin in seeds of the opium poppy. Natural molecules serving as ligands bind to a receptor or target molecule and form a complex that triggers a biological response - this time the effect is required in cancer cells.

PPAR- γ is a ligand-activated transcription factor that regulates genes, which are important in cell differentiation and various metabolic processes, especially lipid and glucose homeostasis. After interaction with the specific ligand, nuclear receptor is translocated to the nucleus, where it changes its structure and regulates gene transcription. [5] The intricate relationship between PPAR- γ and cancer involves pathway modulation, especially cell proliferation, differentiation, apoptosis, angiogenesis, and inflammation. Although, its involvement in cancer is multifaceted, with both pro- and antitumorigenic effects reported in various contexts, the demand for PPAR- γ agonists as chemotherapeutic agents persists. [6][7][8]

B. Thiazolidinediones

Thiazolidinediones (TZDs) are synthesized PPAR- γ agonists, primarily utilized in the treatment of type 2 diabetes as insulin sensitizers. In comparison to other ligands, including mainly fatty acids, TZDs bind to PPAR- γ more firmly, with an affinity in the range of 40–200 nM. [2][9] Over a few TZDs generations, ciglitazone, troglitazone, rosiglitazone and pioglitazone were introduced. Some of them were withdrawn from the market due to risk of hepatotoxicity or cardiovascular incidents. Recently, TZDs proved also anti-cancer effects, dependently or independently of PPAR- γ activation, in monotherapy and in combined treatment with chemotherapeutics, both, on the transcriptional and the protein level. [2] The goal of the work is to find a natural molecule that will have better properties than standards from TZDs family.

C. Principle of gene expression

PPAR- γ belongs to a subset of nuclear receptors that form heterodimers with the retinoid X receptor (RXR), greatly enhancing the ability of the receptor to bind specific DNA sequences in target genes. The DNA sequences recognized by the PPAR–RXR heterodimer are referred to as PPAR-response elements (PPREs). PPAR–RXR heterodimers bind to PPREs in the absence of ligand. After that, ligand is bound, which leads to a PPAR- γ conformational change that results in activation of transcription of the target gene. [9] The Fig. 1 depicts mechanism of PPAR- γ in the function of transcription factor.



Fig. 1. Mechanism of PPAR- γ in role of transcription factor

III. COMPUTER AIDED DRUG DESIGN

Bringing a new drug into the market is a costly process in terms of money, manpower, and time. Conventional approach of drug discovery and development takes an average of 10-15 years with an approximate cost of 800 million to 1.8 billion US dollars. Over the last few decades, CADD has emerged as a powerful technique playing a crucial role in the development of new drug molecules. In biomedical area, CADD is being utilized to accelerate and aid hit identification, hit-to-lead selection, optimize the absorption, distribution, metabolism, excretion, and avoid safety issues. [10][11][12]

A. Target and ligands Selection and Preperation

The 2PRG protein structure was downloaded from the RCSB Protein Data Bank. The designation 2PRG stands for ligand-binding domain of human peroxisome proliferator activated receptor gamma. The downloaded receptor was incorporated into the Maestro environment, followed by protein preparation. The preparation consists of mainly H-bond optimization and water molecules minimization, what results in biologically relevant and energetically favourable protein structure state. Then receptor grid generation was required, what is a three-dimensional space generated around a target protein's binding site, where the ligand is expected to bind.

The next step included the selection and preparation of the ligand libraries to perform reliable virtual screening later. Molecules were downloaded from six different databases, namely: natural molecules from Interbioscreen, NP Atlas and Enamine, synthetic molecules from Lead-like and Drug-like database (Swiss similarity) and kinases from ChemDiv database. Standards Troglitazone and Rosiglitazone were downloaded from PubChem.

B. Structure-based Virtual Screening

Virtual screening narrows down the vast chemical space and prioritize compounds for experimental testing in the early stages of drug discovery. It helps to identify lead compounds that have a higher likelihood of binding to the target of interest. Virtual screening in that way saves time and expenses. Within the virtual screening, more than 7000 ligands and 2 standards were docked to the 2PRG receptor. Predicting ligand-protein binding affinities were determined by using extra precision docking mode with post-docking MMGBSA. Docked complexes were then ranked according to their scores and the top-ranked compound from each library was selected as potential drug candidate. From Interbioscreen natural compound database and NP Atlas database, instead of one top hit, two top hits were used in further validation. The Table 1 depicts glide gscore of hits and standard molecules, where lower gscore indicates a more favourable binding interaction. Most of hits provided better binding affinity score than standard ligands. The best glide gscore was obtained by molecule C21H16FNO3S (-14,791) from Drug-like database. Very similar results have been shown also by molecules C28H39NO6 (-14,577) and C24H20O10 (-14,38) from NP Atlas and by molecule C25H28N2O6 (-14,099) from Interbioscreen natural compound database. Important is also the representation of amino acids of the target protein, which form certain bonds with the given ligand. E.g. top-ranked ligand C21H16FNO3S from Drug-like library forms bonds with following amino acids: Gln 286, Tyr 473, Hie 323, Ser 289.

Compound	Library	Virtual Screening	
		Glide gscore	Amino acids
Rosiglitazone	PubChem	-11,644	Gln 286, Ser 289, Tyr 473, Hie 323
Troglitazone	PubChem	-12,377	Gln 286, Tyr 473, Ser 289, Hie 323
C16H13NO3S	Lead-like	-12,526	Gln 286, Tyr 473, Ser 289, Hie 323, Tyr 327
C21H16FNO3S	Drug-like	-14,791	Gln 286, Tyr 473, Hie 323, Ser 289
C15H16NO4	Enamine	-11,009	Hie 323, Hie 449, Gln 286
C25H28N2O6	IBS - top1	-14,099	Ser 342, Hie 323, Ser 289 , Hie 449
C17H19NO5	IBS - top2	-14,030	Hie 449, Gln 286, Tyr 473, Hie 323
C12H7Cl2F2N2O2	ChemDiv	-10,814	Hie 449, Gln 286, Ser 289, Hie 323
C28H39NO6	NPatlas - top1	-14,577	Leu 340, Hie 449, Hie 323, Ser 289
C24H20O10	NPatlas - top2	-14,380	Gly 284, Gln 286, Ser 289, Hie 323 , Hie 449

C. Induced Fit Docking

Traditional docking methods, which virtual screening is a part of, assume a rigid structure of receptor. Unlike that, induced fit docking (IFD) considers the flexibility of the protein as well as the ligand. It provides a more realistic representation of the ligand-protein interaction by considering the dynamic nature of both components. As shown in Fig. 2, the docking procedure is typically iterative - potential binding positions are generated.



Fig. 2. Rigid vs flexible molecular docking

For IFD were used hits from virtual screening. For each hit were generated several poses and the Table 2 depicts results of the best pose for each molecule. The rule is applied again, lower score suggests a more stable protein-ligand complex. Standard troglitazone provided the best docking (-14,818) and IFD (-1236,876) score. Second best docking score was provided by molecule C21H16FNO3S (-14,329) from Drug-like library until the second best IFD score has been shown by molecule C25H28N2O6 (-1234,578) from Interbioscreen database.

TABLE 2 - SCORES OF BEST POSES FROM INDUCED FIT DOCKING

Compound	Library	Induced Fit Docking (XP)			
		Docking score	IFD score	Amino acids	
Rosiglitazone	PubChem	-11,727	-1228,766	Ser 342, Arg 288, Ser 289, Tyr 473, Hie 323	
Troglitazone	PubChem	-14,818	-1236,876	Gin 286, Hie 323, Ser 289	
C16H13NO3S	Lead-like	-6,235	-1224,805	Ser 289, Hie 449	
C21H16FNO3S	Drug-like	-14,329	-1229,005	Gin 286, Ser 289, Hie 323	
C15H16NO4	Enamine	-8,907	-791,632	Lys 367, Arg 288, Ser 289, Hie 323, Cys 285	
C25H28N2O6	IBS - top1	-13,866	- 1234,578	Ser 342, Hie 449, Tyr 327	
C17H19NO5	IBS - top2	-10,675	-1228,579	Ser 342, Hie 323, Ser 289	
C12H7Cl2F2N2O2	ChemDiv	-8,873	-1224,914	Tyr 473, Hie 449, Ser 289	
C28H39NO6	NPatlas - top1	-13,135	-1234,540	Cys 285, Lys 367, Ser 342	
C24H20O10	NPatlas - top2	-13,507	-1226,110	Tyr 327, Tyr 473, Hie 323, Cys 285, Arg 280, Ser 342	

Again, the amino acids involved in the bonds between the protein and the ligands were described. The best pose of standard troglitazone forms bonds with 3 amino acids: Gln 286, Hie 323 and Ser 289. The same amino acids are visible in case of best pose of ligand C21H16FNO3S from Drug-like library and best pose of molecule C25H28N2O6 from Interbioscreen natural compound database forms bonds with totally different amino acids: Ser 342, Hie 449 and Tyr 327.

D. Dynamic Simulation

Molecular dynamic simulation (MDS) as a powerful computational technique was used to simulate the behaviour of molecular systems over time under specified conditions (300K,100nsec). MDSs were applied on the best poses from induced fit docking. As the result, information about molecular trajectories, conformations, interactions, and properties such as energy were obtained. Fig. 3 depicts set of plots catching Root Mean Square Deviation (RMSD) of receptor PPAR- γ (left Yaxis) and four ligands (right Y-axis). Monitoring the RMSD of the protein gives insights into its structural conformation and system equilibration. Desirable is when simulation converges the RMSD values stabilize around a fixed value, maximally changes of the order of 1-3 Å are acceptable. Ligand RMSD indicates how stable the ligand is with respect to the protein and its binding pocket. If the values are significantly larger than the RMSD of the protein, likely is that the ligand will diffuse away.

If we compare the plots representing the RMSD of standard rosiglitazone (1st) and troglitazone (2nd), the difference in the overlap of the curves is clearly visible. In case of troglitazone, which provided better binding affinity results already after induced fit docking, we can see larger overlapping, so complex is much more stable. Next two plots represent RMSD of molecules from Interbioscreen library. Ligand C25H28N2O6

(3rd) showed little lower overlapping compared to standard troglitazone (2nd), however, system seems to be equilibrated, the RMSD values stabilized around a fixed value in the end of simulation. The best overlapping status was provided by ligand C17H19NO5 (4th). Within all plots, the values of protein and ligand RMSD are very similar, what indicates that the ligands will probably not diffuse away from theirs initial binding sites.



Fig. 3. RMSD of receptor PPAR- γ (left Y-axis) and ligands (right Y-axis, from top to bottom: Rosiglitazone, Troglitazone, C25H28N2O6, C17H19NO5)

Another observed phenomenon within MDS is the presence of bonds formed between ligand and protein amino acids at the
time of simulation. Bonds are categorized into four types: Hydrogen Bonds, Hydrophobic, Ionic and Water Bridges. The stacked bar charts were normalized over the course of the trajectory. For example, a value of 0.7 suggests that 70% of the simulation time the specific interaction was maintained. Values over 1.0 are possible as some protein residue made multiple contacts of same subtype with the ligand.



Fig. 4. Protein-Ligand Contact: C17H19NO5

The Fig. 4 shows results specifically for ligand C17H19NO5. We can see the presence of all types of bonds. E.g. significant H-bonds that occurred more than 30% of the simulation time were formed with residue Ser 289, Ser 342 and Tyr 473 until precious ionic interaction was observed only in case of Lys 367.

The principal component analysis (PCA) and Dynamic Cross-Correlation Matrix (DCCM), as powerful statistical methods, were also used in essential dynamics analysis. PCA illustrates the dynamic behaviour of the receptor molecule upon binding of ligand compounds. In Fig. 5, the first component (PC1) represents the large-scale conformational changes in the biomolecular system as the protein domain movements, ligand binding and unbinding events. The second component (PC2) corresponds to secondary motions or correlated fluctuations within a particular region of the system and reveals localized conformational changes that are not captured by PC1. Each point on the PCA plot corresponds to the projection of the biomolecular conformation onto the space defined by the first two principal components. Points that cluster closely together represent similar conformations, indicating that the biomolecule remained relatively stable or underwent only subtle changes during that simulation period. As the legend indicates, the results for individual molecules are colourfully differentiated.



Fig. 5. Principal component analysis

DCCM diagrams are displayed as colour coded matrix of Pearson correlation coefficients. The highly correlated motions are shown with positive values (blue region) while anticorrelation motions are shown by negative values (red region). When a cell in the DCCM heatmap is represented with high intensity colours, typically a mixture of red and blue, it indicates a strong correlation between the motions of the corresponding pair of residues. Lighter shades of red and blue indicate weaker dynamic couplings or less prominent functional interactions between the residues. Fig. 6 depicts heatmaps for two ligands.



Fig. 6. DCCM - Troglitazone (left), C17H19NO5 (right).

E. ADMET

Undesirable pharmacokinetics and toxicity of candidate compounds are the main reasons for the failure of drug development, and it has been widely recognized that absorption, distribution, metabolism, excretion, and toxicity (ADMET) of chemicals should be evaluated as early as possible. In silico ADMET filters are derived from chemical or molecular descriptors and respect that drug must reach the site of action, exert its pharmacological effect, and be eliminated in reasonable timeframe. In total, 7 categories containing 88 properties were considered during the molecule's evaluation. Specifically, it was about physicochemical properties, medicinal chemistry absorption, distribution, metabolism, excretion, and toxicology.



Fig. 7. Pie charts - evaluated properties

Every single property was evaluated according recommended empirical decision rule which determines ranges of acceptance (green - excellent, yellow - medium, red - poor). Pie charts depicted in Fig. 7 were created to compare results for individual ligands. The green pie chart shows the relative percentage of properties that provided an excellent result (the high percentage is desirable) while the red one visualizes relative percentage of properties with poor result (the low percentage is desirable). The Table 3 shows an example of the evaluation of specific ligand properties. PPB stands for plasma protein binding - one of the major mechanisms of drug uptake and distribution is through PPB, thus the binding of a drug to proteins in plasma has a strong influence on its pharmacodynamic behaviour. DILI represent Drug-induced liver injury. Next one, important for the work, is nuclear receptor PPAR-y. Others mentioned are Genotoxic Carcinogenicity Mutagenicity and Carcinogenicity itself, because of their serious effects on human health. Last one included is Lipinski rule.

Range Rosiglitazon PubCherr 92,04% 0,986 0,757 0,673 Accepte Troglitazon 100,31% 0,972 0,9 0 0,485 Accepted C16H13NO35 96,68% 0,977 0.366 0 0,888 Lead-like Accepte C21H16ENO35 99,68% 0,983 0.902 Drug-like 0.946 Accepte C15H16NO4 37,68% 0,973 0,341 0 0,196 Accepted 98,34% 0,311 0,52 C25H28N2O6 IBS - top1 0,663 Accepted C17H19NO5 68,68% 0,82 0,013 0,17 IBS - top2 Accepted C12H7Cl2F2N2O2 96,49% 0,988 0,003 0 0,08 Accepte 97.80% 0,039 0 0,314 C28H39NO6 NPatlas - top1 Accepte C24H20O10 99,94% 0,969 0,702 0,024 NPatlas-top2

TABLE 3 - ADMET EVALUATION

IV. CONCLUSION

After considering all the results so far, from virtual screening, through dynamic simulation to ADMET evaluation, several molecules showed better results than the standard Troglitazone and Rosiglitazone in certain aspects. Predominantly, outstanding results were observed for the naturally occurring small molecule C17H19NO5 – Piperlongumine, derived from the plant Piper longum, found in southern India and southeast Asia. This molecule seems worthy of further research.

ACKNOWLEDGMENT

I would like to thank my supervisor Sudeep Roy, Ph.D. for his invaluable advice and support. My big thanks also go to Pankaj Vaishali, M. Tech. and Bhogal Inderjeet, Mgr. for providing technical support and data for work implementation.

REFERENCES

- Cancer, 2022. Online. In: World Health Organiztaion. Dostupné z: https://www.who.int/NEWS-ROOM/FACT-SHEETS/DETAIL/CANCER.
- [2] CHI, Tiange; WANG, Mina; WANG, Xu; YANG, Ke; XIE, Feiyu et al., 2021. PPAR-γ Modulators as Current and Potential Cancer Treatments. Online. Frontiers in Oncology. S. 11. ISSN 2234-943X. Dostupné z: https://doi.org/10.3389/fonc.2021.737776.
- [3] CVEK, Jakub a HALÁMKA, Magdalena, 2023. Onkologie pro neonkology. Grada. ISBN 978-80-271-3090-0.
- [4] PINTO, Gaspar P.; HENDRIKSE, Natalie M.; STOURAC, Jan; DAMBORSKY, Jiri a BEDNAR, David, 2022. Virtual screening of potential anticancer drugs based on microbial products. Online. Seminars in Cancer Biology. Č. 86. ISSN 1044-579X. Dostupné z: https://doi.org/10.1016/j.semcancer.2021.07.012.
- [5] GRYGIEL-GÓRNIAK, Bogna, 2014. Peroxisome proliferator-activated receptors and their ligands: nutritional and clinical implications. Online. Nutrition Journal. Č. 13. Dostupné z: https://doi.org/10.1186/1475-2891-13-17.
- [6] HOUSEKNECHT, Karen L; COLE, Bridget M a STEELE, Pamela J, 2002. Peroxisome proliferator-activated receptor gamma (PPARγ) and its ligands: A review. Online. Domestic Animal Endocrinology. S. 1-23. ISSN 0739-7240. Dostupné z: https://doi.org/10.1016/S0739-7240(01)00117-5.
- [7] TACHIBANA, Keisuke; YAMASAKI, Daisuke; ISHIMOTO, Kenji a DOI, Takefumi, 2008. The Role of PPARs in Cancer. Online. National Center for Biotechnology Information. Dostupné z: https://doi.org/10.1155/2008/102737.
- [8] HARTLEY, Andrew a AHMAD, Imran, 2022. The role of PPARy in prostate cancer development and progression. Online. British Journal of Cancer volume. Dostupné z: https://doi.org/10.1038/s41416-022-02096-8.
- [9] REGINATO, Mauricio J. a LAZAR, Mitchell A., 1999. Mechanisms by which Thiazolidinediones Enhance Insulin Action. Online. Trends in Endocrinology & Metabolism. S. 10. Dostupné z: https://doi.org/10.1016/S1043-2760(98)00110- 6.
- [10] MACALINO, Stephani Joy Y.; GOSU, Vijayakumar; HONG, Sunhye a CHOI, Sun, 2015. Role of computer-aided drug design in modern drug discovery. Online. Archives of Pharmacal Research. Dostupné z: https://doi.org/10.1007/s12272-015-0640-5.
- [11] BAIG, M. H.; AHMAD, K.; ROY, S.; ASHRAF, J.M.; ADIL, M. et al., 2016. Computer Aided Drug Design: Success and Limitations. Online. Current Pharmaceutical Design. ISSN 1873-4286. Dostupné z: https://doi.org/10.2174/1381612822666151125000550.
- [12] BAIG, M. H.; AHMAD, K.; ROY, S.; ASHRAF, J.M.; ADIL, M. et al., 2016. Computer Aided Drug Design: Success and Limitations. Online. Current Pharmaceutical Design. ISSN 1873-4286. Dostupné z: https://doi.org/10.2174/1381612822666151125000550.

Modernization of Tenryu Pick and Place Machine

Bc. Jan Drška Faculty of Electrical Engineering and Communication Brno University of Technology Brno, The Czech Republic xdrska04@vutbr.cz

Abstract-This article discusses the possibility of extending the life of an older Tenryu automatic picking machine using modern technology with the possibility of connecting the machine to higher-level production planning and inventory control systems. The article describes the state of the machine before refurbishment and the design of a new system for its control, including the selection of control software.

Keywords - Automation, LinuxCNC, Mesa 5i24, OpenPnP, Pick and Place

I. INTRODUCTION

This article deals with the modernization of the Tenryu MT5530LQ setup machine, which is located at the headquarters of NNM Electric s.r.o. This machine is equipped with a now obsolete control system without the possibility of compatibility with superior systems and its maintenance is almost impossible due to the unavailability of spare parts. This setting machine is a full-size four-head placing machine with conveyor system in its original design. The conveyor allows very easy integration into the production line. The machine also includes an external computer with special software to create programs for the assembly process. Using this software, the machine can be controlled and the movement of all the assembly heads can be also controlled, including speeds, nozzle types and the order of the components to be assembled. However, the actual preparation of the setup program is very time consuming, and even the need to use an old MS-DOS computer is problematic nowadays. The original system has a big problem with fitting more complex or transparent components. The correct rotation of the component is sensed using special laser micrometers, which are now almost unobtainable and unreliable due to their age. These parts also need to be replaced by other technology and allow the recognition of all kinds of standard components.

The main task is to evaluate the current state of the machine, design and subsequently implement the modernization of the key parts of the machine. After such a modernization, the machine should be working again, and it should be possible to connect it with superior systems that will help to facilitate and streamline production in the company. These systems include, for example, database systems for parts inventory management, which allow the machine to be used efficiently and to have sufficient information on the status of the numbers of each type of parts. And to warn operators in time in case of stock shortages. The main advantage of upgrading should be compatibility with the spare parts available today and simpler operation including data preparation for the fitting cycle [3].



Fig. 1. Pick and Place Machine Tenryu MT5530LQ

II. PRINCIPLE OF THE PLACING PROCESS

The principle of assembly is that components are removed from the packages and then placed in a precise position on a printed circuit board (PCB) using a three-axis manipulator. The assembly process itself can be divided into four phases.

The first phase is data preparation. By data preparation we mean the actual design of the PCB, from which data is then generated for the production of these boards. The next output from the design software is a file that carries information about the positions and number of components. This is basically a table where basic information about the location, rotation and type of component is given for each component. This data is then imported into the fitting software that is part of the machine.

The second phase is the preparation of the machine, importing data and setting up the whole machine for mounting. First, a file with data about components and positions is imported. After the import, the feeders with the individual parts must be inserted into the machine and these feeders must also be set up. There are many different types of feeders, with the most basic ones being plastic or paper tape feeders, matrix feeders or vibratory feeders for parts in plastic tubes. After the feeders are set up, there is still the preparation for clamping the conveyor and the PCB clamp, where the parameters of the so-called crosssections are set. These crosspieces are the PCBs assembled into the matrix, where the program needs to specify where and how many PCBs to expect and which to place with components.

Once the machine is ready, the setup cycle can begin. At the beginning of this cycle, the machine will check the nozzle rack and calibrate the nozzles. Next, a camera is used to check the aiming points on the PCB, giving the possibility to adjust the machine coordinates and angle according to the physical state. After the check, the coordinates are recalculated and thanks to this we are able to fit the components very accurately regardless of the inaccuracies caused by clamping or milling the PCB during its manufacture. Once the coordinates have been successfully corrected, the machine starts to place the parts. During this process, the assembly head slides over the feeder and picks up the part using a vacuum nozzle. This component must first be oriented by the machine, where orientation is most often done using a camera. The part is moved over the camera and the dimension and rotation measurements are taken. The part is aligned and then the nozzle places it in the desired position in the pre-applied solder paste, where the part sticks.

The last stage is the soldering process. Soldering can be carried out in several ways, often using hot air furnaces or infrared furnaces. It is also possible to solder using wave or selective wave. After soldering, the assembly process is complete, and the PCB can start to come to life [3].

III. SETTING MACHINE MT-5530LQ

The Tenryu Automatic Placing Machine is a machine that is designed to be built into a production line. The main parts are a three-axis manipulator, a PCB conveyor, a quadruple setup head and 120 positions for feeders.

This machine has its own control computer that collects data from the process and controls all the actuators. The user interface includes a keyboard and two monitors. One monitor displays the direct output from one of the two cameras and the other monitor serves as the graphical interface to the automation system. Using this interface, the user can control the entire system. The system also includes a communication interface with an external computer running a program under MS DOS, which makes it easier to prepare the fitting data.

The machine has four set-up heads to handle up to 120 part feeders. Each head consists of two motors that take care of the lift and rotation of the parts. The machine can carry 4 nozzles directly in the heads and another 12 in the automatic nozzle changer. The nozzles have different shapes and diameters because the SMDs are different shapes and sizes, and a different nozzle is needed for each housing. There is also a laser micrometer on each head to measure the components. These micrometers are not used as much today and are being replaced by camera systems. There is a big problem with these micrometers on this machine. The micrometers are already very old and due to the aging of the materials their characteristics change and the control system of the machine does not allow to adjust or calibrate these sensors to their new values. Replacement micrometers are almost impossible to find and even then, they are very old pieces that may not work properly.

Another major disadvantage is the extreme time required to prepare the data and the entire machine. All the preparation has to be done manually, all the components and their positions with rotation angles have to be entered manually. For large batches this time is relatively small, but if the production needs to be changed frequently the impact is substantial and any small change leads to higher costs.

The machine is designed in such a way that if an error occurs during the setting cycle, the machine releases the PCB, and it is no longer possible to continue where the cycle left off. In the case of large batches this is negligible, but in small batches it is a significant loss also in terms of environmental friendliness. Massive throwing away of easily repairable parts is very common in large manufacturers, and yet often it would be enough to just fix a few things where the material and energy spent on production could be recovered instead of being thrown away.

The last disadvantage is the lack of network connectivity, including the possibility of integrating the machine into a business management system or inventory management system. The machine consumes on average thousands of parts per setup cycle and so it is necessary to know the stock status [3].

IV. MODERNIZATION PROPOSAL

The assembly of the entire machine requires a system that can handle the processing of camera signals, component measurements, reading and writing to digital and analog inputs. In addition, the system should be able to communicate with the parent database system and provide remote user access.

Due to the complexity of reassembling the entire system, it was decided to use ready-made software components to speed up the modernization. The machine needs to control the drives, actuators and perform the fitting algorithm. For the new system a new control computer has to be provided on which the control software will run. A standard DELL Workstation personal computer was selected for these needs, running the Linux operating system with Debian distribution. This configuration was chosen because of the intended use of the LinuxCNC software.

LinuxCNC is open-source software that allows you to control a large number of motor types from stepper motors to servomotors. This software is often deployed on computer-controlled machine tools. A special PCI slot card is required to control the drives and actuators. The Mesa 5i24 card was chosen for its parameters and compatibility with LinuxCNC. This card has an FPGA on it that can be configured for various applications from stepper generator for stepper motors or as an RS422 communication interface. This card will be used to control servo drives and collect data from the end position sensors of each axis of the machine [1][4].

However, to manage the assembly cycle, software is required to prepare the data for assembly and also to manage the component database and the ability to monitor and effectively manage the assembly cycle. After a survey of available solutions, the open-source software OpenPnP was selected. This software is designed to control automatic nozzle assembly machines for 1 to 4 heads with the possibility of future expansion. Furthermore, this system allows importing data from PCB design programs without the need to rewrite coordinates, etc. Furthermore, this software is able to work better with the setup cycle, where it can be freely paused and restarted from where the user needs it, leading to faster and better production with minimized losses. The software also supports many other things such as interfacing with database systems using scripts and other methods. The software is designed to run on both Linux and Windows operating system and the interfacing with LinuxCNC is done using network communication inside the computer through Telnet interface [1][2].

This software will also help to get rid of the dependence on the lack of spare parts for laser micrometers. These micrometers are replaced here by a camera that scans the parts optically and the image obtained is then evaluated and the part is placed in position with the corrections recalculated. Most of today's automatic machines already use this technology and laser measurement is being abandoned. The field of image processing is very attractive today and thanks to the use of OpenCV we are able to recognize shapes, measure their dimensions, detect anomalies and with these advances we are taking the accuracy and reliability of machines to a whole new level [3].

V. MODERNISATION OF SELECTED PARTS OF THE MACHINE

After examining the available solutions and the needs of the machine, the following design was decided. First of all, it is necessary to remove the old and discarded parts of the machine. Those parts that can be reused will be retained or modified to ensure their compatibility with the new system. Those parts that require communication to operate are taken out of service and replaced with new parts due to the lack of documentation of communication protocols.

It is possible to keep the whole machine chassis with the basic power circuits, which are the main supply, the switchboard with circuit breakers, the 230 V / 110 V transformer and the power supplies for the 5 and 12 V control voltages. In addition, the motors including their drivers will be retained. Two CRT monitors will be removed from the original machine and replaced by two new monitors. The video computer and the micrometer signal processing unit will be removed along with the main computer. Removing the old and unnecessary parts will leave plenty of space that can be used to incorporate the new systems. The original machine includes an extensive pneumatic system that will be retained in its entirety and will be expanded to include additional valves to control the pneumatic component feeders.

A PCB for controlling the drivers for the motors, a new PCB for controlling the mounting heads and their movement, and also a PCB for communicating with the control computer need to be made again. All these cards will be connected by an RS485 industrial communication interface, where the information needed to control the entire machine system will be exchanged using the Modbus protocol. Figure 2 shows the hierarchy and symbolic connection of the control cards to each other, including the different control levels. The interconnection and structure of the new control system can be clearly seen. The main system is the OpenPnP program, which controls all machine operations. The movements are handled by LinuxCNC, which is controlled by the OpenPnP program. Data acquisition from sensors and

control of actuators is handled by OpenPnP itself, which sends write or read commands over USB. OpenPnP sends the request, the new cards process it and return the requested value or confirm the execution of the command. The other cards will periodically communicate with each other and pass status information to each other, thus providing all the necessary information from the process needed to control the further operation of the machine.



Fig. 2. Diagram of the whole system

So far, 3 different surface connections have been designed for the new system. The first is the control electronics to control the drivers for the actuators that move all 10 axes. These boards are designed to be able to handle the signal from the Mesa 5i24 card, which is controlled by the LinuxCNC software. This software sends motion commands in the form of step and DIR signals. On this board, the conversion to CW and CCW occurs in differential form. In addition, this board contains a processor with RS485 communication interface, several current-amplified transistor outputs, and the ability to connect the drive encoder directly to the control software. This designed board is capable of controlling 2 motors and therefore there will be a total of 5 motors in the machine. Each mounting head has two drives and so will have one board, which means 4 boards for the heads and the last board will be for the X and Y axis drives.

The next control board will be the electronics to control the actuators and sensors for the mounting heads. Actuator means a set of pneumatic valves to control the suction or blowing of air through the nozzle. Sensorics means collecting data from the end positions of the mounting heads as well as sensing vacuum to check that the component is properly picked onto the nozzle. This board also has an RS485 communication interface to communicate with another boards. There is also a camera on the mounting head which is used for aiming the aiming points and for calibrating the positions of the components in the feeders.

This camera needs good lighting to function properly. The original lighting of the machine will be retained, but its control will be replaced by switching drivers. This board will also include two channels for controlling the camera backlight.

The last board produced is the communication interface between the OpenPnP control program and the machine, and also between LinuxCNC and the end sensor states. This board will contain several communication interfaces and a direct connection to the Mesa card. The card will have two independent RS485 lines for communication with the motor control and for communication with the mounting head board. There will also be a conversion to a USB interface through which the OpenPnP program communicates with the machine hardware. The machine also includes a static bottom camera, which is used to aim the picked parts. There will be two channels on this board for this camera to control the camera lighting in the same way as there is for the head control board.

There are a large number of other actuators and sensors in the machine that need to be serviced. A programmable logic controller (PLC), developed by the company that owns the machine, will be used for these purposes. This PLC will take care of the operation of the conveyor system, the operation of the automatic nozzle changer, the operation of the new feeders for the parts and for the data acquisition from other sensors. Such as the pressure sensor in the pneumatic system or the detection of the presence of PCBs on the conveyor. This system will also include ensuring the functional safety of the machine, including safeguards to protect against accidents. The machine contains fast moving parts and therefore operator protection needs to be ensured. There are covers on all sides of the machine with end position sensors as well as emergency stop buttons. All these features will be restored and will be part of the new machine system.

The final part of the upgrade is to increase the number of positions for the parts feeder. The control of the feeders is purely mechanical. The mounting head has a small hammer on it that, when driven into position, strikes the mechanical feeder slider where it moves the part to the take-up position and the nozzle picks up the part. This hammer is located only on the front of the head and is not located on the rear. For this reason, it is not possible to put these mechanical feeders on the back side because they cannot be operated by the hammer. With the new system and the added pneumatic elements, it will now be possible to put the frequently used pneumatic feeders for larger parts also on the rear positions, which were previously only used for vibratory feeders, which are not used that much. The ability to use the rear positions of the machine will also make it possible to fit boards with more types of components. This will make the machine more efficient and there will be no need to switch feeders between the machine and the stand so often, which will lead to simpler machine operation and more efficient production [3].

VI. CONCLUSION

Thanks to the deployment of the previously mentioned system and functions, the machine can continue to operate. In its original configuration, the machine was able to assemble approximately 3,000 components in one hour of operation. The new system is a little weaker and the first tests have confirmed the possibility of seeding something between 2500-2800 components per hour of machine operation. However, the new system is far superior to the old one mainly because of its flexibility and compatibility of parts and software solutions with today's technologies and trends. Selected LinuxCNC and OpenPnP systems are still under development and thus some support for the future is guaranteed. Another huge advantage is the openness of the systems and the possibility of easy modifications in the software and hardware solutions.

Due to the large number of parts to be assembled, it is necessary to introduce an automatic reading of the removed parts from the warehouse. For this purpose, company uses the Inventree program, which offers the possibility of subtracting parts from the warehouse thanks to a common API. Thanks to the new system, the company will already be able to plan production better and the data in the system will automatically change according to the current status and there will no longer be time delays due to lack of information or low stock.

Another great advantage of the modernization is the versatility of the solution used, where the system can be deployed on other machines owned by the company. The company owns two more machines, one the same Tenryu and then another Heeb. All these machines are already out of service due to obsolescence and unavailability of spare parts for the original systems. Thanks to this solution, a total of 3 machines can be made operational and can continue to perform their task.

ACKNOWLEDGMENT

This is an internal project of NNM Electric s.r.o., where the machine is located. This project is conducted as a thesis. This solution uses a number of open-source solutions that are freely available on Github or their own websites. The project is further created with the help of the company's managing director Mr. Jan Navrátil and the thesis supervisor doc. Jan Mikulka. I would like to thank all of the above mentioned for their guidance and help in the project of upgrading the embedding machine.

REFERENCES

- [1] LinuxCNC, c2023. Online. www.linuxcnc.org. [cited 2024-3-8].
- [2] OpenPnP, c2023. Online. openpnp.org/. [cited 2024-2-8].
- [3] DRŠKA, Jan. Modernization of Tenryu Pick & Place machine [online]. Brno, 2024 [cit. 2024-03-10].
- [4] MESA 5124, c2023. 5124 ANYTHING I/O MANUAL. Online. MESA. www.mesanet.com/pdf/parallel/5i24man.pdf. [cited 2024-3-8].

Naše práce je věda

V našem brněnském technologickém centru vznikají špičkové elektronové mikroskopy a spektrometry, které dodáváme do celého světa. Studují se jimi viry, vznikají díky nim vakcíny, vyvíjí se lepší materiály i elektronika. Pracujeme se špičkovými technologiemi, které posouvají lidské poznání. Najdeš mezi námi odborníky na fyziku, elektroniku, software, mechanickou konstrukci nebo logistiku. Chceš toho být součástí?

thermofisher.jobs.cz

Thermo Fisher



Automated testing device for turbojet ECU

Matúš Halgoš Department of Control and Instrumentation Brno University of Technology Brno, Czech Republic 211420@vut.cz

Abstract—This article presents the research and development behind the automated testing device tailored for electronic control unit (ECU) E040, which controls turbojet engine TJ40 to enhance efficiency and reliability in the testing phase of ECU. The proposed device offers significant advantages over traditional manual testing methods, including reduced testing time and improved accuracy. This paper will give the reader necessary information about TJ40 engine and E040 ECU, then will provide the overall concept design behind automated testing device.

Index Terms—Testing device, ECU, turbojet

I. INTRODUCTION

During the production of the device, a defect may occur in one or more production processes, which may affect the resulting function of the device. The device that contains such a defect must be detected and the cause of the defect must be identified. The systematic process by which these defects can be detected is called testing and consists of a set of actions designed to match the real function of the device as closely as possible. Testing is one of the key aspects of production, especially when it comes to serial production, in which it is necessary to observe a high level of standard. In serial production, it is important that product testing takes as little time as possible, and at the same time all functional elements are reliably tested. One option for performing such testing is testing using automated testing equipment that performs tests independently with minimal operator intervention. Manual testing of the E040 ECU can take up to 6 hours and can be easily susceptible to human error. Automated process of testing same ECU can be performed in 25 minutes. The main advantages of automated testing are mainly speed and reliability. [1]

II. TURBOJET ENGINE TJ40

The TJ40 is a small turbo-jet engine developed by PBS Velká Bíteš, a. s.. The engine, with its parameters, weight of 3.8 kg and dimensions of 147 x 304 mm, belongs to the category of small jet engines, making it an ideal solution for various applications eg. unmanned aerial vehicles (UAVs), training targets or it is also popular between RC modelers. Depending on the version, it can generate a thrust of 395 N or 425 N, thanks to which it can develop a speed of up to 0.8 M. The motor is equipped with a BLDC starter-generator, which guarantees the start of the motor from any position and after start, starte-generator generates electricity for on-board systems [2]



Fig. 1. Turbojet engine TJ40 [2]

III. ECU E040

Electronic control unit E040 is complex electronic device developed by division of aerospace and advanced control of the company Unis a.s. Its main purpose is to control TJ40 using data gathered from sensors and pilot. E040 consists of two modules called BER and BEM. This paper will be focused mainly on BEM. The functions of BEM are: monitoring temperature of intake air with sensor PT100, monitoring temperature of exhaust gasses with K type thermocouple, reading and processing turbine revolutions, evaluation of overspeed protection, controlling BLDC starter generator, maintaining desired on-board voltage level and controls 3 fuel valves. BEM is built of two printed circuit boards which one of them is control (BEMC) and the other is power (BEMP). [3]

IV. CONCEPT DESIGN OF AUTOMATED TESTING DEVICE

Designing device which will be able to detect errors during manufacturing or handling is composed of multiple interrelated processes. A difference is made between faults that cause an immediate malfunction of the device, for example a short circuit of the power supply to ground or a malfunction of some key components such as the power supply for example. The second type of faults are component faults that cause the malfunction of individual functional blocks of the device and the device is consequently unable to perform its function properly. The first type of faults is relatively easy to detect by measuring the continuity of the circuit or by measuring the voltage, on the other hand, for faults that affect the function



Fig. 2. 3D model of PCB BEMC(left) and BEMP (right) [3]

of the functional blocks and may occur during the operation of the device, it is necessary to make a set of tests that cover the overall function of the block. At first it is necessary to understand how the function blocks of the unit works and find non destructive yet sufficient way spot mistake.

V. TEST CASES

Hardware test cases form a document that specifies, in bullet points, the exact steps to create a test procedure. It is created for a specific device, describing the set of parameters and cases that are needed to create an automated test device. The test cases are intended to cover the tests associated with PCB bring-up as well as functional tests that test the correct function of the device. The tests are based on the technical specification of the device and designed to cover the function under test in its entirety. One test case is created for each function under test. The test cases are in the form of a table to make their description clear and easy to read. [1]

VI. TEST PROCEDURES

In order to properly test the equipment, each test must be approached consistently. The test procedures are based on the test cases and the individual cases are described in such a way that it is possible to clearly follow the description. The main purpose for this is, that the person creating software for test device doesn't need to know the hardware part of the test device or the device under test. The procedures describe the entire test procedure, define the instruments used, and specify the connection sequence and setup of each instrument. The test procedures describe in detail the testing sequence as well as the prescribed input/output peripherals to be used for the execution of the test. The exact device registers that are required for value verification or possible calibration are also listed. [1]

VII. DESIGN OF TEST CASES PROCEDURES

When designing test cases and procedures, it is essential to gain a deeper understanding of the device's function and then identify the most appropriate testing method for each functional block. It is important that the test method is both simple enough to implement and comprehensive enough to be able to cover all possible situations that may occur in a real environment. For example, when verifying the functionality of a thermocouple, it is not necessary to use a real thermocouple that needs to be complexly heated to a precise temperature. It is sufficient to use an isolated voltage source that can simulate both sides of the measured range and thus provide relevant input data for testing. Similarly, when testing a speed processing circuit, it is not necessary to use a real motor with a speed sensor. It is simpler to use a generator that can simulate the output from a speed sensor. These approaches allow efficient and simple testing without unnecessary complexity and cost, which increases the efficiency of the testing process and reduces the risk of errors.

VIII. AUTOMATED TESTING DEVICE TJP1

The device consists of two printed circuit boards. The first is named BTES and second is named BINT which are stored in a manual test fixture to make it easier for the operator of the device to perform testing.



Fig. 3. 3D Model of testing device TJP1 (front view)

In figure 3 is shown placement of PCBs in the testing device. This interchangeable cartridge will be placed in testing fixture. In order to test thoroughly, it is necessary to divide the testing into two parts - production and functional tests. During testing sequence, the PCB BEMC and BEMP production tests are performed first to ensure that both printed circuit boards were correctly manufactured. All tests are performed by PCB BTES. The manufacturing tests include:

- Open short circuit voltage circuits: This test will be performed before the tested device is connected to power, to ensure that both testing and tested device will not be damaged due short circuit
- Open short circuit of power MOSFETs: This test will ensure that there will be no short circuit to ground due to failed MOSFET.
- Inverter temperature sensor
- Intermediate circuit capacitance

- CAN bus test: This test procedure will provide information about correct impedance of can bus protection. If the impedance is too high the communuication may not work properly.
- Inverter temperature sensor
- Current sensor test: Is provided to find out if ECU inerter current readings are correct
- · Voltage source test

After the production tests have been performed, the individual signals on the BEMP and BEMC PCBs are interconnected using the BINT PCB and the software can be loaded into the control unit and all the functional tests of the control unit can be tested. Functional tests include:

- Voltage measurement and comparison with the values measured by the unit
- Measurement and calibration of intake air and exhaust gas temperatures
- Measurement and calibration of phase voltages, speed measurement test and check of over-speed protection function
- Switching of fuel valves
- Checking the function of the inverter
- Engine start sequence

After the testing is accomplished the two PCBs are solder into each other. This type of testing ensures that when a small error occurs it will be possible to replace broken component. Otherwise after soldering PCBs even this small error would be unable to fix



Fig. 4. Architecture of testing device TJP1

In the figure 4 is shown architecture of the device. The PCB BINT consists of relays and connectors which are connecting both signals of PCBs together to perform function tests. All the signals from both tested PCBs are connected trough interface block to the PCB BTES. BTES provides all the voltage levels and voltage references to provide exact measurements. The two measurement cards based on NI LabView platform control sequence of tests. In test sequence all the tests are performed and measured data are stored into test protocol.



Fig. 5. 3D model of testing PCB BTES

In the figure 5 is shown final 3D model of layout PCB BTES. Dimensions of the board are 290x190 and it is a 4 layer printed circuit board. The two inside layers consists of copper planes to which they are assigned power and ground signals and also are used for an insulator layers between signals. In the top and bottom layers are routed all other signals.

ACKNOWLEDGMENT

I would like to thank Ing. Tomáš Benešl for his guidance and helpful thoughts during the writing of this paper and also my thank belongs to Unis a.s. company for providing access to its infrastructure and resources during this research.

References

- M. Halgoš, Tester pro řídicí jednotku [online]. Brno, 2024 [cit. 2024-03-02]. Available from: https://www.vut.cz/studenti/zavprace/detail/155531. Supervisor Tomáš Benešl.
- [2] PBS Velká Bíteš, a. s. Turbojet engine PBS TJ40-G1. Online. Available from: https://www.pbs.cz/cz/Letectvi/Letecke-motory/ Proudovy-motor-PBS-TJ40-G1. [cit. 2024-03-03].
- [3] Unis, a. s. Technical specification of electronic control unit E040 Internal document of the company Unis, a. s.

Semiautomatic crimping machine

Jakub Nosek Brno University of Technology FEEC Department of control and instrumentation Brno, Czech Republic jakub_nosek@outlook.com

Abstract—This document deals with explanation of the design, creation and testing of a new automated machine for the production of ureteral catheters for a company producing medical equipment. The work is a continuation of the diploma thesis of a student of the College of Polytechnics Jihlava, whose work only concerned the design of the mechanical part of the machine and the analysis of the actual manual production process.

Keywords—PLC, TIA Portal, HMI, machine safety, EPLAN

I. INTRODUCTION

This work deals with the creation of a machine for the automation of the production process of ureteral catheters. The machine is divided into two parts. A manual process shortens ureteral catheters to the desired length. An automated process assembles the support wire designed to transport the catheter. This wire is also called a stylet. Part of the assignment for the this work is the creation of an electrical project with the selection of electrical components, the creation of a risk assessment according to applicable legislation, the design of software equipment together with visualization, the practical connection of visualization and control software and, finally, the verification of the entire solution in operation [1].

A. Urological catheter

A ureteral catheter is used in patients when it is not possible to drain urine naturally from the body. If we talk about draining urine from the body in a way other than natural, then we talk about urine derivation. The catheter is most often made from a polymer, specifically from polyvinyl chloride known under the abbreviation PVC. In appearance, it is a hollow thin long tube. The tip may include drainage holes for power derivation, or depending on the application, it may include a differently shaped tip [1].

When introduced into the patient's body, the catheter may include a support wire for transport. This is to prevent deformation and bending, as the catheter is very flexible. This support wire is called a stylet and includes a crimped end cap. This end is located on the side of the catheter that the healthcare professional holds during insertion. The stylet passes through the catheter and must not protrude from the catheter on the side for application into the body [1]. Radek Štohl Brno University of Technology FEEC Department of control and instrumentation Brno, Czech Republic stohl@vut.cz



Fig. 1. A stylet produced by the company commissioning the production of the machine.

B. Purpose

The original production process is only manual and divided into several steps, where each production operator is in charge of one part of the process.

Among the biggest suggestion for building the machine was the process of shortening the wires for the stylets, as this is a possible safety hazard for the operator. When cutting, small parts of the wire can fly off and injure the operator. At the same time, the automation of the process apart from safety will also solve time savings, increase accuracy and reduce costs. By combining the process of shortening and crimping the stylet, there will again be savings due to the efficiency of the material flow, which again leads to a reduction in costs [1].

C. The concept of creating a new machine

The machine was developed to include an automatic and a manual part. The whole machine works as a whole, so it has common control and safety. The safety of the machine works in such a way that if there is a breach in one of the workplaces, the whole machine stops and is brought to a safe state. But if a process failure occurs, when, for example, a step of the state machine is not completed in the required time, the machine control will set a failure at the given station, which will stop only the part of the machine where the given station is located. This means that if there is a malfunction, for example, on the automatic part, the operator on the manual part can continue and her work is not affected.



Fig. 2. Overview of the machine design.

D. Automatic part of the machine

The automatic part of the machine has the task of completing the stylet. The task of this part is to remove the wire from the magazine, shorten the wire to the length according to the recipe and then crimp the end. The finished piece will be thrown into the sink with the finished pieces, where the operator will then pick them up. The machine can also recognize some poorly made pieces. When a bad piece is detected, the process stops and must be removed by the operator. The machine has no sinkhole for bad products. The automatic part is divided into four independent stations with possibility exchange information with each other.

E. Manual part of the machine

This part of the machine is used to manually shorten ureteral catheters to the required length. In the layout of the machine by stations, it is station 5 for manual shortening of catheters. The station is located in front of the automatic part. The process begins by inserting the catheter into the trough and waiting for an activated optical sensor to detect the end of the catheter with a shaped tip. The sensor is located in a small house, which is attached to the board of the manual workplace. There are pre-prepared mounting positions for the houses, where the position is directly prepared for a specific length to be shortened. The operator therefore does not have to measure the distance for shortening. After switching on the sensor, there is a wait for the shortening process to be activated. The process can be started by pressing a button or pressing a foot pedal.

II. WORK PROCESSING

A. Electro project

1) In general

An electrical project was created using the EPLAN 2023 tool and at the same time suitable components were selected. Since the selection of elements is huge, only the most important parts will be listed.

2) Control system

A Siemens SIMATIC S7-1214FC DC/DC/RLY PLC was chosen to control the machine. The reason for choosing this PLC with Fail-safe function was the possibility of programming machine control and safety in one program, when variables can communicate with each other. This is a more convenient solution for the programmer. At the same time, the PLC can communicate using the standard PROFINET communication protocol as well as its PROFISafe superstructure, which enables communication with safety elements at the PLe level according to the ČSN EN ISO 13849-1:2023 standard. The S7-1500 PLC was not used, as it is a smaller machine and the chosen PLC is sufficient for the needs.

3) Visualization

HMI TX110-00VPST from Turck was selected for visualization purposes. Visualization programming takes place via the TX VisuPro development tool, which is available free of charge. The panel has similar properties to the Siemens HMI of the higher Comfort series at the price of the lower Basic series panel, which is why it was chosen. The price for a similar equivalent would be approximately twice as much [2].

4) Elements for the safety of machinery

Machine safety is controlled using a safety PLC and safety modules, where the greatest risk to human health is the mechanical movements of the pneumatic cylinders and the shearing point in the manual shortening workplace. Data exchange between the safety modules and the PLC is ensured using the PROFIsafe communication bus, which is an extension of the PROFINET network. This network is used both for the transmission of normal data and additionally for the transmission of safety data, when it achieves certification up to the PLe level according to ČSN EN ISO 13849-1:2023.

To prevent access to the machine, all doors are equipped with an RFID Guardmaster 440G-LZ security lock from Allen-Bradley. The exception is the rear wing door, where one wing is guarded by the aforementioned security lock and the other by the SI-RFDT-LP8 SI-RFDT-LP8 security RFID switch from Banner company. The leaf with a safety lock contains an aluminum angle for mechanically preventing the leaf with a safety switch from opening in the event of a lock. If the door is locked and the leaf with the switch remains open, the control system will recognize this and prevent the machine from starting. It will be possible to close only after the locks are reopened.

In order to monitor the cutting site at the manual workplace, a safety RFID switch SI-RFDT-LP8 is placed on the cover, which is active only in the case of an active manual process for shortening the catheter.

There are three emergency stop buttons in total. One is located on the control panel, the other on the left side of the machine near the wire magazine and the last one near the manual workstation. The last emergency stop is equipped with a cover against accidental pressing and the option to lock the LOTO (LockoutTagout) system, which is used by machine maintenance when working on machinery. All emergency stop buttons are Harmony XB5 RU 22mm detents from Schneider Electric firm.

In the event of a security breach, all safety elements for controlling the action elements will be switched off. It is a BLDC motor to which two STO (Safe Torque Off) signals from the safety output module are connected. In the case of compressed air, on the FRL unit for the treatment of compressed air there are two VP546-5DZ1-M-D safety solenoid valves (SMC company) with gate position monitoring, which are again controlled by the safety output module. Furthermore, all valves for controlling pneumatic cylinders are controlled by the EX260-FPS1 valve terminal (SMC company) with PROFIsafe communication protocol. In the event of a fault, the power part of the valve terminal is disconnected via a command sent via the PROFIsafe protocol from the safety PLC and all valves return to the middle closed position. Only the safety PLC takes care of turning off all these mentioned safety devices in the event of a malfunction. All physical electrical connections are two-channel with a pulsating signal to diagnose the preservation of the safety function in the circuit.

5) Wire unwinding motor

A BLDC (brushless direct current) motor BG 66x25 dPro PN 24 V with integrated control and encoder was chosen for measuring the wire. The motor is controlled via the PROFINET communication protocol. The motor is powered of 24 VDC and has a rated current of 8.03 A. The encoder is the incremental type including 4096 pulses per revolution [3].

6) Pneumatic components

All pneumatic components and associated equipment will be from SMC company.

B. Safety assessment

1) Safety automation builder

Safety automation builder (SAB) is a tool from Rockwell Automation. It has been used for the design and analysis of security systems. For the purposes of the work, a project was created in this tool, where safety risks were defined and the required PLr was subsequently calculated according to the EN ISO 13849-1:2023 standard. Protective measures were applied to all risks and it was subsequently determined whether the risk was reduced to an acceptable level. The risks were taken into account in terms of the type of danger, the operation performed and the level of operator training [4].

2) Sistema

The Sistema tool enables safety assessment according to the EN ISO 13849-1:2023 standard. It was used to calculate individual safety functions based on individual device functional data obtained automatically using libraries supplied directly from industrial device manufacturers or by manual input [5].

3) Documents for EU declaration of conformity

The default revision was created to meet the conditions and requirements of EN 60204-1:2018. For these purposes, a test protocol was drawn up, where electrical equipment was measured to see if it complied with the above-mentioned standard. The measuring device METREL MI 3100 SE and ILLKO REVEX plus were used for this. At the same time, an overall verification of the requirements for meeting the standard was developed. From measurements and development of protocols, it was found that the default revision meets the requirements.

As part of the analysis, there was a risk assessment according to EN ISO 12100:2010, and EN ISO 13849-1:2023. The assessment is recorded in the risk assessment report. The assessment includes mention of the use of some standards, types of users, machine modes, possible performed actions, ranges of movements, effects of energies and many others.

All requirements for issuing an EU declaration of conformity for machinery have been met. The declaration contains all the standards used in the creation of machinery.

C. Software and visualization design

1) General operation of the program

The program will be written in TIA Portal V17 from Siemens. According to the proposal, the program will be divided into a control part for the joint operation of the program, a safety part for servicing the safety of the machine and a part for individual stations. Each station will contain a state machine. Individual stations will exchange information with each other for the overall functioning of the machine.

2) Controlling part of the program

This part of the program will be in charge of general tasks for the operation of the machine. Alarm handling, sensor data processing, control of individual valves and actuators, visualization data handling and various other auxiliary tasks will mainly be carried out here. In general, it takes care of the control of the common parts for the machine independently of the individual station.

3) Safety part of the program

This is where all the machine safety staff will be. Only certified functions from Siemens will be used for work. There is a limitation of programming instructions as it is a safety part. Furthermore, specially created data blocks will be used to share data between the regular and security programs. Only one function from the safety organizational block will be called to handle the instructions, triggered by a cyclic interrupt with the highest priority compared to the standard program.

Any modification to this part of the program will result in a change to the checksum, and normal playback to the PLC cannot be performed without consent to the change and a change to a new checksum with a signature.

4) General operation of visualization

In general, a common background (template) will be created for the entire visualization, which will include the possibility of user management, language change, an overview of active modes, error indications, and the opening of a pop-up window, an overview of the status of security elements, and the current date and time. The visualization will be divided into four main parts.

The first group of screens will serve the main information about the stations and the operation of the machine. Here will be screens for the main overview and control of individual stations. In this group of screens, for example, you can change the recipe, set the batch for production, monitor the statuses of individual stations, change the maximum times of individual steps, start the step mode of the state automaton at each station, confirm bad pieces and turn on the lighting.

The second group will be used to control the machine in manual mode. Furthermore, the setting of the parameters of individual actuators and sensors will be included here. In total, there will be seven screens for manual control and adjustment of pneumatic cylinders, valves, BLDC wire gauge motor, camera sensor, vibrating feeder, ITV proportional valve, and sensors.

The third group will be intended for setting up, and monitoring the entire machine. It will be divided into six parts. The first part will be for recipe management. The second one includes the monitoring the inputs and outputs of the machine. The third part will contain the machine monitoring data, where there will be information about the machine operation, calculated OEE (overall equipment effectiveness) of the machine, and compressed air consumption as a function of time. The fourth section contains the settings and information on predictive maintenance of the actuators. The fifth part is the audit trail, where it will be monitored that user made which change. The last sixth part will be the detailed setting of the bed at the station 3 to calculate the values for guiding the wire.

The last fourth group is alarms. Here will be the current alarms, and the total alarm history with filtering.

D. Electrical wiring

The electrical connection was implemented based on the design, and wiring diagram created before implementation. Thanks to the prepared scheme, and board layout, a lot of time was saved.



Fig. 3. Connecting the cabling in the switchboard.

E. Software creation and visualization

The software was created on the design before actual implementation. Thanks to the design, a lot of time was saved and the operation was efficient. The program itself is logically organized and custom functions are used for almost the entire program. The visualization was created based on the design. The goal was to create it in such a way that it was clear and organized for the best possible orientation and detection of all machine states. Emphasis was also placed on the size of text and images. Everything was consulted with potential users.

F. Testing and production rollout

The machine has been successfully tested and debugged to all states. All users have been trained. The production test, validation was successful and now it is waiting to be released to production.

III. CONCLUSION AND EVALUATION

The machine was successfully designed, manufactured, documented and handed over to the client. Now it is ready to start the production for which it was made. A lot of experience was gained from building the machine, when it was possible to build a large machine from start to finish. It was also necessary to think about the future, manage the costs for construction and work. Unfortunately, this work could not be mentioned here in its entirety, as it is very extensive and everything cannot be described here.



Fig. 4. A machine made and prepared at a designated location in production.

References

- KUŘÁTKO, Petr. Development of a new automated process for the production of ureteral catheters. Diploma thesis, supervisor Lukáš Vaculík. Jihlava: Jihlava Polytechnic University. Also available from: <u>https://isz.vspj.cz/bp/get-bp/student/71744/thema/9834</u>
- Product TX110-00VPST data sheet. Online. 2023. Turkey. Available from: <u>https://www.turck.cz/datasheet/_cz/edb_100002313_ces_cz.pdf</u>.
- [3] BG 66 dFor CO/IO/PN/EC/EI. Online. 2023. Dunkermotoren. Available from: https://www.dunkermotoren.de//media/project/oneweb/oneweb/dunkerm otoren/downloads/pdf/downloads/manuals/brushless-dcmotors/24_bg_66dpro.pdf?la=pt-br&revision=43c382ba-752e-4510-8896-a7de5834518d.
- Safety Automation Builder. Online. Rockwell Automation. Available from:https://www.rockwellautomation.com/en-us/capabilities/industrialsafety-solutions/safety-automation-builder.html.
- [5] SISTEMA Software Assistant. Online. German Social Accident Insurance (DGUV). Available from: <u>https://www.rockwellautomation.com/en-us/capabilities/industrial-safety-solutions/safety-automation-builder.html</u>

Design of Miniaturized Logarithmic-Periodic Antenna

Matěj Podaný Department of Radio Electronics Brno University of Technology Brno, Czech Republic xpodan00@vut.cz Jaroslav Láčík Department of Radio Electronics Brno University of Technology Brno, Czech Republic lacik@vut.cz

Abstract—This paper presents a design of logarithmic-periodic antenna for frequency band 500 MHz to 4000 MHz. In order to decrease antenna dimensions, the longest radiation elements are bent. The antenna in the whole frequency band has a voltage standing wave ratio lower than 2 and achieves a gain exceeding 8 dBi, while its width was reduced by 7.89 %.

Keywords—logarithmic-periodic antenna, broadband antenna, antenna miniaturization, impedance matching

I. INTRODUCTION

Logarithmic-periodic antenna, specifically the logarithmicperiodic dipole array, is a type of wideband antennas. High gain and low voltage standing wave ratio (VSWR) is achievable across the whole desired frequency band. At the same time, they are light weight, low-cost, and easy to manufacture. The ability to deviate from certain theoretical parameters can also be advantageous. For example, by utilizing elements of the same diameter rather than elements with many different diameters, similar performance can be achieved. Typical use for logarithmic-periodic antenna is in broadband applications due to its wide bandwidth possibilities while maintaining relatively high gain [1].

This contribution presents design of miniaturized logarithmic-periodic antenna. Firstly, the basic description of logarithmic-periodic dipole array is provided. Then, a simplified design procedure is described including some findings, followed by final antenna optimization. Subsequently, antenna miniaturization technique is proposed. Finally, the design performance is presented.

II. ANTENNA DESIGN

A logarithmic-periodic dipole antenna is an array of parallel dipoles (Fig. 1) that have specific lengths, spacings between them, and diameters. With lower frequency, the dimensions increase by the inverse of the geometric ratio τ or with higher frequency, the dimensions decrease by τ . Another design parameter is spacing factor σ . Those two parameters are connected by:

$$\alpha = \tan^{-1}((1-\tau)/(4\sigma)), \tag{1}$$

where α is angle according to the Fig. 1 [1].



Fig. 1. Logarithmic-periodic antenna with 5 elements and constant t_n

A. Initial Desing

To find basic dimensions of the antenna, the initial design was carried according to steps described in [1]. For the desired antenna directivity of 7 decibels relative to isotropic (dBi), initial directivity D_0 was selected to 7.65 dBi. Extra initial directivity reserve was added in order to make sure that directivity will be higher than 7 dBi in whole frequency band. According to the nomograph in [1], two important design parameters were selected: $\tau = 0.876$, $\sigma = 0.16$. Then, angle α was calculated from (1). To determine all necessary dimensions of the logarithmicperiodic dipole array, step by step calculation was made according to [1].

• At first, the initial design was simulated in program ANSYS HFFS with a discrete port. All parts of the antenna model were assigned as perfect electrical conductors (PEC). The results were not perfect. There were narrowband increases in VSWR, and narrowband directivity drops, both at lower frequencies. After this, some changes to the design were made. Element diameters were simplified to four groups. To facilitate production, all dimensions were rounded to whole millimeters. At the same time, the presence of antenna cover was considered. For that reason, the expected realized gain was increased to 8 dBi to meet the requirement for at least 7 dBi for the final antenna with cover. Following this, the antenna was successfully optimized by extending the spacings between elements *s*

and adjusting the distance *h* from Fig. 2. The change of the targeted realized gain was not significant, and the antenna met the new requirements despite the fact that the initial directivity was selected to 7.65 dBi. Also, no change was made to the parameters τ and σ .

- Secondly, the discrete port was replaced with a 50-ohm coaxial cable and wave port at the end of the cable was applied. The results of the simulation were similar in terms of directivity. However, VSWR was significantly worse, as expected. To ensure better impedance matching, distances h and L_z were adjusted. Distance L_z should be about one-eighth of the maximum wavelength. In this specific design, it was found out that higher values of L_z gave better results in a combination with slightly higher h, comparing with the initial values. The effect of previous optimalization was global decrease of VSWR and small global increase of directivity. Eventually, the antenna was partially optimized, and the best combination of h and L_z was found. At this point, realized gain was over 8 dBi for low and mid frequencies and VSWR was under 2 for low frequencies.
- Finally, it is important to successfully perform impedance matching in order to optimize VSWR. To achieve this, the booms may be extended from the feeding point excluding the feeding itself [2]. This situation is shown in Fig. 2 and the important parameter is *e*. This means that the inner conductor of coaxial cable must lead through the booms. One possible way to do this is to drill a hole through the booms (Fig. 3). It was discovered that size of the hole does not significantly impact antenna performance.

Then, the optimization of parameter e is needed. In this specific scenario, the extension e is only few millimeters after the feeding point so that the miniaturization of the antenna is not violated. Extending e too much can lead to significant impedance mismatch, with the mid frequencies being most impacted. The effect of extending or shortening the feeding point was explored. It was discovered that changing the feeding



Fig. 2. Matching parameters of the logarithmic-periodic antenna while $t_n = 0$



Fig. 3. Extended U-shape aluminium profiles (booms) with the feeding point

point by a few millimeters does not significantly alter the performance of the antenna, however the performance may be shifted in frequency. It is recommended that the distance h is also optimized together with parameter e [2]. However, during this specific design optimization, it was found that the previous (global) optimization of parameters h and L_z is sufficient. No change of parameter h was needed during the optimization of e.

B. Antenna Miniaturization

There is a continuing need for miniaturization and for that reason, aluminium strips were chosen. They may be bent to decrease the antenna width. Of course, resonance lengths of radiators were preserved. Fig. 4 depicts a proposed design of logarithmic-periodic antenna.

Considering the production of this antenna should be uncomplicated, the elements are bent only once each end. Additional bending at the end of the dipole (Fig. 5) may reduce the antenna width by up to 50 %. However, due to the high complexity of the production, the impact on the antenna is not investigated. Elements are bent on the other side than their boom is and there are two reasons for that. Firstly, the gap between the feeding lines provides more space for longer bend. Secondly, this design achieves the connection between the booms and the elements by two aluminium rivets. To ensure effortless production, the rivets need to be accessible. Bending the elements on the other side allows for easy access to the rivets with a tool. Rotating the elements by 180 degrees (upside down) would make even more space for longer bending, although according to the simulation, the performance was decreased.

III. DESIGN PERFORMANCE

The proposed logarithmic-periodic antenna design was simulated using the dimensions provided in Table 1. Total antenna length L_c is 6 millimeters longer than the sum of the parameters L, L_z , and e. According to the Fig. 2, the short is negligibly thin, but practically, it must have some width due to mechanical reasons. In this specific situation, the short width is 6 millimeters and other sizes match the booms dimensions (Fig. 4c). For clarity, parameter l_{21} is width of the largest element (total antenna width; Fig. 1), parameter h is center to center



Fig. 4. Proposed design of miniaturized logarithmic-periodic antenna with 21 elements: a) perspective, b) top view, c) front view

distance between the booms (Fig. 2). Also, the conducting materials, previously perfect electrical conductors, were assigned to aluminium (elements, booms and the short) and copper (inner and outer conductor of coaxial cable).

Fig. 6 contains the radiation pattern of the proposed antenna design at 2 GHz. The maximum is located at the tip of the antenna (from feeding point to the direction of extension e), as expected. Fig. 7 contains VSWR and realized gain.

a)

10.00



Fig. 5. Additional bending of antenna elements



Fig. 6. Radiation pattern at 2 GHz (red: E-plane, green: H-plane)



Fig. 7. Antenna performance: a) VSWR, b) realized gain



TABLE I. BASIC ANTENNA DIMENSIONS

Variable	Size [mm]	Variable	Size [mm]
Lc	910	Lz	178
L	718	е	8
h	14	l_{21}	280

IV. CONCLUSION

According to the simulations, both VSWR and realized gain have met the requirements while maintaining simplicity and miniaturization of the antenna. In case of realized gain, the results exceeded the expectations. For the whole frequency band, average realized gain is slightly under 9 dBi, so the presence of antenna cover should not decrease the performance under the requirements. In case of antenna miniaturization, the antenna width was reduced by 7.89 %.

An interesting finding is that the bending of elements slightly increased the performance of the antenna. That might be caused by a change in the antenna impedance, which in this case better matches the impedance of the coaxial cable.

The effective bandwidth has been slightly increased to cover the range from 450 MHz to 4050 MHz, thanks to the realized gain exceeding 8 dBi across this whole band.

ACKNOWLEDGMENT

This work was supported by the Internal Grant Agency of the Brno University of Technology under project no. FEKT-S-23-8191.

REFERENCES

- C. A. Balanis, Antenna theory: analysis and design, 3rd ed. Hoboken, NJ: John Wiley, c2005.
- [2] Q. Zhao and X. Yin, "The influence of feed tube changes on performance of log periodic dipole antenna," 2012 International Conference on Microwave and Millimeter Wave Technology (ICMMT), Shenzhen, China, 2012, pp. 1-4, doi: 10.1109/ICMMT.2012.6230340

Mapping and analyzing of signal coverage of 4G/5G mobile networks

1st Michal Baranek Department of Telecommunications Brno University of Technology Brno, Czech Republic xbaran17@vut.cz 2nd Ladislav Polak Department of Radio Electronics Brno University of Technology Brno, Czech Republic polakl@vut.cz 3rd Jan Kufa Department of Radio Electronics Brno University of Technology Brno, Czech Republic kufa@vut.cz

Abstract—This paper addresses the enhanced measurement of signal coverage, capacity, and reliability in mobile networks, particularly with the growing prevalence of 4G and 5G technologies. Given the escalating importance of these networks in everyday activities, there arises a demand for open-source solutions to evaluate and enhance their performance effectively. The objective of this research is to analyze gathered data to pinpoint areas necessitating network enhancements and to develop opensource software and hardware solutions for extracting essential performance metrics (KPIs) from 4G/5G networks. The proposed system offers an interface for assessing network performance and signal coverage, enabling cost-efficient measurements across diverse environments.

Index Terms—coverage mapping, 4G, 5G, key performance indicators, mobile networks, optimization, machine learning, interpolation, data analysis

I. INTRODUCTION

The increasing dependence on '4G' and '5G' networks highlights the crucial need to comprehend signal coverage, capacity, and reliability for uninterrupted connectivity. This paper introduces open-source software and hardware solutions designed to extract key performance indicators (KPIs) from these networks. We conducted static measurements for seven days at the BUT campus in block A04 and dynamic measurements throughout March 2024 in Brno City, concentrating specifically on LTE Band 3 (1800 MHz). Our findings unveil fluctuating KPIs even in static scenarios, with discernible timeof-day influences. The outcomes of dynamic measurements are leveraged to construct coverage maps.

II. METHODOLOGY

This paper primarily focuses on driving conducted by a bicycle or a car to measure long distances or walking tests to collect substantial data. The measurement methodology, as depicted in Fig. 1, involves collecting data from the mobile site using a measuring device, storing it in log files (preferably on the cloud but now on the device's internal memory), and retrieving it for visualization. This aligns with practices used by most Mobile Network Operators (MNOs) or companies specializing in benchmarking mobile networks. They employ drive tests categorized by environment and dynamic measurements, utilizing sophisticated equipment and vehicles for realworld scenario capture. These measurements provide crucial



Fig. 1. Measurement methodology.

input for optimization and planning to enhance overall network performance.

III. MEASUREMENT DEVICE

Currently, various applications and devices, such as "G-NetTrack" and "TEMS Pocket," support mobile network measurements, as seen in papers [1] and [2]. This paper introduces an effort to design a novel measurement device and explores the selection between 4G and 5G modules integrated with GNSS technology.

Current module availability narrowed the choice to Simcom SIM8200EA-M2 and Quectel 5G RM502Q, both supporting high data speeds in both directions. Both manufacturers offer development boards, easing module integration. Both modules share a similar price range and have USB 3.1 for potential high-speed tests.

Despite the advantages of Quectel in [3], where both modules were tested for DL/UL speeds and latency using tools OpenSpeedTest, LibreSpeed, iPerf3, RTT, and ping. The Waveshare 5G module was selected for the measurement due to difficulties encountered in making the Quectel module operational.

Further detailed information about the Waveshare 5G module is provided in Table 1. However, it is noteworthy that this module only supports the n78 band (3600 MHz) in 5G NSA, which is currently deployed in the Czech Republic.

A. Processing Unit

After evaluation, Raspberry Pi (RPi) 4 Model B was chosen for its cost-effectiveness and high performance. Its Quad-core Cortex-A72 64-bit SoC @ 1.8 GHz and two USB 3.0 ports with 5 Gbps transfer capacity ensure efficient throughput. The

Key aspects	SIM8200EA-M2
5G Category	5G NSA/SA
5G Frequency Bands (NSA)	n78
4G Category	Cat-20
4G Frequency Bands	All Supported
Peak Download Rate	4 Gbps
Peak Upload Rate	500 Mbps
Max. RF Output Power	23 dBm
Satellite Systems	Supported
USB	USB 3.1
Price	300€

TABLE I SIM8200EA-M2 SPECIFICATIONS.

RPi's popularity and affordability make it practical, supporting Python programming and essential packages. The final setup includes a 40000 mAh battery power bank, Waveshare 5G module, RPi, and GNSS antenna.

IV. DATA VISUALIZATION

Inspired by the 'TEMS Discovery' by InfoVista, our visualization tool offers a map interface for exploring measured areas with color-coded signal quality samples. Utilizing Python tools, specifically *Dash* and *Plotly*, the application provides interactive data visualization, meeting project requirements effectively.

The *Dash* application comprises a dropdown data selector, map display using OpenStreetMap API, waveform graphs, and filtering capabilities. Structured into two parts, it runs analytics on a high-processing device and measurement settings directly on the Raspberry Pi using Gunicorn. Users can access settings via a browser on a connected device to the Pi hotspot, facilitating seamless parameter customization.

V. KEY PERFORMANCE INDICATORS

In terms of KPIs, critical parameters are analyzed to understand mobile network properties such as coverage, capacity, quality, reliability, etc. These KPIs play a vital role in improving Quality of Experience (QoE) and Quality of Service (QoS).

- **RSSI (Received Signal Strength Indicator)** Received Signal Strength Indicator (RSSI) is defined as the linear average of the aggregate received power, measured in watts. It is observed within the configured OFDM symbol and across the measurement bandwidth, encompassing N number of RBs [4].
- **RSRP** (**Reference Signal Received Power**) RSRP is characterized as the linear average of the power contributions (measured in watts) from the RE responsible for carrying the CRS conducted in designated measurement frequency bandwidth, as specified in [4].
- **RSRQ** (Reference Signal Received Quality) RSRQ is expressed as the ratio of N times the RSRP to the E-UTRA carrier RSSI, where N represents the number of RBs within the E-UTRA carrier RSSI measurement bandwidth. Both the numerator (N times RSRP) and the

denominator (E-UTRA carrier RSSI) measurements are conducted over the identical set of RBs [4].

• **RSSNR** (Reference Signal Signal to Noise Ratio) -RSSNR is the power ratio of usable signals measured from CRS to the average noise within the measurement bandwidth.

VI. MEASURED DATA

A. Static Signal Measurement

A static indoor measurement was conducted at Pod Palackeho Vrchem dormitory on the BUT campus from 9:40 a.m. on Friday, March 22nd to 9:40 a.m. on Friday, March 29th. The measurement device was positioned under a table in a secondfloor room in block A04, aligning with the base station located at Brno - Královo Pole, Kolejní 2905/2, block A04 of the BUT dormitory, with CelIID 790278. Despite being on the opposite side of the building from the cell, there was no Line Of Sight with the cell.

Fig. 2 illustrates all measured samples and 15-minute averages of all measured signal KPIs.

The measurement was conducted during normal college operation, resulting in significant user presence on weekdays, while fewer users were present during the weekend.

The measurement utilized the proposed solution discussed in this work within the operational LTE network of the Czech MNO T-Mobile. Signal KPIs (RSRP, RSRQ, RSSI, RSSNR) were recorded at a frequency of one sample per second using the provided hardware and software.

The measuring module was configured to the specified cell 790278 in LTE band 3 (1800 MHz, FDD (Frequency Division Duplex)).

As depicted in Fig. 2, variability is evident even in the static measurement of cellular network signal metrics. Notably, an intriguing pattern emerges, revealing periodic behavior in the RSRP and RSSNR metrics corresponding to the time of day. During peak working hours, RSRP and RSSNR values tend to decrease due to high cell load. Conversely, in the early



Fig. 2. Static Measurement Conducted at Pod Palackeho Vrchem dormitory on the BUT campus from March 22nd to March 29th.

morning hours, typically between 1:00 a.m. to 6:00 a.m., these values increase, peaking around three to four o'clock in the morning daily.

The relationship between RSRP and RSSI demonstrates a significant negative correlation, with a correlation coefficient of -0.8.

Concurrently, notable deviations are observed during peak cell load, typically occurring between 8:00 p.m. to midnight. During these periods, the RSRP data exhibits significant downward spikes, corroborating the bimodal distribution evident in the histogram of RSRP values in Fig. 3. This distribution highlights two prominent values, predominantly around -80 dBm and, during peak load hours, approximately -110 dBm.

The RSRP demonstrates a substantial standard deviation of 15.145 dBm, indicating significant variability in signal strength even during static measurements. In contrast, RSRQ, RSSI, and RSSNR exhibit low standard deviations, suggesting a consistently clean and strong signal when measured in a static environment. Fig. 4 presents more detailed distributions of individual metrics through boxplots.

B. Dynamic Signal Measurement

In the dynamic measurement phase, multiple walk and public transport tests were conducted to assess signal coverage across the city of Brno. The measuring device, locked to LTE band 3, was securely placed in a backpack. These measurements were conducted throughout March. The measured route in the form of RSRP points is depicted in Fig. 5.

Fig. 6 illustrates individual distributions depicted through boxplots of respective measurements. The median values of RSRP and RSSI exhibit greater variability across different data sets compared to RSRQ and RSSNR.

The empirical cumulative distribution function (CDF) of the RSRP sample from individually measured sets in Fig. 7 reveals noteworthy insights. On Friday, March 29, approximately 50% of the sample fell below the -141 dBm limit. Similarly, on



Fig. 3. Histogram of Static Measurement Conducted at Pod Palackeho Vrchem dormitory on the BUT campus from March 22nd to March 29th. Bin size = 2 dBm.



Fig. 4. Boxplots of a Static Measurement Conducted at Pod Palackeho Vrchem dormitory on the BUT campus from March 22nd to March 29th.

Thursday, March 14, with about 50% probability, the sample is below the -135 dBm limit. In contrast, the last dataset displays more favorable results, with only 16% of the sample below the -140 dBm threshold on March 6 and 30% on March 30. The better performance observed on March 6 is likely attributed to it being conducted as a walk test, while the others were conducted in public transport vehicles, resulting in significant attenuation.

C. Interpolation of Measured Data

For approximating coverage maps python package *scipy.interpolate*, is used. Three interpolation methods were applied to determine the most suitable coverage for a given area, as can be seen in Fig. 8 using a designed visualization tool.

- Bilinear Interpolation calculates a weighted average of the four nearest points.
- The Nearest Neighbor preserves input values without alteration, producing blocky results.



Fig. 5. RSRP Points of a Dynamic Measurement Conducted at Brno city throughout March 2024 on LTE Band 3 (1800 MHz).



Fig. 6. Boxplots of Signal Metrics in a Dynamic Measurement Conducted at Brno City Throughout March 2024 on LTE Band 3 (1800 MHz).



Fig. 7. Empirical Cumulative Distribution Functions of RSRP Metric from Individual Measured Datasets on LTE Band 3 (1800 MHz).

• Cubic Convolution analyzes the 16 closest points to interpolate with a smooth curve, excelling in smoothing continuous data.

The bilinear method was identified as the optimal approach, displaying an acceptable distribution of interpolated samples. Conversely, the nearest neighbor method appeared visually complex, while enhancing the accuracy of the cubic method mandates additional track measurements within a specified area to achieve improved interpolation.

VII. CONCLUSION

In this study, we introduced a methodology for gathering KPIs from both 4G and 5G networks. By employing the SIM8200EA-M2 module and RPi 4 Model B alongside Python packages such as *Plotly* and *Dash*, we devised a comprehen-



Fig. 8. RSRP Coverage Map via Bilinear, Nearest Neighbor, and Cubic method.

sive approach for data collection, visualization, and analysis¹. Through static measurements conducted over 7 days within LTE band 3, we amassed and scrutinized over 600,000 KPI samples, unveiling notable fluctuations in cellular network signal metrics. Daily patterns emerged in metrics like RSRP and RSSNR, with peak load hours exhibiting diminished values and a robust negative correlation observed between RSRP and RSSI. Furthermore, dynamic measurements conducted across the city of Brno throughout March underscored the variability across datasets, with RSRO and RSSNR maintaining stability. The empirical CDF depicted divergent performance levels, with walk tests yielding superior results compared to measurements conducted within public transport vehicles. Additionally, we employed three interpolation methods to construct an LTE band 3 coverage map, determining the bilinear method to be the most optimal. Plans include extending measurements to encompass 5G and other bands within LTE technology.

REFERENCES

- [1] El-Saleh, Ayman and Abdullah Al Jahdhami, Majan and Alhammadi, Abdulraqeb and Shamsan, Z.A. and Shayea, Ibraheem and Hassan, Wan, "Measurements and Analyses of 4G/5G Mobile Broadband Networks: An Overview and a Case Study," Wireless Communications and Mobile Computing, vol. 2023, 04 2023.
- [2] Minovski, Dimitar and Ögren, Niclas and Mitra, Karan and Åhlund, Christer, "Throughput Prediction Using Machine Learning in LTE and 5G Networks," IEEE Transactions on Mobile Computing, vol. 22, no. 3, pp. 1825–1840, 2023.
- [3] Lackner, Thorge and Hermann, Julian and Dietrich, Fabian and Kuhn, Christian and Angos-Mediavilla, Mario and Jooste, Wyhan and Palm, Daniel, "Measurement and comparison of data rate and time delay of end-devices in licensed sub-6 ghz 5g standalone non-public networks," Procedia CIRP, vol. 107, pp. 1132–1137, 01 2022.
- [4] 3GPP "LTE: Evolved Universal Terrestrial Radio Access (3GPP TS (E-UTRA); Physical layer; Measurements 14.2.0 Release 14)." Online, April 2017. 36.214 version https://www.etsi.org/deliver/etsi_ts/136200_136299/136214/14.02.00_60 /ts_136214v140200p.pdf

¹Data visualization tool: https://github.com/michalizn/coverage-analysis



ABB v Brně Lepší svět začíná s vámi

Jsme součástí globálního technologického lídra působícího v oblastí elektrotechniky, robotiky, automatizace a pohonů. Z Brna dodáváme rozváděče vysokého i nízkého napětí, modulární a digitální systémy, přístrojové transformátory a senzory. Tyto produkty míří do nejnáročnějších projektů ve více než 100 zemích světa. Najdete ke mimo jiné v nejvyšší budově světa Burdž Chalífa, v londýnském metru nebo v datovém centru Facebooku. Svou pobočku zde má i Operační centrum.



Přibližně 2100

zaměstnanců



Tržby 10 mld. Kč



Významné centrum výzkumu a vývoje



Nejnáročnější projekty ve 100 zemích světa

Kontaktujte naši náborářku:

lucie.vasatova@cz.abb.com

4 výrobobní závody, centrum digitálních řešení, R&D, servis a operační centrum

Studentům a absolventům nabízíme:

- Placené brigády
- Odborné praxe
- Stáže
- Trainee program
- Vedení balakářské nebo diplomové práce
- Hlavní pracovní poměr po ukončení studia



Lucie Vašatová

+420 704 997 141

Virtuální prohlídka



Webové stránky Operační centrum







Rozváděče vysokého napětí

Výroba rozváděčů vysokého napětí se v našem závodě těší dlouhé tradici a pyšní se statusem největšího závodu svého druhu v Evropě. Co do počtu zaměstnanců, tak i plochy výroby, se jedná o největši jednotku v rámci brněnského závodu. Své produkty dodává do přibližně 100 zemí světa.





Rozváděče nízkého napětí

Rozváděče nízkého napětí fungují v režimu tzv. virtuálního závodu a sdílí svou kapacitu se závodem v polské Bielsko-Białe. Opět se jedná o jeden z největších závodů svého druhu v Evropě. Jeho produkty jsou dodávány do přibližně 50 zemí světa.

