



ISSN 1310-8271

JOURNAL

**OF THE TECHNICAL UNIVERSITY - SOFIA
PLOVDIV BRANCH, BULGARIA**

**FUNDAMENTAL SCIENCES
AND
APPLICATIONS**

Volume 24, 2018

**Journal of the Technical University - Sofia
Plovdiv branch, Bulgaria
“Fundamental Sciences and Applications”**

EDITORIAL BOARD

EDITOR-in-chief

Prof. Michail Petrov, PhD

Honorary EDITOR-in-chief

Prof. Marin Nenchev, DSc Eng., DSc Phys.

Journal Editorial Board

Prof. Valyo Nikolov, PhD

Prof. Andon Topalov, PhD

Prof. Veselka Boeva, PhD

Prof. Galidia Petrova, PhD

Prof. Grisha Spasov, PhD

Prof. Angel Zumbilev, PhD

Prof. Dobrin Seizinski, PhD

Prof. Teofil Iamboliev, PhD

Programme Editorial Board

Prof. Venelin Zhivkov, DSc

Prof. Georgi Andreev, DSc

Prof. Georgi Totkov, DSc

Prof. Emil Nikolov, DSc

Prof. Ivan Iachev, DSc

Prof. Marin Hristov, PhD

Prof. Ognian Nakov, PhD

Prof. Marc Himbert DSc

Prof. Tinko Eftimov DSc

Acad. Yuriy Kuznietsov DSc

Journal Scientific Secretary: Sevil Ahmed, PhD



CONTENTS

1. PETR MARCON, FRANTISEK ZEZULKA, ZDENEK BRADAC	7
TERMINOLOGY OF INDUSTRY 4.0	
2. ZDENEK BRADAC, FRANTISEK ZEZULKA, PETR MARCON	13
TECHNICAL AND THEORETICAL BASIS OF INDUSTRY 4.0 IMPLEMENTATION	
3. KRASIMIRA STOILOVA, TODOR STOILOV, MIROSLAV VLADIMIROV	19
RESOURCE ALLOCATION BY PORTFOLIO OPTIMIZATION	
4. RADOSLAV HRISCHEV	25
PLANNING AND IMPLEMENTATION OF THE ERP SYSTEM IN PACKAGING PRODUCTION. PRACTICAL ASPECTS	
5. BORISLAV RUSENOV, ALBENA TANEVA, IVAN GANCHEV, MICHAIL PETROV	29
SYSTEM DEVELOPMENT FOR MACHINE VISION INTELLIGENT QUALITY CONTROL OF LOW VOLTAGE MINIATURE CIRCUIT BREAKERS	
6. STOITCHO PENKOV, ALBENA TANEVA, MICHAIL PETROV, KIRIL HRISTOZOV, ROBERT KAZALA	37
CONTROL SYSTEM APPLICATION USING LORAWAN	
7. AHMET SENPINAR	43
THE IMPORTANCE OF SOLAR ANGLES IN MPPT SYSTEMS	
8. MARGARITA DENEVA, PEPA UZUNOVA, VALKO KAZAKOV, VANIA PLACHKOVA, KAMEN IVANOV, MARIN NENCHEV AND ELENA STOYKOVA	49
USE OF INTERFERENCE WEDGED STRUCTURE AS AN ATTRACTIVE, SIMPLEST, LIGHT POWER DIVIDING ELEMENT	
9. ONUR KAYA , CETIN AKINCI, EMEL ONAL	55
CALCULATION OF DIELECTRIC DISSIPATION FACTOR AT VARIABLE FREQUENCIES OF MODEL TRANSFORMER	
10. SARANG PATIL	61
TWO NEW TYPES OF COMPACT ULTRA-WIDE BANDWIDTH ANTENNAS FOR 3-AXIS RF FIELD STRENGTH MEASUREMENT.	
11. SVETLIN STOYANOV	67
AN EXPERIMENTAL SETUP FOR DETERMINATION OF THE RESONANT FREQUENCIES OF A MECHANICAL FRAME STRUCTURE	
12. NIKOLA GEORGIEV	71
STUDYING AN AXIAL GENERATOR WITH ROTATING MAGNETS IN ITS STATOR WINDINGS	

13. VASIL SPASOV, IVAN KOSTOV, IVAN HADZHIEV	75
IMPLEMENTATION OF A NOVEL FORCE COMPUTATION METHOD IN THE FEMM SOFTWARE	
14. IVAN HADZHIEV, DIAN MALAMOV, VASIL SPASOV, DIMITAR NEDYALKOV	81
STUDYING THE ELECTRICAL AND THERMAL FIELDS, PRODUCED BY THE CURRENT IN THE INSULATION OF A MIDDLE VOLTAGE CABLE	
15. ERDİNÇ UZUN, TARIK YERLİKAYA, OĞUZ KIRAT	87
COMPARISON OF PYTHON LIBRARIES USED FOR WEB DATA EXTRACTION	
16. ERDİNÇ UZUN, TARIK YERLİKAYA, OĞUZ KIRAT	93
OBJECT-BASED ENTITY RELATIONSHIP DIAGRAM DRAWING LIBRARY: ENTRELJS	
17. GÜNGÖR YILDIRIM, YETKIN TATAR.....	99
AN ALTERNATIVE EXECUTION MODEL FOR OPTIMUM BIG DATA HANDLING IN IoT-WSN CLOUD SYSTEMS	
18. TIHOMIR TENEV, DIMITAR BIROV	105
SECURITY PATTERNS FOR MICROSERVICES LOCATED ON DIFFERENT VENDORS	
19. MEHMET VURAL, MUHARREM TUNCAY GENÇOĞLU.....	109
EMBEDDED AUDIO CODING USING LAPLACE TRANSFORM FOR TURKISH LETTERS	
20. VENETA ALEKSIEVA, SVETOSLAV SLAVOV.....	117
MANAGED ACTIVE DIRECTORY IN DIRECTORY-AS-A-SERVICE	
21. TONY KARAVASILEV, ELENA SOMOVA.....	123
OVERCOMING THE SECURITY ISSUES OF NOSQL DATABASES	
22. DANIEL TRIFONOV, HRISTO VALCHANOV	129
VIRTUALIZATION AND CONTAINERIZATION SYSTEMS FOR BIG DATA	
23. DIMITAR GARNEVSKI, PETYA PAVLOVA	133
PERFORMANCE ESTIMATION OF PARALLEL APPLICATION FOR SOLAR IMAGES PROCESSING	
24. VLADIMIR DIMITROV, DENITSA TSONINA	137
PROGRAMMING TOOLS FOR RELIABLE USER AUTHENTICATION ACROSS MULTIPLE SERVERS	
25. ATANAS KOSTADINOV	143
TESTING AND DIAGNOSTICS OF COMPUTER SYSTEMS	
26. TEODORA HRISTEVA, MARIA MARINOVA.....	147
USING GRAPHIC PROCESSING UNITS FOR IMPROVING PERFORMANCE OF DEEP LEARNING PROCESSES	
27. GEORGI ILIEV	151
MOBILE CROWDSOURCING FOR VIDEO STREAMING	
28. DIMITRE KROMICHEV	157
FAST GAUSSIAN FILTERING FOR SPEED FOCUSED FPGA BASED CANNY EDGE DETECTION COMPUTATIONS	

29. DIMITRE KROMICHEV	163
EFFICIENT COMPUTATION OF ORTHOGONAL GRADIENTS TO BE USED IN SPEED FOCUSED FPGA BASED CANNY EDGE DETECTION	
30. GEORGI PAZHEV	167
SURVEY OF METHODS AND TECHNOLOGIES FOR BUILDING OF SMART HOMES	
31. HATICE POLAT, KUBRA CELIK, HALUK EREN	175
NATURAL BARRIER WALL DESIGN FOR RADIATION AND NOISE PROBLEM ON THE ECOLOGICAL AIRPORTS	
32. ZHIVKO ILIEV, GEORGI DINEV	181
POSSIBILITIES FOR A COMPARATIVE STUDY OF THE VIBRATIONS IN A COMPLEX PENDULUM JAW CRUSHER	
33. TONI MIHOVA, VALENTINA NIKOLOVA-ALEXIEVA, MINA ANGELOVA	187
CONCEPTUAL IMPACT MODEL OF PROCESS MANAGEMENT ON THE MEAT INDUSTRY ENTERPRISES IN BULGARIA	
34. TODOR TODOROV, GEORGI TSANEV	193
ANALYSIS OF THE MEAN FIELD APPROXIMATION FOR TRAINING THE DEEP BOLTZMANN MACHINE	
35. ANGELINA POPOVA	199
CORROSION PROTECTION WITH INHIBITORS QUATERNARY AMMONIUM BROMIDES	
36. KALINA KAMARSKA	203
STUDY OF CORROSION BEHAVIOUR OF ALUMINIUM ALLOYS EN AW-2011 AND EN AW-2024	
37. GEORGI PASKALEV	207
VARIATIONAL PRINCIPLE FOR A CLASS OF NONLOCAL BOUNDARY VALUE PROBLEMS	

TERMINOLOGY OF INDUSTRY 4.0

PETR MARCON, FRANTISEK ZEZULKA, ZDENEK BRADAC

Abstract: *Industry 4.0 is set to be a fourth industrial and scientific – technologic revolution. However, due to the recent development, it will be more very sophisticated and very rapid evolution of scientific – technological background of existing and recent state of the art of existing industrial production, organization of work, new business models and business praxis in high-developed countries all over the world. The benefit of the Industry 4.0 to the existing technological development is not only in new IT technologies, but also in the organization of work in a massive implementation of new materials, life cycle management, and quality of work, safety and security features of industrial, technological and other production. Paper deals with terminology, which is necessary for understanding of goals, procedures, aspects and implementation of Industry 4.0 principle. Paper introduces in technologies, aspects, habits, sources, standards and theories and their application in a systematic and standardized way.*

Key words: *cyber physical system, digitalization, Industry 4.0, RAMI model*

1. Introduction

When the Industry 4.0 (I4.0) system of production are to be successfully implemented and are to bring expected and asked support in concurrency with economies without I4.0 features, it is necessary to understand, enhance, implement and integrate into modern enterprises of the future following technologies, aspects, habits, sources, standards and theories, their application and some others in a systematic and standardized way. Let us introduce you in the background terminology of the I4.0, hence in such terms and their content as Digitization, OPC UA, UML, RFID, Cloud and edge computing and control, Cyber - physical systems, Vertical and horizontal integration of control, IEC 62443, IEC 62264/IEC 61512, IEC 62890, and Standards for I4.0 in preparation.

Authors are persuaded that the first step in I4.0 future is a good understanding of the I4.0 process, state of the art of initial I4.0 ideas and following the step by step development works in correction of initial ideas and in standardization activities of the most appropriate procedures of the I4.0 systems implementation. Let us know to begin with short specification of above mentioned areas, their contents and description.

2. Industry 4.0 terminology

The following subchapters describe the basic terminology of I4.0.

2.1. Industry 4.0

Industry 4.0 are all activities related to new style of industrial production in smart future factories. It is said to be the 4th technical revolution. However, in the reality, the I4.0 is a very rapid evolution in many aspects of human style of living, particularly in activities connected with industrial production.

In the Europe, the Internet of Things (IoT) is sorted into the CIoT (Commercial Internet of Things) and the IIoT (Industrial Internet of Things). The CIoT abbreviation is not frequently used and the IoT is used for Internet of whichever things. But in the United States of America technical terminology, the IoT represents the all issue which are in the Europe covered by the I4.0 activities [1], [2].

2.2. Digitalization

Digitalization is the process of converting information into a digital (i.e. computer-readable) format, in which the information is organized into bits. The result is the representation of an object, image, sound, document or signal (usually an analog signal) by generating a series of numbers that describe a discrete set of its points or samples. The result is called digital representation or, more specifically, a digital image, for the object, and digital form, for the signal. In modern practice, the digitized data is in the form of binary numbers, which facilitate computer processing and other operations, but, strictly speaking, digitizing simply

means the conversion of analog source material into a numerical format; the decimal or any other number system that can be used instead [3,4].

From the I4.0 point of view, digitalization is the crucial topic. Digitalization of process variables, market and economy of industrial production has been started evolutionary thanks to digital control systems (PLC, DCS, embedded control systems, information and IT technologies, economy of production – MES and ERP system). In the new generation of industrial production, digitalization of information appears in the whole human activities. Therefore, it is also very important for new production and market models and production.

The I4.0 initial principle and idea goes out from very comprehensive digitalization of data from the whole production chain. It is non – systematically provided already in the existing production. The digitalization for the future fabric needs to be not only comprehensive more than the existing one, but it has to be provided in the systematic way to be used in the most appropriate way when it is needed, in the real time and without failure, non- correct interpretation and should be obtained (measured) only one time for more applications and use. A good example how to realize such a goal is in implementation existing and future generation of MESs (Manufacturing Execution System) and ERP systems.

2.3. OPC Unified Architecture (OPC UA)

The OPC UA is a machine-to-machine communication protocol for industrial automation developed by the OPC Foundation. Shortly OPC UA is an open standardized SW interface on highest communication levels in production control systems.

The Foundation's goal for OPC-UA was to provide a path forward from the original OPC communications model (namely the Microsoft Windows-only process exchange COM/DCOM) that would better meet the emerging needs of industrial automation. The original OPC is named OLE for Process Control, whereas OLE is Object Linking and Embedding. The original OPC is applied in different technologies such as in building automation, discrete manufacturing, process control and many others and is no more intended for the Microsoft Windows OS only, but it enables to include other data transportation technologies including Microsoft's .NET Framework, XML, and even the OPC Foundation's binary-encoded TCP format [5].

On the other hand the OPC UA differs significantly from its predecessor, Open Platform Communications (OPC). OPC UA better meets the emerging needs of industrial automation [1].

OPC UA shows distinguishing characteristics are:

- Focus on communicating with industrial equipment and systems for data collection and control.
- Open - freely available and implementable without restrictions or fees.
- Cross-platform - not tied to one operating system or programming language.
- Service-oriented architecture (SOA).
- Robust security.
- Integral information model, which is the foundation of the infrastructure necessary for information integration where vendors and organizations can model their complex data into an OPC UA namespace take advantage of the rich service-oriented architecture of OPC UA. There are over 35 collaborations with the OPC Foundation currently. Key industries include pharmaceutical, oil and gas, building automation, industrial robotics, security, manufacturing and process control [5].

Even for above mentioned features, OPC UA is very convenient for the Industry 4.0 information and communication infrastructure. It enables free, open, rapid, safety and security and at least soft real – time communication.

The first version of the Unified Architecture was released in 2006. The current version of the specification is on 1.03 (10 Oct 2015). The new version of OPC UA now has added publish/subscribe in addition to the client/server communications infrastructure [5].

2.4. UML

The Unified Modeling Language (UML) is a general-purpose, developmental, modeling language in the field of software engineering, that is intended to provide a standard way to visualize the design of a system. In 1997 UML was adopted as a standard by the Object Management Group (OMG), and has been managed by this organization ever since. In 2005 UML was also published by the International Organization for Standardization (ISO) as an approved ISO standard [6]. Since then the standard has been periodically revised to cover the latest revision of UML [3]. It is originally based on the notations of the Booch method, the object-

modeling technique (OMT) and object-oriented software engineering (OOSE), which it has integrated into a single language [4], [7], [8].

UML is not a development method by itself; however, it was designed to be compatible with the leading object-oriented software development methods of its time, for example OMT, Booch method, Objectory and especially RUP that it was originally intended to be used with when work began at Rational Software [7].

It is important to distinguish between the UML model and the set of diagrams of a system. A diagram is a partial graphic representation of a system's model. The set of diagrams need not completely cover the model and deleting a diagram does not change the model. The model may also contain documentation that drives the model elements and diagrams (such as written use cases) [8].

There are many type of diagrams sorted in:

Structural UML diagrams: Class diagram, Component diagram, Composite structure diagram, Deployment diagram, Object diagram, Package diagram, Profile diagram.

Behavioural UML diagrams: Activity diagram, Communication diagram, Interaction overview diagram, Sequence diagram, State diagram, Timing diagram, Use case diagram.

In UML, one of the key tools for behaviour modelling is the use-case model, caused by OOSE. Use cases are a way of specifying required usages of a system. Typically, they are used to capture the requirements of a system, that is, what a system is supposed to do. Simply, the Use case diagram – shows possible kinds of the use, it serves to the specification of users requirements in the analytical period of a system design [7], [8].

2.5. RFID [6]

Radio-frequency identification (RFID) uses electromagnetic fields to automatically identify and track tags attached to objects. The tags contain electronically stored information. Passive tags collect energy from a nearby RFID reader's interrogating radio waves. Active tags have a local power source (such as a battery) and may operate hundreds of meters from the RFID reader. Unlike a barcode, the tag need not be within the line of sight of the reader, so it may be embedded in the tracked object. RFID is one method for Automatic Identification and Data Capture (AIDC) [1].

RFID tags are used in many industries, for example, an RFID tag attached to an automobile during production can be used to track its progress through the assembly line; RFID-tagged pharmaceuticals can be tracked through warehouses; and implanting RFID microchips in livestock and pets allows for positive identification of animals.

Since RFID tags can be attached to cash, clothing, and possessions, or implanted in animals and people, the possibility of reading personally-linked information without consent has raised serious privacy concerns [6]. These concerns resulted in standard specifications development addressing privacy and security issues. ISO/IEC 18000 and ISO/IEC 29167 use on-chip cryptography methods for untraceability, tag and reader authentication, and over-the-air privacy. ISO/IEC 20248 specifies a digital signature data structure for RFID and barcodes providing data, source and read method authenticity. This work is done within ISO/IEC JTC 1/SC 31 Automatic identification and data capture techniques. Tags can also be used in shops to expedite checkout, and to prevent theft by customers and employees.

In 2014, the world RFID market was worth US\$8.89 billion, up from US\$7.77 billion in 2013 and US\$6.96 billion in 2012. This figure includes tags, readers, and software/services for RFID cards, labels, fobs, and all other form factors. The market value is expected to rise to US\$18.68 billion by 2026, [6].

In the Industry 4.0 environment the RFID chips will create very important information storages and very decentralized control elements. The Asset Administration Shell (AAS), the crucial element of the future industrial production, will be in many applications placed in a RFID chip. The RFID chip enables not only to carry initial information about what should be done with the production component, but it can carry the all information of the component during the whole production time.

2.6. Cloud and edge computing and control

Cloud computing is an information technology (IT) paradigm that enables ubiquitous access to shared pools of configurable system resources and higher-level services that can be rapidly provisioned with minimal management effort, often over the Internet. Cloud computing relies on sharing of resources to achieve coherence and economies of scale, similar to a public utility [6].

Third-party clouds enable organizations to focus on their core businesses instead of expending resources on computer infrastructure and maintenance [1]. Advocates note that cloud computing allows companies to avoid or minimize up-front IT infrastructure costs. Proponents also claim that cloud computing allows enterprises to get their applications up and running faster, with improved manageability and less maintenance, and that it enables IT teams to more rapidly adjust resources to meet fluctuating and unpredictable demand [1-3]. Cloud providers typically use a "pay-as-you-go" model, which can lead to unexpected operating expenses if administrators are not familiarized with cloud-pricing models [4].

Since the launch of Amazon EC2 in 2006, the availability of high-capacity networks, low-cost computers and storage devices as well as the widespread adoption of hardware virtualization, service-oriented architecture, and autonomic and utility computing has led to growth in cloud computing.

Edge computing is a way how to remove delay and overload of communication infrastructure in Industry 4.0 factories of future. It is clear, after 2 – 3 years history of first I 4.0 case studies, that Industry 4.0 components will not be mostly equipped by very powerful distributed computer systems (situated in AASs), as it was expected by the initial ideas of the Industry 4.0 infrastructure. It seems to be more efficient to provide computing for production purposes on the edge among physical and virtual domains (in a fog). The edge computing will be provided on the level of servers (private or public) and only one part of computing will be done in clouds.

2.7. Cyber - physical systems

A cyber-physical (also styled cyber physical) system (CPS) is a mechanism that is controlled or monitored by computer-based algorithms, tightly integrated with the Internet and its users. In cyber-physical systems, physical and software components are deeply intertwined, each operating on different spatial and temporal scales, exhibiting multiple and distinct behavioural modalities, and interacting with each other in a myriad of ways that change with context [1]. Examples of CPS include smart grid, autonomous automobile systems, medical monitoring, process control systems, robotics systems, and automatic pilot avionics [6], [10].

CPS involves transdisciplinary approaches, merging theory of cybernetics, mechatronics, design and process science [3],[11] The process control is

often referred to as embedded systems. In embedded systems, the emphasis tends to be more on the computational elements, and less on an intense link between the computational and physical elements. CPS is also similar to the Internet of Things (IoT), sharing the same basic architecture; nevertheless, CPS presents a higher combination and coordination between physical and computational elements [6].

Precursors of cyber-physical systems can be found in areas as diverse as aerospace, automotive, chemical processes, civil infrastructure, energy, healthcare, manufacturing, transportation, entertainment, and consumer appliances [6].