Přístrojové transformátory a senzory

Jednotka přístrojových transformátorů a senzorů je 2. největší v rámci brněnského závodu. Zároveň se jedná o největší závod svého druhu na celém světě. Vůbec první transformátor byl v závodě vyroben již v roce 1919 a jejich výroba je tak založena na více než stoletém know-how.





Modulární systémy

Přestože se jedná o mladou jednotku, má už za sebou několik významných milníků. Své modulární rozvodny, dodávané obvykle v kontejnerovém řešení, dodává do projektů napříč odvětvími. Nejvýznamnější z nich jsou data centra a projekty v oblasti těžby a distribuce ropy a zemního plynu.

Centrum digitálních řešení

Cílem jednotky Digital Solution Center je realizace inteligentních řešení v oblasti monitoringu, sběru a archivace dat a hlavně vzdáleného řízení distribuce energie. Zákazníkům nabízí řešení na míru, flexibilitu a vysokou přidanou hodnotu.





Servis přístrojů vysokého a nízkého napětí

Jednotka servisu spolupracuje se zákazníky v rámci oprav, údržby a vylepšování jejich zařízení po celou dobu jejich životního cyklu. U svých zákazníků v přibližně 100 zemích světa stráví technici asi 40 tisíc hodin ročně. Jednotka využívá nejmodernější nástroje včetně tréninkového centra pro zákazníky nebo rozšířené reality.

Technologické centrum

Technologické centrum se zabývá výzkumem a vývojem produktů. V rámci ABB se celosvětově jedná o velmi významnou výzkumnou jednotku. K dispozici má špičkově vybavenou laboratoř umožňující provádění nejrůznějších simulací a výrobu a testování prototypů.





Operační centrum Evropa

Operační centrum Evropa (EUOPC) je globálním centrem ABB pro průmyslovou automatizaci a elektrifikaci. Zákazníky má asi v 50 zemích. Z EUOPC pocházejí řídicí systémy pro nové generace ropných plošin. Zajímavé jsou projekty přispívající ke stabilizaci el. sítě při přechodu na generování a využití zeleného vodíku nebo elektrifikaci lodní přepravy.

Platform for Digital Predistortion Based on RFSoC

Tomáš Kříčka Brno University of Technology FEEC Brno, Czech Republic <u>xkrick00@vutbr.cz</u>

Signal predistortion occurs in modern communication systems and in other systems requiring linear signal amplification. The objective is to learn about signal predistortion, the QORVO 1800 MHz FDD radio front-end and to implement the front-end control using RFSoC and Matlab computer program. I focus on controlling individual attenuators and amplifiers on the QORVO front-end with a program written in C. Using a single board computer Raspberry Pi 4B running Linux. Furthermore, the work involves the application of a SCPI server that operates based on TCP packets. As part of the introduction to the front-end, the power characteristic of the power amplifier is measured.

Keywords— signal distortion, signal predistortion, DPD, power amplifier, PA, RFSoC, radio frequency, RF, frontend, QORVO 1800 MHz, Raspberry, Linux, GPIO, SPI, C, SCPI, TCP, Matlab

I. INTRODUCTION

To communicate wirelessly, it is necessary to propagate the signal using a wave. This is achieved by using radio frequency (RF) communication, which takes place between two or more users. To transmit and receive the radio signal, power amplifiers must be used because the signal from the generator does not have enough power.

When amplifying a PA (power amplifier) signal, it is inefficient to move only within its linear part of the conversion characteristic. When amplifying above the one-decibel compression point, distortion occurs, which limits the signal (the PA is saturated) and introduces unwanted components into the spectrum.

Signal distortion can be avoided by pre-distortion. This method can be used to operate the PA in the saturation region. Pre-distortion takes place in the Matlab scripts, where it is adjust based on the observation of the output signal of the amplifier. This is done to compensate for the amplification characteristics of the PA. The result is then a linear amplification of the signal in the saturation.

II. POWER AMPLIFIER

One of the main components of wireless communication is RF power amplifier (PA). Since the modulated signal doesn't have enough power to be transmitted into free space via an antenna, it must be amplified to a predefined value, which is often expressed in dBm. Jan Král Masaryk University Faculty of Informatic Brno, Czech Republic 2403998@mail.muni.cz

The negative behaviour of PA is non-linearity in amplified signal. This nonlinearity is due final supply voltage and real onchip structure. The linear area of PA is taken under 1 dB compression. Beyond this point amplified signal starts to saturate, this is called "clipping" or "cut off" [1]. Higher harmonics (second and more order) begin to form in the signal Fig 1.

These higher harmonics are too close to the useful signal, to be filtered out by a low-pass filter.



Fig. 1. Distorted signal simulation.

III. DIGITAL SIGNAL PREDISTORTION

The digital predistortion (DPD) is used for example in modern communications and others critical applications. DPD is allowing us to use PA more efficiently however it makes more demand on processing power on the side of transmitter.

DPD systems are basically two nonlinear systems, which form one linear system [2]. In the area, where the amplified signal starts to saturate, the predisposition signal starts to bend in the opposite direction. The resulting signal reaches much more linearity, then basic signal Fig 2.



Fig. 2. The principle of a linear signal. [2]

IV. QORVO 1800 MHZ

For DPD we use power amplifier QPA9903, which is part of small cell RF front-end for RFSoC called QORVO 1800 MHz. This platform was chosen because of the suitable PA characteristics. PA has gain +32 dB (point of 1 dB compression is about +28 dB on the output) [3]. This PA distorts signal which is then followed by DPD.

The one-way in font-end is composed of filters, low noise amplifiers (LNA), PAs and digital step attenuators (DSA). All these active components are controlled from external device Fig 3.



V. APPLICATION

To implement digital signal predistortion, a control unit is required to set active components on front-end. In this solution, a Raspberry pi 3B is used. The program in Raspberry has two functions, one is to control the components and the second one works as a SCPI server for receiving commands from PC. RFSoC holds the position of generating and receiving signal amplified from QORVO Fig 4, 5.

Dgital predistortion is made on PC using Matlab scripts. These scripts set the values for QORVO and the signal parameters for RFSoC. Matlab commands are transferred via TCP packets..



Fig. 4. Application block diagram.



Fig. 5. Measuring workspace.

VI. MEMORY POLYNOMINAL MODEL

For signal predistortion a MP model is used, which is usually derived from Volterra series. MP model can be simply achieved by linearized signal, via predistortion signal from PA [4][5]. The discrete MP model is given (1):

$$y[n] = \sum_{k=1}^{K} \sum_{q=1}^{Q} b_{k,q} x[n-q] |x[n-q]|^{k-1} .$$
(1)

The maximum nonlinearity order and memory length is represented via K and Q. x represents the MP input and $b_{(k,q)}$ the coefficients.

VII. MEASUREMENTS

The measurement verified all proposed application (PC, raspberry, RFSoC and QORVO). The transmit way in QORVO was set with 1 dB attenuation and power was changed in DAC on RFSoC. When the PA output is 22 dBm, the signal started to be distorted. This value is approximately 6 dBm as declared in the datasheet. After this output value, the PA began to saturate Fig. 6.

Measurement was performed using a script in Matlab and all predefined scenario were executed automatically.



In Fig. 7 shows the spectrum of the distorted signal, for demonstration QAM modulation with input power -25 dBm is used. This modulation has a broader bandwidth than, for example, a sin modulated signal. The main difference is that the spectrum has, in addition to the useful signal, a sidebar descend of power.



Fig. 7. Spectrum of distorted signal.

CONCLUSION

The result of this paper is an RFSoC-based signal predistortion platform. This application uses RF front-end Qorvo 1800 MHz to amplify the transmitted signal. Qorvo is controlled by single board computer Raspberry pi, which sets attenuations values and enables amplifiers. The program is written in C language and in addition of controlling Qorvo it can work as SCPI server for Matlab client on PC. Maltab also evaluates received signal from RFSoC and then based on its distortion, adjusts the parameters of the transmitted signal.

ACKNOWLEDGMENT

I would like to express thanks to my supervisor Jan Král for his guidance. And to Josef Vychodil for his support with RFSoC.

REFERENCES

- Anristu: Intermodulation Distortion (IMD) Measurements: Using the 37300 Series Vector Network Analyzer. 9 2000, Available online: <u>https://reld.phys.strath.ac.uk/local/manuals/Anritsu37xxxVNAintermod.pdf.</u> [cit. 2024.3.5].
- [2] ROWE, M.: How DPD improves power amplifier efficiency. web, ANALOG IC TIPS, 1 2022, Available online: <u>https://www.analogictips.com/how-dpd-improves-power-amplifier-efficiency/</u>, [cit. 2024.3.9].
- [3] AVNET: Qorvo 2-Channel RF Front-end 1.8 GHz Card Hardware User's Guide. AVNET, 2020, Available online: <u>Qorvo 1800 MHz Card</u> <u>Hardware UG (avnet.com)</u> [cit. 2024.3.5].
- [4] KRÁL, J.: Digital Predistorters with Low-Complexity Adaptation. Doctoral thesis, Brno University of Technology, Faculty of Electrical Engineering and Communication, Department of Radio Electronics, 2022, Available online: <u>vut.cz/www_base/zav_prace_soubor_verejne.php?file_id=239465</u> [cit. 2024.3.9].
- [5] GOTTHANS, T., Baudoin, G., MBAYE, A., Comparison of modeling techniques for power amplifiers," 2013 23rd International Conference Radioelektronika (RADIOELEKTRONIKA), Pardubice, Czech Republic, 2013, pp. 232-235, doi: 10.1109/RadioElek.2013.6530922, Available online: <u>https://ieeexplore.ieee.org/document/6530922?denied=</u> [cit. 2024.3.9].

Utilizing Dynamic Analysis for Web Application Penetration Testing

1st Patrik Pis Department of Telecommunications FEEC, Brno University of Technology Brno, Czechia xpispa00@vut.cz

Abstract—This paper presents the design and implementation of a new modular tool, called PtWebDA, for dynamic analysis of web applications as one of the techniques used in penetration testing. Compared to other available tools and their limitations, our solution enables efficient rate limiting while also allowing testing of HTTP headers, cookie attributes, and content security policy directives. To verify its effectiveness in supporting manual web application penetration testing, we performed experimental testing in a controlled environment. The results of testing the presented tool PtWebDA are discussed in detail and highlight the key contributions of our solution.

Index Terms—cybersecurity, dynamic analysis, penetration testing, rate limiting, cookies, CSP directives, HTTP headers

I. INTRODUCTION

Information security practice, despite a gradual increase in awareness, is still lagging behind the rapid development of complex applications and systems. Security testing is still not a standard procedure during application development. The gradual increase of complexity of applications results in increasing number of opportunities for new errors and vulnerabilities. Most modern companies present themselves to the public using websites, and internal business processes and their management are gradually moving to the cloud in the style of complex web applications and services. The focus on secure development and operation of web services is thus becoming increasingly desired. With threats on the rise, testing web application security is becoming essential for secure operation of web services and reduction of potential negative impact on customers, employees, or the entire company.

Web applications have become a popular way for sharing content with users or providing services in various industries such as online banking, e-commerce, social networks, and many others. Web applications make it easier for users to access services without the need to install proprietary software, and thus become the easiest way for companies to provide content or service to a wide range of customers. However, due to their accessibility to the public, web applications are not only exposed to potential customers, but also to potential attackers who can cause their unavailability. For this reason, the overall security of web applications is critically important.

The research described in this paper was financially supported by the Ministry of the Interior of the Czech Republic, project No. VK01030019.

2nd Willi Lazarov Department of Telecommunications FEEC, Brno University of Technology Brno, Czechia lazarov@vut.cz

The impact of cyberattacks on web applications can vary widely and is largely influenced by the technologies used and the nature of the web application itself. However, if vulnerabilities found in a web application are exploited, not only will be the confidentiality, integrity, or availability of the web application operator's assets and infrastructure compromised, but sensitive data belonging to users may also be leaked, which can have an even deeper impact [1].

The main contribution of this paper is the design and development of a new tool PtWebDA for dynamic analysis of web applications that aims to support manual penetration testing. PtWebDA can test rate limiting, HTTP headers, cookie attributes, and Content Security Policy (CSP) directives.

The paper is divided into 5 sections. Section II discusses background and related work. Section III focuses on the description of the proposed dynamic analysis tool, including its key features and implementation details. Section IV presents the experimental testing of the tool developed, including a discussion of the results. The last Section V highlights the presented solution and indicates possible future work.

II. BACKGROUND AND RELATED WORK

Web applications are constantly evolving, and it is extremely important to pay adequate attention to the study of their security, as the increase in the complexity of web services can lead to new vulnerabilities, new tactics of threat actors, or new methods of circumventing the already existing defense mechanisms of web applications. To assess the security of a web application, it is crucial to point out the most critical vulnerabilities that are present in modern web applications. The OWASP Top 10 is a document of 10 categories that lists the most common vulnerabilities present in web applications. This document is put together and maintained by an organization called the Open Web Application Security Project (OWASP) [2]. OWASP is a non-profit organization dedicated to improving the security of web applications. The OWASP Top 10 is updated regularly to reflect current trends and threats in application security. It includes various types of vulnerabilities (e.g., broken access control) or security misconfigurations, from authentication and authorization flaws to vulnerabilities caused by incorrect processing of user input and exploitation of publicly known vulnerabilities.

One of the methods that we can utilize to assess the security of web applications is dynamic analysis. Dynamic analysis can be understood as testing or evaluating the functionality or security of an application during its runtime. It can be utilized in a wide array of scenarios, as it is non-exclusive to security testing. Such scenarios include software quality testing [3] and malware analysis [4]. This type of analysis also requires active execution of the tested application. During dynamic analysis, the individual tools do not have access to the source code as in static analysis, which does not require the web application to run, and its process and techniques are thus different from dynamic analysis. These tools simulate end users and have the same access to the application's sources as its potential users and their devices [5], [6].

Dynamic analysis of web applications can be performed manually or automatically. Manual and automated analysis of web applications is applicable in many test scenarios, but the choice of the right approach is strongly dependent on the nature and specification of the test. However, both approaches require a black-box approach without any insight into the internal structure of the web application. Dynamic analysis and penetration testing are two very closely related concepts. Principles of black-box penetration testing of web applications imply that dynamic analysis is a kind of building block of the whole security assessment. To perform a black-box penetration test, a given web application must be running, either in a test or production environment and the tester examines its reactions to various stimuli in real time [7].

To perform dynamic analysis, there already exist publicly available open source or paid commercial tools. Although the terms dynamic analysis, dynamic application security testing (DAST), and vulnerability scanning are different in definition, in practice these terms are often used interchangeably. In the context of the presented solution, we therefore focus on tools for dynamic analysis of web applications. The comparison of selected tools is shown in Table I.

 TABLE I

 Comparison of selected tools for dynamic analysis

Tool	Fuzzing	Manual	Automated	Open source
Burp Suite [8]	1	1	✓ ✓	×
OWASP ZAP [9]	1	1	1	1
ffuf [10]	1	1	X	1
Invicti [11]	×	X	1	X
Acunetix [12]	1	X	1	X
w3af [13]	1	X	1	1
Nikto [14]	X	X	1	1
AppSpider [15]	X	X	1	X
WebInspect [16]	X	X	✓	X

III. IMPLEMENTATION OF TOOL FOR DYNAMIC ANALYSIS

As shown in Table I, many publicly available tools lack support for a manual approach to testing. The biggest problem associated with automated dynamic analysis is the fact that in the current state-of-the-art, it is not capable of detecting all vulnerabilities. Certain more complex vulnerabilities, or chains of vulnerabilities, still require human thinking and creativity. Although automated dynamic analysis is a powerful tool for security testing, the manual approach is still more successful in finding vulnerabilities and eliminating false positives. The tool we are proposing in this paper should serve as the best of both worlds. The tool still automates the important tasks, but the tester has full access to the process and can adjust the testing parameters based on the observed behavior.

Compared to other tools, PtWebDA only tests for a limited set of security issues. This is done by design and should be considered only as an addition to a penetration tester's toolkit during a manual penetration testing engagement, not as a standalone automated security testing solution. The highlevel diagram of PtWebDA is shown in Figure 1.



Fig. 1. High-level diagram of PtWebDA

To show the difference and key contributions of our solution, we compared PtWebDA with already mentioned tools, focusing mainly on implemented modules for dynamic analysis. The results of the comparison are listed in Table II, where the rate-limiting test does not include any of the tools compared. This result was expected given that most automated tools care about covering as much as possible predefined tests and not the operation security¹.

¹Secure testing, i.e. limiting the speed of tools or invasive tests in order to avoid detection of the activity in question.

Tool	Rate limit	HTTP headers	Cookies	CSP
Burp Suite	X	✓ ✓	1	X
OWASP ZAP	×	✓	1	1
ffuf	×	✓	1	X
Netsparker	×	✓	✓ ✓	1
Acunetix	×	✓	✓ ✓	1
w3af	×	✓	 ✓ 	X
Nikto	×	✓	1	X
AppSpider	×	✓	 ✓ 	1
WebInspect	×	✓	1	X
PtWebDA	1	✓ ✓	1	1

 TABLE II

 Comparison of tools functions with the proposed tool

The tool was developed in Python programming language and, as of the time of writing this paper, consists of four modules representing the four aforementioned tests; HTTP headers (general), cookie attributes, CSP directives, and rate limiting, which are described in the following sections.

A. HTTP Headers

Testing for the presence of HTTP headers involves capturing the HTTP response to a request to a web server. Using the Python requests library [17], the test executes an HTTP request with a GET method to the specified URL that was passed to the test as a parameter. While this test is inherently very simple and easy to implement, the power of the findings it can provide can have a significant impact on the overall security of a web application or web server. Additionally, HTTP headers can often be missed during manual inspection of responses, and the tool helps to highlight potential issues to the tester.

B. Cookie Attributes

The principle behind cookie attributes and CSP directives testing is practically the same as in HTTP headers involving capturing the HTTP responses to requests to a web server. However, this time the tool focuses more on the contents of these HTTP headers than on their presence alone.

The cookie attributes module in the current version of PtWebDA focuses mainly on these cookie attributes: *Secure*, *HttpOnly*, and *SameSite*. These attributes play a crucial part in lowering the potential impact of vulnerabilities such as cross-site scripting (XSS) or cross-site request forgery (CSRF), since they directly contribute to the definition of how and when entities may access the cookies themselves. The module checks these attributes, analyzes them, and informs the tester of possible misconfigurations.

C. CSP Directives

Content Security Policy (CSP) defines the rules on how a web browser should load and execute content provided by the web application. It requires a careful setup and precise definition. If CSP is defined, it has a significant impact on the way web browsers render web pages. CSP helps detect and prevent a wide range of attacks, including cross-site scripting, other cross-site attacks, and other attack vectors that lead to compromise users' privacy. This module analyzes the CSP directives defined in the *Content-Security-Policy* HTTP header and searches for insecure definitions and misconfigurations. The important thing during CSP analysis is that there are multiple ways to execute client-side JavaScript, not just direct code definitions (e.g. in <script> tags). JavaScript code can be executed from CSS style definitions, which are often neglected during CSP configuration and may lead to data exfiltration. This module focuses on informing the tester about these possible misconfigurations and highlights possible dangerous CSP directives.

D. Rate Limiting

Testing rate limiting in dynamic analysis involves building a model to evaluate the response of a web application to a significant number of requests in a specified amount of time. However, the actual fact that a web application or a web server has a mechanism for limiting user requests may not be immediately obvious. It is a significant challenge for a tool evaluating rate limiting to differentiate between the slowdown or complete termination of requests by the actual rate limiting mechanism and the insufficient computational or memory resources of the tested system.

Mathematically, we can describe the rate limiting with the following equation:

$$\int_{t_1}^{t_2} R(t) \, dt \le L \times (t_2 - t_1), \tag{1}$$

where

R(t) represents the rate of HTTP responses in time t, L represents the maximal number of requests per minute, t_1 , t_2 represent the time interval.

Analysis of *best practices* from companies such as Cloud-Fare [18] shows that relatively low values of the maximum number of requests per specified time are used for the effective application of rate limiting. Thus, an effective defense against attacks on web application login forms can be a rate limit of, for example, 4–10 requests per minute. PtWebDA performs an analysis of this HTTP traffic and determines the number of requests necessary to trigger the rate limiting mechanism. The detection of rate limiting and the number of failed and successful requests provide sufficient information to solve the inequality defined in Eq. 1 and thus allow us to compute the theoretical value of the maximum number of requests per minute for a given time interval.

From the penetration tester's perspective, this provides important information on the target's security status and shows how testers can alter her/his approach by trying to circumvent the rate limit, e.g. by limiting the number of requests just below the detection threshold, allowing the tester to continue with testing uninterrupted. In addition, the tool does not perform automatic scans without the tester's knowledge, which makes it suitable for manual penetration testing, where the tester has the entire tool under control.

IV. EXPERIMENTAL TESTING

As described in section III, the tool consists of four working modules in its first version. We conducted a series of tests to prove the effectiveness and reliability of each module before testing it out in the production environment. To test it, we created a controlled environment using virtual machines. The virtual machine was running on the Ubuntu 22.04.3 LTS operating system and used Apache 2.4.58 as its web server. To create a controlled sandbox environment, we created a misconfigured web application written in Python using a Flask web framework, containing all the flaws the tool is expected to successfully detect. The web application consisted of a total of three endpoints; "/login", "/menu", and "/nolimit".

The login endpoint had a rate limit configured as low as 4 requests per minute, the menu endpoint had a rate limit configured as 30 requests per minute, and the final endpoint had no rate limit configured. At the same time, the different endpoints of the web application returned the content with different set of HTTP headers and their contents. For example, the login endpoint did not return the *Content-Security-Policy* header, but the menu endpoint did. A summary of results of the experimental testing conducted in our controlled environment is shown in Table III.

TABLE III Results of experimental testing

Tested module	Findings (found/total)	Duration
HTTP headers	8/8	< 1s
Cookie attributes	3/3	< 1s
CSP directives	4/4	< 1s
Rate limit	1/1	18.46s

The main goal of the experimental testing was to ensure that the proposed tool was capable of detecting and effectively highlighting important insights in web server and web application configurations and to try to identify potential bottlenecks and areas for future improvement.

The controlled environment provided us with a strong position to interpret the output of individual module's test runs and compare them with the intended misconfigurations. Since all modules, except rate limiting, are very basic in their core principles, the results of the experimental testing were very accurate. The tool was effectively able to identify all missing headers and provided output regarding headers that reveal information. Regarding the rate limiting module, the results were accurate as well with a few exceptions. The rate limiting module uses mutlithreaded approach and can estimate the web application's rate limit settings with an accuracy of approx. 5 requests margin. The experimental testing of rate limiting provided us with insights on how to handle some edge cases, which can alter the tool's output and produce false information. Some of the edge cases include; the tool is too fast and the rate limit is quite small, or the rate limit is very permissive and the tool is not fast enough to trigger the rate limit, or the tool does not send enough requests to trigger the rate limit defined by the tested application.

V. CONCLUSION AND FUTURE WORK

To summarize our paper, we reviewed the current state of dynamic analysis with a primary focus on web application penetration testing, including a comparison of available solutions. Based on the limitations of existing tools, we designed and developed a new modular tool, PtWebDA, to test the rate limiting, HTTP headers, cookie attributes, and CSP directives of web applications. The tool is developed to be fully under the control of the tester, which makes it particularly suitable for manual penetration testing.

Based on the experimental testing, the tool can be improved in its efficiency and reliability. Additionally, dynamic analysis of web applications is not limited to rate limiting, HTTP headers, cookie attributes, and CSP directives. The area of dynamic analysis is wide, and the modular design of the PtWebDA allows us to develop it even further and transform it into a more complex and versatile tool. Future work and development of this tool will include the development of more modules covering more testing scenarios included in the OWASP Top 10 methodology.

REFERENCES

- M. A. Kunda and I. Alsmadi. "Practical web security testing: Evolution of web application modules and open source testing tools," in 2022 International Conference on Intelligent Data Science Technologies and Applications (IDSTA), 2022, pp. 152-155.
- [2] "OWASP Top Ten," OWASP. https://owasp.org/www-project-top-ten/ (accessed Feb. 10, 2024).
- [3] G. J. Myers, T. Badgett, and C. Sandler, *The Art of Software Testing*. John Wiley & Sons, 2011.
- [4] A. Afianian, S. Niksefat, B. Sadeghiyan and D. Baptiste. "Malware Dynamic Analysis Evasion Techniques: A Survey," ACM Computing Surveys (CSUR), vol. 52, no. 6, pp. 1-28, Nov. 2019.
- [5] R. Baloch. *Ethical hacking and penetration testing guide*. CRC Press: Taylor & Francis Group, 2014.
- [6] T. A. Nidecki. "Vulnerability Assessment and Penetration Testing of Web Application," https://www.acunetix.com/blog/web-security-zone/ dynamic-static-code-analysis-web-security/ (accessed Feb. 17, 2024).
- "Dynamic Application Security Testing (DAST)," PortSwigger. https: //portswigger.net/burp/application-security-testing/dast (accessed Feb. 17, 2024).
- [8] "Burp Suite Community Edition," portswigger.net. https://portswigger. net/burp/communitydownload (accessed Feb. 17, 2024).
- [9] "Zed Attack proxy (ZAP)," zaproxy.org. https://www.zaproxy.org/ (accessed Feb. 17, 2024).
- [10] "Fast web fuzzer written in Go," github.com. https://github.com/ffuf/ffuf (accessed Feb. 17, 2024).
- [11] "Application security with zero noise," invicti.com. https://www.invicti. com/product/ (accessed Feb. 17, 2024).
- [12] "Web Application Security Scanner," acunetix.com. https://www.acunet ix.com/product/ (accessed Feb. 17, 2024).
- [13] "Introduction," docs.w3af.org. https://docs.w3af.org/en/latest/phases.h tml (accessed Feb. 17, 2024).
- [14] "Nikto," github.com. https://github.com/sullo/nikto (accessed Feb. 17, 2024).
- [15] "Welcome to AppSpider," docs.rapid7.com https://docs.rapid7.com/ap pspider/ (accessed Feb. 17, 2024).
- [16] "OpenText Fortify WebInspect," opentext.com. https://www.opentext.c om/products/fortify-webinspect (accessed Feb. 17, 2024).
- [17] "Requests 2.31.0," pypi.org. https://pypi.org/project/requests/ (accessed Feb. 28, 2024).
- [18] "Rate limiting best practices," developers.cloudflare.com. https://develo pers.cloudflare.com/waf/rate-limiting-rules/best-practices/ (accessed Feb. 28, 2024).

Navigation of UAV in GNSS denied area

1st Marco Pintér Department of Control and Instrumentation Brno University of Technology Brno, Czech Republic xpinte06@vutbr.cz 2nd Petr Marcoň

Dept. of Theoretical and Experimental Electrical Engineering Brno University of Technology Brno, Czech republic marcon@vut.cz

Abstract—This paper examines the concept of navigation of of Unmanned aerial vehicle (UAV) in three-dimensional space using visual odometry. In the near future navigation of the UAV without GNSS is becoming a critical part of autonomous navigation systems, using information from on-board cameras to estimate the UAV's movement and position. In the paper, different types of visual odometry, sensors for visual odometry, components of the implementation, and scenarios of usage. For the development and future application we utilize widely used Robotic Operating System (ROS).

Index Terms-drone, visual-odometry, automation, navigation

I. INTRODUCTION

This study presents the development of UAV navigation system using visual data in GNSS-denied areas. In more detail, we are focusing on the use of visual odometry (VO) for UAV position estimation. Today, it is becoming key part of the navigational system of autonomous cars, UAVs, and robots. VO is a technique that compares the currently captured image with the previous one, seeking differences in the displacement of selected features or reference points using the optical flow method. The new pose is obtained by adding the estimated position vector relative to the previous pose.

For the VO position estimation, we utilize Intel Realsense depth camera. We are currently employing the RTAB-Map algorithm for VO position estimation. As a companion computer for the control of UAV and image processing, we utilize the Intel NUC. We also evaluated the difference in accuracy between VO and GPS. For the flying during the development process, we utilize Gazebo simulation alongside PX4-SITL.

II. VISUAL ODOMETRY

Visual odometry is the process of position estimation and camera motion from frame sequence. VO methods can by divided into Relative Visual Odometry (RVL), and Absolute Visual Odometry (AVL).

A. Relative visual odometry

This method employs optical flow techniques to estimate position changes by detecting edges, corners, and changes in pixel brightness. Feature matching between frames identifies key points, from which position changes are calculated and camera pose is updated. However, the method faces drift over time due to errors in recursive pose calculations, which accumulate over time. Solutions include Loop Closure Detection, Simultaneous Localization and Mapping (SLAM), fusion with Inertial Measurement Unit (IMU), or reinitialization to mitigate drift.

The main advantage of this method is that it does not require knowledge of the environment or some additional data about the environment for functionality. It is also suitable for use in indoor applications as it is independent of the Global Navigation Satellite System (GNSS). [1] [2]

B. Absolute visual odometry

This method fundamentally differs from RVL in its approach to feature detection and assignment. It utilizes precollected reference data, assuming accurate georeferencing prior to use in localization. These data can include aerial images captured and georeferenced using UAV GNSS or sourced from platforms like Google Earth.

An advantage over RVL is its immunity to error accumulation due to its use of georeferenced data. However, its effectiveness relies on the quality of the dataset it interacts with. Freely available data may lack uniform lighting conditions or fail to account for dynamic terrain changes, such as vehicle movement on roads, leading to significant errors.

Another key difference compared to RVL is that AVL does not operate with consecutive frames, as it compares the captured image during localization with a database and seeks matches within the image. One option is the use of template matching, where the captured image is directly compared with a known database, and the similarity of the images is evaluated. Another possibility is to employ feature matching, similar to RVL, but with the distinction that the search for distinctive points in the captured image is compared with the points in the dataset.

III. INERTIAL ODOMETRY

The method utilizes motion sensors like accelerometers and gyroscopes to determine the position, orientation, and speed of a moving object without external reference, known as dead reckoning. All inertial navigation systems encounter an issue known as integration error. The problem arises because even the slightest error in measuring acceleration or speed will gradually accumulate over time. Even the best accelerometers on the market, on average, accumulate a deviation of around 50 meters over 17 minutes. This error also manifests itself in situations where the device remains stationary, as the system will, after a certain period, declare a change relative to the initial pose.

Inertial odometry is primarily designed for measuring short distances. For measuring longer distances, it is necessary to combine this technology with other methods of position determination, such as GPS, Light Detection and Ranging (LiDAR), or another source of location tracking. When measuring greater distances, is crucial to eliminate the integration error of this method by resetting the position from a different source. This combination allows for a system capable of determining position with higher precision. [4]

The primary limitation of the VIO approach arises during rapid changes in position where there is a lack of continuity between consecutive frames. This leads to a loss of the current position estimation, requiring VIO to be restarted.

IV. ALGORITHMS FOR IMAGE PROCESSING

Nowadays, many visual odometry algorithms are available, such as ORB-SLAM, OpenVINS, RTAB-Map and many more. In the experiment described in this paper, RTAB-Map algorithm will be utilized.

A. RTAB-Map

RTAB-Map (Real-Time Appearance-Based Mapping) is a SLAM algorithm for trajectory acquisition, which is based on processing RGB-D, stereo, or lidar data. This algorithm relies on an incremental image matching detector with the detection of previous features. This principle is practically based on relative visual odometry, with the difference being the detection of matches with previous frames. This way, the problem of accumulating position estimation errors over time is eliminated. Loop closure detection utilizes the bag-ofwords (BoW) method to determine the probability of a match between the new image and the previous location, or if it represents a new location. In the case of loop closure detection, which involves recognizing patterns in visual elements to identify previous locations, the algorithm retroactively corrects accumulated trajectory errors. To perform this detection, a larger number of previous frames need to be stored. Without limiting the number of stored frames, the computational complexity would increase over time, as it would be necessary to process a larger amount of data. To address this issue, RTAB-Map also implements memory management, which restricts the number of locations designated for loop closure detection. This ensures that the computational performance remains sufficient for realtime comparison. Additionally, RTAB-Map allows for fusion with inertial odometry. Another advantage is its compatibility with ROS 1 and ROS 2.

RTAB-Map implements both AVL and RVL techniques in its visual odometry process. While AVL focuses on determining the camera's absolute position by referencing maps or landmarks using loop closure detection, RVL estimates its position relative to previous locations based on incremental visual changes. By integrating both AVL and RVL approaches, RTAB-Map enhances its ability to achieve robust localization and mapping results, particularly in dynamic environments. Additionally, RTAB-Map also utilizes Inertial Odometry to further improve localization accuracy and robustness. [5] [?] [3]

V. HARDWARE SELECTION

In the following chapter, suitable companion computer platforms and cameras for processing visual and inertial data on UAV are described. For the practical experiments, we utilize Holybro X500 airframe.

A. Companion Computer

When making a selection, it's essential to consider appropriate dimensions, weight, and sufficient computational power. While platforms like Raspberry Pi or NVIDIA Jetson Nano are commonly employed for most current image processing applications due to their balance of affordability, suitable dimensions, and adequate computational power, we opted for the Intel NUC (Next Unit of Computing) to achieve the maximum frame rate of visual odometry.

B. Camera

Intel provides a variety of stereo depth cameras, each differing in sensor type, supported resolution, and frame rate, with some models including an integrated IMU unit. A notable advantage is the support for ROS, as Intel offers a ROS wrapper for processing both visual and inertial data. From the Intel RealSense family, the most well-known depth camera series is the D400 series, which utilizes infrared projectors and sensors to capture depth in the image, enabling 3D mapping of space. Additionally, these cameras feature a high-resolution color sensor and a wide field of view, enabling the capture of a large area in a single frame. For our experiments, we decided to use the Intel RealSense D455 depth camera.

In the figure 1, the connection between the UAV and the Intel NUC companion computer is shown. For the connection of the RealSense camera to the companion PC, we used a high-speed USB 3.0 bus to achieve the fastest data transfer. The flight controller is connected to a computer via a USB TTL converter.

VI. UAV CONTROL

For communication with the UAV, a ROS node for MAVLink communication was implemented. The FC with PX4 firmware allows control of the UAV in different flight modes such as HOLD, MISSION, ALTCTL, POSCTL, MAN-UAL, STABILIZED, and more. Most of the mentioned flight modes utilize GNSS positions for position control and stabilization in space. For commanding the UAV without GNSS, only OFFBOARD mode allows control of the UAV from an external companion computer without using GNSS.

The main goal is to implement the functionality of returning to the launch position when GNSS positioning is interrupted or unavailable. There are two scenarios for the return to launch position. The first involves a direct approach to the launch position, which means computing the appropriate yaw



Fig. 1. Hardware connection between UAV and companion computer

direction towards the launch position. The second approach utilizes the previous flight path to return to the launch position. Both scenarios are shown in the figure 2. Subsequently, the



Fig. 2. Scenarios for return to launch position

mechanism for controlling the UAV to return to the launch position can also be used for flying to a desired relative position in space.

VII. SOFTWARE IMPLEMENTATION

For the fast development, we utilized widely used ROS (Roboting operating system), Humble version and Gazebo simulation alongside PX4-SITL. We decided to use Ubuntu 22.04 as the operating system for the companion computer due to its compatibility with the chosen version of ROS. For visual and inertial data from Realsense camera we used ROS wrapper developed by Intel. From the RTABmap library we used appropriate nodes for processing visual and inertial data from Realsense camera. We utilized Rviz for data visualization

purposes. We also employed RTAB-Map Viz to visualize the process of visual odometry. To compute the relative position from the WGS84 GPS format, we implemented custom nodes in our system.

For communication with the flight controller (FC), we utilized the MAVSDK library for MAVLink commands communication and telemetry processing from the UAV. A ROS node **uav_control** for commanding the UAV in offboard mode and telemetry communication was implemented. In the figure 3 the connection of topics and nodes of the system is shown.



Fig. 3. Software infrastructure of ROS nodes and topics connection

VIII. SIMULATION OF UAV

We utilized the Gazebo simulation alongside PX4-SITL for testing the **uav_control** node. The simulation allows us to verify proper communication, control, and telemetry communication with the FC. In the figure 4, the UAV during flight in Gazebo simulation is shown.

IX. EXPERIMENT OF UAV POSITION ESTIMATION

The experiment was done using created ROS nodes, during which all data on the UAV were recorded. A rosbag was



Fig. 4. Drone in Gazebo simulation environment

recorded during the experiment, which allows for subsequent reconstruction of the flight. Using this rosbag, it is possible to develop various other applications with UAV data, but also to evaluate the measured data. During our experiment, we focused on comparing the precision of positioning between visual odometry and GPS.

After launching the UAV, it was necessary to wait for a sufficient number of GPS satellites to determine the GPS position. Subsequently, all nodes for communication with the flight controller, RealSense camera, and visual-inertial odometry needed to be started. After launching all the necessary nodes, it was possible to start recording all the published topics into the rosbag. After the initialization of the system, the experiment flight was taken. During the experiment, the GPS system had access to approximately 8 satellites, and RTK (Real Time Kinematic) GPS was not utilized. The data stored in the rosbag was visualized afterward, allowing for comprehensive analysis and evaluation of the recorded information. The experimental flight was executed at an altitude of approximately 2 meters.

A. Experiment evaluation

The measured data was processed using the created ROS nodes for error evaluation. The resulting trajectory from visual-odometry and the GPS trajetory can be visualized using the rviz environment. In figure 5, the visualization of both trajectories in the rviz environment is shown. The blue trajectory represents the GPS trajectory, and the green one represents the visual odometry trajectory

As we can see in the trajectory image, they have approximately the same shape with certain deviations. During the experiment, a straight flight over the sidewalk was taken, followed by a return to the starting point. The output data from the ROS node was processed in the spreadsheet editor. From the measured data we can say that the biggest deviation was recorded in y-axis. The average deviation in the x-axis was 0.294 m, in the y-axis was -0.431 m, and in the z-axis was 0.357 m. In the evaluation of the data, no consideration was given to the displacement of the camera and GPS module. This fact may lead to distortion in the y-axis deviation. As the average accuracy of GPS typically ranges from 2 to 5 meters, our experiment demonstrates that the deviation between visual odometry and GPS was minimal.



Fig. 5. Visualization of both trajectories in the rviz environment

X. CONCLUSION

Based on the experiment, we can conclude that visualinertial odometry achieves very good results compared to GPS. From the data, we can see that visual odometry exhibits less oscillations in position estimation compared to the GPS used. The maximum absolute deviation was observed in the Y-axis, with a value of 3.474 m, while in the X-axis, the maximum absolute deviation was 1.895m. Experiments have proven that visual odometry can replace GPS for monitoring the position of UAVs. The main advantage of visual odometry over GPS is its independence from satellites or other signal sources for positioning. This fact is crucial in most current applications, as the GPS signal can be interfered with or may not be available at all.

In the future, it would be interesting to measure the precision of both technologies against reference points in the terrain, as our current experiment evaluates the accuracy of visual odometry against GPS.

Based on the previous experiment, we implemented the architecture for communication with the FC, along with a ROS node for commanding the UAV based on estimated position with VO. The ROS node created for UAV commanding was tested in the Gazebo simulation. In the future, it will be essential to make tests on a real drone.

REFERENCES

- Y. Bai, B. Zhang, N. Xu, J. Zhou, J. Shi, and Z. Diao, "Visionbased navigation and guidance for agricultural autonomous vehicles and robots, *Computers and Electronics in Agriculture*, vol. 205, 2023.
- [2] A. Couturier and M. A. Akhloufi, "A review on absolute visual localization for UAV", *Robotics and Autonomous Systems*, vol. 135, 2021.
- [3] M. Labbé and F. Michaud, "RTAB-Map as an open-source lidar and visual simultaneous localization and mapping library for large-scale and long-term online operation", *Journal of Field Robotics*, vol. 36, no. 2, pp. 416-446, 2019.
- [4] R. Acharya, "Introduction to Navigation", in Understanding Satellite Navigation, Elsevier, 2014, pp. 1-26.
- [5] M. Labbe and F. Michaud, "Appearance-Based Loop Closure Detection for Online Large-Scale and Long-Term Operation", *IEEE Transactions* on Robotics, vol. 29, no. 3, pp. 734-745, 2013.

Educational PocketQube Satellite Demonstrator

Jiří Veverka Department of Radio Electronics FEEC, Brno University of Techonology Technická 12, Brno, Czech Republic E-mail: xvever12@vutbr.cz Aleš Povalač Department of Radio Electronics FEEC, Brno University of Techonology Technická 12, Brno, Czech Republic E-mail: povalac@vutbr.cz Kamil Jaššo Faculty of Military Technology University of Defense Brno, Czech Republic E-mail: kamil.jasso@unob.cz

Abstract—This paper presents the design of a 2P $(50 \times 50 \times 50 \text{mm})$ PocketQube satellite which will be used in STEM laboratory classes. The proposed system was designed with emphasis on modularity, ease of replacement using commercial off-the-shelf components, and easy understanding of the subsystems one can find in a satellite of this type. It also encourages students to come up with their own ideas and modules. The satellite itself is made out of four main modules: On-Board Computer, Communication module, Electric Power System module, and Payload module. These four modules are encapsulated in a structure made entirely out of PCBs and together form the basis of the satellite.

Index Terms—On-board Computer (OBC), Electrical Power System (EPS), Communication system (COM), PocketQube, PCB, STM32, LoRa

I. INTRODUCTION

The PocketQube is a type of microsatellite with a base unit of 1P ($50 \times 50 \times 50$ mm), and its structure is standardized by Alba Orbital, Delft University of Technology, and GAUSS Srl [1]. The aim of this paper is to describe the design of a 2P-sized PocketQube satellite for STEM educational purposes in a laboratory setup. In the classes, students will be encouraged to contribute by designing their own modules compatible with the proposed system bus pin header and satellite architecture. Through this approach, students will gain invaluable practical experience in firmware, electrical and printed circuit board (PCB) design, problem troubleshooting, functional analysis, and system engineering.

The satellite's architecture is designed according to the PocketQube standard [1] and its structure is purely made out of PCBs. To ensure easy satellite modularity, the design incorporates easily replaceable components through standard commercial off-the-shelf (COTS) parts and is partially influenced by the design of already existing 1U $(100 \times 100 \times 100 \text{mm})$ CubeSat microsatellite [2] in the student's laboratory to ensure some level of compatibility between those two systems.

II. POCKETQUBE SYSTEMS

Just like CubeSats [3], PocketQubes are composed of several subsystems that, when integrated into one unit, ensure satellite operations. Among the most commonly used subsystems in PocketQube satellites are:

- On-board Computer (OBC) Satellite's central processing unit. The OBC takes care of mission and telemetry data handling and utilizing available interfaces (I²C, SPI), executes commands, and operates the mission autonomously.
- Electrical Power System (EPS) This system acquires electrical energy through photovoltaic cells and then stores it in an on-board battery/ies. It also regulates and provides that energy to other subsystems and the rest of the satellite through the system bus, ensuring the operational capacity of the entire satellite.
- Communication System (COM) This subsystem ensures communication between space and Earth, utilizing advanced transmission technologies and an antenna for receiving commands and sending satellite data back to Earth. The parameters of the COM, such as bandwidth, frequency, or transmission speed, are typically determined by mission requirements.
- Payload This is the main mission objective. It typically involves some scientific experiment, technological demonstrator, Earth observation, etc.
- Attitude Determination and Control System (ADCS) This subsystem is in charge of positioning the satellite based on its relative position to the stars or the Earth which is obtained from various input sensors all across the satellite. The positioning itself is achieved using e.g. reaction wheels (one wheel per axis) or magnetorquers locked to the Earth's magnetic field. This system is not considered to be implemented in this work.

A. OBC Design

The On-Board Computer (OBC) is implemented using a separate module, as depicted in Figure 1. It is driven by the STM32F44RET microcontroller with an external 8MHz clock, which is directly soldered onto the PCB module.

The module further includes a connector for the SWD interface for firmware uploading and debugging purposes, external EEPROM memory for data logging connected via the I^2C interface, a debug UART interface, and a system bus.

The system bus is a connector linking the OBC to the rest of the satellite and is standardized across all modules. Through the system bus, the OBC is capable of acquiring
data from the EPS or processing and transmitting data through the COM module. All other unused pins on the MCU are left unconnected. A general description of satellite software operations can be seen in Figure 2.

The OBC module mainly utilizes the SPI and I^2C interfaces (as depicted in Figure 3), with an extra I^2C_CSP interface, which is isolated from the first I^2C interface and is destined to be used in the future, specifically designated for communication over the CubeSat Space Protocol (CSP) [4].

For the proper functioning of the satellite, the OBC module also serves as the control unit for the EPS, monitoring the available battery capacity to determine whether to continue operations according to current settings or to suspend the satellite's activities and put itself into a sleep mode.



Fig. 1. 3D model of OBC.

B. COM Design

As the heart of the communication (COM) module was used the RFM98W RF Transceiver from HopeRF which has communication speed up to 300kbit/s, integrated CRC check, maximum output power of +20 dBm (100mW), and supports a wide range of modulations (FSK, GFSK, LoRa, OOK) [5]. This module operates at a selectable frequency in the range from 410 to 525MHz and uses the LoRa modulation which is employed for transmitting telemetry data from the satellite and receiving commands from the ground station. The transceiver module is directly soldered onto the PCB along with an SMA connector for an antenna connection. One of the main reasons why the transceiver is not integrated with the OBC into a single board is to ensure that the resulting system remains as modular as possible, allowing students in the future to experiment with different configurations and communication types. The entire COM module is connected via a system bus to the rest of the satellite utilizing the SPI interface for communication between the OBC and the RFM98W. Apart from the SPI interface, the



Fig. 2. Block diagram of the satellite's software.



Fig. 3. Block diagram of the OBC.

reset and interrupt pins are also connected, so the OBC can effectively control the module and receive an interrupt when incoming communication for processing is received.

C. EPS Design

The energy from 20 solar cells $(45 \times 15 \times 2.1 \text{mm})$ of maximum peak power 123mW per cell and cell efficiency of 25% typically [6] is stored in one Li-Ion 18650 battery with a capacity of 3350mAh. These solar cells are connected in 10p2s configuration with Schottky diodes on each parallel set for short circuit protection and soldered directly onto the satellite's construction PCBs (see chapter III) and their input is monitored from OBC via ADC lines connected to ADC



Fig. 4. Block diagram of the COM.

expander. This is to ensure that the energy keeps flowing into the system even if not all of the solar cells are illuminated. The battery itself is protected against a short circuit with 1.5A recoverable fuses and the input and output current is measured by two INA219 modules on its ends. These two modules are connected via the I²C interface to the OBC which then regulates the power consumption itself. To turn the satellite on, a kill switch is connected to the main power rail with a P-channel MOSFET that disconnects the power rail if the switch is in the off position. The EPS main power rail provides only 3.3V of voltage power since the satellite has no needs (and space) for anything else. The block diagram of the EPS can be seen in Figure 5.



Fig. 5. Block diagram of the EPS.

III. SATELLITE CONSTRUCTION

The proposed PocketQube consists of two parts – an external 2P-sized structure ($114 \times 50 \times 50$ mm without the backplate) and a so-called "sliding backplate" ($58 \times 128 \times 1.6$ mm) which forms one of the sides and serves for securing the satellite into the deployer (as can be seen in Figure 6). The entire structure is implemented using six PCBs with a thickness of 1.6mm, which are assembled by soldering and screws. Solar cells that are directly connected to the EPS are also soldered on these PCBs. Thanks to this solution, more volume and weight budgets are available in the final construction for the internal modules. In the PocketQube standard, we are limited



Fig. 6. 3D model of the structure with example layout of inside components.

to a maximum mass of 250g/1P, however, this fact was not taken into account, as it is only an educational model.



Fig. 7. Block diagram of the Ground station and PC assembly.

IV. GROUND STATION DESIGN

The Ground station was designed similarly to the communication module for the satellite – utilizing a COTS RFM98W module operating at 433.92MHz frequency and an SMA connector for an antenna. The module is designed as a shield with Arduino Uno V3 connector spacing which is then connected to the STMicroelectronics Nucleo-F030R8 development board (as can be seen in Figure 7). This configuration is then connected via an USB connector to a computer, thus enabling communication with the satellite for data reception and command transmission. The ground station communicates with the satellite via a wireless, LoRa-modulated connection. The processes of data reception and command transmission are managed in an application executed on the computer to which the ground station is connected.

V. CONCLUSION

In this paper, we described the design of a 2P PocketQube for educational purposes with emphasis on modularity, ease of understanding, replaceability, and hands-on learning experience. The outcomes were finalized schematics, diagrams, PCB designs, and structural designs ready to be manufactured. In this configuration, the OBC is the sole computing unit that controls the COM and EPS modules although that can be changed in the future, thanks to the integrated support of the CubeSat Space Protocol.

The designed satellite is completely modifiable and serves as the basis for future student projects, laboratory tasks, or entirely new modules which makes it an ideal educational tool.

ACKNOWLEDGMENT

Research described in this paper was supported by the Internal Grant Agency of the Brno University of Technology under project no. FEKT-S-23-8191 and by the institutional support of the Ministry of Defence of the Czech Republic (VAROPS).

REFERENCES

- [1] The PocketQube Standard [online]. [cit. 2024-03-10]. Available at: http://www.albaorbital.com/pocketqube-standard
- [2] KOŠÚT, Martin. Návrh a realizace výukového CubeSatu. Brno, 2022. Dostupné také z: https://www.vut.cz/studenti/zav-prace/detail/141541. Diplomová práce. Vysoké učení technické v Brně, Fakulta elektrotechniky a komunikačních technologií, Ústav radioelektroniky.
- [3] CAPPELLETTI, Chantal, Simone BATTISTINI a Benjamin K. MALPHRUS. CubeSat handbook: from mission design to operations. San Diego, CA, United States: Academic Press is an imprint of Elsevier, [2021]. ISBN 978-0-12-817884-3.
- [4] libcsp. The Cubesat Space Protocol (Source code) [online]. [cit. 2024-05-04]. Available at: https://github.com/libcsp/libcsp/
- [5] HopeRF: RFM98W RF Transceiver Module (Datatsheet) [online]. [cit. 2024-03-10]. Available at: https://www.hoperf.com/modules/lora/RFM98W.html
- [6] SM141K04LV ANYSOLAR High Efficiency SolarMD (Datasheet) [online]. [cit. 2024-03-10]. Available at: https://ixapps.ixys.com/DataSheet/SM141K04LV.pdf

SUPER HOT CHIPS

A COMPANY OF THE SWATCH GROUP

Naše srdce planou pro navrhování čipů! Hledáme další nadšence, co by rozšířili naše řady!

***JOIN US!**

ASICENTRUM PRAHA / BRNO

WE LOVE TO DESIGN THEM.

ČIPY KTERÉ JINDE NENAJDETE

- Navrhujeme ULTRA LOW POWER integrované obvody
- Naše čipy mají vždy něco jedinečného, co konkurence neumí
- Bluetooth, RFID NFC, RFID UHF, Sensor hubs, Optical Sensors
- Na tom všem se můžete podílet

ČESKÁ FIRMA SE ŠVÝCARSKÝM ZÁZEMÍM

- ASICentrum je středně velká firma s rodinnou atmosférou
- Jsou zde špičkoví odborníci, od kterých se můžete hodně naučit
- Švýcarské zázemí díky mateřské firmě EM Microelectronic
- V centrále firmy EM Microelectronic ve Švýcarsku máme vlastní továrnu na výrobu integrovaných obvodů

PŘIJĎ NAVRHOVAT ČIPY

- Nabízíme brigády, diplomky a zaměstnání
- Potřebujeme návrháře, verifikátory i softwarové inženýry
- Studenty naučíme moderní techniky návrhu integrovaných obvodů
- Zaměstnanci získají obrovské zkušenosti s celým procesem vývoje čipů
- Skvělé platové ohodnocení a benefity
- Co očekáváme zájem, chuť do práce a nezbytné základní znalosti ze školy

HLEDÁ SE KOLEGA HOMO TECHNICUS

• Digital IC Designer

ARTIU

- Digital IC Verification Engineer
- Embedded Software Engineer
- Digital Backend Engineer
- Analog IC Designer





Electrical characterization of graphene sensors

1st Patrik Staroň Department of Physics Brno University of Technology Brno, Czech Republic patrikstaron@vutbr.cz

Abstract-Graphene has shown to have great electrical, thermal and chemical properties. These qualities are suitable for construction of a chemical sensor with graphene as its sensing material. These graphene-based sensors are currently in the field of research and they are not vet viable for mass commercial application. Studies in graphene-based sensors span its very high sensitivity, selectivity, functionalization, manufacturability, stability and material configuration. We focus on the noise characterization in pristine graphene sensors, with different active channel size and layout of the electrodes. The topic of sensor characterization is difficult. Very often researchers omit the measurement setup and information about sample preparation. This article is focused on design of a sample, that is suitable for various graphene sensor experiments. We have come up with a successful design of highly configurable sample board with mounted graphene sensor, that are easy to work with. These samples are to be used for noise analysis of a graphene sensor, that is put in different gas environments.

Index Terms—graphene, 2D sensor, measurement, sample preparation

I. INTRODUCTION

Multiple industries, such as manufacturing, healthcare and automotive, have a need for reliable measurement of chemical composition in their environment. These environments include places where people work with dangerous substances and the environments need to comply with safety requirements, or areas with fire hazard or low air circulation. Measurement and detection of volatile substances in the air can be achieved only by gas sensing devices. This demand is one of the reasons why gas sensors are a topic of research up to this day.

As the research in material science progresses, there are new methods and options that can be tried, which promise advancement in design of a sensor with better properties and manufacturability [5].

While there are many commercial implementations of gas sensors, they have problems with accuracy, measurement precision, selectivity to certain molecule, cost and so on. [6].

Gas sensors are a subcategory of chemical sensors. While chemical sensor is meant to be used in generally any environment where detection of chemicals can be performed, gas sensors are meant to be placed in a gaseous environment, although a gas sensor can also to some extent work as a general chemical sensor [5] [7]. Further, the chemical (or gas) sensor can be classified as one based on 2D sensing material. These materials include graphene and its variants, phosphorene, transition metal dichalcogenides, hexagonal boron nitride, Molybdenum Disulfide and others.

The article focuses on graphene 2D chemical/gas sensor. Graphene is a material with spectacular properties. It has high carrier mobility [9] and thermal conductivity, has high mechanical strength [8] and provides various optical and geometrical advantages over all other materials.

This is a premise that lead to the research of how well this material can be used in detecting various gas molecules. It has been found, that graphene might be a suitable material to detect various gases at extreme resolution [10].

In reality, as it is a common case, graphene has problems not only with detection of certain gases, but also with interpreting the measured information from the sensor. As the resolution of the sensor goes up, the signal to noise ratio goes down.

The noise characteristics of graphene sensors is not well understood [11]. Making a sample which is easy to work with and is suitable for noise analysis is therefore very important. Some researchers provide some information about their experiment configuration, such as [12] or [13], little details are given regarding their sample mounting and configuration.

We are focusing on creating and testing a recipe for 2D sensor mounting which is documented in this paper, and later on this configuration will be used for sensor measurement and data analysis, mainly in spectral domain in a low-noise environment.

II. SENSOR MEASUREMENT

Generally there are multiple methods of data acquisition from sensors. Some sensors change their capacitance based on the gas concentration. Most of the sensors change their resistance. Theoretically, even optical state of the sensor and acoustic wave propagation differences can be used to detect gas on the sensing material. Most of the sensors contain electrodes that interface the sensing material and the electrodes provide means for measurement of the sensing material physical properties. These measurements may be more or less straight forward. Chemiresistive gas sensor measurement is discussed in this chapter. Graphene sensors fall into this category and measurement of the sensors brings its own challenges. Especially with graphene sensors, which have potential to be very sensitive. Measurement can be performed using time and frequency domain. Time domain measurement is simple detection of resistance at particular time. Data processing is many times required to interpret the resistance as a concentration of particular gas. Frequency domain measurement consists of measuring the impedance of the sensing material at particular frequency and a range of frequencies is chosen. This can be beneficial, as some sensing materials have predictable resonance frequency and its shift is used to detect the concentration of the gas. Noise is a factor that influences the sensing performance too. Reaction process at the sensing material can be so slow, that the noise of the system may overtake and measurement can be compromised. This kind of noise is low-frequency noise (also known as 1/f noise). It is discussed in semiconductor based sensors but it may be also applicable to graphene sensors, as graphene provides its own source of noise. In research facilities, gas sensors are measured with devices such as 8753ES [1]. Network analyzer can measure resonant frequency by measuring resonator's return loss. Otherwise, resistance measurement with a good supply of current suffices. Sensors contain multiple sources of noise. Most dominant are two. Thermal noise and lowfrequency noise. Thermal noise is present in all devices. In thermal sensors the thermal noise is caused by vibration of the atomic structures that produce current inconsistencies. This noise has spectral density

$$S = 4 \cdot k \cdot T \cdot R \tag{1}$$

where k is Boltzmann constant, T is temperature of the material in kelvin and R is the resistance of the material. This noise is constant at all measurement frequencies and this means it can be isolated. The other source of noise is the low-frequency noise and this noise is dominant at frequencies usually up to 100kHz [2]. The spectral density of this noise is roughly 1/f. This noise decreases with increasing frequency forming a knee called f0. At this knee the noise decreased to the point where it equals the thermal noise [2]. Another source of noise is noise from material inconsistencies. These inconsistencies come from metallic interfaces for example. Within graphene sensors, the grain boundaries contribute to noise. Grain boundary noise of graphene contributes to 1/f noise [3]. After the data from the sensor is acquired, the data can be post processed to increase the accuracy of the detected gas type and concentration. It is known that sensors tend to have long detection and recovery time and this poses problem when sensor is continuously measuring gas concentration. That is why advanced methods of gas sensor data processing have been developed, such as back propagation neural networks [4].

III. MEASUREMENT SETUP

We have several pristine graphene sensors available with geometry 1 and 2. These sensors are commercially available and many times require separation from a die and cleaning. These two steps are crucial for subsequent measurements with these sensors and they will be discussed later. The pattern used in the star electrode arrangement is called 'Hall Bar Geometry'. This electrode arrangement is done purposefully to enable the measurement of hall-effect [14]. However, the electrodes could be used for potential injection experiments. An ideal case would be to manufacture a silicon-based 3D structure with integrated analog amplifier, with graphene layer transferred directly onto the silicon oxide, with electrodes running several micrometers to the analog amplifier, that would be electrically well isolated and whose output would be available via bonding pads. This sensor must be connected to the amplifier via longer route and for good workability requires bonding onto a plate, which can be better worked with. The path should be low-noise and configurable at the same time. To efficiently test these sensors, we came with a solution with the following recipe: A standard FR-4 board with 25µm of copper was pattern-etched using photolitography. Sodiumhydroxide with 30% hydrogen-peroxide was used for copper etching and sodium-carbonate was used for developing. This board was coated using ENIG process, which includes chemical plating of nickel onto copper (4.5µm) and chemical plating of gold onto nickel $(0.05\mu \text{ m})$. This is a 'soft' gold, while 'hard' gold is electroplated. Soft gold is cheaper to manufacture and there was a concern if the golden bonding wires would attach well onto 'soft' gold. Two-sensors bundle was cut off of a die and this bundle was cleaned and glued directly onto the PCB, using conductive silver-based polymer paste with baking time of 1.5h at 80°C. After the glue was set, the sensor was wedge-bonded by 25µ m Au wire. The bond onto the sensor pad was performed at 340kHz needle vibration during 230mS at 350mN force. This bonding was reproduced across all pads. The bonds onto the golden pad on the PCB were more difficult. The recipe is 440kHz during 320mS at 430mN force and the parameters vary based on the nonhomogenity of the golden layer. Electroplated gold might be more homogenous. The Figures 4 and 5 show a visualization of the traces layout. To enable various configurations during testing, a matrix connection PCB was created, to enable a connection of any connector pin to any sensor pad. These structures enable to measure the two sensors at the same time. This setup is portable and the sensor is well-fixed onto the PCB. U.FL connectors were introduced for better HFfriendly path and easy connection and disconnection of the sensors. It has been proven to us that investing this little money and time into preparation of the samples using these connectors simplifies the work with the samples and is worth it. Connecting and disconnecting the sensors for measuring in a noise-reduced environment is therefore easy. The traces have been gold-plated to minimize changes of the material from the pads to the amplifier, although there is a concern that the inter-metallic phase of CU-NI-AU planar stack will contribute to the noise more, than a single hard side-by-side transition between AU and CU. This means that ultra-high sensitivity experiments may not be suitable with this solution.

Figure 1 requires three extra bonds to interconnect the con-

nectors to the sensors and this is better suited to measure the hall-effect. These boards can be mass-produced in a factory and the primary bottleneck is the bonding, which too could be automated.

Sensor is then connected to a circuit as shown in Figure 3, where chemiresistive properties can be measured. Current from a battery source is recommended, as batteries provide steady and very low-noise current. Voltage on the resistor is then amplified and measured in the spectral domain.



Fig. 1. Star electrode arrangement, acquired using scanning electron microscope.



Fig. 2. Two-electrode arrangement, acquired using scanning electron micro-scope.



Fig. 3. Electrical diagram for mesuring the grapene sensor.

V. CONCLUSION

There is a great potential in graphene sensor technology. But making measurements on the graphene samples is not straightforward. We have proposed a solution on sensor preparation including its dicing and cleaning. Experiments with PCB substrate preparation, gold protection and sensor binding with U.FL sensor terminals were performed and proposed solution is portable and ready for various measurements. We have also proposed a noise measurement setup. Our results can be used for direct reproduction for reader's concrete application, and if the application is outside the usecase defined in this paper, the results and processes provided can be used for inspiration. The sample we have shown here will be used for noise analysis while changing the chemical environment around the sensor.

IV. CUTTING AND CLEANING

Sensor was cut using Laser-Dicer Oxford. Results after cutting are visible in the Figures 1 and 2. Whole sensor was coated by flake-like several microns thick deposition layer, which can be removed by acetone and incressed temperature. We used the following recipe that should theoretically clean the sensor: Preheat the sensor to 80° C and pour acetone directly on it. Then transfer the sensor to ultrasonic cleaner and clean the sensor for 5 minutes in acetone. While cleaning, the sensor cools down to room temperature. Then clean the sensor using isopropylalcohol to get rid of the dissolved molecules from the previous step. Also use an ultrasonic cleaner for 5 minutes at room temperature.

The result was that the graphene layer was indeed cleaned, but after raman spectroscopy, we observed defects of several microns in the graphene layer. Results are shown in the Figure 6.



Fig. 4. PCB for the sensor interface with direct connections.



Fig. 6. The red area is missing graphene. Image is from a Witec Alpha 300R Raman imager. Red mesh is edited in for clarity.

ACKNOWLEDGMENT

First of all, I thank Lord Most High for creating such an interesting universe that we study until now and humbly still don't understand a lot, and I thank Him for giving me the ability to gain wisdom, learn and understand concepts I write about in this paper. I thank to my great supervisor doc. Mgr. Dinara Sobola PhD. for her help in overcoming obstacles that occurred while progressing in this research. Thanks to Brno University of Technology for providing the necessary equipment for the research, Gatema company for the help in manufacturing of necessary PCB parts, and several colleagues for their advices and practical help.

REFERENCES

- S. CHOPRA, A. PHAM, J. GAILLARD, A. PARKER, A.M. RAO, Carbon-nanotube-based resonant-circuit sensor for ammonia, Appl. Phys. Lett. Vol. 80, pages 4632-4636, June 2002.
- [2] A. A. BALANDIN Low-frequency 1/f noise in graphene devices Nature Nanotechnology, pages 549-555, 2013.
- [3] W. SHIN, S. HONG, Y. JEONG, G. JUNG, J. PARK, D. KIM, K. CHOI, H. SHIN, R. KOO, J. KIM, J. LEE, *Low-frequency noise in gas sensors:* A review, Sensors and Actuators B: Chemical, Volume 383, 2023.
- [4] M. SRIYUDTHSAK, A. TEERAMONGKOLRASASMEE, T. MORI-IZUMI:Radial basis neural networks for identification of volatile organic compounds, Sensors and Actuators B: Chemical, vol. 65, no. 1, pages 358-360, 2005.
- [5] Z. MENG, R. M. STOLZ, L. MENDECKI, K. A. MIRICA: *Electrically-Transduced Chemical Sensors Based on Two-Dimensional Nanomaterials*, Chemical Reviews vol. 119, no. 1, pages 478-598, 2019.
- [6] J. H. CHOI, J. LEE, M. BYEON, T. E. HONG, H. PARK: Graphene-Based Gas Sensors with High Sensitivity and Minimal Sensor-to-Sensor Variation, ACS Applied Nano Materials vol. 119, no. 1, pages 2257-2265, 2020.
- [7] F. G. Banica: Chemical Sensors and Biosensors: Fundamentals and Applications, John Willey & Sons, Ltd. ISBN: 9780470710661, 2012.



Fig. 5. PCB for the sensor interface with matrix interconnection capability.

- [8] A. H. Castro Neto, F. Guinea, N. M. R. Peres, K. S. Novoselov and A. K. Geim: *The electronic properties of graphene*, Reviews Of Modern Physics, vol. 81, no. 1, pages 109-162, 2009.
- [9] D.S.L. Abergel , V. Apalkov , J. Berashevich , K. Ziegler & Tapash Chakraborty: *Properties of graphene: a theoretical perspective*, Advances in Physics, vol. 59, no. 4, pages 261-482, 2010.
- [10] F. Schedin, A. K. Geim, S. V. Morozov, E. W. Hill, P. Blake, M. I. Katsnelson, K. S. Novoselov: *Detection of Individual Gas Molecules Adsorbed on Graphene*, Nature Materials, vol. 6, no. 1, pages 652-655, 2007.
- [11] Y-M. Lin, P. Avouris: Strong suppression of electrical noise in bilayer graphene nanodevices, Nano Letters, vol. 8, no. 8, pages 2119-2125, 2008.
- [12] J. Ma, M. Zhang, L. Dong, Y. Sun, Y. Su, Z. Xue, Z. Di: Gas sensor based on defective graphene/pristine graphene hybrid towards high sensitivity detection of NO₂, AIP Advances, vol. 9, 075207, 2019.
- [13] S. Novikova, N. Lebedevaa, A. Satrapinskib, J. Waldenc, V. Davydovd, A. Lebedevd: *Graphene based sensor for environmental monitoring of* NO₂, Sensors and Actuators B: Chemical, vol. 236, no. 1, pages 1054-1060, 2016.
- [14] K. R. Amin, A. Bid: Graphene as a sensor, Currenet Science, Vol. 107, No. 3, pages 430-436, 2014.

Development of Beta-NMR Detection Electronics for VITO Beamline at ISOLDE Facility at CERN

Daniel Havránek Brno University of Technology, FEEC, CERN Brno, Czech Republic <u>221061@vut.cz</u> Michael Pešek CERN, ISOLDE Geneva, Switzerland michael.pesek@cern.ch Daniel Paulitsch University of Innsbruck, CERN Geneva, Switzerland paulitsch.d@gmail.com

Magdalena Kowalska CERN, ISOLDE UNIGE, DPNC Geneva, Switzerland magdalena.kowalska@cern.ch

Mark Lloyd Bissel CERN, ISOLDE Geneva, Switzerland mark.lloyd.bissell@cern.ch Jan Král Masaryk University Brno, Czech Republic jan.kral@fi.muni.cz

Beta detected NMR (Nuclear Magnetic Resonance) is a method to determine the magnetic moments of short-lived isotopes with very high precision using beta particles emitted when the isotopes decay. This paper focuses on designing electronics for a beta particle detector to be used at the beta-NMR beamline VITO at ISOLDE facility at CERN. In the detector beta particles lead to the creation of eV photons which are detected with use of SiPMs (Silicon PhotoMultipliers). Key aspects of the design are linearity, vacuum compatibility, magnetic field compatibility and small size, due to space constraints.

Keywords—beta-particle detection, Beta-detected-NMR, SiPM, PCB, AD8001, AFBR-S4N66P024M, CERN, ISOLDE, VITO

I. INTRODUCTION

 β detected NMR can be used to determine the magnetic moments of short-lived isotopes with higher sensitivity than conventional NMR. The next isotope to be investigated with β detected NMR at VITO is ¹¹Be. For such measurements the β detectors are crucial. Their purpose is to detect asymmetric emission of β -particles at 0 and 180 degrees to a magnetic field direction from the hyperpolarized decaying isotopes. The β detection setup consists of two detectors – the Big that can be seen in Fig. 12 and the Small detector. [1]

II. NMR

A. Conventional NMR

Nuclear magnetic resonance is an analytical technique employed in various fields. Atomic nuclei with nuclear magnetic moments form discrete spin orientation energy levels in external magnetic fields. These energy levels correspond to electromagnetic frequencies and after an application of a weak oscillating magnetic field, electromagnetic waves emitted by the nuclei are detected and analyzed. The resonance frequency is dependent on the electron shielding of individual atoms, making it possible to identify structures of molecules. [2]

Research funded by CERN and ERC

B. Beta-NMR

 β -NMR is type of NMR that is based on beta decay of a polarized ensemble. It extends the capabilities of conventional NMR to study of magnetic properties of radioactive nuclei. The NMR resonances can be observed as changes in the beta decay asymmetry which significantly increases the signal strength. β -NMR relies on a beam of hyperpolarized short-lived β -decaying nuclei. [3]

III. THE EXPERIMENTAL SETUP

The experimental setup for beta detection part for β -NMR consists of two β detectors. They consist of a plastic material that emits visible-light photons under impact from β particles and several SiPMs. The last part of the detection setup is a DAQ (Data Acquisition System) used for processing the signals from detectors. The block diagram of the setup can be seen in Fig. 1.

The DAQ system consists, among others, of attenuator/amplifier, oscilloscope and host PC for controlling the DAQ system. The oscilloscope that is used is NI PXIe-5170R. It's input voltage has fixed range of -5 V to 5 V. This needs to be considered in the design not to overshoot this voltage limit. Optionally the signals can be filtered before entering the DAQ system.

Both detectors are made of the same components, only the dimensions are different. A summing card is the same for both detectors and its purpose is just to sum signals from all detection PCBs (Printed Circuit Board) in the specific detector.

Each detector consists of eight detection PCBs. PCBs for the Big and Small detector have the same schematic, but dimensions of the PCBs and their layout are different due to space constraints. The design of the detection PCBs was inspired by a board designed at DPNC (Department of Nuclear and Particle Physics) at the University of Geneva.



Fig. 1. The experimental setup for beta detection

IV. BETA PARTICLE DETECTION

The detection of beta particles starts with a scintillator. The scintillator is a material that emits (fluorescence) photons, when a β particle passes through it, because the latter excites the electrons in the atoms of the material. [4]

For the photon detection we used SiPMs (Silicon PhotoMultipliers). SiPM is a component similar to a photodiode. It is a photodetector that generates short current pulses with very steep rising edge as a response to photon absorption. SiPMs are highly sensitive with very low time jitter that allows them to capture single photons. Between the scintillator and the SiPMs, we introduced a plastic diffuser that allows a similar fraction of emitted photons from the scintillator to reach the SiPMs, independent of the position of the beta particle interaction. In our design we are using AFBR-S4N66P024M SiPM. This SiPM was selected by a PhD student on the project, Daniel Paulitsch according to its good parameters and low price. [5,6]

V. DESIGN AND SIMULATION

The new β -NMR measuring system should be able to also measure the energy of the β particles. The energy of the particles is directly proportional to the number of emitted photons. For measuring the energy, the linearity of the detection PCB is strongly advantageous.

The whole design is also placed in a strong magnetic field of 4.7 T that is necessary for the β -NMR measurement. The design needs to be magnetic field compatible thus all components need to be non-magnetic or very weakly magnetic, not to be affected by the magnetic field and also not to disturb the magnetic field. This is especially crucial for big metallic parts, such as connectors.

The whole detector system will be placed in a high vacuum that is also necessary for the β -NMR measurement. Thus, the whole design needs to be vacuum compatible. This is important mostly for plastic parts of the design which can reduce the vacuum quality because of outgassing (when a small amount of gas escapes from the plastic material). This is especially crucial for big plastic parts such as cables.

A. Detection PCB

One of the key aspects of the design is linearity. We succeeded in designing such a circuit and we simulated it in LTspice. The simplified schematic of the circuit can be seen in Fig. 2.

AFBR-S4N66P024M



Fig. 2. Simplified schematic of detection PCB

All sixteen SiPMs are connected in parallel, and their output goes into a summing amplifier. For the summing amplifier AD8001AR is used. It is the same amplifier that was used in the University of Geneva design. They are using it for the same purpose, and they selected carefully due to its parameters such as bandwidth, speed and others. For saving the time we used the same amplifier. We checked its magnetic properties with strong magnetic field and it is only weakly magnetic. The model of the SiPM for simulation was given to us by colleagues at University of Geneva. We have just changed the model to our specific SiPM. An example of the output signal can be seen in Fig. 3. The results of the simulation can be seen in Fig. 4 and 5. In the graph there is a peak voltage values for specific number of photons and cells.



Fig. 3. Simulated output pulse for 116 photons and 16 cells



Fig. 4. Simulation results – Voltage output dependence on number of absorbed photons



Fig. 5. Simulation results - Output voltage dependence on number of cells

As we can see in Fig. 4, the design is linear, up to roughly 230 photons absorbed by all sixteen SiPMs. This range is sufficient for the first planned isotope ¹¹Be. The expected maximum number of photons for ¹¹Be beta particles is 116, according to the results of simulation (using GEANT4 software) performed by Daniel Paulitsch. At some point the circuit stops being linear. It is caused by the saturation of the amplifier. When the SiPMs detect too many photons, their output signal becomes too large and the amplifier is saturated. But this can be easily solved by changing the gain of the amplifier.

The Small and Big detector PCB needs to have specific dimensions to fit inside of the detector. For a better detection, SiPMs had to be placed at the corners of the PCB and no components between them on the same side of the PCB. For the Big PCB this means that all components need to be on the other side of the PCB. The Small PCB is extended so the connectors

can be reachable. Thus, some small components are allowed on the other side. This also limits the design to SMD (Surface Mount Device) components only.

Another important aspect of the PCB design is the length matching of the signals. Each trace needs to be the same length, otherwise the signal from each SiPM arrives at a different time, causing the resulting signal to be distorted.

Optionally the PCBs can be cooled to low temperatures (down to -40 °C) to lower the noise. We do not expect for the electronics to be negatively affected by the low temperature. The PCBs need to be designed with option of cooling them. The temperature needs to be the same for all SiPMs to achieve low noise and good stability. For this cooling pads were added on the other side of the PCBs bellow the SiPMs. For the cooling a custom aluminum heatsink will be used with commercial chiller. The PCB design of the Small PCB can be seen in Fig. 6 and 7 and for the Big PCB in Fig. 8 and 9.



Fig. 6. Small detection PCB - bottom



Fig. 7. Small detection PCB - top



Fig. 8. Big detection PCB - bottom



Fig. 9. Big detection PCB - top

B. Summing card

The summing card combines signals from eight detection PCBs to one output. All eight inputs go again to a summing amplifier with a specific gain not to overshoot the oscilloscope's range. The PCB itself was not simulated. All of the PCBs were designed in KiCAD and we had no PCB simulator available. But we tried to replicate this effect in LTspice by adding additional capacitance. According to this the amplifier should be stable up to 4 pF. The layout of the summing PCB is similar to that of the Small and Big detector PCBs. The summing card can be seen in Fig. 10.



Fig. 10. Summing card - top

C. Power distrubution PCB and Filter

Power distribution PCB serves only for power distribution to the detection PCBs and Summing card. The whole design is sealed in vacuum tight detector. The number of cables going in and out is limited by the space available for electrical feedthroughs. Due to these reasons, it is better to have a lower number of cables going in and then split them. The Power distribution PCB can be seen in Fig. 11.



Fig. 11. Power distribution PCB - top

The signal from detectors can be filtered before it is processed by DAQ. For the filtering, a simple first order high-pass filter can be used. If it is used, the filter removes the low frequencies from the signal causing faster return to zero and shortening the tail of peaks.



Fig. 12. The Big detector

CONCLUSION

All necessary PCBs for the β -NMR experiment at CERN-ISOLDE were designed. The key parameters for the design, such as magnetic field and vacuum compatibility, linearity and dimensions were obtained. The next steps is to order the PCBs and conduct the first measurement. The last step is to finish the machining of both scintillators and test everything in real-life operation.

ACKNOWLEDGMENT

I would like to extend my sincere thanks to my supervisor Michael Pešek. I would also like to thank to prof. Federico Sanchez Nieto, Dr. Yannick Favre and their team DPNC at the University of Geneva for sharing their original design and for assisting us with our design.

References

- PAULITSCH, Daniel. Development of a new β Detector Setup for the VITO Beamline. Online. 2023. Available online: <u>https://indico.cern.ch/event/1316940/contributions/5639658</u>. [cit. 2024-01-31].
- [2] RAJA, Pavan V. M. and BARRON, Andrew R. NMR Spectroscopy. Online. 2016. Avalable online: https://chem.libretexts.org/Bookshelves/Analytical Chemistry/Physical Methods in Chemistry and Nano Science (Barron)/04%3A Chemica <u>1 Speciation/4.07%3A NMR Spectroscopy</u>. [cit. 2024-03-06].
- [3] CERN. Beta-NMR. Online. 2017. Available online: <u>https://espace.cern.ch/ISOLDE-SSP/SitePages/B-NMR.aspx</u>. [cit. 2024-03-06].
- [4] MITSUBISHI CHEMICAL GROUP. An easy to understand the scintillator. Online. 2019. Available online: <u>https://www.mchemical.co.jp/en/products/departments/mcc/ledmat/tech/1203825_7554</u> <u>.html</u>. [cit. 2024-01-31].
- [5] PIATEK, Slawomir. What is an SiPM and how does it work? Online. 2016. Available online: <u>https://hub.hamamatsu.com/us/en/technicalnotes/mppc-sipms/what-is-an-SiPM-and-how-does-it-work.html</u>. [cit. 2024-01-31].
- [6] BROADCOM. AFBR-S4N66P024M: 2x1 NUV-MT Silicon Photomultiplier Array. Online. Available online: <u>https://docs.broadcom.com/doc/AFBR-S4N66P024M-DS</u>. [cit. 2024-01-31]. Datasheet.

Transparent materials for planar microelectrode arrays

Bc. Jaromír Jarušek Department of Microelectronics Brno University of Technology Brno, Czech Republic Jaromir.Jarusek@vutbr.cz Ing. Jan Brodský Department of Microelectronics Brno University of Technology Brno, Czech Republic Jan.Brodsky@vutbr.cz Ing. Imrich Gablech, Ph.D. Department of Microelectronics Brno University of Technology Brno, Czech Republic Imrich.Gablech@vutbr.cz

Abstract—This paper is focused on indium tin oxide (ITO), ultrathin gold, titanium nitride (TiN) and diamond-like carbon (DLC) thin films that can be used for transparent or semitransparent electrodes and interconnects for planar microelectrode arrays. Paper further describes methods of thin film deposition of selected materials, microfabrication of planar microelectrode arrays, their preparation for electrochemical measurement and results of electrochemical impedance spectroscopy and cyclic voltammetry.

Keywords— planar microelectrode array (pMEA) indium tin oxide (ITO), titanium nitride, diamond like carbon (DLC), ultrathin gold thin film, electrochemical impedance spectroscopy

I. INTRODUCTION

Planar microelectrode arrays (pMEAs) are devices for *in vitro* research of electrogenic cells, i.e. electrically active cells. Those include neurons, heart cells (cardiomyocytes), retinal cells and smooth muscle cells. They are important tools for neural electrophysiology and investigations of neural network signaling and neural plasticity, more general applications include toxicity testing and drug screening. [1], [2], [3] *In vitro* drug screening is especially important for pharmacological applications because some drugs for disease treatment, unrelated to cardiac diseases, can cause life-threatening cardiac rhythm disturbances (arrhythmias). The usage of pMEAs for this application presents an easier and efficient method in comparison to approaches such as isolated Purkinje fibers or *ex vivo* Langendorff heart. [4]

For cell studies inverted optical microscope is preferred, since observation from top may not be possible, because of presence of culturing medium. The usage of opaque electrodes and interconnections limits the view, thus usage of transparent or semitransparent electrodes is desirable for such applications. [5].

In the field of commercial pMEAs it is common that electrode materials and even materials of interconnects are opaque. [2] Some of the examples are MCS 60MEA with opaque TiN and advertised impedance of $< 100 \text{ k}\Omega$ for 30 µm electrode diameter or MED64 with platinum black electrodes with advertised impedance of $10 \text{ k}\Omega$ for 30 µm electrode diameter. [6], [7]

The only commercial pMEA, with advertised transparent electrodes, is MCS 120tMEA with advertised impedance of $< 250 \text{ k}\Omega$ and utilization of semitransparent TiN. Specified

diameter of electrodes is 100 μ m. [8] Further information about used TiN layer, such as thickness or optical transmittance, is not known. It could also be pointed out that if electrode diameter 100 μ m is not an error in documentation, guaranteed impedance of < 250 k Ω is relatively high.

II. MATERIALS FOR MICROFABRICATION OF PMEAS

Two approaches were employed to fabricate semitransparent electrodes and interconnects. The former utilizes ultrathin ≈ 9 nm gold layer as a main conductor, and latter uses highly transparent ITO layer with low sheet resistance as a main conductor and different material which may have significantly higher sheet resistance, but has better chemical, electrochemical, or other properties. As a source of ITO, commercially available 4" ITO coated glass wafers with advertised sheet resistance $< 10 \ \Omega \cdot \Box^{-1}$ were used. The value of sheet resistance has been verified to be in range from $6.5 \ \Omega \cdot \Box^{-1}$ to $7.8 \ \Omega \cdot \Box^{-1}$ for different wafers. Besides ITO other materials were deposited in vacuum systems in CEITEC Nano cleanrooms.

A. Deposition of thin semitransparent conductive films

Prior to depositions, all substrates were cleaned in ultrasonic bath using acetone and isopropyl alcohol with subsequent rinse with deionized water. Further cleaning was performed using oxygen plasma in Diener plasma cleaner and finally *in situ* Ar ions precleaning was utilized just prior to deposition in respective deposition systems.