Unlike more traditional embedded systems, a full-fledged CPS is typically designed as a network of interacting elements with physical input and output instead of as standalone devices [7]. The notion is closely tied to concepts of robotics and sensor networks with intelligence mechanisms proper of computational intelligence leading the pathway. Ongoing advances in science and engineering will improve the link between computational and physical elements by means of intelligent mechanisms, dramatically increasing the adaptability, autonomy, efficiency, functionality, reliability, safety, and usability of cyber-physical systems.[8] This will broaden the potential of cyber-physical systems in several dimensions, including: intervention (e.g., collision avoidance); precision (e.g., robotic surgery and nano-level manufacturing); operation in dangerous or inaccessible environments (e.g., search and rescue, firefighting, and deep-sea exploration; coordination (e.g., air traffic control, war fighting); efficiency (e.g., zero-net energy buildings); and augmentation of human capabilities (e.g., healthcare monitoring and delivery) [5].

2.8. Vertical and horizontal integration of control

2.8.1. Vertical integration (VI)

In microeconomics and management, vertical integration is an arrangement in which the supply chain of a company is owned by that company [12]. Usually each member of the supply chain produces a different product or (market-specific) service, and the products combine to satisfy a common need. It is contrasted with horizontal integration, wherein a company produces several items which are related to one another. Vertical integration has also described management styles that bring large portions of the supply chain not only under a common ownership, but also into one corporation (as in the 1920s when the Ford

River Rouge Complex began making much of its own steel rather than buying it from suppliers) [12].

Vertical integration and expansion is desired because it secures the supplies needed by the firm to produce its product and the market needed to sell the product. Vertical integration and expansion can become undesirable when its actions become anti-competitive and impede free competition in an open marketplace. Vertical integration is one method of avoiding the hold-up problem. A monopoly produced through vertical integration is called a "vertical monopoly".

2.8.2. Horizontal Integration (HI)

HI is the process of a company increasing production of goods or services at the same part of the supply chain. A company may do this via internal expansion, acquisition or merger [1], [3], [6].

The process can lead to monopoly if a company captures the vast majority of the market for that product or service [3].

Horizontal integration contrasts with vertical integration, where companies integrate multiple stages of production of a small number of production units.

Benefits of horizontal integration to both the firm and society may include economies of scale and economies of scope. For the firm, horizontal integration may provide a strengthened presence in the reference market. It may also allow the horizontally integrated firm to engage in monopoly pricing, which is disadvantageous to society as a whole and which may cause regulators to ban or constrain horizontal integration [13].

2.9. Standards for I4.0

IEC 62264 is an international standard for enterprise-control system integration. This standard is based upon ANSI/ISA-95.

ANSI/ISA-95, or ISA-95 as it is more commonly referred, is an international standard from the International Society of Automation for developing an automated interface between enterprise and control systems. This standard has been developed for global manufacturers. It was developed to be applied in all industries, and in all sorts of processes, like batch processes, continuous and repetitive processes.

The objectives of ISA-95 are to provide consistent terminology that is a foundation for supplier and manufacturer communications provide consistent information models, and to provide

consistent operations models which is a foundation for clarifying application functionality and how information is to be used [1], [6], [14],

ANSI/ISA-95.00.01-2000, Enterprise-Control System Integration Part 1: Models and Terminology consists of standard terminology and object models, which can be used to decide which information, should be exchanged.

3. Conclusion

Paper brings a short comprehensive overview of terminology and specification of basic items concerned and available links among them to enable specialists from industry to understand what is for implementation of I 4.0 principles necessary. Author are persuade that the first step in Industry 4.0 implementation should be done in a very good understanding of I 4.0 terminology That's why the paper utilizes many information sources including internet accessible papers and vocabularies and use them as modules in the engineering construction of the paper. Authors believe, that such a papers are in time of starting phases of Industry 4.0 principles, technologies, procedures implementation the most useful for the consequent development and implementation I4.0 principles into industrial praxis.

ACKNOWLEDGEMENT

The research was carried out under support of Technology Agency of the Czech Republic (TF04000074). The authors also gratefully acknowledge financial support from the Ministry of Education, Youth and Sports under projects No. CZ.02.2.69/0.0/0.0/16_027/0008371 and LO1210 - "Energy for Sustainable Development (EN-PUR)" solved in the Centre for Research and Utilization of Renewable Energy).

REFERENCES

1. Manzei, Ch., Schlepner, L., Heinze R. (2016). *Industrie 4.0 im internationalen Kontext*. VDE Verl. GmbH, Beuth Verlag, Berlin.
2. Zezulka, F., Marcon, P., Vesely, I., Sajdl, O. (2016). Industry 4.0 – An Introduction in the phenomenon. IFAC-PapersOnLine, 49(25), pp. 8-12.
3. Internet of Things. In: *Wikipedia: the free encyclopedia* [online]. San Francisco (CA): Wikimedia Foundation, 2001[cit. 2018-02-22]. http://en.wikipedia.org/wiki/Internet_of_things
4. Digitalization. In: *Wikipedia: the free encyclopedia* [online]. San Francisco (CA): Wikimedia Foundation, 2001[cit. 2018-02-22]. <http://en.wikipedia.org/wiki/Digitization>

5. OPC Unified Architecture. In: Wikipedia: the free encyclopedia [online]. San Francisco (CA): Wikimedia Foundation, 2001[cit. 2018-02 22]. https://en.wikipedia.org/wiki/OPC_Unified_Architecture
6. Time-Sensitive Networking. In: Wikipedia: the free encyclopedia [online]. San Francisco (CA): Wikimedia Foundation, 2001[cit. 2018-02 22]. https://en.wikipedia.org/wiki/Time-Sensitive_Networking
7. Unified Modeling Language. In: Wikipedia: the free encyclopedia [online]. San Francisco (CA): Wikimedia Foundation, 2001[cit. 2018-02 22]. https://en.wikipedia.org/wiki/Unified_Modeling_Language
8. Cernohorsky J. (2013), Řídicí systémy s počítači. Učební text a návody do cvičení. Skriptum, VSB – TU Ostrava, pp. 77 – 107.
9. Cloud computing. In: Wikipedia: the free encyclopedia [online]. San Francisco (CA): Wikimedia Foundation, 2001[cit. 2018-02 22]. https://cs.wikipedia.org/wiki/Cloud_computing
10. Cyber-physical system. In: Wikipedia: the free encyclopedia [online]. San Francisco (CA): Wikimedia Foundation, 2001[cit. 2018-02 22]. https://en.wikipedia.org/wiki/Cyber-physical_system
11. Industrie 4.0 (2016). Open Automation, Special Issue to the Hannover Industriemesse 2016, VDE Verlag, pp. 2 – 4
12. Vertical integration. In: Wikipedia: the free encyclopedia [online]. San Francisco (CA): Wikimedia Foundation, 2001[cit. 2018-02 22]. https://en.wikipedia.org/wiki/Vertical_integration
13. Horizontal integration. In: Wikipedia: the free encyclopedia [online]. San Francisco (CA): Wikimedia Foundation, 2001[cit. 2018-02 22]. https://en.wikipedia.org/wiki/Horizontal_integration
14. Industrie 4.0 – Technical Assets (Basic terminology concepts, life cycles and administration models, Status Report, ZVEI, VDE/VDI, 2016

Contacts:

Department of Control and Instrumentation
 Faculty of Electrical Engineering and
 Communication
 Brno University of Technology
 Technická 12
 61600 Brno
 Czech Republic
 E-mail: marcon@feec.vutbr.cz
 E-mail: zezulka@feec.vutbr.cz
 E-mail: bradac@feec.vutbr.cz

TECHNICAL AND THEORETICAL BASIS OF INDUSTRY 4.0 IMPLEMENTATION

ZDENEK BRADAC, FRANTISEK ZEZULKA, PETR MARCON

Abstract: *This paper uses specified terminology to explain theoretical bases of the Industry 4.0 (I4.0). Consequently, paper deals with both the RAMI and the I4.0 component models, which create a theoretical basis of I4.0 principles and their implementation in case studies and next in their implementation into the industrial praxis. Paper uses the German way to develop and implement I4.0 principles into different case studies. Paper's topics give stress on communication systems of the I4.0 and specifies IoT for I4.0 purposes. It deals also with the most recent communication system for purposes of all control levels and makes attention to Time Sensitive Networks. The last part of the paper is an introduction of the creation of the "electronic rucksack", hence the Asset Administration Shell (AAS). Authors introduces readers in this crucial non – simple topic which enables virtualization and modeling of the all production chain.*

Keywords: *Asset Administration Shell, IIoT, Industry 4.0, RAMI, TSN*

1. Introduction

The Industry 4.0 (I4.0) begins and ends with very huge communication activity among components of the production. It is enabled with already existing communication and digitization of information. However, the I4.0 will need still greater information flow than it needs existing state of the art of the industrial production. Next, digitalization and communication have to be more systemic, more rapid, more ordered, more cooperative, more economic. The technical background for it is in recent development of already existing digitalization of information from controlled process, cyber physical-systems and information from the all-technical – economical – marketing chain. Have a look in up to date systems, which will enable this development, hence the TSN, IIoT, RAMI model, I4.0 component model and particularly into the crucial I4.0 item, the Asset Administration Shell (AAS).

2. Communication systems for purposes of Industry 4.0

This chapter deals with the Industrial Internet of Things and Time Sensitive Networking.

2.1. The Industrial Internet of Things (IIoT)

The Internet of things (IoT) is the network of physical devices, vehicles, home appliances and other items embedded with electronics, software, sensors, actuators, and connectivity, which enables these objects to connect and exchange data. Each

thing is uniquely identifiable through its embedded computing system but is able to inter-operate within the existing Internet infrastructure. [1-4].

Experts estimate that the IoT will consist of about 30 billion objects by 2020 [5]. It is also estimated that the global market value of IoT will reach \$7.1 trillion by 2020 [6].

The IoT allows objects to be sensed or controlled remotely across existing network infrastructure, creating opportunities for more direct integration of the physical world into computer-based systems, and resulting in improved efficiency, accuracy and economic benefit in addition to reduced human intervention.

When IoT is augmented with sensors and actuators, the technology becomes an instance of the more general class of cyber-physical systems, which also encompasses technologies such as smart grids, virtual power plants, smart homes, intelligent transportation and smart cities.

"Things," in the IoT sense, can refer to a wide variety of devices such as heart monitoring implants, biochip transponders on farm animals, cameras streaming live feeds of wild animals in coastal waters, automobiles with built-in sensors, DNA analysis devices for environmental, food, pathogen monitoring or field operation devices that assist firefighters in search and rescue operations [7-9]. Legal scholars suggest regarding "things" as an "inextricable mixture of hardware, software, data and service" [10].

These devices collect useful data with the help of various existing technologies and then autonomously flow the data between other devices.

It is useful to separate IoT into two systems. The CIoT (Commercial IoT) and the IIoT (Industrial IoT). The two systems differ in performance and in applications. While the CIoTs give less stress to the hard – real time communication, the IIoT enables communication in near real time parameters. Next the CIoT is performed to be a standard commercial homogenous Ethernet based networks, the IIoT is based on heterogeneous industrial networks based on industrial Ethernet standards. Therefore, in the IIoT networks has to be solved gateways among different communication protocols. It deals with strongly oriented issues – application of internet technologies and networks in industry and for purposes of information exchange among components of industrial production.

The commercial CIoT are intended for more commercial issues and activities such as Smart Building, Smart home, Infotainment Systems, Connected Cars, Smart TV and utilize for connection purposes cloud and big data as well as homogenous TCP/IP networks. The IIoT on the other hand are used for communication in Smart Grids, Smart Cities, Smart Factories, in the all activities of the I4.0 and uses heterogeneous Industrial Ethernet as well as the industrial Fieldbuses and lower industrial networks and protocols.

2.2. Time-Sensitive Networking

The Time-Sensitive Networking (TSN) is a set of standards under development by the Time-Sensitive Networking task group of the IEEE 802.1 working group [1]. The TSN task group was formed at November 2012 by renaming the existing Audio / Video Bridging Task Group (see [11]) and continuing its work. The name changed because of extension of the working area of the standardization group. The standards define mechanisms for the time-sensitive transmission of data over Ethernet networks.

The majority of projects define extensions to the IEEE 802.1Q – Virtual LANs [3]. These extensions in particular address the transmission of very low transmission latency and high availability. Possible applications include converged networks with real time Audio/Video Streaming and real-time control streams, which are used in automotive or industrial control facilities.

Work is also currently being carried out in AVnu Alliance's specially created Industrial group to define Compliance & Interoperability requirements for TSN networked elements [11].

Time sensitive networks are to be general communication tools for communication in the I4.0 environment. They have to fulfill real time requirements on the larger process area then do that industrial Ethernet standards (IE) such as Profinet, PowerLink, Ethernet/IP, EtherCAT and other IEC 61588 standards for real time communication among control systems, operator level, sensors and actuators in the industrial automation systems. The TSN are under development, but the success of the I4.0 implementation is dependent on their standardization. A close cooperation of IEC 61588 standards and development of the standardization process of TSNs is expected. The reason of the TSN topic stems from importance of real – time topic in the Industry 4.0 production, which differs from the existing industrial communication networks in the huge amount of links, entities, data, conditions, distances, heterogeneity of components and business models in smart factories of the future.

3. Models of Industry 4.0 principles, procedures, technologies

The following subchapters describe two I4.0 models, namely RAMI 4.0 and I4.0 component model.

3.1. Industry 4.0: RAMI 4.0

The Authors of the RAMI 4.0 (Reference Architecture Model Industry 4.0) model are BITCOM, VDMA and ZVEI. They decided to develop a 3D model because the model should represent all different manually interconnected features of the technical – economical properties. The model SGAM, which was developed for purposes of communication in networks of renewable energy sources, seemed to be as an appropriate model for the Industry 4.0 applications as well. The RAMI 4.0 is a small modification of the SGAM (Smart Grid Architecture Model).

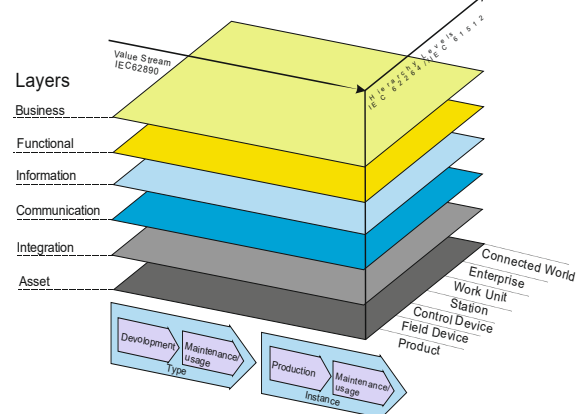


Fig. 1. RAMI 4.0 model

Because into the SGAM as well as into the RAMI 4.0 enter approximately 15 industrial

branches, the RAMI 4.0 model enables looks from different aspects. That's why layers in the vertical axis represent the look from different aspects (a look from the market aspect, a look from a perspective of functions, information, communication, a look from an integration ability of the components) [1], [12].

Very important criterion in the modern engineering is the product life cycle with the value stream, which it contains. The left – hand horizontal axis displays this feature. There are expressed e.g. constant data acquisition throughout the life cycle. Even the totally digitization of the whole development – market chain offers great potential for improvement of products, machines, and other layers of the I4.0 architecture throw-out the all life cycle. This look corresponds well with the IEC 62890 draft standard.

The next model axis (right in the horizontal level) describes function position of the components in the I4.0. In this axis, there is specified the functionality of the components, no any specification for implementation but the function assignment only. The axis respects both IEC 6224 and the 61512 standards. However, the IEC 6224 and the 61512 standards are intended for specification of components in a position in one enterprise or works unit only. Therefore, the highest level in the axis horizontal right is the connected world.

3.2. Industry 4.0: Component Model

The second very important model for purposes of the I4.0 that has been developed by BITCOM, VDMA and ZVEI during the last year is the I4.0 components model (see Fig. 2). It is intended to help producers and system integrators to create HW and SW components for the I4.0. It is the first and the only (in July 2016) specific model which goes out from the RAMI 4.0 model.

It enables better description of cyber – physical features and enables description of communication among virtual and cyber – physical objects and processes. The HW and SW components of future production will be able to fulfil requested tasks by means of implemented features specified in the I4.0 components model.

The most important feature is the communication ability among the virtual objects and processes with real object and processes of production while this model specifies the conform communication. Physical realization of it is that any component of the I4.0 system takes an electronic container (shell) of secured data during the all life cycle. The data are available to all entities of the technical – production chain.

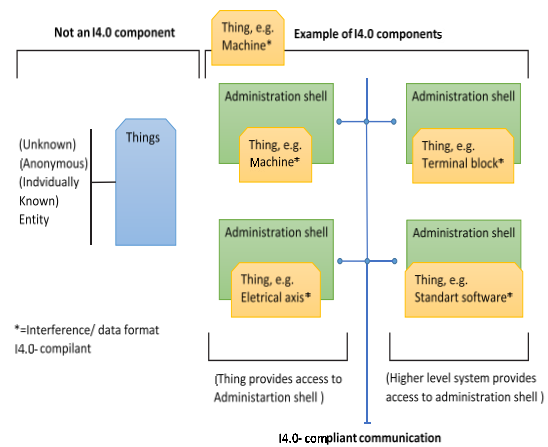


Fig. 2. Industry 4.0 components model [14]

Therefore this model goes out from a standardized, secure and safety real time communication of all components of production. The electronic container (shell) of data and the all Industry 4.0 component model is specified in the Fig. 3 [14]. The main important part of I4.0 components is the AAS, see the Fig. 3.

4. Asset Administration Shell

The most important development has been done in specification of the Asset Administration Shell (AAS). The AAS is the crucial item in the all Industry 4.0 idea. It creates an interface between the physical and virtual production steps. AAS is a virtual digital and active representation of an I4.0 component in the I4.0 system [15,16].

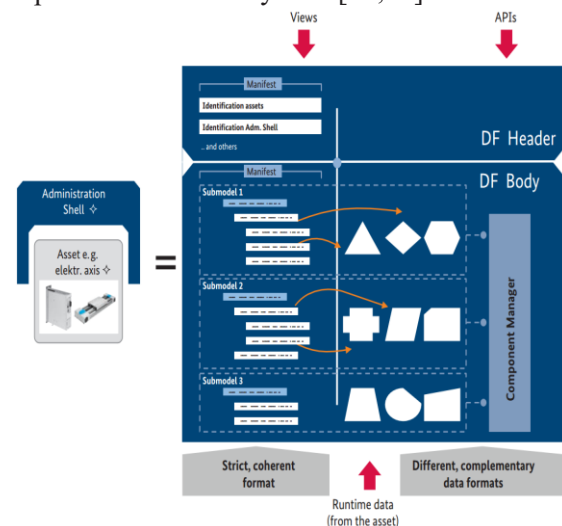


Fig. 3. Asset Administration Shell [6]

Any component of production in the I4.0 environment has to have an administrative shell. The structure of the AAS has to fulfill requirements of different aspects of production and has to enable functionality of I4.0 components from all basic perspective: market, construction, power, function,

positioning, security, communication ability, understandability.

The AAS is composed from a body and a header. The header contains identifying details regarding the asset administration shell and the represented asset. The body contains a certain number of submodels for an asset-specific characterization of the asset administration shell [6].

As it can be seen in the Fig. 3, the AAS is made of a series of submodels. These represent different aspects of the asset concerned. For example, they may contain a description relating to safety or security, but could also outline various process capabilities such as drilling or installation. Possible submodels of an AAS are pictured in the Fig. 4.

Administration Shell IEC TR 62794 & IEC 62832 Digital factory	
Submodels	Standards
Identification	ISO 29005 or URI unique ID
Communication	IEC 61784 Fieldbus profiles
Engineering	IEC 61360/ISO13584 Standard data elem. IEC 61987 Data structures and elements Ecl@ss Database with product classes
Configuration	IEC 61804 EDDL IEC 62453 FDT
Safety (SIL)	EN ISO 13849 EN/IEC 61508 Functional safety discrete EN/IEC 61511 Functional safety process EN/IEC 62061 Safety of machinery
Security	IEC 62443 Network and system security
Lifecycle status	IEC 62890 Lifecycle
Energy Efficiency	ISO/IEC 20140-5
Condition monitoring	VDMA 24582 Condition monitoring
Examples of AAS using	Drilling, Milling, Deep drawing, Clamping, Welding, Painting, Mounting, Inspecting, Validating ...

Fig. 4. Possible submodels of an asset administration shell [16]

The aim is that to standardize only one submodel for each aspect. Thus it will be possible to search for e.g. a welding machine with searching for an AAS containing “welding” with appropriate properties. Second submodel in the example e.g. “energy efficiency” could ensure that the welding

stand can be able to save electricity when it is not in operation mode.

Each submodel contains a structured quantity of properties that can refer to data and functions. A standardized format based on IEC 61360 is required for the properties. Data and functions may be available in various, complementary formats.

The properties of all the submodels therefore result in a constantly readable directory of the key information or, as it were, the Manifest of the asset administration shell and thus of the I4.0 components. To enable binding semantics, asset administration shells, assets, submodels and properties must be clearly identified. Permitted global identifiers are ISO 29002 – 5 (e.g. eCl@ss and IEC Common Data Directories) and URIs (Unique Resource Identifiers, e.g. for ontologies).

Next Fig. 5 [16] shows how an interaction pattern is directed towards the domain specific submodels in the asset administration shell. It is shown on a possible example from a discrete manufacturing process.