For deposition of semitransparent gold layer, a borofloat wafer was used. To enhance adhesion to the substrate, ≈ 1 nm of titanium was sputtered before gold in the same process and the same deposition parameters, without breaking the vacuum. The semitransparent gold layer with thickness of ≈ 9 nm was deposited using ion-beam sputtering technique. The thickness of ≈ 9 nm was chosen as a compromise between sheet resistance, which steeply rises below 10 nm, and optical transmittance, which is around 50 % at 10 nm. [9] The deposition was performed in system equipped with two Kaufman & Robinson RFICP40 ion-beam sources. Energy of ion beam was 600 eV and current 44 mA.

Similarly to the semitransparent gold layer, 1 nm of titanium was sputtered as adhesion layer under TiN. TiN layer with thickness of ≈ 20 nm was deposited using ion-beam assisted deposition using both ion-beam sources. This process was described in more detail in previous publications. Deposition was carried out in temperatures < 100 °C with process

parameters as follows. Energy of primary beam was 600 eV, primary beam current was 44 mA, primary beam gas supply was a mixture of argon and nitrogen with mass flows 3.6 sccm and 3 sccm respectively. Energy of secondary beam was 26 eV and secondary beam current was 15 mA.[10], [11]

DLC layer with thickness of $\approx 150 \text{ nm}$ was deposited directly on precleaned ITO layer by plasma-enhanced chemical vapor deposition (PECVD) using Oxford Instruments PlasmaPro 80 RIE system. Layer was deposited from CH₄ plasma. Deposition was carried out at an elevated temperature of $\approx 60 \,^{\circ}\text{C}$ and working pressure of $\approx 53.3 \,^{\circ}\text{Pa}$ at the corresponding deposition rate of $\approx 8 \,\text{nm}\cdot\text{min}^{-1}$. [12]

B. Optical transmittance of used films

Optical transmittance was measured using three-channel optical spectroscope Ocean Insight JAZ3 and extracted results are shown in Fig. 1. It is apparent that ITO has the highest optical transmittance in the visible range out of all investigated conductive layers. An interesting result is the DLC layer on ITO, because decrease in transmittance compared to ITO layer alone is only small. DLC on ITO layers could with combination of good electrical and electrochemical properties present very interesting material combination for transparent pMEA. This is in contrast to TiN with average transmittance ≈ 45 % and ultrathin gold thin film with average transmittance ≈ 55 %, that can serve only as semitransparent conductors. For further comparison DLC-ITO layer has transmittance on par with other optical parts of pMEA like encapsulation consisted of 50 nm AlN and 3 μ m of Parylen-C.



Fig. 1. Graph of optical transmittances of prepared thin films.

III. MICROFABRICATION OF PMEAS

A. Description of used design

Design consisting of three photolithographic masks was used for microfabrication of pMEAs. Illustration of pMEA design is shown in Fig. 2. First mask contains a pattern for transparent electrodes and underlying interconnects. The second mask contains a pattern for gold metallization of pads and interconnects outside of the active diameter which is comparable with field of view of 10x microscope objectives. Pattern of the third mask allows the creation of opening in deposited encapsulation for pads and transparent electrodes.



Fig. 2. Design of pMEA chip with pitch of 100 μm Parts conductive pattern are highlighted with colors.

Used masks contain four variants of pMEA design, where pitches between microelectrodes are 100 μ m and 200 μ m and diameters of microelectrodes are 30 μ m and 50 μ m. The variant with 100 μ m pitch and 30 μ m is illustrated in Fig. 2.

B. Photolithographic process

pMEAs were fabricated at CEITEC Nano class 100 cleanrooms. Photolithography was performed using semiautomated coating system SÜSS MicroTec RCD8 and a mask aligner SÜSS MicroTec MA8 Gen3. All etching processes were performed using plasma processes, specifically by ion beam etching using Scia Systems Coat 200 and reactive ion etching using Oxford Instruments Plasma Technology PlasmaPro 100 and PlasmaPro 80. For etching using ion beam, AZ MIR 701 type photoresist spincoated at 4000 rpm was used. Ion beam etching was used for gold and DLC layers.

Workflow of fabrication of conductive pattern had varied for different materials. In case of ultrathin film gold electrodes, semitransparent thin film was deposited after pad and interconnection metallization. In case of TiN, underlying ITO pattern was etched first and subsequently TiN thin film was deposited. DLC on ITO was etched using same lithography using ion beam etching and subsequently reactive ion etching.

After completion of conductive pattern, encapsulation layers were deposited. First, a 50 nm aluminum nitride (AlN) thin film was sputtered using two ion beam assisted deposition, then 3 μ m of Parylene-C was deposited using SCS Parylene Deposition System. For etching of opening in encapsulation for pads and

microelectrodes, thick photoresist AZ 1518 spincoated at 1000 rpm was used. Parylene-C was etched using oxygen plasma and AlN was etched using mixture of BCl_3 and Cl_2 plasma.

Before the wafers dicing a photoresist was spincoated on the wafer for protection during dicing process. An example of the results of manufactured pMEA is in Fig. 3.



Fig. 3. Example of manufactured pMEA under optical microscope in transmitted light.

IV. ELECTROCHEMICAL CHARACTERIZATION

A. pMEAs preparation

In order to use the fabricated chips for measurements, several preparatory steps were required. First, the chips were cleaned using the following cleaning procedure. The chip held by the tweezers was rinsed with acetone using a washing bottle to wash away the debris from dicing. Then, the chip was moved to a beaker with acetone, where it was agitated in an inclined position for better removal of the protective photoresist. To completely clean the residual impurities, the chip was left in pure acetone for 5 minutes, another 5 minutes in isopropyl alcohol, followed by rinse with pure isopropyl alcohol by agitating in a beaker, and dried with compressed air. Proper results of cleaning procedure were verified using a microscope.

Second step involves fixing of the chip in the open cavity of ceramic leadless chip carrier (LCC) package using a small amount of molten wax. The pads of pMEAs were then electrically connected to the package using wirebonding method. Wirebonding was done using TPT HB16 wire bonder in wedge bonding setup with gold wire.

The last step involves attaching a small well to the chip for holding electrolyte or cell culture medium. Wells were made using fused filament fabrication (FFF) 3D printing from polyethylene terephthalate glycol (PETG) filament. Wells have internal dimensions $\approx (3.3 \times 3.3) \text{ mm}^2$, wall thickness of \approx 0.45 mm, and height ≈ 6 mm. Attachment was done manually by applying tiny amount of polydimethylsiloxane (PDMS) on base edge of well and then placing it on pMEA using tweezers. Applied PDMS was partially cured in oven at temperatures around 80 °C. Another dose of PDMS was applied to enhance watertightness and attachment robustness. PDMS was applied on the rest of chip outside the well, including pads and wire bonds, and then cured in oven at temperatures around 80 °C for 30 minutes. Illustration of pMEA prepared for electrochemical measurements is in Fig. 4.

A different approach was used for ITO-only microelectrodes chips, where issues with wirebonding reliability were encountered. This limited amount of measurement results.



Fig. 4. Cross-section of pMEA prepared for electrochemical measurements.

B. Electrochemical measurements

Prepared pMEAs, as depicted in Fig. 4, were placed into platform with connectors and socket for the LCC package. As an electrolyte for electrochemical measurement, a 0.1M KCl solution was used. The solution was applied using micropipette. In case of ITO chip, pMEA was rinsed with ethanol before application of the 0.1M KCl solution. A silver wire with electrodeposited AgCl tip was used as a pseudoreference electrode.

Electrochemical measurements were performed using Metrohm Autolab III potentiostat. Measurements consisted of 5 cycles of electrochemical impedance spectroscopy interlaced with 10 cycles of cylic voltammetry in range (-0.3 + 0.5) V relative to the Ag/AgCl pseudoreference electrode. For the measurements, pMEAs with 50 µm microelectrode diameter and, in most cases, with 100 µm pitch were selected.

V. RESULTS OF ELECTROCHEMICAL MEASUREMENTS

From Fig. 5, it can be concluded that 9 nm gold, 20 nm TiN on ITO and 150 nm DLC on ITO exhibit similar impedance magnitudes down to 10 Hz, whereas the ultrathin gold layer deviates to higher magnitudes at lower frequencies. Bare ITO electrodes exhibit significantly higher impedance across the entire frequency spectrum. From Fig. 6, it can be further observed that double layer capacitance is so low that series resistance of electrochemical system is not clearly apparent even at high frequencies.

These characteristics can also be seen in cyclic voltammograms in Fig. 7. Both TiN and DLC have similar capacitive curves with a larger area under the curves, that suggests higher double layer capacitance than gold and bare ITO. Impedances for ultrathin gold, semitransparent TiN and DLC were found to be $\approx 700 \text{ k}\Omega$.



Fig. 5. Statistical results of electrochemical impedance spectroscopy impedance magnitude.



Fig. 6. Statistical results of electrochemical impedance spectroscopy - impedance phase.



Fig. 7. Averaged steady-state cyclic voltammograms.

VI. CONCLUSION

From measured optical transmittances and electrochemical measurements, it can be concluded that DLC on ITO layer has very interesting properties in comparison to other investigated materials. Further investigation of DLC layers is important to better evaluate the potential for practical use as material with superior properties for pMEAS. Further optimization of TiN and DLC layers is desirable for more competitive performance.

ACKNOWLEDGMENT

CzechNanoLab project LM2023051 funded by MEYS CR is gratefully acknowledged for the financial support of the measurements/sample fabrication at CEITEC Nano Research Infrastructure.

REFERENCES

- H. S. Jeong, S. Hwang, K. S. Min, and S. B. Jun, "Fabrication of Planar Microelectrode Array Using Laser-Patterned ITO and SU-8", *Micromachines*, vol. 12, no. 11, 2021.
- [2] M.-G. Liu, X.-F. Chen, T. He, Z. Li, and J. Chen, "Use of multi-electrode array recordings in studies of network synaptic plasticity in both time and space", *Neuroscience Bulletin*, vol. 28, no. 4, pp. 409-422, 2012.
- [3] A. Koklu, R. Atmaramani, A. Hammack, A. Beskok, J. J. Pancrazio, B. E. Gnade, and B. J. Black, "Gold nanostructure microelectrode arrays for in vitro recording and stimulation from neuronal networks", *Nanotechnology*, vol. 30, no. 23, Jun. 2019.
- [4] A. Stett, U. Egert, E. Guenther, F. Hofmann, T. Meyer, W. Nisch, and H. Haemmerle, "Biological application of microelectrode arrays in drug discovery and basic research", *Analytical and Bioanalytical Chemistry*, vol. 377, no. 3, pp. 486-495, Oct. 2003.
- [5] T. Ryynänen, R. Mzezewa, E. Meriläinen, T. Hyvärinen, J. Lekkala, S. Narkilahti, and P. Kallio, "Transparent Microelectrode Arrays Fabricated by Ion Beam Assisted Deposition for Neuronal Cell In Vitro Recordings", *Micromachines*, vol. 11, no. 5, 2020.
- [6] "MED Probe (MEA)" https://www.med64.com/products/med-probemea/ (accessed Apr. 01, 2024)
- "60StandardMEA", multichannelsystems.com. https://www.multichannelsystems.com/sites/multichannelsystems.com/fi les/MCS_60StandardMEA_Layout.pdf (accessed Apr. 1AD)
- [8] "120tMEA100/30iR-ITO," 2019. https://www.multichannelsystems.com/sites/multichannelsystems.com/fi les/documents/data_sheets/120tMEA100-30-ITO_Layout.pdf (accessed Apr. 01, 2024)
- [9] S. Wilken, T. Hoffmann, E. von Hauff, H. Borchert, and J. Parisi, "ITOfree inverted polymer/fullerene solar cells: Interface effects and comparison of different semi-transparent front contacts", *Solar Energy Materials and Solar Cells*, vol. 96, pp. 141-147, 2012.
- [10] J. Jarušek, "Naprašování nitridových vrstev pro bioelektronické aplikace pomocí Kaufmanova iontového zdroje", Bakalářská práce, Brno, 2022.
- [11] I. Gablech, L. Migliaccio, J. Brodský, M. Havlíček, P. Podešva, R. Hrdý, J. Ehlich, M. Gryszel, and E. D. Głowacki, "High-Conductivity Stoichiometric Titanium Nitride for Bioelectronics", Advanced Electronic Materials, vol. 9, no. 4, 2023.
- [12] O. Sharifahmadian, A. Pakseresht, S. Mirzaei, M. Eliáš, and D. Galusek, "Mechanically robust hydrophobic fluorine-doped diamond-like carbon film on glass substrate", *Diamond and Related Materials*, vol. 138, 2023.

Estimating the Equivalent Circuit of Lithium-Ion Batteries During Operation

Jakub Vašíček Department of Electrical and Electronic Technology Brno University of Technology Brno, Czech Republic xvasic33@vut.cz

Abstract— This paper deals with the issue of estimating the equivalent circuit of lithium-ion cells during normal operation. The study explores fundamental battery models, such as the simplified Randles circuit and the full Randles circuit. Conventional methods for model estimation, including Spectroscopy Electrochemical Impedance (EIS) and Galvanostatic Intermittent Titration Technique (GITT) are examined. The research extends into the real-world data analysis, gained from a home photovoltaic system. Using these data, an algorithm to estimate the equivalent circuit model during battery operation is evaluated. The study investigates the correlation between the proposed algorithm's outcomes and results obtained through EIS and GITT methods.

Keywords— lithium-ion batteries, battery cell model, Randles circuit, EIS, GITT, SoC, SoH, BMS

I. INTRODUCTION

Lithium-ion batteries are finding applications in an increasing number of scenarios today, from small batteries designed to power consumer electronics to large stationary storage systems containing thousands of cells. Recently, there has been a growing demand for these batteries, primarily driven by the development of electromobility and the increasing popularity of home photovoltaic systems [1]. This progress is linked with the requirements for precise monitoring of individual cells.

The main goal of this work is to explain the equivalent circuit models of lithium-ion cells, analyze conventional algorithms for their monitoring, and propose an algorithm capable of estimating the equivalent circuit model during battery operation.

The equivalent model has the potential to be used for State of Charge (SoC) estimation precision enhancement. Furthermore, the model might be a powerful diagnostic tool for monitoring the State of Health (SoH), capable of exposing internal defects and abnormalities, including issues that extend beyond the cell case, such as increased contact resistance of the connections between cells.

II. EQUIVALENT BATTERY CIRCUITS

Equivalent battery circuits models describe the battery response to the load transients. The model parameters may be affected by the battery SoC, SoH and by the interconnections between battery cells. Petr Vyroubal Department of Electrical and Electronic Technology Brno University of Technology Brno, Czech Republic vyroubal@vut.cz

A. Simplified Randles Circuit

Simplified Randles Circuit consists of three elements. R_s – the resistance of cell electrodes. The value of R_s can also be affected by the cell interconnection contact resistance. R_p – charge transfer resistance. C_p represents the double layer capacitance, which appears at the interface between cell electrode and electrolyte [2]. The C_p appearance greatly reduces the cell impedance for AC ripple current.



Fig. 1. Simplified Randles Circuit [2]

B. Full Randles Circuit

Full Randles Circuit including the Warburg diffusion impedance Z_{Wd} in addition to the simplified variant. The impedance Z_{Wd} becomes dominant at low frequencies [3].



Fig. 2. Full Randles Circuit [2]

III. CONVENTIONAL METHODS FOR EQUIVALENT BATTERY CIRCUIT ESTIMATION

Conventional methods usually require the battery to be disconnected from the load and measured by special equipment. Consequently, they are impractical for real-world application usage. However, these methods can provide results suitable for reference measurements. All the measurements in this chapter were done using the BioLogic BCS-815 system [4].

A. Galvanostatic Intermittent Titration Technique (GITT)

The GITT method involves measuring the cell's transient response to current pulses, which can be either charging or discharging. This method is primarily suitable for estimating the simplified Randles circuit [5]. Figure 3 displays the transient response of the Winston 60 Ah cell to a 10 A current pulse.



Fig. 3. Transient response of Winston 60 Ah cell, SoC = 75%

To determine the equivalent circuit parameters, three values from the figure must be subtracted: ΔI – current difference, $V_{\rm s}$ – voltage drop on $R_{\rm s}$, $V_{\rm p}$ – voltage drop on $R_{\rm p}$, τ – time constant defined by $R_{\rm p}C_{\rm p}$

The V_s value is obtained as the voltage drop occurring immediately in reaction to the current pulse. The V_p voltage drop can be measured after double layer capacitance C_p charge reaches the steady-state value. The voltage settling time also determines the time constant τ of R_pC_p , which is defined as the time when the voltage reaches 1 - e⁻¹ $\approx 63.2\%$ of the steady state value [6]. The parameters of the equivalent circuit can be calculated using the following equations:

$$R_i = \frac{u_i}{\Delta i} \tag{1}$$

$$R_p = \frac{u_p}{\Delta i} \tag{2}$$

$$C_p = \frac{\tau}{R_p} \tag{3}$$

The following table shows measured and calculated values for the tested cell Winston 60 Ah [7]:

TABLE I. SIMPLIFIED RANDLES CIRCUIT VALUES OBTAINED BY GITT METHOD

Value	Result	Unit	Value	Result	Unit
$V_{\rm i}$	10	V	R _s	1	mΩ
$V_{\rm p}$	20	V	R_p	2	mΩ
ΔI	10	А	$C_{\rm p}$	50	kF
τ	100	S			

B. Electro impedance Spectroscopy (EIS)

The EIS method involves measuring the cell's response to a voltage or current harmonic signal of various frequencies. To ensure a linear response and prevent irreversible changes to the electrochemical system, signals with a small amplitude are used. Based on the measured response, it is possible to calculate the impedance of the cell using equation 4.

$$Z = \frac{v(t)}{i(t)} = \frac{V_0 \sin(\omega t)}{I_0 \sin(\omega t + \varphi)} = Z_0 \frac{\sin(\omega t)}{\sin(\omega t + \varphi)}$$
(4)

Where: v(t) – test voltage, i(t) – test current

 V_0 , I_0 , Z_0 – voltage, current, impedance amplitudes

 ω – Angular frequency

 φ – phase shift

When the measured impedance values are represented as a Nyquist plot, it becomes possible to derive advanced battery models such as the full Randles circuit, by finding correlations between the model and measured data [8]. Various iteration methods are used for this purpose. Figure 4 displays the Nyquist plot for the Winston 60 Ah cell measured on different SoC levels. The line Zfit represents the characteristic of the full Randles circuit estimated by the EIS method.



Fig. 4. Nyquist plot for the Winston 60 Ah cell and estimated model

The plot shows that the cell model is not significantly affected by SoC level. The equivalent circuit was obtained by EC-Lab software, provided with the measuring system, and its parameters are shown in the following table [9]:

TABLE II. FULL RANDLES CIRCUIT VALUES OBTAINED BY EIS METHOD

Value	Result	Unit
R_s	0,93	mΩ
R_p	0,5	mΩ
C_p	100	F
R_d	5	mΩ
t_d	1 000	S

The elements of the full Randles circuit are not equivalent to the simplified variant except for R_s . Therefore, R_s is the only element which can be directly compared with the previous GITT data. The results for R_s are very similar for both methods.

IV. EQUIVALENT BATTERY CIRCUIT ESTIMATION UNDER LOAD

A. Calculating VoC using simplified Randles Circuit

Once the parameters of the equivalent model are determined, it becomes possible to calculate the theoretical V_{OC} of the loaded battery. For simplicity, the simplified Randles Circuit was chosen for initial experiments. To determine the V_{OC} , the voltage drops across the Randles Circuit elements must be calculated.

Voltage drop across the R_s resistance can be calculated simply using the ohms law:

$$V_s(n) = I(n) \cdot R_s \tag{5}$$

Where: *n* is the index of the sample

To calculate voltage drop across R_sC_p , the circuit's differential equation must be solved, resulting in equation 6 [10].

$$V_{p}(n) = R_{p}I(n-1)\left(1 - e^{-\frac{\Delta t}{R_{p}C_{p}}}\right) + e^{-\frac{\Delta t}{R_{p}C_{p}}} \cdot V_{p}(n-1) \quad (6)$$

Where: *n* is the index of the sample

 Δt time difference between two samples

According to Kirchhoff's laws, V_{OC} is simply the sum of all voltages obtained:

$$V_{OC}(n) = V_B(n) + V_S(n) + V_p(n)$$
(7)

Where: V_{OC} is the calculated open circuit voltage

 $V_{\rm B}$ is the measured cell voltage

 $V_{\rm s}$ is the voltage drop across the $R_{\rm s}$ element

 $V_{\rm B}$ is the voltage drop across the $R_{\rm p}C_{\rm p}$ element

To verify if the V_{OC} calculation using the simplified Randles circuit works as expected, the cell with known model parameters obtained by the GITT method was discharged using 30 A current pulses. The discharging process is illustrated in figure 5.



Fig. 5. Calculated V_{OC} during the cell pulse discharging

The measured curves demonstrate that the measured cell voltage $V_{\rm B}$ is significantly influenced by the load current, primarily due to the voltage drop across the equivalent circuit elements. In contrast, the calculated $V_{\rm OC}$ closely simulates an ideal constant current discharge curve.

The analogous experiment was performed after integrating the cell into a real photovoltaic system. Data were gathered over a span of 2 months using a self-built BMS [11]. Figure 6 depicts the cell's behavior collected on November 1, 2023.



Fig. 6. Calculated Voc in real-world conditions

The compensation using the equivalent model worked as expected, even for real-world data. The only imperfection can be seen in the marked area 1. In this area, the load was completely disconnected due to the cell's low SoC. Consequently, the cell voltage took more than 6 hours to settle to its steady state value. With the use of the more accurate full Randles circuit model, the compensated V_{OC} should appear as a straight line in area 1.

B. Estimating simplified Randles Circuit using calculated Voc

Analysis of real-world data shows that, with the appropriate equivalent circuit, the calculated V_{OC} value remains unaffected by the load current and can be regarded as a constant within a short time interval. This situation is illustrated in Figure 7, which displays a 15-minute interval extracted from Figure 6.



Fig. 7. Ripple of the calculated $V_{\rm OC}$

This observation can be used for equivalent circuit estimation. The value of R_s can be obtained based on the immediate voltage reaction to the load change, while the values of R_p and C_p can be obtained using iteration methods.

Figure 8 illustrates the relationship between the V_{OC} ripple voltage V_{PP} and the values of the elements R_p and R_pC_p time constant τ . The same set of data as shown in Figure 7 was used for this analysis.



Fig. 8. $V_{\rm OC}$ ripple voltage $V_{\rm PP}$ vs $R_{\rm p}$ and $R_{\rm p}C_{\rm p}$ time constant τ

The minimum V_{pp} value represents the most accurate battery model. As depicted in Figure 8, the function exhibits a single minimum with no local extremes, making it highly suitable for iteration methods. In the Matlab simulation, the fminsearch function was utilized to find the optimal result. This function employs the simplex search method and has demonstrated highly promising outcomes [12]. Even when initialized with values $\tau = 400$ s and $R_p = 10$ m Ω , which are significantly distant from the expected result as depicted in Figure 8, the optimization process converges quickly, typically within 100 iterations. This efficiency makes the method suitable for computational tasks on average microcontrollers with enough memory like the STM32. In this case, the computation was performed over 180 samples with sample rate of 0.2 S/s, resulting 15-minute measurement time. However, the optimal measurement time needs further evaluation through additional experiments to ensure robustness and accuracy in real-world applications.

V. CONCLUSION

The fundamental battery models, such as the simplified Randles circuit and the full Randles circuit, were described, followed by conventional methods for their estimation, including Electrochemical Impedance Spectroscopy (EIS) and Galvanostatic Intermittent Titration Technique (GITT).

Real-world data analysis, particularly from a home photovoltaic system, has provided practical insights into the behavior of lithium-ion cells in operational environments. Based on this data, an algorithm to estimate the simplified Randles circuit model, utilizing battery transient response during normal operation, was proposed. The algorithm makes use of natural load transients, such as load switching. This makes the algorithm suitable for integration into virtually any battery management system (BMS) without the need for specialized hardware. Basic voltage and current measurement along with sufficient computational power are adequate. However, optimal conditions for obtaining the most accurate results from the algorithm, such as the measurement time, need to be evaluated.

The equivalent circuit model holds the potential to enhance State of Charge (SoC) estimation precision, as it can provide battery voltage values independent of load current [13]. Additionally, the model can facilitate the monitoring of State of Health (SoH) and the detection of internal defects and abnormalities within battery systems, such as poor surface contact between cell connections.

ACKNOWLEDGMENT

This work was supported by the specific graduate research of the Brno University of Technology No. FEKT-S-23-8286.

References

- WAKEFIELD, Faith. Top 25 Solar Energy Statistics for 2024. Online. In: EcoWatch. 2023. Available at: <u>https://www.ecowatch.com/solar-energy-statistics.html</u>.
- [2] LEGRAND, N.; RAËL, S.; KNOSP, B.; HINAJE, M.; DESPREZ, P. et al., 2014. Including double-layer capacitance in lithium-ion battery mathematical models. Online. *Journal of Power Sources*. vol. 251, s. 370-378. Available at: <u>https://doi.org/10.1016/j.jpowsour.2013.11.044</u>.
- [3] HAEVERBEKE, Maxime Van; STOCK, Michiel and DE BAETS, Bernard, 2022. Equivalent Electrical Circuits and Their Use Across Electrochemical Impedance Spectroscopy Application Domains. Online. *IEEE Access.* vol. 10, s. 51363-51379. Available at: <u>https://doi.org/10.1109/ACCESS.2022.3174067</u>.
- [4] BCS-800 battery cycler series. Online. In: BioLogic. Available at: <u>https://www.biologic.net/products/bcs-800/</u>.
- [5] LEBEL, Félix-A.; MESSIER, Pascal; SARI, Ali and TROVÃO, João Pedro F., 2022. Lithium-ion cell equivalent circuit model identification by galvanostatic intermittent titration technique. Online. *Journal of Energy Storage*. vol. 54. Available at: <u>https://doi.org/10.1016/j.est.2022.105303</u>.
- [6] RL Circuit Time Constant / Universal Time Constant Curve, Electrical Academia. Online. In: Available at: <u>https://electricalacademia.com/basic-electrical/rl-circuit-time-constanttime-constant-of-rl-circuit/.</u>
- [7] www.gwl.eu Winston LFP060AHA. Online. In: GWL. Available at: https://files.gwl.eu/inc/_doc/attach/StoItem/1575/Winston_DS_LFP060 AHA.pdf.
- [8] MESSING, Marvin; SHOA, Tina and HABIBI, Saeid, 2021. Estimating battery state of health using electrochemical impedance spectroscopy and the relaxation effect. Online. *Journal of Energy Storage*. vol. 43. Available at: <u>https://doi.org/10.1016/j.est.2021.103210</u>.
- [9] EC-Lab®, 2023. Online. In: BioLogic. Available at: <u>https://www.biologic.net/topics/ec-lab/</u>.
- [10] SHIN, Donghoon; YOON, Beomjin and YOO, Seungryeol, 2021. Compensation Method for Estimating the State of Charge of Li-Polymer Batteries Using Multiple Long Short-Term Memory Networks Based on the Extended Kalman Filter. Online. *Energies*. vol. 14, č. 2. Available at: <u>https://doi.org/10.3390/en14020349</u>.
- [11] VAŠÍČEK, Jakub. Intelligent BMS for lithium batteries. Brno, 2022. Available at: <u>https://www.vut.cz/studenti/zav-prace/detail/142484</u>. Bakalářská práce. Faculty of Electrical Engineering and Communication, Brno University of Technology. Supervisor Petr Vyroubal.
- [12] Fminsearch: Find minimum of unconstrained multivariable function using derivative-free method. Online. In: MathWorks. Available at: <u>https://nl.mathworks.com/help/matlab/ref/fminsearch.html</u>.
- [13] HUA, Yin; XU, Min; LI, Mian; MA, Chengbin and ZHAO, Chen, 2015. Estimation of State of Charge for Two Types of Lithium-Ion Batteries by Nonlinear Predictive Filter for Electric Vehicles. Online. *Energies*. vol. 8, č. 5, s. 3556-3577. Available at: <u>https://doi.org/10.3390/en8053556</u>





TRANSFORMATIVE TECHNOLOGY IS IN OUR NATURE

DISCOVER THE JOB OPENINGS AT GARRETT MOTION R&D CENTER IN BRNO

Contact: kariera@garrettmotion.com









Corrosion potential analysis of iron-magnesium alloys

1st Silvia Bátorová Department of Electrical and Electronic Technology, FEEC Brno University of Technology Brno, Czech Republic xbator05@vutbr.cz

Abstract—This paper analyses the corrosion of Fe-Mg alloys meant to be used as prosthetic implants. Four samples with varying amounts of polystyrene added in were measured and analyzed using the Tafel extrapolation method of the polarization curves over the course of three months.

Keywords—corrosion, aqueous solutions, iron, magnesium

I. INTRODUCTION

The corrosion of metals inside the human body (*in vivo*) is an area of great potential. These metals would form prosthetic implants which could potentially dissolve within the body and therefore the need for invasive surgery in order to remove said implant is avoided. Currently, one of the most studied elements for this purpose is iron. It is often used either alone or as an alloy together with a different metal.

The elements chosen in this paper are iron with magnesium. Four samples of different iron-magnesium-polystyrene makeups were submerged in a 9g/L NaCl aqueous solution meant to roughly substitute the conditions *in vivo*. Their corrosion potentials were studied using the Tafel analysis and the results are documented in this paper.

II. CORROSION OF METALS IN AQUEOUS SOLUTIONS

A. Corrosion of iron

The corrosion of iron in aqueous solutions is described by two reactions, an anodic (1) and a cathodic reaction (2).

$$Fe \rightarrow Fe^{2+} + 2e^{-}$$
 (1)

$$2H_2O + O_2 + 4e^- \to 4OH^- \tag{2}$$

The rate of corrosion is controlled by the speed of the cathodic reaction, which has numerous forms based on the pH and oxygen saturation levels of the solution. Equation (2) concerns solutions with high pH and high oxygen saturation [1][2].

The rate of iron corrosion is increased under higher temperatures, lower pH, higher oxygen concentrations, higher microbial presence and higher velocities of the solution. An additional notable factor influencing the corrosion rate of iron is the concentration of sodium chloride in the solution. At first, corrosion rate increases until the maximum point of 3%, after which it decreases linearly, eventually dropping under the corrosion rate of distilled water [3]. 2nd Miroslav Zatloukal Department of Electrical and Electronic Technology, FEEC Brno University of Technology Brno, Czech Republic zatloukal@vut.cz

B. Corrosion of magnesium

Magnesium is a metal highly prone to corrosion, which is only exacerbated when it is present in alloys. The layers formed on the surface of the metal during corrosion do not protect from further corrosion [4]. The most major factors increasing its corrosion rates are higher concentrations of ions (mainly chloride and phosphate ions), lower pH levels and higher solution flow velocities [5].

III. MEASURING THE CORROSION POTENTIAL OF PREPARED PROBES

A. Probe characteristics

The measurements were performed on four probes named 1, 2, 3 and 4. The probes consisted of Fe-Mg alloys with varying amounts of polystyrene added in for increased porousness of the material. The specific contents of the probes are listed in the table below:

TABLE I. THE CHEMICAL CONTENTS OF THE FOUR PROBES

Probe number	Contents of probe	
1	9 g Fe + 1 g Mg + 1.5 g PS	
2	9 g Fe + 1 g Mg + 1 g PS	
3	9 g Fe + 1 g Mg + 0.5 g PS	
4	9 g Fe + 1 g Mg + 2 g PS	

These probes were sintered under high temperatures and after sintering were submerged in a 9 g/L NaCl aqueous solution. The samples were kept in a temperature-controlled container at 37 °C. The measurements were performed periodically over the course of three months, starting on 22nd of November and ending on 23rd of February.

B. Measuring method used

The method used to evaluate the corrosive processes of the probes was the Tafel extrapolation method of the polarization curves. The measurements were performed using a threeelectrode system of a referential saturated calomel electrode, a platinum counter electrode and a paraffin-impregnated graphite electrode (PIGE), which was the measuring electrode.

This work was supported by Grant "FEKT-S-23-8286 Materials and technology for electrical engineering V" from Brno University of Technology.

In the beginning of the measurement, the probes were taken out of their containers and rinsed in deionized water. After they were dry, samples used in the measurement were scratched off of the probes' surface at different places and the resulting powder was then pressed onto the top of the measuring electrode. The electrode was subsequently submerged into a 9 g/L NaCl solution and the measurement would begin with the use of μ AutoLab Type II potentiostat.

C. Results

The results of the measurement can be found in Table 2.

TABLE II.	RESULTS OF CORROSION POTENTIAL AND
CURREN	T DENSITY OF PROBES OVER TIME

No. of	Date	No. of	Ecorr	icorr
meas.		sample	[V]	[µA/cm ²]
		1	-0.631	5.55
1	22.11.	2	-0.547	4.12
	2023	3	-0.627	5.09
		4	-0.63	3.22
		1	-0.605	4.35
2	29.11.	2	-0.604	3.51
	2023	3	-0.63	3.78
		4	-0.616	2.69
		1	-0.639	3.13
3	6.12.	2	-0.598	3.73
	2023	3	-0.513	3.15
		4	-0.554	3.32
		1	-0.606	3.41
4	14.12.	2	-0.624	3.96
	2023	3	-0.639	5.06
		4	-0.637	3.98
		1	-0.615	4.90
5	19.1.	2	-0.501	3.43
	2024	3	-0.65	4.02
		4	-0.648	3.77
		1	-0.583	3.79
6	13.2.	2	-0.572	3.32
	2024	3	-0.614	3.89
		4	-0.631	4.75
		1	-0.567	3.76
7	23.2.	2	-0.512	3.52
	2024	3	-0.594	3.53
		4	-0.501	3.48

The deviations between specific measurements are depicted in the 4 graphs below. The graphs show how the values of corrosion current density, i_{corr} , and corrosion potential, E_{corr} , changed with time.



Fig. 1. Change of characteristics of probe 1 (y – function equation, s – standard deviation of the fit, v – variation coefficient)



Fig. 2. Change of characteristics of probe 2 (y – function equation, s – standard deviation of the fit, v – variation coefficient)



Fig. 3. Change of characteristics of probe 3 (y – function equation, s – standard deviation of the fit, v – variation coefficient)



Fig. 4. Change of characteristics of probe 4 (y – function equation, s – standard deviation of the fit, v – variation coefficient)

IV. DISCUSSION

The graphs show that the values of corrosion current density kept slightly decreasing with time, while the values of corrosion potential kept slightly increasing. The one exception is the probe 4, which has the biggest amount of polystyrene added in. For probe 4, the corrosion current density increased, which could be due to higher porousness of the probe. The results suggest that the rate of corrosion is slowing down with time. This could be due to changes in the upper layers of the probes or due to presence of corrosion products in the solution. Another effect taking place could be a different rate of corrosion of iron and magnesium. Overall, the behavior of the corrosive processes of the probes is a very complex matter and will require further detailed study.

ACKNOWLEDGMENT

This work was supported by Grant "FEKT-S-23-8286 Materials and technology for electrical engineering V" from Brno University of Technology.

References

- T. M. Devine, "Corrosion of iron-base alloys," in *Embrittlement of* engineering alloys, vol. 25, C. L. Briant and S. K. Banerji, Eds., New York, NY, USA: Academic Press, Inc., 1983, pp.201–234.
- [2] M. Sedlaříková, M. Zatloukal, J. Kuchařík, P. Čudek, G. Fafilek, and E. Doleželová, "Corrosion processes of sintered materials based on Fe," *Journal of Physics: Conference Series*, vol. 2382, Nov. 2022, doi: https://doi.org/10.1088/1742-6596/2382/1/012019.
- [3] R. W. Revie and H. H. Uhlig, "Iron and steel," in *Corrosion and corrosion control: An introduction to corrosion science and engineering*, 4th ed. Hoboken, NJ, USA: John Wiley & Sons, Inc., 2008, pp. 115–133.
- [4] A. Atrens et al., "Review of Mg alloy corrosion rates," Journal of Magnesium and Alloys, vol. 8, no. 4, pp. 989–998, Dec. 2020, doi: https://doi.org/10.1016/j.jma.2020.08.002.
- [5] A. Atrens, S. Johnston, Z. Shi, and M. S. Dargusch, "Viewpoint -Understanding Mg corrosion in the body for biodegradable medical implants," *Scripta Materialia*, vol. 154, pp. 92–100, Sep. 2018, doi: https://doi.org/10.1016/j.scriptamat.2018.05.021.

LCA of different types of cars in Czech Republic

Bc. Václav Kutnar Faculty of Electrical Engineering and Communication Brno University of Technology Brno, Czech Republic xkutna02@vutbr.cz

To reduce greenhouse gas emissions, countries have focused on (among other things) automobile transport. However, there are different car powertrains, and each produces a different amount of emissions from different parts of their lifecycles, which also vary from state to state. Knowing exact, comparable values is necessary for deciding how to reduce these emissions most efficiently and economically possible. Existing studies often focus on only one part of the life cycle and/or are not localized to a particular state. This study analyzes the entire life cycle of a vehicle and its fuel or energy for a vehicle manufactured, operated, and recycled in the Czech Republic. It draws on existing studies, supplemented by data from state agencies, vehicle manufacturers, and our calculations. It was found that in the Czech Republic, compared to a vehicle with a petrol internal combustion engine (ICEVg), an electric vehicle (BEV) has 33% lower total GHG emissions, and a vehicle with hydrogen fuel cells (FCEV) 19% lower. However, an internal combustion engine vehicle has lower emissions from its production and as a result, its cumulative emissions are lower than an electric vehicle until 70,000 km driven and a hydrogen fuel cell vehicle until 111,000 km.

Keywords - car, greenhouse gas, lifecycle, LCA, well-to-wheel

I. INTRODUCTION

Excessive greenhouse gas emissions (GHGs) are key driver of global warming. While these gases are essential for atmospheric and climatic functions, their overabundance can disrupt climate processes. The most prevalent and impactful greenhouse gas is water vapor, followed by CO_2 , CH_4 , NOx, CFCs, and ozone. The EU contributes 7% to global GHG emissions, with transport being the second largest emitter (22% of EU's total), while road transport accounting for 43% of transport emissions.

LCA - LIFE CYCLE ANALYSIS

In order to assess the environmental impact of a product as accurately as possible, all stages of its life cycle must be analysed accordingly. The life cycle assessment method makes it possible to identify, quantify and assess the impact of the entire life cycle in all relevant areas, such as gas emissions, use of natural resources, waste production, etc. [2]. The aim of this work is to compare the total greenhouse gas emissions produced over the entire life cycle of vehicles with different types of propulsion. The following factors will be considered: Ing. Kamil Jaššo, Ph.D. Faculty of Electrical Engineering and Communication Brno University of Technology Faculty of military technology University of Defense Brno, Czech Republic xjasso00@vutbr.cz, kamil.jasso@unob.cz

- Vehicle production
- Fuel life cycle
- Vehicle recycling
- The powertrains analysed are ICEV, BEV and FCEV.
- Production, operation and recycling will take place in the Czech Republic, including production of fossil fuels, electricity and hydrogen
- Future development of the recycling industry and the expected reduction of the electricity or hydrogen emission factor

What is not included in this analysis:

- Consumables and minor repairs
- Replacing the traction battery of BEV or FCEV
- Effect of driving style or extreme weather conditions

II. METHODOLOGY

Every vehicle is made up of a large number of components. Ideally, an assessment would be made of the production of each of these components, but due to limitations imposed primarily by the lack of adequately accurate data, most of the work to date has made some approximations. The most common case is that of splitting the vehicle into glider and powertrain [3][4][5][6][7]. Glider is usually considered identical for all types of propulsion and only powertrain varies accordingly. There are other, more thorough methods with analysing almost every car part and manufacturing procedure in detail. Yet, they are not very common, and often are based on accurate data directly from the manufacturer or use their own simulated cars for which they have all the necessary data [8][9][10].

For each vehicle part, process or energy/fuel, an emission factor was determined, based on evaluation of already completed studies and/or author's own calculations.

A. Glider production

Existing LCAs have taken different approaches to the determination of CO_{2eq} emissions from glider production. Various already existing studies were evaluated [3][6][7][9] [11][12][13][14][15]. Values from Buberger et al. (2022) [16] were used in this paper, namely because they utilize data directly from European car manufacturer located in country neighboring Czechia.

This study presents one average value of emission factor EF_{glider} for almost all types of propulsion is used in the calculations,

namely 4.56 kg $CO_{2 eq}$./kg_{veh}. Only exception is FCEV vehicle, with value of 5.34 kg $CO_{2 eq}$ /kg_{veh}.

B. Battery production

The data describing the environmental impact of EV battery production varies widely across different studies [6][7][9][10] [11][12][13][14][17][18][19]. Values in study by Emilsson, Dallöf (2017) for the Swedish Environmental Research Institute IVL was chosen for this paper [19]. It uses its own research supplemented by data from other studies, presenting a value of 61 to 106 kg CO_{2 eq}/kWh _{batt}, with the exact value based on the carbon intensity of the energy used in battery production. Low value corresponds to the use of electricity with a carbon intensity of 0 g CO_{2 eq}/kWh, whereas high value corresponds to 1 kg CO_{2 eq}/kWh. In the case of Czech Republic, this translates to EF_{batt} 83.5 kg CO_{2 eq}/kWh_{batt}.

C. Fossil fuels production

For fossil fuels production in the Czech Republic, paper by the Transport Research Centre made for the Ministry of the Environment on the LCA of fossil motor fuels and biofuels was used [20]. According to calculations, burning 1 liter of petrol generates 2.629 kg of CO_{2eq} , while 11% is from production and transportation and 89% is from combustion.

D. Electricity production

In Czech Republic, 44% of electricity is made from coal, over 40% from nuclear, 9% from natural gas and 7% from renewable energy sources - photovoltaic, hydro and wind. Losses in generation, transmission and transport networks represent approximately 12% and must be considered in the calculations [21]. Based on these data, an emission factor EF_{el} for the Czech energy mix was calculated.

$$EF_{el} = \frac{GHG \ produced}{Production \cdot efficiency} = (1)$$

$$\frac{37,819,476,500}{84,527,500 \cdot 0.88} = 508.43[gCO_{2eq}/kWh]$$

The Czech Ministry of Industry and Trade states that the emission factor of the Czech energy mix has been decreasing at an average rate of 1.7% per year from 1990 to 2022 [22]. Given the EU commitments to reduce GHG emissions and the goal to be emission neutral in 2050, the planned phase-out of coal-fired power plants and the planned construction of up to 4 new nuclear units, this trend can be expected to continue.

E. Hydrogen production

There are several ways to produce hydrogen. These can be fossil-based processes such as steam reforming or pyrolysis [22], or hydrogen can be produced by electrolysis. Here we then distinguish what electricity has been used for its production. Generally, hydrogen is referred to by a "colour", which indicates the method of its production and its approximate emission factor [23][24][25][26][27]. The Czech Hydrogen Concept prepared by the Ministry of Industry and Trade indicates a total emission factor of 16.356 kg CO₂/kg_H, with majority being produced by pyrolysis [28]. It should be noted,

however, that this refers only to carbon dioxide, not all greenhouse gases. Regarding future developments, the strategy does not foresee significant changes in the technologies used until 2036-2039.

Losses occurring during pressurization and transportation of hydrogen correspond to approximately 2.43 kWh/kg_H of electricity consumed. The tanker consumption results in an average of 10.8 gCO₂/kg_H. These factors EF of Czech hydrogen to 17.602 kgCO₂/kg_H. [20][29][30][31]

F. End of the life

The average lifetime of a vehicle in the EU ranges from 8 to 35 years, with a median of 21.7 years, compared to 18.1 years in Western European countries and 28.4 years in Eastern European countries [32]. Data on average mileage before recycling then varies, ranging between 205,000 [33] and 250,000 km in most cases [15][34]. The 2021 ICCT report gives an average vehicle life expectancy of 17 to 18 years in Germany, 19 years in France and 20 years in Poland, with an average mileage of 13,500 km per year (KM_{annual}) for a lower mid-range vehicle [35].

Among already existing studies, [5][6][9][16][17][36][37] [38][39], two different approaches occurs. Some studies, like the Swedish IVL [19], determines emission factor from recycling (EF_{recyc}). on the basis of energy consumed during the recycling process. Other works use a credit system resulting in negative EF_{recyc} . From the work of Buberger et al. [16], emission factor of glider recycling (R_{glider}) 2.93 kg CO_{2 eq}./kg_{veh} for vehicle without battery was used in this study. As for traction battery, approach of Mohr et al. [40] was used, where recycling bonus is presented as percentual reduction of manufacturing emissions. From this paper, battery recycling coefficient (R_{batt}) of 29% was used, for it was considered more accurate and better reflecting future development in recycling industry and is in accordance with other studies [41][42].

III. ANALYSED VEHICLES

The vehicles selected to represent BEVs and ICEVs were the upper mid-range (D-segment) BMW i4 eDrive35 (BEV) and the petrol-powered BMW 430i Gran Coupe (ICEVg). These vehicles were selected because they share a common platform,

	BMW i4 eDrive35	BMW 430i Coupe	Toyota Mirai MKII
Fuel	Electricity	Gasoline	Hydrogen
Power	210 kW	180kW	134 kW
Curb weight	2065 kg	1733 kg	2415 kg
(m _{veh})	2005 Kg	1755 Kg	2415 Kg
Batt. capacity	70,2/67 kWh		1,24 kWh,
(Q _{batt})	NMC 811	-	NiMH
Batt.weight	550 kg		44,5 kg
(m _{batt})	550 Kg	-	
Consumption	18,7-	7,6-	0.81 kg/100 km
(KPLaver)	15,8 kWh/100 km	6,8 l/100 km	0,01 Kg/100 Kill
Range	406-482 km	771-867 km	650 km

TABLE I. ANALYZED VEHICLES [43][44][45]

a number of parts and part of the supply chain. Due to the low penetration of FCEVs, the Toyota Mirai MK2 was selected as the vehicle most similar to the BMW 4 Series. It is also the only mid-size sedan/coupe body FCEV currently sold. At the time, no European car company offered a hydrogen car.

A. Charging losses

In the case of ICEV, filling of the tank usually happens with negligible losses. That is not the case with BEV, where amount of energy charged to the battery is lesser that amount of energy delivered to the charger. These losses are called charging losses (ChL). They consist of conversion losses (converting AC to DC), heat losses and losses occurring while transfering electrical energy into battery. Average charging efficiency ranges from 78 to 91%, (15 - 80% SoC) [46][47][48][49]. Another major factor is temperature, as very low temperatures can reduce charging efficiency down to 59%. Charging beyond 80 % SoC increases charging losses almost twice the normal value [50]. On the other hand, high temperatures seem to have negligible effect on losses, as does using "fast" DC chargers with power exceeding 80 kW[50]. For this study, an average charging losses value of 15% was chosen.

B. Battery self-discharge

The EV spontaneously loses power even when not in use. This is caused by self-discharge (SD) of the battery itself and by the consumption of on-board appliances. The exact value varies, but common values of SD are 0.5 to 2% per 24 hours in optimal climatic conditions, 2 to 4% in sub-zero temperatures, and in some cases 5% or even more (sentry mode in Tesla vehicles, preheating in winter, etc.)[51][52][53]. Conversely, these losses can be reduced by putting the vehicle into hibernation, where these losses can be reduced to as much as 1-2% per month. It should be noted, however, that there are scarcely any scientific papers on this issue and therefore these values are taken from manufacturer's vehicle manual or non-scientific literature. An average self-discharge value of 2% was chosen for this study.

IV. CALCULATION

Calculation of lifecycle emissions is divided in to three modules – production, operation and end-of-the-life.

BEV production module:

$$E_{prod} = \left((m_{veh} - m_{batt}) \cdot EF_{glider} \right) + (Q_{batt} \cdot EF_{batt})$$
(2)

BEV operation module (first year):

$$E_{oper} = \begin{pmatrix} (KPL_{aver} \cdot KM_{annual}) + \\ (Q_{batt} \cdot SD \cdot 365) \end{pmatrix} \cdot ChL \cdot EF_{el}$$
(3)

BEV EoL module:

$$E_{EoL} = \left(R_{glider} \cdot (m_{veh} - m_{batt}) \right) + \left(\left(Q_{batt} \cdot EF_{batt} \right) \cdot R_{batt} \right)$$
(4)



Fig. 1. Total lifecycle emissions

V. CONCLUSION

Total lifetime GHG emissions are lowest for BEVs, by 33% compared to ICEVs. In the case of FCEVs, the reduction is 19%. However, these values change significantly during lifetimes of cars. BEVs and FCEVs have higher production emissions - 61% for BEVs and 63% for FCEVs compared to ICEVg. Thus, the ICEVg has lower cumulative emissions in the first years of operation, and break-even occurs after approximately five years / 70,000 km of operation in the case of the BEV or eight years / 111,000 km in the case of the FCEV. Fig. 1 shows the evolution of GHG emissions over the entire vehicle life cycle, including the production and recycling phases.

Of the three types of powertrains, the emissions from vehicle operation are highest for the ICEVg, followed by the FCEV and lowest for the BEV. For all vehicle types, these emissions can be reduced by reducing vehicle fuel or energy consumption. In addition, for BEVs these emissions can be further lowered by reducing the emission factor of the electricity used and for FCEVs by using hydrogen production methods with lower emission factor. In the case of ICEVs, the production process and the chemical composition of the gasoline is fixed, and further reductions in GHG emissions from operation can only be achieved by changing the fuel itself.

It also needs to be noted, that this study evaluates only GHG emissions. Other environmental impacts, such as acidification, abiotic depletion, ecological toxicity potential and others were omitted from this work.

VI. ACKNOWLEDGEMENTS

"This work was supported by the specific graduate research of the Brno University of Technology No. FEKT-S-23-8286 and by the institutional support of the Ministry of Defence of the Czech Republic (VAROPS)."

REFERENCES

- [1] Climate change. Meteorological dictionary [online]. Czech Meteorological Society [cit. 2022-01-11]. Available from: Meteorological Dictionary (cmes.cz)J. Clerk Maxwell, A Treatise on Electricity and Magnetism, 3rd ed., vol. 2. Oxford: Clarendon, 1892, pp.68–73.
- [2] MURALIKRISHNA, Iyyanki V. a Valli MANICKAM. Life Cycle Assessment. Environmental Management [online]. Elsevier, 2017, 2017, 57-75 [cit. 2023-03-30]. ISBN 9780128119891. Dostupné z: doi:10.1016/B978-0-12-811989-1.00005-1
- [3] Determining the environmental impacts of conventional and alternatively fuelled vehicles through LCA: Final Report for the European Commission, DG Climate Action [online]. Luxembourg: Publications Office of the European Union, 2020 [cit. 2023-03-12]. ISBN 978-92-76-20301-8. Dostupné z: <u>https://climate.ec.europa.eu/system/files/2020-09/2020 study main report en.pdfpa.eu%2Fsystem%2Ffiles%2F2020-09%2F2020 study main report en.pdf&usg=AOvVaw3fF4djld86tWc L4q4uZUpc</u>
- [4] DE SOUZA, Lidiane La Picirelli, Electo Eduardo Silva LORA, José Carlos Escobar PALACIO, Mateus Henrique ROCHA, Maria Luiza Grillo RENÓ a Osvaldo José VENTURINI. Comparative environmental life cycle assessment of conventional vehicles with different fuel options, plug-in hybrid and electric vehicles for a sustainable transportation system in Brazil. *Journal of Cleaner Production* [online]. 2018, 203, 444-468 [cit. 2023-03-12]. ISSN 09596526. Dostupné z: doi:10.1016/j.jclepro.2018.08.236
- [5] TAGLIAFERRI, Carla, Sara EVANGELISTI, Federica ACCONCIA, Teresa DOMENECH, Paul EKINS, Diego BARLETTA a Paola LETTIERI. Life cycle assessment of future electric and hybrid vehicles: A cradle-to-grave systems engineering approach. *Chemical Engineering Research and Design* [online]. 2016, **112**, 298-309 [cit. 2023-03-12]. ISSN 02638762. Dostupné z: doi:10.1016/j.cherd.2016.07.003
- [6] LOMBARDI, Lidia, Laura TRIBIOLI, Raffaello COZZOLINO a Gino BELLA. Comparative environmental assessment of conventional, electric, hybrid, and fuel cell powertrains based on LCA. *The International Journal of Life Cycle Assessment* [online]. 2017, 22(12), 1989-2006 [cit. 2023-03-12]. ISSN 0948-3349. Dostupné z: doi:10.1007/s11367-017-1294-y
- [7] HELMERS, Eckard, Johannes DIETZ a Martin WEISS. Sensitivity Analysis in the Life-Cycle Assessment of Electric vs. Combustion Engine Cars under Approximate Real-World Conditions. *Sustainability* [online]. 2020, **12**(3) [cit. 2023-03-12]. ISSN 2071-1050. Dostupné z: doi:10.3390/su12031241
- [8] DELOGU, Massimo, Francesco DEL PERO, Laura ZANCHI, Marcos IERIDES, Violeta FERNANDEZ, Kristian SEIDEL, Dinesh THIRUNAVUKKARASU a Thilo BEIN. Lightweight Automobiles ALLIANCE Project: First Results of Environmental and Economic Assessment from a Life-Cycle Perspective [online]. 2018-05-30, - [cit. 2023-03-12]. Dostupné z: doi:10.4271/2018-37-0027
- [9] PERO, Francesco Del, Massimo DELOGU a Marco PIERINI. Life Cycle Assessment in the automotive sector: a comparative case study of Internal Combustion Engine (ICE) and electric car. Procedia Structural Integrity [online]. 2018, 12, 521-537 [cit. 2023-03-12]. ISSN 24523216. Dostupné z: doi:10.1016/j.prostr.2018.11.066
- [10] KARAASLAN, Enes, Yang ZHAO a Omer TATARI. Comparative life cycle assessment of sport utility vehicles with different fuel options. The International Journal of Life Cycle Assessment [online]. 2018, 23(2), 333-347 [cit. 2023-03-12]. ISSN 0948-3349. Dostupné z: doi:10.1007/s11367-017-1315-x
- [11] KAWAMOTO, Ryuji; MOCHIZUKI, Hideo; MORIGUCHI, Yoshihisa; NAKANO, Takahiro; MOTOHASHI, Masayuki et al. Estimation of CO2 Emissions of Internal Combustion Engine Vehicle and Battery Electric Vehicle Using LCA. Online. *Sustainability*. 2019, roč. 11, č. 9. ISSN 2071-1050. Dostupné z: <u>https://doi.org/10.3390/su11092690</u>. [cit. 2023-10-05].
- [12] QIAO, Qinyu; ZHAO, Fuquan; LIU, Zongwei; JIANG, Shuhua a HAO, Han. Comparative Study on Life Cycle CO 2 Emissions from the Production of Electric and Conventional Vehicles in China. Online. Energy Procedia. 2017, roč. 105, s. 3584-3595. ISSN 18766102.

Dostupné z: https://doi.org/10.1016/j.egypro.2017.03.827. [cit. 2023-10-05].

- [13] Hao, H., Qiao, Q., Liu, Z. et al. Comparing the life cycle Greenhouse Gas emissions from vehicle production in China and the USA: implications for targeting the reduction opportunities. *Clean Techn Environ Policy* 19, 1509–1522 (2017). <u>https://doi.org/10.1007/s10098-016-1325-6</u>
- [14] SATO, Fernando Enzo Kenta a NAKATA, Toshihiko. Energy Consumption Analysis for Vehicle Production through a Material Flow Approach. Online. Energies. 2020, roč. 13, č. 9. ISSN 1996-1073. Dostupné z: <u>https://doi.org/10.3390/en13092396</u>. [cit. 2023-10-05].
- [15] European Commission, Joint Research Centre, Castellani, V., Fantoni, M., Cristòbal, J. et al., Consumer footprint – Basket of products indicator on mobility, Publications Office, 2017, https://data.europa.eu/doi/10.2760/539712
- [16] BUBERGER, Johannes, Anton KERSTEN, Manuel KUDER, Richard ECKERLE, Thomas WEYH a Torbjörn THIRINGER. Total CO2equivalent life-cycle emissions from commercially available passenger cars. *Renewable and Sustainable Energy Reviews* [online]. 2022, **159** [cit. 2023-10-05]. ISSN 13640321. Dostupné z: doi:10.1016/j.rser.2022.112158
- [17] VAN MIERLO, Joeri; MESSAGIE, Maarten a RANGARAJU, Surendraprabu. Comparative environmental assessment of alternative fueled vehicles using a life cycle assessment. Online. Transportation Research Procedia. 2017, roč. 25, s. 3435-3445. ISSN 23521465. Dostupné z: <u>https://doi.org/10.1016/j.trpro.2017.05.244</u>. [cit. 2023-10-06].
- [18] PETERS, Jens F.; BAUMANN, Manuel; ZIMMERMANN, Benedikt; BRAUN, Jessica a WEIL, Marcel. The environmental impact of Li-Ion batteries and the role of key parameters – A review. Online. *Renewable* and Sustainable Energy Reviews. 2017, roč. 67, s. 491-506. ISSN 13640321. Dostupné z: <u>https://doi.org/10.1016/j.rser.2016.08.039</u>. [cit. 2023-10-09].
- [19] EMILSSON, Erik a DALLÖF, Lisbeth. Lithium-Ion Vehicle Battery Production: Status 2019 on Energy Use, CO2 Emissions, Use of Metals, Products Environmental Footprint, and Recycling. Swedish Environmental Research Institute IVL, 2017.
- [20] JEDLIČKA, Jiří. Life cycle analysis of fossil motor fuels and biofuels for the development of concept documents for the introduction of CO2 tax in the field of mobile sources of pollution [online]. Ministry of the Environment; Transport Research Centre, 2010 [cited 2021-11-21]. Available from: <u>https://www.kverulant.org/upload/kc/files/bioomyl_SPII4i1_33_07_extr</u> <u>akt.pdf</u>
- [21] QUARTERLY REPORT ON THE OPERATION OF THE ELECTRICITY SYSTEM OF THE CZECH REPUBLIC FOR IV. QUARTER 2022 [online]. Prague: Energy Regulatory Office, 2023 [cit. 2023-04-05]. Available from: <u>https://www.eru.cz/ctvrtletni-zprava-oprovozu-elektrizacni-soustavy-cr-za-iv-ctvrtleti-2022</u>
- [22] MINISTRY OF INDUSTRY AND TRADE. CO2 emission factor from electricity generation for the years 2010-2022. Online. 2023. Available from: https://www.mpo.cz/cz/energetika/statistika/elektrina-ateplo/emisni-faktor-co2-z-vyroby-elektriny-za-leta-2010_2022--273197/. [cited 2024-02-24].
- [23] YOUNAS, Muhammad, Sumeer SHAFIQUE, Ainy HAFEEZ, Fahad JAVED a Fahad REHMAN. An Overview of Hydrogen Production: Current Status, Potential, and Challenges. *Fuel* [online]. 2022, **316** [cit. 2023-04-11]. ISSN 00162361. Dostupné z: doi:10.1016/j.fuel.2022.123317
- [24] Basic information on hydrogen [online]. HYTEP Czech Hydrogen Technology Platform, 2023 [cited 2023-04-11]. Available from:: https://www.hytep.cz/o-vodiku/ve-zkratce
- [25] The many colors and applications of hydrogen: Four things you need to know about hydrogen use in industry [online]. ELSEVIER, 2022 [cit. 2023-04-11]. Dostupné z: https://www.elsevier.com/connect/the-manycolors-and-applications-of-hydrogen
- [26] SOAM, Shveta a BÖRJESSON, Pål. Considerations on Potentials, Greenhouse Gas, and Energy Performance of Biofuels Based on Forest Residues for Heavy-Duty Road Transport in Sweden. Online. Energies. 2020, roč. 13, č. 24, s. 21. ISSN 1996-1073. Dostupné z: https://doi.org/10.3390/en13246701. [cit. 2023-10-26].

- [27] IEA, Comparison of the emissions intensity of different hydrogen production routes, 2021, IEA, Paris https://www.iea.org/data-andstatistics/charts/comparison-of-the-emissions-intensity-of-differenthydrogen-production-routes-2021, IEA.
- [28] MINISTRY OF INDUSTRY AND TRADE. HYDROGEN STRATEGY OF THE CZECH REPUBLIC. Online. 2021. Available from: https://www.mpo.cz/assets/cz/rozcestnik/pro-media/tiskovezpravy/2021/7/3VL-03-Vodikova-strategie_v030b.pdf. [cited 2024-02-24].
- [29] BAUER, Artur, Thomas MAYER, Malte SEMMEL, Martin Alberto GUERRERO MORALES and Joerg WIND. Energetic evaluation of hydrogen refueling stations with liquid or gaseous stored hydrogen. International Journal of Hydrogen Energy [online]. 2019, 44(13), 6795-6812 [cited 2023-04-11]. ISSN 03603199. Available from: doi:10.1016/j.ijhydene.2019.01.087
- [30] How are hydrogen filling stations supplied with hydrogen? [online]. Prague 9: APT, 2021 [cited 2023-04-11]. Available from: https://www.apt.cz/q11
- [31] What are the basic functional parts of a hydrogen filling station for road vehicles? [online]. Prague 9: APT, 2021 [cited 2023-04-11]. Available from: <u>https://www.apt.cz/newspaper/q1z</u>
- [32] HELD, Maximilian, Nicolas ROSAT, Gil GEORGES, Hermann PENGG a Konstantinos BOULOUCHOS. Lifespans of passenger cars in Europe: empirical modelling of fleet turnover dynamics. *European Transport Research Review* [online]. 2021, **13**(1) [cit. 2023-04-25]. ISSN 1867-0717. Dostupné z: doi:10.1186/s12544-020-00464-0
- [33] WEYMAR, Elisabeth a Matthias FINKBEINER. Statistical analysis of empirical lifetime mileage data for automotive LCA. *The International Journal of Life Cycle Assessment* [online]. 2016, 21(2), 215-223 [cit. 2023-04-25]. ISSN 0948-3349. Dostupné z: doi:10.1007/s11367-015-1020-6
- [34] European LCA Results: European Life Cycle Assessment Results and Fact Sheets [online]. 2023 [cit. 2023-04-25]. Dostupné z: <u>https://www.greenncap.com/european-lca-results/</u>
- [35] How Many Miles Does a Car Last? [online]. Car and Driver [cit. 2023-04-25]. Dostupné z: https://www.caranddriver.com/research/a32758625/how-many-milesdoes-a-car-last/
- [36] HAO, Han, Qinyu QIAO, Zongwei LIU a Fuquan ZHAO. Impact of recycling on energy consumption and greenhouse gas emissions from electric vehicle production: The China 2025 case. *Resources*, *Conservation and Recycling* [online]. 2017, **122** [cit. 2023-04-25]. ISSN 09213449. Dostupné z: doi:10.1016/j.resconrec.2017.02.005
- [37] KOROMA, Michael Samsu, Nils BROWN, Giuseppe CARDELLINI a Maarten MESSAGIE. Prospective Environmental Impacts of Passenger Cars under Different Energy and Steel Production Scenarios. *Energies* [online]. 2020, 13(23) [cit. 2023-04-25]. ISSN 1996-1073. Dostupné z: doi:10.3390/en13236236
- [38] ZENG, Dan, Yan DONG, Huajun CAO, Yuke LI, Jia WANG, Zhenbiao LI a Michael Zwicky HAUSCHILD. Are the electric vehicles more sustainable than the conventional ones? Influences of the assumptions and modeling approaches in the case of typical cars in China. *Resources, Conservation and Recycling* [online]. 2021, **167** [cit. 2023-04-25]. ISSN 09213449. Dostupné z: doi:10.1016/j.resconrec.2020.105210
- [39] SATO, Fernando Enzo Kenta a NAKATA, Toshihiko. Energy Consumption Analysis for Vehicle Production through a Material Flow Approach. Online. Energies. 2020, roč. 13, č. 9. ISSN 1996-1073. Dostupné z: <u>https://doi.org/10.3390/en13092396</u>. [cit. 2023-12-14].
- [40] MOHR, Marit; PETERS, Jens F.; BAUMANN, Manuel a WEIL, Marcel. Toward a cell-chemistry specific life cycle assessment of lithium-ion

battery recycling processes. Online. Journal of Industrial Ecology. 2020, roč. 24, č. 6, s. 1310-1322. ISSN 1088-1980. Dostupné z: https://doi.org/10.1111/jiec.13021. [cit. 2023-12-15].

- [41] SUN, Xin; LUO, Xiaoli; ZHANG, Zhan; MENG, Fanran a YANG, Jianxin. Life cycle assessment of lithium nickel cobalt manganese oxide (NCM) batteries for electric passenger vehicles. Online. Journal of Cleaner Production. 2020, roč. 273. ISSN 09596526. Dostupné z: <u>https://doi.org/10.1016/j.jclepro.2020.123006</u>. [cit. 2023-12-13].
- [42] CIEZ, Rebecca E. a WHITACRE, J. F. Examining different recycling processes for lithium-ion batteries. Online. Nature Sustainability. 2019, roč. 2, č. 2, s. 148-156. ISSN 2398-9629. Dostupné z: <u>https://doi.org/10.1038/s41893-019-0222-5</u>. [cit. 2023-12-15].
- [43] PIELECHA, Ireneusz, Andrzej SZAŁEK a Grzegorz TCHOREK. Two Generations of Hydrogen Powertrain—An Analysis of the Operational Indicators in Real Driving Conditions (RDC). *Energies* [online]. 2022, **15**(13) [cit. 2023-03-13]. ISSN 1996-1073. Dostupné z: doi:10.3390/en15134734
- [44] BMW GMBH. BMW. Online. 2023. Available from: https://www.bmw.cz/cs/all-models.html. [cited 2023-10-17].
- [45] TOYOTA. Mirai mk2 tech specs. Online. 2021. Available from: https://media.toyota.co.uk/wp-content/uploads/sites/5/pdf/220203M-Mirai-Tech-Spec.pdf. [cited 2023-10-17].
- [46] SEVDARI, Kristian; CALEARO, Lisa; BAKKEN, Bjørn Harald; ANDERSEN, Peter Bach a MARINELLI, Mattia. Experimental validation of onboard electric vehicle chargers to improve the efficiency of smart charging operation. Online. Sustainable Energy Technologies and Assessments. 2023, roč. 60. ISSN 22131388. Dostupné z: https://doi.org/10.1016/j.seta.2023.103512. [cit. 2024-02-27].
- [47] KOSTOPOULOS, Emmanouil D.; SPYROPOULOS, George C. a KALDELLIS, John K. Real-world study for the optimal charging of electric vehicles. Online. Energy Reports. 2020, roč. 6, s. 418-426. ISSN 23524847. Dostupné z: https://doi.org/10.1016/j.egyr.2019.12.008. [cit. 2024-02-27].
- [48] APOSTOLAKI-IOSIFIDOU, Elpiniki; CODANI, Paul a KEMPTON, Willett. Measurement of power loss during electric vehicle charging and discharging. Online. Energy. 2017, roč. 127, s. 730-742. ISSN 03605442. Dostupné z: https://doi.org/10.1016/j.energy.2017.03.015. [cit. 2024-02-27].
- [49] REICK, Benedikt; KONZEPT, Anja; KAUFMANN, André; STETTER, Ralf a ENGELMANN, Danilo. Influence of Charging Losses on Energy Consumption and CO2 Emissions of Battery-Electric Vehicles. Online. Vehicles. 2021, roč. 3, č. 4, s. 736-748. ISSN 2624-8921. Dostupné z: https://doi.org/10.3390/vehicles3040043. [cit. 2024-02-27].
- [50] TRENTADUE, Germana; LUCAS, Alexandre; OTURA, Marcos; PLIAKOSTATHIS, Konstantinos; ZANNI, Marco et al. Evaluation of Fast Charging Efficiency under Extreme Temperatures. Online. Energies. 2018, roč. 11, č. 8. ISSN 1996-1073. Dostupné z: https://doi.org/10.3390/en11081937. [cit. 2024-02-27].
- [51] TESLA. Tesla model 3 owner's manual. Online. 2024. Dostupné z: https://www.tesla.com/ownersmanual/model3/en_us/GUID-7FE78D73-0A17-47C4-B21B-54F641FFAEF4.html. [cit. 2024-02-28].
- [52] MAXWELL, Alfred. How long can an electric car sit without charging? Online. 2022. Dostupné z: https://topcharger.co.uk/how-long-can-anelectric-car-sit-without-charging/#:~:text=in%20the%20process.-,Energy%20loss%20in%20electric%20car%20batteries,how%20cold%2 Otemperatures%20affect%20batteries).. [cit. 2024-02-28].
- [53] NEIGHBOR. How Long Can an Electric Car Sit Without Charging? Online. 2024. Dostupné z: https://www.neighbor.com/storage-blog/howlong-can-electric-car-sit-without-charging/. [cit. 2024-02-28].

Calculation of Bone Mineral Density from Dual-energy CT and its Application on Patient with Multiple Myeloma

Michal Nohel Department of Biomedical Engineering Brno University of Technology, FEEC Brno, Czech Republic xnohel04@vutbr.cz Jiri Chmelik Department of Biomedical Engineering Brno University of Technology, FEEC Brno, Czech Republic chmelikj@vutbr.cz

Abstract—This article presents the results of the calculation of bone mineral density in the spine of a patient with multiple myeloma and lytic lesions. The findings indicate that the average value for a healthy vertebra falls within the physiological range. In the case of the patient with myeloma, a low value was measured in the area of the lytic lesion, suggesting a high risk of pathological fractures. The research also revealed lower values in areas without lytic lesions. These results emphasize the importance of precise evaluation of mineral density in the diagnosis of spinal diseases.

Index Terms-BMD, multiple myeloma, dual-energy CT, spine

I. INTRODUCTION

Multiple myeloma (MM) is a hematologic disorder marked by the clonal expansion of plasma cells within the bone marrow. It is commonly linked to skeletal issues, with osteolytic bone lesions serving as a key diagnostic indicator for disease progression and being incorporated into diagnostic criteria [1].

The main manifestations of MM are succinctly captured by the acronym CRAB, which represents hypercalcemia (C), renal impairment (R), anemia (A) and bone disease (B). Among these, bone disease emerges as the prevailing symptom, impacting more than 80 % of all patients [2], [3].

Identifying osteolytic lesions, a common manifestation of the condition, is crucial for implementing therapy immediately. In contemporary diagnostics, the integration of low-dose computed tomography (CT), along with magnetic resonance imaging (MRI) and hybrid imaging methods (particularly PET/CT), has become indispensable [4].

Dual-energy CT (DECT) is currently gaining prominence in the field of medical imaging. This technique involves the use of two different X-ray energies for imaging, allowing for energy decomposition and enhanced differentiation of signals. Unlike conventional CT (cCT) examinations, DECT enables the distinction between photons with different energy levels. Spectral CT (sCT) can also use multi-energy decomposition. Manufacturers employ different technical configurations, including two different X-ray energies or two-layer detectors with varying sensitivities to different X-ray energies, to achieve this capability. This functionality enables the utilization of postprocessing software to generate multiple parametric maps, as well as what is known as virtual monoenergetic images (VMI) [5], [6].

The bone mineral density (BMD) plays a crucial role in assessing bone strength and resilience. This metric serves as a key indicator of bone integrity and can provide essential information on the risk of osteoporosis and fractures associated with bone weakening. Measurement of BMD is often performed using various diagnostic techniques, such as dualenergy X-ray absorptiometry (DXA) or Quantitative Computed Tomography (QCT). DXA, compared to other BMD estimation methods, has a low cost and low absorbed dose to the patient. However, DXA also has its disadvantages, such as the bias in values caused by the summation of soft tissue values with bone tissue. It is also a 2D imaging method, so the resulting area BMD estimate is the summation of the superficial cortical bone with the more metabolically active internal trabecular bone. Therefore, multi-energy Xray computed tomography methods with three-dimensional output are used for a more detailed evaluation of mineral density distribution in trabecular bone. Most QCT methods are limited to the estimation of BMD of cortical and trabecular bone in the presence of a phantom without more extensive estimation of the partial volumes of the individual elemental components of trabecular bone; bone minerals, collagen, water, bone marrow, and fat components. The fat component reduces the CT number value, and collagen has the opposite effect. Therefore, it is reasonable to quantify the volumes of these components for proper calculation of BMD estimation [7].

Regarding the measurement by QCT, a BMD value between 80 and 120 mg/cm³ characterizes osteopenia, while a BMD value below 80 mg/cm³ defines osteoporosis [9]. Thus, these diseases can be diagnosed and distinguished based on BMD measurements.

This paper focuses on initial experiments to calculate and analyze BMD utilization in a database of 10 patients, 5 with confirmed MM and 5 without spinal pathology.

II. MATERIALS AND DATA

In this study, an anonymized database of ten patients was utilized, consisting of five oncological cases with multiple myeloma presenting lytic lesions in the spine, and five patients with spine images showing a pathological-free condition. Data were acquired with the approval of the ethics committee under the application registration number NU23J-08-00027, and all patients provided informed consent. The information was obtained through Philips Healthcare IQon spectral CT in collaboration with the University Hospital Brno, Department of Radiology and Nuclear Medicine.

The scanning acquisition parameters included a peak tube voltage of 100 kV, tube current of 10 mA, matrix size of 512×512 , and slice thickness of 0.9 mm using a sharp reconstruction kernel and hybrid iterative reconstruction technique (iDose4, set to level 4). Scans covered from the head to the knees with the upper limbs crossed over the abdomen. Finally, the scans were reviewed using a specialized workstation (Intellispace Portal version 12.1; Philips Healthcare) by two independent readers. The diagnosis of multiple myeloma was established based on elevated levels of monoclonal immunoglobulin in the blood and an increased count of plasma cells in the bone marrow. Spectral CT with a low-dose protocol was performed for the initial staging of the disease, following recommendations of the International Myeloma Working Group (IMWG).

Raw SBI format spectral CT data were available for each patient. Spectral CT enabled reconstruction of various parametric images including conventional CT, virtual monoenergetic images at different energies, calcium suppression images, and others using a dedicated workstation (Intellispace Portal version 12.1; Philips Healthcare).

For the purpose of this paper, conventional CT images and virtual monoenergetic images at 40, 80, and 120 keV were reconstructed and used. An example of the available data is shown in Fig. 1.



Fig. 1. Example of available data with lytic lesions (multiple myeloma disease) in different parametric images (from left convention CT, VMI at 40 keV, VMI at 80 keV, and VMI at 120 keV.

Moreover, for each patient, a segmented spine mask was available, generated by the nnU-Net machine learning model presented in [8]. Manual segmentation of lytic lesions was also available for patients with multiple myeloma (see Fig. 2).



Fig. 2. 2D Visualisation of 3D segmentation of lytic lesion in VMI at 40 keV. Left original image, right overlay of original CT image with segmentation masks. Each color represents an individual lesion.

III. METHODS

For the calculation of BMD in this paper, the methodology presented in [10] was used. The mass attenuation coefficient for an absorber comprising a blend of elements can be determined through the following formula:

$$\left(\frac{\mu}{\rho}\right)_T = \sum_{i=1}^N \left(\frac{\mu}{\rho}\right)_i \cdot W_j \tag{1}$$

where $\left(\frac{\mu}{\rho}\right)_i$ is the mass attenuation coefficient of individual elements at photon energy *i*, W_j represents the fractional weight of each element, μ represents the attenuation coefficient, ρ represents the mass density, *N* represents the number of elements, and $\left(\frac{\mu}{\rho}\right)_T$ represents the total mass attenuation coefficient.

With the use of this relationship and the standard photonattenuation tables for the elements, the mass and linear attenuation coefficients for any substance can be determined. The trabecular part of a vertebra consists mainly of five distinct materials: bone mineral, collagen matrix, water, red marrow, and adipose tissue. The mineral content within the bone consists primarily of poorly crystallized calcium hydroxyapatite $Ca_{10}(P0_4)_6(OH)_2$, distributed throughout the collagen matrix.

The mass attenuation coefficients of the components found in trabecular bone are connected to the measured CT numbers in Hounsfield units through the following equation:

CT number =
$$K\left[\sum_{i=1}^{N} \frac{\left(\frac{\mu}{\rho}\right)_{\rho_i}}{\left(\frac{\mu}{\rho}\right)_{H_2O}} - 1\right]$$
 (2)

where K is a constant, approximately equal to 1000 for most CT scanners. This relationship can be reformulated as follows:

CT number =
$$\alpha \rho_{BM} + \eta \rho_C + \omega \rho_W + \beta \rho_F + \theta \rho_M + \delta + \epsilon$$
 (3)

where α (calcium hydroxyapatite), η (collagen), ω (water), β (adipose tissue), and θ (red marrow) are attenuation coefficients dependent on photon energy, $\delta = -1000 \, \text{HU}, \epsilon$ is the number of offset of water, and ρ_{BM} , ρ_{C} , ρ_{W} , ρ_{F} , and ρ_{M} are concentrations of bone mineral, collagen matrix, water, adipose tissue, and red marrow, respectively, in grams per cubic centimeter (g/cm³).

The equation that connects measured CT numbers to concentrations of different substances in the cancellous bone of the spine can be simplified by considering its structure. The trabecular compartment is intricately woven with a complex network of collagen matrix that houses the bone mineral. The collagen matrix has minimal water content in older individuals, who are the most common subjects for bone mineral measurement. The remaining space surrounding the collagen matrix is mainly occupied by red marrow and adipose tissue in various proportions. The water in this part of the trabecular space is sufficiently similar in density and photonattenuation properties to that in the red marrow, allowing them to be combined as non-adipose tissue (ρ_T) , as expressed in the following equation:

$$\omega \rho_W + \theta \rho_M = \gamma \left(\rho_W + \rho_M \right) \equiv \gamma \rho_T \tag{4}$$

Moreover, the ratio of bone minerals to collagen in the matrix remains relatively consistent in the majority of elderly individuals. Decalcified matrix is only found in specific cases (e.g., osteomalacia). On the contrary, osteoporosis leads to an overall reduction in bone mass within the trabecular space. Despite this decline, the ratio of bone minerals to collagen in the remaining matrix material remains essentially unchanged.

Based on published data [11], the mean density of the matrix material (bone mineral plus collagen) (ρ_{TB}) is approximately 1.92 g/cm³ in older individuals. The density of collagen Calone is 1.38 g/cm³, and the density of bone mineral (l) alone is 3.06 g/cm³. Therefore, the expression can be written as follows:

$$\rho_{TB} = \frac{lV_{BM} + CV_c}{CV_{BM} + V_c} \tag{5}$$

where V_{BM} and V_c are volumes occupied by bone mineral and collagen, respectively, per cubic centimeter. This equation can be rewritten as follows:

$$V_c = \frac{(l - \rho_{TB})}{(\rho_{TB} - C)} V_{BM} = \lambda V_{BM}$$
(6)

where

$$\lambda = \frac{l - \rho_{TB}}{\rho_{TB} - C} \tag{7}$$

In addition, the total volume (V_{TB}) occupied by the matrix material (bone mineral plus collagen) must be equal to the sum of its parts, which can be expressed as follows: $V_{TB} =$ $V_{BM} + V_c = V_{BM} + \lambda V_{BM}$ and $V_{TB} = (1 + \lambda) V_{BM}$.

As λ is a constant, the volume occupied by bone minerals and collagen can be represented as a proportion of the total volume taken up by the matrix material (V_{TB}) , given by:

$$V_{BM} = \frac{V_{TB}}{1+\lambda} \tag{8}$$

and

$$V_c = \frac{\lambda V_{TB}}{1+\lambda} \tag{9}$$

The densities of collagen and bone minerals in trabecular bone tissue are known, as are the densities of fat-free tissue and adipose tissue. The value of the density of fat-free tissue, typically denoted g, is reported as 1.02 g/cm³, and the density of adipose tissue, denoted t, is 0.92 g/cm³ [11]. These densities can be expressed as $\rho_T = gV_T$ and $\rho_F = tV_F$, where V_T and V_F are the volumes of fat-free tissue and adipose tissue, respectively.

When using modified partial equations and relationships from (3), the number of unknown variables is reduced to three partial volumes of bone minerals and collagen (V_{TB}) and volumes of tissue with fat (V_F) and without fat (V_T) . Since the total volume must be equal to 1 cm^3 , (3) simplifies to two unknown variables: $V_T = 1 - V_{TB} - V_F$.

The relationship of the CT number to the fractional volumes reduces to the following equation:

CT number =
$$\mu V_{TB} + \beta t V_F + \gamma g (1 - V_{TB} - V_F) + \delta + \epsilon$$
 (10)
where

V

$$\mu = \frac{\alpha l + \eta C \lambda}{1 + \lambda} \tag{11}$$

Equation (10), involving two unknowns V_F and V_{TB} , can then be easily solved as a system of two equations using Xray radiation with two distinct energies in a dual-energy CT system:

CT number =
$$(\mu - \gamma g)V_{TB} + (\beta t - \gamma g)V_F + \gamma g + \delta + \epsilon$$
 (12)
and

CT number' =
$$(\mu' - \gamma'g)V_{TB} + (\beta't - \gamma'g)V_F + \gamma'g + \delta + \epsilon'$$
(13)

Where (12) is the equation for higher energy radiation, and (13) is the equation for lower energy radiation.

From the values of V_{TB} and V_F , many other variables that are of interest to the clinician can be determined, as follows: $V_T = 1 - V_{TB} - V_F$ = non-adipose fractional volume per cubic centimeter, and

$$\rho_{BM} = \frac{l + V_{TB}}{1 + \lambda} \tag{14}$$

where (14) expresses the density of bone mineral (BMD) in g/cm^3 .

The total density of the trabecular bone in the vertebral body, in g/cm^3 , can be calculated as follows:

$$\rho_{TBV} = \rho_{BM} + \frac{C\lambda V_{TB}}{1+\lambda} + tV_F + gV_T \tag{15}$$

The density of hydroxyapatite (ρ_{BM}) is represented by calcium, accounting for 39.9%; $\rho_{Ca} = 0.399 \rho_{BM}$.