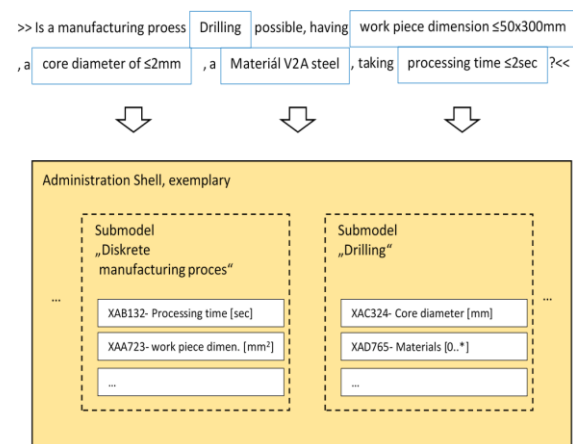


Fig. 5. Interaction pattern directed towards the domain – specific submodels in the AAS [16]

According to the language for Industry 4.0 the Fig. 6 shows an approach to this idea from the sub-working standardization group [16]:

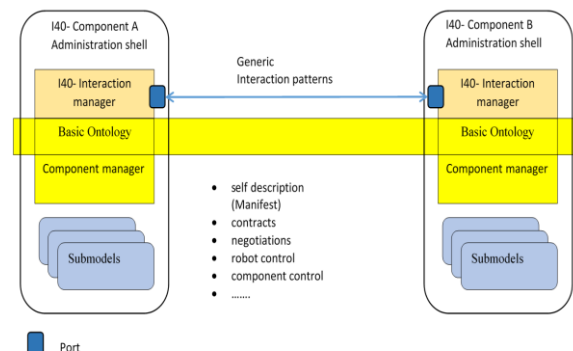


Fig. 6. An approach to the topic “languages of I4.0”

For such purposes is by any component of I 4.0 one interaction manager, which is responsible of processing of interaction patterns in the network. A domain – independent basic ontology safeguards the connection with the domain – specific submodels in the AAS.

5. Conclusion

Paper goes out from the contribution Marcon P., Zezulka F., Bradac Z.: Terminology of Industry 4.0, Proc. of TechSys 2018, Plovdiv 2018 [17] and uses specified terminology to explain theoretical bases of the I4.0.

Paper's topics give stress on communication systems of the I4.0 and specifies IoT for I 4.0 purposes. It deals also with the most up to date communication system for purposes of all control levels and make attention to Time Sensitive Networks.

Consequently, paper deals with both the RAMI and the I4.0 component models, which create a theoretical basis of I4.0 principles and their implementation in case studies and next in their implementation into the industrial praxis. Paper uses the German way to develop and implement I4.0 principles into different case studies.

ACKNOWLEDGEMENT

The authors gratefully acknowledge financial support from the Ministry of Education, Youth and Sports under projects No. LO1210 - "Energy for Sustainable Development (EN-PUR)" solved in the Centre for Research and Utilization of Renewable Energy). The research was carried out under support of Technology Agency of the Czech Republic (TF04000074).

REFERENCES

- Manzei, Ch., Schlepner, L., Heinze R. (2016). *Industrie 4.0 im internationalen Kontext*. VDE Verl. GmbH, Beuth Verlag, Berlin.
- Zezulka, F., Marcon, P., Vesely, I., Sajdl, O. (2016). Industry 4.0 – An Introduction in the phenomenon. IFAC-PapersOnLine, 49(25), pp. 8-12.
- Marcon, P. et al. (2017), Communication technology for industry 4.0. *Progress In Electromagnetics Research Symposium - Spring (PIERS)*, St. Petersburg, pp. 1694-1697, doi: 10.1109/PIERS.2017.8262021
- Internet of Things. In: *Wikipedia: the free encyclopedia* [online]. San Francisco (CA): Wikimedia Foundation, 2001[cit. 2018-02-22]. http://en.wikipedia.org/wiki/Internet_of_things
- Digitalization. In: *Wikipedia: the free encyclopedia* [online]. San Francisco (CA): Wikimedia Foundation, 2001[cit. 2018-02-22]. <http://en.wikipedia.org/wiki/Digitization>
- Special Issue to the Hannover Industriemesse 2016, VDE Verlag, pp. 2 – 4.
- Cloud computing. In: *Wikipedia: the free encyclopedia* [online]. San Francisco (CA): Wikimedia Foundation, 2001[cit. 2018-02 22]. https://cs.wikipedia.org/wiki/Cloud_computing
- Computer and Automation. Fachmedium der Automatisierungstechnik, 9-2017
- Cyber-physical system. In: *Wikipedia: the free encyclopedia* [online]. San Francisco (CA): Wikimedia Foundation, 2001[cit. 2018-02 22]. https://en.wikipedia.org/wiki/Cyber-physical_system
- Vertical integration. In: *Wikipedia: the free encyclopedia* [online]. San Francisco (CA): Wikimedia Foundation, 2001[cit. 2018-02 22]. https://en.wikipedia.org/wiki/Vertical_integration
- Time-Sensitive Networking. In: *Wikipedia: the free encyclopedia* [online]. San Francisco (CA): Wikimedia Foundation, 2001[cit. 2018-02 22]. https://en.wikipedia.org/wiki/Time-Sensitive_Networking
- Industrie 4.0 (2016). Open Automation, Special Issue to the Hannover Industriemesse 2016, VDE Verlag, pp. 2 – 4.
- VDI/VDE. (March 2016). Gesellschaft Mess- und Automatisierungstechnik. Status report Industrie 4.0 – Technical Assets. Basic terminology concepts.
- VDI/VDE. (2015). Gesellschaft Mess- und Automatisierungstechnik. Status report. Reference Architecture Model Industrie 4.0 (RAMI 4.0).
- ZVEI, VDI/VDE. (March 2016). Industrie 4.0. Technical Assets, Status Report.
- ZVEI (April 2017). Examples of Assets Administration Shell for Industrie 4.0 Components- Basic Part.
- Walter K.D. (2016). Die „Thinks“ im Nebel, Computer and Automation, vol. 8, pp. 12-14.
- Marcon, P., Zezulka, F. and Bradac Z. (2018), Terminology of Industry 4.0. *In Proc. Of TechSys 2018*, Plovdiv.

Contacts:

Department of Control and Instrumentation
Faculty of Electrical Engineering and
Communication, Brno University of
Technolgy
Technicka 12, 61600 Brno, Czech Republic
E-mail: zezulka@feec.vutbr.cz
E-mail: marcon@feec.vutbr.cz
E-mail: bradac@feec.vutbr.cz

RESOURCE ALLOCATION BY PORTFOLIO OPTIMIZATION

KRASIMIRA STOILOVA, TODOR STOILOV, MIROSLAV VLADIMIROV

Abstract: An optimal allocation of financial resources based on the Portfolio theory is worked out. Standard portfolio optimization problem is defined and solved on the base of statistical data of four financial indices, available at the Bulgarian stock exchange. The optimization problem of linear-quadratic programming is solved for different values of the coefficient, formalizing the investor's preferences for having risk during the investment process. Assessments and comparisons of the investment solutions are presented. Graphical illustrations of the problems' solutions are given.

Key words: *Portfolio theory, investments, optimization*

1. Introduction

The allocation of investments for a set of financial assets (securities, bonds), bought or sold from the stock exchange is a practical task, which is a current activity for many investors. The investor's goal is to invest today capital in financial assets in order to obtain later return by selling their assets. Each asset for the investor has a potential for future income. The set of securities is called "portfolio". The investor wants to know which is the best combination of the securities in the portfolio in order to receive better return. It means that there are different opportunities to allocate the investments for a set of financial assets, bought from the stock exchange. Which asset to be chosen and what amount of them to be bought concerns the decision making of the investor. A formal model for supporting such decision making is the portfolio theory [1]. The formalization of the portfolio theory results in definition and solution of a portfolio optimization problem. It gives as solution the optimal allocation of financial resources for trading financial assets. For the investment process the target is to maximize the return while the investment risk has to be minimal [2,3,4,5]. The problem of portfolio optimization targets the optimal resource allocation in investment process [5,7,8].

2. Portfolio Optimization Problem

The analytical relations between the portfolio risk V_p , portfolio return E_p and the

values of the investment per type of assets x_i according to the portfolio theory are [5]

$$E_p = \sum_{i=1}^n E_i x_i = E^T x \quad (1)$$

$$V_p = \sum_j \sum_{i=1}^n x_i x_j \text{cov}(i, j) = x^T \text{cov}(\cdot) x \quad (2)$$

where

E_i - the average value of the return of asset i ;

$E^T = (E_1, \dots, E_n)^T$ - a vector with dimension $1 \times n$;

$\text{cov}(i, j)$ - the covariation coefficient between the assets i and j . The covariation is calculated from available historical statistical data for the returns of assets i and j . The matrix $\text{cov}(\cdot)$ by definition is a symmetric one, or $\text{cov}(i, j) = \text{cov}(j, i)$.

Relation (1) is the quantitative evaluation of the portfolio return. Relation (2) formalizes the quantitative assessment of the portfolio risk. The portfolio problem solutions $x_i, i=1, n$ determine the relative amounts of the investment per security i . The covariation is determined from previously available statistical data of the returns of assets i and j and it represents a symmetrical matrix:

$$\text{cov}(\cdot) = \begin{pmatrix} \text{cov}(1,1) & \text{cov}(1,2) & \dots & \text{cov}(1,n) \\ \text{cov}(2,1) & \text{cov}(2,2) & \dots & \text{cov}(2,n) \\ \vdots & & & \\ \text{cov}(n,1) & \text{cov}(n,2) & \dots & \text{cov}(n,n) \end{pmatrix}_{n \times n}$$

The components $\text{cov}(i,j)$ are evaluated from the profits of assets i and j $R_i^{(1)}, R_i^2, \dots, R_i^{(N)}$ and $R_j^{(1)}, R_j^2, \dots, R_j^{(N)}$ for discrete time moments (1), (2), ..., (N). The covariation's component between assets i and j is calculated as

$$\text{cov}(i,j) = \frac{1}{N} \left[(R_i^{(1)} - E_i)(R_j^{(1)} - E_j) + (R_i^{(2)} - E_i)(R_j^{(2)} - E_j) + \dots + (R_i^{(N)} - E_i)(R_j^{(N)} - E_j) \right]$$

where

$$E_i = \frac{1}{N} [R_i^{(1)} + R_i^{(2)} + \dots + R_i^{(N)}]$$

$$E_j = \frac{1}{N} [R_j^{(1)} + R_j^{(2)} + \dots + R_j^{(N)}]$$

are the average profits of the assets i and j for the period $T = [1, 2, \dots, N]$. Particularly, the value

$\text{cov}(i,i) = \tau_i^2$ gives the variation of the return of asset i . The portfolio theory defines the so called standard optimization problem as [2]

$$\begin{aligned} \min_x & \left[\frac{1}{2} x^T \text{cov}(\cdot)x - \sigma E^T x \right] \\ & x^T \cdot \mathbf{1} = 1, \end{aligned} \quad (3)$$

where $\text{cov}(\cdot)$ – is a symmetrical positively defined square matrix $n \times n$,

E – is a $(n \times 1)$ vector of the average profits of the assets for the period of time $T = [1, 2, \dots, N]$;

σ – is a parameter of the investor's preferences to undertake risk in the investment process.

$$\mathbf{1} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}, \text{ is a unity vector, } n \times 1;$$

The constraint of the optimization problem presents the equation $x_1 + x_2 + \dots + x_n = 1$, which is a formalization of the assumption that the full amount of the resources are devoted for the investments. If the right side of the constraint is less than 1, this means that all amount of the investments are not effectively used. The investment per different assets has to be performed for the total amount of the available investment resources, numerically presented as relative value of 1. The solutions $x_i, i=1, n$ give the relative values of the investment, which are allocated for the assets $i, i=1, n$.

The component of the goal function $V_p = x^T \text{cov}(\cdot)x$ is the quantitative assessment of the portfolio risk. The component $E_p = E^T x$ is the quantitative value of the portfolio return. The goal function of problem (3) targets the minimization of the portfolio risk V_p as well as the maximization of its return E_p . The parameter σ formalizes the investor's ability to undertake risk and it has numerical value in the range $[0, +\infty]$. When $\sigma=0$ the investor is very cautious (even coward) and his main goal is to decrease the risk of the investment, $\min_x [x^T \text{cov}(\cdot)x]$. For the case $\sigma=+\infty$ the investor has forgotten the existence of risk in the investments. His target is to obtain a maximal return from the investment. For that case the relative weight of the return in the goal function is most weighted, and then the optimization problem has an analytical form:

$$\min_x [-\sigma E^T x] \equiv \max_x [E^T x] .$$

Thus, in the portfolio problem it is introduced a new unknown parameter σ , which assesses the investor's preferences for undertaking risk in decision making. This parameter influences the portfolio problem, making it a parametric one. Respectively, for a new value of σ , the portfolio problem (3) has to be solved again. The trivial case when σ is not properly estimated the optimization problem has to be solved for a set of σ . For practical reasons, the portfolio problem has to be multiple solved with a set of values for the coefficient of the investor's preferences σ to undertake risk.

The numerical assessment of σ parameter is a subjective task of the financial analyzer. This coefficient strongly influences the definition and respectively the solutions of the portfolio problem. Respectively, σ changes also the final investment decision.

The portfolio theory uses the relation risk-return $V_p = V_p(E_p)$ for the assessment of the portfolio characteristics, which result from the combinations of the assets, used in the portfolio. The investors have to choose optimal portfolios from this relation $V_p = V_p(E_p)$, which is titled "efficiency frontier". This "efficiency frontier" is not evidently found. Points from this curve can be found by solving the portfolio optimization problem with different values of the parameter σ . The "efficiency frontier" is evaluated point after point according to the iterative numerical procedure:

1. An initial value of σ for the investor's preferences is chosen, for instance $\sigma=0$. This corresponds to investor who is far from risky decisions.

2. The portfolio problem is solved with the stated σ

$$\min_x \left[\frac{1}{2} x^T \text{cov}(\cdot) x - \sigma E^T x \right]$$

$$x^T \times 1 = 1$$

and the optimal solution $x(\sigma)$ is found.

3. Evaluation of the portfolio risk and portfolio return according to (1) – (2):

$$V_p = x^T(\sigma) \text{cov}(\cdot) x^T(\sigma), \quad E_p = E^T x(\sigma).$$

These values give one point of the relation $V_p = V_p(E_p)$, which belongs to the efficiency frontier.

4. New value $\sigma_{\text{new}} = \sigma_{\text{old}} + \Delta$ is chosen, where Δ is determined by considerations for completeness by considering the set of $\sigma = [0, +\infty]$. Then, go to point 2.

For each solution of the portfolio optimization problem one point of the space $V_p = V_p(E_p)$, belonging to the curve of the efficiency frontier is found, Fig.1.

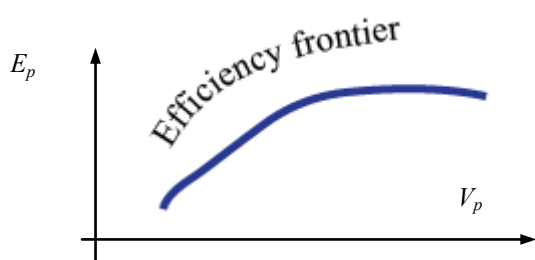


Fig.1. Efficiency frontier of the portfolio optimization

For practical cases of individual investor, problem (1) is solved with a set of values of σ according to the expert's experience [9]. Having a set of solutions $x(\sigma)$ the final value of σ^* for that investor is empirically estimated, which gives also the final optimal portfolio solution $x(\sigma^*)$ as well.

3. Assessment of Bulgarian financial indices according to the Portfolio theory

This research illustrate the application of the portfolio theory for the Bulgarian stock exchange. Currently, this market is assessed by evaluation of four indices: SOFIX, BGBX40, BGTR30 and BGREIT. This research assumes that each index represents a particular security. Applying the portfolio theory by definition and solution of a portfolio optimization problem with these securities will give information about which part of the stock exchange is preferable for investment.

For the definition of the portfolio problem it is necessary to evaluate the correlations, respectively the covariation matrix for these four indices. For the evaluations of the covariation matrix it has been used public data of the Bulgarian stock exchange, available from [6]. Each of these four indices, SOFIX, BGBX40, BGTR30 and BGREIT, evaluated for the Bulgarian stock exchange is calculated according to predefined rules of usage of characteristics of established set of securities. This research analyzes the behavior of these indices and as a result it recommends to the investors which part of the stock exchange, assessed by the corresponding index is preferable for investment receiving better return. The statistical data about the behavior of the four indices for a period of one year are given in Table 1 [6]

Table 1 Statistical data of 4 Bulgarian indices

	SOFIX	BGBX40	BGTR30	BGREIT
30.1.2018	712,73	138,23	571,59	115,41
29.12.2017	677,45	132	555,98	116,1
30.11.2017	665,03	130,49	547,89	113,99
31.10.2017	671,41	131,19	547,08	115,88
29.9.2017	688,11	134,34	559,26	114,88
31.8.2017	705,44	134,85	553,9	115,33
31.7.2017	715,21	135,52	548,7	115,01
30.6.2017	703,46	134,22	535,47	113,78
31.5.2017	661,23	130,61	516,72	111,12
28.4.2017	657,29	130,25	519,92	108,39
30.3.2017	633,04	124,9	502,24	108,78
28.2.2017	611,12	120,56	486,06	107,83
28.1.2017	610,1	117,72	471,76	107,61

The covariant matrix according to the financial data is in the form

	SOFIX	BGBX40	BGTR30	BGREIT
SOFIX	cov(1,1)	cov(1,2)	cov(1,3)	cov(1,4)
BGBX40	cov(2,1)	cov(2,2)	cov(2,3)	cov(2,4)
BGTR30	cov(3,1)	cov(3,2)	cov(3,3)	cov(3,4)
BGREIT	cov(4,1)	cov(4,2)	cov(4,3)	cov(4,4)

All the components of the covariant matrix using the data of the Bulgarian indices consisting the portfolio are calculated. The covariant matrix according to Table 1 is

Cov() =

	SOFIX	BGBX40	BGTR30	BGREIT
SOFIX	1178,053	190,601	896,482	93,672
BGBX40	190,601	32,965	157,193	15,373
BGTR30	896,482	157,193	856,949	87,683
BGREIT	93,672	15,373	87,683	10,372

Using the data from Table 1, the average returns of the four indices are:

$$E^T = [670,12 \quad 130,38 \quad 532,04 \quad 112,62].$$

Having the parameters for Cov() and E[.], the optimization problem (3) is defined up to the value of the coefficient σ . The optimization problem (3) is solved several times by using different values of the coefficient of the investor's preferences to undertake risk σ . The steps of the above sequence of calculations gives as solutions the allocation of the investment per securities $x(\sigma^*)$. Additionally, using relations (1) and (2) the Return(σ^*) and Risk (σ^*) are evaluated for different values of σ . It means that for each σ the portfolio optimization problem (3) has been solved. The calculation environment for the problem's solution is the popular software application of Excel – Solver.

For the different values of σ the portfolio optimization problem gives solutions x_1, x_2, x_3, x_4 , given in Table 2. The values of Risk and Return for each σ are given in the right two columns of Table 2.

The results of the optimization solutions show that for small values of σ (0; 0,1) the investments should be done to BGREIT. When σ varies between 0,2 and 0,8 the investments to BGREIT should decrease and the investments to SOFIX should increase. For $1 < \sigma < 2$ the investments to SOFIX and BGTR30 increase and the investment to BGREIT decrease. For $\sigma > 2$ the optimization results recommend the investments to be done only to SOFIX.

Table 2 Portfolio optimization problem's solutions

σ	SOFIX	BG BX 40	BG TR 30	BG REIT	Risk	Return
	x_1	x_2	x_3	x_4		
0	0	0	0	1	10,37	112,62
0,1	0	0	0	1	10,37	112,62
0,2	0,0282	0	0	0,9718	15,86	128,33
0,25	0,0560	0	0	0,9439	22,84	143,85
0,3	0,0839	0	0	0,9161	31,38	159,38
0,35	0,1117	0	0	0,8883	41,47	174,90
0,4	0,1396	0	0	0,8604	53,12	190,42
0,45	0,1952	0	0	0,8047	81,06	221,47
0,5	0,1952	0	0	0,8047	81,06	221,47
0,6	0,2509	0	0	0,7490	115,21	252,52
0,8	0,3623	0	0	0,6376	202,14	314,61
1	0,4737	0	0	0,5263	313,91	376,71
1,5	0,7253	0	0,037	0,2377	703,48	532,51
2	0,9776	0	0,022	0	1165,59	667,04
2,5	1	0	0	0	1178,05	670,12
3	1	0	0	0	1178,05	670,12
3,5	1	0	0	0	1178,05	670,12

The variations of the Risk and Return according to σ are presented in Figures 2 and 3. The variation of Risk to Return is presented in Fig.4. This curve represents the efficiency frontier. The investor chooses a point of this curve to allocate the investment resources in optimal manner. For example, the investor can choose Risk 100 and the corresponding value of Return is about 240. It means that the investments should be allocated to two financial funds: 75% to BGREIT and about 25% to SOFIX, Fig.5. If the investor prefers better Return, greater than 670, it means that the resources must be allocated only to one fund - SOFIX.

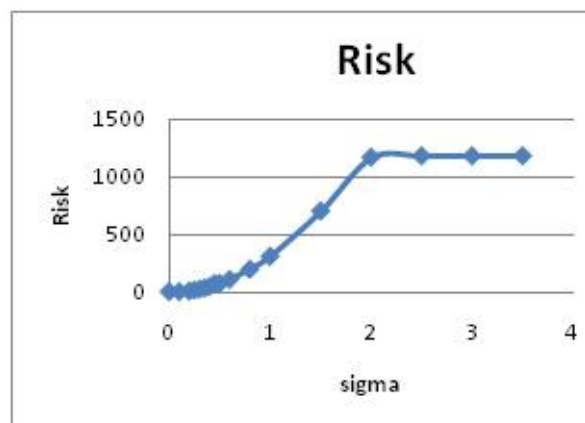


Fig.2. The portfolio Risk according to σ

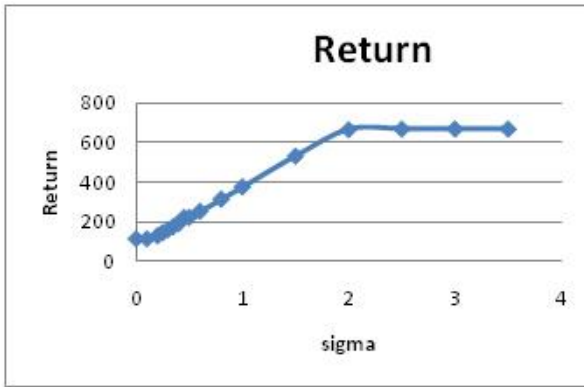


Fig.3. The portfolio Return according to σ

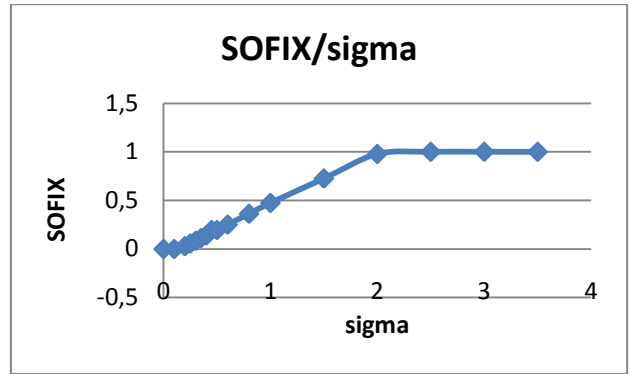


Fig.6 Optimal solutions changes SOFIX /sigma

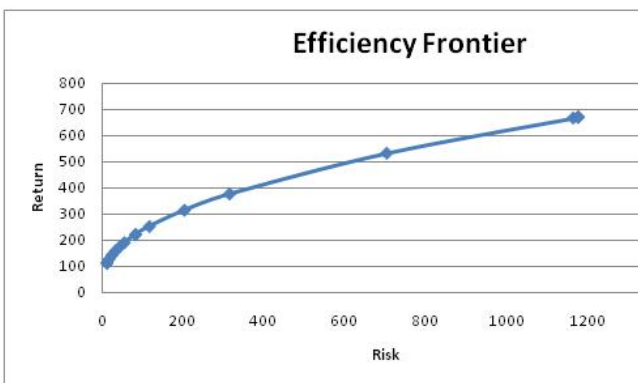


Fig.4 The Efficiency frontier (Return/Risk)

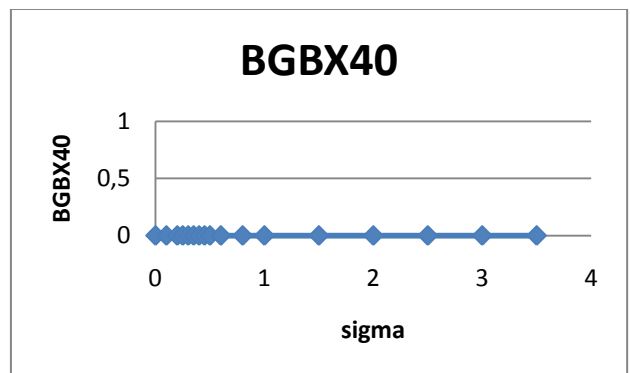


Fig.7 Optimal solutions changes BGBX40 /sigma

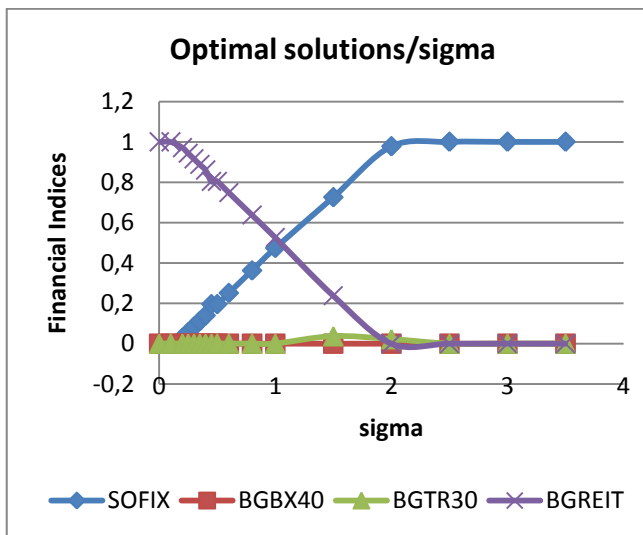


Fig.5 Optimal solutions changes/sigma

The variation of each of the fourth indices are presented in Fig.6 – Fig.9.

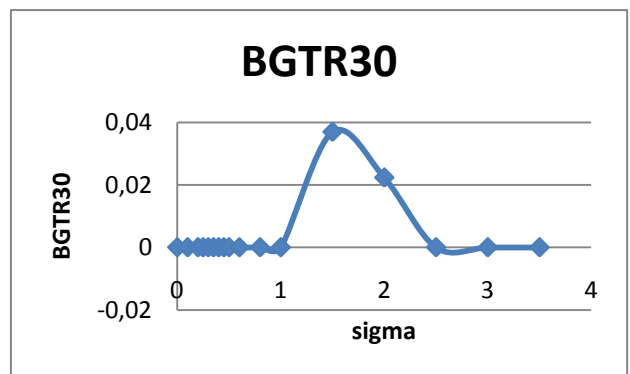


Fig.8 Optimal solutions changes BGTR30 /sigma

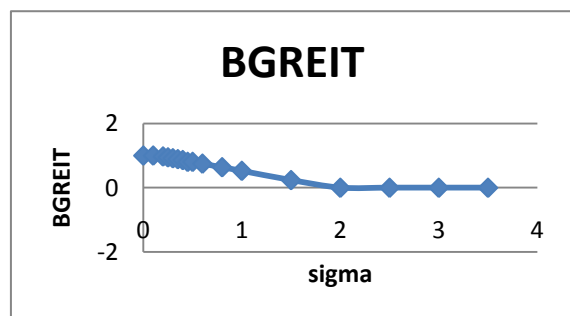


Fig.9 Optimal solutions changes BGREIT /sigma

4. Conclusions

The Portfolio theory is a strong theoretical support for financial resource allocation. The mathematical formalization of the financial process is on the base of solving linear-quadratic optimization problem so called standard portfolio optimization problem. The peculiarity is that this problem depends on a coefficient σ which represents the investor's relation to take a risk. That is why the portfolio optimization problem and its corresponding solutions depend on the value of σ . The standard portfolio optimization problem for different values of σ is solved. The most popular Bulgarian financial available indices are used for the assessment of the financial investments. The main characteristics of the portfolio optimization problem – Return and Risk, which determine the main portfolio characteristic - efficiency frontier are obtained in results of problem's solutions. Analysis of solutions variations and assessment of their behavior is proposed. This research provides an actual analysis for the recent behavior of the Bulgarian stock exchange. The results obtained provide advices to the investors which part of the Bulgarian stock exchange to be used for investment. The added value of this research concerns the definition and solution of a portfolio problem, which addresses the dynamical behavior of the Bulgarian stock exchange. As a result this research provides a comparative analysis for real assets which are under trade on the Bulgarian stock exchange.

REFERENCES

1. Bodie, Z., Kane, A., Marcus, A. (2000). *Investments*. Naturela, Sofia, 906 p.
2. Campbell J., Chacko G., Rodriguez, J., Viceira, L. (2002) Strategic asset allocation in a continuous-time VAR model, Harvard Institute of Economic Research, Harvard University Cambridge, Massachusetts, 1-21.
3. Kohlmann, M., Tang, S. (2003) Minimization of risk and linear quadratic optimal control

theory. *SIAM J. Control optim*, Vol. 42, No. 3, 1118–1142

4. Magiera P., Karbowski, A. (2001) Dynamic portfolio optimization with expected value-variance criteria. Preprints of the 9th IFAC Symposium on LSS'01, Bucharest, Romania, 308-313.
5. Sharpe W. (2000). Portfolio theory & Capital markets, Mc Grow Hill, No4.
6. <http://infostock.bg>
7. Sharpe, W. Adaptive asset allocation policies.- J. Financial Analysts. Vol. 66 issue 5, 2010, pp.45-49.
8. Schulmerich, M., Leporcher, Y.M., EU C.H. Applied asset and risk management. A guide to Modern Portfolio Management and behavior-driven markets. 2015, XVII, 476 p, ISBN 978-3-642-55443-8.
9. Vatchova B. "Logical method for knowledge discovery based on real data sets" Proceedings of the IADIS European Conference Data Mining 2011, Rome, Italy, ISBN: 978-972-8939-53-3 © 2011 IADIS, pp.203-207.

Acknowledgement

This work has been partly supported by project H12/8, 14.07.2017 of the Bulgarian National Science fund: Integrated bi-level optimization in information service for portfolio optimization, contract ДН12/10, 20.12.2017

Contacts

Krasimira Stoilova

Todor Stoilov

Institute of Information and Communication Technologies – Bulgarian Academy of Sciences

Acad. G.Bonchev str., bl.2, Sofia 1113

Tel. +359 2 979 27 74

E-mail: todor@hsi.iccs.bas.bg

k.stoilova@hsi.iccs.bas.bg

Miroslav Vladimirov

University of Economics - Varna

77, Boris I bul., Varna 9002

E-mail: vladimirov@ue-varna.bg

PLANNING AND IMPLEMENTATION OF THE ERP SYSTEM IN PACKAGING PRODUCTION. PRACTICAL ASPECTS

RADOSLAV HRISCHEV

Abstract: *The article presents practical aspects - necessary conditions, model of planning and implementation of a specialized ERP system for packaging production at a plant, part of multinational company. Are explored main problems in implementing of such systems.*

Key words: *ERP, packaging production, Kiwiplan*

1. Introduction

Evolutionary automation of production in industrial plants goes through the building of a sustainable IT infrastructure, the development and implementation of simple specialized IT systems / financial, human resources management, warehouse, manufacturing, etc./ to come to the need of ERP /Enterprise Resource Planning/ systems.

ERP systems cover all /or almost/ all information flows and provide the needed information to employees and managers to quickly make efficient decisions, make business processes more effective and help to reduce costs and increase revenue in organization. In a complex ERP system, software solutions are becoming more flexible and user-friendly. However, ERP systems are expensive and time-consuming investment that requires serious and professional planning.

In the specialized literature [2], the implementation of the ERP systems is defined as a key element of FoF - Factory of Future building, together with the renovation of the production facilities - technologies, machines and equipment.

2. Implementation of ERP system in a packing plant

In this article is explored an example of implementation of specialized ERP system for the managing of plant for packaging production – boxes from corrugated board, part of a large multinational company.

The basic requirements for implementing of ERP system:

- Comply with group policies - it should be compatible with the existing systems in the main company, preferably to be part of corporate ERP system. This leads to a significant reduction of deployment costs,

based on shared experience and the cost of the needed IT infrastructure - equipment, licenses, administration and maintenance.

- To be integrated with existing corporate systems - financial, BI, CRM and others. As a rule, multinational companies are using unified financial systems - mainly SAP.
- Be flexible and scalable; to consist of separate standalone integrated modules. This allows only those modules that are in line with business processes to be deployed. On the other way, with the expansion of the business can be easily added new modules if it's necessary.
- To have an upgrade and maintenance policy and rules.

One of the most popular packaging plant management ERP systems for cardboard and corrugated packaging production is Kiwiplan - <http://www.kiwiplan.com> [1]. The system is widely distributed around the world and has more than 600 corporate customers with more than 680 covered plants.



Fig. 1. *Kiwiplan in the world*

The system is built on a modular principle and is extensively flexible. All modules are presented on the diagram /Fig. 2/ and are connected with embedded interfaces.

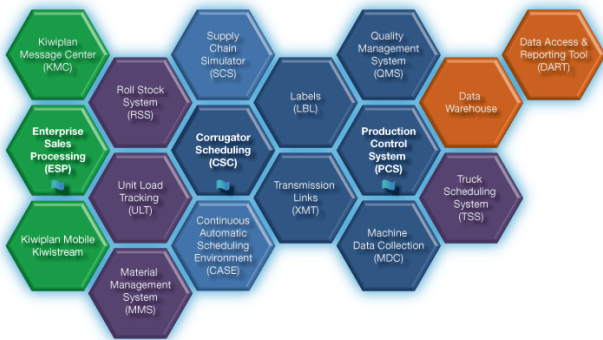


Fig. 2. Structure scheme of Kiwiplan

Main modules of Kiwiplan:

- Enterprise Sales Processing /ESP/ - managing of the sales
- Roll Stock System /RSS/ - roll paper warehouse system
- Corrugator Scheduling /CSC/ - planning of corrugated and boxes production
- Production Control System /PCS/ - managing of the production
- Data Warehouse – finish good store management
- Truck Scheduling System /TSS/ - expedition and loading

The system is bidirectional linked to the control systems of the corrugated aggregates and converting /for the boxes/ production machines. This means that in real-time production systems are exchanging information with the ERP system. Operators of the machines on a special terminals production machines monitored the required order parameters and outputs.

3. Example of implementation of ERP system Kiwiplan in a packaging plant

The specificities of this implementation are existing already implemented modules SAP R3 - FI /Finance/, CO /Controlling/ MM /Material Management/. Additionally, in the plant was in use own process management IS, developed by the local IT specialists of the plant by request of users, accurately reflecting the workflows of the units of the plant. This has resulted in resistance of implementation at various levels, from machine operators up to management. On the other hand, the availability of information in a structured electronic Data Bases facilitated the migration of data at the ERP implementation stage.

3.1. Customization of the system for current structure of the plant

The structure of the implemented system shown on the Figure 3.

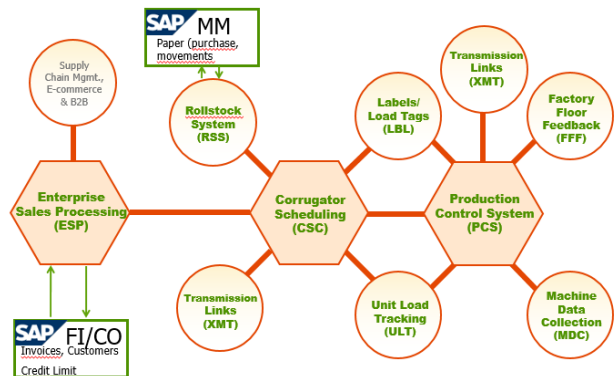


Fig. 3. Structure scheme of concrete implementation of Kiwiplan

In this implementation process, not all the system modules have been used. As an example, the quality control is done through another system because of special characteristics of the process and client needs. This has reduced the expenses/costs of the system implementation. On the figure /Fig.3/ shows also interfaces to existing SAP modules. In fact, this is done by BCS /Business Connect Servers/.

3.2. Logical model of implementation

For the project realization, a model of implementation - ERP Kiwiplan has been created.

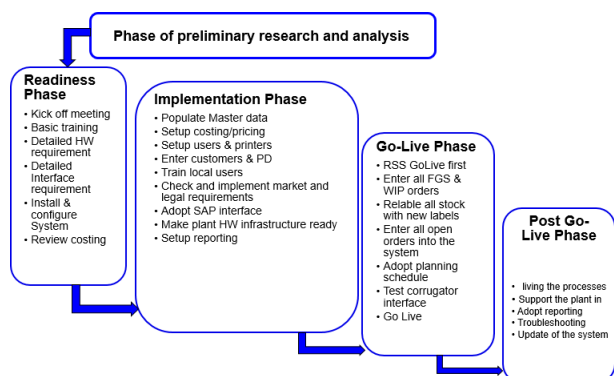


Fig. 4. Model of implementation

It includes the whole process in five basic steps:

- Analysis Phase

- Presentation of current plant status - Local Management
- Analyze standard process flow
- Detailed analysis / process area
- Presentation of analyzed plant status - Process Management
- Discussion of results Managers Board approval meeting

- **Readiness Phase**
 - Kick off meeting
 - Basic training
 - Detailed HW requirement
 - Detailed Interface requirement
 - Install & configure System
 - Review costing
- **Implementation Phase**
 - Populate Master data
 - Setup costing/pricing
 - Setup users & printers
 - Enter customers & Project designs
 - Train local users
 - Check and implement market and legal requirements

- Adopt SAP interface
- Make plant HW infrastructure ready
- Setup reporting
- **Go-Live Phase**
 - RSS GoLive first
 - Enter all FGS & WIP orders
 - Relabel all stock with new labels
 - Enter all open orders into the system
 - Adopt planning schedule
 - Test corrugator interface

- **Go Live**
 - Post Go-Live Phase
 - Support the plant in the processes
 - Adopt reporting
 - Troubleshooting
 - Update the system

All steps have their meaning.

Analysis Phase - this stage is a time of a complete detailed review of the production processes, of the existing information systems, outlines the possible gaps and advantages of the new ERP system. This is the time to make decisions – about time of implementation, budget, responsibility.

Readiness Phase - preparatory phase. Proper planning and building of infrastructure /hardware, communications systems, licenses, interfaces/ is the result of the preliminary investigation and determine the scope of the system. Key users of the system training also predetermines the implementation of the system as deadlines.

Implementation Phase - Implementation of the system. At this stage, the infrastructure is "aliving", the system is adapted according to the legal requirements according to the location of the plant. Specific reports are also being developed, for example required by the public authorities. In this phase starts the training of users in a test ERP environment.

Go-Live Phase - start of the project. Enter all data from the old system to the new ERP, re-label the available product in stock finished product. This stage is a critical. Incorrect input of previously prepared information may result in huge losses from misplaced queries, erroneous pricing and availability. Interfaces to the existing external systems are tested and real work is start.

Post Go-Live – This step is predicted to find the possible problems and mistakes and gives the opportunity to make some fine settings on the system.

3.3. Distribution of responsibilities

The main factor for success of the project is the correct distribution of responsibilities [3]. On the following table can be seen an example of possible distribution.

Table 1. Distribution of responsivities

Project Success Factors
• Management Board - commitment, trust, openness and leadership
• Local Team - expertise, enthusiasm and professionalism
• Local Organization - willingness and ability to implement change
• Process Application Team - expertise, professionalism and leadership
• Project Management - quality, methodology and control
• IT Team - involvement, infrastructure ownership and cooperation
• Software Partner - responsibility, professionalism and support
• Complexity - of the local organization (Products, Processes)

The success of the implementation depends on the tempo of the weakest chain in the future system. According to statistics, only 30% to 60% of the implemented ERP systems are successful. In big and complicated structures, it can be stopped or blocked on a lot of levels for a variety of reasons. For a successful implementation, the following key steps are necessary:

- Fixing the budget;
- Creating a working team of specialists at all levels;
- Choosing the right team leader. He must have organizational experience and full authority;

- Permanent control over the steps and the preparation of the system implementation according to the projects phases.

The main factor is planning and building a stable IT structure, according to the specific requirements of the system – Wi-Fi covering of the stores and production areas, installation of specialized equipment – MDC /machine data connectors/, forklift terminals, scanners, etc.

4. Results of the implementation

The result of the implementation is difficult to be defined uniquely because it's a complex amount from the work of all units, but can be depicted as a sum of the cost savings on the one hand and the realized benefits on the other /Fig. 5/.



Fig. 5. Bidirectional benefit generation

Implementing of ERP system doesn't mean automatically generating benefits. According to recent research, a significant amount of implementations not only doesn't bring revenues, but can also generate losses.

The reasons can be externally, for example, market situation, but also an internal one, due to improper planning of implementation, too short deadlines, inability of staff to work with the system, system or IT infrastructure instability. For example, in wrong built communication connectivity, downtimes and customer problems are possible.

Of course, with a careful analysis of business processes, detailed planning, the creation of an efficient team with sufficient budget and implementation time, the ERP system provides huge competitive advantages for the business.

In this example, the implementation of the Kiwiplan ERP system led to a steady growth of volumes and profits of 3-12% on an annual basis.

Time deployment of Kiwiplan in a plant produced more than 100000 m² of packaging per year was 14 months with a budget of 800,000 €.

For more than five-year exploitation downtime of factory due to ERP system and IT infrastructure problems on an annual basis does not exceed 20 minutes.

5. Conclusion

To conclude, ERP-system implementation does not have an alternative on the industry, because it allows flexible and effective management of the processes. Especially for the medium and large companies.

The availability of information system covering all information processes enables tracking and managing key business activities, processing much more information, and providing access to quality analytics and queries. This allow managers different levels and business executives make the right decisions within the required timeframe. It is a system which unites all the processes in a big industrial factory, gives an opportunity to control all the key activities in the company, and makes access to analyses and references easier.