Fig. 3. Illustration of the computed BMD map for the first lumbar vertebra in a patient with multiple myeloma: on the left, the conventional CT image with the lesion highlighted by the green contour, and the blue contour indicating the ROI of the trabecular tissue. In the middle, there is a VMI at 40 keV. On the right is an illustration showing the derived BMD map.



Fig. 4. Illustration of the BMD map for the first lumbar vertebra in a patient without spinal pathologies: the left displays a conventional CT image with the blue contour marking the ROI of trabecular tissue. In the middle, a VMI at 40 keV is shown, while the right illustrates the derived BMD map.

IV. RESULTS AND DISCUSSIONS

This initial study used virtual monoenergetic images at 80 and 120 keV to estimate BMD. BMD and other substances were calculated only in the spine region using the available spine segmentation mask. An example of the available data in the conventional CT bone radiology window, along with a VMI at 40 keV and the calculated BMD map for the vertebra of a patient with multiple myeloma and a marked lytic lesion and region of interest (ROI) of trabecular tissue, can be seen in Fig. 3. On the contrary, Fig. 4 presents an example of the same images for the vertebra of a patient without spinal pathology. It can be seen from the images that the contrast between the lesion and the surrounding tissue is much better on the 40 keV VMI than on conventional CT. Even better contrast is evident on the BMD map, where the resulting image is also much smoother.

To evaluate the differences in calculated BMD values between healthy vertebrae and vertebrae affected by lytic lesions, ROIs of trabecular tissue of lumbar vertebrae L1, L2, and L3 of the same patient with multiple myeloma were selected for demonstration, where vertebra L1 was affected by a lytic lesion marked by a radiologist. Furthermore, the L1 and L2 vertebrae of the patient without pathology in the spine were selected. Box-and-whisker plots were created from ROIs of conventional CT data (see Fig. 5) and from the BMD map (see Fig. 6). The examples clearly illustrate that the disparity between the lytic lesion, the vertebrae of the patient with multiple myeloma, and those of the healthy spine is significantly greater on the BMD map compared to conventional CT.

From the ROIs, the mean value and standard deviation were calculated. Tab.I showing mean values (with standard deviations) extracted from ROIs on conventional CT, VMI at 40 keV, and BMD map. MM denotes a patient with multiple myeloma, while H denotes a healthy patient. ROI refers to regions of interest derived from trabecular tissue.

Based on the results obtained and references to the literature, it can be concluded that the calculated BMD results are in line with expectations. The mean BMD values of 222 and 217 g/cm³ for healthy vertebrae fall within the physiological



Fig. 5. The box-and-whisker plot displays ROIs from conventional CT. The lytic lesion is highlighted in red, ROIs from the patient with multiple myeloma are in purple, and ROIs from healthy vertebrae are in green.



Fig. 6. The box-and-whisker plot displays ROIs from BMD map. The lytic lesion is highlighted in red, ROIs from the patient with multiple myeloma are in purple, and ROIs from healthy vertebrae are in green.

range for healthy tissue. On the contrary, in a patient (female, 66 years) with multiple myeloma and lytic lesions in the spine, a mean BMD value of 62 g/cm³ for lytic lesions indicates a high risk of pathological fractures at this site. Regarding the ROI of the same vertebra, the calculated mean BMD value of 152 g/cm³ suggests a moderate impact of osteoporosis on the rest of the vertebra. When examining other vertebrae from the same patient without lytic lesions, the mean BMD values of 132 and 156 g/cm³ suggest that although there are no significant lytic lesions, the vertebrae exhibit lower BMD values than would be physiological. This could be due to

TABLE I AVERAGE VALUES AND STANDARD DEVIATIONS OF THE ROIS

Region of interest	Conv CT [HU]	VMI [HU]	BMD [g/cm ³]
MM - L1 - lesion	84 (38)	160 (56)	62 (17)
MM - L1 - ROI	166 (53)	398 (81)	152 (25)
MM - L2 - ROI	149 (60)	401 (110)	156 (36)
MM - L3 - ROI	140 (46)	341 (83)	132 (28)
H - L1 - ROI	278 (32)	622 (42)	222 (13)
H - L2 - ROI	270 (38)	604 (54)	217 (17)

factors such as patient age, diffuse MM infiltration, or incipient osteopenia/osteoporosis.

V. CONCLUSION

This article presents the results of the calculation of bone mineral density (BMD) in the vertebrae of a patient with multiple myeloma and lytic lesions in the spine. The findings indicate that the average BMD value for a healthy vertebra falls within the physiological range. In the case of the patient with myeloma, a low BMD value was measured in the area of the lytic lesion, suggesting a high risk of pathological fractures. The research also revealed that BMD was lower than physiological even in areas without lytic lesions. The results are promising for future research, indicating that BMD maps could serve as valuable parametric input maps for segmenting lytic lesions in the spine and for further analysis.

ACKNOWLEDGMENT

The paper and the research were supported by Philips Healthcare and Brno University Hospital, Department of Radiology and Nuclear Medicine.

REFERENCES

- R. Silbermann and G. D. Roodman, "Myeloma bone disease: Pathophysiology and management", *Journal of Bone Oncology*, vol. 2, no. 2, pp. 59-69, 2013.
- [2] B. Jamet, C. Bailly, T. Carlier, C. Touzeau, A. -V. Michaud, M. Bourgeois, P. Moreau, C. Bodet-Milin, and F. Kraeber-Bodere, "Imaging of Monoclonal Gammapathy of Undetermined Significance and Smoldering Multiple Myeloma", *Cancers*, vol. 12, no. 2, 2020.
- [3] B. R. Madhira, V. M. Konala, S. Adapa, S. Naramala, P. M. Ravella, K. Parikh, and T. C. Gentile, "Recent Advances in the Management of Smoldering Multiple Myeloma", *World Journal of Oncology*, vol. 11, no. 2, pp. 45-54, 2020.
- [4] J. Hillengass, S. Usmani, S. V. Rajkumar, B. G. M. Durie, M. -V. Mateos, S. Lonial, C. Joao, K. C. Anderson, R. García-Sanz, E. Riva, J. Du, N. van de Donk, J. G. Berdeja, E. Terpos, E. Zamagni, R. A. Kyle, J. San Miguel, H. Goldschmidt, S. Giralt, S. Kumar, N. Raje, H. Ludwig, E. Ocio, R. Schots, H. Einsele, F. Schjesvold, W. -M. Chen, N. Abildgaard, B. C. Lipe, D. Dytfeld, B. M. Wirk, M. Drake, M. Cavo, J. J. Lahuerta, and S. Lentzsch, "International myeloma working group consensus recommendations on imaging in monoclonal plasma cell disorders", *The Lancet Oncology*, vol. 20, no. 6, pp. e302-e312, 2019.
- [5] R. Forghani, B. De Man, and R. Gupta, "Dual-Energy Computed Tomography", *Neuroimaging Clinics of North America*, vol. 27, no. 3, pp. 371-384, 2017.
- [6] N. Rassouli, M. Etesami, A. Dhanantwari, and P. Rajiah, "Detectorbased spectral CT with a novel dual-layer technology: principles and applications", *Insights into Imaging*, vol. 8, no. 6, pp. 589-598, 2017.
- [7] A. D. Brett and J. K. Brown, "Quantitative computed tomography and opportunistic bone density screening by dual use of computed tomography scans", in *Journal of Orthopaedic Translation*, 2015, vol. 3, no. 4, pp. 178-184.
- [8] M. Nohel, R. Jakubicek, L. Blazkova, V. Valek, M. Dostal, P. Ourednicek, and J. Chmelik, "Comparison of Spine Segmentation Algorithms on Clinical Data from Spectral CT of Patients with Multiple Myeloma", in *MEDICON-23 and CMBEBIH-23*, 2024, pp. 309-317.
- [9] "ACR—SPR—SSR Practice Parameter for the Performance of Quantitative Computed Tomography (QCT) Bone Densitometry.", 02023.
- [10] E. L. Nickoloff, F. Feldman, and J. V. Atherton, "Bone mineral assessment: new dual-energy CT approach", in *Radiology*, 1988, vol. 168, no. 1, pp. 223-228.
- [11] R. B. Mazess, "Errors in measuring trabecular bone by computed tomography due to marrow and bone composition", in *Calcified Tissue International*, 1983, vol. 35, no. 1, pp. 148-152.

Assessing Diversity in Predictive Equations for Body Compartment Estimation

Dávid Kampo

Department of Biomedical Engineering, Faculty of Electrical Engineering and Communication Brno University of Technology Brno, Czech Republic xkampo00@vutbr.cz

Abstract—Abstract:

This paper evaluates the validity of prediction equations for estimating total body water (TBW), extracellular water (ECW), fat-free mass (FFM), and body cell mass (BCM) using bioelectrical impedance analysis (BIA). The study focuses on a sample of 10 Czech individuals of European ethnicity. Prediction equations were selected based on similarity to the study population and were compared to reference ranges for accuracy.

Findings show promising outcomes for TBW estimation, with low relative errors (RE) of 0.11% and -3.26% for equations by Deurenberg et al. and Kotler et al. respectively. Matias et al.'s equation for ECW estimation demonstrated the most accurate results with an RE of 2.13%. For FFM estimation, equations by Lukaski et al. and Deurenberg et al. showed favorable outcomes with RE values of 4.4% and -1.27% respectively. However, none of the BCM prediction equations provided satisfactory accuracy.

Further research with larger sample sizes is needed for more accurate validation. Nonetheless, this study offers valuable insights into selecting appropriate prediction equations for BIAbased body composition analysis.

Index Terms—body composition, bioelectrical impedance analysis, multi-frequency BIA, total body water, extracellular water, fat free mass, body cell mass, prediction equations

I. INTRODUCTION

Bioelectrical impedance analysis (BIA) serves as a noninvasive, safe, straightforward, cost-effective, rapid, and portable technological method extensively employed by healthcare practitioners in the medical field, sports experts, individuals engaged in physical activity within fitness environments, or by everyday individuals for at-home health monitoring. Additionaly, BIA has been seamlessly incorporated into smart scales or devices, facilitating the advancement of e-health applications for self-evaluation [2] [3].

The simplicity of BIA measurements stems from administering a low-intensity alternating electric current (AC), denoted in micro-amperes (μA), across the human body at both low (ranging from 1 to 30 kHz) and high frequencies (exceeding 50 kHz) via electrodes positioned on the skin. For the purpose of body composition analysis, impedance (Z, Ohm) is obtained as the ratio of voltage (U, Volt) and current (I, Ampere):

$$Z = \frac{U}{I} \,. \tag{1}$$

Impedance is a complex function and its magnitude can be expressed by the modulus of Z, moreover, the other parameters can be calculated, for instance, resistance (R) as a real part

of impedance and reactance (X) as an imaginary part of impedance. It is relevant to note that the use of direct current (DC) influences the imaginary part of impedance which, subsequently, equals zero [2] [4] [5].

In terms of potential uses of BIA, there is considerable opportunity for this technique to be utilized not only by the general public but also within clinical settings. This potential is supported by extensive research conducted in this field, which has focused on various applications such as staging lung cancer, monitoring pulmonary edema, assessing hydration status and hyponatremia, estimating dry weight in kidney failure, evaluating neural system diseases, monitoring muscular activity, and providing an overall assessment of health status [2] [3] [6] [11].

BIA estimates human body composition based on tissuespecific electrical properties. This paper aims to estimate body compartments using mathematical models and prediction equations with empirical variables. Regression models are used to correlate raw bioimpedance data and body composition against validated reference methods such as Dual-Energy Xray Absorptiometry (DEXA) and isotopes dilution in order to obtain empirical variables. Reference methods that may be used as alternatives to BIA, on the other hand, represent negatives like radiation exposure, higher costs, or complexity [2] [7] [11].

According to the number of frequencies, BIA can be divided into [2]:

- Single-frequency BIA (SF-BIA) uses the frequency of 50 kHz.
- Multiple-frequency BIA (MF-BIA) uses at least one low and one high frequency.
- Bioimpedance spectroscopy (BIS) uses a broad range of frequencies.

Furthermore, there exist two methodologies for conducting BIA measurements concerning the segments on which the analysis is performed. Initially, there is the whole-body BIA (WH-BIA), which evaluates total body impedance, typically with electrodes positioned between the wrist and ankle. Secondly, there is segmental BIA (S-BIA), which assesses body impedance in specific body segments such as the arms, trunk, or legs [2].

Despite BIA's advantages and wide-ranging applications, numerous published studies suffer from methodological in-
consistencies, hindering reproducibility. Many research papers lack critical details such as device selection, electrode placement, and calibration procedures. A 2016 survey by Charlotte Stoddart found that among researchers in the study, 70% couldn't replicate their peers' experiments, and over half failed to reproduce their own [10]. This inadequate reporting may compromise study validity and lead to misinterpretations in future research [8] [9].

Additionally, prediction equations are often tailored to specific samples, highlighting the importance of ensuring the similarity of the analyzed sample to those in the referenced studies in terms of ethnicity, age, or weight. Adherence to standardized protocols and precise measurement of anthropometric data are crucial to guarantee result accuracy; otherwise, variations in measured impedance may arise, potentially resulting in overor underestimations of body compartments. When considering the application of general prediction equations across different ethnic populations, prior testing and validation become essential [8] [11] [12] [13].

The primary objective of this paper is to compare, statistically analyze, and validate five prediction equations for estimating total body water (TBW), extracellular water (ECW), fat-free mass (FFM), and three prediction equations for estimating body cell mass (BCM) utilizing a sample of 10 individuals of European ethnicity. Prediction equations were derived from a diverse range of studies focusing on cohorts analogous to the sample population studied in this paper. Each participant's data were processed using individual prediction equations, and the relative errors (RE) between measured and predicted reference values were computed. The accuracy and validity of each prediction equation were then evaluated based on the RE analysis.

II. METHODS

A. Subjects and Antropometric Measurements

For the purpose of testing and validating prediction equations, measurements were undertaken on a sample consisting of 10 individuals of Czech nationality, representing the European Caucasian ethnicity. The sample comprises 6 females and 4 males, aged between 24 and 49 years. None of the subjects reported any health concerns or regular use of medications, including diuretics, which could potentially impact the measured data. This study obtained approval from the Ethics Committee of Faculty of Medicine, University of Ostrava, denoted by the reference number 21/2020, signifying compliance with ethical standards and protocols.

 TABLE I

 Descriptive Statistics of Height, Weight, and Age for the

 Studied Sample (n=10)

Parameter	Average	Standard deviation [SD]
Height [cm]	173.7	9.89
Weight [kg]	68.3	12.58
Age [years]	40.6	6.81

Before conducting the measurement, anthropometric parameters such as height and weight were assessed to guarantee the precision of the findings presented. Each participant was weighed using a Tefal Premiss 2 PP1401V0 scale, accurate to the nearest 0.2 kg, while wearing light clothing. Height was measured in the standing position, without shoes, to the nearest 0.5 cm using a soft measuring tape. Subsequently, both height and weight were used in prediction equations to determine the relevant body compartments needed for this study.

Standard deviation (SD) and mean values of height, weight, and age for the examined sample are presented in Table I.

B. Bioimpedance analysis

Data were acquired using a bioimpedance analyzer provided by the Institute of Scientific Instruments of the Czech Academy of Sciences, focusing on the left side of the human body due to negligible differences concluded by available studies between the left and right sides [9]. For this study, the MF-BIA approach was employed, utilizing a four-electrode configuration system for data collection. In this electrode composition, one pair (CC - current carrying electrodes) is utilized to administer a constant current into the tissue, while the other pair (PU - voltage pick-up electrodes) detects voltage changes resulting from varying tissue conductivity. Subsequently, equation (1) is used to calculate impedance.

KendallTM H34SG ECG electrodes were utilized during the measurements, applied in whole-body BIA. Voltage electrodes were positioned on the wrist (a. radialis) and ankle (a. tibialis), while current electrodes were placed distally, 5 cm from the positions of the voltage electrodes, to prevent any potential unwanted inter-electrode interactions that could lead to elevated impedance values.

With the objective of acquiring the necessary parameters incorporated into the prediction equations for estimation, impedance was measured at 13 frequencies ranging from 1 kHz to 1000 kHz with an applied electric current of 0.08 mA. However, only impedance, resistance, and reactance measured at 5 kHz, 50 kHz, and 100 kHz were utilized for predicting body composition via prediction equations. The standard deviation and average values of impedance (Z,Ohm), resistance (R,Ohm), and reactance (X,Ohm) measured at these frequencies are summarized in Table II.

TABLE II SUMMARY OF IMPEDANCE, RESISTANCE, AND REACTANCE AT SELECTED FREQUENCIES FOR BODY COMPOSITION ESTIMATION

Parameter	Average	SD
R_{5kHz} [Ohm]	349.95	55.77
Z_{50kHz} [Ohm]	280.68	49.33
R_{50kHz} [Ohm]	269.29	47.21
X_{50kHz} [Ohm]	-78.96	15.77
R_{100kHz} [Ohm]	228.12	38.79
Z_{100kHz} [Ohm]	246.57	42.79

Measurements were conducted at 9 am in a room with a stable temperature of 22 °C. Prior to each test, the analyzer

underwent calibration, with successful calibration confirmed if the resistance reference value reached 500 Ohms.

All data were obtained under standardized conditions for MF-BIA, and participants were instructed to adhere to the protocol and conditions for BIA analysis to ensure maximum precision in measurements. The protocol includes the following instructions [9]:

- Refraining from eating for a minimum of 4 hours before measurements.
- Ensuring normal hydration levels.
- Emptying the bladder, preferably before measurements.
- Avoiding any physical activity for at least 4 hours before measurements.

During the measurements in the upright position, participants were instructed to refrain from talking and to stay relaxed, in order to achieve reliable results and avoid electrode detachment or changes in the acquired signal.

C. Prediction Equations

This paper focuses on assessing and examining the validity of five prediction equations for estimating TBW, ECW, and FFM, as well as three prediction equations for predicting BCM. The selection of prediction equations was made considering the similarity between the sample of subjects in previous studies and the sample studied in this paper.

D. Total Body Volume (TBW)

For the assessment of TBW (in liters), the first two equations studied were developed by Deurenberg et al. in 1995. Their research involved a sample of 139 individuals from the Netherlands and Italy, with an average age of 25.35 years for all participants. These prediction equations include parameters such as height (in cm), weight (in kg), impedance (Ohm) at 50 kHz and 100 kHz, age (in years), and sex (coded as 1 for male and 0 for female) [1] [13]. The equations by Deurenberg et al. are as follows:

$$TBW = 6.69 + 0.34573 \cdot \frac{height^2}{Z_{100kHz}} + 0.17065 \cdot weight - 0.11 \cdot age + 2.66 \cdot sex [1],$$
(2)

$$TBW = 6.53 + 0.3674 \cdot \frac{height^2}{Z_{100kHz}} + 0.17531 \cdot weight - 0.11 \cdot age + 2.83 \cdot sex [1].$$
(3)

The third prediction equation was formulated by Heitmann et al., also based on a sample of 139 individuals in 1990. The study was conducted on the Danish population aged between 35 and 65 years. The developed equation incorporates height (in cm), weight (in kg), resistance (in Ohms) at 50 kHz, age (in years), and sex (coded as 1 for male, 0 for female) [2]. The equation by Heitmann et al. is as follows:

$$TBW = -17.58 + 0.24 \cdot \frac{height^2}{R_{50kHz}} - 0.172 \cdot weight + 0.04 \cdot age \cdot weight + 0.165 \cdot height [l].$$

$$(4)$$

The fourth prediction equation was formulated by Kotler et al. based on a sample of 332 individuals comprising White, Black, and Hispanic participants, with a mean age of 41 years. Kotler's study involved the creation of sex-specific prediction equations using parameters such as height (in cm), impedance (in Ohms) at 50 kHz, and weight (in kg) [11] [16]. The equations by Kotler et al. are as follows:

$$TBW_{male} = 0.58 \cdot \frac{height^{1.62}}{Z_{50kHz}^{0.7}} \cdot \frac{1}{1.35} + 0.32 \cdot weight - 3.66 [1],$$
(5)

$$TBW_{female} = 0.76 \cdot \frac{height^{1.99}}{Z_{50kHz}^{0.58}} \cdot \frac{1}{18.91} + 0.14 \cdot weight - 0.86 [l].$$
(6)

The final prediction equation for TBW estimation examined in this paper is the formula outlined by Schoeller et al., validated on a sample of 125 Caucasians aged between 14 and 53 years. Schoeller developed equations incorporating parameters such as height (in cm), resistance (in Ohms) at 50 kHz, and weight (in kg) [7]. The equation by Schoeller et al. is as follows:

$$TBW = 0.499 \cdot \frac{height^2}{R_{50kHz}} + 0.08 \cdot weight + 2.9 [1].$$
(7)

E. Extracellular Water (ECW)

As the initial equation employed for predicting the ECW compartment (in liters), was utilized the one formulated by Matias et al. This equation was derived from a sample of 208 individuals, all of whom were athletic with a mean age of 21 years. It includes parameters such as height (measured in centimeters), resistance (R) and reactance (X), both measured in Ohms at 50 kHz, weight (measured in kilograms), and sex (coded as 1 for male and 0 for female) [11]. The equation by Matias et al. is as follows:

$$ECW = 1.579 + 0.055 \cdot \frac{height^2}{R_{50kHz}} + 0.127 \cdot weight + 0.006 \cdot \frac{height^2}{X_{50kHz}} + 0.932 \cdot sex [l].$$
(8)

Next in order, was the equation designed by Sergi et al. on a sample of 40 Caucasians ranging in age from 21 to 81. The equation incorporates height (cm), resistance, and reactance (Ohm) at 50 kHz, weight (kg), and sex (0 for male, 1 for female) [11] [14]. The equation by Sergi et al. is formulated as:

$$ECW = -5.22 + 0.2 \cdot \frac{height^2}{R_{50kHz}} + 0.08 \cdot weight + 0.005 \cdot \frac{height^2}{X_{50kHz}} + 1.86 \cdot sex + 1.9 [l].$$
(9)

The third equation considered for assessing ECW is the formulation by Lukaski et al. This equation was established based on a sample of 110 individuals, consisting of White and African American participants aged between 20 and 73 years. It includes parameters such as height (in centimeters),

resistance, and reactance (in Ohms) at 50 kHz, along with weight (in kilograms) [11] [14]. The equation by Lukaski et al. is as follows:

$$ECW = 0.189 \cdot \frac{height^2}{R_{50kHz}} + 0.052 \cdot weight - 0.0002 \cdot \frac{height^2}{X_{50kHz}} + 1.03 [l].$$
(10)

The last two equations were applied in the reference study involving a sample of 169 elderly individuals aged over 60, demonstrating accurate results. The equation by Segal et al. includes the variables of height (measured in cm), resistance (measured in Ohms) at 5 kHz, and weight (measured in kg), and is as follows [3]:

$$ECW = -6.1 + 0.284 \cdot \frac{height^2}{R_{5kHz}} + 0.112 \cdot weight [1].$$
(11)

The final equations, also from the same study, were developed by Visser et al. and are sex-specific, incorporating the variables of height (measured in cm), resistance (measured in Ohms) at 5 kHz, and weight (measured in kg), and are as follows [3]:

$$ECW_{male} = 4.8 + 0.225 \cdot \frac{height^2}{R_{5kHz}} [1],$$
 (12)

$$ECW_{female} = 1.7 + 0.2 \cdot \frac{height^2}{R_{5kHz}} + 0.057 \cdot weight [l].$$
(13)

F. Fat Free Mass (FFM)

Concerning FFM estimation (in kilograms), the first evaluated prediction equation was formulated by Deurenberg et al. based on a sample of 827 individuals aged between 16 and 83, of unspecified ethnicity. The variables in the equation include height (measured in cm), resistance (measured in Ohms) at 50 kHz, weight (measured in kg), and age (measured in years), and the equation is as follows [17]:

$$FFM = -12.44 + 0.34 \cdot \frac{height^2}{R_{50kHz}} + 0.1534 \cdot height + 0.1534 \cdot height + 0.273 \cdot weight - 0.127 \cdot age [kg].$$
(14)

Kyle et al. conducted a study involving 343 White individuals aged 22 to 94. The variables used in their equation comprise height (measured in cm), resistance, and reactance (both measured in Ohms) at 50 kHz, as well as weight (measured in kg), and the equation is [17]:

$$FFM = -4.104 + 0.518 \cdot \frac{height^2}{R_{50kHz}} + 0.231 \cdot weight + 0.13 \cdot X_{50kHz} [kg].$$
(15)

The third tested equation was developed by Lukaski et al. using a sample of 114 subjects of White and Black ethnicity aged 18 to 50. The equation incorporates variables such as height (measured in cm), resistance, and reactance (both measured in Ohms) at 50 kHz, and weight (measured in kg), and the equation is [15]:

$$FFM = 0.756 \cdot \frac{height^2}{R_{50kHz}} + 0.11 \cdot weight + 0.107 \cdot X_{50kHz} - 5.463 \, \text{[kg]} \,.$$
(16)

Bedogni et al. employed and tested the prediction equation for FFM on a sample of 35 anorexic women from Italy, without specifying the age range of the individuals in the study. The equation from the study includes variables such as height (measured in cm), impedance (measured in Ohms) at 50 kHz, and weight (measured in kg), and the equation is [17]:

$$FFM = 0.6 \cdot \frac{height^2}{Z_{50kHz}} + 0.2 \cdot weight + 3.3 \,[kg] \,.$$
(17)

The last prediction equation utilized in the paper was proposed by Sun et al. The study involved a sample of 1304 individuals of Black and White ethnicity aged 12 to 94. The study describes sex-specific prediction equations that involve variables such as height (measured in cm), resistance (measured in Ohms), and weight (measured in kg), and the equations are [15]:

$$FFM_{male} = -10.68 + 0.65 \cdot \frac{height^2}{R_{50kHz}} + 0.26 \cdot weight + 0.02 \cdot R_{50kHz} \text{ [kg]}, \qquad (18)$$

$$FFM_{female} = -9.53 + 0.69 \cdot \frac{height^2}{R_{50kHz}} + 0.17 \cdot weight + 0.02 \cdot R_{50kHz} \text{ [kg]}.$$
(19)

G. Body Cell Mass (BCM)

Prediction equations for BCM (in kilograms) were sourced from the study led by M. Dittmar and H. Reber, which outlines three prediction equations suitable for estimating BCM. In pursuit of the main goal of this paper, all three equations from this study were utilized. The study was conducted on 160 German participants aged 60 to 90 years. This decision was made due to the lack of alternative equations in the literature that exhibit some minimal similarities with the studied sample in this paper. The three equations incorporate parameters such as height (measured in cm), reactance and resistance (measured in Ohms) at 5 kHz, 50 kHz, and 100 kHz, weight (measured in kg), and sex (coded as 1 for male and 0 for female), and are as follows [6]:

$$BCM = 1.898 \cdot \frac{height^2}{X_{50kHz} + \frac{R_{50kHz}^2}{X_{50kHz}}} - 0.051 \cdot weight + 4.180 \cdot sex + 15.496 \, [kg],$$
(20)

$$BCM = 1.118 \cdot \frac{height^2}{\frac{R_{5kHz} \cdot R_{50kHz}}{R_{5kHz} - R_{50kHz}}} + 4.250 \cdot sex + 14.457 \, [kg] \,,$$
(21)

$$BCM = 0.822 \cdot \frac{height^2}{\frac{R_{5kHz} \cdot R_{100kHz}}{R_{5kHz} - R_{100kHz}}} + 4.158 \cdot sex + 14.096 \, [\text{kg}] \,.$$
(22)

III. RESULTS

To validate the measured data, anticipated values were computed for each individual by considering the normal percentage ranges specified in the literature for various body compartments. Typically, total body water (TBW) accounts for 60% of total body weight, fluctuating between 45% to 75%, while extracellular water (ECW) generally constitutes around 20% of total body weight in individuals with normal hydration levels. Fat-free mass (FFM) typically ranges from 68% to 90% of total body weight, and body cell mass (BCM) is expected to represent between 30% to 45% of total body weight [18] [19] [20] [21].

By contrasting the calculated values from the prediction equations outlined in this paper with these reference ranges, the accuracy of each prediction equation was evaluated for Czech individuals of European ethnicity within the study cohort. The measured estimations, alongside the computed reference bands based on the percentage ranges of analyzed body compartments, are illustrated in Figures 1 and 2.

The evaluation of TBW estimation reveals promising outcomes across various prediction equations, presenting low relative errors. Deurenberg et al.'s equation (2) stands out with an RE of 0.11%, followed closely by Kotler et al.'s equations (5) and (6), displaying a slight underestimation with an RE of -3.26%. Deurenberg et al.'s Equation 3 also yields acceptable results, demonstrating an RE of 4.98%. However, Schoeller et al.'s equation (7) notably overestimates TBW with an RE of 10.47%, while Heitmann et al.'s equation (4) significantly underestimates TBW, recording an RE of -22.4%. Therefore, the equations by Kotler et al. and Deurenberg et al. exhibit accurate and precise estimations suitable for the middle-European population.

For ECW estimation, equations by Sergi et al. (9), Lukaski et al. (10), Segal et al. (11), and Visser et al., (12) and (13), consistently overestimate ECW, exhibiting an average relative error (RE) of 59.22% across all four equations. In contrast, the prediction equation by Matias et al. provides accurate ECW estimation with an RE of 2.13%, demonstrating its validity for the sample in this paper. http

In the subsequent analysis, the prediction equations for estimating FFM exhibit varying degrees of accuracy. The equation developed by Sun et al., (18) and (19), demonstrates the highest level of inaccuracy, with an RE of 22.87%. This is followed by the prediction equation designed by Bedogni et al. (17), which performs slightly better with an RE of 18.11%. The prediction equation proposed by Kyle et al. (15) yields an RE of 7.87%, indicating a moderate level of accuracy. Notably, the prediction equations by Lukaski et al. (16) and Deurenberg et al. (14) produce more favorable results, with RE values of 4.4% and -1.27% respectively. This suggests that these two prediction equations can be considered valid for estimating FFM, with Deurenberg's equation exhibiting a slight underestimation.

Finally, three prediction equations for estimating BCM were applied, all developed by Dittmer et al. (20,21,22). Equation

(20) notably underestimates the BCM compartment, with an RE of -272.41. Equation (21) exhibits a similar trend, demonstrating an RE of 58.24%. Although Equation (22) shows the most favorable outcome among the three equations, it still falls short of providing satisfactory accuracy, with an RE of 35.15%. Consequently, none of the utilized equations appear to be suitable for obtaining accurate and valid estimations of BCM.

All relative errors for the prediction equations are summarized in Table III.

TABLE III Relative Error (%) for Prediction Equations Estimating TBW, ECW, FFM, and BCM

Total Body Water	RE [%]	Extracellular Water	RE [%]
Eq. (2)	0.11	Eq. (8)	2.13
Eq. (3)	4.98	Eq. (9)	53.25
Eq. (4)	-22.4	Eq. (10)	67.73
Eq. (5) and (6)	-3.26	Eq. (11)	63.83
Eq. (7)	10.47	Eq. (12) and (13)	52.06
Fat Free Mass	RE [%]	Body Cell Mass	RE [%]
Eq. (14)	-1.27	Eq. (20)	-272.41
Eq. (15)	7.87	Eq. (21)	35.15
Eq. (16)	4.4	Eq. (22)	58.24
Eq. (17)	18.11	-	-
Eq. (18) and (19)	22.87	-	-

IV. DISCUSSION

The prediction equations defined in this paper were employed to compute estimates of TBW, ECW, FFM, and BCM for each participant. Consequently, each participant had five estimates for TBW, ECW, and FFM, and three estimates for BCM. Given the absence of studies focusing on the middle-European population, specifically nationalities such as Slovakian, Czech, Hungarian, or Polish, prediction equations developed on comparable study groups under similar conditions and relatively similar age ranges were selected for this study. However, this approach could not be universally applied due to variability in the studies and a lack of reporting, making it challenging to replicate the calculations. Furthermore, the availability of prediction equations developed on similar samples to that of this study was limited for certain body compartments such as BCM. Hence, all equations were chosen based on the best available alternatives and compromises that could align with the sample group and measurement procedure.

Equation (4) by Heitmann et al. for TBW estimation markedly underestimates TBW volume, potentially attributed to the older age range (35 - 65 years) of the study group compared to the average age of 40.6 years of participants in this study, given that TBW typically decreases with age [22].

The majority of prediction equations for ECW estimation in this paper yield unsatisfactory results. Equation (10) by Lukaski et al., developed on a sample of African Americans, exhibits a notably high RE of 67.73%, likely due to ethnic disparities, as prediction equations often display high ethnicity



Fig. 1. TBW and ECW Estimations for Studied Subjects



Fig. 2. BCM and FFM Estimations for Studied Subjects

specificity. Similarly, equation (11) by Segal et al. demonstrates a high RE of 63.83% as it was derived from a study involving individuals over 60 years old, outside the age range of participants in this study.

Overall, almost all tested prediction equations for FFM estimation yield satisfactory results, except equations (18) and (19) by Sun et al., which involved a mixture of Black and White participants, potentially explaining inaccuracies. Similarly, equation (17) by Bedogni et al. fails to provide valid results, possibly due to its inclusion of anorexic Italian women, diverging from the normally weighted study group in this paper.

Regarding BCM estimation, all three equations, (20), (21) and (22) yield erroneous results as they were derived from a study by M. Dittmar and H. Reber, which included participants aged 60 - 90 years, significantly older than those in this study, likely resulting in invalid outcomes.

Concluding, one of the limitations of this study is its relatively small study group. Further research with a larger number of participants is needed to gain data that could more accurately assess the validity of the prediction equations presented in this paper. Nevertheless, the findings here can help identify which prediction equations hold promise for accurate estimation and which do not. Additionally, the utilization of different bioimpedance analyzer in this paper compared to those in the analyzed studies, as well as variations in participants' adherence to standardized measurement protocols and conditions before measurements, may also introduce some level of error in estimations.

V. CONCLUSION

In conclusion, this research highlights the importance of selecting appropriate prediction equations tailored to the specific population and context for accurate assessments of body composition using BIA. Notably, the prediction equations by Deurenberg et al. (0.11% RE for TBW), Matias et al. (2.13% RE for ECW), Lukaski et al. (4.4% RE for FFM), and Schoeller et al. (-1.27% RE for FFM) exhibit satisfactory accuracy with relative errors below 5%. These findings emphasize the need for further research to develop and validate prediction equations that account for diverse populations and clinical settings.

ACKNOWLEDGMENT

I would like to thank Dr. Ing. Vlastimil Vondra for his supervision and guidance, as well as Mendel University in Brno for providing participants for this study.

REFERENCES

- KHALIL, Sami; MOHKTAR, Mas and IBRAHIM, Fatimah, 2014. The Theory and Fundamentals of Bioimpedance Analysis in Clinical Status Monitoring and Diagnosis of Diseases. Online. Sensors. Roč. 14, č. 6, s. 10895-10928. Available at: https://doi.org/10.3390/s140610895.
- [2] EL DIMASSI, Sali; GAUTIER, Julien; ZALC, Vincent; BOUDAOUD, Sofiane and ISTRATE, Dan. Mathematical issues in body water volume estimation using bio impedance analysis in eHealth. Online. S. 1-3. Available at: https://hal.science/hal-04220658/document.

- [3] RITZ, P., 2001. Bioelectrical Impedance Analysis Estimation of Water Compartments in Elderly Diseased Patients: The Source Study. Online. The Journals of Gerontology Series A: Biological Sciences and Medical Sciences. 2001-06-01, roč. 56, č. 6, s. M344-M348. Available at: https://doi.org/10.1093/gerona/56.6.M344.
- [4] BERA, Tushar Kanti, 2014. Bioelectrical Impedance Methods for Noninvasive Health Monitoring: A Review. Online. Journal of Medical Engineering. 2014-06-17, roč. 2014, s. 1-28. Available at: https://doi.org/10.1155/2014/381251.
- [5] AMINI, M.; HISDAL, J. and KALVØY, H., 2018. Applications of bioimpedance measurement techniques in tissue engineering. Online. Journal of Electrical Bioimpedance. 2018-12-31, roč. 9, č. 1, s. 142-158. Available at: https://doi.org/10.2478/joeb-2018-0019.
- [6] DITTMAR, Manuela and REBER, Helmut, 2001. New equations for estimating body cell mass from bioimpedance parallel models in healthy older Germans. Online. American Journal of Physiology-Endocrinology and Metabolism. 2001-11-01, roč. 281, č. 5, s. E1005-E1014. Available at: https://doi.org/10.1152/ajpendo.2001.281.5.E1005.
- [7] SCHOELLER, D. A. and LUKE, A., 2000. Bioelectrical Impedance Analysis Prediction Equations Differ between African Americans and Caucasians, but It Is Not Clear Why. Online. Annals of the New York Academy of Sciences. Roč. 904, č. 1, s. 225-226. Available at: https://doi.org/10.1111/j.1749-6632.2000.tb06456.x.
- [8] KYLE, Ursula G.; BOSAEUS, Ingvar; DE LORENZO, Antonio D.; DEURENBERG, Paul; ELIA, Marinos et al., 2004. Bioelectrical impedance analysis—part II: utilization in clinical practice. Online. Clinical Nutrition. Roč. 23, č. 6, s. 1430-1453. Available at: https://doi.org/10.1016/j.clnu.2004.09.012.
- [9] BRANTLOV, Steven; JØDAL, Lars; LANGE, Aksel; RITTIG, Søren and WARD, Leigh C., 2017. Standardisation of bioelectrical impedance analysis for the estimation of body composition in healthy paediatric populations: a systematic review. Online. Journal of Medical Engineering & Technology. 2017-08-08, roč. 41, č. 6, s. 460-479. Available at: https://doi.org/10.1080/03091902.2017.1333165.
- [10] STODDART, Charlotte. Is there a reproducibility crisis in science? Online. Nature. S. d41586-019-00067-3. Available at: https://doi.org/10.1038/d41586-019-00067-3.
- [11] CORATELLA, Giuseppe; CAMPA, Francesco; MATIAS, Catarina N.; TOSELLI, Stefania; KOURY, Josely C. et al., 2021. Generalized bioelectric impedance-based equations underestimate body fluids in athletes. Online. Scandinavian Journal of Medicine & Science in Sports. Roč. 31, č. 11, s. 2123-2132. Available at: https://doi.org/10.1111/sms.14033.
- [12] MATIAS, Catarina N.; SANTOS, Diana A.; JÚDICE, Pedro B.; MA-GALHÃES, João P.; MINDERICO, Cláudia S. et al., 2016. Estimation of total body water and extracellular water with bioimpedance in athletes: A need for athlete-specific prediction models. Online. Clinical Nutrition. Roč. 35, č. 2, s. 468-474. Available at: https://doi.org/10.1016/j.clnu.2015.03.013.
- [13] DEURENBERG, Paul; TAGLIABUE, Anna and SCHOUTEN, Frans J. M., 1995. Multi-frequency impedance for the prediction of extracellular water and total body water. Online. British Journal of Nutrition. Roč. 73, č. 3, s. 349-358. Available at: https://doi.org/10.1079/BJN19950038.
- [14] MATIAS, Catarina N.; SANTOS, Diana A.; JÚDICE, Pedro B.; MA-GALHÃES, João P.; MINDERICO, Cláudia S. et al., 2016. Estimation of total body water and extracellular water with bioimpedance in athletes: A need for athlete-specific prediction models. Online. Clinical Nutrition. Roč. 35, č. 2, s. 468-474. Available at: https://doi.org/10.1016/j.clnu.2015.03.013.
- [15] HOFSTEENGE, Geesje H.; CHINAPAW, Mai JM and WEIJS, Peter JM, 2015. Fat-free mass prediction equations for bioelectric impedance analysis compared to dual energy X-ray absorptiometry in obese adolescents: a validation study. Online. BMC Pediatrics. Roč. 15, č. 1. Available at: https://doi.org/10.1186/s12887-015-0476-7.
- [16] KOTLER, DP; BURASTERO, S; WANG, J and PIERSON, RN, 1996. Prediction of body cell mass, fat-free mass, and total body water with bioelectrical impedance analysis: effects of race, sex, and disease. Online. The American Journal of Clinical Nutrition. Roč. 64, č. 3, s. 489S-497S. Available at: https://doi.org/10.1093/ajcn/64.3.489S.
- [17] COËFFIER, Moise; GÂTÉ, Mathilde; RIMBERT, Agnès; PETIT, André; FOLOPE, Vanessa et al., 2020. Validity of Bioimpedance Equations to Evaluate Fat-Free Mass and Muscle Mass in Severely Malnourished Anorectic Patients. Online. Journal of Clinical Medicine. Roč. 9, č. 11. Available at: https://doi.org/10.3390/jcm9113664.

- [18] What is the average percentage of water in the human body?, 2020. Online. Medical News Today. Available at: https://www.medicalnewstoday.com/articles/what-percentage-of-thehuman-body-is-water.
- What is body composition?, 2019. Online. Metabolic Igniters. Available at: https://www.tntmetabolicigniters.ca/2019/03/05/body-compositionwhat-is-it-how-to-measure-it/.
- [20] Physiology, Body Fluids, 2023. Online. E. BRINKMAN, Joshua; DO-RIUS, Bradley and SHARMA, Sandeep. National Library of Medicine. Available at: https://www.ncbi.nlm.nih.gov/books/NBK482447/.
- [21] MUURAHAINEN, Norma, 1998. What This Means: Have Mass, Will Travel. Online. Available at: https://www.poz.com/article/What-This-Means-Have-Mass-Will-Travel-4859-5861.
- [22] LU, Hong; AYERS, Eric; PATEL, Pragnesh and MATTOO, Tej K., 2023. Body water percentage from childhood to old age. Online. Kidney Research and Clinical Practice. Roč. 42, č. 3, s. 340-348. Available at: https://doi.org/10.23876/j.krcp.22.062.

NXP SEMICONDUCTORS CZECH REPUBLIC



NXP Semiconductors se zabývá vývojem a výrobou mikroprocesorů a analogových součástek pro automobilový a spotřební průmysl, počítačové sítě a bezdrátové technologie. V České republice je společnost NXP Semiconductors zastoupena vývojovými centry v **Rožnově pod Radhoštěm** a v **Brně**. Nově jsme otevřeli i kanceláře na UTB Zlín a VŠB TU v Ostravě.

NXP je světová špička v oblasti **NFC** technologií, **Wireless** charging či **Automotive Electronics**. Kariéra ve firmě NXP Czech Republic znamená příležitost spolupracovat s těmi nejlepšími odborníky v **přátelském R&D prostředí**.

V moderních laboratořích pracuje mezinárodní tým, čítající více než 500 embedded software odborníků **(75 studentů)** z oblastí **Machine Learning**, Security, Wireless charging, Connectivity, **Motor Control**, Functional safety, **Embedded Systems**, Automotive Ethernet, **Signal processing**, Software Validation and Verification, Technical Marketing a v neposlední řadě zákaznické centrum spolu s technicko-obchodním oddělením.





Hledáme talentované studenty

Studentům nabízíme

- Spolupráci na vývoji nejmodernějších technologií
- Práci na reálných projektech, technický mentoring od světových kapacit
- Výborné mzdové ohodnocení a časovou flexibilitu
- Pomoc s bakalářskou, diplomovou či disertační prací
- Studentské akce a soutěže NXP CUP, Campus Engage – meet the experts, Women in NXP, atd.

Jak se zapojit?

- Přihlaste se na studentské či absolventské pozice skrze QR kód výše
- Kontaktujte nás na: josef.halamik@nxp.com

Optimizing of pre-processing analysis for Illumina RNA-Seq data in *Arabidopsis thaliana*

1st Jana Schwarzerová Department of Biomedical Engineering Faculty of Electrical Engineering and Communication, Brno University of Technology Brno, Czech Republic Department of Molecular and Clinical Pathology and Medical Genetics, University Hospital Ostrava, Ostrava, Czech Republic Molecular System Biology (MOSYS), University of Vienna, Vienna, Austria Jana.Schwarzerova@vut.cz

^{2nd} Patrícia Janigová Department of Biomedical Engineering Faculty of Electrical Engineering and Communication, Brno University of Technology Brno, Czech Republic <u>231027@vut.cz</u>

^{3rd} Martina Dvořáčková *Mendel Centre for Plant Genomics and Proteomics, CEITEC, Masaryk University,* Brno, Czech Republic <u>martina.dvorackova@ceitec.muni</u> .cz ^{4th} Wolfram Weckwerth Molecular Systems Biology (MOSYS), University of Vienna. Vienna, Austria Vienna Metabolomics Center (VIME), University of Vienna, Vienna, Austria wolfram.weckwerth@univie.ac.at

Abstract — Gene expression analysis through RNA sequencing (RNA-Seq) has revolutionized molecular biology, providing profound insights into the intricate transcriptional landscapes of organisms. Arabidopsis thaliana, a widely studied model plant, serves as a cornerstone for investigating fundamental biological and ecology processes. However, accurate interpretation of RNA-Seq data hinges on meticulous pre-processing methods to ensure data integrity and trustworthiness, especially in the context of Illumina sequencing. In this research, we present a comprehensive framework for optimizing pre-processing analysis tailored specifically for Arabidopsis thaliana RNA-Seq datasets generated through Illumina sequencing. Our approach encompasses rigorous quality control, precise read alignment, transcript quantification, and normalization procedures crucial for subsequent differential expression analysis. Additionally, we address unique considerations and challenges inherent to Arabidopsis thaliana datasets, providing valuable insights for researchers in the field.

Keywords — Gene expression analysis, Quality control, Arabidopsis thaliana, Transcriptomics

I. INTRODUCTION

Gene expression analysis has become a cornerstone of modern molecular biology research, offering invaluable insights into the complex regulatory mechanisms governing cellular processes [1]. Among the various techniques employed for gene expression profiling, RNA sequencing (RNA-Seq) has emerged as a powerful tool due to its high sensitivity, wide dynamic range, and ability to provide comprehensive transcriptome information. With Illumina sequencing having established itself as one of the primary RNA sequencing platforms in genomics, it continues to be a dominant force. There is an increasing imperative to consistently optimize and discuss its outputs, especially for non-traditional organisms such as various plants utilized in ecological research, which are gaining prominence. Plant science, primarily focused on the study of *Arabidopsis thaliana*, is particularly prominent within the realm of ecology.

Arabidopsis thaliana [2], a small flowering plant native to Eurasia, has long served as a model organism for studying plant genomics. Its relatively simple genome, short life cycle, and amenability to genetic manipulation make it an ideal system for investigating fundamental biological processes [3].

Currently, there is a continuous push for advancements in sequencing methodologies to outpace techniques like microarrays. Despite the availability of numerous algorithms for processing microarray data, this method proves to be more costly and complex to perform in wet lab processes, potentially leading to inaccuracies in the data obtained. However, with precise and effective pre-processing analysis, RNA-Seq data has the potential to provide equal or even superior informational value compared to microarray data [4].

Accurate interpretation of RNA-Seq data necessitates rigorous pre-processing steps to ensure data quality and reliability. Quality control measures are crucial for identifying and mitigating potential biases and artifacts introduced during library preparation and sequencing. Furthermore, read alignment to the reference genome, transcript quantification, and normalization are essential steps in the pre-processing pipeline, enabling accurate and reproducible analysis of gene expression patterns.

In this study, we focus on elucidating the pre-processing methods tailored specifically for RNA-Seq data analysis in Arabidopsis thaliana. We aim to provide a comprehensive overview of the key steps involved in ensuring the robustness and accuracy of gene expression analysis. Additionally, we discuss specific considerations and challenges unique to Arabidopsis thaliana datasets, particularly in the context of systems biology approaches aimed at unraveling the intricacies of plant molecular networks based on our previous studies [5], [6]. By elucidating the gold standard of preprocessing protocols, this study aims to contribute to the advancement of research in Arabidopsis thaliana and provide valuable insights into the broader field of plant systems biology.

II. MATERIALS AND METHODS

A. Dataset – Arabidopsis thaliana

Our study analysed on the utilization of *Arabidopsis thaliana*, a model plant species renowned for its pivotal role in molecular biology research. Specifically, we focused on the Columbia (Col-0) ecotype, a widely studied genetic background, along with specific mutants obtained from the Arabidopsis Biological Resource Center. These mutant lines were selected based on their relevance to the biological processes under investigation, thus enriching the diversity of our experimental material.

The RNA-seq data utilized in our research were sourced from publicly available repositories, specifically the NCBI Gene Expression Omnibus (GEO) database. The dataset was deposited under the GEO Series accession number GSE188493 [7], facilitating transparency and reproducibility in our analyses. The data originated from a study conducted by Zhong et al [7]. All Arabidopsis samples utilized in this investigation belong to the Columbia ecotype (Col-0) and were cultivated at 22°C under LD conditions (16 hours of light, 8 hours of darkness) [7].

The analysed Arabidopsis samples includes Col-0 (SRR16892914, SRR16892916, SRR16892917); and mutants: asf1a (SRR16892918, SRR16892919, SRR16892920); hira-1 (SRR16892921, SRR16892923); SRR16892922, fas1 (SRR16892924, SRR16892925, SRR16892926); fas2 asf1b (SRR16892927, SRR16892928, SRR16892929); (SRR16892930, SRR16892931, SRR16892932) and fwa (SRR16892933, SRR16892934, SRR16892935). [7]

B. Pre-processed RNA-Seq pipeline

The gold standard preprocessing workflow [8] was applied to RNA-Seq data from *Arabidopsis thaliana*. FastQC [9] was utilized to evaluate overall sequence quality, encompassing GC percentage distribution and the presence of overrepresented sequences. The initial step involves merging data for pair-end sequencing, which is optional depending on the data type. In our case, as we utilized pair-end data, we utilized the PEAR tool [10] for the "Merged data" step. Subsequently, conducting additional rRNA filtering is advisable, despite standard wet lab protocols typically including rRNA removal. For this purpose, the SortMeRNA tool [11] proved to be useful (see Fig 1, part 1*).

Quality in base pairs is influenced by their position in the read, leading to lower average quality in later cycles of the sequencing process. To enhance read mapping rates, a common strategy involves removing low-quality bases through quality trimming. The Trimmomatic tool [12] is commonly employed for this task. Following quality trimming and adapter removal, FastQC is rerun to validate the improvements.

The last step involves aligning reads (see Fig 1. part 2^*) to a reference genome – TAIR10 [13]. The choice of aligner depends on the type of reference available, with STAR [14]

being recommended for genome-based alignment of RNA-Seq data. This alignment process requires a prepared index genome. Once prepared, sample reads can be aligned to it, resulting in the creation of a SAM and BAM format.



Fig. 1. Visualization of RNA-Seq Data Pre-processing Pipeline: This visual guide illustrates the pre-processing pipeline implemented in the study. The pipeline is divided into three color-coded sections. The orange section represents the foundational steps of the pipeline, the blue section indicates potential follow-up analyses, and the yellow section focuses on the main aspect of this study, which is the differential expression analysis. In general, the entire pipeline includes steps such as quality control, read trimming, alignment to a reference genome, transcript quantification, and normalization.

The Fig 1 outlines the sequential steps involved in preparing raw RNA-Seq data for downstream analysis. The main focus of the presented pipeline is the generation of a count table containing information on differential gene expression.

III. RESTULS AND DISCUSSION

In this section, we delve into a pivotal stage within our preprocessing pipeline, which significantly contributes to the robustness and trustworthiness of our RNA-Seq data analysis in *Arabidopsis thaliana*. This intermediate step plays a crucial role in ensuring the accuracy and integrity of the subsequent analytical procedures.

Our pre-processing pipeline meticulously handles a comprehensive set of 21 raw datasets, each comprising a total of 7 samples, with each sample having 3 replicates. To fortify our analyses and account for variability, each sample is meticulously replicated three times. This approach not only enhances the statistical power of our study but also enables us to confidently discern genuine biological signals from potential artifacts or noise. By meticulously curating and standardizing our datasets through this pre-processing pipeline, we establish a solid foundation for subsequent analyses, ensuring that our findings are robust, reproducible, and reflective of the true biological phenomena underlying gene expression in *Arabidopsis thaliana*.

A. Quality control before and after pre-processing

Firstly, the study present rigorous quality control measures before and after pre-processing to ensure the reliability and integrity of the RNA-Seq data. Quality control assessments, such as those conducted with FastQC, enabled us to identify potential issues such as overrepresented sequences and sequence quality of GC distribution, see Fig 2 (before), and Fig 3 (after).



Fig. 2. Quality Control (QG) report represented GC distribution before preprocessing pipeline



Fig. 3. Quality Control (QC) report represented GC distribution after preprocessing pipeline

After the application of SortMeRNA and Trimmomatic tool, we observed an approximation of the GC distribution curve towards a normal distribution. This shift suggests an improvement in the quality of the RNA-Seq data following preprocessing.

The initial GC distribution, as depicted in Fig. 2, exhibited deviations from the expected normal distribution, indicating potential biases or artifacts in the sequencing data. However, post-processing, as illustrated in Fig. 3, the GC distribution

curve appears to align more closely with the expected normal distribution, reflecting a reduction in sequencing errors and an enhancement in data quality. This normalization of the GC distribution is indicative of the effectiveness of the preprocessing steps, particularly in mitigating biases and improving the reliability of the RNA-Seq data for downstream analyses.

B. Comparison of filtering rRNA using defualt database and special A. thaliana database

Subsequently, we applied various pre-processing techniques, including rRNA filtering using SortMeRNA ($1^* - 1.1$). based on default rRNA eukaryote database offers SortMeRNA and 1.2. based on special rRNA database for *A. thaliana* available on RNACetnral database [15]) and quality trimming with Trimmomatic, to enhance the quality of the sequencing data.

These steps resulted in improved sequence quality metrics and enhanced the overall reliability of the data for downstream analyses.

C. Final Count Table: Transcript Abundance Analysis

The final output was generated using featureCount [16], and R\DESeq normalization [17] was applied to it. In our analysis of transcript abundance, we focused on the raw count data for three specific samples associated with Col-0, as depicted in Table I. These counts, obtained from replicates 1, 2, and 3, offer insight into the variability of transcript abundance across different samples.

TABLE I. EXAMPLE OF RESULTS TRANSCRIPT ABUNDANCE FOR COL-0

Sample ID	Merged data	Forward	Reverse
SRR16892914	33207850	17151179	16056671
SRR16892916	43639467	22380486	21258981
SRR16892917	31546058	16114941	15431117

Following this, we proceeded to assess transcript abundance before and after normalization for three selected samples – see Table II and Table III. In Table II, we present the transcript abundance after normalization for the first replicates. The normalized counts provide a more accurate representation of gene expression levels, facilitating robust downstream analyses.

 TABLE II.
 EXAMPLE OF RESULTS: TRANSCRIPT ABUNDANCE

 BEFORE NORMALIZATION – FIRST REPLICATES

	SRR16892916	SRR16892917	SRR16892918
AT1G01010	148	146	104
AT1G01020	358	340	225
AT1G03987	4	3	1
AT1G01030	126	115	82
AT1G01040	1944	1835	1932

Following normalization, which is crucial for removing systematic biases and enabling fair comparisons across samples,

we proceeded to assess transcript abundance for three selected samples. Normalization adjusts the raw count data to account for differences in sequencing depth and other technical factors, ensuring that gene expression levels are accurately represented. In Table II, we present the transcript abundance before normalization for the first replicates. These raw counts provide an initial insight into gene expression levels but may be influenced by technical variation.

Subsequently, in Table III, we present example of the transcript abundance after normalization for the first replicates. The normalized counts provide a more accurate representation of gene expression levels, as they have been adjusted to account for differences in sequencing depth and other technical factors. This normalization step enables robust downstream analyses, allowing for meaningful comparisons of gene expression levels between samples.

 TABLE III.
 EXAMPLE OF RESULTS: TRANSCRIPT ABUNDANCE

 AFTER NORMALIZATION – FIRST REPLICATES

	SRR16892916	SRR16892917	SRR16892918
AT1G01010	7.065	7.084	7.038
AT1G01020	8.314	8.314	8.207
AT1G03987	1.052	1.045	1.031
AT1G01030	6.783	6.764	6.723
AT1G01040	10.719	10.715	10.960

IV. CONCLUSSION

In conclusion, our study has successfully implemented a robust pre-processing pipeline specifically tailored for RNA-Seq data analysis in *Arabidopsis thaliana*. By incorporating stringent quality control measures and leveraging specialized databases for rRNA filtering, we have fortified the reliability and fidelity of the sequencing data. These outcomes underscore the pivotal importance of effective pre-processing methodologies in facilitating precise interpretation of RNA-Seq data and propelling advancements in plant molecular biology research.

Moreover, our application of this gold standard approach to transcriptomic data holds significant promise in bridging the gaps between genomic and metabolomic analyses within the panOMICs platform. This integrated approach not only enhances our understanding of the complex regulatory networks underlying plant biology, but also offers a novel insight into the interconnectedness of various molecular processes.

Ultimately, our research endeavors to contribute to the broader landscape of systems biology and pave the way for transformative discoveries in plant science. By elucidating the intricate molecular mechanisms at play in *Arabidopsis thaliana*, we aim to provide valuable insights that will inform future studies and drive innovation in the field of plant molecular biology.

ACKNOWLEDGMENT

This work has been supported by grant project FEKT/FIT-J23-8274.

LM2015042 and the CERIT Scientific Cloud LM2015085, provided under the programme "Projects of Large Research, Development, and Innovations Infrastructures.

REFERENCES

- Rapaport, F., Khanin, R., Liang, Y., Pirun, M., Krek, A., Zumbo, P., Mason, C.E., Socci, N.D. and Betel, D., 2013. Comprehensive evaluation of differential gene expression analysis methods for RNA-seq data. Genome biology, 14, pp.1-13.
- [2] Meinke, D.W., Cherry, J.M., Dean, C., Rounsley, S.D. and Koornneef, M., 1998. Arabidopsis thaliana: a model plant for genome analysis. Science, 282(5389), pp.662-682.
- [3] SCHWARZEROVÁ, J. Metabolite Genome-wide Association Studies of Arabidopsis Thaliana. In Proceedings of the 27th Conference STUDENT EEICT 2021 selected papers. 1. Brno: Brno University of Technology, Faculty of Electrical Engineering and Communication, 2021. s. 41-44. ISBN: 978-80-214-5943-4.
- [4] SCHWARZEROVÁ, J. Reproducible analytical pipeline for using raw RNA-Seq data from non-model organisms. Proceedings of the 26th Conference STUDENT EEICT 2020. Brno: Vysoké učení technické v Brně, Fakulta elektrotechniky a komunikačních technologií, 2020. s. 225-228. ISBN: 978-80-214-5867-3.
- [5] Weiszmann, Jakob, et al. "Metabolome plasticity in 241 Arabidopsis thaliana accessions reveals evolutionary cold adaptation processes." Plant Physiology 193.2 (2023): 980-1000.
- [6] SCHWARZEROVÁ, Jana; BARTOŇ, Vojtěch; WALTHER, Dirk a WECKWERTH, Wolfram. Comprehensive analysis of putrescine metabolism in A thaliana using GWAS, genetic risk score, metabolic modelling and data mining. Online. In: Proceedings II of the 29st Conference STUDENT EEICT 2023: Selected papers. 2023, p.151-155. ISBN 978-80-214-6154-3. ISSN 2788-1334.
- [7] Zhong, Z., Wang, Y., Wang, M., Yang, F., Thomas, Q.A., Xue, Y., Zhang, Y., Liu, W., Jami-Alahmadi, Y., Xu, L. and Feng, S., 2022. Histone chaperone ASF1 mediates H3. 3-H4 deposition in Arabidopsis. Nature communications, 13(1), p.6970.
- [8] Delhomme, N., Mähler, N., Schiffthaler, B., Sundell, D., Mannapperuma, C., Hvidsten, T.R. and Street, N.R., 2014. Guidelines for RNA-Seq data analysis. Epigenesys Protoc, 67(1-738), p.24.
- [9] Wingett, S.W. and Andrews, S., 2018. FastQ Screen: A tool for multigenome mapping and quality control. F1000Research, 7.
- [10] Zhang, J., Kobert, K., Flouri, T. and Stamatakis, A., 2014. PEAR: a fast and accurate Illumina Paired-End reAd mergeR. Bioinformatics, 30(5), pp.614-620
- [11] Kopylova, E., Noé, L. and Touzet, H., 2012. SortMeRNA: fast and accurate filtering of ribosomal RNAs in metatranscriptomic data. Bioinformatics, 28(24), pp.3211-3217.
- [12] Bolger, A.M., Lohse, M. and Usadel, B., 2014. Trimmomatic: a flexible trimmer for Illumina sequence data. Bioinformatics, 30(15), pp.2114-2120.
- [13] Lamesch, P., Berardini, T.Z., Li, D., Swarbreck, D., Wilks, C., Sasidharan, R., Muller, R., Dreher, K., Alexander, D.L., Garcia-Hernandez, M. and Karthikeyan, A.S., 2012. The Arabidopsis Information Resource (TAIR): improved gene annotation and new tools. Nucleic acids research, 40(D1), pp.D1202-D1210.
- [14] Dobin, A., Davis, C.A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P., Chaisson, M. and Gingeras, T.R., 2013. STAR: ultrafast universal RNA-seq aligner. Bioinformatics, 29(1), pp.15-21.
- [15] "RNAcentral: an international database of ncRNA sequences." Nucleic acids research 43, no. D1 (2015): D123-D129.
- [16] Liao, Y., Smyth, G.K. and Shi, W., 2014. featureCounts: an efficient general purpose program for assigning sequence reads to genomic features. Bioinformatics, 30(7), pp.923-930.
- [17] Anders, S. and Huber, W., 2012. Differential expression of RNA-Seq data at the gene level–the DESeq package. Heidelberg, Germany: European Molecular Biology Laboratory (EMBL), 10, p.f1000research.

Quantitative Analysis of Vocal Tract Resonances in Patients with Parkinson's Disease

1st Daniel Kováč Department of Telecommunications Brno University of Technology Brno, Czech Republic xkovac41@vut.cz

Abstract-Parkinson's disease (PD) is often associated with hypokinetic dysarthria, impacting speech-motor function and resulting in articulatory deficits such as reduced speech intelligibility due to stiffness in articulatory structures. This study aims to introduce novel speech metrics to assess articulation impairments and differentiate between healthy individuals and those with PD. Resonant frequency attenuation (RFA) was confirmed as a valid measure, and two additional acoustic features based on supraglottic resonances were proposed. Their efficacy in distinguishing between groups was evaluated using a dataset comprising 19 healthy controls (HC) and 28 PD patients. Significantly divergent values between PD and HC were observed for all three features via Student's t-test. Employing logistic regression with leave-one-out cross-validation, PD classification achieved 96% sensitivity and 100% specificity. This preliminary study suggests that supraglottic resonances could be pivotal in aiding PD diagnosis through acoustic speech analysis.

Index Terms—Resonance Frequency Attenuation, Articulatory Decay, Parkinson's Disease, Hypokinetic Dysarthria

I. INTRODUCTION

Parkinson's disease (PD) is a chronic neurodegenerative disorder characterized by progressive degeneration of dopaminergic neurons in the substantia nigra [1]. The dopamine deficiency in this part of the midbrain leads to difficulties associated with both motor and non-motor functions, such as sleep disturbances, speech impairments, emotional changes, and a range of other deficits [2]. The risk of PD increases with age, with other significant factors including gender and genetic predisposition. It is also evident that the number of affected individuals increases with the growing population and life expectancy [3]. Early detection and initiation of treatment significantly impact the disease's progression and the quality of life of patients.

According to Hoehn and Yahr, the progression of PD manifests in five stages, with speech abnormalities appearing as early as the second phase of the disease [4]. These speech disorders, collectively known as hypokinetic dysarthria (HD), can manifest in phonation, articulation, prosody, and respiration, impacting up to 90% of PD patients [5].

This work analyses articulation disorders manifested as resonance frequency attenuation (RFA) in the spectral envelope of the vocal tract of PD patients. Rusz et al. highlighted the acoustic speech feature RFA in their publication, which proves to be highly effective in detecting PD in both early-onset (PD onset before the age of 50) and late-onset (PD onset after the age of 70) patients, demonstrating overall independence from the speaker's age. This feature is defined as the difference between the second formant region's maximum and the local valley's minimum, known as the antiformant [6]. In another publication, this feature shows a significant difference between healthy controls (HC) and PD patients. Logistic regression modelling achieved an area under the receiver operating characteristic curve (AUC-ROC) of 0.86 for men and 0.93 for women [7]. Tykalová et al. statistically analyzed this feature in two subtypes of PD. According to the results of this study, the feature did not exhibit a significant difference between patients with postural instability/gait difficulties (PIGD) and those with dominant tremor (TD). It differed from healthy controls only in patients with PIGD [8].

This work aims to verify the results of the mentioned studies, expand the concept with new features, and subsequently test and analyze them in detail.

II. METHOD

A. Testing Database

The effectiveness of the features in discriminating patients from HC was tested on a database of Italian speakers, consisting of 47 speech recordings [9]. The database contains HC and PD patients on medication (ON state); their distribution in the dataset can be observed in Figure 1.



The average score on the Unified Parkinson's Disease Rating Scale (UPDRS-III) for speech in PD patients is $\bar{x} = 1.0$

This work was supported by the Quality Internal Grants of BUT (project KInG, reg. no. CZ.02.2.69/0.0/0.0/19_073/0016948; financed from OP VVV), grant no. NU20-04-00294 of the Czech Ministry of Health, and EU – Next Generation EU (project no. LX22NPO5107 (MEYS)). Many thanks to Dr Mekyska, head of the Brain Diseases Analysis Laboratory (BDALab).

with a standard deviation of $\sigma = 1.2$, where 0 indicates normal speech, and 4 indicates unintelligible speech. According to the Mann-Whitney U test, the difference in age is not significant. The average age is 67 for the HC group and 66 for PD patients. The speaker's task was to read a text, with the average recording length being approximately 1 minute.

B. Design of Acoustic Features

The algorithm is implemented in MATLAB. First, a signal in the .wav format is loaded. If the signal was not originally sampled at $f_s = 16$ kHz, it is resampled to this value. Then, the signal offset is removed using the internal function detrend. Subsequently, the signal is segmented with a window length of 25 ms and an overlap of 10 ms.

1) Linear Prediction Coefficients: From each segment of the signal, linear prediction coefficients a_p are computed using forward linear prediction:

$$\widetilde{x}[n] = \sum_{i=1}^{p} a_{\mathbf{p}}[i] \cdot x[n-i], \qquad (1)$$

where $\tilde{x}[n]$ is the predicted value of the signal at discrete time n. The variable i denotes the delay, and p is the order of linear prediction, set to p = 12 in this case. The forward linear prediction error $f_p[n]$ is then the difference between the current value and the predicted value:

$$f_{\mathbf{p}}[n] = x[n] - \widetilde{x}[n], \qquad (2)$$

and the mean squared error ϵ_p^f is given by:

$$\epsilon_{\rm p}^f = \frac{1}{M} \sum_{n=0}^{M-1} f_{\rm p}^2[n].$$
 (3)

The variable M denotes the length of the signal in samples. The goal is to find linear prediction coefficients that minimize this mean squared error:

$$\frac{\partial \epsilon_{\rm p}^f}{\partial a_{\rm p}[i]} = 0. \tag{4}$$

This leads to a system of normal equations described as an autocorrelation Toeplitz matrix solved using the Levinson-Durbin recursive method [10]. To compute the coefficients in MATLAB, the lpc function from the publicly available COVAREP repository (v1.4.1) was utilized [11].

2) *Transfer Function of the Human Vocal Tract:* From the obtained linear prediction coefficients, we can model the transfer function of the human vocal tract as follows:

$$H(z) = \frac{1}{\sum_{i=0}^{p} a_p[i] \cdot z^{-i}},$$
(5)

where $a_p[0] = 1$ and $z = e^{j2\pi \mathbf{f_d}}$. The frequency vector $\mathbf{f_d}$ is the result of discretizing the sequence:

$$\mathbf{f} = 1, 2, 3..., \frac{f_s}{2}$$
 (6)

with the sampling frequency f_s . The denominator of the transfer function is a polynomial, and the solutions of the

polynomial equation are roots, which we call poles; they have a complex form and are associated with resonances. These resonances arise in the supra-glottal cavities and correspond to speech formants. In MATLAB, the roots function is used to find the roots of the equation. Only poles with positive imaginary parts are considered (resonances mirrored across half of the sampling frequency are not considered). The magnitude is given by the absolute value of this complex number, and the frequency at which the resonance occurs is obtained from the phase. An example of pole distribution in the unit circle and the corresponding magnitude frequency response of the vocal tract for a specific segment is shown in Figure 2.



Fig. 2. Unit Circle with Pole-Zero Plot (top) and Amplitude Frequency Response of the Vocal Tract (bottom).

3) *Feature Extraction:* In the amplitude frequency response of the vocal tract, local maxima and minima are found using the findpeaks function. Then, the computation of the acoustic features is performed:

• RFA1

The difference between the amplitude frequency response value at the second formant and the value at the first local minimum (between the first and second formants). In Figure 2 (bottom), this feature corresponds to the vertical distance between the second formant and the first local minimum. This feature was mentioned in the introduction of the article.

• RFA2

The difference between the amplitude frequency response value at the second formant and the value at the second local minimum (between the second and third formants). In Figure 2 (bottom), this feature corresponds to the vertical distance between the second formant and the second local minimum.

#locMAX

The number of local maxima in the amplitude frequency

response of the vocal tract.

These features are computed for each voiced segment, and their average values are considered.

4) Voiced Segment Identification: Features are extracted only from voiced segments of speech. Voicing is verified by the following conditions:

- The fundamental speech frequency is between 75 Hz and 400 Hz. This frequency is computed using PRAAT software [12].
- The frequency of the second formant is below 3 kHz.
- There exists at least one local minimum.
- The magnitude of the first formant is higher than that of the third.
- The lowest magnitude among the first three formants is higher than the highest magnitude among the remaining formants.

The feature extraction algorithm is available online: https://github.com/BDALab/Articulatory_decay.

C. Statistical Analysis

The extracted acoustic features underwent a comprehensive statistical analysis. The Shapiro-Wilk test was used to verify the normal distribution of features in HC group. Since the features exhibited a normal distribution, z-scores were subsequently calculated to determine the number of patients deviating from the norm established by HC (values of PD patients were normalized based on the mean and standard deviation of HC):

$$\boldsymbol{z}_{\rm HC} = \frac{\boldsymbol{p}_{\rm HC} - \bar{\boldsymbol{x}}_{\rm HC}}{\sigma_{\rm HC}},\tag{7}$$

$$\boldsymbol{z}_{\rm PD} = \frac{\boldsymbol{p}_{\rm PD} - \bar{\boldsymbol{x}}_{\rm HC}}{\sigma_{\rm HC}}.$$
 (8)

The vectors p_{HC} and p_{PD} represent the feature values of HC and PD, respectively, and the vectors z_{HC} and z_{PD} represent the resulting z-scores. The number of PD patients whose feature z-scores lie more than two standard deviations from the mean of HC indicates the strength of the feature to discriminate and thus detect PD. The Student's t-test then determines whether the values of features in HC and PD patients come from the same distribution. Pearson's correlation coefficient was used to assess the dependence of features on the age of HC and the degree of correlation between the PD patient's features and their clinical data (UPDRS-III).

D. Logistic Regression

Mathematical modelling of the features was performed using the Python programming language. Acoustic features were used as input independent variables for training machine learning models. Logistic regression with L1 regularization was chosen as the classifier. Training data were balanced using the Synthetic Minority Over-sampling Technique (SMOTE). Using grid search, hyperparameters were tuned to achieve the best possible balanced accuracy, and leave-one-out crossvalidation was used to validate the models. Evaluation metrics for assessing the model's performance include the Matthew's correlation coefficient (MCC), accuracy (ACC), sensitivity (SEN), and specificity (SPE).

III. RESULTS

A. Statistical Analysis

Figure 3 depicts the probability distributions of the z-scores for each feature across the entire dataset of speakers (divided into HC and PD speakers). The red dots highlight patients who deviate from HC, with their proportion to the total number of patients indicated.

Table I describes the changes in the average value \bar{x} , median \tilde{x} , and standard deviation σ of the feature in PD speakers compared to HC.

TABLE I FEATURES' DESCRIPTIVE STATISTICS.

	\bar{x}_{HC}	\bar{x}_{PD}	\tilde{x}_{HC}	\tilde{x}_{PD}	σ_{HC}	σ_{PD}
RFA1 [dB]	3.36	1.64	3.17	1.71	0.87	0.89
RFA2 [dB]	12.75	10.92	12.69	10.72	1.90	2.17
#locMAX	4.05	3.69	4.03	3.74	0.16	0.21

In Table II, the results of the statistical tests are presented. The Shapiro-Wilk test checks the null hypothesis that the acoustic features extracted from the speech signal of HC originate from a normal probability distribution. The null hypothesis of the Student's t-test is that the features of PD patients and HC have insignificantly different means.

TABLE II Results of Statistical Tests.

	Shapiro-Wilk Test	Student's T-test
	P-value	P-value
RFA1	0.18	< 0.01
RFA2	0.7	0.01
#locMAX	0.35	< 0.01

Table III describes the results of correlation tests. The null hypothesis states that there is no linear relationship between the HC features and their age. The null hypothesis of the second test is set such that there is no linear relationship between the features of PD speakers and their scores on the UPDRS-III scale (speech).

TABLE III RESULTS OF CORRELATION TESTS.

	HC A	.ge	UPDRSIII	(speech)
	Coefficient	P-value	Coefficient	P-value
RFA1	0.21	0.38	-0.57	< 0.01
RFA2	0.27	0.27	-0.18	0.37
#locMAX	0.26	0.28	0.05	0.81

B. Logistic Regression

The performance of the model in classifying speakers as HC or PD based on individual features is listed in Table IV. The table also includes the model's performance when all independent features were contained. Figure 4 displays the ROC curve of this model, which expresses the relationship



Fig. 3. Probability Distributions of Acoustic Features and Deviations of PD Patients from Norms given by Healthy Controls.

between sensitivity and specificity. The area under the curve (AUC = 0.97) indicates the overall classifier's performance. The bottom graph also shows the probability of individual samples classified by the model as HC or PD speakers. The decision threshold is optimal for a p = 0.5 probability.

TABLE IV MODELS' PERFORMANCES.

	MCC	ACC [%]	SEN [%]	SPE [%]
RFA1	0.86	93	89	96
RFA2	0.43	71	75	68
#locMAX	0.79	89	89	89
together	0.96	98	96	100

IV. DISCUSSION

The RFA acoustic feature was verified, and two additional quantifying articulation disorders in PD patients were proposed. Subsequently, they were statistically analyzed, and their sensitivity/specificity in classifying PD patients in the Italian dataset was determined.

The results of statistical analyses indicate significant differences between features extracted from HC and PD recordings. Both the mean and median of all features were lower in PD patients. They exhibited a resonance attenuation at the second formant, which is directly related to the position of the speaker's tongue, suggesting increased stiffness [13]. Reduced values of the #locMAX feature indicate an overall attenuation of all speech formants. The RFA1 feature shows a linear relationship (negative correlation) with PD patients' UPDRS-III (speech) scores with a significance level of p = 0.01. Hence, this feature might be associated with speech intelligibility. In the second test, the hypothesis that features depend on the age of HC was not rejected. This is consistent with the results published by Rusz et al. [6]. According to the Shapiro-Wilk test, all three features of HC originate from a normal distribution. This allows us to normalize data to z-scores and



Fig. 4. ROC Curve (top) and Classification Graph; (TN = True Negative; FN = False Negative; TP = True Positive).

directly identify patients deviating from norms. For the RFA1 feature, it is 46% of PD patients below two standard deviations from the HC mean. With the RFA2 feature, there are 11% of PD patients outside the norm, and with #locMAX, it is 41%. An interesting phenomenon occurs with the RFA2, where one person deviates from the HC norm in the opposite direction. This can be explained by the fact that in individuals with PD, hypokinetic dysarthria may manifest only in some dimensions of speech [14]. This is also supported by the fact that the standard deviation increased for all features in PD patients compared to HC. Moreover, the Student's t-test results show that HC and PD features originate from different probability distributions. For the above reasons, all three can be considered suitable for quantifying articulation disorders in PD patients.

The classification performance is highest for the RFA1 feature, with the model achieving an accuracy of 93%, which is increased to 98% when all three were employed in the logistic regression.

V. CONCLUSION

This article focuses on quantifying articulation disorders based on the resonance frequency attenuation in the spectral envelope of the vocal tract. Three acoustic features based on the linear predictive coding (RFA1, RFA2, and #locMAX) were proposed and evaluated on a database of Italian-speaking HC and PD patients. All three features exhibited normal distribution in the HC group and independence from the speaker's age. Additionally, the RFA1 feature correlated with PD patients' UPDRS-III (speech) scores. Based on statistical analysis, all three features appear suitable for the automatic assessment of articulation disorders associated with PD. A machine learning model based on logistic regression achieved a classification accuracy of 93% using the RFA1 feature, which increased to 98% after adding the remaining two features as inputs to the model.

The study's drawback is the limited number of speakers. Integrating L2 regularization alongside L1 could mitigate overfitting and improve the model's ability to generalize. Additionally, expanding the dataset to include more clinical data could enhance the robustness of the findings. While the research primarily concentrates on attenuation of the second speech formant, broadening the scope to analyze formants across the entire speech spectrum could yield deeper insights into articulation challenges among PD patients.

REFERENCES

- O. Hornykiewicz, "Biochemical aspects of parkinson's disease," *Neurology*, vol. 51, no. 2_suppl_2, pp. S2–S9, 1998.
- [2] H. Narabayashi, "The neural mechanisms and progressive nature of symptoms of parkinson's disease—based on clinical, neurophysiological and morphological studies," *Journal of neural transmission-Parkinson's disease and dementia section*, vol. 10, pp. 63–75, 1995.
- [3] G. DeMaagd and A. Philip, "Parkinson's disease and its management: part 1: disease entity, risk factors, pathophysiology, clinical presentation, and diagnosis," *Pharmacy and therapeutics*, vol. 40, no. 8, p. 504, 2015.
- [4] R. Bhidayasiri, D. Tarsy, R. Bhidayasiri, and D. Tarsy, "Parkinson's disease: Hoehn and yahr scale," *Movement disorders: a video atlas: a video atlas*, pp. 4–5, 2012.
- [5] A. K. Ho, R. Iansek, C. Marigliani, J. L. Bradshaw, and S. Gates, "Speech impairment in a large sample of patients with parkinson's disease," *Behavioural neurology*, vol. 11, no. 3, pp. 131–137, 1998.
- [6] J. Rusz, T. Tykalova, M. Novotny, E. Ruzicka, and P. Dusek, "Distinct patterns of speech disorder in early-onset and late-onset de-novo parkinson's disease," *npj Parkinson's Disease*, vol. 7, no. 1, p. 98, 2021.
- [7] J. Rusz, T. Tykalova, M. Novotny, D. Zogala, E. Ruzicka, and P. Dusek, "Automated speech analysis in early untreated parkinson's disease: Relation to gender and dopaminergic transporter imaging," *European journal of neurology*, vol. 29, no. 1, pp. 81–90, 2022.
- [8] T. Tykalová, J. Rusz, J. Švihlík, S. Bancone, A. Spezia, and M. T. Pellecchia, "Speech disorder and vocal tremor in postural instability/gait difficulty and tremor dominant subtypes of parkinson's disease," *Journal of Neural Transmission*, vol. 127, no. 9, pp. 1295–1304, 2020.
- [9] G. Dimauro, F. Girardi *et al.*, "Italian parkinson's voice and speech," 2019.
- [10] Z. Smékal, Analog and Digital Signal Processing in Examples and Programs. Brno University of Technology, VUTIUM Press, 2021.
- [11] G. Degottex, J. Kane, T. Drugman, T. Raitio, and S. Scherer, "Covarep—a collaborative voice analysis repository for speech technologies," in 2014 ieee international conference on acoustics, speech and signal processing (icassp). IEEE, 2014, pp. 960–964.
- [12] P. Boersma, "Praat, a system for doing phonetics by computer," in *Glot International 5:9/10, 341-345.*
- [13] J. Lee, S. Shaiman, and G. Weismer, "Relationship between tongue positions and formant frequencies in female speakers," *The Journal of the Acoustical Society of America*, vol. 139, no. 1, pp. 426–440, 2016.
- [14] J. Rusz, T. Tykalova, M. Novotny, D. Zogala, K. Sonka, E. Ruzicka, and P. Dusek, "Defining speech subtypes in de novo parkinson disease: response to long-term levodopa therapy," *Neurology*, vol. 97, no. 21, pp. e2124–e2135, 2021.

Long-Term Effect of Repetitive Transcranial Magnetic Stimulation on Parkinson's Disease Patients with Different Severity of Hypokinetic Dysarthria

Krystof Novotny Department of Telecommunications Brno University of Technology Brno, Czech Republic ORCID: 0009-0005-7232-0841

Abstract—The prevalence of Parkinson's disease (PD), severe neurodegenerative disorder, has steadily increased. Among the symptoms of PD is hypokinetic dysarthria (HD), a motor speech disorder, characterised by respiratory, articulatory, prosodic and phonatory impairments. It has been demonstrated that both motor and non-motor symptoms of PD can be improved using the repetitive Transcranial Magnetic Stimulation (rTMS). This study analyses acoustic speech characteristics of 19 participants diagnosed with PD before (one pre-stimulus) and after (four poststimulus) evaluation sessions of rTMS treatment. The participants were divided into two groups - receiving either rTMS or sham stimulation (1:1 randomization). Based on the prestimulus subresults of the Test 3F, participants were stratified into two cohorts, according to their possible HD severity level. Speech recordings were also taken during each evaluation session. The outcome of the follow-up acoustic analysis resulted in 16 parameters for each of those sessions. Their evaluation demonstrated the dependence of the effect of rTMS treatment on the severity level. The actively stimulated group of the first cohort showed consistent improvement in articulation and prosody (sham did not) while the actively stimulated group of the second stratified cohort showed consistent improvement in phonation (sham did not). The study provides early preliminary insights into the benefits of rTMS for the alleviation of HD manifestations (symptomatic treatment of PD). In addition, it provides new insights into the possible relationship between the effectiveness of rTMS and the degree of severity in HD.