The presented case in the article is a general example and can be used as a model for implementing different ERP systems in another branches of the industry.

REFERENCES

1. Kiwiplan – official site of the company www.kiwiplan.com
2. Moutaz Haddaraab, Ahmed Elragala (2015). The Readiness of ERP Systems for the Factory of the Future Publisher, Conference on Enterprise Information Systems, HCist 2015 October 7-9, 2015
3. Heechun Yang (2016). Project team right-sizing for the successful ERP implementation, Information Technology and Quantitative Management (ITQM 2016)

Authors' contacts

PhD, Eng. Radoslav Hrishev

Technical University–Sofia, Branch Plovdiv
25 Tsanko Dystabanov St.

4000 Plovdiv, Bulgaria

E-mail: hrishev@tu-plovdiv.bg

SYSTEM DEVELOPMENT FOR MACHINE VISION INTELLIGENT QUALITY CONTROL OF LOW VOLTAGE MINIATURE CIRCUIT BREAKERS

BORISLAV RUSENOV, ALBENA TANEVA, IVAN GANCHEV, MICHAEL PETROV

Abstract: *An intelligent system for automated quality control of manufacturing process applications, based on machine vision is presented in this paper. The quality of many produced parts in manufacturing processes depend on dimensions and surface features. The presented automated machine vision system analyzes those geometric and surface features and decides about the quality by utilizing statistical analysis. Refined methods for geometric and surface features extraction are presented also. The efficiency of processing algorithms and the usage of an advanced analysis as a substitution of human visual quality control are investigated and confirmed.*

Key words: *quality control, machine vision, intelligent systems*

1. Introduction

Many industrial processes use or require visual inspection in quality control as an integrated part of their production stages. Such processes are based on visual perception principles to successfully determine levels of product quality by quantifying its visual appearance in general and some specific visual features, respectively [1]. A visual inspection system is based on machine vision principles by using acquisition cameras and also, one or more industrial computers. The main motivation for machine vision implementation is economic factors, which constantly require less production costs.

One of processes that use machine vision for product quality control is the production in mechanical manufacturing processes [2]. The production phases are more or less automated. The exception is quality control stage with mostly human vision inspection. Some production lines still use human vision in quality control. The main reason lies in complexity of this task. Human resources are used because the visual quality control process is very complex and highly demanding and often should be on-line adaptive on changeable quality requests in classification stage of production. Because of human features limitations as controlling element in production line, man becomes one of the weakest and unreliable links. By replacing the human with machine, the whole process should have better production yield and could be more efficient [3].

Industrial control system (ICS) is a general term that encompasses several types of control systems used in industrial production, including

supervisory control and data acquisition (SCADA) systems, distributed control systems (DCS), and other smaller control system configurations such as programmable logic controllers (PLC) often found in the industrial sectors and critical infrastructures. Industrial automation is a discipline that includes knowledge and expertise from various branches of engineering including electrical, electronics, chemical, mechanical, communications and more recently computer and software engineering.

In order to stay on top of a competitive market, companies have to keep their production costs as low as possible. One element of their strategy is to collect production and control data, analyze it to find improvements, and incorporate those improvements in each new plant.

The role of including quality control in other control systems may lie in control and information flow of plants, in integrating processing machines, in the Manufacturing Execution Systems (MES) [7] that monitor the processes, and in the data-based Enterprise Resource Planning (ERP) [8] system that provides decision support.

Process quality control system is a novel computer-aided process quality control system, which integrates hardware and software. The system could realize quality data collection, transmission, storage, quality monitoring and quality statistical analysis for spare parts production process. It could accomplish the collection and monitoring of quality data automatically in field. Once the production process has problems, it can give an alarm and begin to analyze, providing a basis for process quality control. Furthermore, it can also carry out offline quality statistical analysis of

the quality data derived from the machining field, guaranteeing after-process control of processing quality.

2. Quality Control System

The complete system block diagram in Figure 1 presents the role of the SCADA system; the quality measurement data is collected from the final product and stored in special registers inside the controller/ Remote Terminal Unit (RTU) which in this case is a programmable logic controller (PLC). This data is transferred to the SCADA node using industrial network which could be a local or remote network, this data is analyzed and a control decision to tune the controller if necessary to ensure that the product is within the bounds of required quality.

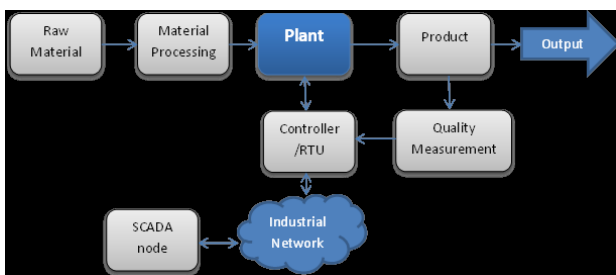


Fig. 1. Quality control system block diagram

Process quality data acquisition and controlling contents include the monitoring of parameters process product quality testing. Process parameters monitoring is realized by measuring the relevant parameters on product quality characteristics. Process product quality testing is achieved by testing products' quality feature in the machining processes or machining process interval. Common processing measurement parameters for example could include cutting force, temperature, spindle motor current changes, vibration and noise signals. Process quality control should establish the correlation between process parameters and the final product quality characteristics, and ensure the quality of the final product by the adjustment to parameters.

Role of SCADA in Quality Management

A production process includes the quality assurance testing of samples from each product lot. The test data is used to produce a certificate of conformance report for such a lot. The data is collected from test equipment, then sometimes manually entered into a customized database form, and then formatted to produce the certificate of conformance. Manual data entry is time-consuming, error-prone, and repetitive. Here is the challenge to introduce a supervisory control system to automate that process and integrate data collection with its

other manufacturing systems. With the automation of the quality assurance process by electronically collecting data from the measurement tools, the operator does not need to write test measurements into a form and then into a computerized spreadsheet. Dozens of samples with up to hundreds of measurements are displayed; with manual work, the system can only handle a limited number of measurements. In being automated, it can be upgraded to manage much more. After data collection, a system can retrieve the expected measurement values from database for the samples of that lot.

A summary screen might be immediately displayed to the operator indicating, for each sample, whether all measurements were within control limits, within specification limits, or outside specification limits. The operator can also view details about each sample and adjust the data manually. Any data modifications are stored in an audit trail, and the initial raw data is kept for historical records. When the operator is satisfied that the data is correct, the system sends this validated data to its database, and it is possible to produce the certificate of conformance. Time saving for these procedures is significant: it could easily take an operator longer to record the measurement data than it takes the tool to create it. A supervisory control system brings the time for the whole process, from measurement to a report, down to a matter of seconds.

3. Description of an application system.

A typical machine vision and quality control system has the structure shown on fig.2.

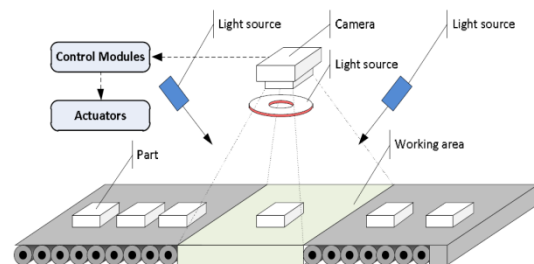


Fig. 2. Structure of a typical machine vision system

The system includes: camera, light sources, a transport system for moving the inspected products, control modules and actuators. Typically, the transport system is a part of the production process. The cameras are positioned so that the inspected parts or finished article falls into the camera's work area. In general, the control module serves to process the camera signal and to communicate with the executive actuator. If it is

provided, it serves with the quality measurement system and the SCADA system of the factory as well, where the inspected product is manufactured. The light source allows accurately determine the right amount of light flux to properly capture the scene.

3.1. Machine vision system of low voltage miniature circuit breakers. In the manufacture of miniature circuit breakers, a number of parameters are considered for their reliable operation. One of these is the distance between the bimetal plate and a specific design hole. For each type of automatic circuit breaker this parameter must be within specified limits. The subject of this article is the development of a system for qualifying this parameter. Figure 3 shows the required Quality Score.

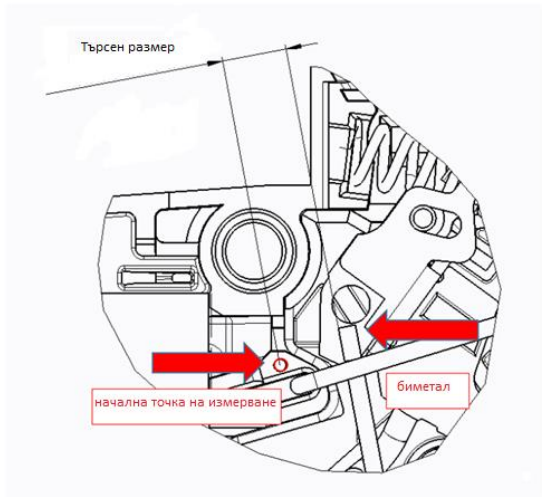


Fig. 3. Parameter for inspection

The task of identifying and separating the inappropriate incoming product from the incoming good one is accomplished by the laboratory machine vision system for testing purposes shown in Fig. 4.

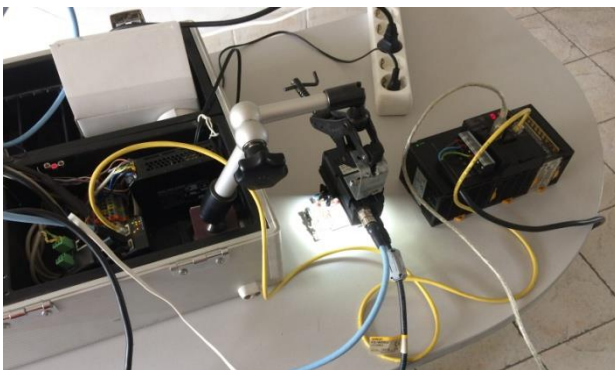


Fig. 4. Laboratory machine vision system

The system of machine vision and quality control of low voltage miniature circuit breakers is constructed schematically according to Fig. 5.

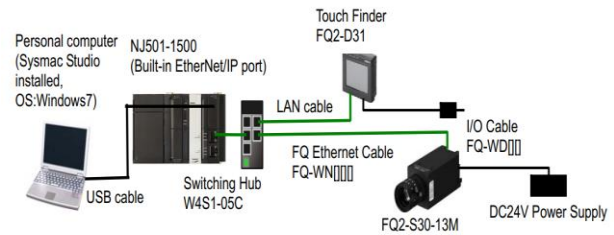


Fig. 5. Structure of machine vision and quality control system

Few words about the part. The size of the inspected part is about 30X30mm. The inspected parameter (distance) is from 2.0 to 6.0 mm. To choose the right resolution is important to know the size of the smallest feature to detect. To have a stable detection is necessary to have at least 3 pixels covering that area. It is not always possible to detect a feature by interpreting a single pixel. Several pixels across the feature can ensure it is not an aberration. In our task the smallest part is a hole with diameter of 0.8 mm. We measure a distance between center of this hole and the bimetal. That means our smallest part is 0.4 mm. If for stable detection we must have at least 3 pixels for this part, then the size of the pixel would be at least 0,13mm.

The smart camera model FQ2 [4] of the manufacturer for industrial automation company Omron was used. The equipment includes a visual inspection camera with embedded setup controller with predefined different quality parameters and embedded light source with controller. This camera already has had an EtherNet/IP™ interface for digital output control. The available inspection camera for this task is model FQ2-S45100N. The field of view is shown on fig.6.

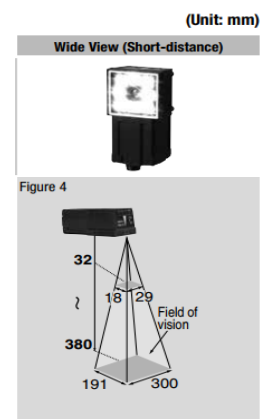


Fig. 6. Field of view of the inspection camera FQ2

The resolution of described camera is 350,000 pixels (752X480). The size of the image element is ½ inch color CMOS sensor (6.4mm X 4.8mm). That means the size of the pixel is 8µm. This is enough for stable detection of the smallest part of the low voltage miniature circuit breakers.

The controller of the visual inspection camera allows operation in industry most frequently recognition modes. Камерата предлага 3 режима за оценяване на геометрични размери: „Edge pitch”, „Edge position” and „Edge width”.

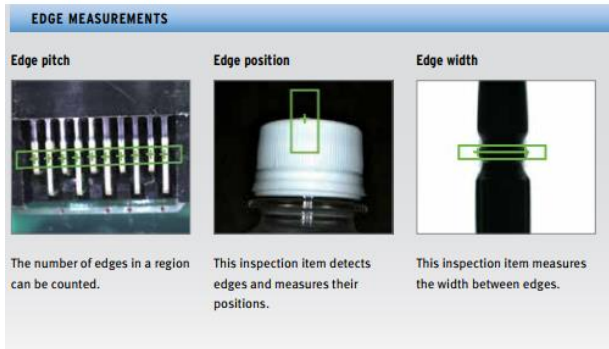


Fig. 7. Edge measurements of the inspection camera FQ2

In this work, the "Edge width" criterion is used. The smart sensor has created group settings files for each particular miniature circuit breaker type. Each group setting includes a shutter speed, white balance, a set of filters, and a function for offsetting the motion position of the scanned object. With these tools, it is possible to highlight the desired characteristic feature in order to work the Edge width criterion. In our case this is the distance between the bimetal plate and the technological opening shown in Fig.3 For the purposes of the manufacturing process, the camera controller also sets a tolerance for the qualified parameter. A tolerance is defined within which the parameter defines the product as fit or the exit beyond these limits defines it as unfit. This is achieved with the Expression function, shown on fig.8.



Fig. 8. Function "Expression" of the inspection camera FQ2

The embedded digital outputs of the intelligent camera are used to send signal to the programmable logic controller PLC of the low voltage miniature circuit breakers manufacturing machine. When an unsecured low voltage miniature circuit breaker is detected, the cycle stops and the actuator removes the product. After that well trained worker adjusts the desired distance between the bimetal and the described in fig.3 hole. When the adjustment procedure is done, the worker returns the circuit breaker for new inspection.

In order to be able to derive full information about the quantity of inferior production which is evident and what is deviation in relation to the set parameters, it is necessary to use more digital outputs. Intelligent camera FQ2 has a datalog function. This feature allows keeping a file with a table of values for each criterion measurement. For a modern production system, this data exchange is not fast and effective enough. It requires human intervention to be processed, which means it can not be included in the enterprise's SCADA system.

When using such a camera included in a SCADA system it is no longer a problem. The high-speed EtherNet/IP™ interface allows to be achieved the necessary data exchange.

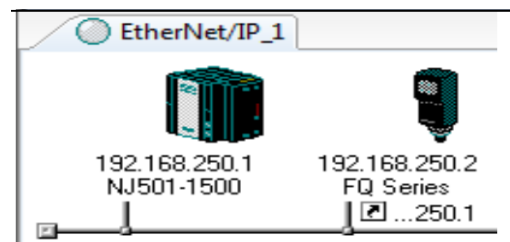


Fig. 9. EtherNet/IP™ data exchange between the inspection camera FQ2 and machine controller.

Significant improvements are also in recognition algorithms. This generation of intelligent cameras have the high-tech Shape Search III algorithm [3], [6]. It allows work with complex objects in light interference conditions and poor background. Several objects entering the camera's work area can be inspected simultaneously, even if objects are poorly or partially illuminated or rotated and overlapped. It is possible to measure distances between the outline of an object, detection of defects such as cracks, for example. There is compensation for displacement and rotation of the object. The ability to recognize text, different types of bar code and 2D code is built in. All of these features make it possible to expand the range of possible smart camera application areas. In addition to the enhanced recognition algorithms, a significant difference is also noted in the camera

Date	Time	Measuren	Scene No.	Judge	I0,JG	I0,C	I0,CR0	I0,X0	I0,Y0
2018/4/4	09:43:46	0	0	1	-1	1	5,8	970,5	502
2018/4/4	09:43:47	1	0	1	-1	1	5,8	975,5	555
2018/4/4	09:43:48	2	0	0	-1	1	5,2	958,5	546
2018/4/4	09:43:48	3	0	0	-1	1	5	882,5	528
2018/4/4	09:43:49	4	0	0	-1	1	4,8	825,5	570
2018/4/4	09:43:49	5	0	0	-1	1	3,9	895,5	581
2018/4/4	09:43:50	6	0	0	-1	1	4	869,5	572
2018/4/4	09:43:51	7	0	0	-1	1	4,5	868,5	665
2018/4/4	09:43:51	8	0	0	-1	1	3,8	878,5	716
2018/4/4	09:43:54	9	0	0	-1	1	4,8	988,5	678
2018/4/4	09:43:54	10	0	1	-1	1	6,5	988,5	683
2018/4/4	09:43:55	11	0	0	-1	1	6,7	988,5	681
2018/4/4	09:43:55	12	0	1	-1	1	6,5	988,5	679
2018/4/4	09:43:56	13	0	0	-1	1	6,6	988,5	684
2018/4/4	09:43:57	14	0	0	-1	1	6,6	988,5	685
2018/4/4	09:43:57	15	0	0	-1	1	6,6	988,5	681
2018/4/4	09:44:01	16	0	0	-1	1	8	988,5	688
2018/4/4	09:44:54	17	0	1	-1	1	6,5	988,5	685
2018/4/4	09:44:55	18	0	0	-1	1	6,6	988,5	681
2018/4/4	09:44:56	19	0	1	-1	1	6,4	988,5	684
2018/4/4	09:44:56	20	0	1	-1	1	5,9	988,5	680
2018/4/4	09:44:57	21	0	1	-1	1	6,1	988,5	680
2018/4/4	09:44:58	22	0	1	-1	1	6,5	988,5	686
2018/4/4	09:44:58	23	0	1	-1	1	6	988,5	678
2018/4/4	09:44:59	24	0	1	-1	1	6,2	988,5	684
2018/4/4	09:44:59	25	0	0	-1	1	6,8	988,5	687
2018/4/4	09:45:00	26	0	1	-1	1	6,5	988,5	687
2018/4/4	09:45:01	27	0	1	-1	1	6,3	988,5	678
2018/4/4	09:45:01	28	0	0	-1	1	6,6	988,5	681
2018/4/4	09:45:28	29	0	0	-1	1	6,7	988,5	688
2018/4/4	09:45:29	30	0	1	-1	1	5,7	988,5	684
2018/4/4	09:45:32	31	0	0	-1	1	6,7	375,5	398
2018/4/4	09:45:33	32	0	0	-1	1	4,5	320,5	403
2018/4/4	09:45:34	33	0	0	-1	1	5,5	988,5	679

Fig. 13. Datalog file of low voltage miniature circuit breaker Type 1

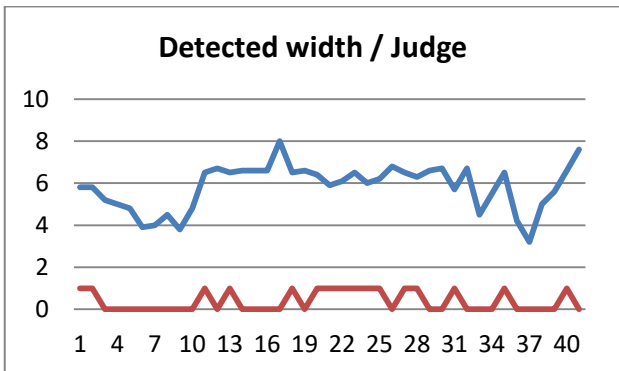


Fig. 14. Time series plot of low voltage miniature circuit breaker Type 1

Type 2:



Fig. 15. Low voltage miniature circuit breaker

Measurement ID	Scene No.	Judge	I0,JG	I0,C	I0,CR0	I0,X0	I0,Y0
1	0		0	0	0	0	0
2	1		0	0	0	0	0
3	2	1	1	84,1643	500,0332	486,9069	5,5982
4	3	1	1	81,8463	499,7947	487,136	5,5171
5	4	1	1	83,5629	499,6835	487,8125	5,625
6	5	1	1	84,4523	499,8132	487,9536	5,625
7	6	1	1	83,2155	500,0086	488,9676	5,2807
8	7	1	1	83,9125	499,9061	488,7804	4,8485
9	8	0	1	84,2345	499,8651	488,6979	4,5143
10	9	1	0	0	0	0	5,5143
11	10	1	0	0	0	0	4,625
12	11	1	1	89,3993	527,212	489,2409	5,625
13	12	0	1	90,8239	595,0901	491,7836	6,3547
14	13	1	1	87,5479	661,4408	494,0627	5,5392
15	14	0	0	0	0	0	6,5392
16	15	1	1	80,7649	716,2281	511,0236	5,7229
17	16	0	1	88,2181	669,6731	511,9482	4,1979
18	17	1	1	92,2955	623,8528	504,8731	5,1942
19	18	1	1	93,8163	577,7789	500,9427	5,625
20	19	1	1	87,5411	543,1352	560,2128	5,2149
21	20	0	1	84,0912	525,9865	577,5245	3,8232
22	21	1	1	90,7844	560,1119	558,8035	5,0437
23	22	1	1	93,5776	564,0857	557,8314	5,625
24	23	1	1	92,4632	564,5972	557,0912	5,3184
25	24	1	1	92,842	564,1059	558,0998	5,625
26	25	1	1	93,4857	564,0561	558,3625	5,625
27	26	1	1	92,6834	564,505	557,7792	5,1804
28	27	1	1	89,0224	563,7449	560,5371	5,0594
29	28	1	1	87,3246	563,695	560,6186	5,0012
30	29	1	1	88,3407	563,7019	560,6701	4,8651
31	30	1	1	86,9342	563,6206	560,1287	5,0764
32	31	1	1	90,3787	563,0682	558,9451	5,625
33	32	1	1	84,2756	511,528	558,9456	5,625
34	33	1	1	85,5408	509,728	556,0991	5,0913

Fig. 16. Datalog file of low voltage miniature circuit breaker Type 2

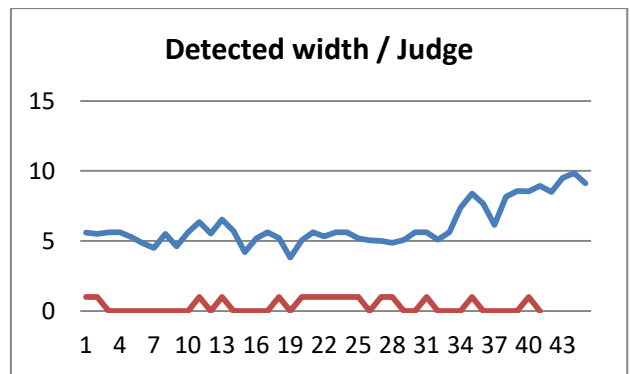


Fig. 17. Time series plot of low voltage miniature circuit breaker Type 2

Thus, the accumulated and evaluated data from the quality control of the low voltage miniature circuit breaker can be used not only to enhance and permanently stabilize the quality. They are also to take preventative measures in servicing the machinery and equipment, involved in the production of the materials and elements involved in the final product.