Index Terms—Parkinson's disease, hypokinetic dysarthria, repetitive transcranial magnetic stimulation, acoustic analysis, digital speech and voice biomarkers

I. INTRODUCTION

Parkinson's disease (PD) is the second most common neurodegenerative disorder worldwide with an age-standardised prevalence rate of 106.28 per 100,000 cases in 2019 (94.09 for Central Europe and 126.01 for Western Europe). Moreover, the number is following an increasing trend in most parts of the world [1]. The prevalence rates are then significantly higher for the older population [2].

The disease is caused by a malfunction of the motor loop of the basal ganglia and is manifested by tremors, muscle stiffness and bradykinesia [3]. Thus, it is a motor disorder. This also has a significant impact on the speech and voice of patients. Up to 90 % of them have a motor speech disorder known as hypokinetic dysarthria (HD) [4]. Its manifestations include monotony of the fundamental pitch and loudness of speech, a blurred hoarse voice, unnatural and inconsistent intonation, incorrect phrasing, inappropriate pauses and sudden changes in speech rate [5]. These impairments have an undeniable impact on patients' quality of life.

Apart from pharmacological and surgical treatments [6], one of the possible solutions is the repetitive Transcranial Magnetic Stimulation (rTMS). Studies on patients with PD have shown that the superior temporal gyrus (STG) is the part of the brain that is responsible for projecting the motor aspects of speech onto the voice during its production. It also implements prosody and manages auditory feedback processing [7]. Excitation of neurons in this area by rapid magnetic field changes (the principle of functioning of rTMS based on magnetic induction) can bring improvements in the named areas of speech [8]. The advantage of this approach is that it is one of the non-invasive brain stimulation methods it does not involve surgical procedures and has relatively mild side effects. Its effectiveness has already been demonstrated by multiple studies on the primary symptoms of PD [9], [10]. There are fewer relevant studies in the context of dysarthric speech [11].

The short-term effects of rTMS applied to the primary motor cortex on speech were investigated by Dias et al. The results of the speech task of the sustained vowel [a:] showed an improvement in both the fundamental pitch and the intensity of voice [12]. Excitation of the same brain region was also investigated by Hartelius et al. but the studied features quantifying sustained fricative, prolonged vowel phonation, diadochokinetic rate and sentence intelligibility did not yield clear results. In addition, a strong placebo effect was also mapped [13]. More positive results were obtained by Eliasova et al. when rTMS was performed on an adjacent region of the cerebral cortex – the primary sensorimotor cortex. The best results here were observed for the five sustained vowel tasks. There was a general improvement in voice quality and intensity, speech rate and tongue movement. This was also a study focused on short-term effects [14]. Brabenec et al. were then the first to provide evidence of improved articulation after STG excitation as part of a short-term effect. Significant improvement was achieved in parameters monitoring speech formants characteristic for tongue and jaw movements [15].

Even fewer studies have addressed the long-term effects of rTMS on speech. Moreover, it should be noted that they are all based on the same data collection (the excited region was STG). Brabenec et al. conducted the first of these. Speech is quantified based on the Test 3F. The study focuses mainly on the phonation region, which is the only one for which they observed a noticeable effect of rTMS. The improvement over baseline is long-lasting in the active group. However, so is the significant placebo effect in the sham group [16]. Gomez-Rodellar et al. involve multiple speech biomarkers and observe long-term improvement in jitter, cepstral peak prominence and selected features quantifying the tremor of vocal cord tension (divided into frequency bands corresponding to those standardly used in electroencephalography) [17]. Next, Brabenec et al. report noticeable difference in the active group (compared to sham) for parameters describing the left anterior arcuate fasciculus. This is the region connecting the auditory feedback area with the motor regions involved in articulation. Furthermore, the values correlate with the time evolution of the phonetic subscore of the Test 3F [18]. Subsequently, Gomez-Rodellar et al. expand knowledge in the relationship between biomechanical correlates of the vocal cord tension and rTMS therapy [19]. However, all longitudinal studies have yielded mixed results. Improvement is not uniform across the active group and the sham group often shows a placebo effect.

It is hypothesized that the effectiveness of rTMS therapy might be related to the varying severity of HD. Steurer et al. reported different behaviors of dysarthric speech symptoms in a cohort of 83 PD patients. Three different cohorts showing distinct speech characteristics were examined: a group with no HD, a group with mild HD, and a group with moderate HD. Furthermore, within these, they also documented correlations of the digital speech and voice biomarkers with the results of other clinical tests and screenings [20].

The aim of this paper is to stratify PD patients based on the Test 3F into different cohorts possibly representing the severity level of HD. Subsequently, digital speech and voice biomarkers should be calculated from the recordings of each session of all patients. Using one session completed by healthy controls (HC), the relationship of the calculated parameters to normative values should be determined. The goal is to investigate the evolution of each feature over time after active treatment (or sham stimulation), taking into account the severity level of HD (the outcome of the stratification).

As a result, this study aims to be the first ever to explore the possible influence of severity level of HD on rTMS therapy outcomes. This could shed light on new findings in the treatment of HD.

II. MATERIALS AND METHODS

A. Database

The input data for this work were audio recordings from the HIDI database, which contains a total of 19 PD patients (14 females/5 males, mean age 71.38 ± 7.43). The group was divided into two parts. The first group (10 participants, mean age 71.83 ± 6.69) received active stimulation (labeled "STG") according to the protocol described in [16] during ten sessions spread over two weeks at the Central European Institute of Technology, Masaryk University. The second group (9 participants, mean age 70.87 \pm 8.57), underwent a similar process using the same device, but here no magnetic field was applied (labelled "SHAM"). Audio recordings of each participant were acquired during five sessions. Pre-stimulus T0 before undergoing the actual stimulation, post-stimulus T1 two weeks after, post-stimulus T2 six weeks after, poststimulus T3 10 weeks after and post-stimulus T4 14 weeks after the stimulation. Among other examinations, patients also completed the Test 3F at each session. The database was created under the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement no. 734718 (CoBeN) and under the grant of the Czech Ministry of Health no. 16-30805A. All patients were recorded in the ON state.

At the same time, a dataset of HC was used to enable the calculation of normative values. It contains a total of 31 individuals (15 females/16 males, mean age 67.10 \pm 6.05) who completed the same speech tasks as the representatives from the patient database within one session. The recordings were obtained within the framework of the Ministry of Health project no. NU20-04-00294. The studies mentioned above were approved by the local Ethics Committee and all participants signed informed consent.

B. Speech processing

Two speech tasks were investigated – sustained phonation of the vowel [a:] (as long as possible), and a free speech monologue (at least 90 s long).

Maximum phonation time (MPT, the total length of the phonation) was calculated from the first task. The following parameters were also calculated from the first task, but after its adjustment – only the section of the recording starting at 2 s of its time and ending at 4 s was treated. Such parameters are: relative standard deviation of the first (relF1SD) and second (relF2SD) formant, relative standard deviation of the fundamental frequency (relF0SD), jitter (PPQ), shimmer (APQ) and harmonics-to-noise ratio (HNR).

For the second mentioned task, the relative standard deviation of the first (relF1SD) and second (relF2SD) formant and the relative standard deviation of the fundamental frequency (relF0SD) were extracted. The following features were also calculated from the second task, but after its modification by removing the silent regions that exceeded 250 ms. From the pauses removed, the median of their duration (DurMED) and the mean absolute deviation of their duration (DurMAD) were obtained, and the ratio of their number to the duration of the whole task (SPIR) was also computed. Similarly, after removing regions of silence exceeding 50 ms, the parameters relative standard deviation of speech energy (relSEOSD), energy evolution of speech (EEVOL) and percentual pause ratio (PPR) were extracted.

The values of all calculated parameters were adjusted so that the worsening of the manifestation of a particular symptom was tracked by increasing the parameter value. Parameters with the opposite manifestation (an increase in value meant an improvement) were adjusted by multiplying by -1.

In the next step, the values of all parameters (dimension a) across all sessions of all participants (dimension b) were collectively processed, including the results of one session of HC. The data prepared in this way were adjusted using linear regression. The effect of the age of the speakers at the time of each session and the effect of gender were removed.

From the set of values of HC, norms were established for each parameter – the median $(MED_{\rm HC})$ and the value corresponding to the 95th percentile $(PR_{95,\rm HC})$ were determined. Using these values, it was then possible to relate to the norms all calculated values of PD patients from all sessions for each parameter according to the following relationship:

$$DIST_{\text{PD},k,T} = \frac{VAL_{\text{PD},k,T} - MED_{\text{HC}}}{PR_{95,\text{HC}} - MED_{\text{HC}}}.$$
 (1)

The result measures the relative distance $(DIST_{PD})$ of the investigated feature value (VAL_{PD}) for a particular session T of the selected patient k from the norm determined by the HC.

C. Stratification

Based on the hypothesis of different effects of rTMS on HD according to its severity, participants were stratified. Input values for the k-means cluster analysis were the three Test 3F subscores obtained in T0. Using the silhouette method, it was decided that there were most likely to be two distinct cohorts in the dataset. Thus, the k-means algorithm identified 8 participants (4 STG and 4 SHAM) as the first group (labeled "0") and 11 participants (6 STG and 5 SHAM) as the second group (labeled "1"). The groups were then compared using the Mann–Whitney U (MWU) test applied to the individual Test 3F subscores. In general, group 0 values exceeded group 1 in all subscores. In addition, the test evaluated the differences as significant in all three cases (see Fig. 1).

D. Evaluation

The main part of the analysis relied on plotting the stacked grouped bar graphs. For each calculated parameter, one graph was plotted showing four clusters of results. Each cluster belongs to one cohort (STG 0, STG 1, SHAM 0, SHAM 1). Four results are then displayed within each cluster. These are percentages of cases where the cohort improved/deteriorated for that parameter compared to T0. Each column then considers the more strictly selected values compared to the previous one. Column one – improvement over T0 was achieved at least in one post-stimulus session, column two – improvement was



Fig. 1. Distribution of the Test 3F subscores within stratified groups and the results of the MWU test.

achieved at least in two post-stimulus sessions, column three – improvement was achieved at least in three post-stimulus sessions and column four – improvement over T0 was achieved in all post-stimulus sessions. The examined values are the data generated by equation (1) for adjustment with respect to norms.

Additionally, the median percentage improvement is also observed for each case. It is calculated from the set of all values where there was an improvement in relation to the norm compared to T0, for all patients meeting the minimum number of improvements condition.

III. RESULTS

The best result for stratified group 0 was yielded by relF1SD computed from the monologue task. In STG 0, there was an improvement in 50% of cases in all four post-stimulus sessions and in another 25% in three sessions (see Fig. 2). In SHAM 0, then, there was no improvement in relF1SD of the monologue in two or more sessions (only 75% in one session). Additional results of improvement (and discrimination from SHAM) for this group were then provided by the parameters PPR, SPIR and relF2SD calculated from the monologue task.

The HNR parameter delivered the best results for stratification group 1, with STG 1 showing improvement in all four post-stimulus sessions in 67 % of cases and in two sessions in the remaining 33 %. In contrast, the improvement among all sessions in SHAM group 1 was not achieved at all (only 20 % in three sessions and another 20 % in one session). Similar results with slight variations were obtained within group 1 for PPQ, APQ, MPT and relF0SD calculated from extended phonation.

IV. DISCUSSION

There is not much significant difference in the distribution of values of individual parameters in prestimulus sessions within the stratified groups. However, longitudinal analysis showed a difference in the effect of rTMS on a particular cohort.

The first two parameters for which the most consistent results in improvement were observed for group 0 both deal with the range of values of the selected formant in terms of the longer monologue. A larger range of values (expressed







Fig. 2. Resulting graphs for selected parameters.

in relative standard deviation) indicates a higher articulation ability. For the relF1SD parameter, this indirectly refers to the size of the pharyngeal cavity and thus to the mobility of the tongue root. For the relF2SD parameter (which does not yield as clear results as relF1SD, but similar nonetheless), on the other hand, it refers to the oral cavity, i.e. the openness of the lips and the position of the tongue in the mouth. The other two parameters (PPR and SPIR) with consistent improvement results in group 0 both relate to pauses in speech. In all four cases, for STG 0, there was a steady improvement in the parameter values in the majority of cases, whereas for SHAM 0 the improvements were rather sporadic. Both STG 1 and SHAM 1 gave inconsistent results in these cases.

The HNR parameter (monitoring the increase of the noise

in the voice) provides the clearest discrimination for group 1. Again, for this cohort it provides a clear differentiation between STG and SHAM, with STG showing improvement over prestimulus in most sessions for the majority of patients and SHAM showing rather isolated cases of improvement. In contrast, for group 0, it yields inconsistent non-differentiating outcomes. Similarly performing in favor of group 1 stratification are the parameters PPQ (microperturbations in frequency), APQ (microperturbations in amplitude), MPT (monitoring airflow insufficiency), and relFOSD calculated from the sustained vowel (irregularity of vocal fold vibration). All four of these features show similar patterns to HNR, with improvements for STG 1, rather sporadic improvements for SHAM 1, and mixed results for all of group 0.

Looking at the relationship between the selected parameters, a possible link emerges. The parameters describing the difference between STG and SHAM group 0 can be described as articulatory and prosodic, and the parameters describing the differences between STG and SHAM group 1 can all be seen as members of the phonation category. Thus, this longitudinal study suggests that the first stratified group can be viewed as a cohort of people in whom the influence of rTMS can lead to an improvement in articulation and prosody abilities, and the second stratified group as a cohort in which rTMS has a positive effect on phonation abilities.

Both short-term improvements in voice quality and tongue movement after the application of rTMS were observed by Eliasova et al. in [14]. Therefore, this result is consistent with the behaviour of both our groups 0 and 1. However, it should be mentioned here that in their research a different part of the brain was stimulated. Brabenec et al. then provided positive results regarding the short-term effect on speech formants in [15]. This matches the behavior of our group 0. In their case, the STG area was stimulated (this is a dataset partially identical to the one used here). Further research then provides consistent evidence from longitudinal studies of a positive effect of rTMS in the phonation domain (voice quality, vocal cord tremor...) [16], [17], [19]. This corresponds to the behaviour of our group 1.

A significant phenomenon present in this work is the placebo effect. A consistent long-term improvement is observed in a number of parameters (relSEOSD, EEVOL) not only for STG but also for SHAM. The hypothesis of a likely placebo effect is supported by the aforementioned studies that also dealt with it. In addition, an alternative explanation may be the phenomenon described in [21] and [22]. The studies focused on comparing patients' general communicative speech with speech where patients focused on being understood as clearly as possible. Noticeable differences were observed between the two styles of speech, indicating that PD patients may be able to produce more intelligible speech with sufficient concentration and thus minimize the classic features of dysarthric speech. Therefore, the results of our research may also be influenced by fluctuations in momentary concentration and patients' efforts to achieve higher intelligibility across different sessions.

However, the objectively biggest weakness of this work is the limited database. Moreover, when dealing with stratification and further subdividing the dataset into smaller cohorts, the low number of research participants is a major drawback.

V. CONCLUSION

This work had three main objectives. The first aim was to stratify the dataset of PD patients based on their speech. This resulted in two cohorts differing in the subscores of the Test 3F at T0.

The second aim was to obtain values of the digital biomarkers of speech and voice that would allow the results of the speech tasks of individual patients to be compared with each other. Thus, 16 different features were calculated for each session of each research participant. In addition, these values were then related to the norms obtained from the HC.

The last and the major aim of the whole study was to investigate whether the effect of rTMS on the different stratified groups differs in any way. The findings suggest that rTMS may have varying effects on speech impairments. Depending on the severity level of HD, there was an improvement in a specific group of parameters characterizing a particular speech area (articulation and prosody versus phonation).

The research offers novel observations regarding the advantages of using rTMS to ease the speech impairments in PD patients. Furthermore, it sheds light on the potential correlation between the effectiveness of rTMS and various HD severity levels in PD patients.

The database of research participants offers even more unexplored data in relation to the effects of rTMS on HD according to its severity. Thus, follow-up research is likely. For example, studies observing behavior within each stratified group in search of relationships between changes in digital speech and voice biomarkers and the development of other clinical testing outcomes are suggested.

ACKNOWLEDGMENT

This study was supported by EU – Next Generation EU (project no. LX22NPO5107 (MEYS)). I would like to express my sincere gratitude to my supervisor Jiri Mekyska for his guidance, inspirational suggestions and knowledge sharing throughout this research project.

REFERENCES

- [1] Z. Ou, J. Pan, S. Tang, D. Duan, D. Yu, H. Nong, and Z. Wang, "Global trends in the incidence, prevalence, and years lived with disability of parkinson's disease in 204 countries/territories from 1990 to 2019," *Frontiers in public health*, vol. 9, pp. 776 847–776 847, 2021.
- [2] M. Balaz, J. Buril, P. Burilova, A. Pokorna, and I. Kovacova, "Representation of parkinson's disease and atypical parkinson's syndromes in the czech republic-a nationwide retrospective study," *PloS one*, vol. 16, no. 2, pp. e0246342–e0246342, 2021.
- [3] R. B. Postuma, D. Berg, M. Stern, W. Poewe, C. W. Olanow, W. Oertel, J. Obeso, K. Marek, I. Litvan, A. E. Lang, G. Halliday, C. G. Goetz, T. Gasser, B. Dubois, P. Chan, B. R. Bloem, C. H. Adler, and G. Deuschl, "Mds clinical diagnostic criteria for parkinson's disease," *Movement disorders*, vol. 30, no. 12, pp. 1591–1601, 2015.
- [4] G. Moya-Galé and E. S. Levy, "Parkinson's disease-associated dysarthria: prevalence, impact and management strategies," *Research* and Reviews in Parkinsonism, pp. 9–16, 2019.

- [5] I. Rektorova, M. Mikl, J. Barrett, R. Marecek, I. Rektor, and T. Paus, "Functional neuroanatomy of vocalization in patients with parkinson's disease," *Journal of the neurological sciences*, vol. 313, no. 1, pp. 7–12, 2012.
- [6] D. B. Freed, Motor speech disorders: Diagnosis and treatment, 3rd ed. San Diego, CA, US: Plural Publishing Inc, 2020.
- [7] N. Elfmarková, M. Gajdoš, M. Mračková, J. Mekyska, M. Mikl, and I. Rektorová, "Impact of parkinson's disease and levodopa on resting state functional connectivity related to speech prosody control," *Parkin-sonism & related disorders*, vol. 22, pp. S52–S55, 2016.
- [8] B. Murdoch and C. Barwood, "Non-invasive brain stimulation: A new frontier in the treatment of neurogenic speech-language disorders*," *International journal of speech-language pathology*, vol. 15, 2012/12/17.
- [9] A. W. Shukla, J. J. Shuster, J. W. Chung, D. E. Vaillancourt, C. Patten, J. Ostrem, and M. S. Okun, "Repetitive transcranial magnetic stimulation (rtms) therapy in parkinson disease," *PM & R*, vol. 8, no. 4, pp. 356–366, 2016.
- [10] R. Li, Y. He, W. Qin, Z. Zhang, J. Su, Q. Guan, Y. Chen, and L. Jin, "Effects of repetitive transcranial magnetic stimulation on motor symptoms in parkinson's disease," *Neurorehabilitation and Neural Repair*, vol. 36, no. 7, pp. 395–404, 2022.
- [11] L. Brabenec, J. Mekyska, Z. Galaz, and I. Rektorova, "Speech disorders in parkinson's disease," *Journal of Neural Transmission*, vol. 124, no. 3, pp. 303–334, 2017.
- [12] A. E. Dias, E. R. Barbosa, K. Coracini, F. Maia, M. A. Marcolin, and F. Fregni, "Effects of repetitive transcranial magnetic stimulation on voice and speech in parkinson's disease," *Acta neurologica Scandinavica*, vol. 113, no. 2, pp. 92–99, 2006.
- [13] L. Hartelius, P. Svantesson, A. Hedlund, B. Holmberg, D. Revesz, and T. Thorlin, "Short-term effects of repetitive transcranial magnetic stimulation on speech and voice in individuals with parkinson's disease," *Folia phoniatrica et logopaedica*, vol. 62, no. 3, pp. 104–109, 2010.
- [14] I. Eliasova, J. Mekyska, M. Kostalova, R. Marecek, Z. Smekal, and I. Rektorova, "Acoustic evaluation of short-term effects of repetitive transcranial magnetic stimulation on motor aspects of speech in parkinson's disease," *Journal of Neural Transmission*, vol. 120, no. 4, pp. 597–605, 2013.
- [15] L. Brabenec, P. Klobusiakova, M. Barton, J. Mekyska, Z. Galaz, V. Zvoncak, T. Kiska, J. Mucha, Z. Smekal, M. Kostalova, and I. Rektorova, "Non-invasive stimulation of the auditory feedback area for improved articulation in parkinson's disease," *Parkinsonism & related disorders*, vol. 61, pp. 187–192, 2019.
- [16] L. Brabenec, P. Klobusiakova, P. Simko, M. Kostalova, J. Mekyska, and I. Rektorova, "Non-invasive brain stimulation for speech in parkinson's disease," *Brain stimulation*, vol. 14, no. 3, pp. 571–578, 2021.
- [17] A. Gómez-rodellar, J. Mekyska, P. Gómez-vilda, L. Brabenec, P. Simko, and I. Rektorova, "Evaluation of tms effects on the phonation of parkinson's disease patients," in *Artificial Intelligence in Neuroscience*. Cham: Springer International Publishing, 2022, pp. 199–208.
- [18] L. Brabenec, P. Simko, A. S. Minsterova, M. Kostalova, and I. Rektorova, "Repetitive transcranial magnetic stimulation for hypokinetic dysarthria in parkinson's disease enhances white matter integrity of the auditory-motor loop," *European journal of neurology*, vol. 30, no. 4, pp. 881–886, 2023.
- [19] A. Gómez-Rodellar, J. Mekyska, P. Gómez-Vilda, L. Brabenec, P. Šimko, and I. Rektorová, "A pilot study on the functional stability of phonation in eeg bands after repetitive transcranial magnetic stimulation in parkinson's disease," *International Journal of Neural Systems*, vol. 33, no. 06, p. 2350028, 2023/03/24. [Online]. Available: https://doi.org/10.1142/S0129065723500284
- [20] H. Steurer, E. Schalling, E. Franzén, and F. Albrecht, "Characterization of mild and moderate dysarthria in parkinson's disease: Behavioral measures and neural correlates," *Frontiers in Aging Neuroscience*, vol. 14, 2022. [Online]. Available: https://www.frontiersin.org/articles/10.3389/fnagi.2022.870998
- [21] K. L. Stipancic, F. van Brenk, A. Kain, G. Wilding, and K. Tjaden, "Clear speech variants: An investigation of intelligibility and speaker effort in speakers with parkinson's disease," *American journal of speechlanguage pathology*, vol. 31, no. 6, pp. 2789–2805, 2022.
- [22] A. M. Goberman and L. W. Elmer, "Acoustic analysis of clear versus conversational speech in individuals with parkinson disease," *Journal of communication disorders*, vol. 38, no. 3, pp. 215–230, 2005.

Workplace buildup for measurement of lowfrequency absorption of acoustic structures

Š. Skvaril

Faculty of electrical engineering and communication Brno University of Technology Brno, Czech Republic xskvar08@vut.cz

Abstract— This work deals with realisation of laboratory instrument focused on measurement of sound absorption coefficient of low frequency acoustic resonating systems. Construction based on the mechanism of impedance tube using transfer-function method is chosen with operating frequency chosen between 30 Hz - 200 Hz. Proposed impedance tube was constructed, and its functionality was verified using a set of Helmholtz resonators. Vibration properties of the tube were tested by the means of additional measurement using accelerometers.

Keywords—Sound waves, impedance tube, acoustics, resonator, transfer function method

I. INTRODUCTION

For certain rooms it is crucial that the sound behaves in controlled way. Using acoustic treatment to control reverberation time and low-order reflections are brought to a level, that matches the use of the room, whether it is a concert hall or a recording studio. The most common type of acoustic treatment is utilization of porous material. However, their functionality depends on their thickness. For treatment of lower frequencies, a thicker material is needed. In acoustics the frequencies below 150 Hz are the hardest to treat as the required thickness of porous material is too large, exceeding thicknesses of 0.5 m [1]. Additionally, when the dimensions of the room start are close to wavelengths of specific frequencies, uneven distribution of sound field is caused as a result of standing waves. Example of such uneven sound field distribution for frequency of 100 Hz can be seen in Fig.1 that shows result of FDTD simulation for empty room with dimensions of 4 x 5 x 3 m. This phenomenon is called 'Room modes' and is very common in small listening spaces. [1]



Fig. 1. FDTD simulation of spatial distribution of sound field for 100 Hz in room with dimensions of 5 x 4 x 3 m

Attenuation of those frequencies is achieved using resonant systems, such as Helmholtz resonators or membrane resonators [1] which principle is based on behaviour of mass-spring system. However, design and testing of such devices is not very feasible. The basic property describing the acoustic properties of a material is the sound absorption coefficient, which can take values from 0 to 1, where 0 means absolute sound reflection and 1 means absolute sound absorption [1]. The two most common methods for measuring sound absorption coefficient are the reverberation room method or the impedance tube method [2] [3]. The disadvantage of reverberation room is in the required size of sample area which is several meters squared. It significantly increases costs of testing experimental materials and constructions. On the other side, the impedance tube method uses samples of material that need to be too small. For example, Brüel & Kjær impedance tube type 4206, which is one of the most widely used impedance tube systems, allows to measure only samples in shape of circle with 100 mm diameter.

The aim of research project that is described in the paper is to design, build and verify a laboratory setup, able to use a sample of sufficient size and measure sound absorption coefficient in frequency range from 30 Hz to 200 Hz.

II. DESIGN OF ENCLOSURE

A. Principle of operation

Impedance tube consists of straight tube with loudspeaker attached to one side and highly reflective surface on the other. As illustrated in Fig.2.



Fig. 2 Diagram of impedance tube that uses transfer function method

Between there are two positions of microphones that pick up acoustic pressure in time. Then transfer function between those two positions is calculated using (1)

$$H_{12} = \frac{p_2 p_1^*}{p_1 p_1^*} \tag{1}$$

where p_1 and p_2 are Fourier transforms of acoustic pressures acting on the microphones. The * symbol stands for complex conjugate.[3] H_{12} is then used to compute the reflection coefficient using (2)

$$r = \frac{H_{12} - e^{-jk_0(x_1 - x_2)}}{e^{jk_0(x_1 - x_2)} - H_{12}} e^{2jk_0x_1}$$
(2)

where x_1, x_2 are distances of microphones from the measured sample, k_0 is the wavenumber. The reflection coefficient is then used to compute absorption coefficient using (3) [3]

$$\alpha = 1 - |r|^2 \tag{3}$$

B. Impedance tube design

As the proposed impedance tube should be able to measure larger samples of acoustic elements, cross-section was chosen to be shape of square with length of the side to be 30 cm. This makes possible to measure samples of sufficient size. Equations (4), (5), (6) are specified in [3], to be used for calculation of dimensions of the tube.

$$d < 0.5 \frac{c_0}{f_u} \tag{4}$$

$$f_u \cdot s < 0.45c_0 \tag{5}$$

$$s > 0.05 \frac{c_0}{f_l} \tag{6}$$

where *d* is largest cross-sectional distance, c_0 is speed of sound in the air which corresponds to 343 m/s at 20°C. f_u and f_l are selected upper and lower limit of frequency range and *s* is the distance between microphones. The microphone mounted closer to the loudspeaker should be at least one largest cross-sectional distance away from the loudspeaker. The same applies to the distance of the microphone closer to the sample. [3]

To fulfil the conditions mentioned above with f_u and f_l to be chosen 200 Hz and 30 Hz the distance between microphones should be at least 0.57 m but no more than 0.77 m, so the distance was chosen to be 0.6 m. With cross-section with square shape of 0.3 x 0.3 m is *d* 0.42 m. i.e. the distance of microphones from the loudspeaker and from the sample was chosen 0.5 m both with total length of 1,6 m.

The ISO standard [3] does not assume tube of such dimensions, as it recommends walls to be made of metal with thickness to be 10% of cross-section dimension. That would lead to 30 mm thick metal walls, which would be very difficult and expensive to build and work with. For this reason, the medium density fibreboard (MDF) with thickness of 18 mm was chosen to be the supporting material for 3 mm thick metal sheets which increases mass and rigidity of the structure. The connection of materials was made using epoxy resin. A 20 mm thick metal panel was used as the reflective surface at the end of the tube. The cross section of the wall structure is shown in Fig.3.



Fig. 3 Detail of wall cross-section of proposed impedance tube

All tube walls were joined together using wooden dowels and wood glue as seen in Fig.4



Fig. 4 Impedance tube assembly in progress

Completed tube is shown in Fig.5. It contains four microphone openings in distances of 50 cm, 65 cm, 85 cm and 110 cm from the sample. It makes possible to measure in multiple combinations. When the openings are not in use, they are sealed with fitting cap.



Fig. 5 Completed impedance tube with microphone inserted in position with distance of 50 cm from the sample.

End of tube with loudspeaker mounted is embedded with 35 cm thick absorbent material made of polyester fibre wool with density of 50 kg/m³. It is used to damp the standing wave inside the tube along with reflections from the wall that contains the loudspeaker. It ensures, that the microphones pick up only the progressive wave from the loudspeaker and the wave reflected from the sample.

III. VERIFICATION MEASUREMENTS

A. Measurement using accelerometers

Measurement was performed on the assembled impedance tube with Bruel & Kjaer type 4507-B-004 accelerometers. This measurement was carried out in the frequency range from 30 Hz to 400 Hz in order to determine the inherent mechanical vibrations of the tube walls. The reason for this was to investigate their possible influence on the measurement of the absorption coefficients.

The accelerometers were evenly spaced at a total of nine different positions on the front and 24 positions on the side wall of the tube.

The measurements from the centre of the front wall shows the resonant frequency of 230 Hz, see Fig.6.



Fig. 6 Measurement of displacement from center point of front wall

The measurement from the middle of the side wall of the tube shows that resonances at lower frequencies are excited due to the larger wall dimensions. These resonances start at 75 Hz and then increase in intensity with a relatively constant step of 75 Hz. The measurement result can be seen in Fig.7



Fig. 7 Measurement of displacement from center point of side wall

In order to eliminate the effect of mode, where the accelerometer could be placed in a node and thus not detect any displacement, all measurements from the sidewall were averaged. The result of this averaging shows a main resonant frequency of 230 Hz as in the case of the front-side measurements with additional vibrations above 250 Hz, see Fig.8.



Fig. 8 Average of displacement values from side wall measurements

The interpolated data from the displacement measurements for 230 Hz on the side wall is shown in Fig.9.



Fig. 9 Interpolated measurement for resonant frequency of the tube at 230 Hz

The side of the wall where the displacement reaches the highest values is closer to the loudspeaker, still in the part of tube before the internal damping of the impedance tube.

B. Measurement of sound absorption coefficients

Furthermore, a verification measurement of absorption coefficients was carried out. In the required band from 30 Hz to 200 Hz, the measurement of empty tube results in absorption coefficient less than 0.1. The detail of this measurement is shown in Fig.10. The slight increase at lower frequencies is probably due to sound penetration through the walls of the tube. The reason for the increase in noise-like distortion at higher frequencies is likely due to the resonances of the impedance tube walls as presented in previous chapter. However, because the resonances occur only above the desired operating frequency band, their effect on the measurement results is minimal.



Fig. 10 Detail of sound absorption coefficient of empty impedance tube

Above the working frequency band, the tube wall resonances result in a ripple of measured values around the resonance frequency of 230 Hz, with a shift to negative values, see Fig.11. Measurements in this region are therefore not possible. Large distortion above 260 Hz is caused by microphone spacing which is too large, according to (5).



Fig. 11 Detail of measured sound absorption coefficient above working frequency range of the impedance tube

For the verification of the measurement of absorption coefficient Helmholtz resonators were constructed. They have shape of a cube with internal dimensions of 550 cm³ and 680 cm³ and the possibility of changing the openings. The individual openings were created using 3D printer. Photography of these resonators and available openings is in Fig.12. Their absorption coefficients were measured both individually and in combination when both were placed inside the tube. Those measurements are illustrated in Fig.13 and Fig.14. The results were compared with theoretically calculated resonance frequencies using lumped parameter model and the deviation is less than 10 Hz. However, this difference in values is not due to the impedance tube measurement inaccuracy, but to the difference between the actual Helmholtz resonator and its theoretical model.



Fig. 12 Photography of test Helmholtz resonators with available openings



Fig. 13 Measurement of one combination of cavity and opening



Fig. 14 Measurement of both cavities with different openings both placed inside the impedance tube

IV. CONCLUSION

Impedance tube for measurements of sound absorption coefficients in the range from 30 to 200 Hz was designed for a maximum sample size is 0.3 x 0.3 m.

A working prototype has been successfully built according to this design. Accelerometer measurements were carried out, to observe the range of displacement caused by possible resonance of the impedance tube structure. It was found that the main resonating frequency is 230 Hz, which is above the required working frequency band. Absorption coefficient of the empty tube does not exceed value of 0.1 in the desired frequency range. The functionality of the impedance tube was verified on measurement of sound absorption coefficient of set of Helmholtz resonators.

This impedance tube has great potential both in further investigation and development of acoustic metamaterials and low frequency absorbers and in the possibilities of extending the teaching by incorporating it into laboratory tasks.

REFERENCES

- COX, Trevor J. a D'ANTONIO, Peter. Acoustic absorbers and diffusers: theory, design and application. Third edition. Boca Raton: CRC Press, Taylor & Francis Group, [2017]. ISBN 978-1-4987-4099-9.
- [2] ISO 354:1985, Acoustics Measurement of sound absorption in a reverberation room
- [3] ISO 10534:1996, Acoustics Determination of sound absorption coefficiend in impedance in impedance tubes







A taky Vašek, Ondra, Lukáš, Eva, Olga... Ti všichni táhnou za jeden provaz, kterým přitahují novinky do automotive světa. Za každou naší inovací stojí konkrétní lidé, kterým Valeo vděčí za svůj úspěch. Díky nim mohou být auta autonomní, elektrická, komfortní i bezpečná. Postavte se také za svou myšlenku, kterou za pár let řidiči ocení v provozu.

Koukni na valeo.jobs.cz nebo naskenuj QR kód.

Hardware Design of Mobile Probe for Validating 5G-IoT Technologies

1st Pavel Palurik Department of Telecommunications Brno University of Technology Brno, Czech Republic xpalur01@vut.cz 2nd Radim Dvorak Department of Telecommuniations Brno University of Technology Brno. Czech Republic xdvora2g@vut.cz

Abstract—The evolution of 5G-IoT systems (5G - Internet of Things) has had a significant impact on how modern IoT devices are viewed and has opened numerous new use-cases for such technologies. As the new technologies become widely available, more IoT-technology-driven devices are being adopted in many branches of the industry and more of the bussiness subjects opt for the new IoT-enabled devices over the standard devices. For some businesses, choosing the correct technology may proove challenging as their knowledge in the particular topic may be limited.

In this paper we present the hardware design and initial test cases of a modular Mobile Probe for validating 5G-IoT technologies. The probe is intended to help bussinesses to make informed decisions for the particular use-case based on the autonomous measurements of various 5G-IoT technologies parameters. Due to it's simplicity the probe could also serve as an educational platform for students. After the description of the hardware design, initial testing results for NB-IoT and LTE Cat-M technologies are presented and discussed in terms of the probe capabilities and compatibility with given technologies.

Index Terms—IoT, 5G, hardware design, LTE Cat-M, NB-IoT, ESP32

I. INTRODUCTION

As the market with IoT technologies grows, many new IoTenabled devices and use-cases appear on the market constantly. This has steered the industry for faster adoption of the IoT revolution by incorporating such devices into their portfolio/inventory. While some businesses choose to adopt these devices on their own some are required to do so by legislation especially in the distribution industry.

In Czech Republic, a legislation that requires all new energy meters, installed from 2024, to be remotely accessible, was passed. This has left the distributors in a difficult position as some of the use-cases could not be covered by standard broadband technologies due to variety of the environment and coverage issues (metal electric distribution boxes, cellars, remote areas etc). With little to no expertise they often turned to third party evaluation to perform measurements and seek advice, effectively costing a lot of money and resources in the process for testing, as the distributor technicians have to be present on-site during the measurements.

To tackle this issue we created a modular, battery-powered testing probe for 5G-IoT technologies that helps the businesses evaluate the communication parameters for the given use-case

and make an educated decision based on the measurement results. Furthermore the device saves the businesses significant amount of costs and human resources as the device can be left unattended by the technician while performing the measurement.

This paper focuses on description and validation of the hardware design of the test probe as well as validating the measurement capability by measuring the border parameters and throughput for two CIoT (Cellular IoT) technologies NB-IoT and LTE Cat-M.

The key outputs of our work can be summarized as follows:

- The testing probe is capable of succesful validation of a given CIoT technologies and is suitable for bussiness use as well as a education platform.
- The modularity of the platform ensures long usability for technologies to come (for example RedCap Reduced Capability).

The rest of the paper is organised as follows. In Section II we give an overview of the design process of the probe. In Section III brief description of tested technologies NB-IoT and LTE Cat-M is given as well as result of the hardware testing. Finally, in Section IV, we conclude the paper.

II. DESIGN PROCESS

At the start of the design several must-have requirements have been set based on our prior experience with such devices and intended technologies for testing.

Sufficient computing power - For this task the ESP32-S3 module from Espressif system was selected. It features a dualcore XTensa LX7 microcontroller with maximum clock rate of 240 MHz supplemented with high-speed octal SPI flash memory. Additionally, it includes integrated 2.4 GHz, 802.11 b/g/n Wi-Fi, and Bluetooth 5 (LE) for local connectivity. This module was chosen based on a sufficient number of GPIO pins, price, availability and connectivity options.

Battery power - To ensure the testing capability in enclosed environments the battery power was a necessity for such a tester. This task has been met by implementing a single Li+ battery with USB-C connector and appropriate charging and protection circuits. **Local storage** - To store the results of a performed measurement or user configuration the device was equipped with a SD card interface.

Modularity - To allow the device to be used with various types of technologies, modular interface had to be included. The available options were Mini PCI Express (mPCIe) or M.2, formerly Next Generation Form Factor (NGFF). The first one mentioned has larger pin spacing including the dimensions of entire connector and a fixed structure of pulled out signals. In contrast, M.2 is a specification that is divided into several subgroups based on the position of the cutout, each with a different choice of buses that are pulled out and used on the connector. The connector surfaces have smaller dimensions, allowing more of the pins to fit on a connector. For this task the M.2 interface was selected.

User accesible SIM card connector - To allow users to swap the SIM card without dissassembling the device's enclosure a SIM card slot holder had to be mounted on the main board. This feature allows user to test the CIoT technologies with various connectivity providers.

User Interface - The user interface is a crucial part of the whole device design. As the device is meant for use without prior technology knowledge it has to be intuitive and user friendly.

A. Motherboard PCB

This is the device's main printed circuit board. It contains an ESP32 microprocessor module, an external connector for a communication module adapter, a charging connector, integrated circuits that deliver suitable voltage levels with sufficient current capacity to the device, connectors for local, removable storage, a SIM card slot, and user interface connectors, including those for probe programming and debugging, as shown in Fig. 1.

1) Power Section: The device features a USB-C (Universal Serial Bus - C) connector supporting USB 2.0 data communication and power supply, including battery charging. The connector supplies USB standard signals such as VBUS, DP, DN, CC1, CC2, and GND to the board. The input supply voltage is applied to the autonomous charger MAX77751 IC (Integrated circuit) intended for 1-cell Li+ batteries. This IC has integrated USB Type-C CC detection and reverse boost compatibility. It also offers fast-charge current selection and top-off current threshold setting using externally connected resistors. The IC implements the Adaptive Input Current Limit (AICL) function to limit the maximum input current and regulate the input voltage to prevent the supply voltage from collapsing due to a possible soft source. The device features a Smart Power Selector and a battery true-disconnect FET (Field Effect Transistor) to regulate charging and discharging or isolate it in case of a failure [1].

To ensure data communication with the ESP32, a USB DPDT switch is included in the design. The ESP32 controls and switches the switch after the initial setting of the maximum possible charging currents that the adapter can deliver. The LiPol battery selected for the probe serves as the power source



Fig. 1. Block diagram describing device composition.

in regular mobile operation. The step-down converter is used to reduce the voltage to the required supply level of 3.3 V. Two converters from Renesas, the ISL9120IR and ISL91107IR, are connected in parallel, and their connection to the circuit can be selectively changed using resistors with zero resistivity. Both chips require few external components and have a wide input voltage range of 1.8V to 5.5V and a conversion efficiency of at least 85%. The weaker IC can deliver a maximum output current of 0.8A, while the stronger IC can deliver up to 2A [2], [3]. This parallel connection serves for future current draw testing.

The resulting 3.3V voltage level powers all peripherals, including the Low Drop-Out (LDO) 1.8V voltage regulator MCP1802T, which can supply up to 300mA. The power rails are equipped with sufficient filter capacitors, including a TVS diode PESD3V3L1BA,115, and LED signaling of the present voltage [4].

2) Expansion boards: The M.2 connector is utilized for expansion boards of various communication module adapters. The connection of this connector partially adheres to the reference design presented in the congatec document. The E-key version is used, meaning pins numbered 24 to 31 are omitted. This text describes multiple deviations from the congatec reference specification. As a result, the connection falls under the category of proprietary applications of the M.2 E-key connector [5].

The connector incorporates multiple interfaces, including I2C, SPI, UART, USIM, and USB, along with various generic signals. The USIM interface signals are safeguarded by TVS diodes against overvoltage conditions and connected to the push-push type housing for Micro SIM cards. In the initial version of the probe, the USB data are used for control and FW local update of used comm. module is connected to the secondary debug USB-C connector, which may not be included in future versions. The Texas Instrument's TXB0102DCUR chip is adopted as a level-shifting IC, fully satisfying the

requirements for possible NB-IoT and Cat-M modules. Future prototypes may employ the TSX0108E IC, which can provide up to eight possible level-shifting signals for I2C, SPI, and UART in addition to the existing architecture.

3) Local Storage: The micro SD, figuring as local storage, card was connected using a Molex 47352-1001 connector with a PUSH-PUSH type mechanism, utilizing the SPI bus. All signals are equipped with 10k Ohm pull-up resistors and are safeguarded by a TVS diode array, specifically the SP0503BAHTG model, which incorporates three protective diodes in its body.

4) Placement of components on the PCB: A preview of the component layout on the probe's motherboard is shown in Figure 2. Here you can see two USB-C connectors, a USB switch, a circuit ensuring charging with voltage transmission from the battery to two step-down converters connected in parallel and a 1.8V voltage source. In addition, we can find a slot for a Micro SIM card with the necessary ESD protection, a Micro SD card, the necessary header connectors, an M.2 slot for expansion adapter boards and an ESP32 process module.



Fig. 2. Motherboard PCB in Altium with 3D render.

B. User Interface PCB

The second board is used to complete the device. It contains a mechanically fixed color TFT (Thin Film Transistor) display, a group of control buttons, and a connector for connection to the lower baseboard. Additionally, there is space on the bottom of this board for soldering two other unconnected slots, which can be used for local storage of additional micro SIM cards for possible connection to the networks of other potential operators. The probe device has a total of three slots for SIM cards, allowing for testing of major domestic operators such as Vodafone, T-Mobile, and O2 in the Czech Republic. Each button connected to the connector in its closed state acts as a 100 k Ω to 10 k Ω voltage divider connection, i.e., a 10:1 conversion ratio. The signal itself is equipped with a TVS diode to protect against possible accidental transients that could, for example, cause faulty actuation.

III. MEASUREMENT CAPABILITY TESTING

To evaluate the devices testing capability for various technologies we have opted for the M2 key add-on board with LPWA (Low Power Wide Area) module Quectel BG77.

BG77 is a dual RAT (Radio Access Technology) modem enabling connectivity via two CIoT technologies LTE Cat-NB (Narrowband IoT) and LTE Cat-M. NB-IoT and LTE Cat-M are a part of the 5G-ecosystem technologies for IoT applications and their utilization has been on the rise in the Czech Republic due to the growing demand for IoT connectivity and legislative requirements especially in the smart metering field.

A. NB-IoT and LTE Cat-M

Both technologies, described as a part of 3GPP Release 13, are derived from the existing LTE standard which ensures the interoperability with the existing public mobile network infrastructure and easy deployment for the mobile network operators. To better suit the IoT applications and to enable massive deployment, on the scale of millions, while keeping the cost-per-modem low, both technologies significantly reduce the hardware complexity (up to 80%), feature capability (such as MIMO) and also system bandwidth to achieve better coverage, enabling the connectivity for broader spectrum of devices. The LTE Cat-M1 technology reduces the system bandwidth to 1.4 MHz¹ and 5 MHz² respectively, while the NB-IoT reduces it even further down to 200 kHz (180+20 kHz). This results in a significant reduction of throughput for both technologies as the maximum theoretical throughput for LTE Cat-M1 is up to 1 Mbps and up to 3 Mbps for the Release 14 LTE Cat-M2, while the Release 13 NB-IoT achieves up to 62.5 kbps and up to 159 kbps in Release 14 respectively [6].

As both technologies are focused on IoT battery powered devices with goal of achieving 10+ years of battery life, they both implement standard power saving features such as PSM (Power Saving Mode) and eDRX (extended Discontinuous Reception) which, in addition with the reduced hardware complexity and reduced capability, help to reduce the power consumption of the technology-enabled devices [6].

To achieve better coverage both technologies implement Coverage Enhancement (CE) which modifies the radio parameters such as used modulation, number of repetitions or TX power to achieve more robust transmission in harsh radio conditions. The Coverage Enhancement level is assigned to the device by the operator's network and is based on well defined thresholds, based on the measured values of RSRP (Reference Signal Received Power) and SINR (Signal to Interference Ratio). The set threshold can vary from operator to operator as the operator is able to set these thresholds as required to ensure the best network performance. The NB-IoT features three Coverage Enhancement levels also known as ECL (Enhanced Coverage Level) classes 0,1 and 2 while the LTE Cat-M features Coverage Enhancement Modes A and B [6].

¹3GPP Release 13

²3GPP Release 14

 TABLE I

 Comparison of the CIOT technologies defined in 3GPP Releases

 13 and 14

Technology	NB-IoT (NB1)	NB-IoT (NB2)	LTE Cat M1	LTE Cat M2	
3GPP	Release 13	Release 14	Release 13	Release 14	
Frequency		700-21	100 MHz	•	
Bandwidth	200	kHz	1.4 MHz	5 MHz	
Link budget	164	dB	155	.7 dB	
Max. EIRP		23 dBm			
Max. payload	1600 B		8188 B		
III data mata	0.3-62.5 Kbps 0.3-159 Kbps	0.3.150 Khos	HD: 375 Kbps;	HD: 2.625 Kbps;	
OL uata Tate		FD: 1 Mbps	FD: 7 Mbps		
DI data rata	0.5.27.2 Khns	0.5.127 Khns	HD: 300 Kbps;	HD: 2.35 Mbps;	
DL uata Tate	0.5-27.2 Kops	0.5-127 Kops	FD: 0.8 Mbps	FD: 4 Mbps	
	Tx: 240 mA	Tx: 240 mA	Tx: 360 mA	Tx: 360 mA	
Consumption	Rx: 12 mA	Rx: 46 mA	Rx: 46 mA	Rx: 70 mA	
	PSM: <1uA	PSM: <3uA	PSM: <3uA	PSM: <8uA	
Security	LTE security				

B. Testing Scenarios

The tester was placed in a RF shield box R&S CMW-Z10 and connected through a step attenuator directly to the local commercial RRU (Remote Radio Head) via coaxial cable (Figure 3). This setup enabled us to directly control the RSRP (Reference Symbol Received Power) and SINR (Signal-to-Interference-plus-Noise Ratio) values read by the module. The local RRU is a production RRU of the Vodafone Czech Republic operator therefore any results measured on our testbed should directly match the results and parameters of a public network.



Fig. 3. Measurement testbed.

For the evaluation, several tests were performed for both technologies available (NB-IoT and LTE Cat-M) to test the probe's capability of testing the technologies to their full extent. As the design is targeted for the technology evaluation, validation and user testing, the test scenarios were selected to fit the scheme. The selected tests include:

- Border communication parameters measurement.
- Throughput measurement.

C. Border communication parameters measurement

This test was designed to evaluate the HW testing capabilities in terms of the border radio parameters for a given technology.

1) Measurement Scenario: For this measurement we utilized the variable attenuator to lower the value of RSRP read by the device. In each iteration, we attempted a registration to the network and evaluated whether the registration was succesful. For each iteration we also performed a basic data exchange utilizing the *ping* feature and monitored the CE level for each step. During the tests, all modem< - >MCU communication was logged and saved onto the SD card and later evaluated.

2) Measurement Results: From the measurement the NB-IoT technology was able to perform a succesful registration to the network at the threshold of -129 dBm RSRP. After the -129 dBm threshold, we were not able to succesfuly register to the network. However, after already being registered to the network the device was able to successfuly exchange data for down to -132 dBm RSRP. For the CE Level, we have reached the ECL 0 at -111 dBm RSRP and ECL 2 at -127 dBm RSRP due to the SINR exceeding the threshold value of -3 dB.

For LTE Cat-M, here the technology was able to successfuly register to the network for down to **-127 dBm RSRP** with the last data exchange being measured at **-129 dBm RSRP**.



Fig. 4. NB-IoT and LTE Cat-M1 measured sensitivity.

D. Throughput measurement

The goal of this measurement was to evaluate the design of the communication lines, in other words, whether the device is capable of saturating the communication interface and fully exploit the limits of the given technology. This measurement also allowed us to verify the technology's limits as seen in table I.

1) Measurement Scenario: For testing throughput, PPPoS (Point to Point over Serial) connection was estabilished to the BG77 module. We utilized the *iperf* application, in UDP mode and uplink direction, which is available for the ESP32 via the its development framework. We started our measurement at -60 dBm RSRP and using the variable attenuator we lowered the RSRP value in granular steps based on the current RSRP measured by the module. In each iteration we ran the *iperf* and saved the results on a SD card.

2) Throughput measurement results: After the measurements we removed the SD card from the device and evaluated the results saved on the SD card.

For the LTE Cat-M technology the throughput of maximum of **1003 kbps** was achieved. This is the maximum throughput the LTE Cat-M technology can achieve as seen in table I. From this result we can conclude that the communication interface is designed correctly as it is able to saturate the technology communication channel. The LTE Cat-M maintains stable maximum throughput up to -104 dBm RSRP where the sudden decline can be seen in figure 5. The decline appears to be linear with the average decline of 52 kbps per dB of RSRP. In comparison to the Border limits measurement we were not



Fig. 5. LTE Cat-M1 measured throughput.

able to achieve any data throughput beyond the -122 dBm of RSRP which is probably due to the *iperf* application not being able to estabilish a connection due to harsh conditions and heavy packet loss.



Fig. 6. NB-IoT measured throughput.

For the NB-IoT case, the throughput was stable for up to -107 dBm RSRP until the data throughput decline can be seen as in comparison with the LTE Cat-M. The maximum achieved throughput was measured as 28 kbps which does not match the comparison in table I, however we believe the value to be correct for the given setup and to be the result of network limitations as the tester was able to saturate the LTE Cat-M technology at 1 Mbps without further issues. These results also correlate with our previous measurements of the NB-IoT technology.

For the NB-IoT data rate decline, two "step-like" deviations from the linear decline pattern can be seen in Fig. 6 as the throughput declines. These deviations are due to changes of the ECL class from 0 to 1 and from 1 to 2 respectively, which significantly improved the channel realiability for the given RSRP value.

Here, as well as with the LTE Cat-M we were not able to get any throughput beyond -128 dBm due to *iperf* issues.

IV. CONCLUSION

With the growing demand for IoT technologies in Czech Republic in mind, especially in the smart metering field, the target was set to create a battery-powered probe that can be used directly by the distributor or other subjects, without third-party contractor involvement, to evaluate the technology for the specific use-case, effectively saving costs and human resources required for such validation.

In this paper, the Hardware design of the Mobile Probe for validating 5G-IoT Technologies was presented as well as tests performed to validate the correct functionality of the device.

The probe was designed as a modular device platform, enabling the user to test and validate different technologies or module manufacturers which is achieved via the M.2 Ekey interface where the external add-ons can be mounted.

For the hardware evaluation the BG77 M.2 E-Key add-on board was selected. During the evaluation, various aspects of the hardware design were tested (such as USB and UART speed, SIM card slot, sufficient SD card interface speed for logging) as well as the basic comparison for two uprising CIoT technologies NB-IoT and LTE Cat-M was presented.

While the LTE Cat-M technology is capable of reaching speeds up to 1003 kbps the NB-IoT is able to provide more coverage especially in environments with bad radio conditions (RSRP < -125 dBm).

From the initial testing, we conclude that the probe design is sufficient for testing 5G-IoT technologies as the capabilities and limits of both technologies could be achieved with the design. This leads to the conclusion that the **device's capabilities are suitable for use by companies as a verification device** or as an education tool for students.

A. Future Work

In the future work we will focus on developing the software for validating different aspects of both technologies and higher layer protocols to further demonstrate the technology limits and features to the user of the probe.

The future development includes:

- Software for autonomous testing of various technology aspects and features (throughput, PSM, eDRX, Packet loss and more).
- Sufficient User Interface for the tester.
- Layer 4 and Layer 7 protocols testing features. (TCP, UDP, MQTT, CoAP and more).

REFERENCES

- MAX77751, 3.15A USB-C Autonomous Charger for 1-Cell Li+ Batteries, Maxim Integrated, November 2020, Rev. 4.
- [2] ISL9120IR, Compact High Efficiency Low Power Buck-Boost Regulator, Renesas, February 2016, Rev. 1.00.
- [3] ISL91107IR, High Efficiency Buck-Boost Regulator with 4.1A Switches, Renesas, March 2015, Rev. 0.00.
- [4] MCP1802, 300 mA, High PSRR, Low Quiescent Current LDO, Microchip, October 2010, Rev. C.
- [5] *M.2TM Pinout Descriptions and Reference Designs*, congatec, January 2020, First release; Rev. 1.0.
- [6] O. Liberg, M. Sundberg, E. Wang, J. Bergman, J. Sachs, and G. Wikström, Cellular internet of things: from massive deployments to critical 5G applications. Academic Press, 2019.

Strategic Capacity Planning and Optimization in Communication Networks: A Case Study

1st Aneta Kolackova Department of Telecommunications FEEC, Brno University of Technology Brno, Czech Republic aneta.kolackova@vut.cz

Abstract—This paper explores the complex task of capacity planning for IP-based communication sites, emphasizing the importance of strategic resource management and infrastructure optimization. It highlights the benefits of the methodology developed recently by the Czech Telecommunications Office (CTU) for both large and small providers. The study examines the implementation of an upspeed enhancement in communication sites operated by a major provider. Using CTU's methodology, cost-effective adjustments were proposed, revealing an immediate bandwidth requirement of 1500 Mbps, closely aligning with the pre-upspeed value of 1000 Mbps (in one of considered scenarios). The analysis also demonstrated the feasibility of adding an additional Network Termination Point (NTP) within the existing 3000 Mbps limitation. However, to accommodate a fifth NTP, a bandwidth upgrade to 3500 Mbps would be necessary.

Index Terms—Capacity Planning, Network Optimization, Upspeed, Bandwidth, Utilization of network

I. INTRODUCTION

Network capacity is generally planned to meet the requirements of target networks in the foreseeable future based on the prediction of future user and service development, as well as to resolve congestion problems that occur in the real IP network. Therefore, current and future service requirements must be considered during capacity planning. After a capacity plan is implemented, user and service potential of the live network can be unleashed, and network capacity can be expanded in phases to meet the requirements of the future growth of the user base and services.

Following the implementation of a capacity plan, rigorous testing becomes paramount. Validating the network's performance against projected requirements ensures that the planned enhancements effectively address congestion issues and accommodate future growth. Robust testing allows us to fine-tune the network, optimize resource utilization, and maintain seamless service delivery.

When testing both the mobile and fixed electronic communications networks, it is advisable to based on a coherent methodology. In the Czech Republic, the Czech Telecommunications Office (CTU) carries out measurements of data parameters of mobile networks and fixed networks within its competences. In March 2021, the CTU issued a new methodology for measuring mobile networks [1]. The methodology defines the guidelines for measuring mobile networks, which

2nd Jan Jerabek Department of Telecommunications FEEC, Brno University of Technology Brno, Czech Republic jerabekj@vut.cz

the CTU follows in control measurements of mobile networks, which may be stationary or under movement (typically drivetests). The measurement methods in the CTU methodology are based on documents issued by the Association of European Regulators for Electronic Communications BEREC (The Body of European Regulators for Electronic Communications) [1].

This article deals with the new methodology [2], which was published in November 2023, and which follows the abovementioned mobile methodology and [3].

Although the research measurements originate from an authentic mobile network, it is essential to clarify that they do not constitute end-user data. The focus of investigation lies specifically in the network backhaul – the link connecting cell sites to the core network. Given this context, the recently developed methodology proves particularly suited for validating realistic network configurations and operational scenarios.

Furthermore, the CTU has introduced an auxiliary tool known as 'netcalculator' [4]. This tool facilitates the assessment of electronic communications network capacity impact. In the context of this article, netcalculator will play a pivotal role in the verification process.

II. BACKGROUND AND RELATED WORKS

The evaluation of the impact of the capacity of electronic communications networks, specifically in the distribution or connection segment of the electronic communications network, on the performance of Internet access services is based on the theory of bulk service, with the assumption that the elementary input data flow of bulk service requests is stationary, regular and incrementally independent.

Stationarity is the number of data streams - network termination points (NTPs) that arrive at the bulk service system in time Δt depend only on the length of this interval and does not depend on its position on the timeline. Regularity is the probability of more than one data flow occurring in a sufficiently small interval of length Δt is negligibly small. Increment independence is the number of data streams that occur in one time interval, independent of the number of data streams in other intervals [2].

The capacity of the network, whether it is a backhaul network or a distribution network, has a significant impact on the performance of the Internet access services provided. Insufficient network capacity will cause an increase (deterioration) in the QoS parameters that fall within the set of extended data parameters according to [3].

The performance of Internet access services is characterised by Actual Achieved Speed (SDR^1) , which can be imagined as the transmission rate corresponding to the transport layer of the ISO/OSI model (L4). The resulting SDR at a particular NTP can be determined through a measurement process specified by [3], however, for a global assessment of the impact of network capacity as a function of the number of NTPs on the resulting decrease in SDR, it is necessary to use this new methodology, which is designed additionally to be used for activation, planning or similar activities.

A. Network termination points (NTPs)

For NTP, the network capacity is most commonly specified by the NBR bit rate value corresponding to the physical layer (L1), possibly the IR rate value corresponding to the link layer (L2), or the IP TR throughput value, i.e., the rate value corresponding to the network layer (L3).

The actual network capacity value can be verified by measuring the bandwidth in the downlink or uplink direction using the method specified in [6]. According to [6], it is also possible to assess the impact of the network capacity with respect to a monthly time frame-based on the result of network traffic monitoring. This will be exploited later in the paper.

B. Impact of network Utlization Factor (UF) by NTPs

Commonly used methods to identify the utilization rate of the network, include the determination of the link utilization (LU) parameter of the network topology [4]. The LU parameter is defined as the ratio of the nominal speed value over time (the current NBR bit rate value) and the corresponding capacity. As a rule of thumb, if the average value of the LU parameter reaches 0.4 to 0.5, it is recommended to increase the capacity of the data link.

The nature of typical NTP usage of an electronic communications network in a given location and traffic type can also be described by the aggregate utilization factor (UF) [4], which is determined as the ratio of the average speed value (average bit rate) over a certain time frame and the nominal bit rate value at peak operating times. The UF does not depend directly on the capacity value of a given data link (connection network), or on the actual value of the LU usage parameter at a given time, but on the behaviour of NTP depending on the location and nature of the data traffic.

C. Performance of Actual Achieved Speed (SDR)

In assessing the impact of network capacity on the resulting performance of Internet access services, represented by the SDR at the NTP location, it is necessary to rely on the knowledge of network capacity itself, just as in assessing the impact on the average number of NTPs. It is also necessary to determine the percentage value of the probability at which a certain number of traffic flows generated simultaneously by a specified number of NTPs will arrive at the system. It is generally recommended that the probability value should not be less than 90% [2].

Researchers deal with a similar topic in many articles with different approaches. Thanks to the knowledge of historical data, they often work with time series forecasting, such as here [7] and [8]. However, current data are not always available. So, a mathematical approach that corresponds more to our used solution is found [6]. Although this paper focuses on wireless networks, it proposes a mathematical framework for the planning capacity of a 5G and beyond wireless networks.

Nevertheless, the surveyed literature does not present an analogous opportunity for expedient and straightforward validation, thereby we decided to use the method proposed by CTU.

III. AUXILIARY TOOL NETCALCULATOR

CTU created an auxiliary tool to provide a software capable to calculate and explain the issues of evaluating the impact of the capacity of IP communication networks in the context of this methodology. The application is written in python and works on the principle of recalculating values based on input data. Netcalutator has a graphical interface, which is shown in Fig. 1 (unfortunately, only Czech version is available).

The calculator has 4 parts. In part A, the "Input parameters of the electronic communications network" are entered. These are the mandatory parameters of the network under consideration, which correspond to the identified bottleneck of the network infrastructure. The important value is A.1, which identifies the capacity of site bottleneck. The MTU in field A.2 is left at the maximum value, 1514. A.3 specifies whether IPv6 or IPv4 is used. And entry A.4 indicates the number NTPs affected by the assessed value of the bottleneck capacity.

Part B should be filled with the results of the monitoring of the screen traffic (ideally for at least a month) at the site under assessment. B.1 indicates the Maximum Bit Rate (NBR_{MAX}) at the bottleneck location. To this value, the value of the usage parameter LU_{MAX} can be determined, which corresponds to the ratio of the maximum bit rate (field B.1) and the bottleneck capacity (field A.1). Field B.2 gives the average bit rate value (NBR_{AVG}) from network traffic monitoring (over the last month) at the bottleneck location. In B.3, the value of the utilization factor is calculated to reflect the end-user real behavior in the considered network and with respect to bottleneck average utilization.

In part C, the input parameter of the Poisson process, or probability, is determined, which is closely related to the number of NTPs generating simultaneous data flows to the bulk system (the bottleneck under consideration) at time t. A probability of 90% is recommended for rural areas and 95% for urban areas.

In part D, the calculator addresses the impact of aggregation on the bottleneck of communications network. With this section, we are able to analyze the impact of the desired SDR at the active connection point, the NTP endpoint, on the characteristics of the identified bottleneck. Field D.1 is

¹Because of the netcalculator usage, that has Czech shortcuts.
Český tele	komunikační	úřad			
Metodika pro na výkon služ	vyhodnocování d šeb přístupu k inte	opadu kapacity sítí e rmetu, verze 1.0	lektronických ko	omunikací	Nápověda
. Vstupní parametry sítě elektronic	kých komunikaci				A
A.1 Kapacita sítě elektronických kon	nunikací (L1)	3 000,0	Mb/s		1
A.2 MTU (Maximum Transport Unit)		1,500 🔹 B			2
A.3 Velikost IP záhlaví		40 (IPv6) 🗸 B			3
A.4 Agregační poměr (počet posuzo	vaných NTP)	1: 3 🔹			4
Výsledky monitoringu síťového pr	ovozu (měsíc)				В
B.1 Maximální bitová rychlost NBR _m	_{ax} (L1)	1,000.0	Mb/s	LU _{max} 0,33	1
B.2 Průměrná bitová rychlost NBR _{avc}	, (L1)	500.0	Mb/s	LU _{avg} 0,17	2
B.3 Faktor využití UF		0,5			3
Vstupní parametr Poissonova proc	esu				С
C.1 Pravděpodobnost		95 (městské obla	esti) 🗸 %		
Dopad agregace na přípojku					D
D.1 Výsledná SDR (L4)	936,3	Mb/s (bez UF)	1 404,4	Mb/s (s UF)	
Dopad agregace na úzké hrdlo sítě	é elektronických	komunikací			E
E.1 Požadovaná SDR (L4) přípojky	100.0	Mb/s			1
E.2 Průměrný počet přípojek (NTP)	582	(bez UF)	1 748	(s UF)	2
E.3 Pokles výkonu služby na přípojce	26,1	% (bez UF)	16,4	% (s UF)	3
E.4 Potřebná šířka pásma (L3) úzkéh	o hrdla	301,7	Mb/s		4

Fig. 1. Example of particular values entered to netcalculator for calculating the impact of IP communication network capacity (available only in Czech)

auxiliary and indicates the input value of the average SDR to be considered with respect to the properties of the identified network infrastructure bottleneck in the Poisson process for the purpose of aggregation impact analysis. With and without taking UF into account.

In Section E, the calculator addresses the impact of aggregation on the bottleneck of the communication network on NTP. With this section, we are able to analyze the impact of the desired SDR rate at the active connection point, i.e., the NTP endpoint, on the characteristics of the identified bottleneck.

Item E.1 is optional and indicates the input value of the average SDR to be considered with respect to the properties of the identified network infrastructure bottleneck in the Poisson process for the purpose of aggregation impact analysis. Then, in E.2, the it calculates the average number of NTP without and with UF. Without UF, this will take into account the required SDR in E.1 and the identified network capacity at bottleneck entered to A.1 with the current assumption of full utilization. Whilst the resulting average NTP with UF will also take into account the typical NTP behaviour as per network traffic monitoring B.3 at the site.

Field E.3 describes the performance drop in the situation

where NTP users meet at the same time, again from the perspective with and without UF.

The last fields E.4 and E.5 are similar, but E.5 additionally takes into account the UF, calculated by the calculator in field B.3. E.4 describes the resulting average bandwidth value at the identified bottleneck location corresponding to the number of NTPs, i.e. field A.4, at the desired SDR rate (E.1) at the NTP network endpoint location. E.5 is the same, but in addition to E.1 it also takes into account the typical NTP behaviour according to B.3.

IV. DATA EXAMINED AND MONITORING OF NETWORK TRAFFIC

In this chapter we will focus on the description of three real 4G+5G cell sites where the transmission speed limit in the bottleneck has been recently increased. Specifically, from a transfer rate of 900 Mbits to 3000 Mbits, In all three cases these are macro cell sites, situated in a urban area in cell site shelters. The precise locations of the base stations have been intentionally omitted to preserve the commercial benefits of the operator. Each site has 3 baseband units (BBUs), which we can consider as NTP under certain conditions, see in Fig. 2.



Fig. 2. A diagram of one cell site with three BBUs connected by optical fiber to a router that is connected towards the Internet.

Over the span of 58 days, each BBU from all 3 sites has been monitored and records consisting of hourly observations of mobile data usage have been collected. The monitoring took place 29 days before the change and 29 days after the bottleneck parameters change and router reconfiguration to higher throughput. When considering the data, we always consider it from the perspective of a single site containing 3 BBUs.

All calculations work with traffic in the download direction as an example. In these scenarios, download is the most critical. For the purposes of the calculator and according to the methodology, data from monitoring were averaged to get the NBR_{AVG} and their NBR_{MAX} was found. The NBR_{MAX} is taken from the average value per hour, because monitoring has no less granularity. The individual values for a given cell site are given in the Table I.

The daily comparison of individual sites reveals discernible variations in load patterns, as depicted in Fig. 3, 4 and 5. Despite the observed increase in upspeed, the NBR behavior remains consistent in the graphical representations, even though the capacity has expanded threefold.

The graphical representations provide evidence that the increase in upspeed has not yielded any discernible alterations.

B.2 NBR_{AVG} [Mbps] **B.1** NBR $_{MAX}$ [Mbps] Before After Before After Site A 300.0 419.6 157.6 159.7 422.0 104.2 Site B 316.2 104.1472.1 262.9 Site C 461.3 261.3

 TABLE I

 MONTHLY MONITORING OF THE SITES - BEFORE AND AFTER THE UPSPEED

Instead, it primarily pertains to strategic capacity planning for future requirements, rather than an immediate necessity. In the subsequent chapter, we shall employ the netcalculator to ascertain the maximum capacity and assess the rationality of

the proposed planning.





Fig. 3. Site A - Daily comparison before and after the upspeed

Fig. 4. Site B - Daily comparison before and after the upspeed

V. PLANNING OPTIMIZATION USING NETCALCULATOR

In the preceding section, we expounded upon the scrutinized data and meticulously depicted the NBR both prior to and subsequent to the upspeed. The operator is presently remunerating for a line that remains underutilized, as evidenced by



Fig. 5. Site C - Daily comparison before and after the upspeed

the graphical representation. Leveraging the netcalculator, we can computationally ascertain the theoretical upper bounds of speeds and meticulously strategize and engineer the uplink to optimally exploit the site's resources, thereby resulting in cost savings for the operator.

Table II, describes the calculated UF and the resulting SDR. Only the values of NBR_{*MAX*} and NBR_{*AVG*} are relevant for the UF calculation, so the upspeed, i.e. the increase capacity of site bottleneck (A.1), does not affect this calculation. Only Site A has improved, in the other cases UF has deteriorated. However, it should be noted that NBR_{*MAX*} is calculated from the averaged hourly peak, not the total peak on the link.

In the case of SDR, the value of capacity of site bottleneck (A.1) has an effect and after upspeed the values in resulted SDR (D.1) increased, see Table II. This calculation confirms that we are able to increase the capacity of each NTP and thus design an ideal ratio, with the help of part E in netcalculator.

TABLE II

CALCULATING UF AND SDR USING THE NETCALCULATOR

	B.3 Util	ization Factor (UF)	D.1 Resu (without [M	llted SDR /with UF) bps]
	Before	After	Before	After
Site A	0.53	0.38	281.1/421.6	936.9/2380.3
Site B	0.25	0.33	281.1/421.6	936.9/2810.6
Site C	0.56	0.57	281.1/421.6	936.9/2165.1

We first considered the required SDR of the connection (NTP) at 1000 Mbps on L1, so 936.3 Mbps assuming the use of IPv6 and the transport protocol header that is present at L4. This is the speed at which the NTP connections are now set. The calculated values can be seen in Table III. At the same time, a hypothetical value has been calculated assuming that the required SDR is halved, given the traffic that flows through the network in all 3 sites, see Chapter IV.

In column E.1 - desired SDR, the value of 936.3 Mbps, explained above, was chosen, and also the half value of 468.2 Mbps was chosen for the hypothetical site connection

TABLE III CALCULATING PART E, IMPACT OF AGGREGATION ON THE BOTTLENECK REGARDING TO NTP

	E.1 [Mbps]	E.2 w/o UF [-]	E.2 w/UF [-]	E.3 w/o UF [%]	E.3 w/UF [%]	E.4 [Mbps]	E.5 [Mbps]
Site A	936.3	4	29	54.5	30.6	2824.2	1091.6
H_Site A	468.2	19	141	45.3	23.5	1412.3	771.9
Site B	936.3	4	38	54.5	27.2	2824.2	947.6
H_Site B	468.2	19	187	45.3	20.7	1412.3	670.1
Site C	936.3	4	26	54.5	32.1	2824.2	1144.5
H_Site C	468.2	19	128	45.3	22.7	1412.3	809.3

H_Site A, B or C.

For column E.2 - number of NTP - without (w/o) UF and with (w/) UF, values of 4 and 19 were calculated for all three sites. This means that even in the case of using a 1000 Mbps link, it would be possible to connect 1 more NTP, i.e. 4, compared to the current 3. Taking UF into account, we even get to much higher NTP numbers.

Columns E.3 - performance drop - w/o and w/ UF expresses the performance drop when all NTP users load the bottleneck at the same. When considering UF, lower values are obtained.

In Table III, we see that if we consider the value of 936.3 Mbps in desired SDR (E.1), the bandwidth needed really corresponds to the set value of 3000 Mbps at the bottleneck, see column E.4 (resulting average bandwidth value at the identified bottleneck). Which also corresponds to a hypothetical situation where the NBR_{MAX} would increase to a maximum of 3000 Mbps. However, if we take the real traffic that is now at the sites, counting desired SDR (E.1) with half value of 468.2 Mbps, we could reduce the capacity at all three sites to 1500 Mbps. Moreover, if we take UF into account, see column E.5 (resulting average bandwidth value at the identified bottleneck with respect to NTP behaviour), the value could be set to 1000 Mbps at the bottleneck. For better display, these values have also been plotted in the graph in Fig. 6.



Fig. 6. Comparison of site design, for NTP SDR 936.3 Mbps and 468.2 Mbps.

Conversely, if we were to consider a situation where the individual NTPs would have desired SDR (E.1) greater than 1000 Mbps, e.g. 1500 Mbps, the required capacity of site bottleneck (A.1) would have to be increased, in our case to 4500 Mbps.

Furthermore, we then confirmed that the NBR_{AVG} does not affect the required bottleneck capacity, until it exceeds the NBR_{MAX} value. In other words, if the UF is equal to 1, it corresponds to the calculation in E.4, i.e. to the full load (however, it is recommended to keep the UF in the values 0.4 to 0.5). In the situation when the NBR_{MAX} exceeds the value of the set capacity on the bottleneck, the capacity needs to be logically increased.

VI. CONCLUSION

Capacity planning for communication sites is a multifaceted endeavor that demands strategic resource management and infrastructure optimization. Providers, both large and small, grapple with this intricate task, employing distinct procedures and tailored solutions. Notably, smaller companies can benefit from the methodology developed by some authority such as CTU.

In this article, we scrutinized the involvement of communication sites operated by a major provider, both prior to and following an upspeed enhancement. Leveraging CTU's methodology, we proposed cost-effective adjustments. Assuming differential pricing for a 3000 Mbps line versus a 1000 Mbps line, our calculations indicated an immediate bandwidth requirement of 1500 Mbps. Accounting for utilization factor (UF), this aligns closely with the pre-upspeed value of 1000 Mbps.

Furthermore, if the operator intends to augment the number of NTPs at these sites, our analysis—guided by CTU's methodology—revealed the feasibility of adding one additional NTP within the existing 3000 Mbps limitation. To accommodate a fifth NTP, however, a bandwidth upgrade to 3500 Mbps would be necessary.

REFERENCES

- Czech Telecommunications Office (CTU). Metodika pro měření a vyhodnocení datových parametrů mobilních sítí elektonických komunikací [online]. 1. 3. 2021, verze 2.3 [cit. 27.2.2024]. Available at https://shorturl.at/anDGJ
- [2] Czech Telecommunications Office (CTU). Metodika pro vyhodnocování dopadu kapacity sítí elektronických komunikací na výkon služeb přístupu k internetu [online]. 1.1. 2023, verze 1.0 [cit. 27.2.2024].Available at https://shorturl.at/hoBDI
- [3] Czech Telecommunications Office (CTU).Metodika pro měření a vyhodnocení datových parametrů pevných sítí elektronických komunikací [online]. 1. 10. 2018, verze 2.0 [cit. 27.2.2024]. Available at https://shorturl.at/cfpyM
- [4] Czech Telecommunications Office (CTU). Netcalculator (Kalkulačka pro výpočet dopadu kapacity sítě elektronických komunikací).[cit. 27.2.2024]. Available at https://shorturl.at/cerxy
- [5] Huawei. WTTx Capacity White Paper [online]. 18.3. 2016, issue 01 [cit. 28.2.2024]. Available at https://carrier.huawei.com/en/technicaltopics/wireless-network/WTTx/WTTx-white-paper
- [6] X. Ge, S. Tu, G. Mao, V. K. N. Lau and L. Pan, "Cost Efficiency Optimization of 5G Wireless Backhaul Networks," in IEEE Transactions on Mobile Computing, vol. 18, no. 12, pp. 2796-2810, 1 Dec. 2019, doi: 10.1109/TMC.2018.2886897.
- [7] Mahmood, A., Kiah, M.L.M., Azzuhri, S.R. et al. Wireless backhaul network's capacity optimization using time series forecasting approach. J Ambient Intell Human Comput 12, 1407–1418 (2021). https://doi.org/10.1007/s12652-020-02209-2
- [8] A. Mahmood, L. Binti Mat Kiah, S. R. B. Azzuhri and A. N. Qureshi, "Wireless Backhaul Network Optimization Using Automated KPIs Monitoring System Based on Time Series Forecasting," 2018 IEEE World Symposium on Communication Engineering (WSCE), Singapore, 2018, pp. 1-6, doi: 10.1109/WSCE.2018.8690543.

Optimizing components for Dilithium and Kyber unified hardware implementation

1st Patrik Dobiáš Department of Telecommunications Brno University of Technology Brno, Czech Republic xdobia13@vut.cz 2nd Lukáš Malina Department of Telecommunications Brno University of Technology Brno, Czech Republic malina@vut.cz

Abstract—As the ongoing standardization process of postquantum schemes yields initial outcomes, it is important not only to focus on optimization of standalone implementations but also to explore the possibilities of combining multiple schemes into one unified architecture. In this paper, we focus on exploring the combination of two NIST selected schemes, namely the CRYSTALS-Dilithium digital signature scheme and the CRYSTALS-Kyber key encapsulation scheme. We present optimized designs of unified hardware components that can be used as building blocks for these schemes. All implemented components outperform state-of-the-art implementations in both hardware utilization and performance.

Index Terms—CRYSTALS-Dilithium, CRYSTALS-Kyber, FPGA

I. INTRODUCTION

In 2022, the National Institute of Standards and Technologies (NIST) selected four post-quantum cryptography (PQC) schemes for near-term standardization as a result of its standardization process. CRYSTALS-Kyber (which we refer to as 'Kyber' for the rest of the paper) for key encapsulation, and CRYSTALS-Dilithium (which we refer to as 'Dilithium' for the rest of the paper), Falcon and SPHINCS+ for digital signatures.

From the beginning of this post-quantum standardization process many works proposed hardware implementation of PQC schemes focusing on high-performance implementations [1], [2], [3], low-area implementations [4], [5] or comparison between implementations [6]. However, most of these works focus on a standalone implementation for a specific security level. Given that securing communication in a real-world scenario typically requires both key encapsulation and digital signature schemes on the device, it would be advantageous to integrate at least two of these schemes into a single hardware implementation. Otherwise, we would have to use multiple hardware implementations, leading to a large area overhead. Additionally, for server-side applications, we usually want to support a wider variety of security levels to provide different levels of security based on client's options.

Due to their similarities, merging Dilithium and Kyber appears to be a good choice, and this work focuses on exploring possibilities to combine individual blocks of these schemes.

A. Related Work

To receive real-world results, some researchers focus on the unification of multiple schemes, supporting all of their security levels. Authors in [7] implement a unified polynomial multiplier for Dilithium and Saber. They propose to use same Dilithium's Number Theoretic Transform (NTT) multiplier also for Saber multiplication, which would lead to negligible probability of producing wrong result. Following this approach, in [8], authors implement a complete unified architecture for Dilithium and Saber. Even though Saber is not considered in the NIST standardization process anymore, their work is a valuable insight into the combination of seemingly incompatible schemes.

On the unification of multiple signature schemes focused authors in [9]. They unify the Dilithium signature scheme, supporting all security levels, with Falcon verification phase for specific use cases.

In [10], authors propose the first implementation of a unified architecture for Dilithium and Kyber. They identify the main components that would have a significant impact, such as polynomial multiplication and Keccak with rejection samplers, and focus on their optimized sharing. The result of their work is single implementation of all security levels of Dilithium and Kyber, while being comparable with state-of-the-art standalone implementations. Following this work, authors in [11] focused on further optimizations of the unified NTT multiplication for Dilithium and Kyber. They also examined the trade-off between area consumption and performance when changing the number of butterfly units. Their results promise lower area usage with better performance than [10].

B. Contribution

The contribution of the paper is as follows:

- We design main components needed for the implementation of unified architecture for Dilithium and Kyber schemes.
- We propose a novel way of decoding polynomial coefficients that allows for the use of different bit lengths at little cost.
- We implemented and verified all proposed designs on the Zynq Ultrascale+ platform. The results show that

	Input Processing		Output Processing	
Input	Head Reorder	Butterflies	Tail Reorder	Output
0 1 2 3	0 1 2 3			
128 129 130 131	2 3 130 131	0 128 1 129	0 1 128 129	
4 5 6 7	4 5 6 7	2 130 3 131	128 129 130 131	0 1 2 3
132 133 134 135	6 7 134 135	4 132 5 133	4 5 132 133	128 129 130 131

Fig. 1. Coefficients flow during first iteration of NTT.

our implementations are superior in both utilization and performance compared to existing implementations.

C. Organization

The remainder of this paper is organized as follows. in Section II, we describe in short post-quantum schemes Dilithium and Kyber focusing on polynomial multiplication, sampling, and decoding/encoding operations. Section III focuses on the unified implementation of the selected components. in Section IV we present and discuss achieved results and compare them with existing implementations, and finally, in Section V, we conclude this paper.

II. BACKGROUND

In this section, we provide a brief overview of two postquantum schemes, which were selected by NIST for standardization, Dilithium and Kyber, highlighting the key operations that were selected for unification.

A. Dilithium

Dilithium is a lattice-based digital signature scheme, whose security is based on Module Learning With Errors (M-LWE) and Shortest Integer Solution (SIS) problems. For the key generation, signing and verification phases, multiple different operations are needed. Here, we will describe only the operations that we decided to unify with Kyber due to shared characteristics.

1) Polynomial Arithmetic: Algebraic operations are performed over the polynomial ring $R_{8380417} = \mathbb{Z}_{8380417}[X]/(X^{256} + 1)$. Because Dilithium has 512^{th} root of unity, polynomial multiplication is performed using complete Number Theoretic Transform (NTT).

2) Coefficients Sampling: Dilithium uses two forms of coefficients sampling. The first occurs during the generation of the matrix A, when a rejection is applied to the output of the SHAKE-128 hash function. The second during the generation of secret key and error vectors, when rejections are applied to the output of the SHAKE-256 hash function.

3) Coefficients Unpacking/Packing: The process of unpacking and packing coefficients is used to reduce the size of keys and signatures. This is done by using only a certain number of bits of the coefficients. It is required to support unpacking/packing of 20-bit, 18-bit, 13-bit, 10-bit, 6-bit, 4-bit and 3-bit long coefficients.