5. Conclusion

An intelligent system for automated quality control of manufacturing process applications based on machine vision is presented in this paper. The quality of many produced parts in manufacturing processes depends on dimensions and surface features. The presented automated machine vision system analyzes those geometric and surface features and decides about tile quality by utilizing statistical analysis. Refined methods for geometric and surface features extraction are presented also. The efficiency of processing algorithms and the usage of an advanced analysis as a substitution of human visual quality control are investigated and confirmed.

The presented machine vision system has purpose to replace a human vision quality controller in the manufacture of low voltage miniature circuit breaker. The presented system consists of at least two cameras for image registration and one computer. After the image acquisition and trimming basic optic parameters, the geometric and surface analysis is performed. Geometric analysis relies on contour tracing method where several geometric inspection methods are united into one. Using this method, the complete time for analysis is reduced to acceptable limits suited for real time operations.

REFERENCES

- [1] Batchelor B.G., Whealan P.F. (2002). *Intelligent vision systems for industry*.
- [2] Forsyth D.A., Ponce J. (2003)., *Computer vision. A modern approach*. Whiliams publishing house.
- [3] Gonzalez R., Woods R. (2002)., *Digital image processing*. Whiliams publishing house
- [4]<http://www.omron.com/technology/core/thinkAndSee/>
- [5] Said Ibrahim Abu Al-Roos *SCADA Online Product Quality Control*.(2013)
- [6] OMRON Corporation, *Technology overview Pattern Matching Algorithm Shape search III*
- [7] Jürgen Kletti (Ed.) *Manufacturing Execution Systems – MES. Springer-Verlag Berlin Heidelberg 2007. p.271*
- [8] Monk, Ellen; Wagner, Bret (2006). *Concepts in Enterprise Resource Planning (Second ed.)*. Boston: Thomson Course Technology.

Authors' contacts

Organization:

Technical University Sofia, Plovdiv branch

Address: 25, Tsanko Dyustabanov Str.,
4000 Plovdiv, Bulgaria

E-mail: borislav.rusenov@abv.bg

E-mail: altaneva@tu-plovdiv.bg

E-mail: ganchev@tu-plovdiv.bg

E-mail: mpetrov@tu-plovdiv.bg

CONTROL SYSTEM APPLICATION USING LORAWAN

STOITCHO PENKOV, ALBENA TANEVA, MICHAIL PETROV,
KIRIL HRISTOZOV, ROBERT KAZALA

Abstract: *The goal of this work is to apply LoRaWAN to the system control development. A designed network with low cost equipment is used for industrial application. The solution is related to WSN using network protocols and standards summarized in this article. The main goal is focused on WSN to collect data and net operation by combining low energy protocol and transmission method. In the developed network project for industrial data exchange MQTT protocol and LoRaWAN are used. An application for control system purpose is presented. For the system implementation a LoRa gateway is used, combining iC880A LoRaWAN[®] Concentrator board and Rpi 3 as a host controller. The control algorithm, for a loop system operating, is developed. It is implemented in open source programming environment. Some advantages of the developed network with MQTT and LoRa are given. The preliminary investigations for verification are conducted. The real test with sending and receiving data between the connected nodes are made. The recent work can be viewed as an example related to the so called low cost automation.*

Key words: *Industrial Network, LoRaWAN, WSN, Low Cost Automation (LCA), IoT*

1. Introduction

In recent years, many different network specifications are applied to the industrial systems [1]. The communications respectively industrial networks are the key technologies of the present and the near future. Using network in industrial applications is usual, complex, responsible, and sometimes dangerous tasks. More complex problem is when there is no power supply or electricity. There are many different ways to solve, depending of cases. It is case of highly dynamic phenomena a solution of such problem is to use low cost nodes equipped with relevant networked sensors for data collection. Several nodes can be organized and formed a grid which will bring more complexity, and we will gather all needed information. The paper presents wireless networked sensors (WSN) for data collection in control system application. Each node has to be equipped with sensor and communication hardware. If the goal is to measure air pollution environment, or to gather info under the sky, it could be use nodes with GPS to get exact location, but considering power plan that is not energy efficient. Where is needed can be used mobiles and flying robots for deploying these nodes, and at moment when nodes are deployed, it can be marked GPS position. Thus are covered large areas and different surfaces. One of the main factors for

such system development is the implementation cost. In order to achieve it is appropriate to use low-power nodes, which are in the frame of LoRaWAN.

This paper is focused on configuration and development of the LoRa network gate and nodes with temperature and light monitoring system. This work summarizes solutions related to the network protocols and standards for predefined application. Proposition of node based sensor network with specified low energy protocols are presented. The control system application with feedback is planned to be further development.

2. Layer protocols

Regarding to the OSI model is important to discuss the way of data transmission through layers. In this point of view, the transmission method will be summarized.

2.1 LoRa[™] is a proprietary spread spectrum modulation scheme that is derivative of Chirp Spread Spectrum modulation (CSS). It trades data rate for sensitivity within a fixed channel bandwidth. It implements a variable data rate utilizing orthogonal spreading factors, which allows the system designer to trade data rate for range or power. Furthermore, the network performance is optimized in a constant bandwidth. LoRa[™] is a PHY layer implementation and is agnostic with to higher-layer implementations. This allows LoRa[™]

to coexist and interoperate with existing network architectures. This application note explains some of the basic concepts of LoRa™. Modulation and the advantages of the scheme can provide when deploying both fixed and mobile low-power real-world communications networks [6].

2.2 LoRaWAN™ /layer2/ is a Low Power Wide Area Network (LPWAN) specification intended for wireless battery operated Things in a regional, national or global network LoRaWAN would correspond to the [Media Access Control \(MAC\)](#) layer. LoRaWAN targets key requirements of Internet of Things (IoT) such as secure bi-directional communication, mobility and localization services. Its specification provides seamless interoperability among smart Things without the need of complex local installations and gives back the freedom to the user, developer, businesses enabling the roll out of IoT. The choice of transport protocol, when internet connectivity is needed basically, is reduced to two options TCP and UDP. The Protocols allow multiple devices to communicate effectively using the Internet. However, the determination way for different data types, how they are divided and stored in frames, are required. In case of a system design for collecting data in order to reduce the workload, related to the data exchange organization, it is possible to use application-layer protocols. There is a variety of options. One is to use industrial automation protocols. In problem definition arises high cost of implementation. By analyzing network protocols, extra attention should be paid to the model they use for data exchange. Many of the technologies used in modern computer systems use a data exchange model referred to as Request-Response. However, when you try to use such a model for data exchange in the sensor network you can encounter some difficulties.

The Internet of Things networking technology cheat sheet 1.0

Network:	Sigfox	LoRa	NB-IoT (cat NB1)	LTE-M (cat M1)	LTE Cat 0	LTE Cat 1
Type:	PLWAN	PLWAN	DSSS modulation	LTE (cellular)	LTE (cellular)	LTE (cellular)
Low Power:	++++	++++	++++	+++	++	++
Throughput Kbit/s:	0.1	50	100	375	1000	10.000
Bandwidth:	Ultra-narrowband	Narrowband	Narrowband	Low	High	High
Latency:	1 – 30s	Based on profile	1.6 – 10s	10 – 15ms	Unknown	50 – 100ms
Standard:	Proprietary	Proprietary	3GPP Rel. 13	3GPP Rel. 13	3GPP Rel. 12	3GPP Rel. 8
Availability world-wide:	++	+++	++	++	++++	++++
Spectrum:	Unlicensed ISM	Unlicensed ISM	Licensed LTE	Licensed LTE	Licensed LTE	Licensed LTE
Complexity:	Very low	Low	Very low	Low / medium	High	High
Coverage / range:	Medium / high	Medium / high	High	High	High	High
Battery life:	Very high	Very high / high	High	Medium / high	Low	Low
Gateway needed:	Yes	Yes	No, but optional	Optional	Optional	Optional
Signal penetration:	High	Medium / high	Medium / high	Medium / high	Low	Low
Security:	+++	+++	+++	++++	++++	++++
Future proof:	+++	+++	++++	++++	+++	+++

Fig. 1. Physical layer parameters of some network technology

A possible solution is using the Publish-Subscribe method. In this method, the data publishing modules send it to a server called the broker, which then sends the data to clients subscribed to certain information. Using these methods of data exchange allows the clients to receive not all the information sent by the node, but only the data that interests them. There is also no need to constantly call the modules that generate information about the data. On Fig. 1 is presented physical layer parameters of different Low Energy protocols, applicable to the *IoT*. A better option for recent purposes is LoRaWAN. Some features show that it has advantages: low power and complexity, high battery life, low cost implementation, easy maintenance, use of free bandwidth regarding to the listed licensed network types.

2.3 MQTT protocol

Among different options described one of the most appropriate is MQTT protocol (Message Queue Telemetry Transport), details in [5]. It was designed in 1999 for transferring data from telemetry devices. The main goal of the designers was to create an efficient protocol to transfer data from devices with limited hardware resources, which is equipped with a low-performance microprocessor and a small amount of memory. Also expected to work in networks with severely limited bandwidth for data transmission. The protocol uses a publish-subscribe method and transmits the data over TCP/IP or UDP. In its implementation requires a special computer called a message broker. The task of the broker is to collect messages and send them to devices interested in specific information. Fig. 2 shows the organization diagram for data exchanging between Publishers and Subscribers by MQTT Broker.

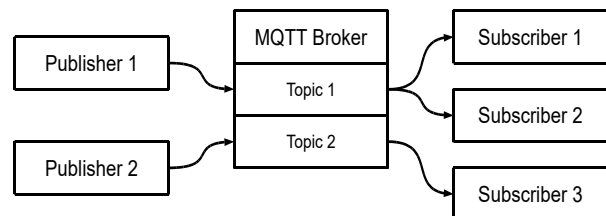


Fig. 2. Organization of data exchange in MQTT.

MQTT protocol messages are assigned to names that are topics. In context of the client and the broker, there is no need to configure the topic. The client sends a message to a specific topic. If there is a particular topic the broker will update its data, in the absence a new topic will be created automatically, to which will be assigned the information transmitted in the message. Topics may be organized in a hierarchical manner using the

separator in the form of a forward slash (/). This allows us to organize data in a broker in a manner similar to the file system. Example topic for networked grid nodes may have the following form:

LoRaNet1/NODE1/sensor3/DATA

An important feature of the MQTT protocol is the ability to manage the quality of service by implementing QoS (Quality of Service). It allows you to manage the way to deliver a message and confirmation of its receipt. LoRaWAN has several different classes of end-point devices to address the different needs reflected in the wide range of applications.

3. System design and configuration

The developed solution for industrial system combines LoRa and MQTT and follows from *Infrastructure Overview* [2], fig.3. Therefore could be achieved LPWAN based LoRaWAN, including Gates (G), Nodes (N), connected to each G, in the frame of the things network (TTN), presented on fig.4. It is evident there is variety of tasks: many sensor types, devices through Gateway to the many user defined applications.

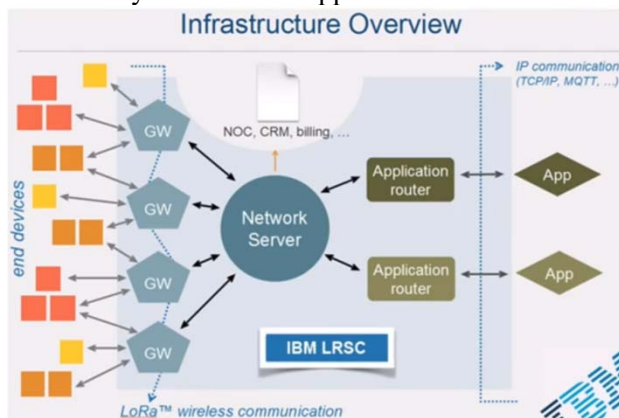


Fig. 3. Overview of the LoRa network

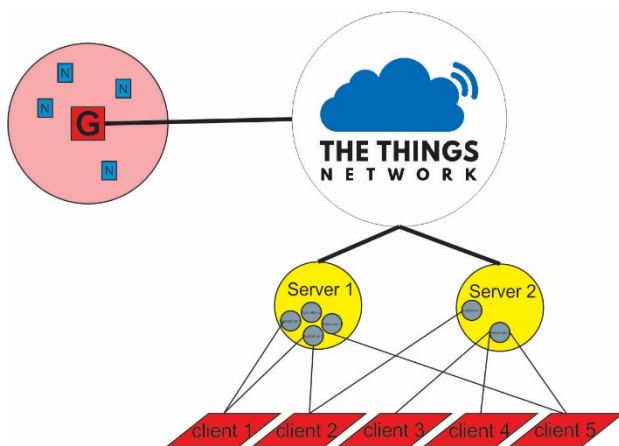


Fig. 4. Overview of the developed network

The network can consist of thousands nodes. In case of new LoRaWAN project is need a Gateway in case with no LoRa coverage. The proposed solution includes a Gateway. This can be viewed as an advantage of the application, fig.4. LoRaWAN uses licence-free spectrum, usually ISM (Industrial, Scientific, Medical) bands to communicate over the air. In Europe, ETSI regulates the ISM band access on the 868MHz and 433MHz bands. The usage of these bands is submitted to limitations: The output power (EIRP) of the transmitter shall not exceed 14dBm or 25mW, and the duty cycle imposed in Europe by ETSI is limited to 1% (for devices) or 10% (for gateways) depending on the used sub-band.

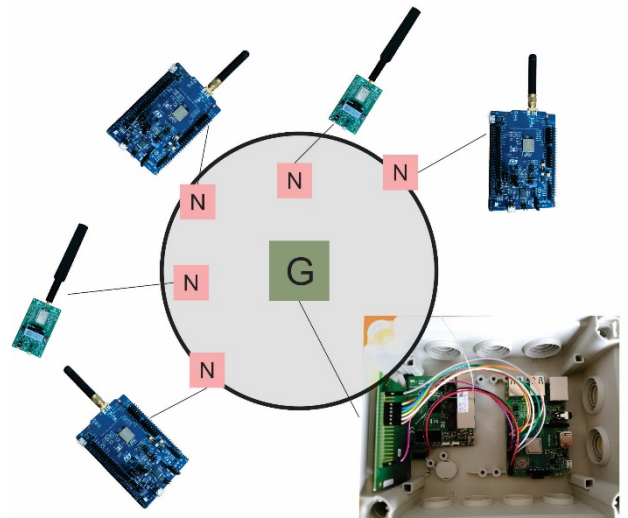


Fig. 5. Designed LPWAN and Assembled Gateway

On the next fig.5 is shown the used and connected LoRa node. In this way is developed and obtained network based on LoRaWAN. The node is equipped with thermometer and light sensor, hence could be used as feedback in industrial control system application with low cost, as well as in standalone security and/or fire alarm system, battery operated in places with no electricity, as well as part of home automation system, etc..

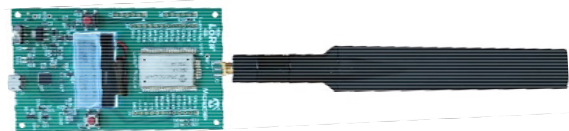


Fig. 6. A single node in the developed network.

LoRa gateway packet forwarder is built by combining iC880A LoRaWAN® Concentrator board at 868 MHz and a Rpi 3 for a host controller, interconnected via SPI bus. Gateway forwards received from Air packets to a public LoRa server, in our case *TheThingsNetwork* using Semtech packet-forwarding protocol. Then MQTT protocol is used to traverse collected packets to our user

application server. It can have multiple gateways to cover wider areas. The Nodes are built by using Murata LoRa SOC type ABZ, which contains SX1278 for radio front-end and STM32L0 as MCU. It is ported an open source LoRa stack implementation by IBM, called LIMC to this particular platform. The nodes are powered from a long live lithium battery CR123A type. It can be combined with different type of sensors connected to analog or digital GPIO's, or using the built-in I2C and SPI buses. For temperature and humidity DHT22 is the cheapest and reliable option, but also have BMP180, HYT261, CCS821. Furthermore each node has 1-Wire iButton compatible bus, a digital thermometer DS18B20 and token ID reader socket for DS1990, which can be used for security authentication. There are multicolor LEDs and a buzzer on node's board for human interaction, buttons and spare GPIOs to wire it in custom user applications. Simple example of and a user application can be environmental real-time monitor and logger. The attached sensors generate feedback data, updating to the network on predefined intervals, or immediately in case of over passing the threshold values. For example temperature, humidity, air pressure, pollution and light sensing are normally changing slowly and will be updated on schedule. If some of the readings change a lot, an immediate update event would be processed. Any external trigger from GPIO (for e.g. PIR sensor or MC) will be transmitted immediately or delayed, depending of the current operational mode – armed, monitoring, status report.

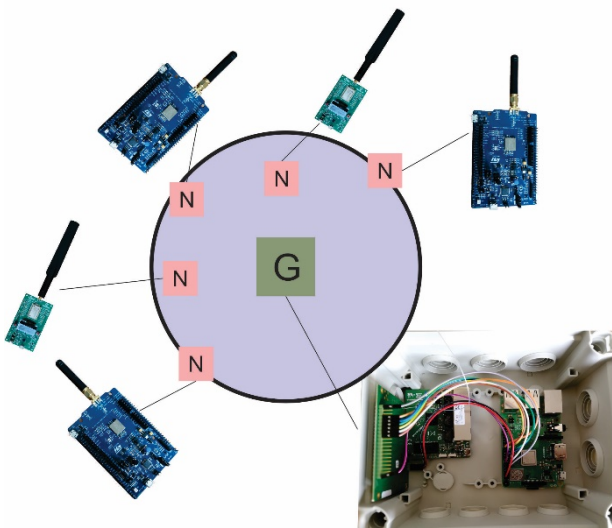


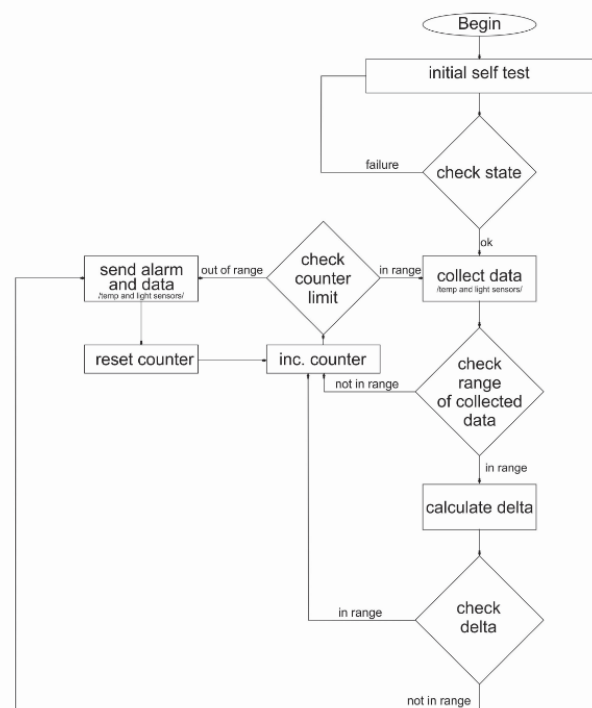
Fig. 7. Components overall view of the developed Network

Each transmission has sequential ID number, and various payload, but always will carry status report for battery level and important circuits (e.g. Temperature loop). The sent information will reach TTN via route from LoRa packet, through the

gateway concentrator, VPN and Internet. Then our application server registered to corresponding MQTT channel and topics will receive the message to process and forward to the user application and front-end which are logged with necessary credentials to have monitoring and/or control authority. All components of the developed Network /Gateway with installed packet forwarder and nodes On the fig.7 the G (iC880A LoRaWAN® Concentrator board), an integrated light and temperature sensors in each N (NODE - Raspberry pi 3) are presented. Hence a MSN with LoRa for control system application is achieved.

4. The developed algorithm

For the WSN operation a simple algorithm is developed, shown below:



The main steps are devoted on checking for state, range of collecting data and delta. A part of program loaded on Raspberry pi 3 as GATE, regarding to the algorithm is:

```

{
  "gateway_conf": {
    "gateway_ID": "B827EBFFFE79235B",
    "servers": [
      {
        "server_address": "router.eu.thethings.network",
        "serv_port_up": 1700,
        "serv_port_down": 1700,
        "serv_enabled": true
      }
    ],
    "ref_latitude": 42.1517592,
    "ref_longitude": 24.7381178,
    "ref_altitude": 17,
    "contact_email": "stoitcho@abv.bg",
    "description": "TTN TU-Plovdiv Gate 001"
  }
}

```


Setting the **network key** and **app key** to the NODE using terminal is given on fig.8. An assigned Gateway ID, to the “router of the things network”, using port 1700, giving the position of our gate to predefined coordinates, is established.

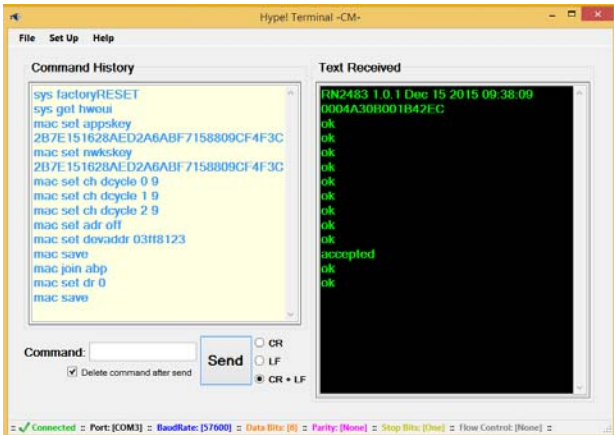


Fig. 8. Programming of the Node

5. Experiments and results

The real test and verification in the presented WSN with sending and receiving data between the connected node (N), (G) and TTN are made. The followed figures show the results of information exchanging (packets and payload) in the frame of the developed network, denoted with TU-Plovdiv Gate 001. Due to the successful registration at TTN the configured and programmed gate starts operate. The *Network Session Key*, *App Session Key* and *Device Address* are set to the nodes, using ABP activation, can be seen on fig. 9, 10. The messages are forwarded between the nodes. It is established real communication between recent nodes in the frame of TTN.

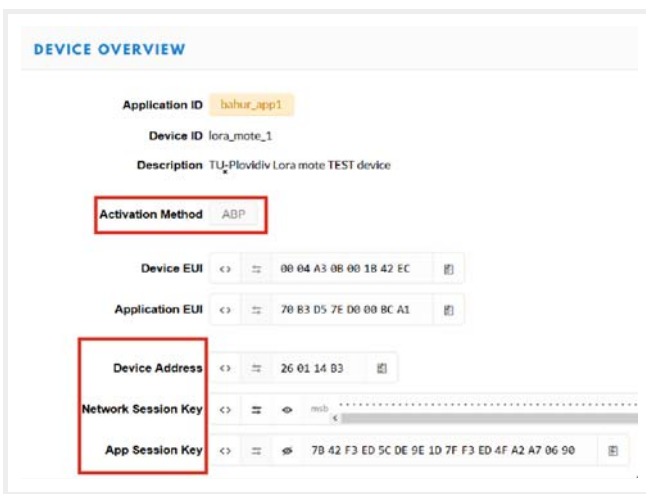


Fig. 9. Registering to TTN

The device parameters are shown on next fig.10. Real received messages are count (average 5382) and can be seen in menu gateway overview,

Hence it can be seen the chosen activation method ABP, device address, NSK and AppSK.

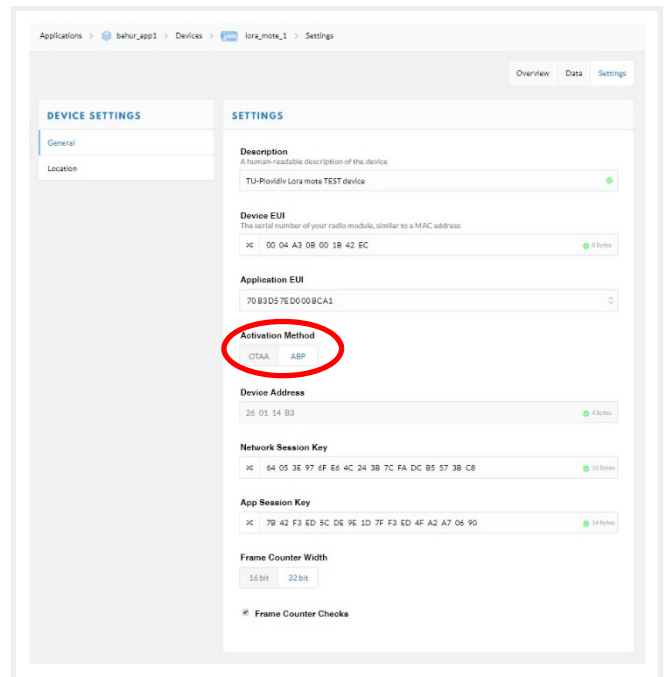


Fig. 10. Setting an activation method

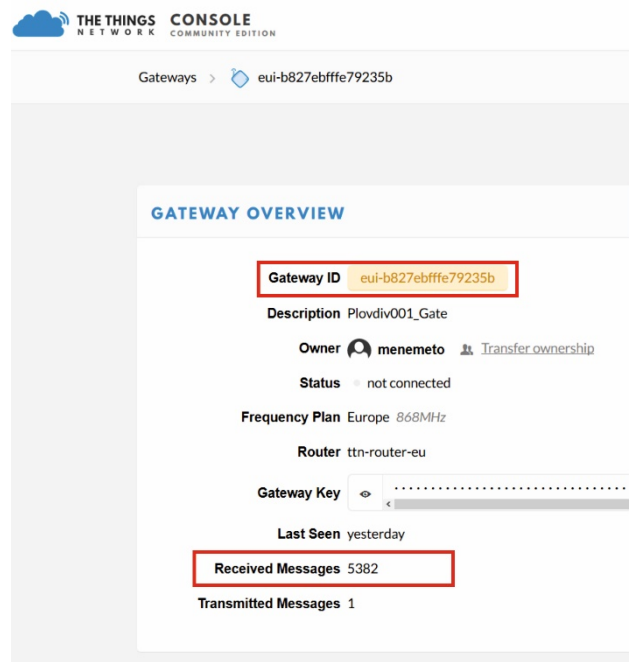


Fig. 11. Overall view of the GATE



Fig. 12. Data of the GATE operating

Single packet data can be seen on fig.13. Payload size, used frequency, etc. can be seen on fig.12. These data could be appropriately representing the temperature and light monitoring readings. The next step of program development is to ensure the proper, useful and readable interface to the users.

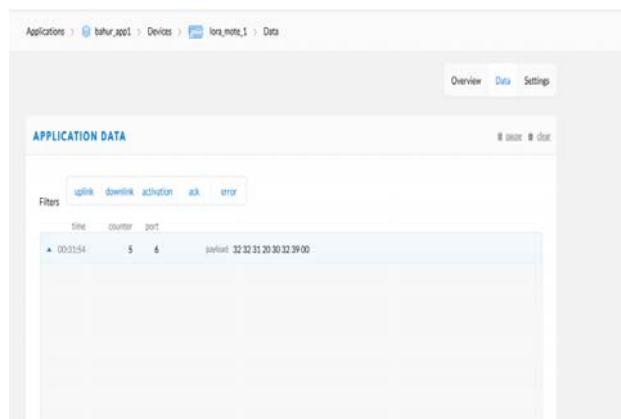


Fig. 13. Confirmation of the application data transmission

6. Conclusion

The work is devoted on developing, configuration and commissioning of the LoRaWAN for system control purpose in the case of missing 220V supply. A solution is found out by combination between the LoRaWAN network and MQTT protocol.

Some advantages of the system are: easy maintenance; low cost; low energy consumption; reliability and control system applicability; easy connection to the Internet, long lifecycle. The number of system nodes can be enlarged with up to more than 1000 nodes connected with one gate (G). In the particular test application for the industrial network development it is obtained promising results. To industrial environment for monitoring process values like temperature, air conditioning systems and etc. Each “node” could present sensor one or many, or robot sensors, connected the main node through the Gateway with LoRa. The paper shows an example of the use of open network protocols in networked mobile robot sensor system. MQTT protocol allowed for the use of standardized methods for data exchange in the sensor network. Also greatly simplified the integration of new nodes, the use of nodes information from other systems and integrating with the Internet. The use of open protocols simplifies software development work, especially when it consists of a large number of independent nodes. Also simplifies the maintenance process, since it is possible to read information about the industrial values and the status of individual sensors without having to use special tools.

The project is made to the stage designed, configured, program developed in the frame of the network and readiness for low cost control system application. The achieved results (fig.9-13) show the promising future work. Next improvement could be devoted on program development to ensure the proper, user interface, regarding to the process control tasks. The future work will be focused on many other applications designing, depicted by separate NODES on fig. 4.

REFERENCES

1. M. Rizzi, P. Ferrari, A.Flammini, E.Sisinni, Evaluation of the IoT LoRAWAN solution for distributed measurement applications, IEEE Transactions on Instrumentation and Measurement September 2017, DOI: 10.1109/TIM.2017.2746378
2. IBM Research - Zurich: Industry & Cloud Solutions, Industry Solutions, Cross-Industry: LRSC, LoRaWAN, www.research.ibm.com - February 7, 2015
3. <https://www.thethingsnetwork.org/>
4. IMST GmbH, iC880A Quick Start Guide – Quick Start Guide Document ID:4100/40140/0078
5. <https://market.thingspark.com/solutions/tracking/ic880a-lorawan-concentrator-868mhz-404802>
6. Comparison of 6LoWPAN and LPWAN for the Internet of Things – Article in Australian Journal of Electrical and Electronics Engineering, December 2017 DOI 10.1080/1448837X.2017.1409920

Department of Control Systems
 Technical University–Sofia, Branch Plovdiv
 25 Tsanko Diustabanov St.
 4000 Plovdiv

Phone : +359 032 659 585
 E-mail: stoitcho@abv.bg
 E-mail: altaneva@tu-plovdiv.bg
 E-mail: mpetrov@tu-plovdiv.bg
 E-mail: khristozov@icloud.com
 E-mail: rkazala@tu.kielce.pl

THE IMPORTANCE OF SOLAR ANGLES IN MPPT SYSTEMS

AHMET SENPINAR

Abstract: *The movement of the world around the solar system and its axis results in the change in the position of the sun according to the world. This causes the different seasons and day lengths in year. One way of the utilizing solar energy is using solar array systems which generates electricity energy as expose to the solar radiation on it. At these systems, to obtain high efficiency from the solar radiation, the fixed and tracking solar array systems are used. Fixed systems are located at certain slope with horizontal. This slope changes with seasons and position of the region. Calculating optimum slope angle depend on the seasons in year at the region, the amount of energy which was obtained the solar radiation is provided increasing.*

Keywords: *Solar Angles, Fixed Array, Slope Angle.*

1. Introduction

Some kind of energy like the sun, wind, geothermal, biomass and wave energy produce a part of alternative/renewable energy sources. Solar energy has significant importance on human health and environment because of its abundance, renewability and pollution free.

Energy obtained from the sun on earth per unit time is known as the solar constant and is represented by GSC. The value of the solar constant as accepted by the World Radiation Center is 1367 W/m^2 ($1.96 \text{ cal/cm}^2 \text{ min}$) [1]. The sun is a gaseous body, with a mass of approximately $2 \cdot 10^{30} \text{ kg}$ and a diameter of $1.39 \cdot 10^9 \text{ m}$. The distance from the sun to earth is approximately $1.49 \cdot 10^{11} \text{ m}$.

There are different using areas of pv systems and these increase day by day on the World. Some of these areas are as follow: house with pv energy, street illumination, cooling, pumping water, traffic signalling, pv plants, hybrid (sun+wind) systems, space systems, and telecommunication systems etc. [2-4].

In some PV systems, motion of the sun is wanted to be tracked in a day. The purpose of these systems is to increase the system's efficiency. Therefore, two different systems which were the fixed and the tracking system were established for applications of the stand-alone PV system [5, 6].

Different sun-tracking systems have been developed to track the sun's movement across the sky [7-19].

The value of the solar angles coming to the surface of the array in any fixed or tracking array varies with the geographical location of the settlement in which the array is located, the date of

that day and the time of day. Fixed systems are systems in which the array of solar cells is placed with a specific fixed slope. The slope angle changes according to the season and region. Tracking arrays follow the sun to maximize the incident beam radiation on their surfaces. Tracking control is based on angles of incidence and surface azimuth angles. Solar tracking systems are more expensive and complex than fixed systems.

PV arrays can be mounted to track the sun, but fixed systems must be maintained at a certain angle to the horizontal to fully exploit available sunlight at the location. If this slope angle is determined well, the amount of insolation and the generated energy increase. To maximize energy, solar panels, such as photovoltaic modules, are usually oriented toward the equator with an optimal slope angle from the horizon, which depends on climatic conditions and site latitude [20-23]. Some solar angles as follows;

1.1. Latitude angle

The angle θ on the earth's surface measured North or South of the equator to a point is its latitude. Latitude values increases toward the poles, with the North pole being 90° , and the South pole -90° .

1.2. Declination angle (δ)

The declination angle is the angular position of the sun at solar noon with respect to the plane of the equator, north positive, $-23,45^\circ \leq \delta \leq 23,45^\circ$. The declination δ can be found from the equation of Cooper:

$$\delta = 23,45 \cdot \sin\left(\frac{360 \cdot (284 + n)}{365}\right) (\text{°}) \quad (1)$$

where n represents the day of the year (n=1, for 1 January) [1].

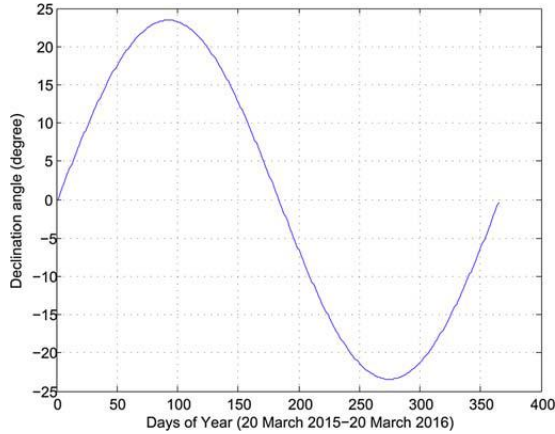


Fig.1. Change of declination angle

1.3. Zenith angle (θ_z)

The zenith angle (θ_z) is the angle between the vertical and the line to the sun and is calculated as follows [1]:

$$\cos\theta_z = \cos\delta \cdot \cos\phi \cdot \cos\omega + \sin\delta \cdot \sin\phi \quad (2)$$

1.4. Solar altitude angle (α_s)

It is the amount of angle that the horizontal direction with the direction of the sun. Since the zenith angle has been completed at 90° , here is the altitude of the sun;

$$\alpha_s + \theta_z = 90^\circ, \quad \alpha_s = 90^\circ - \theta_z \quad (3)$$

1.5. Solar incidence angle (θ)

It is the amount of angle between the light coming directly to a surface and the normal of that surface. This angle represents the angle of incidence of sun. This angle is calculated as follows:

$$\cos\theta = \cos\theta_z \cdot \cos\beta + \sin\theta_z \cdot \sin\beta \cdot \cos(\gamma_s - \gamma) \quad (4)$$

where γ , is the surface azimuth angle.

1.6. Slope angle(β)

This angle is the angle between the plane of the surface in question and the horizontal, The slope angle varies between $0^\circ \leq \beta \leq 180^\circ$.

$$\tan\beta = \tan\theta_z |\cos\gamma_s| \quad (5)$$

2. The tilt angle of the fixed or tracking array

The fixed systems, on the other hand, are installed at a certain tilt angle which varies depending on the geographical location and time of day. The tracking systems are classified into two groups: a single axis and two axes. The rotation motion in practice is usually horizontal east-west, horizontal north-south, vertical, or parallel to the earth's axis. For a plane rotated about a horizontal east-west axis with a single daily adjustment so that beam radiation is normal to the surface at noon each day [1],

$$\cos\theta = \sin^2\delta + \cos^2\delta \cdot \cos\omega \quad (6)$$

and the slope angle of this array

$$\beta = |\phi - \delta| \quad (7)$$

where, θ represents angle of incidence, which is the angle between the beam radiation on a surface and the normal to that surface, ω represents hour angle which is the angular displacement of the sun east or west of the local meridian due to rotation of the earth on its axis at 15° per hour.

Tracking PV systems are classified by their motions. Rotation can be about a single axis (horizontal east-west, horizontal north-south, or parallel to the Earth's axis), or it can be about two axes. For a plane rotated about a horizontal east-west axis with continuous adjustment to minimize the angle of incidence [1],

$$\cos\theta = (1 - \cos^2\delta \cdot \sin^2\omega)^{1/2} \quad (8)$$

the slope angle of this array

$$\tan\beta = \tan\theta_z |\cos\gamma_s| \quad (9)$$

The amount of insolation received at different locations across Turkey varies

according to geographical position and local climatic conditions. Thus, the researcher calculated an optimal slope angle for 9 cities across Turkey using data on insolation levels and meteorological records from 2014.

The meteorological data for the 9 cities are shown in Tables 3 and 4 along with the

average monthly and seasonal optimal slope angles. First, monthly average values for each city were calculated. The annual average slope angle value ($0,9*\emptyset$) degrees was then calculated using the monthly average values [1].

Table 1. Monthly Average Values of Optimum Slope Angles for 9 Different Cities in Turkey

City	January	February	March	April	May	June	July	August	September	October	November	December
Ankara	60,40	52,88	41,94	30,06	20,75	16,48	18,45	26,26	37,56	49,40	58,61	62,65
Mugla	57,96	50,44	39,50	27,62	18,31	14,04	16,01	23,82	35,12	46,96	56,17	60,21
Elazig	59,52	52,00	41,06	29,18	19,87	15,60	17,57	25,38	36,68	48,52	57,73	61,77
Samsun	62,01	54,49	43,55	31,67	22,36	18,09	20,06	27,87	39,17	51,01	60,22	64,26
Istanbul	61,85	54,33	43,39	31,51	22,20	17,93	19,90	27,71	39,01	50,85	60,06	64,10
Malatya	59,05	51,53	40,59	28,71	19,40	15,13	17,10	24,91	36,21	48,05	57,26	61,30
Mersin	57,32	49,80	38,86	26,98	17,67	13,40	15,37	23,18	34,48	46,32	55,53	59,57
Sinop	62,85	55,33	44,39	32,51	23,20	18,93	20,90	28,71	40,01	51,85	61,06	65,10
Canakkale	60,93	53,41	42,47	30,59	21,28	17,01	18,98	26,79	38,09	49,93	59,14	63,18

Table 2. Seasonal and Annual Average Values of Optimum Slope Angles for 9 Different Cities in Turkey

Latitude-Longitude and Seasonal values as degree							
City	Latitude	Longitude	Spring	Summer	Autumn	Winter	Annual values
Ankara	39,56	32,52	30,92	20,40	48,52	58,64	35,60
Mugla	37,12	38,22	28,48	17,96	46,08	56,20	33,40
Elazig	38,68	39,14	30,04	19,52	47,64	57,76	34,82
Samsun	41,17	36,20	32,53	22,01	50,13	60,25	37,05
Istanbul	41,01	28,58	32,37	21,85	49,97	60,09	36,91
Malatya	38,21	38,19	29,57	19,05	47,17	57,29	34,38
Mersin	36,48	34,38	27,84	17,32	45,44	55,56	32,83
Sinop	42,01	35,09	33,37	22,85	50,97	61,09	37,81
Canakkale	40,09	26,24	31,45	20,93	49,05	59,17	36,08

3. Conclusion

The continuity of the sun has advantages because it does not require maintenance costs, it has no adverse effect on the environment, it is established in a short time, and the storage of the energy obtained. In addition, obtaining energy from these sources allows us to use our existing resources for a longer period of time. With the help of the solar array, the fixed arrays are placed in such a way as to have the slope angles varying in MPPT systems according to the geographical position of the region used, in order to utilize the sunlight at high speed while obtaining the electric energy. Monthly, seasonal and annual optimum values of different cities were calculated for optimum slope angle. When these values are examined, it is advantageous to use optimum values which vary according to the seasons instead of annual average values, since there is a big difference between winter and summer. At the end of each season, these slopes are manually adjusted to provide efficient operation of the solar array system.