B. Kyber

Kyber is a lattice-based key encapsulation scheme, whose security is based on the M-LWE problem. The key generation, encapsulation and decapsulation phases require multiple different operations. Here, we will once more specify only operations that we aim to unify with Dilithium.

1) Polynomial Arithmetic: Algebraic operations are performed over the polynomial ring $\mathbb{Z}_{3329}[X]/(X^{256}$ + R_{3329} = 1). Because Kyber does not have 512^{th} root of unity, polynomial multiplication is performed using incomplete Number Theoretic Transform (NTT).

2) Coefficients Sampling: Similar to Dilithium, Kyber employs two types of coefficients sampling. The first occurs during the generation of the matrix *A*, when a rejection is applied to the output of the SHAKE-128 hash function. The second is used for secret key and error vectors sampling, during which only a specific number of bits without rejection are used to sample the output of SHAKE-256.

3) Coefficients Decoding/Encoding: To reduce the size of keys and ciphertexts, decoding and encoding of coefficients is used. This is done by using only a certain number of bits of the coefficients. It is required to support decoding/encoding of 12-bit, 11-bit, 10-bit, 5-bit,4-bit and 1-bit long coefficients.

III. IMPLEMENTATION DETAILS

In this section we describe implementation details of unified hardware components that can be shared between Dilithium and Kyber. When designing these components, our focus is on ensuring high performance while maintaining minimal hardware utilization.

A. Polynomial Multiplication

We implement polynomial multiplication that uses two parallel butterfly unit cores with flexibility for both Dilithium and Kyber as proposed in [10]. These butterfly units process together 4 coefficients for Dilithium and 8 coefficients for Kyber every clock cycle. The polynomial multiplication unit can perform NTT/NTT⁻¹ operations in 512 cycles for Dilithium and 224 cycles for Kyber, and pointwise multiplication in 256 cycles for Dilithium and 128 cycles for Kyber. Achieving the same latency as presented in [10] and [11].

To eliminate the requirement of using a single memory block for each butterfly unit, we propose using head and tail reorder units, as described in [6]. These units use a 96-bit long register to store coefficients that have not been processed yet, to reorder them in a way, so they are always arranged in the same order in memory. To better understand the flow of coefficients in this component, the first iteration of Dilithium's NTT is depicted in Figure 1. In the first cycle, the initial four coefficients are fed on input and stored in the head reorder registers. In the subsequent clock cycle, the first four coefficients from the latter half of the polynomial are fed on input, while the coefficients with indices 0, 1, 128, 129 are directed to butterflies, and the coefficients 2, 3, 130, 131 are stored in head reorder registers and are used in the next clock cycle. The tail reorder operates in a similar way, storing coefficients 0, 1, 128, 129 in the register and outputting 0, 1, 2, 3in the following cycle, while storing 128, 129, 130, 131 for the subsequent cycle. Thus, ensuring the preservation of the coefficient order in memory.

B. Coefficients Sampling

We implement a unified coefficients sampling unit supporting all of the samplings described in Section II. This unit uses the standard AXI stream interface on both input and output. The input data are 64-bit long to align with our Keccak implementation, and output data are 96-bit long representing the size of 4/8 coefficients for Dilithium/Kyber.

The architecture of this unit is depicted in Figure 2. First, incoming bytes are stored in a shift register, since we only need up to 48 bits every clock cycle. Then, these bytes are processed concurrently by all rejection units and stored in registers. Register sharing occurs only among similar rejection units (e.g., matrix A rejection for both schemes), as employing a single register for every unit would necessitate large multiplexor, leading to excessive utilization of lookup tables. Moreover, as the cbd and rej_bounded units produce coefficients with only 3 respectively, 4 valid bits, this approach leads to a very small registers utilization overhead.



Fig. 2. Architecture of the coefficients sampling unit. The orange rectangles represent registers.

Moreover, we have opted not to integrate Keccak with coefficients sampling, deviating from the approach in [10]. By separating Keccak, it can serve various hashing purposes, thereby enhancing the versatility of the designed crypto core. Our Keccak component is derived from previous work, with minor modifications made to its interface.

C. Coefficients Decoding

We propose a novel way for coefficients decoding supporting additional coefficient sizes at a cost of little to no area overhead. Because of that, we were able to combine Unpack operation for Dilithium with Decode operation for Kyber into one component. The component has standard AXI stream input interface with 64-bit data bus that is decoded into 1bit, 3-bit, 4-bit, 5-bit, 10-bit, 12-bit, 13-bit, 18-bit and 20-bit long coefficients. These coefficients are then outputted using a standard memory interface.

Our proposal uses four 48-bit long registers that store input data rotated left by multiple of 16 (e.g. $reg_0 := input[47 : 0]$, $reg_1 := input[31 : 0] || input[63 : 48]$), the registers are enabled based on current state of decoding. These registers are then used in four 4-to-1 multiplexers, where they are rotated left again, this time with multiples of 4 with overlap from neighboring registers. The output of these multiplexers is fed to the final 4-to-1 multiplexer. The output from this register has valid coefficients data aligned from LSB. The example data flow during decoding of 5-bit, 10-bit, 11-bit and 20-bit long coefficients is shown in 3.



Fig. 3. Data flow during decoding of 5-bit, 10-bit and 20-bit long coefficients (left) and 11-bit long coefficients (right).

This approach allows decoding four 3-bit, 4-bit and 10bit long coefficients and two 13-bit, 18-bit and 20-bit long coefficients for Dilithium, and eight 1-bit, 4-bit and 5-bit long coefficients and four 10-bit, 11-bit and 12-bit long coefficients for Kyber.

This approach consumes additional registers compared to solution described in [12], specifically for our settings 192 vs. 104 registers. However, it consumes a significantly smaller number of lookup tables, since we only need to determine when to enable registers compared to the 1-to-N demultiplexor used when storing input in shift registers as in [12]. This can be seen in the results, where our component uses only 1.6x more LUTs while achieving 2-4x throughput and supporting more decoding bit lengths.

D. Coefficients Encoding

When implementing coefficients encoding, we adopt a methodology similar to that proposed in [12]. Nevertheless, to mitigate the complexity in the final 1-to-16 multiplexor, which in our scenario would be even more extensive due to our support for a broader range of encoding bit lengths, we introduce an additional register. This register stores coefficients rotated by multiples of 4 bits, depending on the current encoding state, which makes use of the 1-to-4 multiplexor. Subsequently, on the output side, only another 1-to-5 multiplexor is necessary. Using this approach, we reduce the complexity of the multiplexor network, leading to lower hardware utilization.

IV. RESULTS AND COMPARISON

In this section we present results of implemented components. All of these components were implemented using VHDL and tested using the reference C implementations and python framework cocotb. The results are obtained using synthesis in Vivado 2022.2, and the UltraScale+ ZCU102 platform was chosen as the target to ensure a fair comparison with the designs implemented in [10] and [11], which also use this platform.

A. Utilization and Performance Results

We target all components at a frequency of 400 MHz. The corresponding utilization of hardware resources is presented in Table I. The table illustrates the utilization of hardware resources in implemented hardware components, highlighting the demands posed by each in terms of lookup tables (LUTs), flip-flops (FFs), block RAMs (BRAMs), and DSPs. In particular, the most resource-intensive components are Keccak, which consumes 3052 LUTs and 1697 FFs, and Polynomial Multiplication, which utilizes 2516 LUTs, 1200 FFs, 1 BRAM, and 4 DSPs. The remaining components collectively consume fewer resources compared to these two. Specifically, the Sample unit consumes 720 LUTs and 380 FFs, the Decode unit uses 388 LUTs and 314 FFs, and the Encode unit utilizes 530 LUTs and 216 FFs.

Component	LUT	FF	BRAM	DSP
Polynomial Multiplication	2514	1200	1	4
Keccak	3052	1697	0	0
Sample	720	380	0	0
Decode	388	314	0	0
Encode	530	216	0	0

TABLE I

HARDWARE UTILIZATION OF IMPLEMENTED COMPONENTS.

B. Comparison With Other Unified Designs

We conduct a comparison between our implementations and those presented in [10] and [11] in terms of hardware utilization and operational frequency. We cannot compare our implementation of coefficients sampling unit, because [11] does not include this unit in their design, and [10] reports its performance with Keccak, making a fair comparison difficult. 1) Polynomial multiplication: Comparison of the polynomial multiplication unit is presented in II. Our implementation demonstrates significantly lower utilization in both lookup tables and register usage compared to existing implementations. As mentioned in Section III, our unit maintains the same latency for all operations as the other works while achieving the highest working frequency. Therefore, our component offers the best performance with the lowest hardware utilization.

Work	LUT	FF	BRAM	DSP	Frequency [MHz]		
[10]	3487	1918	1	4	270		
[11]	2893	2356	4.5*	4	342		
This Work	2514	1200	1	4	400		
*Penorted with memories for storing coefficients							

TABLE II COMPARISON OF UNIFIED POLYNOMIAL MULTIPLICATION UNIT.

2) Coefficients Decoding/Encoding: We compare our unified designs of coefficients decoding and encoding units in Table III. Because the authors in [10] did not unify these components in their design, but rather employed separate components for Dilithium and Kyber, we present their results as a sum of those components, as they together perform same operations. We can see that our components are superior in terms of reduced area utilization and higher operational frequency.

Work	LUT	FF	Frequency [MHz]
Decode - [10]	552	362	270
Encode - [10]	1099	371	270
Decode - This Work	388	314	400
Encode - This Work	530	216	400

^aSum of unpack and decode units. ^bSum of pack and encode units. TABLE III

COMPARISON OF COEFFICIENTS DECODING AND ENCODING UNITS.

V. CONCLUSION

In this work, we have presented implementation of main building blocks needed for the unified design of Dilithium and Kyber schemes. These implementations outperform other unified designs in both area utilization and performance. To achieve these results, we proposed novel way of decoding coefficients leading to efficient decoding of all bit sizes needed for Dilithium and Kyber schemes. We also combined ideas of [10] and [6] while implementing polynomial multiplication unit to achieve high performance with low area utilization. For coefficients sampling unit we suggested separation of Keccak from the unit, to offer even more usage options of the unified crypto core. In future, we plan to finish our unified design, by combining all of these components and apply mechanisms to protect it against side-channel and fault-injection attacks.

VI. ACKNOWLEDGMENT

This work is supported by Ministry of the Interior of the Czech Republic under Grant VJ02010010.

References

- A. Ferozpuri and K. Gaj, "High-speed fpga implementation of the nist round 1 rainbow signature scheme," in 2018 International Conference on ReConFigurable Computing and FPGAs (ReConFig). IEEE, 2018, pp. 1–8.
- [2] X. Li, J. Lu, D. Liu, A. Li, S. Yang, and T. Huang, "A high speed postquantum crypto-processor for crystals-dilithium," *IEEE Transactions on Circuits and Systems II: Express Briefs*, 2023.
- [3] M. Schmid, D. Amiet, J. Wendler, P. Zbinden, and T. Wei, "Falcon takes off-a hardware implementation of the falcon signature scheme," *Cryptology ePrint Archive*, 2023.
- [4] M. Andrzejczak, "The low-area fpga design for the post-quantum cryptography proposal round5," in 2019 Federated Conference on Computer Science and Information Systems (FedCSIS). IEEE, 2019, pp. 213–219.
- [5] D. Kales, S. Ramacher, C. Rechberger, R. Walch, and M. Werner, "Efficient fpga implementations of lowmc and picnic," in *Topics in Cryptology–CT-RSA 2020: The Cryptographers' Track at the RSA Conference 2020, San Francisco, CA, USA, February 24–28, 2020, Proceedings.* Springer, 2020, pp. 417–441.
- [6] V. B. Dang, K. Mohajerani, and K. Gaj, "High-speed hardware architectures and fpga benchmarking of crystals-kyber, ntru, and saber," *IEEE Transactions on Computers*, vol. 72, no. 2, pp. 306–320, 2022.
- [7] A. Basso, F. Aydin, D. Dinu, J. Friel, A. Varna, M. Sastry, and S. Ghosh, "Where star wars meets star trek: Saber and dilithium on the same polynomial multiplier," *Cryptology ePrint Archive*, 2021.
- [8] A. C. Mert, D. Jacquemin, A. Das, D. Matthews, S. Ghosh, S. S. Roy et al., "A unified cryptoprocessor for lattice-based signature and keyexchange," *IEEE Transactions on Computers*, 2022.
- [9] P. Karl, J. Schupp, T. Fritzmann, and G. Sigl, "Post-quantum signatures on risc-v with hardware acceleration," ACM Transactions on Embedded Computing Systems, 2023.
- [10] A. Aikata, A. C. Mert, M. Imran, S. Pagliarini, and S. S. Roy, "Kali: A crystal for post-quantum security using kyber and dilithium," *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 70, no. 2, pp. 747–758, 2022.
- [11] S. Mandal and D. B. Roy, "Kid: A hardware design framework targeting unified ntt multiplication for crystals-kyber and crystals-dilithium on fpga," arXiv preprint arXiv:2311.04581, 2023.
- [12] Y. Xing and S. Li, "A compact hardware implementation of cca-secure key exchange mechanism crystals-kyber on fpga," *IACR Transactions on Cryptographic Hardware and Embedded Systems*, pp. 328–356, 2021.

Overdoping effect with Zr and Hf on the oxidation behaviour of FeCrAl-Hf by means of Atom Probe Tomography

Samer I. Daradkeh Central European Institute of Technology Brno University of Technology Brno, Czech Republic 252679@vutbr.cz Oscar Recalde Department of Materials- and Earth Science Technical University of Darmstadt Darmstadt, Hessen, Germany o.recalde@aem.tu-darmstadt.de

Marwan S. Mousa Department of Physics Mu'tah University Al-Karak, Jordan marwansmousa@yahoo.com

Dinara Sobola Institute of Physics of Materials, Department of Physics Czech Academy of Sciences, Faculty of Electrical Engineering and Communication Brno, Czech Republic sobola@vut.cz

Abstract—The study investigated the oxidation behaviour and grain boundary diffusion of minor/major elements of FeCrAl alloys, doped with over-critical concentrations of reactive elements (REs) Zr and Hf. While the formation of $\alpha - Al_2O_3$ scale on these alloys is conventionally attributed to inward oxygen transport along grain boundaries, this research proposes that metal ion outward diffusion also contributes to the development of oxide scales and their microstructural characteristics. Samples were analyzed after thermal exposure at 1100 °C using scanning electron microscopy (SEM), transmission electron microscopy (TEM), and atom probe tomography (APT). Results revealed increased oxide growth, deeper internal oxidation, and RE-oxide formation near and at oxide grain boundaries due to enhanced inward and outward diffusion resulting from overdoping. The impact of overdoping varied with RE type and concentration, influenced by solubility, ionic size, and electronic structure of alumina. Notably, Zr-doped samples maintained alumina adhesion to the alloy after thermal exposure, whereas severe spallation occurred in Hf-doped samples.

Index Terms—Overdoping, Reactive element, Grain boundary, Alumina

I. INTRODUCTION

FeCrAl alloys are renowned for their exceptional resistance to corrosion and oxidation, making them ideal for hightemperature applications like heating elements and furnace linings [1], [2]. At temperatures of 1100 °C, these alloys develop a protective $\alpha - Al_2O_3$ scale, crucial for their outstanding oxidation resistance. Alumina, the primary component of this scale, offers remarkable thermodynamic stability and low defect concentration due to its high point defect formation energy, especially for charged defects [3]. However, extended exposure to high temperatures can lead to breakaway corrosion

Identify applicable funding agency here. If none, delete this.

as aluminum reservoirs in the alloy are depleted over time. Extensive lifetime modelling has substantiated this phenomenon.

The role of Reactive Elements (REs) in modifying the oxidation process and oxide layer formation has been extensively studied [4], [5]. REs enhance scale adhesion to FeCrAl alloy, reduce parabolic growth, and impede outward diffusion of cation ions. Their benefits are attributed to their large ion size and high oxygen affinity, with distribution and concentration determining their effectiveness [6]. Research suggests that RE concentration should remain below the solubility limit in the alloy to avoid adverse effects like the formation of secondphase compounds such as yttrium aluminum oxide (YAG) and Y-rich oxide particles at alumina Grain Boundaries (GBs), which compromise scale adhesion [7].

While alumina primarily grows through inward diffusion, evidence suggests the involvement of outward diffusion of cations, driven by the gradient in oxygen chemical potential across the scale [8]. The composition of newly formed oxide in surface GBs over time and alumina scale thickness provides quantitative data on GB flux and aluminum diffusivity, where in this study it aligns with Fick's 1st law.

This paper investigates FeCrAl alloys overdoped with REs (specifically, Zr and Hf) to understand elemental outward segregation in thermally grown alumina scales, shedding light on alloy lifetime. Atom probe tomography (APT) was employed for 3-D mapping of atoms in the oxide scale at near-atomic resolution, with Transmission electron microscopy (TEM) used for correlation.

Table 1 Alloys composition as analyzed by ICP-OES/EDS

in the second									
Alloy/Compositions (at. %)	Fe	Al	С	Cr	Mn	Si	Y	Zr	Hf
Fecralloy TM (OES)	67.4	9.65	0.07	22.03	0.19	0.56	0.06	0.06	_
FeCrAl-Zr (OES)	66.5	10.27	0.3	22.45	-	-	-	0.5	-
FeCrAl-Hf (EDS)	66.8	9.64	0.3	22.74	-	-	-	-	0.52

II. EXPERIMENTAL PROCEDURES

A. Sample Preparation and Exposure

The samples were prepared by arc melting technique under a protective Ar gas atmosphere using purity-high elements. Composition analysis was conducted using Spark-Optical Emission Spectrometry (Spark-OES) for FecralloyTM and FeCrAl–Zr, and Energy-Dispersive Spectroscopy (EDS) for FecralloyTM, FeCrAl–Zr, and FeCrAl–Hf. The Zr and Hf levels in both samples exceeded the optimal amount of REs in alloys, typically around 0.05 at.%.

After sectioning with a diamond wire saw and annealing at 1200 °C for 24 hours for homogenization, the samples underwent polishing with abrasive SiC paper and oil-base diamond suspension, followed by ultrasonic cleaning. Initial hightemperature oxidation exposures were conducted in laboratory air at 1100 °C for 200 hours to form an adherent alumina scale.

Taper grinding was performed on one side of the preoxidized specimens at an angle of 0.08°, removing the outer equiaxed layer of scale to reveal the columnar grains of the inner layer. This was followed by further oxidation at 1100 °C to measure the new outward-growing oxide associated with grain boundaries. The optimum conditions to grow measurable



Fig. 1. Schematic illustration of the experimental procedures used to observe the outward growth of Alumina due to Al outward GB diffusion a typical alumina oxide structure after first oxidation. b) After the first oxidation, the sample was polished by removing a portion of the formed oxide layer. c) after Re-oxidation after taper polishing

oxide ridges in the second exposure were determined to be 60 minutes for FecralloyTM, 20 minutes for FeCrAl–Zr, and 40 minutes for FeCrAl–Hf.

B. Analytical Techniques

The investigation of oxide scales on Fecralloy[™] and Fe-CrAl–Zr involved scanning electron microscopy/focused ion beam (SEM/FIB), transmission electron microscopy (TEM), and laser-assisted atom probe tomography (APT). However, after the second exposure, the FeCrAl–Hf sample exhibited notable spallation and cracks in most of the regions of interest. so TEM and APT were not employed in this condition. The SEM/FIB analysis utilized a Zeiss Auriga 60, while the TEM examination employed an FEI Thermo Fisher Themis Z. Laser-assisted atom probe tomography (APT) was conducted using a CAMECA LEAP 4000 XHR equipped with a UV laser (wavelength 355 nm) and a spot size less than 1 μ m. Sample preparation for TEM and APT involved SEM/FIB, with a maximum accelerating voltage of 30 kV for thinning steps and 2 kV to minimize Ga ion damage to the sample.

The initial preparation of TEM and APT samples involved creating lamellae containing alumina grain boundaries and attaching them to Micro-Post/TEM grids. Needle-shaped APT specimens with a shank angle of about 10° and an apex radius of less than 100 nm were prepared using the lift-out technique. These oxide samples containing grain boundary oxide were aligned approximately parallel to the field evaporation direction to achieve a high electric field for ion removal. In the case of materials like alumina with poor electrical conductivity, a combination of ultra-short laser pulses and high electrical voltage was used. APT analysis was conducted at a temperature of 35 K, with a UV-laser pulse frequency of 125-200 kHz, a laser pulse energy of 30-40 pJ, and a detection rate of 0.25%. Reconstruction of the atom probe data was performed using AP Suite 6 software, based on the tip evolution profile. Various analyses, including 3D atom distribution, 1D concentration profile across grain boundaries and interfaces, and proxigram analysis, were employed. The Thermo-Fisher Themis-Z was operated in S/TEM mode at an accelerating voltage of 300 keV, utilizing both brightfield (BF) and high-angle annular dark-field (HAADF) image modes.

Quantitative measures of solute segregation at interfaces, typically evaluated using methods like those proposed by Krakauer and Seidman [9], often simplify characterization to a single value, known as the Gibbsian interfacial excess Γ .

$$\Gamma_i = N_i^{excess} / A = (1/A) N^{vol} [C_i^{vol} - C_i^{\alpha} \xi - C_i^{\beta} (1-\xi)]$$
(1)

where N_i^{excess} is the excess no. of atoms at interface, A is the interface area over which Γ_i is determined, N_i^{α} and N_i^{β} are the no. in α -phase and β -phase adjacent to the interface, ξ is the gibbs dividing surface, and C represents the concentration of an element *i* in phase α or phase β .

C. Grain Boundary Outward Aluminum Flux

Oxide ridge volume measurements were obtained after the second oxidation using SEM/FIB. And it can be calculated using the following relation [10]:

$$J_{GB}^{Al} = N_{GB}^{Al} / L_{GB}.t \tag{2}$$

where J_{GB}^{Al} is the number of Al ions per length and time t, L_{GB} is the lateral length of the GB, N_{GB}^{Al} is the no. of segregated Al atoms along GB. The unit cell volume of $\alpha - Al_2O_3$ ($V_{unitcell}$) is 0.254 nm^3 and with the no. of Al ions per unit cell of $\alpha - Al_2O_3$ ($N_{unitcell}^{Al}$) being 12. This yields:

$$N_{GB}^{Al} = (V_m / V_{unit \, cell}) N_{unit \, cell}^{Al} \tag{3}$$

Then:

$$J_{Al} = 12.A_{ridge}/V_{unit\,cell} \tag{4}$$

where A_{ridge} is the cross-sectional area of ridge.

III. RESULTS

A. Oxide Scale Microstructure (SEM)

After an initial 200 h of oxidation at 1100 °C, the average external oxide scale thickness of FecralloyTM, FeCrAl–Zr, and FeCrAl–Hf was found to be $3.1 \pm 0.3 \ \mu\text{m}$, $5.6 \pm 2.7 \ \mu\text{m}$, and $4.9 \pm 2.2 \ \mu\text{m}$, respectively. Whereas FecralloyTM did not form internal oxides, quite extensive internal oxidation had occurred on both FeCrAl–Zr and FeCrAl–Hf samples to average depths of 79.2 ± 3.2 \ \mu\text{m} and 18.2 ± 2.2 \ \mu\text{m}, respectively.

After bevel polishing and re-oxidation in the second exposure, a severe scale spallation occurred at FeCrAl-Hf preventing further investigation using APT and TEM. It has been found that is transient Fe-rich scale had formed, most of it had spalled extensively to expose partially spalled alumina which exhibited ridges beneath. Distinct areas of Hf oxide were present in some areas from which the alumina had spalled. While Topographies of the scales formed on FecralloyTM and FeCrAl–Zr after the same re-exposure. In addition, new alumina ridges formed over grain boundaries on the alumina scale.

B. Oxide Scale and Composition (TEM+APT)

A TEM and APT investigation was performed on Fecral-loyTM and FeCrAl–Zr after the 2nd thermal exposure. The results from both instruments coincided and confirmed each other.

Figure 2 and 3 show the elemental maps of Fe, Cr, and Zr ions for FecralloyTM. The atomic fraction in the oxide layer for Fe ranges between 4.05×10^{-4} to 3.3×10^{-3} , Cr between 4.53×10^{-5} to 2.1×10^{-3} and Zr between 5.3×10^{-5} to 3×10^{-3} , meaning nearly no segregation to the oxide layer.



Fig. 2. APT reconstruction elemental maps distribution (Al, O, and C, respectively) of a tip-sample from oxide scale grown on FecralloyTM, indicating no Fe, Cr, or Zr ions segregation to alumina scale through GB.

Figure 2 depicts the absence of Fe, Cr, or Zr ion segregation to the alumina scale via grain boundaries, which indicates the feasibility of controlled and limited doped REs-FecralloyTM.

A further APT sample of Fecralloy[™] was taken from a sample of the alumina grain boundary nearest to the oxide ridge (near the oxide/air interface), as in figure 4, showing the spatial distribution of Ag (deposited silver layer to mitigate the charging effect during SEM/FIB work), O, Fe, and Al ions. Fe was found in the grain boundaries of the alumina layer close to the oxide/air interface. Also, there are traces of Cr in the vicinity of the same region. The Gibbsian interfacial

excess of Fe (FecralloyTM sample) at the alumina oxide GB was calculated and it was $\Gamma_{Fe}^{GB} = 0.42 \text{ nm}^{-1}$.



Fig. 3. a) HAADF-STEM image of the oxide layer grown on FecralloyTM after the second exposure. b–d) EDS maps of Fe, Cr, and Zr ions, respectively. EDS maps showing the segregation of Zr ions to the oxide layer



Fig. 4. a) The location from which the APT sample was prepared. b) Reconstructed 3-D APT images showing the atomic distribution of Fe, Al, O, and Ag on FecralloyTM. c and d) Atomic density maps of Fe and Cr

In the case of FeCrAl–Zr, TEM and APT were utilized and a sample containing an oxide-alloy interface and GB were successfully prepared (Figures 5 and 6). Both the TEM and APT results, as depicted in Figures 5 and 6, coincide with each other. The results show the enhanced segregation of Fe, Cr, and Zr at the oxide/alloy interface and at the oxide scale, forming the Cr-rich region (as TEM results show - Figure 5b and c) as an effect of the overdoping. From Figure 5c, a notable rise in Cr concentration is observed (the black arrow in Figure 5a denotes the direction of concentration profile analysis), signifying the diffusion of Cr across the oxide/alloy interface.



Fig. 5. a) The oxide scale formed on FeCrAl–Zr alloy after the 2nd thermal exposure. b) Elemental map shows the distribution of Cr ions. c) 1D concentration profile along the indicated direction in (a)

Figure 6 shows the isosurface map of Fe and 3D atomic distribution of Cr indicating a triple junction point at the oxidealloy interface, representing the intersection point of adjacent



Fig. 6. a) APT reconstructed isosurface maps of Fe in an oxide layer grown on FeCrAl–Zr tip-sample after the 2nd exposure for 20 min. b) atom distribution of Cr



Fig. 7. Proximity histograms (proxigrams) of a) Triple junction point. b) Fe precipitate

oxide grain and the oxide/alloy interface. The triple junction point was confined between the alloy and two oxide grains and was identified by its hump-like location and shape. The emergence of the Fe-rich band and Cr-rich band gives an indication of Cr and Fe ion segregation as a result of overdoping, which coincides with the figure 5b and c (TEM results). Additionally, the presence of precipitates near the grain boundary (Figure 6a) may contribute to the segregation of Fe, Zr and Cr ions along it. Figure 7 shows the concentration of Fe, Cr, and Zr as a function of distance to the isoconcentration surface at the triple junction and precipitates region. These indicate that Cr (3.3 at%) is more prevalent than Fe or Zr at the triple junction region. Fe is relatively the major constituent of the precipitate region. The Gibbsian excess of Zr near the triple junction was calculated, and it was $\Gamma_{Zr}^{GB} = 0.7 nm^{-1}$.

C. Outward Flux of Al Along Alumina GB

The data acquired for the flux of Al along GB is shown in the figure, which is plotted using a double-logarithmic of the oxide thickness against the Al flux along GB. Very good agreement was found with results published by Nychka and Clarke [11]. Also, the results follow Fick's law of diffusion. From the figure, it can be deduced that the flux of Al along GB has enhanced in overdoped FeCrAl with Zr.



Fig. 8. Grain boundary flux of aluminum versus wedge oxide thickness for a) FeCralloyTM and FeCrAl–Zr alloy after re-oxidation for 1 h and 20 min, respectively, and b) Data for several alloys pre-oxidized for 120 h at 1100 °C in air, and re-oxidized for 4 h at 1100 °C in air [11]

IV. DISCUSSION

The sample preparation methodology focused on creating a specific microstructure on the oxide surface through the second polishing step and adjusting the second oxidation time accordingly to form ridges after the 2nd exposure. This microstructure's growth was found to be dependent on both time and oxide thickness. While the second polishing step is considered critical, it also poses a challenge as it may introduce additional dislocations, potentially leading to a higher flux. To mitigate this effect, the second polishing should utilize the smallest possible particle size and be carried out for at least 11 minutes. Alternative milling techniques have been reported to have drawbacks that impact the aluminum flux [12].

The evaluation of the Al flux using Eq. 4 disregards several factors that could affect the results, e.g., low-, and highangle GBs. Also, the leakage of atoms from GBs to adjacent grains was disregarded. However, Eq. 4 can serve its purpose by giving an estimated value of Al flux along GBs, which correlates with the 1st Fick's law of diffusion.

Previous research [13], [14] has shown that the inclusion of RE(III) elements results in a significant increase in grain boundary width as both the ionic radius and concentration of these elements increase. This widening of grain boundaries induces accumulated stress, which accelerates the precipitation of RE oxides and, in certain instances, causes oxide spallation. The most famous and reliable model to demonstrate the effect of REs are poisoned interface model (PIM) [15] proposed that REs and cation impurities segregate to the interface and react with the defects, having an impact on the oxide growth mechanism. However, the PIM mode contradicts the finding in Pint's paper [16], where it overlooks the role of RE doping in the scale and ignores the possibility of small voids at the oxide/alloy interface. According to the grain boundary segregation model, REs and cation impurities follow similar defect routes during oxide scale growth. Cation impurities typically segregate at grain boundaries like aliovalent cations, necessitating vacancies or interstitials to maintain electric neutrality and thereby increasing diffusivity in grain boundaries [17]. The first mechanism in the GB segregation model can be explained as the presence of REs at GBs, which acts as site blockers to other cation impurities and reduces their segregation. The other postulation of that model is the swampingout mechanism [18], supported experimentally using HRTEM and STEM EDS [19]. The recent trend to explain the effect of RE relies on the band structure of Al_2O_3 scales [20]–[22], where the electron and holes play a role in the scaling reaction (i.e., creation and annihilation of vacancies at interfaces). The involvement of electrons or holes stems from the defect structure near the band edge. Doping with reactive elements can modify the electronic structure by segregating and causing migration of grain boundary disconnections. A specific type of grain boundary defect, characterized by a step height (h) and a Burgers vector (b), contains both positively and negatively charged jogs [23]. The higher the concentration of reactive elements, the greater the migration of these disconnections, ultimately resulting in enhanced segregation.

There was an agreement on the practical limit of REs doping in commercial alloys, which is several hundred ppm. These limits can be affected by several factors like the solubility and the presence of C and S atoms. Exceeding the limits led to the formation of 2nd phases or intermetallic compound, spallation, and oxide layer failure [24]. The level of doping of the FeCrAl-Hf alloy led to disastrous degradation of its ability to form a protective alumina scale. The major undesirable effect was scale spallation, apparently promoted by the formation of $H f O_2$ at the alloy/scale interface, presumably due to differences in thermal expansion among the oxides, the lattice mismatch between alumina and hafina, and the increased Pilling-Bedworth ratio. The Hf and Zr considered an example of the element that has an extremely high affinity for oxygen [25], which is one of the reasons for the formation of internal oxidation. Meanwhile, the extent of the internal oxidation depth in both cases (overdoped sample) is governed by the oxygen solubility of oxygen in Hf and Zr, where the solubility in Hf is lower than Zr. The APT reconstruction analyses correlate with the STEM/EDS results, as both indicate the absence of Fe and Cr ion segregation at alumina GBs in the FecralloyTM sample, which reflects the feasibility of the REs acting as site blockers of outward ion segregation.

V. CONCLUSIONS

The results can be summarized as follows:

- The inward diffusion was enhanced in the overdoped samples where the inward diffusion is more enhanced by overdoping with Zr than overdoping with Hf based on the depth of internal oxidation.
- Based on the thickness of the oxide scale formed on the overdoped alloys, the outward diffusion of aluminum was enhanced.
- 3) In the overdoped FeCrAl with Zr, the Fe and Cr ion segregation within oxide grains was observed, which may attributed to the overdoping effect or as a remnant of the transient stage. Additionally, the presence of Zr ions was detected within the oxide grain boundaries of the overdoped sample with Zr.
- 4) The sample overdoped with Hf exhibited severe scale spallation after the second exposure, surpassing the level

observed in the Zr-overdoped sample, rendering further analysis of this alloy unfeasible.

5) The Fecralloy[™] sample showed favourable outcomes, such as reduced outward diffusion of Al and other constituents along alumina grain boundaries, decreased growth rate, and improved adhesion of the alumina scale to the alloy. Conversely, the Hf-overdoped sample experienced significant spallation after the second exposure, exceeding that observed in the Zr-overdoped sample.

ACKNOWLEDGMENT

The authors would like to thank the German Academic Exchange Service (DAAD) scholarship. This work was carried out with the support of the Karlsruhe Nano Micro Facility (KNMFi, www.knmf.kit.edu), a Helmholtz Research Infrastructure at Karlsruhe Institute of Technology (KIT, www.kit.edu) (Project No. ha 2020-024-028730). Also, special thanks to Dr. Dinara Sobola and Prof. Dr.-Ing. Martin Heilmaier for the scientific discussion. The authors would like to thank CEITEC Nano Research Infrastructure for their assistance with this research.

REFERENCES

- Josefsson H, Liu F, Svensson Jalvarsson M, Johansson LOxidation of FeCrAl alloys at 500–900°C in dry O2. Materials and Corrosion. 2005 Nov;56(11):801–805. Available from: http://dx.doi.org/10.1002/maco.200503882.
- [2] Badini C, Laurella F. Oxidation of FeCrAl alloy: influence of temperature and atmosphere on scale growth rate and mechanism. Surface and Coatings Technology. 2001 Jan;135(2–3):291–298. Available from: http://dx.doi.org/10.1016/S0257-8972(00)00989-0.
- [3] Tang C, Jianu A, Steinbrueck M, Grosse M, Weisenburger A, Seifert HJ. Influence of composition and heating schedules on compatibility of FeCrAl alloys with high-temperature steam. Journal of Nuclear Materials. 2018 Dec;511:496–507. Available from: http://dx.doi.org/10.1016/j.jnucmat.2018.09.026.
- [4] Chevalier S. What did we learn on the reactive element effect in chromia scale since Pfeil's patent? Materials and Corrosion. 2013 Dec;65(2):109–115. Available from: http://dx.doi.org/10.1002/maco.201307310.
- [5] Hindam H, Whittle D. PEG FORMATION BY SHORT CIRCUIT DIF-FUSION IN A1203 SCALES CONTAINING OXIDE DISPERSIONS. Journal of The Electrochemical Society. 1982;129(1147).
- [6] Eklund J, Jönsson B, Persdotter A, Liske J, Svensson JE, Jonsson T. The influence of silicon on the corrosion properties of FeCrAl model alloys in oxidizing environments at 600 °C. Corrosion Science. 2018 Nov;144:266–276. Available from: http://dx.doi.org/10.1016/j.corsci.2018.09.004.
- [7] Pint BA. Optimization of Reactive-Element Additions to Improve Oxidation Performance of Alumina-Forming Alloys. Journal of the American Ceramic Society. 2003 Apr;86(4):686–95. Available from: http://dx.doi.org/10.1111/j.1151-2916.2003.tb03358.x.
- [8] Quadakkers WJ, Holzbrecher H, Briefs KG, Beske H. Differences in growth mechanisms of oxide scales formed on ODS and conventional wrought alloys. Oxidation of Metals. 1989 Aug;32(1–2):67–88. Available from: http://dx.doi.org/10.1007/BF00665269.
- [9] Krakauer BW, Seidman DN. Absolute atomic-scale measurements of the Gibbsian interfacial excess of solute at internal interfaces. Physical Review B. 1993 Sep;48(9):6724–6727. Available from: http://dx.doi.org/10.1103/PhysRevB.48.6724.
- [10] Boll T, Unocic KA, Pint BA, Stiller K. Interfaces in Oxides Formed on NiAlCr Doped with Y, Hf, Ti, and B. Microscopy and Microanalysis. 2017 Mar;23(2):396–403. Available from: http://dx.doi.org/10.1017/S1431927617000186.

- [11] Nychka JA, Clarke DR. Quantification of Aluminum Outward Diffusion During Oxidation of FeCrAl Alloys. Oxidation of Metals. 2005 Jun;63(5–6):325–352. Available from: http://dx.doi.org/10.1007/s11085-005-4391-4.
- [12] Boll T, Unocic KA, Pint BA, Mårtensson A, Stiller K. Grain Boundary Chemistry and Transport Through Alumina Scales on NiAl Alloys. Oxidation of Metals. 2017 Jan;88(3–4):469–479. Available from: http://dx.doi.org/10.1007/s11085-016-9697-x.
- [13] Barth TL, Weber PK, Liu T, Xue F, Valenza TC, Backman L, et al. Grain boundary transport through thermally grown alumina scales on NiAl. Corrosion Science. 2022 Dec;209:110798. Available from: http://dx.doi.org/10.1016/j.corsci.2022.110798.
- [14] Babic V, Geers C, Panas I. Reactive Element Effects in High-Temperature Alloys Disentangled. Oxidation of Metals. 2019 Nov;93(1–2):229–245. Available from: http://dx.doi.org/10.1007/s11085-019-09946-6.
- [15] Pieraggi B. Fundamental Aspects of Reactions at the Metal/Scale Interface During Scaling. Materials Science Forum. 1997 Oct;251–254:299–312. Available from: http://dx.doi.org/10.4028/www.scientific.net/MSF.251-254.299.
- [16] Pint BA. Experimental observations in support of the dynamicsegregation theory to explain the reactive-element effect. Oxidation of Metals. 1996 Feb;45(1–2):1–37. Available from: http://dx.doi.org/10.1007/BF01046818.
- [17] Yokoi T, Yoshiya M, Yasuda H. On modeling of grain boundary segregation in aliovalent cation doped ZrO2: Critical factors in site-selective point defect occupancy. Scripta Materialia. 2015 Jun;102:91–94. Available from: http://dx.doi.org/10.1016/j.scriptamat.2015.02.021.
- [18] Nakagawa T, Sakaguchi I, Shibata N, Matsunaga K, Mizoguchi T, Yamamoto T, et al. Yttrium doping effect on oxygen grain boundary diffusion in -Al2O3. Acta Materialia. 2007 Nov;55(19):6627–6633. Available from: http://dx.doi.org/10.1016/j.actamat.2007.08.016.
- [19] Gemming T, Nufer S, Kurtz W, Rühle M. Structure and Chemistry of Symmetrical Tilt Grain Boundaries in -Al2O3: II, Bicrystals with Y at the Interface. Journal of the American Ceramic Society. 2003 Apr;86(4):590–94. Available from: http://dx.doi.org/10.1111/j.1151-2916.2003.tb03345.x.
- [20] WADA M, MATSUDAIRA T, KITAOKA S. Mutual grain-boundary transport of aluminum and oxygen in polycrystalline Al2O3 under oxygen potential gradients at high temperatures. Journal of the Ceramic Society of Japan. 2011;119(1395):832–839. Available from: http://dx.doi.org/10.2109/jcersj2.119.832.
- [21] Matsudaira T, Wada M, Saitoh T, Kitaoka S. Oxygen permeability in cation-doped polycrystalline alumina under oxygen potential gradients at high temperatures. Acta Materialia. 2011 Aug;59(14):5440–5450. Available from: http://dx.doi.org/10.1016/j.actamat.2011.05.018.
- [22] Matsudaira T, Wada M, Saitoh T, Kitaoka S. The effect of lutetium dopant on oxygen permeability of alumina polycrystals under oxygen potential gradients at ultra-high temperatures. Acta Materialia. 2010 Mar;58(5):1544–1553. Available from: http://dx.doi.org/10.1016/j.actamat.2009.10.062.
- [23] Heuer AH, Zahiri Azar M, Guhl H, Foulkes M, Gleeson B, Nakagawa T, et al. The Band Structure of Polycrystalline Al2O3 and Its Influence on Transport Phenomena. Journal of the American Ceramic Society. 2016 Feb;99(3):733–747. Available from: http://dx.doi.org/10.1111/jace.14149.
- [24] Smialek JL. Invited Review Paper in Commemoration of Over 50 Years of Oxidation of Metals: Alumina Scale Adhesion Mechanisms: A Retrospective Assessment. Oxidation of Metals. 2022 Jan;97(1–2):1–50. Available from: http://dx.doi.org/10.1007/s11085-021-10091-2.
- [25] Grigoriev S, Vereschaka A, Uglov V, Milovich F, Tabakov V, Cherenda N, et al. Influence of the tribological properties of the Zr,Hf-(Zr,Hf)N-(Zr,Me,Hf,Al)N coatings (where Me is Mo, Ti, or Cr) with a nanostructured wear-resistant layer on their wear pattern during turning of steel. Wear. 2023 Apr;518–519:204624. Available from: http://dx.doi.org/10.1016/j.wear.2023.204624.

PŘED VÍCE NEŽ 140 LETY VYNALEZL ZAKLADATEL ALEXANDER GRAHAM BELL TELEFON. AT&T MÁ OD TÉ DOBY DĚDICTVÍ OSMI NOBELOVÝCH CEN A VÍCE NEŽ 12 500 PATENTŮ PO CELÉM SVĚTĚ.

NAŠE JEDNOTKA, AT&T LABS NETWORK SYSTEMS, JE ZAMĚŘENA NA TRANSFORMACI PROSTORU SÍTĚ OSS POMOCÍ NOVÉ GENERACE TECHNOLOGIÍ AI, ANALÝZOU VELKÝCH DAT A DALŠÍCH DISTRIBUOVANÝCH SYSTÉMŮ.

PRO NÁŠ TÝM V BRNĚ HLEDÁME ZKUŠENÉ I ČERSTVÉ ABSOLVENTY SOFTWAROVÉHO INŽENÝRSTVÍ, KTEŘÍ CHTĚJÍ BÝT SOUČÁSTÍ DYNAMICKÉHO ROZVOJE V OBLASTI 5G/6G A SDN. STUDENTŮM IT NABÍZÍME INTERNSHIPS, ČÁSTEČNÉ PRACOVNÍ ÚVAZKY, TÉMATA A VEDENÍ BAKALÁŘSKÝCH A MAGISTERSKÝCH PRACÍ.



KONTAKT: KATEŘINA VELKÁ KATERINA.VELKA@INTL.ATT.COM AT&T GLOBAL NETWORK SERVICES CZECH REPUBLIC PALACHOVO NÁMĚSTÍ 2 625 00 BRNO



An enhanced theoretical approach for accurate measurements of the optical and energy characteristics of semiconducting materials

1st Mohammad M. Allaham Central European Institute of Technology Brno University of Technology Institute of Scientific Instruments of Czech Academy of Sciences Brno, the Czech Republic allaham@vutbr.cz 2nd Zuzana Košelová Department of Microelectronics Faculty of Electrical Engineering and Communication Brno University of Technology Institute of Scientific Instruments of Czech Academy of Sciences Brno, the Czech Republic koselova@vutbr.cz 3rd Dinara Sobola

Central European Institute of Technology Brno University of Technology Department of Physics Faculty of Electrical Engineering and Communication Brno University of Technology Brno, the Czech Republic Sobola@vut.cz

4th Zdenka Fohlerová Department of Microelectronics Faculty of Electrical Engineering and communication Brno University of Technology Brno, the Czech Republic fohlerova@vutbr.cz

Abstract—The absorption coefficient is an important optical property in characterizing semiconducting materials. It plays a significant role in studying optical characteristics, electrical structure, energy band structure, and the creation/annihilation of excitons when a semiconducting material absorbs electromagnetic radiation. In this study, an enhanced theoretical model will be introduced and applied to characterize thin films prepared from UPR4 (unsaturated polyester resin) single-component epoxy resin, which is important to study the charge flow at the interface of tungsten-UPR4 composite field emission cathodes.

Index Terms—absorption coefficient, energy gap, Urbach tailing energy, Tauc plot, transmittance, and reflectance.

I. INTRODUCTION

When a transparent semiconducting film is subjected to electromagnetic radiation of sufficient photon energy, valance band electrons absorb the incident photons and excite to the conduction band. This process leaves empty holes in the valance band structure, and the electron-hole pair is known as the exciton. Moreover, some photons will be scattered too along the film thickness. The attenuation coefficient describes the two events (absorption and scattering) as a linear combination of their coefficients, the absorption coefficient (α) and the scattering coefficient (σ) where $\mu = \alpha + \sigma$. In the case where the scattering of photons is neglected ($\sigma \rightarrow 0$), the attenuation coefficient is approximated to the absorption coefficient ($\mu \approx \alpha$) [1]–[4].

the absorption coefficient is a significant optical property of a semiconducting material that helps to study its energy band structure and the creation/annihilation of excitons. In literature, the Bouguer-Beer-Lambart law (BBL-law) is widely used to express the absorption coefficient (α) as a function of the film thickness (x) and the relative transmittance ($T(\lambda)$) of electromagnetic radiation. The BBL-law is given by the following equation [5]:

5th Alexandr Knápek Institute of Scientific Instruments of Czech Academy of Sciences

Brno, the Czech Republic

knapek@isibrno.cz

$$\alpha = \frac{1}{x} \ln\left(\frac{1}{T}\right) \tag{1}$$

BBL-law is considered as a first good approximation for obtaining α because it does not consider the reflectance of the photons at the first air-film interface or even the multi-internal reflections that can occur at the second film-air interface as shown in Fig.1(a).

Considering the case of a single reflection at the first airfilm interface, the BBL-law will take the following form after considering the reflectance (R) of the photons (as presented in Fig.1(b)):

$$\alpha = \frac{1}{x} \ln\left(\frac{1-R}{T}\right) \tag{2}$$

Furthermore, if only one internal reflection was considered at the second film-air interface, the transmitted radiation will be reduced again (as presented in Fig.1(c)), and the absorption coefficient will be given by [6]:

$$\alpha = \frac{1}{x} \ln\left(\frac{(1-R)^2}{T}\right) \tag{3}$$

In this paper, we derive an equation to measure α by considering infinite internal reflections at the two interfaces, which



Fig. 1. Different transmission regimes through a semiconducting film of thickness x. (a) Describes the transmission without reflection, (b) Describes the case with a single reflection at the air-film interface, (c) describes the case with a single reflection at each interface, and (d) describes the case of infinite reflections at each interface.

causes an infinite absorption and scattering of photons while the film is subjected to electromagnetic radiation. Moreover, thin films of the UPR4 single-component epoxy resin were prepared for testing the new methodology and comparing it with the two main forms of BBL-law.

II. METHODOLOGY

A. Theoretical interpretation

To generalize equation 3, we consider the case of infinite internal reflections (as presented in Fig.1(d)). In this case, the

total transmitted radiation is given by:

$$T = (1 - R)^2 \sum_{n=1}^{\infty} R^{2(n-1)} \exp(-\alpha x (2n - 1))$$
 (4)

Solving the infinite series in equation 4 connects T, R, and α and provides the following equation [6]:

$$T = \frac{(1-R)^2 \exp(-\alpha x)}{1 - R^2 \exp(-2\alpha x)}$$
(5)

In equation 5, applying the definition $y = \exp(-\alpha x)$, multiplying both sides by $1 - R^2 y^2$, and following algebraic

arrangements yields to:

$$TR^2y^2 + (1-R)^2y - T = 0$$
 (6)

Equation 6 is a quadratic equation, and it has the solution:

$$y = \frac{-(1-R)^2 \pm \sqrt{(1-R)^4 + 4T^2R^2}}{2TR^2}$$
(7)

The solution with the negative square root in equation 7 must be neglected since it provides a non-defined solution for α . Thus, considering only the solution with the positive square root, substituting back the value of y, applying the natural logarithm, and dividing both sides by -x yields the following form for α :

$$\alpha = \frac{1}{x} \ln \left(\frac{2TR^2}{\sqrt{(1-R)^4 + 4T^2R^2} - (1-R)^2} \right)$$
(8)

In this paper, equation 8 is referred to as the ASKequation (Allaham-Sobola-Knápek equation) for measuring α by considering infinite internal reflections from transparent semiconducting films. The obtained results of α from the ASK-equation can then be used in the Tauc methodology to study the energy gap of the tested films. Moreover, it can be used in Urbach methodology to study the defects in the structure at low density of states. The following equations are used to determine the energy gap (E_g) and the Urbach tailing energy (E_U) for the tested films [1], [7]:

$$(\alpha h\nu)^n = A(h\nu - E_g) \tag{9}$$

$$\ln(\alpha) = \left[\ln(\alpha_0) - \frac{E_g}{E_U}\right] + \frac{h\nu}{E_U}$$
(10)

In equation 9, A is a proportional parameter, n = 2 is related to the allowed direct transitions, and n = 1/2 is related to the allowed indirect transitions. Moreover, $E_{\rm U} = 1/{\rm slope}$ of equation 10.

B. Materials and experimental setup

To prepare the UPR4 flat thin films, a glass substrate was coated by the epoxy resin and replaced inside the furnace at 180 °C for 4 hours for the epoxy resin to cure. The samples were titled at 45° to prepare flat epoxy films preventing the shrinking of the epoxy layer. The prepared films have an average thickness of 0.1 mm, which was measured by KINEX micrometer from K-MET (KINEX measuring) (Prague, the Czech Republic).

Raw T and R measurements were obtained from the Ocean Optics JAZ-3 and NIR-QUEST (Orlando, USA) spectrometer, where the tested wavelength range was 200 - 1000 nm.

III. RESULTS AND DISCUSSION

The raw T and R spectra are presented in Fig.2(a and b). The obtained results were used to obtain the absorption coefficient from equations 1, 3, and 8 (The results are presented in Fig.2(c)). The allowed direct energy gap values were obtained from the three curves in Fig.2(d), where the results report allowed direct energy gap values of $E_{\rm g}^{\rm dir} = 3.240$ eV as measured by equation 1, while it was $E_{\rm g}^{\rm dir}=3.245~{\rm eV}$ as measured from equations 3 and 8.

Moreover, the curves in Fig.2(e) were used to calculate the allowed indirect energy gap values, where it was found to have values of $E_{\rm g}^{\rm dir} = 3.199$ eV when using equation 1, but $E_{\rm g}^{\rm dir} = 3.174$ eV when using equations 3 and 8. The UPR4 epoxy resin has an amorphous structure and is

The UPR4 epoxy resin has an amorphous structure and is cured thermally, which causes some defects in its energy band structure at a very low density of states. These defects appear as delocalized energy states located inside the material's energy gap. The difference between the modified top of the valance band and the bottom of the conduction band at a very low density of states is known as the Urbach band energy tailing.

The Urbach curves are presented in Fig.2(f) for different α values as obtained from equations 1, 3, and 8. The slopes of the 2 curves are 3.484 eV⁻¹ (for the curve that was obtained from equation 1) and 3.559 eV⁻¹ (for the curves that were obtained from equations 3 and 8). These slope values correspond to Urbach tailing energies of 0.287 eV and 0.281 eV respectively. The summary of the obtained results is presented in Table I.

 TABLE I

 MEASUREMENTS OF THE ENERGY GAP AND URBACH TAILING ENERGIES

 FOR α VALUES AS OBTAINED FROM EQUATIONS 1, 3, AND 8

Equation	$E_{\rm g}^{\rm dir}$ (eV)	$E_{\rm g}^{\rm indir}$ (eV)	$E_{\rm U}~({\rm eV})$
1	3.240	3.174	0.287
3 and 10	3.245	3.199	0.281
ΔE	0.005	0.025	-0.006

The allowed direct and indirect transitions, along with the energy band structure in the momentum space are presented in Fig.2(g). The allowed direct transition $(a \rightarrow b)$ is related to a transition from the top of the valance band to the bottom of the conduction band without any changes in the linear momentum of the electron.

The allowed indirect transition $(a \rightarrow c \rightarrow d)$ is related to a phonon-assisted transition. In this case, the energy of the photon is not enough to accomplish the direct transition. However, a phonon-electron interaction occurs causing a change in the linear momentum of the transmitted electron is changed and the transition occurs at nonadjacent limits of the energy bands.

The schematic in Fig.2(h) describes the energy bands in the density of states space. In the case of the UPR4 epoxy resin, it shows high defects of in the structure of the energy bands and energy gap, which is related to the amorphous structure of the epoxy resin.

IV. CONCLUSION

In this paper, a new formula for measuring the absorption coefficient from raw transmittance and reflectance data was introduced considering an approximation for the cases where the scattering coefficient can be neglected equation 8.

A practical application for equation 8 was introduced by applying it to the transmittance and reflectance of raw data as obtained from thin films of the UPR4 epoxy resin. The



Fig. 2. The optical characteristics of the UPR4 single component epoxy resin. (a) Transmittance (T) spectrum, (b) Absolute reflectance (R) spectrum. (c) Absorption coefficient. (d) Tauc plot for evaluating the allowed direct optical energy gap. (e) Tauc plot for evaluating the allowed indirect optical energy gap. (f) Urbach plot to evaluate the Urbach band tailing energy. (g) Allowed direct and indirect electric transitions and (h) energy bands tailing in the UPR4 single component epoxy resin

optical characteristics of the prepared thin films reported that the UPR4 epoxy resin has values of 3.245 eV for the allowed direct transition and 3.199 eV for the allowed indirect transition. Moreover, the energy band tailing energy was found to have a value of 0.287 eV.

Moreover, the results of the ASK-equation were compared to values that were obtained from the BBL-law, which showed high consistency between the two methodologies as presented in Table I.

ACKNOWLEDGMENT

CzechNanoLab project LM2018110 funded by MEYS CR is gratefully acknowledged for the financial support of the measurements/sample fabrication at CEITEC Nano Research Infrastructure.

The research described in the paper was financially supported by the Internal Grant Agency of the Brno University of Technology, grant numbers CEITEC VUT/FEKT-J-24-8567.

The research was financially supported by the Czech Academy of Sciences (RVO:68081731).

REFERENCES

- [1] M. Allaham, R. Dallaev, D. Burda, D. Sobola, A. Nebojsa, A. Knápek, M. S. Mousa, V. Kolařík, Energy gap measurements based on enhanced absorption coefficient calculation from transmittance and reflectance raw data, Physica Scripta 99 (2) (2024) 025952. doi:10.1088/ 1402-4896/adlcb8.
- [2] J. Tauc, Optical properties and electronic structure of amorphous ge and si, Materials Research Bulletin 3 (1) (1968) 37-46. doi:https:// doi.org/10.1016/0025-5408(68)90023-8.
- [3] S. Yang, Y. Zhang, Structural, optical and magnetic properties of mndoped zno thin films prepared by sol-gel method, Journal of Magnetism and Magnetic Materials 334 (2013) 52–58. doi:https://doi.org/ 10.1016/j.jmmm.2013.01.026.
- [4] M. Hohmann, B. Lengenfelder, D. Muhr, M. Späth, M. Hauptkorn, F. Klämpfl, M. Schmidt, Direct measurement of the scattering coefficient, Biomed. Opt. Express 12 (1) (2021) 320–335. doi:10.1364/BOE. 410248.
- [5] T. G. Mayerhöfer, S. Pahlow, J. Popp, The bouguer-beer-lambert law: Shining light on the obscure, ChemPhysChem 21 (18) (2020) 2029–2046. doi:https://doi.org/10.1002/cphc.202000464.
- [6] J. I. Pankove, Optical processes in semiconductors, Dover publications, 1975.
- [7] O. V. Rambadey, A. Kumar, A. Sati, P. R. Sagdeo, Exploring the interrelation between urbach energy and dielectric constant in hf-substituted batio3, ACS Omega 6 (47) (2021) 32231–32238. doi:10.1021/ acsomega.lc05057.

A Review of the Li-ion Battery in-situ Experiments in Scanning Electron Microscope

David Trochta Dept. of Electrical and Electronic Technology Dept. of Electrical and Electronic Technology Dept. of Electrical and Electronic Technology Brno University of Technology Brno, Czech Republic 199998@vut.cz

Ondřej Klvač Brno University of Technology Brno, Czech Republic xklvac02@vut.cz

Tomáš Kazda Brno University of Technology Brno, Czech Republic kazda@vut.cz

Abstract-This paper focuses on the description of in-situ experiments with lithium-ion batteries in a scanning electron microscope during cycling. With these experiments, we are able to better understand the working principles and internal processes of Li-ion batteries and improve their efficiency, reliability, and safety based on this knowledge. However, conducting these experiments poses several challenges, which are described in the paper, along with ways in which they can be partially mitigated.

Keywords—Li-ion, battery, in-situ, SEM, real-time observation

INTRODUCTION I.

Lithium-ion batteries, due to their high energy density and long cycle life, have become an essential component of modern technology, powering a wide range of devices, from smartphones to electric vehicles. Nevertheless, despite their widespread use, there are still challenges related to their efficiency, safety, and lifespan that need to be addressed [1], [2]. However, it is not only about improving currently known materials and optimizing them, but also about introducing new so-called advanced Li-ion batteries. These batteries are experiencing significant advancements in materials and design. Solid-state batteries, employing a solid electrolyte, promise improved safety and energy density [3]. High-energy cathodes, including nickel-rich and lithium-rich variants, enhance overall performance [4]. Silicon anodes are being explored to boost energy density, albeit with challenges like volume changes [5]. To further optimize and improve all these Li-ion systems, it is crucial to understand the internal processes and mechanisms of these batteries, which requires detailed characterization [6].

This is where scanning electron microscopy (SEM) can be used among many other characterization techniques. SEM, which uses a focused beam of electrons to create detailed images of materials at a microscopic level, provides valuable insights into the structure and composition of the battery materials [7]. Taking this a step further, in-situ SEM allows researchers to observe and study these batteries under operating conditions, providing real-time information about the dynamic processes occurring within the batteries. This has opened new opportunities for improving the performance and safety of lithium-ion batteries [8], [9].

Although SEM has proven to be an indispensable technique for characterization, there are several problems. Achieving high resolution during battery operation can be limited, and the introduction of experimental conditions may lead to sample artifacts that affect accuracy. The field of view is often constrained, making it challenging to capture a comprehensive picture of the entire battery system. Continuous exposure to the electron beam can result in sample damage, impacting the observed properties. Additionally, the complexity of in-situ experiments, the necessity for specialized equipment, and the requirement to closely mimic actual operating conditions pose practical challenges [8], [9]. Also, the interpretation of the results demands a deep understanding of both electrochemistry and electron microscopy.

This paper explores advancements, both recent and foundational, highlighting emerging technologies that show promise in overcoming in-situ SEM-related limitations. By examining the spectrum of progress in the field, from recent innovations to longstanding methodologies, we aim to provide a comprehensive perspective on the in-situ analysis of Li-ion batteries in SEM. Through this complex examination, we endeavor to contribute significantly to the ongoing discourse surrounding battery characterization, fostering dialogue and innovation aimed at enhancing efficiency, safety, and longevity.

II. ADVANTAGES OF THE IN-SITU TECHNIQUES

In the field of lithium-ion battery research, both ex-situ and in-situ techniques play crucial roles in unraveling the mysteries of battery behavior. Ex-situ techniques involve the removal of battery components for subsequent analysis, providing valuable insights into the composition, structure, and chemical properties of the materials [10]. While ex-situ methods have been instrumental in advancing our understanding of lithium-ion batteries, they lack the ability to capture dynamic processes occurring during actual battery operation [11]. Another problem is that sample preparation for ex-situ analysis usually must take place in a box with a protective atmosphere. This is because the batteries contain materials that are sensitive to humidity, oxygen, or nitrogen. During sample preparation (battery disassembly), samples may also be damaged or degraded due to improper handling or use of excessive force [12].

These limitations are where in-situ techniques shine and offer unique advantages that significantly enhance our understanding of Li-ion batteries. SEM provides exceptional temporal resolution, enabling researchers to precisely capture rapid electrochemical changes at the micro or nanoscale. This high-resolution imaging is crucial for real-time observations of dynamic processes, such as the lithiation and delithiation of electrodes, offering direct insights into the evolution of electrode materials and the formation of solid-electrolyte interfaces (SEI) [9], [13]. Additionally, SEM allows for the direct visualization of interfaces, providing detailed information on the structural and morphological changes critical for optimizing electrode materials and electrolytes. The application of in-situ SEM, with its ability to provide real-time, highresolution imaging, empowers researchers to tailor battery designs and materials for improved efficiency, thereby addressing key challenges in the quest for advanced lithium-ion battery technologies [9], [8].

III. SEM INSTRUMENTATION FOR IN-SITU ANALYSIS

An electron microscope consists of several important parts. At its core, a SEM instrument consists of an electron gun that emits a focused beam generated via thermionic emission from a heated filament (W or LaB₆) or a field emission gun (FEG). The electron beam is then focused into a probe on the surface of the sample using an electromagnetic lens [7], [8]. The use of the FEG electron source proves more suitable for in-situ battery analyses due to its high brightness, coherence, and reduced chromatic aberration [7]. The improved stability and smaller probe size contribute to superior spatial resolution and detailed observations of dynamic processes during battery cycling [8].

Another essential part is the detectors. In addition to the classical Everhart-Thornley detector (ETD), it is advisable to use so-called in-lense detectors. These detectors are able to detect more electrons, but also electrons with lower energy, and thus improve the overall image obtained. Other necessary detectors include an EDS detector for energy dispersive spectroscopy and determining the exact material composition together with a quantitative data [7]. However, conventional EDS detectors are not capable of detecting materials with a proton number lower than 4. Therefore, if EDS analysis of lithium with a proton number of 3 is required, a special EDS detector, referred to as windowless, is needed [14]. Yet another detector that can be used for in-situ analysis of batteries is electron back scattered diffraction detector (EBSD). With this detector, phase changes and orientation changes in the material structure can be observed [8].

Other useful tools that can be used in SEM for in-situ analyses include various micromanipulators, which can be used to assemble an electrochemical cell directly in the microscope chamber or, in the case of suitable shielding and wiring, also as current and voltage probes [15]. There can also be various electrochemical probes to connect the battery to the petentiostat. For the investigation of mechanical properties, it is possible to use nanointenders [16]. Other useful tools and instrumentation for in-situ battery research include, e.g., heating or cooling stages [17], gas injection systems [18], or systems for sample transfer in a protective atmosphere [8]. All of the above equipment can be purchased from vendors and easily configured for the purpose of the experiment. However, special holders and electrochemical cells are also needed for in-situ battery analysis, which are not commonly sold, and scientists are usually left to construct them on their own. Configurations of these electrochemical cells can be divided into two groups: open cells and liquid cells (closed cells). [9]

The open cell configuration is easier to manufacture and assemble but has the disadvantage that the entire battery is exposed to the environmental conditions of the microscope. If a deep vacuum is used for imaging, the use of this cell precludes the use of conventional electrolytes, i.e., lithium salts in organic solvents, which would evaporate in the chamber. Thus, this open configuration is suitable for batteries with a solid-state electrolyte or batteries with an ionic liquid electrolyte that does not evaporate in a vacuum [9]. The main advantage of this configuration is the high resolution of the images. However, the high resolution decreases with the use of low-vac mode, or environmental scanning electron microscope (ESEM), which can be used to reduce sample charging or to get closer to the battery's real operational conditions [19], [20].

On the other hand, the completely closed and sealed design of the liquid cells allows the use of ether- and ester-based electrolytes. This can be used, e.g., to observe the formation and evolution of the SEI layer in these electrolytes or for in-situ analysis of the battery under the real conditions of the commonly used electrolytes. The observation itself is then performed through a silicon nitride (SiN) observation window. The main disadvantage of this configuration then lies in this observation window, which significantly reduces the spatial resolution [9]. This configuration is also more challenging to assemble accurately, and the electrochemical cell cannot be modified after sealing, which is done with a resin epoxy. Another practical problem lies in the electrolyte layer that can form just below the observation window, making it impossible to image the electrode structure itself and due to closed construction, the excess electrolyte cannot be removed [21].

IV. LI-ION IN-SITU EXPERIMENTS

The first in-situ observation of a Li-ion battery was conducted by Braudy and Armand [22] in 1987. The authors delineated the fundamental principles and procedures for sample preparation for in-situ observations. In their experiment, lithium metal served as the anode, and the cathode explored various materials, including titanium dioxide (TiO₂), vanadium oxide (V₆O₁₃), and iron(II) sulfide (FeS). A polymeric electrolyte, consisting of a polyether with an ethylene oxide base and lithium perchlorate (LiClO₄) salt, functioned as the separator. Assembling the battery in a drybox using the hot-pressing method, the authors then transferred it into the electron microscope chamber. The sample orientation enabled the sideview observation of the sandwich structure. The research yielded significant results, particularly in observing the phase changes of TiO₂ and V₆O₁₃ during cycling. It was observed that the TiO₂ structure remained remarkably stable, whereas the V₆O₁₃ structure exhibited cracking. The cathode containing FeS underwent significant changes during the cycling, resulting in a breakdown of its structure.

Some research teams have tried to solve the problem of electrolyte evaporation using the ESEM. One notable experiment from 2006, conducted by Rainmann et al. [23] in ESEM with Ar and pressure of 200 Pa. Authors aimed to observe volume changes and potential mechanical damage on the anode of a Li-ion battery. The anode was a mixture of Sn, Super P, and polyvinylidene fluoride-hexafluoropropylene (PVDF-HFP) as a binder. The whole mixture was coated on a stainless-steel grid to achieve better wetting and also to be able to observe the active material of the anode. Lithium metal served as the counter electrode. Between the electrodes was an unspecified type of separator soaked with electrolyte. Unfortunately, the exact electrolyte details were not specified. However, ethylene carbonate (EC) and propylene carbonate (PC) were mentioned as the used solvents due to their higher boiling points. The researchers designed a custom electrochemical cell made of polypropylene, which can be seen in Figure 1. The main feature is that the sample placed in the holder is covered with a mylar film, which has a very small hole for the passage of electrons. In this way, they combined the advantages of liquid and closed cells. They used electrolyte with common solvents and were also able to maintain high spatial resolution due to the small observing window. Another feature is the system for keeping any vapor from the electrolyte inside the cell from escaping through the viewing hole. Finally, the authors were able to observe high-resolution changes in the anode structure, proving the functionality of their cell design.



Fig. 1. Elchem. cell for ESEM in-situ analysis by Rainamann et al. [23]

Different advance occurred in 2011 when Chen et al. [20] published an in-situ experiment involving a Li-ion battery and ionic liquid electrolyte. The authors used lithium metal as the anode, while the cathode consisted of a mixture of SnO₂, carbon black, and polyvinylidene fluoride (PVDF) coated on a steel grid. The steel grid served as a current collector, making the electrode permeable. To assemble their battery, the authors employed a conventional sample holder used in SEM called a stub. In a protective glove box with an argon atmosphere, they affixed a Cu foil to this aluminum stub to ensure material compatibility with the lithium metal placed on the stub. Subsequently, a Whatman glass separator, filled with ionic liquid, was placed on top of the lithium metal. The prepared sample was then transferred to an electron microscope chamber in a protective atmosphere. In the microscope, the cathode was positioned on the separator using a micromanipulator. The experimental setup can be seen in Figure 2. The advantage of this procedure was the ability to observe the wetting process of the cathode. However, a drawback was that the battery examination only provided a top view, allowing observation solely of the cathode. The experiment revealed irreversible changes on the surface of the cathode during discharge, with no observable changes during charging.



Fig. 2. In-situ experiment setup by Chen et al. [20]

In 2019 Shi at al. [19] introduced interesting experiment. This time, the authors decided to build a half cell and used lithium metal as the counter electrode. A composite of Si, graphite, the conductive additive C65, and a carboxymethyl cellulose (CMC) binder coated on copper foil current collector was then used as the anode. The electrodes were separated by a Whatman glass separator that was soaked in ionic liquid electrolyte (10 wt.% bis(trifluoromethane) sulfonimide lithium salt (LiTFSI) in 1-ethyl-3-methylimidazoliumbis (trifluoroomethylsulfonyl) imide (EMIM TFSI). The aim of the work was to compare two types of working electrodes, namely unstructured S/C composite and 3D-line structured Si/C composite, into which channels were made with a laser. Assembling the cells took place in a glove box with argon atmosphere, and the arrangement was placed in the microscope chamber for side-view observations as shown in Figure 3. The results showed that the 3D-structured S/C electrode is more suitable for several reasons, among the main being better mechanical resistance to prevent separation from the current collector. Moreover, the channels enlarged the contact area of the electrodes and improved the utilization of anode materials. This greatly increased the charge capacity of 3D-line-structured anodes.



Fig. 3. In-situ experiment setup and imaged cross-view by Shi et al. [19]

Another paper from 2019 was published by Tsuda et al. [24]. Authors introduced another in-situ experiment with a different open-cell configuration. The experimental setup featured a glass plate as the foundation, upon which the cathode, lithium cobalt oxide (LiCoO₂), on an aluminum current collector was positioned. Two Whatman glass separators were then layered on the cathode, between which a Ni wire was used as a reference electrode. The cathode, composed of Si particles deposited on a copper grid through electrophoretic deposition (EPD), was subsequently added to the top of the prepared sample. Whole battery structure can be observed in Figure 4. An ionic liquid, [C2mim] [FSA] with 1.0 M Li [TFSA] and [Li(G4)] [TFSA], served as the electrolyte. This innovative three electrode design allowed for unique observations and insights into the behavior of the Li-ion battery components during cycling, even though only the anode could be observed from the top view. The authors were able to observe changes in the morphology and phase changes of Si nanoparticles, which they were able to contrast with the discharge and charge curves. They also presented the strength of the three-electrode in-situ battery measurement technique.



Fig. 4. Exploded view of the sample prepared by Tsuda et al. [24]

In 2020 Kaboli et al. [25] published a study in which they focused on detecting the cause of solid-state Li-ion battery failure. The electrolyte in this case was a solid polymer electrolyte (SPE), which consisted of polyethylene oxide (PEO) and LiTFSI in a molar ratio of 30:1. Lithium nickel manganese cobalt oxides (NMC 622) was used as the cathode, and lithium metal was used as the counter electrode. From the abovementioned components, an electrochemical cell was assembled in a dry room, as indicated in Figure 5. The sample was then enclosed in nonconductive resin. The publication also included an experiment to exclude reactions between the used resin and the electrochemical component at room temperature and at 50 °C. Their results showed that no chemical reactions occurred between the sample and the resin enclosure. After the resin had cured, the sample was transferred to a cryo-microtomy machine, where the cross section was prepared. During the preparation of the cross section and polishing, care was taken to ensure that the surface was not contaminated, and therefore it was constantly flushed with argon. The sample thus prepared was then placed in the chamber of the electron microscope. The sample then had to be fixed in the chamber using nonconducting plastic plates, which also set the pressure on the measured cell, which is important for batteries with SPE. The sample was then heated to 50 °C, and after 24 hours of tempering, electrochemical measurements began.



Fig. 5. Schematic view of the sample prepared by Kaboli et al. [25]

During a cycling, the authors were able to discover the main cause of the failures that occur with these types of batteries. The main cause is the gradual thinning of the electrolyte, which has thinned from the original 23 um to 5 um during the cycling. The authors also processed the same data for the NMC cathode and lithium metal. While NMC changed its thickness periodically depending on charging and discharging, the thickness of metallic lithium was almost constant during the whole cycling. [25]

V. CONCLUSION

In summary, lithium-ion batteries stand as indispensable components in modern technology, owing to their high energy density and prolonged cycle life. However, persistent challenges in efficiency, safety, and lifespan require ongoing exploration and innovation. The evolution toward Li-ion and advanced Li-ion batteries brings both promise and difficulty, requiring a complex understanding of internal mechanisms for meaningful progress.

This paper has delved into the critical role of in-situ scanning electron microscopy (SEM) in unraveling the complexity of lithium-ion batteries. Despite facing challenges like spatial resolution limitations and sample artifacts, in-situ SEM provides invaluable real-time insights into dynamic battery processes during operation. The discussion extends beyond challenges to highlight recent and foundational advancements, emphasizing emerging technologies that exhibit potential for overcoming these limitations.

Moreover, the suitability of individual SEM components for in-situ battery analyses has been evaluated, shedding light on the significance of additional equipment tailored for in-situ SEM analysis. A focal point is the electrochemical cell for in-situ SEM analysis and its configurations, with the open configuration emerging as the most common choice for its versatility in examining various materials with high spatial resolution. However, it's crucial to acknowledge the limitations of this configuration, particularly its inability to accommodate organic solvent-based electrolytes. For this reason, batteries with an ionic liquid or solid-state electrolyte are most often investigated in this type of cell.

This review also confirms that there is a continued need to actively discuss the improvement of these in-situ techniques in SEM and focus on the development of additional equipment for these analyses that could make them easier and more efficient. It is also necessary to develop new tools that can enable a whole new branch of experiments that can reveal more about the internal process of not only Li-ion batteries but also post-Li-ion batteries.

ACKNOWLEDGEMENT

This work was supported by the specific graduate research of the Brno University of Technology No. FEKT-S-23-8286.

REFERENCES

- H. Cho, J. Kim, M. Kim, H. An, K. Min and K. Park, "A review of problems and solutions in Ni-rich cathode-based Li-ion batteries from two research aspects: Experimental studies and computational insights", *Journal of Power Sources*, vol. 597, 2024.
- [2] N. Nitta, F. Wu, J. Lee and G. Yushin, "Li-ion battery materials: present and future", *Materials Today*, vol. 18, no. 5, pp. 252-264, 2015.
- [3] W. Ji, B. Luo, G. Yu, Q. Wang, Z. Zhang, Y. Tian, Z. Liu, W. Ji, Y. Nong, X. Wang and J. Zhang, "A review of challenges and issues concerning interfaces for garnet-type all-solid-state batteries", *Journal of Alloys and Compounds*, vol. 979, 2024.
- [4] A. Bin Abu Sofian, I. Imaduddin, S. Majid, T. Kurniawan, K. Chew, C. Lay and P. Show, "Nickel-rich nickel–cobalt–manganese and nickel– cobalt–aluminum cathodes in lithium-ion batteries: Pathways for performance optimization", *Journal of Cleaner Production*, vol. 435, 2024.
- [5] H. Kang, J. Ko, S. Song and Y. Yoon, "Recent progress in utilizing carbon nanotubes and graphene to relieve volume expansion and increase electrical conductivity of Si-based composite anodes for lithium-ion batteries", *Carbon*, vol. 219, 2024.
- [6] M. Gutierrez, M. Morcrette, L. Monconduit, Y. Oudart, P. Lemaire, C. Davoisne, N. Louvain and R. Janot, "Towards a better understanding of the degradation mechanisms of Li-ion full cells using Si/C composites as anode", *Journal of Power Sources*, vol. 533, 2022.
- [7] T. Kogure, "Electron Microscopy", in *Handbook of Clay Science*, Elsevier, 2013, pp. 275-317.
- [8] J. Wu, M. Fenech, R. Webster, R. Tilley and N. Sharma, "Electron microscopy and its role in advanced lithium-ion battery research", *Sustainable Energy & Fuels*, vol. 3, no. 7, pp. 1623-1646, 2019.
- [9] S. Zhou, K. Liu, Y. Ying, L. Chen, G. Meng, Q. Zheng, S. Sun and H. Liao, "Perspective of operando/in situ scanning electron microscope in rechargeable batteries", *Current Opinion in Electrochemistry*, vol. 41, 2023.
- [10] W. Li, D. Lutz, L. Wang, K. Takeuchi, A. Marschilok and E. Takeuchi, "Peering into Batteries: Electrochemical Insight Through In Situ and Operando Methods over Multiple Length Scales", *Joule*, vol. 5, no. 1, pp. 77-88, 2021.
- [11] P. Paul, E. McShane, A. Colclasure, N. Balsara, D. Brown, C. Cao, B. Chen, P. Chinnam, Y. Cui, E. Dufek, D. Finegan, S. Gillard, W. Huang, Z. Konz, R. Kostecki, F. Liu, S. Lubner, R. Prasher, M. Preefer, J. Qian, M. Rodrigues, M. Schnabel, S. Son, V. Srinivasan, H. Steinrück, T. Tanim, M. Toney, W. Tong, F. Usseglio-Viretta, J. Wan, M. Yusuf, B. McCloskey and J. Nelson Weker, "A Review of Existing and Emerging Methods for Lithium Detection and Characterization in Li-Ion and Li-Metal Batteries", Advanced Energy Materials, vol. 11, no. 17, 2021.
- [12] T. Waldmann, A. Iturrondobeitia, M. Kasper, N. Ghanbari, F. Aguesse, E. Bekaert, L. Daniel, S. Genies, I. Gordon, M. Löble, E. De Vito and M. Wohlfahrt-Mehrens, "Review—Post-Mortem Analysis of Aged Lithium-Ion Batteries: Disassembly Methodology and Physico-Chemical Analysis Techniques", *Journal of The Electrochemical Society*, vol. 163, no. 10, pp. A2149-A2164, 2016.

- [13] Y. Yuan, K. Amine, J. Lu and R. Shahbazian-Yassar, "Understanding materials challenges for rechargeable ion batteries with in situ transmission electron microscopy", *Nature Communications*, vol. 8, no. 1, 2017.
- [14] S. Burgess, H. James, P. Statham and L. Xiaobing, "Using Windowless EDS Analysis of 45-1000eV X-ray Lines to Extend the Boundaries of EDS Nanoanalysis in the SEM", *Microscopy and Microanalysis*, vol. 19, no. 2, pp. 1142-1143, 2013.
- [15] L. Peng, Q. Chen, X. Liang, S. Gao, J. Wang, S. Kleindiek and S. Tai, "Performing probe experiments in the SEM", *Micron*, vol. 35, no. 6, pp. 495-502, 2004.
- [16] C. Fincher, D. Ojeda, Y. Zhang, G. Pharr and M. Pharr, "Mechanical properties of metallic lithium: from nano to bulk scales", *Acta Materialia*, vol. 186, pp. 215-222, 2020.
- [17] L. Mele, S. Konings, P. Dona, F. Evertz, C. Mitterbauer, P. Faber, R. Schampers and J. Jinschek, "A MEMS -based heating holder for the direct imaging of simultaneous in-situ heating and biasing experiments in scanning/transmission electron microscopes", *Microscopy Research and Technique*, vol. 79, no. 4, pp. 239-250, 2016.
- [18] H. Zheng, D. Xiao, X. Li, Y. Liu, Y. Wu, J. Wang, K. Jiang, C. Chen, L. Gu, X. Wei, Y. Hu, Q. Chen and H. Li, "New Insight in Understanding Oxygen Reduction and Evolution in Solid-State Lithium–Oxygen Batteries Using an in Situ Environmental Scanning Electron Microscope", *Nano Letters*, vol. 14, no. 8, pp. 4245-4249, 2014.
- [19] H. Shi, X. Liu, R. Wu, Y. Zheng, Y. Li, X. Cheng, W. Pfleging and Y. Zhang, "In Situ SEM Observation of Structured Si/C Anodes Reactions in an Ionic-Liquid-Based Lithium-Ion Battery", *Applied Sciences*, vol. 9, no. 5, 2019.
- [20] D. Chen, S. Indris, M. Schulz, B. Gamer and R. Mönig, "In situ scanning electron microscopy on lithium-ion battery electrodes using an ionic liquid", *Journal of Power Sources*, vol. 196, no. 15, pp. 6382-6387, 2011.
- [21] Y. Qiu, G. Rong, J. Yang, G. Li, S. Ma, X. Wang, Z. Pan, Y. Hou, M. Liu, F. Ye, W. Li, Z. Seh, X. Tao, H. Yao, N. Liu, R. Zhang, G. Zhou, J. Wang, S. Fan, Y. Cui and Y. Zhang, "Highly Nitridated Graphene– Li 2 S Cathodes with Stable Modulated Cycles", *Advanced Energy Materials*, vol. 5, no. 23, 2015.
- [22] P. Baudry and M. ARMAND, "In situ observation by SEM of positive composite elec- trodes during discharge of polymer lithium batteries", Solid State Ionics, vol. 28-30, pp. 1567-1571, 1988.".
- [23] P. Raimann, N. Hochgatterer, C. Korepp, K. Möller, M. Winter, H. Schröttner, F. Hofer and J. Besenhard, "Monitoring dynamics of electrode reactions in Li-ion batteries by in situ ESEM", *Ionics*, vol. 12, no. 4-5, pp. 253-255, 2006.
- [24] T. Tsuda, K. Hosoya, T. Sano and S. Kuwabata, "In-situ scanning electron microscope observation of electrode reactions related to battery material", *Electrochimica Acta*, vol. 319, pp. 158-163, 2019.
- [25] S. Kaboli, H. Demers, A. Paolella, A. Darwiche, M. Dontigny, D. Clément, A. Guerfi, M. Trudeau, J. Goodenough and K. Zaghib, "Behavior of Solid Electrolyte in Li-Polymer Battery with NMC Cathode via in-Situ Scanning Electron Microscopy", *Nano Letters*, vol. 20, no. 3, pp. 1607-1613, 2020.

Lidar systems testing considerations for field use

Helena Picmausová Fraunhofer Institute of Optronics, System Technologies and Image Exploitation IOSB Ettlingen, Germany; University of Defence Brno, Czech Republic https://orcid.org/0000-0003-3683-7243 Jan Farlík University of Defence Brno, Czech Republic <u>https://orcid.org/0000-0001-</u> 7254-2405 Marc Eichhorn Fraunhofer Institute of Optronics, System Technologies and Image Exploitation IOSB Ettlingen, Germany; Karlsruhe Institute of Technology Karlsruhe, Germany <u>marc.eichhorn@iosb.</u> <u>fraunhofer.de</u> Christelle Kieleck Fraunhofer Institute of Optronics, System Technologies and Image Exploitation IOSB Ettlingen, Germany christelle.kieleck@iosb.fraunhofer .de

Abstract— The aim of this paper is to offer a perspective on testing a commercially available lidar system in order to determine its suitability for various practical tasks including mapping, object recognition, and the potential in its coupling with another sensor, in this case a camera. Several mapping missions were carried out over the course of the experiment, with both the lidar and the camera mounted on an Unmanned Aircraft System. Aside from mapping ordinary objects like trees, vehicles, people, and ground elevation, a standardized test target was designed for the purposes of the experiment, and placed in detection area. Influence of external factors on system performance was evaluated, e.g. atmospheric conditions and material properties of scanned surfaces, especially their reflectivity. Post processing of obtained data was carried out, demonstrating the potential of overlaying multiple sensor data for easier object recognition, and an optimal use case for the system is suggested.

Keywords— Sensors, lidar, Unmanned Aircraft Systems, 3D mapping, camera imaging, object detection

I. INTRODUCTION

In recent years, Lidar (light detection and ranging) has been an object of extensive research and development, with ever increasing prominence among active electro-optic sensors. With fields of application spanning from meteorology and geoscience to object identification, autonomous navigation, to aerospace and defence, the technology has also become increasingly accessible, affordable, and compact [1]-[4]. With the variety of solutions available on the market, it is important to choose an appropriate system for a given application, and keep limitations of technology in mind.

II. LIDAR SYSTEMS CONSIDERATIONS

Several parameters need to be accurately calculated and tested to determine effective range and operational limitations of a lidar system. Usually, the main performance factors are the laser source – its output power, pulse length, and wavelength, update frequency of the positioning system and atmospheric conditions. One of the advantages of lidar is that unlike many passive sensors, it can function in low light conditions, and given an equivalent-sized optical aperture, higher resolution can be achieved thanks to a lower diffraction limit at shorter wavelengths [1], [11].

For lidar applications, among the most popular laser sources are diode pumped solid-state lasers, which allow for nanosecond-class pulse generation using Q switching, [7]. Thanks to their widespread availability and maturity, Nd3+:YAG doped lasers emitting at 1064 nm can be therefore encountered in lidar solutions, despite their relatively lower quantum efficiency [1], [10].

Another increasingly popular option is a diode-pumped fiber laser solution, which generally offers exceptional beam quality and resilience to mechanical shock due to the absence of freespace optics. Furthermore, lasers operating above 1550 nm wavelength are deemed ,eye-safer', and also offer easy availability of components owing to telecommunication applications [5], [7].

Diode lasers operating around the 905 nm range are also often used in lidars, especially in the automotive industry, due to their compact size, low cost and good efficiency, and naturally, resolution is better with shorter wavelengths; but as they cannot store energy, their low pulse energies result in considerably shorter detection ranges [8], [9], [11].

Atmospheric conditions should also be taken into account when evaluating lidar performance results. Unlike atmospheric scattering, atmospheric absorption is wavelength specific for each air component, so a careful choice of wavelength can impact beam power decay on the two-way path to the target and back. Neither the 905 nm, 1064 nm, nor the 1550 nm range offer zero absorption. Atmospheric humidity is therefore a key factor when attempting any lidar measurement [10].

Illuminance is another key factor for lidar detection range, which is significantly reduced with higher ambient illuminance. Performance testing on a sunlit day (approx. 100 klx) reportedly results in a reduction of detection range to around 70% of values obtained during clear night conditions (0 klx), [12], [13].

DOI 10.13164/eeict.2024.192

III. EXPERIMENT

The experiment consisted of testing capabilities of a commercially available automotive-grade lidar system coupled with an RGB mapping camera, both coupled with an Inertial Measurement Unit (IMU) and mounted on a UAS (Unmanned Aircraft System), make DJI Matrice 300 RTK. A test area has been set up allowing for multiple mapping flights at various altitudes, and containing various sized targets and landscape elements such as ground elevation, tents, buildings, vehicles and trees.

The goals of the experiment included verification of the accuracy stated in the lidar datasheet, the influence of flight altitude, detection capabilities of both sensors (the lidar and the camera) of objects with various reflectance and their mutual data agreement.

Given the airborne mapping approach, an accurate positioning solution is crucial in this experiment, and although UAS navigation benefits from the advancement of global positioning system (GPS), the inertial navigation system (INS) measurements are key, making the IMU update frequency an important factor, as well as yaw, pitch, and angle/roll accuracy.

A. Test target

Aside from the above-mentioned terrain elements, test targets were designed, providing a measurement standard for the lidar resolution capability, as accuracy can be greatly improved by pre-acquired metric data, as reported in [6].

The experimental targets were a new original design created for the purpose of this experiment. They were 3D printed from materials with various reflectivity to determine the influence of colour and material as well as dimensions.

Figure 1 and Figure 2 show the two versions of the test target, each consisting of multiple cubes with side lengths of 1, 2, 4, and 8 cm, allowing for resolution testing of horizontal, vertical and depth/distance accuracy.



Fig. 1. Test target no. 1, consisting of low reflectivity black cubes and high reflectivity white cubes, providing high contrast.



Fig. 2. Test target no. 2, combining low reflectivity black cubes and medium reflectivity red cubes, providing lower contrast.

B. Lidar system specifications

The manufacturer-declared specifications include the lidar ranging accuracy of 2-3 cm at a 100 m distance, 905 nm wavelength, 30 W system power, IMU update frequency of 200 Hz and a point return rate over 240,000 pts/s, with a detection range reaching over 400 m at 80% reflectivity, and the operating temperature range from -20° to 50° C. The software used offers flight route planning with flight track overlap and angle options, time-of-flight estimation, and postprocessing tools including an overlay of the acquired lidar data point cloud and photogrammetry images from the 20 MP, 1" CMOS camera.

C. Flight conditions

The test area was scanned during multiple mapping flights at 50 and 100 m altitude. The scanning area overlap was set to 20%



Fig. 3. A photo taken by the RGB camera from flight altitude of 50 m, with a clearly visible white tent, a few people standing on the walkway and a very small but visible target – circled in yellow and red, respectively.

for the 50 m, and at 50% for the 100 m flight. The experiment was carried out during clear weather conditions, calm wind, temperature of 18° C, humidity levels under 30%, in full daylight. A reference photo of target area is shown in Figure 3.

IV. RESULTS

The results using both the lidar point cloud, and an overlay from photogrammetry can be seen below in Figure 4. It is worth noting that the side of the white tent hasn't been detected at all, likely due to the scanning incidence angle, a similar issue to one previously reported in [6], that can be sometimes ameliorated by energy balance normalization and realizing multiple scans with different angles. Figure 5 depicts the same area colorized according to reflectivity levels as perceived by the system in post-processing. The bright red stripes of area do not represent a real increased reflectivity however, but rather areas of scanning flight overlap that has been incorrectly evaluated by the software as high reflectivity due to a denser point cloud.

The test target is practically invisible in both pictures.

A. Mapping flight at 50 m altitude



Fig. 4. Results using both the lidar point cloud and an overlay from photogrammetry at 50 m altitude.



Fig. 5. Colorization according to reflectivity levels obtained from the data.



Fig. 6. Results using both the lidar point cloud and an overlay from photogrammetry at $100\ {\rm m}$ altitude.



Fig. 7. Colorization according to ground point classification obtained from the data.

B. Mapping flight at 100 m altitude

Figure 6 is again an overlap of camera data and lidar point cloud, this time taken from a 100 m altitude. At this range, it becomes nearly impossible to identify the small group of people. Figure 7 shows a different colour mode based on ground-point recognition. The point cloud is colorized in post-process to best differentiate between the ground and all other objects, and sometimes can offer a clearer overview than reflectivity or height mapping.

A different area with a distinct treeline, a road, and several cars was also mapped. Below, in Figure 8 the ground point recognition offers a clear distinction between the forest and cars, it however fails to recognize an asphalt road.

This road is clearly visible in Figure 9, thanks to reflectivity measurement carried out on the same flight. In Figure 9 it is also

worth noting that different reflectivity was detected on different vehicles due to their paint colours.



Fig. 8. Ground point recognition in different scene.



Fig. 9. Colorization according to reflectivity levels obtained from the data.

V. CONCLUSIONS

Commercial lidar solutions are becoming ever more widespread and accessible. At an affordable price-point, with an intuitive interface and efficient handling they offer a powerful tool for terrain mapping, autonomous driving and object detection and classification, especially when paired with other sensors like radar and cameras.

It is important, however, to keep the limitations of each system in mind when choosing a solution for a particular scenario. While a lightweight, automotive lidar can offer a fairly reliable mapping of the surrounding area, with distance its resolution can deteriorate far below the manufacturer stated values.

As expected, surface reflectivity had influence on obtained data point cloud density, with bright metallic surfaces reflecting more rays than dark and plastic materials. This effect could be observed on car paint and various tent and buildings materials, but not on the intended test targets, as those proved too small to map regardless of material properties. During the course of the experiment, the only sensor able to detect the testing target was the RGB camera, with the lidar barely picking up enough cloud points at 50 m distance. Achieved resolution values did not exceed 10 cm in either depth or width.

In general, field-experiment results can rarely achieve the values obtained in a laboratory environment. This experiment was carried out in as near-optimal working conditions for the system as possible given the time of year and local climate. Temperature, humidity levels and light conditions were all within the working range of the system. Wind speed stayed between calm and a light breeze throughout all flights.

From the standpoint of system compatibility, the data obtained from the camera and the lidar highly corresponded, and their overlap in post-processing posed no issues.

Should the 3D data be evaluated by an algorithm, the lidar dataset would likely suffice for object recognition on its own. With a person in-the-loop, however, the colourization provided by the RGB camera proved an added value, making the target area much easier to read to the human eye.

Used lidar system is a typical automotive-grade lidar, therefore mounting it on a UAS may add error to the results, increasing the importance of INS measurements accuracy.

Even so, in the experiment, the lidar system proved effective in mapping the ground elevation and larger objects like trees and vehicles, and at a shorter distance, even people. The interface offered time-effective mission planning and intuitive controls, and required little-to-no additional training for the UAS operator.

For higher resolution object detection and classification, and for larger distances, however, a different, higher energy lidar system would be more appropriate.

ACKNOWLEDGMENT

The authors would like to thank the Exercise Control Team of the NATO Counter – Unmanned Aircraft Systems Technical Interoperability Exercise 2023 (NATO C-UAS TIE 23), for their invaluable support in the realization of the experiment: for providing technical facilities, equipment and a dedicated test area, as well as for their extensive expertise, helpful advice and positive attitude.

REFERENCES

- P. McManamon, "Field Guide to Lidar," in SPIE Field Guides, vol. FG36, SPIE PRESS, Bellingham, Washington USA, 2015. ISBN: 9781628416558.
- [2] L. C. G. David, A. H. Ballado, S. M. Sarte and R. A. Pula, "Mapping inland aquaculture from orthophoto and LiDAR data using object-based image analysis," 2016 IEEE Region 10 Humanitarian Technology Conference (R10-HTC), Agra, India, 2016, pp. 1-5, doi: 10.1109/R10-HTC.2016.7906855.
- [3] E. J. Welton et al., "The NASA Micro Pulse Lidar Network (MPLNET): Early Results from Development of Diurnal Climatologies," 2021 IEEE International Geoscience and Remote Sensing Symposium IGARSS, Brussels, Belgium, 2021, pp. 1432-1433, doi: 10.1109/IGARSS47720.2021.9553788.
- [4] J. -K. Liu, K. -T. Chang, C. Lin and L. -C. Chang, "Accuracy evaluation of ALOS DEM with airborne LiDAR data in Southern Taiwan," 2015 IEEE International Geoscience and Remote Sensing Symposium

(*IGARSS*), Milan, Italy, 2015, pp. 3025-3028, doi: 10.1109/IGARSS.2015.7326453.

- [5] L. Holmen, G. Rustad, and M. Haakestad, "Eye-safe fiber laser for longrange 3D imaging applications," Appl. Opt. 57, 6760-6767, 2018.
- [6] C. Bodensteiner, W. Hübner, K. Jüngling, P. Solbrig and M. Arens, "Monocular Camera Trajectory Optimization using LiDAR data," 2011 IEEE International Conference on Computer Vision Workshops (ICCV Workshops), Barcelona, Spain, 2011, pp. 2018-2025, doi: 10.1109/ICCVW.2011.6130496.
- [7] C. Larat, M. Schwarz, E. Lallier, and E. Durand, "Eye-Safe Q-Switched Er:YAG MOPA Laser System," in Advanced Solid-State Lasers Congress, G. Huber and P. Moulton, eds., OSA Technical Digest (online) (Optica Publishing Group, 2013), paper JTh2A.47.
- [8] A. Knigge, A. Klehr, A. Liero, J. Fricke, A. Maaßdorf, A. Zeghuzi, H. Wenzel, and G. Tränkle, "Wavelength stabilized 905 nm diode lasers in the 100 W class for automotive LiDAR," in 2019 Conference on Lasers and Electro-Optics Europe and European Quantum Electronics Conference, OSA Technical Digest (Optica Publishing Group, 2019), paper cb_5_2.

- [9] A. Laugustin, C. Canal, and O. Rabot, "State-of-the-art laser diode illuminators for automotive LIDAR," in 2019 Conference on Lasers and Electro-Optics Europe and European Quantum Electronics Conference, OSA Technical Digest (Optica Publishing Group, 2019), paper cb_p_15.
- [10] T. F. Refaat *et al.*, "Backscatter 2-µm Lidar Validation for Atmospheric CO2 Differential Absorption Lidar Applications," in *IEEE Transactions* on Geoscience and Remote Sensing, vol. 49, no. 1, pp. 572-580, Jan. 2011, doi: 10.1109/TGRS.2010.2055874.
- [11] M. Eichhorn, "Laser Physics: From Principles to Practical Work in the Lab," Springer, 2014, ISBN : 978-3-319-05127-7.
- [12] C. Kyba, A. Mohar, T. Posch, "How bright is moonlight?" A&G Astronomy and Geophysics, 58, 1, pp. 1.31—1.32., 2017, doi: http://doi.org/10.1093/astrogeo/atx025.
- [13] Livox Avia User Manual, Available online: <u>https://terra-1-g.djicdn.com/65c028cd298f4669a7f0e40e50ba1131/Download/Avia/Livox%20Avia%20User%20Manual%20202204.pdf</u>.

Evaluation of the gamma-ray spectrum transmitted upon alteration in scintillator shading

Tadeáš Zbožínek Department of Microelectronics, FEEC Brno University of Technology Brno, The Czech Republic 214774@vut.cz Michal Jelínek Institute of Scientific Instruments Czech Academy of Sciences Brno, The Czech Republic jelinmi@isibrno.cz

Abstract— This paper deals with determining the properties of optical fibres required for transmitting energy spectra in ionising radiation measurements by scintillation. Specifically, optical fibres that provide the communication link between the scintillator and the photosensitive electronic devices. The primary objective is focused on the basic diameter and numerical aperture values required for efficient transmission of energy spectra. An iris mechanism gradually shaded the scintillator crystal, simulating the optical fibre diameter. The results can be used to calculate the numerical aperture to achieve this goal. This way, we can establish basic properties for further energy spectra measurements with optical fibres.

Keywords—Gamma, spectroscopy, scintillation, ionising radiation, scintillator shading, optical fibres

I. INTRODUCTION

Gamma-ray spectroscopy is a widely used method for enhanced characterisation of ionising radiation. This radiation is often associated with increased caution because exposure to it also means exposure to health hazards. Therefore, it needs to be well-characterized to avoid unnecessary accidents. Applications of this spectroscopy are found in many nuclear industry sectors, such as medicine, power generation, research, security, and other specific areas.

Gamma-ray spectroscopy provides information on decaying elements, such as the type of radioisotope and its quantity. To understand this spectroscopy, it is appropriate to describe gamma radiation. Gamma radiation comprises high-energy photons released during radioactive decay in the atomic nucleus. Gamma photons carry no charge, so their measurement gains additional complexity. For this reason, gamma spectroscopy may be primarily conducted through two methods. Germanium semiconductors are used for direct photon detection, or scintillation materials are used with lower energetic photon detectors. This paper focuses on scintillation materials. [1]

For the accurate resolution of radioisotopes, it is essential, among other factors, to ensure the efficient transmission of a light signal from the scintillator to the photosensitive electronics (photomultiplier). In general connection, it is expected to have a scintillator directly attached to the photomultiplier tube for minimal photon loss. However, when measuring high ionising radiation levels, the connection suffers from a reduced electronics life span.

One of the methods to separate electronics from ionising radiation is the adoption of optical fibres to transport

Břetislav Mikel Institute of Scientific Instruments Czech Academy of Sciences Brno, The Czech Republic mikel@isibrno.cz

scintillation radiation. In this way, it extends the lifetime of electronic devices in long-term high-energy measurements.

We have tested an unusual connection with optical fibre between these parts [2]. Using this type of separation brings signal loss, resulting in difficulty in source resolution with gamma-ray spectroscopy. The paper focuses on preparing minimal start conditions for the first fibre transfer to measure gamma-ray spectra.

II. CURRENT STATUS OF OUR WORK

Optical fibres are frequently used in scintillation ionising radiation measurements without spectra resolution. Yet, they come with a downside hidden behind direct radiation measurement. The silica and plastic optical fibres are used for the measurement. The plastic fibres are generally more sensitive to ionising radiation, which causes their attenuation to increase when exposed to it. Due to their cost, they are typically used for short-term measurements, which can be left to recover between further measurements. In contrast, silica fibres are less sensitive to ionising radiation and are suitable for long-term measurements [3] [4] [5]. Thus, we decided to use optical fibres only as a communication channel between a scintillation material and a photosensitive electronic. They maintain the advantages of silica fibre ionising radiation resistance and avoid low sensitivity during measurement. This allows ionising radiation to be effectively monitored. In contrast to this measurement, detecting the scintillation radiation and transmitting enough scintillation radiation for spectral analysis is essential. The current stage of our measurement is described in Fig. 1.



Fig. 1. The measurement of ionizing radiation using optical fibres [2]

In this paper, we focus on two properties, namely the numerical aperture and effective diameter of transferring fibres, which led to the number of photons transferred. The numerical aperture and diameter proved critical in our activity measurement [2]. Our focus is on transferring gammaray spectra for radioactive source identification. We need to find the minimal part of scintillation light (number of photons) that still allows us to recognise different radioisotopes. To simulate this problem, we shaded the scintillation crystal (shading the photons) to discover the limits of transferring the gamma-ray spectrum.

III. GAMMA-RAY SPECTROSCOPY

Gamma-ray spectroscopy with scintillation materials uses a simple principle. Gamma photons are absorbed or transfer part of their energy to the scintillator, which brings the scintillator to an excited state. After a while, the scintillator deexcites and, through internal conversion, the material emits photons in ultraviolet and visible parts of the light spectrum. Number of emitted photons directly corresponds to the energy of the incident gamma ray. These photons create a form of light pulses, which are easily detectable by photomultiplier. In photomultiplier, emitted photons generate electrons after hitting the electrode. The electric pulse generated in this manner corresponds to the light pulse emitted by the scintillator, indicating the energy of the incident gamma ray. Therefore, we can get energy and activity information about the source. Energy from peak height and activity from the number of pulses.

Every radioisotope has a specific energy spectrum. Moreover, some are mono-energetic, meaning their decay produces mainly gamma rays with the same energy. Knowing the spectra of radioisotopes, we can find their occurrences in our measurement.

The energy spectrum of gamma radiation obtained by scintillation includes several sections, which can be specific for every radioisotope. The spectrum also contains noise, which can make resolution harder. Various scintillation materials exhibit different levels of resolution and background noise, resulting in distinct spectra. The theoretical spectrum and spectrum obtained by measuring scintillation crystal Na(Tl) used are shown in Fig.2.

The primary way to distinguish radioisotopes is to refer to their prominent energy peaks (photopeaks). Photopeaks represent gamma energy directly. They are created during the absorption of gamma photon (photoelectric effect) in the scintillator. Another way is to compare the shape of Compton scattering. During scattering, a photon loses part of its initial energy and gives it to an electron in the scintillator. Transmitted energy depends on changes in the scattering angle. A gamma photon's highest energy, the Compton edge, is visible in the spectrum. This edge corresponds to the maximum angle (180 degrees) during scattering, and it is linked to the energy of the photopeak.

An exciting form of noise is peaks that look like photopeaks. We could identify them as fake peaks, which can be easily distinguished in known environments. In unknown environments, these could be described as photopeaks, thus leading to an error in measurement. They are connected to high-



Fig. 2. The gamma ray theoretical energy spectrum with important regions highlighted (A); Measured gamma-ray spectrum of a scintillator NaI(Tl) (B)

energy photons and the presence of high atomic number (high-Z) materials. Precisely backscatter peak (gamma photon backscatter from shielding), annihilation peak (incident photon from shielding created during positron annihilation), and escape peaks (created during pair production and positron annihilation directly in the scintillator).

The last noise we mention is X-ray, which is visible at lower energies in the spectrum. The primary source of this X-ray is the surrounding high-Z material. [1]

IV. EXPERIMENTAL SET-UP

The gamma-ray spectrum is measured using scintillators that are coupled to a photomultiplier. The photomultiplier can detect visible and ultraviolet radiation from the scintillator and generate pulse signals. A digitiser is utilised to process the signal. Changing the signal from analogue to digital opens the way to evaluate the data. The experimental setup is depicted schematically in Fig. 3.

The experiment uses a Caesium-137 reference source.¹³⁷Cs is a radioisotope commonly employed for detector characterisation due to its distinctive single photopeak. Additionally, nuclear states possess well-defined energies, resulting in the emission of gamma rays that are nearly monoenergetic. The decay process of Caesium begins with beta



Fig. 3. The Schematic Connection of the Experiment

decay, followed by the emission of a gamma-ray with an energy of 662 keV [1]. The decay process is shown in Fig. 4.

The experimental setup can then be divided into parts representing its function: scintillator, scintillation detection, spectrum acquisition and iris.

A. Scintillator

The first part of the whole measurement is the scintillator. The scintillators convert high-energy particles, such as gamma radiation photons, into visible light. This conversion process is facilitated by fluorescence in certain materials, resulting in rapid emission of light photons upon radiation interaction. This generates a brief weak pulse for observation. However, competing processes, especially phosphorescence and slow fluorescence, can reduce the scintillator's efficiency by inducing an afterglow effect. There are two types of scintillation materials: inorganic and organic scintillators. Each scintillator material operates through distinct luminescent mechanisms, resulting in unique advantages. Inorganic materials typically exhibit high light yield, while organic counterparts offer rapid response times.

We used sodium iodide, an inorganic material doped with Thulium NaI(Tl). Generally, this material has good properties compared to the price and is widely used in the nuclear industry. The light emission peak is around 420 nm, the decay time is 250 ns, and the light yield is 38 photons/keV for gamma radiation. Another property of scintillators is energy resolution, which represents their response to a monoenergetic radiation source. A simple formula does the calculation:

$$R = \frac{FWHM}{H} \tag{1}$$

Here, *R* represents resolution, *FWHM* stands for the full width at half maximum, and *H* signifies the height of the peak.

The measurement proceeds with NaI(Tl) crystal in a cylinder 1 inch in diameter. These dimensions consider additional fibre optic involvement. The surfaces are enveloped in high-reflection material to maximise signal gain, leaving one of the circular sides uncovered for light extraction. Furthermore, all enclosed sides are further shielded by an aluminium cover, while the extraction side is protected by glass. The cover is thick enough for shielding from alpha and beta particles. Therefore, this scintillator is designed to detect X-rays and gamma rays. [1] [6]



Fig. 4. The Caesium-137 decay scheme. [1]

B. Scintillation detection

A photomultiplier detects photons emitted by the scintillator. The photomultiplier consists of a photocathode, dynodes, and an anode placed in a vacuum. The photons strike the photocathode and emit electrons. The emitted electrons then hit the dynodes. Due to the high voltage applied to the dynodes, the incident electrons eject secondary electrons, multiplying their number. After multiple dynodes, the multiplied number of electrons now reaches the anode, where the amplified signal is collected. Through this process, the signal is strong enough for later processing. The photomultiplier requires a high-voltage source for proper function. For the low-voltage version, the silicon photomultiplier can be utilised. In simplified terms, the Silicon PMT consists of multiple rows of avalanche diodes arranged on a silicon substrate close to each other. Then, each photon triggers the activation of a single diode, generating an electric impulse.

A comparison of these two photomultipliers will be the subject of a further experiment. In this first measurement, we use a high-voltage photomultiplier tube with a blue–green bialkali sensitive photocathode, magnetic shielding, and an active divider. The high-voltage photomultiplier was chosen because of the general use of these devices in nuclear measurement. [1] [7]

C. Spectrum acquisition

The signal is subsequently transmitted to the digitiser. Specifically, we utilise a 10-bit, 2/1 GS/s DT5751 digitiser manufactured by CAEN S.p.A., incorporating pulse processing firmware. Specifically, a multichannel analyser and a counter are included. The counter facilitates count rate measurements, from which the activity of a radioisotope can be calculated with proper calibration. In our experiment, we prioritise the analysis of energy spectra, hence utilising a multichannel analyser. Analysed signal is brought to an input of the digitiser. Analog signal is then divided into discrete levels represented by bits. This is essential because we need to know the height of a single pulse and, thus, the energy of gamma-ray to calculate the energy spectrum.

D. Iris

We used Iris to determine the minimum parameters required. This involves specifying the needed diameter and NA and then recalculating these values for optical fibres. It is essential first to define these values to set the requirements to allow us to transmit spectra.

V. MEASUREMENT OF GAMMA-RAY SPECTRA

The measurement of spectra occurs within the described setup with a minor variation. In this configuration, the scintillator is not directly attached to the photomultiplier. Instead, an iris aperture is positioned between them to simulate a numerical aperture. Radioactivity etalon is placed next to the scintillator. The measurement lasts for 10 minutes, during which 600 values are recorded (1 value per second). These values represent the number of counts per second. During the measurement, energy levels are converted into a spectrum. Following each measurement, the iris is opened and measured. Through this procedure, we can observe the behaviour of spectra under varying degrees of shading and approximately determine the point at which we can no longer accurately extract information from the spectra.

A. Data evaluation

To discern the readability of the measured gamma-ray spectrum, we shall utilise the theoretical framework outlined in Chapter Two. For the correct assumption, the spectrum needs to include photopeak. The peak provides an accurate energy number through energy calibration, thus identifying the radioisotope. Using high-resolution inorganic scintillators, the identification of photopeaks should be sufficient. However, challenges arise when utilising lower-resolution scintillators or operating within environments characterised by elevated noise levels. This is particularly evident in high-Z materials adjacent to or within the scintillator material itself. There were indications in Chapter Two of possible fake photopeaks that could make a spectrum unreadable. The following spectrum element that can be used is the Compton edge. Compton edge is directly connected to the photopeak trough equation:

$$E_C = \frac{hv}{1 + 2hv/m_0 c^2} \tag{2}$$

where hv refers to the energy of gamma rays, h stands for Planck's constant, and v stands for the frequency. The symbols m and c represent an electron's mass and the speed of light. Then E_c refers to the distance between the photopeak and Compton edge. This energy tends to have a higher energy gamma-ray approach value of 256 MeV [1]. For some scintillators, photopeaks are not displayed within the spectrum, so the radioisotopes have to be calculated based on the character of this edge. Also, the number of Compton edges indicates the number of photopeaks. Following this theory, it is imperative to recognise both the Compton edge and potentially the photopeak to identify radioisotopes from a gamma-ray spectrum.

This method can be applied to the measured data for spectrum recognition. The spectra obtained from the iris measurements are shown in Figure 5. The iris was gradually closed. After the measurements, the diameters were measured for accurate values. A photopeak can be observed in all spectra except the closed one, which tends to add to the Compton continuum with a smaller iris diameter. At 74% iris opening, the Compton continuum begins to disappear. At 55%, the Compton portion of the spectrum distorts to a peak, and source resolution is no longer possible. The peaks cease to be resolvable at 35% opening of the iris, and at 5% (maximum closed iris), we get only a single peak, a combination of the entire energy spectrum. The spectral transmission



Fig. 5. The energy spectra of gamma radiation (Cs137) obtained by shading the scintillation detector NaI(Tl); each curve corresponds to the spectrum of a specific opening of the iris.

limitation could then be set at 74% iris opening. The numerical aperture of the iris was calculated to be 0,988. Applying this knowledge to optical fibres, the fibre should have a diameter of at least 74% of the iris opening, i.e. 8.4 mm.

The second observable effect is a shift in the position of the entire measured spectrum. This could be caused by a lower number of photons (lower pulse signal) thanks to Iris closing.

VI. CONCLUSION

The measurements of the gamma-ray scintillation spectrum were made for different sizes of regions defined by the diameter of the closing aperture. The minimum diameter is 8.4 mm for a 1-inch cylindrical NaI(Tl) crystal. This creates problems for measurements using optical fibres. Multiple fibres can be used to ensure these dimensions instead of single fibres with a large diameter. The various fibres provide us with options for their placement. Due to their numerical aperture, they do not need to be placed closely together, but they have a small spacing, which could reduce the number of fibres required.

Further measurements are needed to confirm these considerations, especially now that we have determined the minimum diameter. We will continue experimenting with fibre shading to find the optimal position and number of the fibres needed.

ACKNOWLEDGEMENT

This publication was created with the state support of the Technology Agency of the Czech Republic within the framework of the National Centers of Competence Programme (project No. TN02000020). The work was supported by the European Regional Development Fund-Project " Interdisciplinary Collaboration in Metrology with Cold Objects Fibre Networks Quantum and (No. CZ.02.1.01/0.0/0.0/16_026/0008460) and CAS by (RVO:68081731).

REFERENCES

- [1] KNOLL, Glenn F. Radiation detection and measurement. 4th ed. Hoboken: Wiley, 2010. ISBN 978-0-470-13148-0.
- [2] JELÍNEK, Michal; MIKEL, Bretislav a ZEMÁNEK, Pavel. Optical fibers forming to ionizing radiation sensors preparation. online. In: 21st Czech-Polish-Slovak Optical Conference on Wave and Quantum Aspects of Contemporary Optics. Lednice: SPIE, 2018, s. 14-. ISBN 9781510626072. ISSN 0277-786X. Dostupné z: https://doi.org/10.1117/12.2518106. [cit. 2022-05-04].
- [3] BILLINGSLEY, John; O'KEEFFE, S.; FITZPATRICK, C.; LEWIS, E. a AL-SHAMMA'A, A.I. A review of optical fibre radiation dosimeters. online. *Sensor Review*. 2008, roč. 28, č. 2, s. 136-142. ISSN 0260-2288. Dostupné z: https://doi.org/10.1108/02602280810856705. [cit. 2022-11-21].
- [4] ZUBAIR, H.T.; BEGUM, Mahfuza; MORADI, Farhad; RAHMAN, A.K.M. Mizanur; MAHDIRAJI, Ghafour A. et al. Recent Advances in Silica Glass Optical Fiber for Dosimetry Applications. online. *IEEE Photonics Journal.* 2020, roč. 12, č. 3, s. 1-25. ISSN 1943-0655. Dostupné z: https://doi.org/10.1109/JPHOT.2020.2985857. [cit. 2022-11-21].
- [5] O'KEEFFE, S; FERNANDEZ FERNANDEZ, A; FITZPATRICK, C; BRICHARD, B a LEWIS, E. Real-time gamma dosimetry using PMMA optical fibres for applications in the sterilization industry. online. *Measurement Science and Technology*. 2007, roč. 18, č. 10, s. 3171-3176. ISSN 0957-0233. Dostupné z: https://doi.org/10.1088/0957-0233/18/10/S19. [cit. 2022-11-21].
- [6] NaI(Tl) Scintillation Crystal: Sodium Iodide. online. In: LUXIUM solutions. Hiram (Ohio): Luxium Solutions, 2022. Dostupné z: https://www.crystals.saint-gobain.com/radiation-detectionscintillators/crystal-scintillators/naitl-scintillation-crystals. [cit. 2023-04-12].
- [7] 9266B Series. online. In: ET Enterprises. Uxbridge (London): ET Enterprises, 2023. Dostupné z: https://etenterprises.com/products/photomultipliers/product/p9266b-series. [cit. 2023-05-17].

Jsme skupina E.ON, těší nás

Jsme jedním z největších energetických koncernů u nás i ve světě. Centrálu máme v Německu, ale najdete nás už v 15 zemích Evropy, včetně České republiky. Zakládáme si na tom, že naše energie můžou šetřit peníze i přírodu a čím dál víc využíváme obnovitelné zdroje.

Kdo patří do naší rodiny

E.ON Energie, a.s.

- Obchoduje s elektřinou a plynem, zajišťuje marketing a komunikaci a stará se i o výrobu energií.
- Pro zákazníky připravuje řešení na míru v oblasti fotovoltaiky, tepelné techniky a elektromobility.

EG.D, a.s.

- Distributor energií, který vlastní a provozuje rozvodnou síť elektřiny zejména na jihu Čech a Moravy a rozvodnou síť plynu na jihu Čech.
- Zajišťuje připojení odběrných míst k síti a stará se o dopravu energií k zákazníkům.

E.ON Česká republika, s. r. o.

 Funguje jako podpora výše uvedeným společnostem. Zajišťuje jim služby, jako je účetnictví, právo nebo HR, a na starost má i zákaznickou péči.

Na přírodě nám záleží



E.ON Energie je jedničkou ve výkupu zelené elektřiny v České republice.



EG.D pomocí bezpečnostních prvků na sloupech vysokého napětí chrání ptáky před úrazem.



V budovách E.ONu i ve všech dobíječkách elektromobilů využíváme zelenou energii.

eon

Nastartuj svou kariéru

#spolujsmeeon

Studentské programy

- Letní energetická akademie poslední týden v červenci
- Letní technické brigády od června do září
- Stipendijní program
- Diplomové/bakalářské práce
- Praxe a brigády

Absolventské programy

(plný úvazek)

- Trainee program od září
- Junior technik od září

eon.cz/kariera



O E.ON Kariéra


Enzyme-Based Impedimetric Biosensor dotted with gold nanoparticles

Zuzana Košelová Department of Microelectronics Faculty of Electrical Engineering and Communication, Brno University of Technology, Brno, 612 00, Czech Republic <u>225675@vut.cz</u>

Abstract— This research delves into the realm of biosensor improvement through the utilization of gold nanoparticles (Au NPs). The primary objective is to assess the impact of differentsized Au NPs on sensor performance, specifically investigating whether 100 nm or 20 nm nanoparticles prove more favourable to enhancement. Moreover, we aim to inspect the biosensor's response to varied concentrations of Au NPs, unravelling the involved collaboration between nanoparticle size, concentration, and overall sensor properties. This modification of commercial electrodes with Au NPs, could be way for enhancing surface area and enzyme immobilization. Notably, the investigation also explores the potential drawbacks associated with increasing nanoparticle concentration and offers insights into optimizing biosensor design. It has been observed that while 20 nm Au NPs slightly decreased impedance values at higher glucose concentrations, 100 nm Au NPs, conversely, exhibited an increase in capacitive behaviour. Equally crucial is the parameter chosen for constructing the calibration curve. From impedance values at low frequencies of alternating voltage, such as 2 Hz, a lower Limit of Detection (LOD) is obtained. However, the analysis of R_{ct} in the case of 20 nm Au NPs reveals a broader range of glucose concentrations falling within the calibration area. Through a comprehensive analysis of electrochemical behaviour, impedance, and charge transfer resistance, we endeavour to provide contributions to the improvement of biosensor technologies.

Keywords—gold nanoparticles, biosensors, enzyme immobilization, impedance spectroscopy, glucose detection, calibration curves

I. INTRODUCTION

Over the past decade, there has been a growing interest in the modification of sensors using nanomaterials, offering advantages such as increased surface area, enhanced restriction, and improved selectivity. These modifications prove beneficial for applications like third-generation sequencing and direct detection events involving affinity partners, such as DNA hybridization or antigen-antibody interactions [1]. Various techniques are employed for creating nanoscale structures, ranging from sophisticated methods like ion beam etching to simpler approaches like anodization [2][3]. Among these methods, the packing of spherical nanoparticles (NPs) in a dense planar arrangement has gained attention [4][5]. By precisely choosing the material and shape of NPs, a versatile system suitable for a wide range of applications can be produced. In this study, we explore and tests a strategy to enhance the Limit of Zdenka Fohlerová Department of Microelectronics Faculty of Electrical Engineering and Communication Brno University of Technology Brno, 612 00, Czech Republic <u>fohlerova@vutbr.cz</u>

Detection (LOD) in biosensors, focusing on enlarging the surface area for enzyme storage. The idea is that a densely covered surface by enzyme will reach saturation later. Expanding the electrode surface with NPs provides more space for immobilization while using the same electrodes, allowing us to test whether this approach increases sensitivity. Even though gold nanoparticles (Au NPs) are primarily utilized in optical biosensors, the widespread use in enzyme-based biosensors is notable [6][7][8]. Au NPs exhibit high conductivity and biocompatibility, potentially able forming strong bonds with organic substances (with enhancement of the right chemical This unique feature creates a groups). suitable microenvironment for enzyme immobilization, significantly enhancing enzyme activity [9]. Enzymes, owing to their numerous functional groups such as carboxylic (-COOH), amino (-NH2), thiol (-SH), etc., can be easily immobilized directly onto nanoparticles [9]. There are numerous combinations in which these nanoparticles have been incorporated into biosensing. However, they are mostly used in their oxidized form, as ions, as it makes them more easily electrochemically detectable. For example, Ilkhani et al. [10] also utilized citrate doped Au NPs, although with a differential pulse voltammetry setup (particle diameter = 32 nm, concentration = $0.268 \ \mu g/mL$, LOD = 50 pg/mL). Additionally, Au NPs on gold electrodes have also been tested primarily using amperometry (particle diameter = 11 nm, LOD = 23 μ M) [11], (particle diameter = 3.4 nm, LOD = 0.5μ M) [12], (particle diameter = 2.6 nm, LOD = 8.2 μ M) [13]. However, we are the first to employ AuNPs in a biamperometric setup for calibration measurements using impedance spectroscopy. It is intriguing to note that the addition of nanoparticles has sometimes led to a decrease in biosensor response. According to the authors, this phenomenon may be attributed to the difference in oxidative states between Au NPs and the active centre of the enzyme, as observed in previous studies (colloid containing AuNPs at 5,5 mM) [14]. In our investigation, we aim to scrutinize the effectiveness of sensor enhancement through the utilization of both larger and smaller Au nanoparticles (NPs). The central focus is to discern how the biosensor will respond to varying concentrations of these NPs, evaluating whether their presence will predominantly favour the overall performance or potentially exhibit contrasting effects. By delving into this exploration, we seek to unravel insights that can contribute to refining the design and application of biosensors, particularly those utilizing Au NPs for surface modification.

II. METHODS AND MATERIALS

A. Chemicals

Gold nanoparticles (Au NPs) of diameter 20 and 100 nm stabilized suspension in citrate buffer, poly-L-lysine hydrobromide (p-lys), glucose oxidase (GOD) from Aspergillus Niger, bovine serum albumin (BSA), The mixture of 16 μ l GOD (8 mg/ml) and 25 μ l BSA (16 mg/ml) was then cross-linked with 2.75 μ l 2% glutaraldehyde solution (GO). All were dissolved in PBS (10 mM; pH 7.3). This microliter mixture was applied to the Au NPs-modified electrodes and allowed to dry. The reference sample without GOD was created with 33 μ l BSA (13 mg/ml) and 2.75 μ l of a 2% glutaraldehyde solution. Additionally, a reference without



Fig. 1. Graph presents Nyquist plot (a) and CV (b) for the electrodes after modification steps observed in the presence of the Fe2+/Fe3+ redox probe. Citrat20 and Citrat100 represent the liquid without NPs, serving as reference points to elucidate electrode reactions unaffected by the influence of gold nanoparticles (Au NPs). And samples labelled Au20 and Au100 are electrodes modified with 20 nm Au NPs and 100 nm Au NPs, respectively.

Hexaammineruthenium(III) chloride (98%) (Ru3+), potassium ferrocyanide and potassium ferricyanide (Fe2+/Fe3+), glutaraldehyde solution (GO) (25%), were purchased from Sigma-Aldrich (Germany). D-(+)-Glucose monohydrate, KOH, H_2O_2 (30%), isopropyl alcohol was purchased from Penta (Czech Republic).

B. Electrodes

In this study, commercial electrodes from a biamperometric setup (www.printed.cz, Czech Republic) consisting of a pair of two identical gold disk electrodes were utilized. These electrodes have an internal diameter of 400 μ m and lack a separate reference or counter electrode. The immobilization of gold nanoparticles (Au NPs) with diameters ≈ 20 nm and concentration 6.54 10¹¹ particles/mL, as well as ≈ 100 nm with concentration 3.8 10⁹ particles/mL was carried out on a positively charged poly-L-lysine (p-lys) thin layer.

Before modifying the electrodes, an "activation" process was performed on the gold electrodes. The electrodes underwent polishing with microcloth (Buehler) and isopropyl alcohol to eliminate residual photoresist and passivation layers from the surface. Subsequently, the electrodes were treated with a solution of 0.5 M KOH and 20% H_2O_2 for 10 min. To functionalize the gold electrode, a 50 µg/mL p-lysine in phosphate-buffered saline (PBS; 10 mM, pH 7.4) was adsorbed onto the electrode surface for at least 30 minutes, providing amino groups on the gold surface. Following washing and drying steps, a drop of negatively charged Mu NPs solution was applied to the positively charged modified electrode for 40 minutes.

NPs was created by cross-linking the GOD enzyme with BSA on p-lysine.

Electrochemical impedance spectra (EIS) and cyclic voltammetry (CV) were recorded using the μ AUTOLAB III / FRA2 (Metrohm Autolab, Netherlands) analyser. For EIS enzymatic measurements, the electrode was placed in a beaker with 2 ml of 5 mM [Ru(NH₃)₆]₃ dissolved in PBS. A potential with an amplitude of 50 mV was applied, with a logarithmic distribution of 20 individual frequencies ranging from 100,000 to 2 Hz. The first two measurements were conducted without glucose, followed by a series of measurements with varying glucose concentrations to establish a calibration curve. Glucose was introduced from a solution made of 0.1 M D-(+)-Glucose monohydrate dissolved in PBS.

CV and EIS for confirmation of the modification layers were performed in 50 mM Fe2+/Fe3+ in 10 mM PBS (7.4 pH) as a redox probe. CV measurements were made in the range of -0.5 to 0.5 V. The cycle was repeated twice for each sample. The scan rate was 0.1 V/s and the step potential was 2.44 mV.

III. RESULTS AND DISCUSSIONS

The initial objective was to confirm the immobilization of Au NPs on the electrodes and its impact on glucose sensing. Impedance spectroscopy and cyclic voltammetry measurements using Fe2+/Fe3+ as a redox probe revealed distinct changes in electrode behavior after modification (Fig. 1). The analysis revealed that the immobilization of P-lysine caused a noticeable decrease in impedance. This phenomenon can be attributed to the presence of amino groups, characterized by a positive



c=0	c=0.1 c=0.5
c=0.75	c=1 c=1.25
c=1.5	— c=1.75 — c=3
c=5	c=10

Fig. 2. Bode plots of Ref without NPs (a), with 20 nm Au NPs and as insert image 100 nm Au NPs (b). The 2x (c) and 4x (d) higher concentration of Au NPs for 20 nm Au NPs and showed as insert image

charge, which effectively attracted the negatively charged redox probe. Consequently, the conductivity near the electrode surface increased, leading to a reduction in impedance. [15]. The immobilization of Au NPs led to evident increase in impedance, primarily attributed to the occupation of positive amino groups from p-Lys by negatively charged Au NPs originating from citrate groups. This newly formed layer of nanoparticles on the electrode created a significantly larger conductive surface. Theoretically, the overall increase of surface area could be similar for both larger and smaller particles. However, the crucial factor lies in the quantity of particles and their ability to attach to the positive charge of Plysine. The smaller 20 nm Au NPs caused a more substantial increase in impedance compared to their 100 nm Au NPs counterparts. This suggests a higher concentration of smaller particles with a negative charge, implying a potentially greater surface area for enzyme immobilization. Subsequent measurements involved samples where electrodes were exposed only to a solution from which nanoparticles had been taken out by centrifugation. This aimed to test whether the impedance curve changes were solely due to alterations in pH or other reactions with the solvent. The behaviour of these electrodes markedly differed from those with NPs, providing substantial evidence of the changes induced by Au NPs immobilization. Additional surface examination methods, such as electron microscopy and profilometry, were employed. However, due to the small particle size and the combination of gold on gold, the results from these methods were inconclusive and are not presented here. Cyclic voltammetry (CV) measurements (Fig. 1b) exhibited similar tendencies in peak shifting, reinforcing the observed changes in electrode behaviour due to NP immobilization.

Subsequently, electrodes with a layer of 20 nm and 100 nm NPs covered with a film of the GOD enzyme, immobilized through the cross-linking method, were measured. To establish a baseline for behaviour without NPs, GOD cross-linked on P-lys without NPs served as a reference (Ref without NPs). Additionally, protein BSA cross-linked on NPs without GOD enzyme was employed to isolate the effects caused by the enzyme, essentially measuring the background response.

Fig. 2 illustrates the Bode plots, showing a notable decrease in impedance, particularly at lower frequencies of alternating voltage, with increasing glucose concentration. From this region, calibration curves were constructed, as illustrated in Fig. 3a for a frequency of 2 Hz. Lower frequencies were

essential to allow the manifestation of the transfer of massive ions, and therefore are better for crating calibration curves. To

TABLE I. This table shows the calibration area, parameters from the calibration curves, and LOD for the R_{ct} and impedance measurements conducted at frequency 2 Hz. The Values of R_{ct} have been logarithmically transformed to achieve linear representation.

Sample	R _{ct}			Z for 2 Hz			
	Calibration area [mM]	a [log(Ω).mM ⁻¹]	LOD [mM]	Calibration area [mM]	a [MΩmM ⁻¹]	LOD [mM]	LOD data er. [mM]
Ref – no NPs	0.1 – 1.9	-1.22 ± 0.09	0.24	0.01-1.6	-0.50 ± 0.03	0.18	0.10
Au NPs 100	0.5 – 2	-1.2 ± 0.1	0.51	0.1-1.7	-0.35 ± 0.03	0.31	0.14
2x Au NPs 100	0.5 - 1.75	-5.7 ± 1.5	0.89	0.1-2	-0.34 ± 0.06	0.56	0.15
4x Au NPs 100	0.5 - 1.75	-1.4 ± 0.3	0.60	0.1-1.7	-0.32 ± 0.02	0.22	0.15
Au NPs 20	0.5 – 1.75	-4.9 ± 1.7	1.1	0.1-1.8	-0.32 ± 0.03	0.33	0.16
2x Au NPs 20	0.1 – 2	-1.33 ± 0.09	0.23	0.1-1.5	-0.43 ± 0.05	0.42	0.12
4x Au NPs 20	0.1 - 2	-1.3 ± 0.1	0.29	0.01-1.5	-0.66 ± 0.07	0.33	0.08



Fig. 3. Graph (a) shows impedance measurement for 2 Hz voltage frequency. Graph (b) shows R_{ct} parameter dependence on glucose concentration. The values were normalized by maximal R_{ct} value for given type of sample. R_{ct} was calculated by fitting Randles equation corresponding to showed circuit (b). 100 nm and 20 nm Au NPs represents diameters and x2/x4 is the multiple of the particle concentration.

further illustrate the electrode behaviour, Nyquist graphs would reveal a gradual formation of more pronounced semicircles with increasing glucose concentrations. For lower concentrations, a distinct diffusional tail would be prominent.

The decline in impedance observed on the electrode surface corresponds to the ongoing reaction. The rise in current is a direct result of the re-oxidation of Ru2+ to Ru3+. Thanks to reaction of glucose oxidation catalysed by glucose oxidase (GOD) the reduction from Ru3+ to Ru2+ is occurring and therefore the free electron can be created for increasement of measured current. To delve deeper into the modifications in

electrochemical behaviour, we focused on the charge transfer resistance (R_{ct}). This parameter signifies the difficulty encountered when an electron undergoes transfer from one atom

or compound to another (Fig. 3b). The values for R_{ct} , along with their standard deviation (σ_c), were determined through data fitting using the Randles equation. For both 20 nm and 100 nm Au NPs, we conducted measurements with concentrations that the deviations of individual measurements. Reference samples without GOD did not have a calibration curve and were consequently not included in the table. Generally, regarding our sensors with 20 nm Au NPs, they tended to behave similarly to the reference sample without NPs. In contrast, those with Au NPs exhibited a smaller decrease in impedance but a larger drop in R_{ct} . This suggests that these gold particles might be causing a hindrance rather than accelerating the electron transport to the electrodes. However, since R_{ct} remains low, we can assume that it serves more as a capacitor function rather than a distinct barrier. For practical applications of these biosensors, consideration of the sensor's behaviour is crucial when selecting an appropriate processing form. Ideally, the chosen form should yield an optimal curve and a lower LOD. Generally, calibration from impedance provides a smaller deviation and thus a better LOD, while calibration from R_{ct} can capture a broader concentration range.

IV. CONCLUSION

It is commonly hypothesized that a larger surface area allows for the immobilization of a greater number of enzymes, potentially enhancing sensitivity. This study investigated this hypothesis using two different-sized Au NPs to monitor their impact on biosensor behaviour. The method of packing spherical NPs in a dense planar arrangement was employed for immobilization, known for its rapid and simple application without the need for sophisticated equipment. It was revealed that the nanoparticle layer significantly influenced the biosensor response in impedance measurement. Larger particles, specifically 100 nm Au NPs, tended to cause blocking, enhancing capacitor behaviour, while 20 nm Au NPs slightly increased electron flow, potentially due to a higher amount of immobilized enzyme available for reaction. However, this effect became noticeable only at an elevated concentration of 20 nm Au NPs. The Limit of Detection (LOD) was lower for calibrations derived from impedance curves. However, we must note that, compared to our reference, the addition of Au NPs rather caused a decrease in LOD, most likely due to further blocking. A decrease in response was also observed in [14]. Further optimization, such different NPs diameter, material or exploring alternative enzyme immobilization techniques, may yield improvements in the calibration range. Continued investigation into the properties of nanopore biosensors could prove valuable, particularly in the development of "point of care" devices.

ACKNOWLEDGMENT

This article was supported by the Czech Academy of Sciences (RVO:68081731) and The Technology Agency of the Czech Republic FW03010504. We acknowledge CzechNanoLab Research Infrastructure supported by The Ministry of Education, Youth and Sports of the Czech Republic (LM2018110), and the project FEKT-S-23-8162 and CEITEC VUT/FEKT-J-24-8567.

References

[1] A. Santos, T. Kumeria, and D. Losic, "Nanoporous anodic aluminum

oxide for chemical sensing and biosensors," *TrAC - Trends in Analytical Chemistry*, vol. 44. Elsevier B.V., pp. 25–38, Mar. 01, 2013. doi: 10.1016/j.trac.2012.11.007.

- [2] Q. Chen and Z. Liu, "Fabrication and Applications of Solid-State Nanopores," *Sensors*, vol. 19, no. 8, p. 1886, Apr. 2019, doi: 10.3390/s19081886.
- [3] S. Manzoor, M. W. Ashraf, S. Tayyaba, and M. K. Hossain, "Recent progress of fabrication, characterization, and applications of anodic aluminum oxide (AAO) membrane: A review," Dec. 2021, Accessed: Mar. 02, 2022. [Online]. Available: http://arxiv.org/abs/2112.08450
- [4] J. Sopoušek, J. Věžník, P. Skládal, and K. Lacina, "Blocking the Nanopores in a Layer of Nonconductive Nanoparticles: Dominant Effects Therein and Challenges for Electrochemical Impedimetric Biosensing," ACS Appl. Mater. Interfaces, vol. 12, no. 12, pp. 14620– 14628, Mar. 2020, doi: 10.1021/acsami.0c02650.
- [5] A. de la Escosura-Muñiz, M. Espinoza-Castañeda, M. Hasegawa, L. Philippe, and A. Merkoçi, "Nanoparticles-based nanochannels assembled on a plastic flexible substrate for label-free immunosensing," *Nano Res.*, vol. 8, no. 4, pp. 1180–1188, Apr. 2015, doi: 10.1007/s12274-014-0598-5.
- [6] P. Jiang, Y. Wang, L. Zhao, C. Ji, D. Chen, and L. Nie, "Applications of Gold Nanoparticles in Non-Optical Biosensors," *Nanomater.* 2018, Vol. 8, Page 977, vol. 8, no. 12, p. 977, Nov. 2018, doi: 10.3390/NANO8120977.
- Z. Hua, T. Yu, D. Liu, and Y. Xianyu, "Recent advances in gold nanoparticles-based biosensors for food safety detection," *Biosens. Bioelectron.*, vol. 179, p. 113076, May 2021, doi: 10.1016/J.BIOS.2021.113076.
- [8] C. Shan, H. Yang, D. Han, Q. Zhang, A. Ivaska, and L. Niu, "Graphene/AuNPs/chitosan nanocomposites film for glucose biosensing," *Biosens. Bioelectron.*, vol. 25, no. 5, pp. 1070–1074, Jan. 2010, doi: 10.1016/J.BIOS.2009.09.024.
- [9] I. S. Kucherenko, O. O. Soldatkin, D. Y. Kucherenko, O. V. Soldatkina, and S. V. Dzyadevych, "Advances in nanomaterial application in enzyme-based electrochemical biosensors: a review," *Nanoscale Advances*, vol. 1, no. 12. Royal Society of Chemistry, pp. 4560–4577, Dec. 03, 2019. doi: 10.1039/c9na00491b.
- [10] H. Ilkhani, M. Sarparast, A. Noori, S. Z. Bathaie, and M. F. Mousavi, "Electrochemical aptamer/antibody based sandwich immunosensor for the detection of EGFR, a cancer biomarker, using gold nanoparticles as a signaling probe," *Biosens. Bioelectron.*, vol. 74, pp. 491–497, Dec. 2015, doi: 10.1016/J.BIOS.2015.06.063.
- [11] S. Zhang, N. Wang, Y. Niu, and C. Sun, "Immobilization of glucose oxidase on gold nanoparticles modified Au electrode for the construction of biosensor," *Sensors Actuators B Chem.*, vol. 109, no. 2, pp. 367–374, Sep. 2005, doi: 10.1016/J.SNB.2005.01.003.
- [12] T. Zhang, J. Ran, C. Ma, and B. Yang, "A Universal Approach to Enhance Glucose Biosensor Performance by Building Blocks of Au Nanoparticles," *Adv. Mater. Interfaces*, vol. 7, no. 12, p. 2000227, Jun. 2020, doi: 10.1002/ADMI.202000227.
- S. Zhang, N. Wang, H. Yu, Y. Niu, and C. Sun, "Covalent attachment of glucose oxidase to an Au electrode modified with gold nanoparticles for use as glucose biosensor," *Bioelectrochemistry*, vol. 67, no. 1, pp. 15–22, Sep. 2005, doi: 10.1016/J.BIOELECHEM.2004.12.002.
- [14] G. A. Valencia, L. C. De Oliveira Vercik, and A. Vercik, "A new conductometric biosensor based on horseradish peroxidase immobilized on chitosan and chitosan/gold nanoparticle films," *J. Polym. Eng.*, vol. 34, no. 7, pp. 633–638, Sep. 2014, doi: 10.1515/POLYENG-2014-
 - 0072/MACHINEREADABLECITATION/RIS.
- [15] K. Lacina, J. Sopoušek, V. Čunderlová, A. Hlaváček, T. Václavek, and V. Lacinová, "Biosensing based on electrochemical impedance spectroscopy: Influence of the often-ignored molecular charge," *Electrochem. commun.*, vol. 93, pp. 183–186, Aug. 2018, doi: 10.1016/j.elecom.2018.07.015.

Deployment of deep learning-based anomaly detection systems: challenges and solutions

Stepan Jezek Dept. of telecommunications, FEEC Brno University of Technology Brno, Czech Republic xjezek16@vutbr.cz

Abstract—Visual anomaly detection systems play an important role in various domains, including surveillance, industrial quality control, and medical imaging. However, the deployment of such systems presents significant challenges due to a wide range of possible scene setups with varying number of devices and high computational requirements of deep learning algorithms. This research paper investigates the challenges encountered during the deployment of visual anomaly detection systems for industrial applications and proposes solutions to address them effectively. We present a model use case scenario from real-world manufacturing quality control and propose an efficient distributed system for deployment of the defect detection methods in manufacturing facilities. The proposed solution aims to provide a general framework for deploying visual defect detection algorithms base on deep neural networks and their high computational requirements. Additionally, the paper addresses challenges related the whole process of automated quality control, which can be performed with varying number of camera devices and it mostly requires interaction with other factory services or workers themselves. We believe the presented framework can contribute to more widespread use of deep learning-based defect detection systems, which may provide valuable feedback for further research and development.

Index Terms—deep learning, defect detection, system design, algorithm deployment, image processing, distributed systems

I. INTRODUCTION

Visual anomaly detection systems have garnered significant attention in recent years owing to their potential to enhance safety, security, and efficiency across various domains. These systems utilize advanced computer vision techniques to automatically identify deviations from normal patterns or behaviors within visual data, enabling early detection of abnormalities in diverse environments such as surveillance videos, industrial processes, and medical imaging [1], [2].

The general problem of anomaly detection is based on labeling data samples that deviate from a defined normal state [6]. The problem of visual industrial anomaly detection is mostly focused on detecting defects in manufactured products. Examples of such defects can be seen in figure 1. Recently, many new defect detection methods have been introduced, mostly exploiting current advancements in deep learning algorithms. While considerable progress has been made in the development of anomaly detection algorithms and models, their effective deployment in real-world scenarios presents 2nd Radim Burget Dept. of telecommunications, FEEC Brno University of Technology Brno, Czech Republic burgetrm@vutbr.cz



Fig. 1. Examples of data samples from industrial defect detection datasets. Normal training samples are displayed in the first row, testing samples with highlighted anomalies are displayed in the second row.

various challenges that have not been sufficiently addressed in the current literature.

The deployment of visual anomaly detection systems entails navigating through a wide range of technical, operational, and practical challenges [15]. The visual anomaly detection process mostly present challenges that stem from the complex nature of real-world environments, which often exhibit dynamic and unpredictable conditions. Factors such as lighting variations, occlusions, noise, and object diversity pose significant hurdles to the reliable operation of anomaly detection systems. Deployment of such deep learning-based algorithms needs to address considerations related to computational resources, scalability, and integration with existing infrastructure [3].

Addressing these challenges requires a holistic approach that encompasses not only algorithmic advancements but also considerations of system design, data preprocessing, model optimization, and deployment strategies. Additionally, the deployment of visual anomaly detection systems necessitates a thorough understanding of the specific requirements and constraints of the target application domain. Real-world deployment scenarios often demand robustness, adaptability, and efficiency to ensure effective anomaly detection in diverse operational settings. Moreover, the cost efficiency of such systems also needs to be addressed in order to enable widespread adoption [15].

In light of these challenges and opportunities, this research



Fig. 2. Examples of defect localization maps from industrial defect detection datasets.

paper aims to explore the intricacies of deploying visual anomaly detection systems in practical settings. We aim to identify key challenges encountered during deployment and propose effective solutions based on distributed microservices architecture that can provide computatinally intensive deep learning based defect detection service to many data collection devices at the same time. We also take into account the user interface of the system that enables workers to use the system without high expertise demands.

The main contribution of the paper can be summarized as follows:

- Identification of Deployment Challenges: We examine challenges encountered in deploying visual anomaly detection systems, including data load issues, computational constraints, and integration considerations.
- **Proposal and Validation of Deployment Strategies:** This paper proposes effective system design solution for deployment of deep learning based defect detection algorithms in manufacturing. We also validate the solution by implementing the prototype of the designed system using modern web and distributed microservices technologies.

The rest of the paper is structured as follows: in the next section, we present current trends in visual anomaly detection and provide an overview of recent literature focused on the deployment of visual deep learning algorithms. The next section describes the model use case of a visual defect detection system. Next, we present a proposed system design and implementation approaches for the defect detection system. In the last section, we discuss the implementation results.

II. RELATED WORK

Traditional defect detection methods often rely on manual inspection or rule-based algorithms, which are labor-intensive, time-consuming, and may lack robustness in handling complex defect patterns. In recent years, deep learning-based approaches have shown promising results in automating defect detection tasks by leveraging the power of convolutional neural networks (CNNs) to learn discriminative features directly from raw data [6]. The main task of the defect detection algorithms usually comprises of detection at the image level and also defect localization which involves creating a defect localization maps as illustrated in the figure 2. Accuracy of the methods is most often evaluated using standard AUROC metric [19], either on the image or the pixel level depending on the image level defect detection or localization task.

In the realm of visual anomaly detection, two prominent paradigms have emerged: reconstruction-based methods and feature extraction-based approaches. Reconstruction-based techniques, rooted mainly in autoencoder architectures, aim to learn a compact representation of normal data and subsequently reconstruct input samples. Anomalies are identified by measuring the reconstruction error, with higher errors indicating potential defects. Examples of recent reconstruction-based methods include Skip-GANomaly [4] or DRAEM [5].

Feature extraction-based methods focus on extracting discriminative features from images using convolutional neural networks (CNNs). These features are then fed into traditional machine learning algorithms, such as k-NN, k-Means or oneclass SVMs, to classify anomalies. According to the recent literature and state of the art, feature extraction-based methods show the highest accuracy on common defect detection benchmarks such as MVTec-AD [8], Kolektor SDD [9], Magnetic Tiles Dataset [10] or MPDD [11]. Recent state of the art feature extraction-based methods for visual anomaly detection include CFlow-AD [13], PatchCore [14] or CFA [12].

With recent advancements in state of the art methods and algorithms, the field of machine learning has recently seen a wide adoption in many practical applications. Therefore, the problems related to deployment of the methods and algorithms in practice was given more focus in the community. General problems related to deployment of machine learning-based algorithms in industrial aplications is published in [15]. It describes common machine learning deployment tasks such as data management, model integration and also law considerations. Other deployment proposals and solutions are presented in [16] and [17]. Even though the above-mentioned publications offer a comprehensive description of the deployment topic, they may lack sufficient implementation details for realworld applications.

Our research is inspired by the publication [18], which introduces a thorough design of a deep learning-based system for Covid-19 detection from X-Ray scans. It describes all system components that ensure modular deployment of different machine learning algorithms, high availability to clients and integration with other service. It is however closely related to the medical use case that involves several system components not present in the industrial settings.

III. PROPOSED SYSTEM DESIGN

In this section, we describe general requirements for deployment of defect detection systems in real-world applications, including the general manufacturing workflows. We delve into common requirements for the detection systems and describe the resulting technical necessities in terms of computation, user interactions and communications with other services.

Next, we describe a general technical framework for the detection system that leverages a distributed architecture for



Fig. 3. Design schematics of a general industrial defect detection system. The three main parts of the system are - data collection, interface and processing.

efficient and scalable processing of visual data from multiple camera devices. The system is based on the distributed microservices architecture that interconnects serveral components, each serving a specific role in the defect detection workflow.

A. Industrial Defect Detection

Manufacturing in general is dependant on delivering high quality products with minimal amount of defects. Real-world industrial production therefore usually deploy a quality control processes to ensure defective products are not shipped to customers. Quality control is typically based on visual evaluation of product features. These tasks are usually performed by specialized workers and recently also deep learning and computer vision. This commonly involves cameras placed on the production line that captures the products during different manufacturing phases.

Image data from the cameras are then transmitted to the central server that provides the detection algorithms services. The server also needs to provide a user interface to workers or other responsible personnel. After processing the captured images the main part of the system is detection of defects via deep learning algorithms. This task can be performed on the same server, but because of high computational requirements of deep neural networks, this approach is not usually suitable with intensive data streams. The deep learning algorithms therefore need to run on a dedicated servers with hardware acceleration (usually GPU cards). Detection results are saved to the database to ensure persistence and history tracking and are the presented to the users via specified interface (e.g. web browser). The general diagram of the described system is shown in the figure 3.

B. General Defect Detection System

Considering the workflow and requirements described in the previous section, we propose a general model of a defect detection system that fits a typical industrial production. The system is divided into three main parts: data collection section, interface section and processing section.

The data collection section is composed of varying number of camera devices. In real-world applications, the devices are usually industrial grade cameras with high reliability and endurance certifications. The output image size is usually in the range of 1 to 20 megapixels and the framerate between 1 and 30 frames per second. This data load requirements are addressed in further sections.

Data from the camera devices is sent to the server node that provides interface for forwarding messages into the main processing server. This forwarding may include image dropping if the system does not need high framerates. The server then uses the camera data to create a request to the main processing server. Requests are forwarded to the message queue broker, which acts as a middleman between the interface and the processing server.

The processing server then subscribes to the messages in the request queue and performs preprocessing and defect detection on the input data in the same order as the requests were pushed to the queue. This approach enables the user interface service to pass the computationally intensive tasks to another machine and still be able to process other user requests. After processing a request, the server pushes the results to a separate queue in order to immediately start with processing other requests. The results queue is processed back by the server forwarding messages for processing. This server will save the results to the database server and also provide a user web interface for presentation of the results to the clients (e.g. quality assurance workers).

IV. EXPERIMENT

Based on the general design described in the previous section, here we present a prototype implementation of the visual defect detection system. You can see the main system building blocks in the figure 4.



Fig. 4. System implementation diagram. It shows all service components of the resulting system architecture.

The system is expected to handle a network of camera devices using the MQTT communication protocol, image data transfer is performed via Eclipse Mosquitto MQTT broker [20]. We expect a system with possibly large number of low-powered camera devices, for which this protocol can be efficiently exploited.

A. Main Web Interface

The core functions for user interface and communication with processing services are implemented using the Django web application framework [21]. This server provides web user interface that is responsible for a management of camera devices, presenting results and forwarding the requests to the processing server (see figure 5). The Django server acts as a queue requests producer for image inference and also processes the final results. The results are consumed from a separate process that receives results from a dedicated results queue and sends the data to the Django web server via REST application interface endpoint [7]. The web server then saves the results to a PostgreSQL database server [22].

B. Processing Queue

The main web interface forwards the defect detection requests to the queue that is implemented using the RabbitMQ message broker [23]. The RabbitMQ can be run as a service on a separate server, it provides interface for creating parallel message queues that can either be used by producers to push new messages (requests) or consumers to process the requests. We define separate message queue for defect detection inference requests and detection results, that are consumed by the main web interface.

C. Main Processing Server

The main server for processing defect detection tasks is implemented as a RabbitMQ message consumer, which loads the deep learning model during intialization phase, registers to the inference and results queues and the waits for incoming requests from the queue. We can run this process on an arbitrary number of servers, the RabbitMQ broker will distribute the load on all of them. Using this approach, we are not restricted to a specific deep learning algorithm or framework. For our testing purposes, we deployed the PatchCore [14] algorithm combined with the background segmentation model GroundingDino [24], both implemented with the PyTorch framework.

D. Database Storage

Once the defect detection analysis is completed, the results are stored in the PostgreSQL database server and presented to the user via the web interface. The interface uses a Websocket protocol to send the results to the web browser without the need of explicit user actions or periodic polling.

V. RESULTS AND DISCUSSION

We tested the system with setup that involved dedicated virtual private servers for the main web server, Mosquitto MQTT broker and the RabbitMQ broker. We compared our method against the simple architecture that involved wrapping the deep learning model inside a single application with the REST interface for sending the inference requests. We used a computationally intensive algorithm combining defect detection with background segmentation as mentioned previously. The processing time on our workstation (nvidia RTX 3070 GPU) per image was 1.9 seconds in average. Using our proposed architecture, we were easily able to double the number of processing time of 1.02 second per image (see table I).

 TABLE I

 Results of the tested deployment architectures

Architecture	Throughput	System scalability
REST API wrapper	1.9 s/img.	vertical
Our distributed system	1.02 s/img.	vertical, horizontal

VI. CONCLUSION

In this paper, we focused on challenges related to the scalable deployment of deep learning-based industrial defect detection algorithms. We described current highlights in current visual defect detection methods and the related literature



Fig. 5. Main web server user interface. The image shows the defect detection inference section that displays the data for a selected camera device and also defect detection results that includes the defect localization maps and a image level detection result.

focused on deploying the detection systems in real-world applications. Compared to the research in new defect detection methods, the problem of deploying the methods has not been widely researched in recent publications.

We proposed a unique architecture for deploying the computationally intensive deep learning algorithms using the distributed microservices design principles. We tested the proposed system by implementing a functional prototype using the Django web framework, Mosquitto MQTT broker and the RabbitMQ message queue broker. We tested our solution against a the commonly used REST API model wrapper and achieved significant improvements by 46.3 % in terms of processing time using our proposed distributed architecture.

REFERENCES

- V. Chandola, A. Banerjee and V. Kumar. "Anomaly detection: A survey," ACM computing surveys (CSUR), 2009.
- [2] G. Pang et al. "Deep learning for anomaly detection: A review," ACM Computing Surveys (CSUR), vol. 54, n. 2, 2021, pp. 1-38
- [3] Yang, R. Xu, Z. Qi, and Y. Shi, "Visual Anomaly Detection for Images: A Systematic Survey," Procedia Computer Science, vol. 199, pp. 471–478, 2022, doi: https://doi.org/10.1016/j.procs.2022.01.057.
- [4] S. Akçay, A. Atapour-Abarghouei and T. P. Breckon, "Skip-ganomaly: Skip connected and adversarially trained encoder-decoder anomaly detection," in 2019 International Joint Conference on Neural Networks (IJCNN), 2019, pp. 1–8.
- [5] V. Zavrtanik, M. Kristan and D. Skocaj, "DRÆM A discriminatively trained reconstruction embedding for surface anomaly detection," in 2021 *IEEE/CVF International Conference on Computer Vision (ICCV)*, Montreal, QC, Canada, 2021 pp. 8310-8319.
- [6] I. Goodfellow, Y. Bengio, and A. Courville, Deep Learning. MIT Press, 2016.
- [7] R. T. Fielding et al., "Reflections on the REST architectural style and 'principled design of the modern web architecture' (impact paper award)," in Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering, 2017, pp. 4–14. doi: 10.1145/3106237.3121282.
- [8] P. Bergmann, M. Fauser, D. Sattleger and C. Steger, "Mvtec ad a comprehensive real-world dataset for unsupervised anomaly detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 9592–9600.
- [9] D. Tabernik, S. Šela, J. Skvarč, and D. Skočaj, "Segmentation-Based Deep-Learning Approach for Surface-Defect Detection," in *Journal of Intelligent Manufacturing*, May 2019.

- [10] Y. Huang, C. Qiu, Y. Guo, X. Wang and K. Yuan, "Surface Defect Saliency of Magnetic Tile," in 2018 IEEE 14th International Conference on Automation Science and Engineering (CASE), 2018, pp. 612-617, doi: 10.1109/COASE.2018.8560423.
- [11] S. Jezek, M. Jonak, R. Burget, P. Dvorak and M. Skotak, "Deep learningbased defect detection of metal parts: evaluating current methods in complex conditions," 2021 13th ICUMT, Brno, Czech Republic, 2021, pp. 66-71, doi: 10.1109/ICUMT54235.2021.9631567.
- [12] S. Lee, S. Lee and B. Song, "CFA: Coupled-hypersphere-based Feature Adaptation for Target-Oriented Anomaly Localization", arXiv.org, 2022. [Online].
- [13] D. Gudovskiy, S. Ishizaka and K. Kozuka, "CFLOW-AD: Real-Time Unsupervised Anomaly Detection with Localization via Conditional Normalizing Flows," in arXiv preprint arXiv:2107.12571, 2021.
- [14] K. Roth and L. Pemula, J. Zepeda, B. Schölkopf, T. Brox and P. Gehler, "Towards Total Recall in Industrial Anomaly Detection," in arXiv:2106.08265, 2021.
- [15] A. Paleyes, R.-G. Urma, and N. D. Lawrence, "Challenges in Deploying Machine Learning: A Survey of Case Studies," ACM Comput. Surv., vol. 55, no. 6, Dec. 2022, doi: 10.1145/3533378.
- [16] A. Luckow, M. Cook, N. Ashcraft, E. Weill, E. Djerekarov and B. Vorster, "Deep learning in the automotive industry: Applications and tools," 2016 IEEE International Conference on Big Data (Big Data), Washington, DC, USA, 2016, pp. 3759-3768, doi: 10.1109/Big-Data.2016.7841045.
- [17] H. Heymann, A. D. Kies, M. Frye, R. H. Schmitt, and A. Boza, "Guideline for Deployment of Machine Learning Models for Predictive Quality in Production," Procedia CIRP, vol. 107, pp. 815–820, 2022, doi: https://doi.org/10.1016/j.procir.2022.05.068.
- [18] V. Myska et al., "CovidStopHospital: e-Health Service for X-Ray-Based COVID-19 Classification and Radiologist-Assisted Dataset Creation," 2023 15th ICUMT, Ghent, Belgium, 2023, pp. 62-67, doi: 10.1109/ICUMT61075.2023.10333292.
- [19] T. Fawcett, "ROC graphs: Notes and practical considerations for researchers," in *Machine learning*, vol. 31(1), pp. 1-38, 2004.
- [20] Eclipse mosquitto. Eclipse Mosquitto. (2018, January 8). https://mosquitto.org/
- [21] Django. Django Project. (n.d.). https://www.djangoproject.com/
- [22] Group, P. G. D. (2024, March 8). PostgreSQL. https://www.postgresql.org/
- [23] Rabbitmq: One broker to queue them all: Rabbitmq. RabbitMQ RSS. (n.d.). https://www.rabbitmq.com/
- [24] Shilong Liu, undefined., et al, "Grounding DINO: Marrying DINO with Grounded Pre-Training for Open-Set Object Detection," 2023.

Amplitude measurement of small displacement using video magnification

1st Dominik Ricanek

dept. of Control and Instrumentation Faculty of Electrical Engineering and Communication, Brno University of Technology, Czechia ORCID: 0000-0001-5031-2481

Abstract—Video magnification algorithms show promising results when used to amplify small vibrations. Measuring these small vibrations is integral in modal analysis of an object and is usually done using specialized vibrometers or accelerometers. Computer vision (CV) systems fall short in this task as the magnitude of the vibrations decreases because of a small SNR (Signal-to-Noise Ratio). In this paper we try to further improve the accuracy of the CV approach by adding video magnification into the image pre-processing stage, allowing the algorithms to measure even imperceptibly small vibrations. For this purpose, we have gather ground truth data with a laser vibrometer in tandem with high speed camera footage of a metal cantilever beam, vibrating in its first mode, and have trained a convolutional neural network for regression.

Index Terms-EVM, phase-based, optical flow, PDV, ResNet18

I. INTRODUCTION

Camera based systems (CBS) remain underutilized in modal analysis of complex objects where transducers and laser vibrometers surpass them in both spatial resolution and frequency range. However, recent advancements in video magnification, specifically the phase-based approach by [12], allow us to improve the spatial resolution of cameras dramatically. This combined with the other benefits of CBS, such as the ease of use and remote sensing, can help them set up a niche as a fast modal shape analysis tool, used to identify places on an object where more precise sensors could be placed.

Using a laser vibrometer, we can measure surface vibrations of the object by scanning. However, this limits us to measurements parallel to the laser's plane and requires expensive equipment. In CBS the entire image plane is used to measure vibrations and, depending on the requirements, the equipment can be affordable. The disadvantage lies in frequency range as a high camera frame rate is naturally going to be more difficult to achieve than high measuring frequency of an accelerometer.

The idea to use cameras to amplify invisible motion dates back to [7] and their Lagrangian magnification. Since then major advancements have been made, starting with the invention of Eulerian magnification by [14] which was then further improved by [12] by using phase differences to allow for greater magnification factors without affecting the SNR. The names of these two algorithms stem from fluid mechanics and allude to their different perspectives; whereas Lagrangian magnification focuses on the motion of individual pixels - as one would particles within a fluid - the Eulerian perspective takes into account the motion of all particles globally through the use of spacial filters.

The method has been experimentally proven to magnify real small motion by [2], where the team found the lowest bound of sub-pixel motion to be 0.00025 pixels. The lower bound can be further improved by using specialized image sensors with greater bit depths or higher contrast.

Attempts have also been made to incorporate phase-based Eulerian video magnification (EVM) into edge devices by implementing it in C++ [3]. These are usually still run on a laptop computer, however. A truly edge computing application is used for example for baby breathing monitor in [5]. Phase-based video amplification nevertheless remains highly computationally demanding [9] and sensitive to larger motion within the picture.

II. COMPONENTS

Phase-based EVM can be thought of as a combination of three smaller algorithms. That is:

A. Phase-based optical flow

Optical flow is a tool used in computer vision to gain motion vectors of objects within analyzed video. The phasebased EVM does this by looking at phase differences between two frames, rather than intensity differences, as in the seminal algorithm.

The method itself has been around for over two decades but remains sparsely used relative to the other computer vision algorithms [2]. Apart from EVM, it is used for example in autonomous navigation, movement detection and tracking.

B. Image pyramid

An image pyramid decomposition is an over-complete transformation which results in a sequence of progressively downscaled versions of the original image called "levels" as can be seen in Fig. 1. When applied to a video, decomposition is performed on each frame individually. This requirement is the main cause behind the high computational demands of video magnifying algorithms. Indeed one of the major advancements in video magnification was the introduction of Riesz pyramids

The completion of this paper was made possible by the grant No. FEKT-S-23-8451 - "Research on advanced methods and technologies in cybernetics, robotics, artificial intelligence, automation, and measurement" financially supported by the Internal science fund of Brno University of Technology.

as directional filters by [13], which led to a major increase in processing speed.

In our implementation we choose the usual approach of down-scaling by a factor of two between each level and using Gaussian filter to pool the pixels. This way we are left with multiple levels of spatially filtered images which helps us magnify motion beyond just the sharpest edges.



Fig. 1. A general image pyramid with n-levels.

C. Complex-steerable filters

By filters here we mean 2D convolution masks; complex because they are applied in the Fourier domain and thus as complex numbers; and steerable because their response is unaffected by the direction they are applied in - we use eight directional filters derived from a common complex-steerable filter bank as seen in Fig. 2.

In essence these function as edge detectors applied on each transformed frame of the video e.i. each pyramid. For this reason the filters themselves have to be pyramid form. These are sometimes called Riesz pyramids for their special ability to create orthogonal quadrature pairs upon application and thus reducing the number of required directional filters by half notice the lack of opposing directional filters in our example.

Another useful property of complex-steerable filters are their non-aliased sub-bands thanks to which we can then spatially move the filter response without artefacts. It should be noted that steer-ability is not exploited. [12] [10]

III. PHASE BASED VIDEO MAGNIFICATION

One of the consequences of the Fourier theorem is that signals undergoing translation will have some of their components change phase and vice versa - changes in phase of a component sine wave correlate to motion of the signal. This can be easily imagined on a 1D example in Fig. 4 where a wavelet moves in the x axis based on the changes to its phase. This process was shown to have promise in video processing application and the team behind [14] utilized it to avoid the computationally demanding calculation of motion vectors of other video magnifying methods and instead directly manipulate the phase of a video through complex-steerable filters.



Fig. 2. The complex steerable filter (above) and its bank of eight directional filters.

In practical application we will start by shooting a video with the Nyquist–Shannon theorem in mind. That is at twice, or more, the frame rate of the expected frequency. Note that methods exist to remedy distortion and aliasing caused by capturing at lower frame rates but those were out of scope of this work.

The video is then then split into individual frames and each frame is transformed into an image pyramid using a Gaussian filter. The complex-steerable filters, which will be applied in the next step, are also transformed. However, since the filter bank does not change, these transformations can be done beforehand to reduce processing time. These filters are also already in the Fourier domain in order to utilize the Convolution theorem and further save computational time.

Each frame is transformed into the Fourier domain, centered, and multiplied by each of the eight complex-steerable filters. The result of this process is a new image pyramid composed of complex matrices from which we can calculate the phase as the angle between the real and imaginary parts.

The first frame of the video is chosen as the reference. It is transformed as stated above and used to calculate the difference in phase between all the other frames. These differences are then filtered temporally to band-pass the desired frequencies and multiplied by the amplification constant α (where $\alpha = 1$ corresponds to 100% motion amplification).

The amplified image pyramids are then collapsed and added to the original frames. This way we end up with a synthetic video with amplified motion, without the need to track movements like other video magnification techniques do [12].

IV. RELATED RESEARCH

Eulerian Video Magnification has progressed a long way since its debut at MIT in 2012 [14] and nowadays these sorts of algorithms achieve accurate real time amplification of tiny



Fig. 3. Full processing chain of the phase video amplification algorithm. Left to right: original image, pyramid decomposition, convolution with complexsteerable filters, phase difference multiplication, temporal filtering, collapsing the pyramid and adding to original frame.



Fig. 4. Translation of a wavelet by changes in phase of its components [11].

motions - be it periodic vibrations or short impulses - in a video. Using spatial phase, it is possible to highly negate the effect of noise within the image on the magnified result [12] and attenuate unwanted movements. This can and has been used to ameliorate the fatal flaw of EVM, large motion [17]. Computer vision systems have been used for precise modal analysis of complex structures such a three storied model of a building [16]. Efforts have also been made to transfer the EVM algorithm into the realm of machine learning, with mixed results [8]. However, to our best knowledge, nobody has yet attempted to use phase-based EVM algorithm to aid a machine learning models in determining the amplitude of the detected vibrations.

When it comes to vibrometry using regression CNNs, the idea has been extensively explored. In [6] it has been used to detect anomalies in rotating machinery and predict bearing faults with an RMSE of 0.97, the authors use 1D convolution of an attached accelerometer in combination with a gated recurrent unit (GRU). Although somewhat unrelated, we mention this here to point out that gathering of raw vibrational data has great potential not only for immediate diagnostic purposes but also for overall quality assessment of the machinery and so a need exists for precise vision based solutions that could replace tactile sensors in the future. A slightly more complex approach has been done by [1] and [18], where 2D CNNs have been used on time-frequency data images for surface roughness estimation and fault diagnosis. In the case of [1] a

classifier is then trained to also differentiate between various types of bearing faults or detect tool wear on a milling machine bit. The authors of [18] use continuous wavelet transform (CWT) instead of short-time Fourier transform (STFT) and a nonlinear auto-regressive neural network to generate more signals to tackle the problem of imbalanced datasets. In [4] vibration data is combined with sound samples to allow for fault detection as well as monitoring of the machine's state using fast Fourier transform and principal component analysis (PCA) to reduce the amount of data. This shows that a camera could be used as a two in one type of sensor and opens the avenue for even more opportunities. In some cases, old archived footage from plant cameras could be used to train a more robust model. In [15] the authors use a high speed industrial camera in a well controlled setting to capture, detect and automatically analyze a rotating body. They use a re-identification network which proved to remove the jitter otherwise associated with frame-to-frame automatic, anchorfree region of interest locating.

V. EXPERIMENTAL SETUP

The dataset in this paper has been gathered as a part of a larger experiment. We use three videos, with increasing amplitude of vibration, of a double sided cantilever beam actuated by a vibrating shaker connected to a function generator Siglent SGD 2042X through an amplifier B&K Type 2732.

The right side of the beam is 30 cm long and has a resonant frequency of 33.8 Hz along the first mode. This frequency was chosen because the part of the broader experiment from which we source our data was intentionally limited to capturing at 120 frames per second. That being said, we used a high speed camera Phantom v1840, again limited to 120 frames per second in this particular case, because we needed a wider range of frame rates to be available. We used Canon EF 20-70mm f/2.8 II USM lens and the entire setup was illuminated by high powered LEDs controlled by the GSVitec MultiLEDG8 driver as can be seen in Fig. 5.

The ground truth measurements at each of the annotated points in Fig. 6 were taken with a laser vibrometer Polytec PDV 100 and data logged using LabVIEW. Note that the sensor at the end of the arm in Fig. 6 is a charge accelerometer and is not relevant to this study.



Fig. 5. Experimental setup. Left to right: shaker with cantilever, laser vibrometer, Phantom v1840. This setup shows measurement of the ground truth for one of the 28 points.

Each point along the arm had first been measured by the laser accelerometer before the vibrating arm had been shot at the aforementioned 120 frames per second by the Phantom V1840 camera in raw format. This way we collected three 5 s clips. The distance of the arm from the camera lens was always 50 cm and the exposure was set to 100 μs .



Fig. 6. Labeled arm points that were measured by the laser vibrometer. Spaced apart by 1 cm.

VI. IMAGE PRE-PROCESSING

As we are dealing with amplitude measurements of a temporal signal, the video must first be transformed into a temporal image. In our case we take a 1 pixel wide vertical slice from each frame around the point to be measured, and stack them side by side to create a time-slice image, as can be seen in Fig. 7 on the right.

Before we settled on using a regression CNN, we tried two different classical CV approaches. All our experiments were compared to ground truth data measured with a laser vibrometer. At first, since we mainly focus on observing the edge of the cantilever, we tried implementing an edge detector to get an amplitude envelope from the time-slice image. However, since the amplitude envelope of our time-slices does not have clear enough edges, tuning an edge detector proved to be too complicated.

The second attempt included applying a 2D Fourier transform on the time-slice images and looking for peaks in the frequency domain. This approach had good results but required manual selection of the frequency peaks which was not possible to automate.

Thus we settled on using a regression CNN. The data further proved to be easy to augment by doing linear interpolation between each pair of the 28 measured points. We gained around 1400 image-amplitude pairs (the distance in Fig. 6 from point 28 to point 1 spans 1400 pixels) from each video, resulting in a dataset containing 4212 samples.



Fig. 7. Comparison of non-magnified and magnified point number 2 from the dataset with the lowest amplitude

VII. CNN TRAINING

We tested several different architectures. First we tried a naive approach and created a simple CNN with two convolutional layers and three fully connected layers, which took in 32x32 images of the cropped time-slice. However, the network was not deep enough to provide good feature selection and ultimately gave random results.

Then we tried training an older deep CNN architecture from scratch. For this purpose we chose VGG16 which accepts 244x244 images, making it potentially usable for lower frequency vibrations. But training the architecture from scratch provided results with very high error.

We retained the 244x244 input size and decided to leverage transfer learning instead of randomizing initial weights. At the same time we decided to try two other architecture, AlexNet and ResNet18, from which we chose the latter as it showed best performance during training.

Trying the more complex versions of ResNet provided only marginal improvements, so we stayed on ResNet18. We tested training with three different final layers of increasing depths, however, the simplest one gave the best results; it consists of a single linear layer, reducing the input 512 features into a single output, and a sigmoid function. We can use the sigmoid function because out data fits into the range of 0 to 1, however, if higher displacement is present in the training dataset, ReLU could be used instead. Finally, we used mean squared error as our loss function to train for regression.



Fig. 8. CNN architecture, ResNet18 adjusted for regression, final layer output is passed through the sigmoid function.

VIII. RESULTS

The CNN performed well without EVM pre-processing for two of our three videos where SNR was low enough for the signal to get picked up by the network. However in the last dataset the noise was too much for our network to handle, as can be seen in Fig. 9.

In Fig. 10 we used phase-based EVM in pre-processing before creating time-slices and passing them to the network. The amplification factor was chosen heuristically to be $\alpha = 100$ and the network outputs were divided by the same number to get the final measurement.

While the result of our CNN and the ground truth of our laser measurements generally correlate, there are still major deviations from the true values. However, the result while using EVM is much better than without using it.



Fig. 9. Measurement of amplitude from video without EVM pre-processing



Fig. 10. Measurement of amplitude from video with the aid of EVM

Point	Laser [µm]	ResNet18 [µm]	$\Delta [\mu m]$			
1	6.353	5.933	0.420			
2	5.790	4.445	1.345			
3	5.231	6.354	1.122			
4	4.765	5.686	0.920			
5	4.448	5.656	1.207	l		
6	3.989	4.723	0.733			
7	3.501	5.448	1.946			
8	3.076	4.356	1.279			
9	2.621	4.911	2.289			
10	2.314	3.211	0.897			
11	1.859	2.041	0.182			
12	1.551	1.876	0.324			
13	1.205	1.342	0.136	l		
14	1.000	1.571	0.571			
15	0.775	1.276	0.501			
16	0.111	0.831	0.719			
17	0.535	0.921	0.386			
18	0.429	1.016	0.586			
19	0.977	0.897	0.080			
20	1.046	1.906	0.860			
21	1.211	0.976	0.235			
22	1.454	1.658	0.204	l		
23	1.617	2.376	0.758			
24	1.809	1.893	0.084			
25	1.925	2.176	0.250	l		
26	2.092	2.442	0.349			
27	2.129	2.453	0.323			
28	2.255	4.050	1.794			
TABLE I						

COMPARISON BETWEEN LASER MEASUREMENT AND CNN REGRESSION FOR LOWEST AMPLITUDE DATASET

IX. CONCLUSION

The goal of this paper was to familiarize the reader with the possibilities of phase-based Eulerian video magnification when applied as an auxiliary algorithm to other computer vision algorithms.

We explained the three major parts of phase-based EVM and how they work together to apply linear magnification of the selected motion within a video.

Furthermore, we trained a ResNet18 as a regression CNN, utilizing transfer learning to measure vibration amplitude of a point from a time-slice image of the video.

We used three videos of a resonating cantilever beam to create a dataset containing 4212 images.

The resulting CNN worked well on two of the three videos but struggled with the third, where vibration amplitude was overwhelmed by noise. However, after applying EVM in the pre-processing step, the network's outputs correlated relatively well with the ground truth data.

In future work we will focus on creating a CV algorithm to extract features of the vibration signal from the video and pass them to a more specialized neural network to achieve better results than the more general oriented CNNs can offer.

REFERENCES

- [1] Han-Yun Chen and Ching-Hung Lee. Deep learning approach for vibration signals applications. *Sensors*, 21(11):3929, 2021.
- [2] Sean Collier and Tyler Dare. Accuracy of phase-based optical flow for vibration extraction. Journal of Sound and Vibration, 535:117112, 2022.
- [3] Nikolaos Giannopoulos. Phasebasedevmcpp. https://github.com/NikolaosGian/PhaseBasedEVMCpp, 2023.
- [4] Seulki Han, Nasir Mannan, Daryl C Stein, Krishna R Pattipati, and George M Bollas. Classification and regression models of audio and vibration signals for machine state monitoring in precision machining systems. *Journal of Manufacturing Systems*, 61:45–53, 2021.
- [5] Luke Hsiao. Cribsense. https://github.com/lukehsiao/CribSense, 2016.
- [6] Kwangsuk Lee, Jae-Kyeong Kim, Jaehyong Kim, Kyeon Hur, and Hagbae Kim. Cnn and gru combination scheme for bearing anomaly detection in rotating machinery health monitoring. In 2018 1st IEEE International conference on knowledge innovation and invention (ICKII), pages 102–105. IEEE, 2018.
- [7] Ce Liu, Antonio Torralba, William Freeman, Frédo Durand, and Edward Adelson. Motion magnification. ACM Trans. Graph., 24:519–526, 07 2005.
- [8] Tae-Hyun Oh, Ronnachai Jaroensri, Changil Kim, Mohamed Elgharib, Fr'edo Durand, William T Freeman, and Wojciech Matusik. Learning-based video motion magnification. In Proceedings of the European Conference on Computer Vision (ECCV), pages 633–648, 2018.
- [9] Karl Pauwels and Marc M. Van Hulle. Realtime phase-based optical flow on the gpu. In 2008 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops, pages 1–8, 2008.
- [10] E.P. Simoncelli, W.T. Freeman, E.H. Adelson, and D.J. Heeger. Shiftable multiscale transforms. *IEEE Transactions on Information Theory*, 38(2):587–607, 1992.
- [11] Neal Wadhwa, Michael Rubinstein, Frédo Durand, and William T. Freeman. Phasebased video motion processing. ACM Trans. Graph., 32(4), jul 2013.
- [12] Neal Wadhwa, Michael Rubinstein, Frédo Durand, and William Freeman. Phasebased video motion processing. ACM Transactions on Graphics (TOG), 32, 07 2013.
- [13] Neal Wadhwa, Michael Rubinstein, Frédo Durand, and William T. Freeman. Riesz pyramids for fast phase-based video magnification. In 2014 IEEE International Conference on Computational Photography (ICCP), pages 1–10, 2014.
- [14] Hao-Yu Wu, Michael Rubinstein, Eugene Shih, John Guttag, Frédo Durand, and William Freeman. Eulerian video magnification for revealing subtle changes in the world. ACM Transactions on Graphics - TOG, 31, 07 2012.
- [15] Rongliang Yang, Sen Wang, Xing Wu, Tao Liu, and Xiaoqin Liu. Using lightweight convolutional neural network to track vibration displacement in rotating body video. *Mechanical Systems and Signal Processing*, 177:109137, 2022.
- [16] Yongchao Yang, Charles Dorn, Tyler Mancini, Zachary Talken, Garrett Kenyon, Charles Farrar, and David Mascareñas. Blind identification of full-field vibration modes from video measurements with phase-based video motion magnification. *Mechanical Systems and Signal Processing*, 85:567–590, 2017.
- [17] Yichao Zhang, Silvia L. Pintea, and Jan C. van Gemert. Video acceleration magnification. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), July 2017.
- [18] Quan Zhou, Yibing Li, Yu Tian, and Li Jiang. A novel method based on nonlinear auto-regression neural network and convolutional neural network for imbalanced fault diagnosis of rotating machinery. *Measurement*, 161:107880, 2020.

The Human-machine interface for UAV ground control station

Jan Klouda UTEE Brno University of Technology Brno, Czechia xkloud04@vutbr.cz Petr Marcoň UTEE Brno University of Technology Brno, Czechia marcon@vutbr.cz

Abstract—This paper describes the design of a Humanmachine interface for UAVs and presents general HMI requirements for the operator (commander). The design was created for a ground control station that is tasked with safely and efficiently controlling several UAVs at once. The goal was to create an environment that provides the operator with an overview of each UAV, mission planning and execution.

Keywords—UAV, Ground control station, design, HMI

I. INTRODUCTION

In today's era of constant technological advancement and the increasing importance of unmanned systems in all areas of human activity, the Human-Machine Interface (HMI) is becoming a key element connecting humans with technology. From the perspective of UAV (Unmanned Vehicle) research and utilization, HMI can fundamentally affect the efficiency, safety and success of missions. Ground Control Station (GCS), as the central point of management and control of UAV operations, represents an important node of connection between the operator and technology. Thus, the selection of an appropriate GCS screen design is one of the important elements of the system. Its efficiency and clarity can either simplify the operator's work and increase efficiency, or on the contrary, represent a source of confusion and risk for errors. HMI is appearing in multiple sectors from industry to avionics. This interface is also used in Industry 4.0, where it presents the state of the environment and sensors Chyba! Nenalezen zdroj odkazů.Chyba! Nenalezen zdroj odkazů.Chyba! Nenalezen zdroj odkazů.

During UAV operations, emphasis is placed not only on the aesthetic aspects of design, but above all on the ability of the HMI to present relevant information to the operator in a clear and intuitive form. Proper GCS screen design can not only make the operator's job easier, but also eliminate unnecessary stress, increase decision-making efficiency, and improve overall mission performance. This paper focuses on the importance and process of GCS screen design for UAV operators and analyzes the key factors that affect its effectiveness **Chyba! Nenalezen zdroj odkazů.Chyba! Nenalezen zdroj odkazů.**

Specifically, a graphical GCS design for a swarm of UAVs operating as a reconnaissance squadron will be described. This swarm of UAVs will be fully automated and will perform terrain reconnaissance and data collection in a variety of environments and conditions. The following section will describe the operational procedures that form the basis for the HMI design. Another important aspect can be considered the way data is presented, which may vary depending on the nature of the mission and the operator's level of UAV knowledge **Chyba!** Nenalezen zdroj odkazů.

II. OPERATION PROCEDURES

The Ground Control Station is the central point of control and management of UAVs operations. Effective execution of these operations requires carefully defined procedures that include several key phases. From the Home Screen to mission planning, pre-flight preparation, and actual mission execution, each phase is critical to the successful and safe conduct of the operation. The following paragraphs will describe each step of the GCS operational procedures in more detail, including their importance and impact on the overall effectiveness and safety of UAV missions.

The GCS is planned as a portable station, with only one 1920x1080 pixel monitor. This design was created for a user who has only received basic knowledge as a UAV operator. Thus, an operator with a broader knowledge and subconscious understanding of UAVs is not considered to control a swarm of drones. The GCS does not support the ability to directly control individual UAVs. The overall flight will be fully automated. Thus, from this perspective, even the mission planner is simplified to the level that the operator only enters high-level information and the system creates individual flight plans for the UAV. The transmission of in-flight telemetry information is also simplified, and symbolic methods are chosen to display the status of individual UAVs.

The HMI design is intended for use in emergency services. These emergency services can use unmanned vehicle technology to gain situational awareness. For the Ambulance Service, the technology would be useful in providing an overview of mass car crashes, or other larger incidents. The military could use the system for reconnaissance of an area and use it as a source of information for intelligence purposes. A similar use could be found for the Police. The design is therefore designed for a wider range of security forces, and it is possible to modify it for different types of use in different forces.



Fig. 1. Block diagram of the mission flow with individual screens

A. Home screen

The Start screen is the operator's entry point into the GCS environment and provides an overview of the basic GCS functions. A well-designed start screen should be intuitive and provide the operator with a quick and easy transition to the next phases of the mission. This section of the operating procedures includes an analysis of the available features. From the start screen, one can proceed to the GCS settings, where one can set the connection to the UAV, change the screen contrast, and other functions (GCS settings). It is also possible to progress from the Home screen to the screen where individual UAVs can be connected to perform sensor calibration and detailed setup and inspection. The final path from the Home screen leads to the mission planning screen where the operator creates a mission for a swarm of UAVs.

B. Mission planning

Mission planning is the step before the actual execution of the operation. The operator can select a pre-planned mission from the database, edit it, or create a new mission. If the operator chooses to edit or create a new mission, the GCS operator uses area delineation and flight parameter determination to achieve the mission objectives. Careful mission planning can minimize risks and increase the efficiency of the operation. Proper GCS screen design for this phase should provide the operator with interactive planning tools, terrain visualization, flight plans, and waypoint lists. The operator enters the area of interest, the lowest flight level, the highest flight level (UFL - upper flight level, LFL- lower flight level) and the information of interest in that area. This can be radiation sensors, photogrammetric map using thermal camera, multispectral camera, etc. The GCS is able to use this information to decide what and how many UAVs are needed to cover the selected area. It will plan a flight plan for each UAV so that no collisions occur during the flight.

C. Pre-flight preparation

Pre-flight preparation for flight safety of UAVs. This phase includes checking and setting up the UAVs hardware, verifying GPS signal availability, calibrating sensors, checking the status of batteries and other important elements. The GCS pre-flight preparation screen should provide the operator with an overview of all the necessary tasks and information that need to be checked and performed prior to launch. A list (Check List) will open where the operator will find all available systems to check. The Check List should be fully automated.

D. Mission execution

The actual execution of the mission is the main phase of the whole operation. During this phase, the operator actively monitors and indirectly controls the UAV in accordance with the defined mission objectives and parameters. During mission execution, the GCS screen provides the operator with relevant real-time information on the UAVs position, battery status, environment and other relevant factors to effectively react to changes and situations during the mission.

III. GRAPHIC DESIGN OF SCREENS

The graphical design of the screens is a Human-Machine Interface (HMI) element that has a major impact on the efficiency of UAV operations. As already mentioned, before the design it is necessary not only to describe the procedure of the expected operations, but also to describe the operator in terms of knowledge and skills. Another aspect is the hardware capabilities of the GCS itself. In this proposal, a GCS with only one screen is assumed. The GCS should be part of the command post, be mobile and convey information clearly and quickly.

A. Selection of icons and symbol

This element of the graphic design of GCS screens has a significant impact on user-friendliness and clarity. Iconography is often used to visually represent functions, states and actions, allowing for quick operator identification and orientation within the interface. When choosing icons and symbols, it is important to consider the principles of simplicity, clarity and consistency. Icons should be intuitive and easy to identify so that the operator can quickly recognize their meaning. A thorough analysis of the user needs and context of use is essential for the proper selection of icons and symbols that effectively communicate the desired information and functions without unnecessary confusion.

The choice of colours and the simplicity and clarity of symbols and symbols are also important, contributing to easy identification and visual organisation of information on the GCS screen. Colours can be important not only from an aesthetic point of view, but also from a functional point of view, for example in indicating statuses and priorities. It is important to choose an appropriate colour scheme and apply it consistently throughout the interface. For the graphic design, dark colours were chosen for the background to add contrast when the sun illuminates the screen. Four primary colours were chosen to convey the information. White for informational messages, green for nominal status indication, orange for warning (alert) indication and red for critical failure indication. Simplicity and clarity of symbols and symbols are key to minimising user error and confusion. Symbols that are too complex and unclear could lead to misinterpretations, resulting in incorrect decisions or inadequate operator response. Therefore, it is imperative that icons and symbols are clearly defined and intuitive to the user, which will contribute to the effective use of GCS screens and the overall success of UAVs operations. For example, in the Fig. 2 you can see the UAV icon from the North Atlantic Treaty Organization (NATO) perspective and the icon selected for this design. For use in military operations, NATO symbology will need to be used. A custom design has been chosen for this design which can also display the UAV's status by changing the background color.



Fig. 2. The UAV icon for NATO (left) and the created icon for this design proposal (right) [4].

B. Design procedure

The process of designing the graphic design of GCS screens requires careful planning and a systematic approach. It starts with an analysis of user needs and identification of the key features that GCS screens must postroad. This is followed by defining the functionality and information layout requirements. Wireframes and prototypes are used to visualize the designs and test them with users. An iterative process involves refining designs based on feedback and testing until the optimal design is achieved.

C. Final design

The final graphic design of GCS screens is based on an analysis of user needs and functionality requirements. It includes visual elements such as colour schemes, typography, icons and graphic elements that are designed to provide simplicity, clarity and ease of orientation for the user. The detailing of each element and its placement on the screen creates a functional design that will trigger the user to perform at their best. The following photos describe the individual screens designed for an operator managing a swarm of UAVs to explore areas.

1) Home screen

The home screen shows the signpost for the operator. It is designed simply and shows only the time, the system logo, and buttons leading to other screens.



Fig. 3. GCS home screen for controlling a swarm of UAVs

2) Mission planning screen

This screen is very variable. It varies depending on the choice of operator. In the first phase, it is possible to select a preplanned mission and either edit, delete or select it. The operator can also create a new mission. The new mission creation screen is shown in Fig. 4.



Fig. 4. Mission planning screen

3) Mission execution screen

It is the most critical and important screen of the entire proposal. It must convey mission progress information and basic telemetry information of the UAVs to the operator. It must also transmit real-time information about the threats detected by the selected sensors and write them directly to the map base. It is possible to monitor the mission progress of the entire swarm of UAVs during the operation, or to monitor the image stream from a selected UAV. This screen is also needed to inform the operator of any messages coming from the UAVs. The Messages window in the right corner is used for this purpose. In the upper right corner is basic telemetry information about each of the UAVs. These are presented by symbols in three colours. Each of the colours represents the status of each system. By clicking on the desired UAV, detailed telemetry data can be displayed and detailed information about the UAV status can be transmitted.



Fig. 5. Screen during the progress of a scheduled mission

High-level telemetry uses four icons. The first icon is the multicopter symbol, which shows the status of the UAV (internal sensor calibration and HW status). The second icon shows the GNSS (Global Navigation Satellite System) status, another signal strength between the UAV and the GCS. The last icon shows the battery status. The icons can be seen in the following Fig. 6.



Fig. 6. Icons for high-level telemetry. From right: UAV status, GNSS, GCS connection strength, battery status

D. Verification of the functionality of the designs

Verifying the functionality of the proposed GCS screen designs is a step before deploying them in real operation.

Testing involves examining the user interface, interaction elements and navigation within simulated environments or realworld situations. The goal is to ensure that the graphical design supports effective use of the GCS and contributes to the safety and success of UAVs operations. These validations are performed iteratively where information and suggestions are gathered from test operators. The wider the diversity of operator focus, the better the potential design flaws can be detected.

E. Limitations

Despite efforts to optimize and improve the graphic design of GCS screens, there may be limitations that can affect its effectiveness and usability. These constraints may include limited screen size, insufficient system resources, ergonomic limitations of the operator's work environment, and other technical and operational factors that must be considered when designing and implementing the graphic design. The proposed graphic design has several of these constraints. The first is the hardware limitation of a single screen, where there is not enough space to display all possible information. Another is the choice of the map background, where a dark contrasting map would be better for outdoor use.

IV. CONCLUSION

This paper has dealt with the graphical design of GCS for automatic swarm of UAVs. The design is based on several iterations. The mission planning process and the various nominal scenarios that this design reflects were described. Custom icons were created to communicate high-level information to the operator and justify some of the main design elements. This design has several limitations and can therefore only be used for the purpose for which it was created. A total of three screens were presented, with more in the design itself, with many additional modifications. This graphic design can be implemented after subsequent additional evaluations and the programming itself. This design can be linked to the back-end to form a complete GCS.

References

- M. Friedrich and M. Vollrath, "Human–Machine Interface Design for Monitoring Safety Risks Associated with Operating Small Unmanned Aircraft Systems in Urban Areas", Aerospace, vol. 8, no. 3, 2021.
- [2] L. Garbarino, N. Genito, G. Di Capua, and R. Rocchio, "Innovative Low-Cost Design of a Ground Control Station for Unmanned Aerial Systems Experimentation", in 2023 IEEE/AIAA 42nd Digital Avionics Systems Conference (DASC), 2023, pp. 1-8.
- [3] P. Marcon, J. Arm, T. Benesl, F. Zezulka, C. Diedrich, T. Schröder, A. Belyaev, P. Dohnal, T. Kriz, and Z. Bradac, "New Approaches to Implementing the SmartJacket into Industry 4.0 ‡", Sensors, vol. 19, no. 7, 2019.
- [4] "Milsymbol MIL-STD-2525C Implementation", Spatialillusions. from: https://doi.org/10.1109/DASC58513.2023.10311145.

SKODA Kariéra



Stáže a závěrečné práce

*RED

Studuješ vysokou školu a zároveň chceš získat pracovní zkušenosti? Zkus odbornou stáž ve ŠKODA AUTO. Pod vedením našich odborníků se zapojíš do inovativních projektů – třeba těch v oblasti elektromobility. Získané znalosti pak můžeš zúročit i ve své závěrečné práci.



Doktorandský program

Aplikuj výsledky svého výzkumu v reálném prostředí, využívej nejnovější technologie a dej své disertační práci nový rozměr. Poskytneme ti flexibilní pracovní dobu a možnost spolupracovat s profesionály ve svém oboru.



Trainee program

Dokončil jsi studia a přemýšlíš kam dál? Trainee program pro absolventy ti pomůže najít obor, který tě bude opravdu bavit a naplňovat. V průběhu jednoho roku budeš na plný úvazek pracovat v mezinárodním týmu, poznáš různá oddělení a vycestuješ na zahraniční stáž.



Volné pozice

Vyber si ze široké nabídky volných pozic, staň se součástí nejprogresivnější české firmy a tvoř s námi budoucnost automobilového průmyslu.





O @WeAreSKODA









Machine Learning-based Fingerprinting Localization in 5G Cellular Networks

Thao Dinh Le Department of Telecommunications FEEC, Brno University of Technology Brno, Czech Republic Dinh.Thao.Le@vut.cz Pavel Masek Department of Telecommunications FEEC, Brno University of Technology Brno, Czech Republic masekpavel@vutbr.cz

Abstract—This study explores the viability of employing machine learning (ML)-based fingerprinting localization in 5G heterogeneous cellular networks. We conducted an extensive measurement campaign to collect data and utilized them to train three ML models: Random Forest (RF), Extreme Gradient Boosting (XGBoost), and Support Vector Machine (SVM). The findings reveal that RF delivers the highest accuracy among the three ML algorithms. Furthermore, the results indicate that 5G New Radio (NR) can benefit the most from this localization method due to the dense deployment of base stations, achieving median localization errors of 17.5 m and 106 m during the validation and testing phases, respectively.

Index Terms—Fingerprinting Localization, Machine Learning, 5G New Radio, NB-IoT, LTE-M, Random Forest, XGBoost, Support Vector Machine

I. INTRODUCTION

Localization is a critical aspect of cellular networks, offering a multitude of benefits across various domains such as emergency services, enhanced user experience, network resource optimization, and security. With the emergence of the Internet of Things (IoT) and Machine-to-Machine (M2M) communication in recent years, localization has become essential for managing connected devices and enabling autonomous operations. Within the 5G standardization, NarrowBand Internet of Things (NB-IoT) and Long-Term Evolution Machine Type Communication (LTE-M) have emerged as notable technologies designed specifically to provide efficient communication for a massive number of simple IoT devices. Given the booming rise of M2M communication and IoT, these technologies are in place to become integral parts of the next-generation cellular networks [1], [2].

Due to the unique characteristics of IoT devices and network constraints, localization of devices in LTE-M and NB-IoT networks faces several limitations. IoT devices are typically required to operate on battery power for extended periods of time, usually up to ten years or more. Therefore, localization techniques must be energy-efficient to minimize the impact on battery life. In addition, most IoT devices have limited processing and hardware capabilities compared to smartphones or other mobile devices. For example, in order to reduce module cost, they are typically equipped with only a single antenna, and may lack Global Navigation Satellite System (GNSS) receivers or advanced inertial sensors. These constraints significantly restrict the range of localization techniques viable for implementation on such devices since the more costly, complex, and power-demanding techniques such as GNSS, angle-based positioning, and time of arrival-based positioning are not suitable for them [2].

This paper aims to address the above issue by exploring the practical applicability of fingerprinting localization technique in combination with machine learning (ML) algorithms. Fingerprinting localization has been gaining popularity in recent years as an efficient localization technique. In the context of IoT, this method is suitable for battery-powered, low-complex end devices, as it does not require additional hardware or complex computational processes on their end. Instead, on the network operator's end, it requires measurement campaigns to establish a fingerprint database of signal characteristics in the targeted area.

The main contributions of the paper are:

- An extensive measurement campaign of the three 3GPP cellular technologies—5G New Radio (NR) in Non-Standalone (NSA) mode, LTE-M and NB-IoT—across the city of Brno, Czech Republic, within the network of Vodafone operator.
- Implementation of the fingerprinting localization method utilizing three ML algorithms: Extreme Gradient Boosting (XGBoost), Random Forest (RF), and Support Vector Machine (SVM) for the aforementioned communication technologies.
- Evaluation of the ML models' performance to determine the practicality of the localization method.

The remainder of the paper is structured as follows: Section II presents a description of fingerprinting localization and an overview of the communication technologies of interest. Section III offers insights into the measurement campaign and analyzes the acquired dataset. The performance of the ML models is presented in Section IV, followed by a comprehensive discussion. Finally, the paper is concluded with Section V.

II. BACKGROUND AND RATIONALE

A. Fingerprinting localization

In recent times, fingerprinting has emerged as a prominent positioning technique in wireless networks in general and mobile networks in particular. Fingerprinting localization determines the location of a mobile device within a cellular network by comparing the received signal characteristics of the device to a pre-existing database of signal fingerprints collected from known locations [3], [4].

Fingerprinting localization involves two main phases offline and online, as illustrated in Fig. 1. In the offline phase, a database of signal fingerprints is created by collecting measurements of signal characteristics from various locations within the network's coverage area. These measurements can include metrics about signal strength such as Received Signal Strength Indicator (RSSI) and Reference Signal Received Power (RSRP), Signal-to-Noise Ratio (SNR), Angle of Arrival (AOA), and others, depending on the available infrastructure and requirements. Each fingerprint in the database is associated with its corresponding geographical location, which is usually acquired by mapping the collected signal characteristics to specific coordinates using techniques such as GNSS or manual surveying. In the online (deployment) phase, when a mobile device needs to be localized, it measures the current signal characteristics from nearby base stations (BS). These measurements are then compared to the fingerprints in the database, and the location corresponding to the closet match is considered to be the estimated location of the device. Various matching algorithms can be used to map the measured signal characteristics to the fingerprints in the database. Examples of such algorithms could be nearest neighbor, weighted averaging, probabilistic methods, and ML-based methods [3], [4].



Fig. 1. Principle of ML-based fingerprinting localization [5]

Fingerprinting localization can achieve good accuracy, especially in environments with complex propagation characteristics, where other localization techniques may struggle. Furthermore, fingerprinting c an u tilize a v ariety o f signal metrics, making it adaptable to different wireless technologies and deployment scenarios. However, this technique also has drawbacks. A notable disadvantage is that fingerprinting localization is sensitive to changes. For instance, changes in network infrastructure, environmental conditions, or user behavior may necessitate updates to the database. Therefore, in order to maintain accuracy over time, fingerprint databases have to be updated regularly. Moreover, creating an accurate fingerprint database can be labor-intensive and may require extensive surveying and measurement campaigns, especially in large and complex areas [3], [4].

B. 5G cellular technologies

1) 5G New Radio: Introduced in 3GPP Release 15, 5G NR is the global standard for 5G wireless communication. It defines the air interface and protocols for 5G networks. NR aims to provide significantly higher data rates, lower latency, and better reliability compared to previous generations of mobile networks. In specific, it strives to deliver peak data rates of 20 Gbps and ultra low-latency of around 1 ms. enabling a variety of new applications such as Virtual Reality (VR) and Augmented Reality (AR). 5G NR operates in a wide range of frequency bands, which include spectrum bands below 6 GHz (sub-6GHz) for broader coverage and better indoor penetration, and millimeter-wave (mmWave) for ultrafast, high-capacity connections in dense urban areas. In 5G NR, advanced antenna techniques such as massive Multiple Input Multiple Output (MIMO) and beamforming are employed to increase spectral efficiency and network capacity [6].

Regarding the rollout of 5G technology, there are currently two primary configurations: Standalone (SA) and Non-Standalone (NSA). SA deployment represents a complete 5G infrastructure, including both the 5G Radio Access Network (RAN) and 5G core network. On the other hand, NSA deployment integrates a 5G RAN with an LTE core, resulting in a partial implementation of 5G NR capabilities. Due to its ease of deployment and utilization of existing LTE infrastructure, NSA has been the preferred choice for initial 5G network deployments, and it remains the sole available solution in current public mobile networks within the Czech Republic [6].

2) Narrowband Internet of Things: Introduced in 3GPP Release 13, NB-IoT is tailored to meet the stringent requirements of massive Machine-Type Communication (mMTC) applications. As the name suggests, NB-IoT operates within a narrow bandwidth of only 180 kHz, which neatly fits into a single LTE Physical Resource Block (PRB), an LTE guard band, and even a 200 kHz Global System for Mobile communication (GSM) carrier. NB-IoT is optimized for applications which only require sporadic transmission of small data packets. However, it can offer extended coverage range, enabling communication over long distances and in challenging environments, such as underground areas or deep inside buildings. Considering the typical deployment of a massive number of User Equipments (UEs) in mMTC scenarios, NB-IoT also strives for maximum simplicity to minimize the cost of communication modules. In specific, devices are only equipped with a single antenna and can only operate in half-duplex mode. Furthermore, NB-IoT is designed to operate with minimal power consumption, making it suitable for battery-powered IoT devices [7].

3) LTE-M: Also introduced in Release 13, LTE-M is a variant of the LTE technology optimized for IoT. Similar to NB-IoT, LTE-M operates within a narrow bandwidth, which is 1.4 MHz for Cat-M1 and 5 MHz for Cat-M2. This narrowband design allows LTE-M to coexist alongside other LTE services while efficiently utilizing available spectrum resources. This technology offers extended coverage range compared to traditional LTE, providing reliable connectivity in challenging environments such as underground and remote rural areas. Compared to NB-IoT, LTE-M offers higher data rates, making it suitable for applications that require moderate throughput, such as firmware updates, asset tracking with frequent location updates, and voice over LTE (VoLTE). In addition, LTE-M also supports mobility, allowing devices to maintain connectivity while on the move. Nevertheless, these benefits are accompanied by compromises, as LTE-M requires greater bandwidth, typically involves higher expenses, offers slightly reduced communication range (in comparison with NB-IoT), and cannot operate in guard-band mode [7].

III. MEASUREMENT CAMPAIGN

A. Measurement setup

The core of the measurement setup was the Rohde & Schwarz TSMA6B autonomous mobile network scanner, operating on a Windows platform. The TSMA6B was controlled by a laptop via Wi-Fi and the Remote Desktop Protocol (RDP). Utilizing its integrated software-defined radio (SDR) capability, this scanner has the ability to analyze various cellular standards, including LTE, 5G NR, NB-IoT, and LTE-M. Within the setup, the TSMA6B primarily served to oversee all connected wireless modules and to obtain precise location data through its GNSS module. In terms of communication modules, two Quectel BG77 modules were used for NB-IoT and LTE-M connectivity, and a Quectel RM520N module was used for 5G NR. All communication modules were controlled via AT commands, linked via USB ports to the TSMA6B analyzer, and connected to the cellular network provided by Vodafone Czech Republic, utilizing the corresponding frequency bands available in the Brno area. The measurement data from the communication modules were sent via debug port to the TSMA6B, where they were decoded and stored on the hard drive of the analyzer. The dataset obtained by the TSMA6B from the measurement campaign is extensive. However, in the context of this paper, the parameters of interest are: i) GPS coordinates, ii) BS identifiers (Cell ID, PCI), and iii) RSRP values. The entire measurement procedure was automated using a measurement script and the predefined autonomous functionality of the communication modules to maintain mobile connectivity.

Two different driving sessions were carried out in the city of Brno, as shown in Fig. 2. The first drive was more extensive and was intended to gather necessary data to train the ML models. Subsequently, a shorter second drive was



Fig. 2. Map of measurement campaign in Brno

conducted to acquire data for testing purposes. In total, the two measurement drives lasted 12 hours and spanned 220 kilometers. For the second driving session, we opted for a route that intersected partially with the first d rive, b ut also included minor unmapped areas. This choice was made to simulate a real-life localization scenario where the UE deviates from the known route but still remains within the confines of the trained model. It is anticipated that the ML models' performance with the testing data from the second drive will decline in comparison to the self-validation data from the first drive. Our objective is to determine whether this reduction in precision falls within an acceptable range for the proposed method to remain viable. It's essential to bear in mind that, for many IoT applications, pinpoint accuracy in the location of the IoT device isn't of vital importance, and an error margin spanning several hundred meters to a few kilometers is deemed acceptable.

B. Dataset analysis

In the obtained dataset, 5G NR exhibited the highest density of BS deployment, with a total number of 1388 BSs identified. LTE-M came in second with a total of 532 BSs observed in the dataset. NB-IoT had the lowest BS density, with only 331 BSs identified. This could be attributed to NB-IoT's design, which aims to provide the broadest coverage among the three communication technologies. It's also worth noting that the difference between 5G and the other two technologies may not solely stem from a higher density of BSs; NSA deployment could also be a contributing factor. In NSA deployment, only downlink communication occurs through 5G cells, while uplink communication is still transmitted via LTE cells. This dual connection approach might further contribute to the observed increase in the number of BSs in the case of 5G NR.

In terms of the size of the dataset, roughly 29000 data points were recorded for each of the technologies during the initial driving session. In the second driving session, more than 7300 data points were collected for 5G NR and NB-IoT each, while roughly 6500 data points were logged for LTE-M.

An important factor that could potentially impact localization accuracy is the distance separating the UE from the serving BS. The shortest average distance between the UE and its serving BS was observed in the case of NB-IoT, measuring at 472 m. Following closely behind, 5G NR ranked second with an average distance of 488 m, while LTE-M exhibited the greatest average separation, reaching a value of 577 m.

IV. PERFORMANCE ANALYSIS

The findings depicted in Fig. 3 and Tab. I indicate that RF offers the most reliable localization capability, with XGBoost ranking second, while SVM exhibits the poorest performance among the three ML algorithms. Additionally, it is demonstrated in Fig. 4 that, overall, 5G NR achieved better results than LTE-M and NB-IoT across all three ML models. This can be attributed to the denser deployment of 5G BS in comparison to the other two technologies.

For 5G NR, RF demonstrates median error values of 17.5 m and 106 m for the validation dataset and the test dataset, respectively. Moreover, the 75th percentile of localization errors for RF remains under 51 m with the validation data and 215 m with the test data. Despite being more complex, the XGBoost algorithm demonstrates a decrease in precision of approximately 260 m (validation set) and 263 m (test set) in median error compared to RF. SVM demonstrates slightly better accuracy than XGBoost in the validation phase, but lags behind in the testing phase, yielding a median error of 783 m.

Regarding LTE-M, as depicted in Fig. 3b, RF continues to exhibit superior performance over the other two ML algorithms, demonstrating a median error of 64 m for the validation dataset and 176 m for the test dataset. Interestingly, SVM surpasses XGBoost significantly in the validation phase, achieving accuracy close to that of RF. However, SVM struggles to maintain its localization accuracy with the test dataset, resulting in a median error of 1101 m.

Lastly, in the case of NB-IoT, a comparable trend to LTE-M is observed, as illustrated in Fig. 3c. However, the performance is marginally inferior across all three ML models. Specifically, RF maintains its superior accuracy in both validation and testing, with a median error of 60 m and 205 m, respectively. SVM continues to surpass XGBoost during the validation phase. However, SVM's accuracy declines significantly with the test dataset, resulting in a median error of 1065 m and a 75th percentile error of 1452 m.

In terms of error distribution, as depicted in Fig. 4a and 4d, RF remains the most reliable algorithm, with the majority of its localization errors falling below 250 m for all three communication technologies, across both the validation and testing datasets. With XGBoost, as seen in Fig. 4b and 4e, the localization errors exhibit much greater dispersion, with the highest concentration around 250 m for 5G NR and NB-IoT, and 750 m for LTE-M. In the case of SVM, the algorithm delivers consistent performance for all three technologies during the validation phase, with the majority of errors concentrated within the area under 500 m. Notably, during the validation phase, SVM produces the best results for LTE-M, slightly exceeding the accuracy of 5G NR and NB-IoT, as shown in Fig. 4c. However, as the trend has been observed above,

TABLE I: Comparison of ML models' performance

Model	R ² score	RMSE	Median	75 th pctl.	
RF val.	99.74%	0.0013	17.51	50.72	
RF test	95.01%	0.0045	106.14	214.62	
XGBoost val.	96.90%	0.0046	277.99	511.95	ЪЯ
XGBoost test	90.72%	0.0061	369.46	728.45	IJ
SVM val.	97.78%	0.0040	216.55	302.54	α,
SVM test	74.51%	0.0099	783.41	1300.72	
RF val.	99.23%	0.0022	64.03	155.03	
RF test	88.84%	0.0065	175.97	391.54	
XGBoost val.	90.70%	0.0080	720.36	1012.91	N-
XGBoost test	73.95%	0.0100	746.07	1150.19	ΞĘ,
SVM val.	97.91%	0.0039	211.98	290.72	Ι
SVM test	66.45%	0.0114	1101.00	1700.34	
RF val.	99.38%	0.0019	59.61	167.17	
RF test	93.71%	0.0049	204.89	531.90	r.
XGBoost val.	95.76%	0.0052	532.46	715.35	LoI
XGBoost test	91.70%	0.0057	511.88	780.13	ġ
SVM val.	98.46%	0.0031	226.15	304.65	Z
SVM test	73.32%	0.0103	1065.48	1451.63	

SVM encounters difficulties with unseen data in the testing phase, evident from the wide distribution of errors for all three technologies, which is concentrated within the range of 500 m to 1200 m.

From the obtained results, there is no obvious correlation between localization accuracy and the distance from the UE to the serving BS. Although NB-IoT exhibited the shortest average distance between the UE and the BS, its localization error is only comparable to that of LTE-M, and notably higher than that of 5G NR. Rather, it is the density of BS deployment specific to each technology that appears to exert a more conspicuous influence on the precision of localization.

V. CONCLUSION

In order to evaluate the real-world feasibility of fingerprinting localization employing ML algorithms in modern cellular networks, we conducted a measurement campaign in Brno, Czech Republic, focusing on the three predominant technologies: 5G NR, LTE-M, and NB-IoT. Subsequently, the collected data were used to train three ML models, namely RF, XGBoost, and SVM. The performance of these models was evaluated to determine the method's practicality and identify the ML algorithm capable of achieving the highest accuracy.

The obtained findings suggest that ML-based fingerprinting localization is applicable for contemporary cellular networks, with the RF algorithm emerging as the most accurate among the three analyzed ML algorithms. The RF model consistently exhibited high performance across all three communication technologies, establishing it as a viable option for real-world implementation. Although XGBoost provided lower accuracy than RF, it demonstrated solid performance across all three technologies in both the validation and testing phases. It is expected that with a larger training dataset, the XGBoost model's performance will certainly improve. On the other hand, although the SVM model performed relatively well in the validation phase, achieving accuracy surpassing XGBoost and nearing that of RF, it delivered poor results with unseen



Fig. 3. Cumulative distribution function of localization errors by technology



Fig. 4. Localization error distribution by ML model

testing data. This suggests that the SVM model lacks reliability, rendering it unsuitable for this application.

It's worth noting that the method's accuracy can be significantly improved with a larger dataset. During our measurement campaign, we tried to conduct the measurements as thorough as possible within our resource constraints. Nevertheless, with greater investment, a much larger dataset could be obtained with ease. Additionally, incorporating additional contextual features like time of day and weather conditions into the dataset could further enhance the ML model's performance.

REFERENCES

 J. A. del Peral-Rosado, R. Raulefs, J. A. López-Salcedo, and G. Seco-Granados, "Survey of cellular mobile radio localization methods: From 1G to 5G," *IEEE Communications Surveys & Tutorials*, vol. 20, no. 2, pp. 1124–1148, 2018.

- [2] F. Khelifi, A. Bradai, A. Benslimane, P. Rawat, and M. Atri, "A survey of localization systems in Internet of Things," *Mobile Networks and Applications*, vol. 24, pp. 761–785, 2019.
- [3] F. Alhomayani and M. H. Mahoor, "Deep learning methods for fingerprint-based indoor positioning: A review," *Journal of Location Based Services*, vol. 14, no. 3, pp. 129–200, 2020.
- [4] S. Subedi, J.-Y. Pyun *et al.*, "Practical fingerprinting localization for indoor positioning system by using beacons," *Journal of Sensors*, vol. 2017, 2017.
- [5] M. Stusek *et al.*, "On applicability of timing advance-based localization in 5G cellular networks," *unpublished*, 2024.
- [6] E. Dahlman, S. Parkvall, and J. Skold, 5G NR: The next generation wireless access technology. Academic Press, 2020.
- [7] O. Liberg, M. Sundberg, E. Wang, J. Bergman, J. Sachs, and G. Wikström, Cellular Internet of Things: From massive deployments to critical 5G applications. Academic Press, 2019.

Comparative Analysis of Gaussian Process Regression Modeling of an Induction Machine: Continuous vs. Mixed-Input Approaches

1st Vladimír Bílek

Department of Power Electrical and Electronic Engineering Brno University of Technology Brno, Czech Republic vladimir.bilek@vut.cz

Abstract—This paper investigates the application of machine learning technique for modeling continuous and mixed-input parameters of electrical machines. The design of electrical machines typically requires the consideration of certain parameters as integer values due to their physical significance, including the number of stator/rotor slots, stator wires, and rotor bars. Traditional machine learning methods, which predominantly treat input parameters as purely continuous, may compromise modeling accuracy for such applications. To address this challenge, models capable of handling mixed-input parameters were used for the case study. Two training datasets were generated: one with purely continuous inputs and another with both continuous inputs and a categorical parameter, specifically, the number of stator conductors. Gaussian process regression was employed to build three models: two with continuous kernels, trained on both datasets, and one with a mixed kernel, trained only on the dataset containing a categorical parameter. A comparative analysis, demonstrated on a 1.5 kW induction machine - though applicable to a wide range of machines - illustrates the differences between the proposed approaches. The results highlight the importance of selecting an appropriate model for the Multi-Objective Bayesian optimization of electrical machines.

Index Terms—Finite element method, Gaussian process regression, Induction machine, Machine learning, Mixed-Input surrogate models, Surrogate modeling

I. INTRODUCTION

Nowadays, the use of numerical methods such as the Finite Element Method (FEM) is very common for the electromagnetic calculation of induction machine (IM) [1]. FEM often provides very accurate results compared to the measured data. The design of an IM leads to a multi-objective problem, where the Multi-Objective Optimization is the most efficient solution to find the optimal design. The leading method in the field of Multi-Objective IM design Optimization is the use of the non-dominated sorting genetic algorithm II (NSGA-II) [2]. However, the time-consuming nature of FEM simulations and the need of the high count of designs evaluations, makes this

This research work has been carried out in the Centre for Research and Utilization of Renewable Energy (CVVOZE). Authors gratefully acknowledge financial support from the Ministry of Education, Youth and Sports under institutional support and BUT specific research programme (project No. FEKT-S-23-8430). Author acknowledges Assoc. Prof. Jan Barta, Ph.D. as my research supervisor for his valuable comments and help while writing this paper.

approach not very efficient. It is often necessary to make compromises in the assembly of the electromagnetic model, thus, reducing the overall accuracy.

One solution that directly addresses this issue is the machine learning deployment, which is extensively used in many industrial fields. However, in the field of multi-objective optimization of electrical machine design utilizing machine learning, there is not much literature on this topic. Only few papers have attempted to use this approach [3], [4]. This work aims to find a suitable machine learning model that will serve as a basis for future work in the area of Multi-objective Optimization of electrical machine design utilizing machine learning.

II. INDUCTION MACHINE MODELLING METHODOLOGY

A. Surrogate Modeling

A surrogate model is a special case of supervised machine learning technique, defined as a Black-box function, that replaces all computationally expensive simulations with approximated functions [5]. Black-box functions generally focus on mapping inputs vs. outputs as accurately as possible; with no focus, on the exact behavior of the original modeled system from the physics point of view. The actual internal processes of black-box functions are often not transparent or known, especially as the complexity of the model increases. The indisputable advantage of surrogate models is the ability to calculate large amounts of data in a very short time.

A surrogate model is built using training data (or observations) and appropriate mathematical approximation techniques (or machine learning technique). The main advantage is the low number of training data, although it is essential to have a proper distribution of these data [6]. A suitable machine learning technique can further reduce the prediction error as well as the required number of training data. Several techniques are available [5], however, in the context of electrical machine modeling, the most appropriate is Gaussian process regression (GPR). GPR can be used with the so-called Bayesian optimization (BO), which is a global optimization technique for timeconsuming calculations utilizing specifically GPR models [7]. Therefore, this paper will focus on GPR models furthermore.



Fig. 1. Example of Gaussian process regression modeling of a simple 1D function.

B. Gaussian Process Regression

GPR (also known as Kriging) is a supervised non-parametric machine learning technique for surrogate modeling. GPR is specified by its mean and covariance function [8]:

$$f(x) \sim \mathcal{N}(\mu, \sigma) \sim \mathcal{N}(\mu, K), \tag{1}$$

where f(x) is the modeled function, \mathcal{N} is the Gaussian (normal) distribution, μ, σ is the mean and covariance function, which is also known as kernel K. For the GPR modeling, the *Matérn* [9] kernel was chosen, as it is able to handle non-smooth functions. It is defined by the smoothness coefficient, with 4 recommended values $(1/2; 3/2; 5/2; \infty)$. The smoothness value for $\nu = 5/2$ was selected, because it produces the best results, based on the authors experience:

$$k(x_{i}, x_{j}) = \sigma_{f}^{2} \left(1 + \frac{\sqrt{5}}{l} d(x_{i}, x_{j}) + \frac{5}{3l} d(x_{i}, x_{j})^{2} \right) \cdot \exp\left(-\frac{\sqrt{5}}{l} d(x_{i}, x_{j})\right)$$
(2)

where $d(x_i, x_j)$ is the distance between two samples, σ_f is the covariance scaling value, l is the lengthscale hyperparameter, and ν determines smoothness of the resulting function ($\nu > 0$). The definition of specific kernel is called a prior distribution. It is an initial belief about the underlying function. Afterwards, the Gaussian likelihood is usually defined:

$$p(y|f(x)) \sim \mathcal{N}(f(x), \sigma_n^2 I), \tag{3}$$

where p likelihood model, y are observations, σ_n is observed noise of the training data, and I identity matrix. Presenting the training data, which consists of inputs and outputs, to the GPR model, the prior distribution is updated to the posterior distribution. This reflects the updated beliefs about the true function. Additionally, the GPR models requires training, which is also called hyperparameter optimization. The described GPR modelling process is illustrated in Fig. 1.

GPR models are popular because of their flexibility, ability to model a wide range of regression tasks with small number of training data [10]. Furthermore, these models have the ability to measure the uncertainty of the prediction with the application of a likelihood function. The main disadvantages of this technique are the smaller training data sets, typically a maximum of 2,000 samples and a moderate number of input parameters, usually up to 20. The presented GPR model will be used for the modeling of selected IM.

III. MODELING OF INDUCTION MACHINE USING CONTINUOUS AND MIXED-INPUT MODELS

A. Selected Case Study Machine

For the case study, an industrial three-phase IM with a diecast aluminium squirrel cage was selected. The sketch of the machine is shown in Fig. 2. The stator and rotor sheet material is M470-50A, where Table I lists its key data.

The electromagnetic model of the selected machine, was created in *Ansys Maxwell Electronics Desktop*, which is FEMbased software. In order to simplify the construction of the electromagnetic model and reduce the total computational



Fig. 2. The sketch of the selected machine showing its design parameters considered for optimization.

TABLE I Key Induction Machine Parameters

Parameter	Unit	Value
Output torque	Nm	10
Torque ripple	%	9.2
Power factor	-	0.768
Electromagnetic efficiency	%	86.55
Speed	rpm	1459.7
Pole pairs	-	2
Nominal frequency	Hz	50
Rated voltage line-to-line	V	400 (Wye)
Flux density fundamental harmonic in the air-gap	Т	0.84

time, the machine was modeled in 2D space. The electromagnetic model was supplied from a pure sinusoidal voltage source, where the machine was loaded with its rated torque. The studied machine had to be analyzed using the transient analysis, due to its nature. To achieve accurate results of the model, a very fine mesh was created in the model, especially in regions like air-gap. The total number of elements in the mesh varied around 18,000 elements. The entire process of electromagnetic calculation was automated using the Python programming language, where the simulation time of a single design took approximately 2 hours. This electromagnetic model was subsequently used for the evaluation of the generated training data for GPR models.

B. Set Up of Gaussian Process Regression Models

A total of two types of datasets were considered for the case study, where the selected input parameters are illustrated in Fig. 2. The first dataset had a purely continuous nature of the input parameters. The second dataset was mixed, containing both continuous input parameters and one categorical input parameter (i.e., only integer values) for the *stator conductors* by rounding their values. Each dataset consisted of 800 samples, where a quasi-random Sobol sequence [11] was used to generate the data for input parameters. This technique has unique characteristics that make them particularly suitable for generating initial training data in scenarios which requires well-distributed sample points across the input space. Table II lists all used input and output parameters for the modeling with their corresponding selected or evaluated ranges.

For both datasets, three GPR models were considered. The first two had purely continuous kernels and were built using both generated training data. The third model was a mixed-input space model, which is specifically designed to model categorical input parameters and was built using only the mixed dataset. The GPytorch was employed for the construction of the GPR models, containing state-of-the-art algorithms in the field of GPR modeling [12]. The *Matérn* kernel with a smoothness coefficient $\nu = 2.5$ was considered for all models, where the models were trained using the gradient-based *L*-*FBGS* optimizer. Furthermore, given the categorical nature of the input parameter *stator conductors*, the *Categorical* kernel was used to process it, while the *Matérn* kernel was used for the other continuous input parameters. This resulted in the following architecture:

$$K((x_1, c_1), (x_2, c_2)) = K_{\text{Cont}_1}(x_1, x_2) + K_{\text{Cat}_1}(c_1, c_2) + K_{\text{Cont}_2}(x_1, x_2) \cdot K_{\text{Cat}_2}(c_1, c_2),$$
(4)

where K_{Cont} and K_{Cat} are kernels for continuous (x) and categorical (c) samples. To verify the accuracy of the GPR models, 10 % of the testing data from the overall dataset were randomly selected, where the other 90 % of the data were used as the training data. Since the two data sets are almost identical apart from the input parameter *stator conductors*, the same test samples were considered (selected by their index position).

 TABLE II

 Evaluated Training Data Ranges for All GPR Models

Variable	Parameter	Unit	Value
	Active length $(l_{\rm Fe})$	mm	100 - 160
Innuta	Air-gap length ($\delta_{\rm Fe}$)	mm	0.1 - 1
Inputs	Stator conductors (Q_s)	-	20 - 50
	Rotor radius (R_{ro})	mm	35 - 45
	Electromagnetic efficiency	%	2.8 - 87.5
Outputs	Power factor	-	0.2 - 0.92
Outputs	Torque ripple	%	3.8 - 162.8
	Flux density fundamental harmonic in the air-gap	Т	0.45 - 1.73

C. Results and Discussion

Four main metrics were considered for the verification of the GPR models accuracy. The first one is a plot of predicted vs. simulated data for training and testing data. If the samples are aligned in a perfect line, these would be socalled error-free models. Fig. 3 displays this plot for all three models. Additionally, a coefficient of determination (R^2) was listed, where its score indicates how well the model fits the data [13]. Samples for the training data exhibits error close to zero, implying well-trained models. The most significant discrepancy is in the testing data, where the mixed-input model (Fig. 3c) shows by far the largest sample deviations for all output parameters. Both models with continuous kernels shows fairly similar results, but the model with mixed dataset (Fig. 3b) is little bit more accurate according to the R^2 score.

The second metric displays the distribution of the relative prediction error for all three models compared to the simulated data, as shown in Fig. 4 for the testing data. Here, a more detailed distribution of relative errors can be seen, where the mixed-input model has a noticeably higher value of error (even more than 10 %) and its variance through all output parameters. The other two models are quite similar again. The model with continuous kernel and mixed dataset has the lowest relative error and its variance.

The third metric was a comparison of the models based on statistical analysis of the predicted testing data. This analysis consisted of evaluation of the three most commonly used statistical coefficients: Mean squared error (MSE), Root mean squared error (RMSE), and Mean absolute error (MAE). These coefficient are generally used to evaluate the accuracy of surrogate models and are useful for comparison of multiple models between each other [14]. The analysis results are presented in Table III. Minor model-to-model nuances are apparent here and the analysis confirms the previously established accuracy results for each model. Overall, the model with continuous kernels and mixed dataset yields noticeably the lowest values of the evaluated metrics, indicating its highest accuracy.

The fourth and final metric is the display of the learning curve for the training data of all three models. A learning curve graphically represents the model's performance over time or across varying sizes of the training dataset. It is particularly useful for identifying how well a model benefits from adding more training data and for diagnosing whether



Fig. 3. Comparison of GPR model accuracy results of all modeled output parameters for (a) continuous model (with continuous datasets), (b) continuous model (with mixed datasets), and (c) mixed-input model, with listed corresponding coefficients of determination (R^2) for training and testing data.



Fig. 4. Comparison of continuous GPR model with continuous dataset (1st model), continuous GPR model with mixed dataset (2nd model), and mixed-input GPR model (3rd model), using evaluated relative errors for testing data, highlighting the error differences for the corresponding value of each output parameter.

the model is suffering from overfitting or underfitting [15]. For the learning curve display, MSE was evaluated on the varying training dataset (randomly selected from the initial training data which varied in size from 10 % to 100%) and the testing dataset (constant 80 samples), with the results shown in Fig. 5. The MSE score for all testing datasets is practically constant, potentionaly meaning that the models does not have enough representative data to capture the true patterns and complexity of the modeled function. In the case of the mixed-input model, it has the highest MSE scores for training and testing data. Moreover, it has higher score for testing data compared to the training data. This indicates high bias of the modeled, which means it is underperforming and leads to poor generalization. In order to fix this, the model would require more training data. The other two models have similar results, where yet again the model with continuous kernel and mixed datasets

TABLE III

EVALUATED STATISTICAL METRICS FOR THE CONTINUOUS MODEL WITH CONTINUOUS DATASETS (CONT. 1), THE CONTINUOUS MODEL WITH MIXED DATASETS (CONT. 2), AND THE MIXED-INPUT MODEL (CAT.)

		Parameter				
Coefficient	Model	Electromagnetic efficiency	Power factor	Torque ripple	Flux density in the air-gap	
	Cont. 1	3.25e-2	6.66e-6	7.22	1.39e-5	
$(\%^2 - \%^2 T^2)$	Cont. 2	2.26e-2	6.77e-6	3.35	1.16e-5	
(,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	Cat.	9.92e-1	1.1e-4	29.35	6.66e-5	
	Cont. 1	1.8e-1	2.58e-3	2.69	3.73e-3	
RMSE (%, -, %, T)	Cont. 2	1.5e-1	2.6e-3	1.83	3.41e-3	
(,0, 1, 10, 1)	Cat.	9.96e-1	1.05e-2	5.42	8.16e-3	
MAE (%, -, %, T)	Cont. 1	1.22e-1	1.65e-3	1.43	2.53e-3	
	Cont. 2	9.58e-2	1.6e-3	1.15	2.55e-3	
	Cat.	6.84e-1	6.44e-3	3.34	4.11e-3	



Fig. 5. Comparison of continuous GPR model with continuous dataset (1st model), continuous GPR model with mixed dataset (2nd model), and mixed-input GPR model (3rd model), for testing on their MSE value with different sizes of training and constant testing datasets.

performs slightly better. The MSE score of this model for the training data stabilizes around the value of the MSE score for the testing data. This indicates a very good fit on the training data, and thus potentially, this model does not need more data. In the case of the model with continuous kernel and continuous dataset, the MSE score for the training data is slightly higher than that for the testing data. This may indicate a good model fit, but also that the model may require more training data, which would need to be verified by adding more samples.

IV. CONCLUSION

In this paper, two approaches for modeling IM utilizing machine learning have been presented: continuous and mixedinput models. GPR with Matérn and categorical kernel were used as surrogate models, which are able to accurately model functions or complex problems using only a few training data samples. The results of three models on two datasets were compared. It was shown that the GPR model with pure continuous kernel performs significantly better compared to the mixed-input GPR model. Moreover, if some of the input parameter is defined as categorical (i.e., integer only), the GPR model may even perform better, as shown in the paper. The obtained results are mainly relevant for multi-objective BO electrical machine design, where selecting the correct GPR model is crucial. A suitable GPR model is able to significantly accelerate and make the whole process of BO more efficient. The next step will be to test the selected model on multiobjective BO of electrical machine design, where this method can be applied to any type of machine.

REFERENCES

- J. Bacher, F. Waldhart, and A. Muetze, "3-d fem calculation of electromagnetic properties of single-phase induction machines," *IEEE Transactions on Energy Conversion*, vol. 30, no. 1, pp. 142–149, 2015.
- [2] T. D. Strous, X. Wang, H. Polinder, and J. A. B. Ferreira, "Finite element based multi-objective optimization of a brushless doubly-fed induction machine," in 2015 IEEE International Electric Machines & Drives Conference (IEMDC), 2015, pp. 1689–1694.
- [3] S. Zhang, S. Li, R. G. Harley, and T. G. Habetler, "An efficient multiobjective bayesian optimization approach for the automated analytical design of switched reluctance machines," in 2018 IEEE Energy Conversion Congress and Exposition (ECCE), 2018, pp. 4290–4295.
- [4] M. C. Huber, M. Fuhrländer, and S. Schöps, "Multi-objective yield optimization for electrical machines using gaussian processes to learn faulty design," *IEEE Transactions on Industry Applications*, vol. 59, no. 2, pp. 1340–1350, 2023.
- [5] K. McBride and K. Sundmacher, "Overview of surrogate modeling in chemical process engineering," *Chemie Ingenieur Technik*, vol. 91, no. 3, pp. 228–239, 2019. [Online]. Available: https://onlinelibrary.wiley.com/doi/abs/10.1002/cite.201800091
- [6] W. A. Pruett and R. L. Hester, "The creation of surrogate models for fast estimation of complex model outcomes," *PLOS ONE*, vol. 11, no. 6, pp. 1–11, 06 2016. [Online]. Available: https://doi.org/10.1371/journal.pone.0156574
- [7] B. Shahriari, K. Swersky, Z. Wang, R. P. Adams, and N. de Freitas, "Taking the human out of the loop: A review of bayesian optimization," *Proceedings of the IEEE*, vol. 104, no. 1, pp. 148–175, 2016.
- [8] V. L. Deringer, A. P. Bartók, N. Bernstein, D. M. Wilkins, M. Ceriotti, and G. Csányi, "Gaussian process regression for materials and molecules," *Chemical Reviews*, vol. 121, no. 16, pp. 10073–10141, 2021, pMID: 34398616. [Online]. Available: https://doi.org/10.1021/acs.chemrev.1c00022
- [9] E. Schulz, M. Speekenbrink, and Α. Krause. "A tutorial on gaussian process regression: Modelling, exploring, exploiting functions," Journal of Mathematical and Psvchology, vol. 85, pp. 1–16, 2018. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0022249617302158
- [10] M. Liu, G. Chowdhary, B. Castra da Silva, S.-Y. Liu, and J. P. How, "Gaussian processes for learning and control: A tutorial with examples," *IEEE Control Systems Magazine*, vol. 38, no. 5, pp. 53–86, 2018.
- [11] W. J. Morokoff and R. E. Caflisch, "Quasi-random sequences and their discrepancies," *SIAM Journal on Scientific Computing*, vol. 15, no. 6, pp. 1251–1279, 1994. [Online]. Available: https://doi.org/10.1137/0915077
- [12] J. R. Gardner, G. Pleiss, D. Bindel, K. Q. Weinberger, and A. G. Wilson, "Gpytorch: Blackbox matrix-matrix gaussian process inference with gpu acceleration," in *Advances in Neural Information Processing Systems*, 2018.
- [13] O. Renaud and M.-P. Victoria-Feser, "A robust coefficient of determination for regression," *Journal of Statistical Planning and Inference*, vol. 140, no. 7, pp. 1852–1862, 2010. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0378375810000194
- [14] A. Botchkarev, "A new typology design of performance metrics to measure errors in machine learning regression algorithms," *Interdisciplinary Journal of Information, Knowledge, and Management*, vol. 14, p. 045–076, 2019. [Online]. Available: http://dx.doi.org/10.28945/4184
- [15] T. Viering and M. Loog, "The shape of learning curves: A review," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 6, pp. 7799–7819, 2023.

Feastibility study and comparative study of air breathing electric propulsion systems operating in very low Earth orbit conditions

Marek Šťastný (1,3), Tomáš Dytrych (3), Vladimír Dániel (5), Kryštof Mrózek (2,3,4), Adam Obrusník (2,3,4)

⁽¹⁾ Department of Theoretical and Experimental Electrical Engineering,

Faculty of Electrical Engineering and Communication,

Brno University of Technology, Brno, Czech Republic

⁽²⁾ Department of Physical Electronics, Faculty of Science, Masaryk University, Brno, Czech Republic

⁽³⁾ SpaceLabEU, Zlatniky, Czech Republic

(4) PlasmaSolve s.r.o., Brno, Czech Republic

⁽⁵⁾ VZLU a.s., Prague, Czech Republic

263124@vut.cz

Air Breathing Electric Propulsion (ABEP) systems offer a promising solution to extending the lifetime of Very Low Earth Orbit (VLEO) missions by using residual atmospheric particles as a propellant. ABEP systems need to provide sufficient thrust to compensate for substantial aerodynamic drag present in VLEO environments. The feasibility of operating a hypothetical ABEP system of a given geometry at a given altitude is assessed via Direct Simulation Monte Carlo (DSMC) based on the following parameters: the mean pressure in the ionization chamber, compression factor and drag force acting upon the surface of the given geometry. Atmospheric models were used for reference to ensure realistic VLEO-like conditions. The comparative study is performed using a Global Plasma Model (GPM). A GPM can calculate the volume-averaged quantities of plasma systems with complex physics and reaction kinetics. The results of GPM are compared with a breadboard model of an Electron Cyclotron Resonance (ECR) ABEP system constructed by the Czech aerospace research institute (VZLU a.s.).

Very Low Earth Orbit (VLEO), Direct Simulation Monte Carlo (DSMC), Air Breathing Electric Propulsion (ABEP), Intake, Global Plasma Model (GPM)

I. INTRODUCTION

Carrying out missions in Very Low Earth Orbit (VLEO, orbits with mean altitude below 450 km) has several benefits over the standard operations in higher altitudes. However, atmospheric conditions at VLEO present challenges to the spacecraft (S/C) in the form of thermospheric winds, unpredictable solar activity and predominantly a substantial aerodynamic drag which must be compensated by a propulsion system to maintain the orbit and to reach desirable mission lifetimes. Air Breathing Electric Propulsion (ABEP) systems present a theoretical solution by using residual atmospheric particles as a propellant, however, neither of the concepts was, to date, demonstrated to work in orbit. A general schematic of a platform fitted with ABEP system presented in the Fig. 1 is described as follows:



Fig. 1. Schematic of a platform with ABEP system [1].

An intake is a device capable of collecting atmospheric particles, compressing them, and driving them into the thruster. Intakes vary in geometry and size, but all concepts share a common end goal: to scatter the atmospheric particles into the accelerator, which contains an ionization chamber capable of ionizing the compressed gas and accelerating it to generate thrust. The energy required for plasma ignition is meant to be supplied from solar arrays, attached to the S/C.

Current state of the art ABEP concepts feature RAMjet Electric Propulsion (RAM-EP) designed by SITAEL which was the first time worldwide experimentally validated system using an on-ground VLEO representative environment. They registered a successful ignition but were unable to compensate the atmospheric drag [2]. The Institut für Raumfahrtsysteme (IRS) in Germany provided an ABEP concept utilizing RF helicon-based Inductive Plasma Thruster (IPT). In early 2021 the IRS managed to achieve a successful ignition with satisfying preliminary results [3]. Other noteworthy concept originates in Japan Aerospace Exploration Agency (JAXA) employing Air Breathing Ion Engine (ABIE) with Electron Cyclotron Resonance (ECR). This system did not undergo any experimental verification but a study from 2012 confirmed the possibility of plasma ignition inside the ECR in 200 km of altitude [4]. Additionally in 2012, a different approach to ABEP was presented by BUSEK, focusing on the application of ABEP in extra-terrestrial atmospheres, mainly on Mars [5]. Recent studies show that the geometry and material of an intake contribute significantly to the results. Both University of Stuttgart [6] and the University of Colorado [7] carried out examinations of a parabolic intake made of specular reflecting material and concluded that it is significantly more efficient than material with diffuse reflection. Furthermore, the geometry of front-fitted parallel grid ducts is also under investigation. SITAEL features a split-ring configuration in their concepts. Meanwhile College of Aerospace Science and Engineering [8] compared the split-ring configuration with a honeycomb configuration and concluded that the latter configuration is more suitable for better intake performance. Therefore, in our own investigation of gas compression, we pick the two of the most prevalent geometrical configurations which appeared in recent years.

This article concerns the feasibility of ABEP systems by comparing an experimental setup with simulated data. Our goal is to determine whether plasma in the experimental model ignites at low pressures which are comparable to VLEO environment. This process is captured in the following chapter. The possibility of plasma ignition gives rise to the chapter called gas compression study, in which we examine whether such pressures can be obtained using our representative intake geometries. The conclusion features discussion regarding the results and possible future directions of the project.

II. PLASMA IGNITION STUDY

As the atmospheric pressure at VLEO is low, it is necessary to turn to plasma sources which can operate in such an environment. One of the viable sources is the ECR plasma source, which is used in the experimental breadboard model constructed by VZLU a.s. ECR is believed to compliment low operational pressures found at VLEO as the prevalent method of energy transfer happens between electrons and the high frequency electrodes. Electrons in the discharge chamber follow helical trajectories under the influence of the external magnetic field from coils. The resonance (i.e. maximum power transfer) occurs when the electron cyclotron frequency arising from their circular motion matches the frequency of the electric field created by the high frequency electrodes.



Fig. 2. A schematic picture of the laboratory model of the ABEP thruster [1].

Schematics of the laboratory model of ABEP thruster seen in the Fig. 2 features a glass tube, a triad of extraction grids and the ECR plasma source. Experimental results were validated against Global Plasma Model (GPM) results using the kinetic system of N2/O2 mixture. We are using results from previously developed GPM of magnetized high frequency plasma source. The following paragraph succinctly describes the model as the full details can be found in [1]

At the crux of the GPM lies a set of balance equations, which are solved for all species except electrons (assuming quasineutrality) and for energy. Theses balance equations involve source terms which determine the temporal change in the key variable. The specie density changes via different mechanisms such as kinetic reactions, inflow, outflow, wall loss and wall gain. The energy balance of electrons and neutral gas particles involve different loss or gain channels. In electron energy equation, energy is introduced to the plasma through oscillating RF electric field and is, at first, transferred only to electrons. Different source terms of the energy balance equation describe the process of energy dissipation into the system. For a successful computation, the Electron Energy Distribution Function (EEDF) must be known. In ECR plasma, the EEDF is not Maxwellian and can be obtained for example through Boltzmann equation solver BOLSIG+ [9]. Thus, we can obtain electron temperature, whereas neutral gas energy equation allows us to obtain temperature of heavy particles. The GPM can be understood as a 0D model since it outputs only volumeaveraged quantities. However, this is balanced out by low computational times so it can serve as a fast and efficient tool for estimating plasma properties and observing relationships among key variables.

As was stated above, GPM is based on a reaction kinetics. For example, the model used in this paper contains more than 600 unique reactions and processes. The input of GPM consists of species composition and initial conditions and parameters such as initial density of all species, initial pressure, energy from source and volume of the ionization chamber. The output of interest is plasma density from which we determine whether the plasma ignition occurs. As can be seen in [1], the GPM exhibits good accordance with the experimental data provided by VZLU a.s. This is best showcased in the Fig. 3 in which we can see the validation of GPM against the experiment.



Fig. 3. Experimental validation of GPM: Two blue curves represent two independent measurements; orange curve is a result of GPM simulation [1].

The GPM in conjunction with experimental data also predicts that the plasma in ECR with reasonable dimensions would ignite only above 5 mPa. With such information in mind, we conducted a gas compression study in order to pinpoint if such pressures are achievable at VLEO altitudes.

III. GAS COMPRESSION STUDY

Most of the theoretical results of the performance of any ABEP systems mentioned in the introduction section were gathered solely via simulations. Similarly, our investigation of ABEP intake systems was conducted using primarily Direct Simulation Monte Carlo (DSMC) method as DSMC corresponds well with the molecular flow regime of particles expected to be found at the intake of the S/C. Our DSMC cases were conducted with dsmcFoam+, a solver which is a part of the OpenFOAM open-source software. The DSMC solver requires precise information about the simulated environment, which is imposed at the model's boundary condition. As the atmosphere in VLEO changes based on a variety of different parameters, an atmospheric empirical model with the ability to provide reliable data must be employed. We utilize the NRLMSISE-00 model which is considered a staple model in space research and simulation [10]. NRLMSISE-00 is used primarily to determine number densities of dominant species in VLEO. Additional benefit of NRLMSISE-00 is its ability to provide information regarding density of highly reactive oxygen generated by photochemical dissociation of ozone molecules. This reactive oxygen proves to be a non-negligible factor while considering the lifespan of ABEP components, as it has high corrosive effect especially contributing detrimentally to extraction grids. A representative graph of number densities of species in VLEO can be seen in the following Fig. 4:



Fig. 4. Number densities of dominant species and pressure as a function of altitude [1].

By employing DSMC method we were able to compute not only atmospheric drag F_D affecting the outer walls of the S/C but also mean pressure p_{mean} in the ionization region. We chose two different geometries, as seen in Fig. 5. As stated in the introduction, we decided on popular intake shapes. Firstly, an S/C consisting of a parabolic shaped intake leading directly to the ionization chamber with a 90% reflecting grid fitted at the end was envisioned. The second geometry is more complex in design, consisting of a split-ring intake, not unlike in concepts of SITAEL, and cone shaped compartment leading to the ionization chamber that is also fitted with a grid with the same reflection coefficient. Both ionization chambers were modelled with the experimental breadboard model in mind to share comparable dimensions.



Fig. 5. General parabolic and split-ring configuration.

The first set of simulations determine at which altitude the selected geometry performs the best. The simulations are performed at five different altitudes between 150 km and 250 km for both geometries.



Fig. 6. Total gas pressure per cell – parabolic configuration at 150 km altitude.

Fig. 6 shows a total gas pressure inside and outside of the S/C. The parameters of interest are F_D , p_{mean} , the ratio between the two and lastly, a compression factor β representing the ratio between mean pressure in the ionization chamber and the pressure outside of the S/C. Considering all these parameters, we can determine the optimal altitude to be at 175 km. Although we barely reach the desirable 5 mPa at this altitude, the drag force remains optimal. Furthermore, both parameters can be further improved by changing the geometry of the S/C or by employing front-fitted ducts or by using different materials. The following Tab. 1 shows the best results for both geometries:

TABLE I. COMPARISON OF BOTH GEOMETRIES AT 175 KM OF ALTITUDE

Configuration	p _{mean} [mPa]	F _D [mN]	$\frac{p_{\text{mean}}}{F_{\text{D}}} [\mathrm{m}^{-2}]$	β
Parabolic	5.719	0.808	7.078	55.77
Split-ring	4.588	1.031	4.450	44.75

The second set of simulations is performed at the optimal altitude by scaling key elements or changing the number of elements in both geometries. The parabolic configuration allows for scalable parabolic intake as well as scalable ionization chamber. The split-ring configuration features variable parameters capable of altering the number of lamellas per ring or the number of rings as well as scalable intake length. For example, the split-ring configuration performed better with shorter intakes, as seen in the following Tab. 2.

TABLE II. PERFORMANCE PARAMETERS FOR DIFFERENT LENGTHS OF SPLIT-RING CONFIGURATIONS

x [cm]	15	22.5	30
p _{mean} [mPa]	5.390	5.084	4.588
$F_{\rm D}$ [mN]	0.972	1.059	1.031
$\frac{p_{\text{mean}}}{F_{\text{D}}} [\text{m}^{-2}]$	5.545	4.801	4.450
β	52.56	49.58	44.75

As far as the parabolic configuration is concerned, the geometrical variations reveals that our base geometry performs the best. This is the one pictured in the Fig. 5. The results are further discussed in the conclusion section.

IV. CONCLUSION AND FUTURE WORK

A. Plasma Ignition Study

Studies involving GPM provided us with an important value below which plasma ignition is no longer achievable. In conjunction with experiment carried out by VZLU a.s. we concluded that 5 mPa is the lowest pressure at which plasma can be ignited in the breadboard model's ionization chamber.

The GPM method is a swift tool capable of offering good estimates of plasma properties. This evokes the idea of coupling the DSMC method with GPM. Implementing such conjunction of both methods would allow us to accurately compute the drag force, calculated using DSCM method, as well as the thrust, calculated by GPM (based on pressure and gas composition from DSMC). Deciding whether the thrust is net positive by directly comparing the two variables should lead to better understanding of the feasibility of the ABEP concept.

B. Gas Compression Study

Through DSMC, a study of performance parameters such as mean pressure in the ionization chamber, drag force acting on the surface of the spacecraft or compression factor at various altitudes was performed. Based on the mean pressure and the ratio between the mean pressure and the drag force, we determined that the optimal altitude at which further studies would be conducted is 175 km. Both mean pressure and the drag force decreased with increasing altitude, which was to be expected. The mean pressure to drag force ratio at this altitude yielded the most favorable results: the mean pressure was high enough for the ignition of plasma in the ionization chamber to still be theoretically possible. Once the optimal altitude was selected, a series of geometrical studies was performed to determine the optimal geometry in terms of best performance parameters. Concerning the parabolic intake, we simulated a geometry with parabolic intakes of various lengths. The results suggested a better performance for longer intakes. Reducing the overall length of the split-ring intake to be comparable with the parabolic configuration, led to better performance as well. Subsequently, we picked the geometrical variations that performed best for both geometries and compared them. Surprisingly, both optimal geometries performed similarly, nevertheless the parabolic configuration showed slightly better compression factor and was less affected by the atmospheric drag.

C. Future Work

Apart from further improvements to the experimental device, our current research gravitates towards the coupling of the DSMC method and GPM, as was stated prior. Additional plans include more in-depth look into the effects of atomic oxygen on the S/C. More specifically, the ABEP systems are usually imagined to be fitted with accelerating grids which pose an unstable element in their implementation. The abundance of highly reactive atomic oxygen proves detrimental to many ABEP electrical components, including the grids. Thus, comprehensive research including experimental data of the grid lifetime is currently in preparation.

Lastly, a more detailed insight into plasma processes and composition is required as the GPM cannot solely comprehend all underlying phenomena. The optical emission spectroscopy arises as a suitable candidate for additional plasma diagnostics which can, as well, be used to validate the GPM results.

ACKNOWLEDGMENT

A portion of the work presented in this paper concerned with the gas compression study was carried out by the main author as an undergraduate at the Faculty of Science, Masaryk University in Brno, Czech Republic. We extend our thanks to Masaryk University as well as to PlasmaSolve s.r.o. for providing necessary software utilized in conducting numerical simulations.

REFERENCES

- K. Mrózek, T. Dytrych, P. Moliš, V. Dániel, and A. Obrusník, "Global plasma modeling of a magnetized high-frequency plasma source in lowpressure nitrogen and oxygen for air-breathing electric propulsion applications," *Plasma Sources Sci. Technol.*, vol. 30, no. 12, 2021, doi: 10.1088/1361-6595/ac36ac.
- [2] T. Andreussi et al., "Development status and way forward of SITAEL's air-breathing electric propulsion," *AIAA Propuls. Energy Forum Expo.* 2019, no. August, pp. 1–22, 2019, doi: 10.2514/6.2019-3995.
- [3] F. Romano, Y. Chan, G. Herdrich, and C. Traub, "Design, Set-Up, and First Ignition of the RF Helicon-based Plasma Thruster," 2021, no. March.
- [4] P. Zheng, J. Wu, Y. Zhang, and B. Wu, "A comprehensive review of atmosphere-breathing electric propulsion systems," *International Journal* of Aerospace Engineering, vol. 2020, no. 4. 2020, doi: 10.1155/2020/8811847.
- [5] N. S. Symposium et al., "Atmospheric Breathing Electric Thruster for Planetary Exploration," pp. 1–14, 2012.
- [6] F. Romano et al., "Intake design for an Atmosphere-Breathing Electric Propulsion System (ABEP)," *Acta Astronaut.*, vol. 187, no. June, pp. 225– 235, 2021, doi: 10.1016/j.actaastro.2021.06.033.
- [7] S. W. Jackson and R. Marshall, "Conceptual design of an air-breathing electric thruster for CubeSat applications," J. Spacecr. Rockets, vol. 55, no. 3, pp. 632–639, 2018, doi: 10.2514/1.A33993

- [8] P. Zheng, J. Wu, Y. Zhang, and Y. Zhao, "Design and Optimization of vacuum Intake for Atmosphere-Breathing electric propulsion (ABEP) system," *Vacuum*, vol. 195, no. October 2021, p. 110652, 2022, doi: 10.1016/j.vacuum.2021.110652.
- [9] G. J. M. Hagelaar and L. C. Pitchford, "Solving the Boltzmann equation to obtain electron transport coefficients and rate coefficients for fluid

models," *Plasma Sources Sci. Technol.*, vol. 14, no. 4, pp. 722–733, 2005, doi: 10.1088/0963-0252/14/4/011.

[10] J. M. Picone, A. E. Hedin, D. P. Drob, and A. C. Aikin, "NRLMSISE-00 empirical model of the atmosphere: Statistical comparisons and scientific issues," J. Geophys. Res. Sp. Phys., vol. 107, no. A12, pp. 1–16, 2002, doi: 10.1029/2002JA009


BRNO UNIVERSITY OF TECHNOLOGY FACULTY OF ELECTRICAL ENGINEERING AND COMMUNICATION TECHNICKÁ 3058/10, 616 00 BRNO, CZECH REPUBLIC



www.eeict.cz



www.fekt.vut.cz