REFERENCES

1. Beckman WA, Duffie JA. *Solar engineering of thermal processes*. 2nd ed. Canada: John Wiley and Sons Inc.; 1991.
2. Chambouleyron, I. (1996). Photovoltaics in the developing World. *Elsevir.Energy*, Vol.21, No.5, p385-394.
3. Green, J.M., Wilson, M., Cawood, W. (2001). Maphephethe rural electrification (photovoltaic) programme: the constraints on the adoption of solar home systems. *Development Southern Africa*, Vol:18, No.1, p19-30.
4. Kuwano Yukinori (1998). Progress of photovoltaic system for houses and buildings in Japan. *Elsevie., Renewable Energy*, Vol:15, p535-540.
5. Helwa, N.H., Bahgat, A.B.G., El Shafee, A.M.R. and El Shenawy, E.T. (2000), Maximum Collectable Solar Energy by Different Solar Tracking Systems, *Taylor and Francis, Energy Sources*, vol.22, pp. 23-34.
6. Şenpınar, Ahmet (2005), The Control of The Stand-Alone Photovoltaic Cell Systems By Computer, PhD Thesis, *Firat University Graduate School of Natural and Applied Sciences*, Elazig, Turkey.
7. Al-Soud MS, Abdallah E, Akayleh A, Abdallah S, Hrayshat ES. A parabolic solar cooker with automatic two axes sun tracking system. *Appl Energy* 2010;87:463–70.
8. Kuo YC, Liang TJ, Chen JF. Novel maximum-power-point-tracking controller for photovoltaic energy conversion system. *IEEE Trans Ind Electron* 2001;48:3.
9. Yi Ma, Guihua Li, Runsheng Tang. Optical performance of vertical axis three azimuth angles tracked solar panels. *Appl Energy* 2011;88(5):1784–91.
10. Tian Pau Chang. The gain of single-axis tracked panel according to extraterrestrial radiation. *Appl Energy* 2009;86(7–8):1074–9.
11. Hadi H, Tokuda S, Rahardjo S. Evaluation of performance photovoltaic system with maximum power point (MPP). *Sol Energy Mater Sol Cells* 2002;2670:1–6.
12. Tian Pau Chang. Output energy of a photovoltaic module mounted on a single axis tracking system. *Appl Energy* 2009;86(10):2071–8.
13. Koutroulis E, Kalaitzakis K, Voulgaris NC. Development of a microcontrollerbased, photovoltaic maximum power point tracking control system. *IEEE Trans Power Electron* 2001;16(1):46–54.
14. Prapas DE, Norton B, Probert SD. Sensor system for aligning a single-axis tracker with direct solar insolation. *Appl Energy* 1986;25(1):1–8.
15. Karimova KhS, Saqib MA, Akhter P, Ahmed MM, Chattha JA, Yousafzai SA. A simple photo-voltaic tracking system. *Sol Energy Mater Sol Cells* 2005;87:49–59.
16. Roth P, Georgiev A, Boudinov H. Design and construction of a system for suntracking. *Renew Energy* 2004;29:393–402.
17. Illanes R, De Francisco A, Torres JL, De Blas M, Appelbaum J. Comparative study by simulation of photovoltaic pumping systems with stationary and polar tracking arrays. *Prog Photovolt: Res Appl* 2003;11:453–65.
18. Tomson T. Discrete two-positional tracking of solar collectors. *Renew Energy* 2008;33(March):400–5.
19. Mohamad AA. Efficiency improvements of photo-voltaic panels using a Suntracking system. *Appl Energy* 2004;79:345–54.

20. Gunerhan H, Hepbasli A. Determination of the optimum tilt angle of solar collectors for building applications. *Building and Environment* 2007;42(2):779e83.
21. Gopinathan KK, Maliehe NB, Mpholo MI. A study on the intercepted insolation as a function of slope and azimuth of the surface. *Energy* 2007;32(3):213e20.
22. Tang R, Wu T. Optimal tilt-angles for solar collectors used in China. *Applied Energy* 2004;79(3):239e48.
23. Benghanem M. Optimization of tilt angle for solar panel: case study for Madinah, Saudi Arabia. *Applied Energy* 2011;88(4):1427e33.

Author:

Ahmet SENPINAR

Firat University, College of Technical
Science, Elazig/Turkey
asenpinar@gmail.com

USE OF INTERFERENCE WEDGED STRUCTURE AS AN ATTRACTIVE, SIMPLEST, LIGHT POWER DIVIDING ELEMENT

MARGARITA DENEVA¹, PEPA UZUNOVA², VALKO KAZAKOV¹, VANIA PLACHKOVA¹,
KAMEN IVANOV¹, MARIN NENCHEV¹ AND ELENA STOYKOVA³

Abstract: *Based on our previous experience in the field of the interference wedged structure (IWS) [2-5], we have considered and developed new competitive applications of such structures. We demonstrate that IWS can be used as an attractive, simplest, light power dividing element for fixed wavelength spatially and spectrally narrow light (laser) beam. This assures: 1) precisely and variably controlled ratio of the reflected and transmitted power by simple translation of the list-like wedged structure (having wedge angle of $\sim 10^{-5}$ rad) in its plane; 2) the separation is practically without energy losses; 3) in the optical scheme with complex geometry of the beam propagation, the power ratio control do not provide beam propagation change. In the work we also show that, as an additional advantage, the power ratio control in some specific cases can be combined with spectral control with noted advantages. We show that the application of the composed tunable interference wedged structures [6] can assure an additional advantage such as guaranteed operation for a single line for complex wavelength combined beams.*

Key words: *wedge interference structure, light power dividing element*

1. Introduction

The well established interferential devices, based on the Mickelson, Fabry-Perot, Mach-Zender type interferometers and Bragg's gratings, have essential practical applications [1]. The applications, presented in the literature, include precise measurements of widths and components of spectral lines, spectral analyzers and filters, fiber optics, use in measurement systems, quality assessment of optical elements, metrology, interferential microscopy and spectral control in laser devices. Here, we consider specific type of interferometer devices, which are not so popular, however present essential potential – Interference Wedged Structures (IWS) [1-6]. Actually, the main and known type of wedged structure is the Fizeau-Interferometer (FI) and its solid version – Interference Wedge (IW). There are many similarities between the FI and the type of Fabry-Perot Interferometer. Fabry-Perot interferometer is a structure built by two transmissive **strongly parallel** reflectors whereas Fizeau Interferometer consists of two partially transmissive reflectors **inclined at a small angle in the order of 10^{-5} rad** with respect to each other. This difference leads to essential difference in the properties [1-3]. The

discussed devices present potential for both instrumentation base for scientific research and instrumentation for measurement in modern industrial and control activities (micrometer displacements, spectral analysis for detection of material composition, healthcare and ecology).

In the previous our work we have developed more detailed theory and cycles of experimental investigation and spectral application of interference wedged structures, especially for attractive laser spectral control [2-5]. We have found a new property of IW – the non-Snellius spectral selective reflection and have introduced new spectral selective element – Reflecting Interference Wedge (RIW), the last being with an essential potential in laser technology [3-5]. In the last our works we have also introduced new elements - Tunable Composed Wedged Interference Structures (CWIS) [6] that present list-like elements with selection of a single, narrow resonance, tunable at the entire length of the structure by simple translation in its plane. These structures represent combinations of layer-mirrors, separated by transparent wedged multi-layers with convenient parameters. The pioneering analysis

and experimental investigations show essential potential of such structures. Nevertheless and actually the CWIS have promising applications, some of which are considered in the work. The noted potentials make necessary ongoing scientific research activity to develop knowledge of interference phenomena in such wedged structures and elements and to propose and develop new ideas concerning the use of its potential.

Actually, the main applications of the IWS relate with use of its spectral properties - spectral selectors and spectral control elements [1,7]. The main aim of the work is, as development of previous our ideas [8], to present and demonstrate one new and competitive application of the interference wedged structures, concerning its potential in solutions of energetic problems, where the laser beam is attractive, simplest, light power dividing element, assuring: precisely and variably controlled ratio, work practically without energy losses and without beam propagation direction change.

2. Basic principle of the use of wedged structures as laser beam splitting loss-less elements with precise division control and without beam propagation disturbance

The IWS's, which we have investigated, constitute each by single or multi separating transparency layers with reflecting layers – mirrors in the separating and the end surfaces of the wedged layers. Due to the compactness and simplicity of employment, such hard structures are preferred for practical use and this is the reason that we have investigated mainly such type of IWS. Also, here, as in pioneering work, concerning presentation of the noted in the Introduction new technique, we will direct our attention to the applications of simplest single wedged layer structure (single angle structure). Such structure represents so-called Interference Wedge (IW), noted above [1-3]. The application of the other, more complex–composite tunable wedged interference structure, which application is more convenient for certain case (high spectral resolution and sensibility with power variation by translation) will be object of other publications.

The typical schematic, with noted in the picture composed elements and the real view of the IW is shown in Fig.1a and 1b. The construction parameters of the IW's, used in the work, are – length of the wedge arms of (2-5) cm, the thickness of the separating wedge layer $e \sim (2-10) \mu\text{m}$ (at the half length of the wedge arms; must be also quartz

plate $\sim 100-800 \mu\text{m}$) and refractive index of $n_1=1.5$; apex angle $\alpha = (2-5) \times 10^{-5}$ rad; reflectivity of consisted multi-dielectric mirrors R_i of (70-90)%. This composition with full thickness of $\sim (7-15) \mu\text{m}$ and for the noted apex angle difference between the thickness of the starting of the structure and its end is (1-2) μm . Thus for the wedge arms' length of 4 cm the structure is practically completely parallel plate (Fig.2). The structure was layered on the ~ 1 mm thickness well parallel glass or quartz plate. Thus all the composition - the wedged structure at the support, was practically completely flat-parallel list-like element.

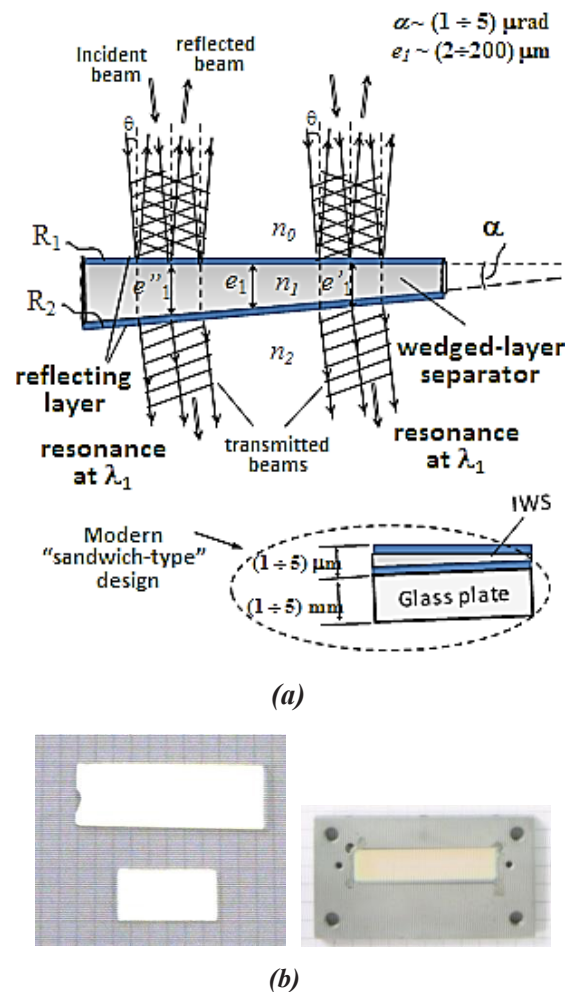


Fig.1. (a) Schematic of the one angle IWS structure with sown two transmission resonances at the same wavelength λ_1 ; **(b)** photograph of two working structures (left, composition of wedged separating layer with $e \sim 5 \mu\text{m}$ with despite two layer-mirrors and at 1 mm glass supports; right – the IW in metallic holder, two mirrors with quartz glass separating wedge).

The property of IW to be spectral filter [1-3], tunable by simple translation of its plane, is illustrated by the schematic in Fig.2(a). The

incident beam is composed by superimposed beams with a series of wavelengths $\lambda_1, \lambda_2, \lambda_3, \dots$, and at a given line, parallel of the wedge angle apex line, the IW transmits only one beam at line resonant λ_2 ; with translation of IW in its plane along the wedge arm, the selection – linear with translation, is at other resonant and respectively IW transmits for the beam with other wavelength - e.g λ_5 , etc.

The real photograph, showing the different line of the transmission for two-wavelength illumination, is shown in Fig.2(b) (the illuminating beam contains the wavelengths 0.6328 μm and 0.5951 μm – red and yellow He-Ne lasers). With sliding the IW, the wavelength transmission resonance is changed. For obtained such wavelength separation the beam must be with small diameter, compared with the distance of the line of equal thickness of the wedge. The typical distance for the given above wedges is 10-15 mm and thus, the beam diameter must be of order of 1-3 mm to coincide well with linewidth of the wedge resonant line.

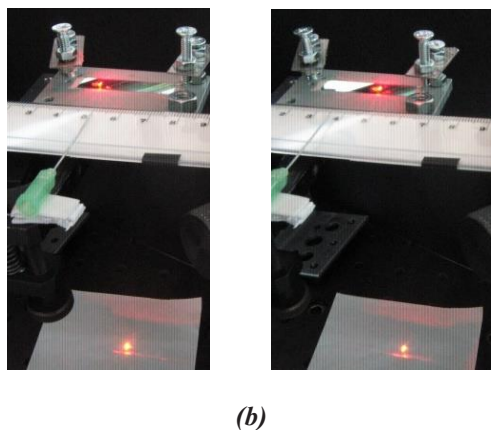
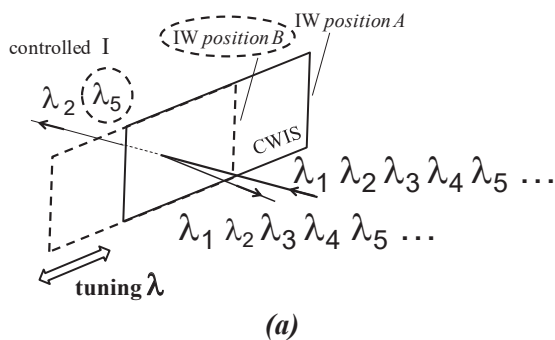


Fig.2. The properties of IW to be spectral filter, tunable by simple translation of its plane. (a) Schematic illustration and (b) Real photographs that show the transmission in two different wavelengths at two places of the interference wedged structure (in the case IW). The change of the place is by translation of the structure in its plane (left for 0.6328 μm and right for 0.5941 μm).

For all experimental investigations, concerning the study of the interference wedged structure (reported here and in series of other) we have composed complex laboratory installation. The installation permits to superimpose in one beam many laser beams (red - 0.6328 μm and yellow - 0.594 μm from He-Ne lasers, green – 0.53 μm from Nd-YAG second harmonic, narrow spectrum and large spectrum semiconductor red \sim 0.63-0.65 μm lasers and variety of semiconductor diodes – red, yellow, green, blue). The illumination is with large parallel beam of 6 cm diameter and with acceptable uniformity of the typically illuminated area of 4 cm or with a small diameter of 1-3 mm. The light powers (of order of milli-watts) were measured via high quality professional “Power meter Thorlabs” (measurement from nW to tens of mW with possibility of registered wavelength adjustment). The graphs of relative intensity distribution of incident and the passed beam profiles were obtained using professional beam scope or by computer treatment of the photographs of the spots using appropriate programming and control for absence of saturation.

Let's we illuminate the IWS (in particular IW) with a narrow spectral line beam (\sim 0.05-1) nm and let's this beam has diameter of order of 1-3 mm. Outside the resonance line the beam will be completely reflected by the structure. When the beam approaches to the resonance line, the IWS starts to be particularly transparent for the beam and with approaching to the exact place (line) of the resonance the transmissivity will increase. The beam will be separated at two parts – transmitted and reflected. The ratio of the two parts will depends on the distance of the beam from the line of exact resonance (maximum transmission). The approaching of the beam to the resonance line can be obtained either by translation of the IWS in its plane or by beam translation. In practice, the interesting case is by translation of the IWS in its plane. Below, are shown the results of our calculations, prepared following the approach in [2]. The calculation gives the wedge transmission as a function of distance from the exact resonance. In the example of calculation we consider IW as one-angle structure (IWS, Fig.1a) with optical thickness of 20 μm , wedge apex angle of 1.2×10^{-5} rad, equal reflectivity of the mirrors $R=0.9$, length of the side planes of the wedge $l=4$ cm. Here, for correct comparison theory-experiment, we will study Gaussian beam with 1 mm Gaussian radius (e.g. Ar-ion laser, He-Ne laser and other). The wavelength is 0.6328 μm (He-Ne laser beam).

The approach for theoretical analysis with result for computer simulation is related to the

analytical treatment of the path of the beam rays in the structure (Fig.3) taking into account the complexity of the optical paths and correspondingly summing [2,3]. It is based on the decomposition of the complex amplitude $g(x)$ of the incident light beam on the plane waves in its plane z , which in paraxial approximation is described by the expression:

$$g(x, z) = \exp[ik(z - z_0)] \int_{-\infty}^{\infty} G_0(a) \times \exp(i2\pi\alpha x) da. \quad (1)$$

where $\exp(i2\pi\alpha)$ is a unit-amplitude plane wave propagating in direction given by the directional cosines $\lambda\alpha$ and $\lambda\gamma = \sqrt{1 - (\lambda\alpha)^2} \cong 1 - (\lambda\alpha)^2/2$, $G_0(\alpha)$ is the Fourier transform of $g(x)$ and represents the angular spectrum of the field $g(x, z)$ at the plane $z = z_0$. The propagation axis of the incident beam, which is also its axis of symmetry, coincides also with the propagation direction of the plane wave with a directional cosine $\lambda\alpha = 0$. The wedge response at an arbitrary point $P(x, z)$ to single component from the beam angular spectrum that falls on the wedge surface at angle $\theta_0 + \eta$ $\eta = \arcsin(\lambda\alpha)$ is

$$\tau(x, z, \eta) = T \sum_{p=1}^{\infty} R^{p-1} \Omega_p \cos \varepsilon_p \times \exp[i2(p-1)\pi] \exp[i\varphi_p(x, z, \eta)] \quad (2)$$

where $R = rr'$ with $r(t)$ and $r'(t')$ being the reflection coefficients of the front and rear wedge surfaces. Using this approach, it can be obtained [2] the transmitted intensity $I_T(x, z, \lambda)$ distribution for a Gaussian from:

$$I_T(x, z, \lambda) = a^2 T^2 (S_1^2 + S_2^2) \dots, \quad (3)$$

where

$$S_1 = \sum_{p=1}^{\infty} R^{p-1} \Omega_p \cos \varepsilon_p; \quad (4)$$

$$S_2 = \sum_{p=1}^{\infty} R^{p-1} \Omega_p \sin \varepsilon_p \quad (5)$$

with Ω_p and ε_p , which are functions of x, z, x_0, z_0 , wavelength, incident angle and the apex angle [3].

Obtained by the calculation graphs of the transmission (noted in the graphs as "Trans.") and reflection ("Refl.") at the resonance as function of the wedge thickness e are presented in Fig.3(a) for spatial finite beam (we accept for limitation 2-times the Gaussian radius - the last of 1 mm). In Fig.3(b)

are plotted the expanded part of Fig.3(a) around the exact resonant position. In this figure is also shown schematically the incident laser beams BP_1 and BP_2 in two positions of incidence points P_1 and P_2 .

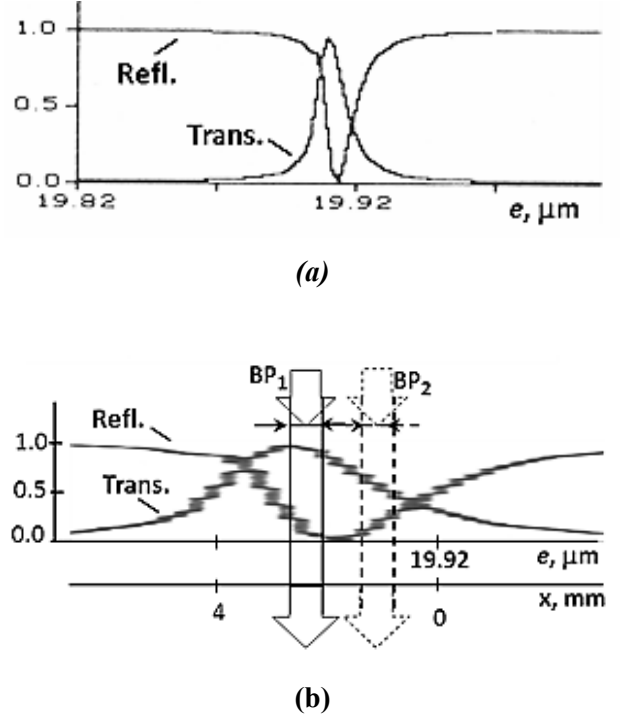


Fig.3. (a) Computer simulation of the transmitted (Trans.) and Reflected (Refl.) relative intensities for the function of the wedged thickness around the exact resonance maximum of the transmission at wavelength 632 nm; **(b)** – the expanded part of Fig.3(a) around the exact maximum of transmission. The transmitted and reflected relative intensities are given as function of thickness e – top scale and X – bottom scale, where X corresponds to the e – length along the wedge arm. In the graph in Fig.3(b) are shown schematically the incident beams BP_1 and BP_2 for two position of incidence.

3. Experiment

Now, as example, for the ~ 1 mm Gaussian diameter He-Ne laser beam ($\lambda = 0.6328 \mu\text{m}$, real diameter of the diaphragm of 1.7 mm, and passed power of $\sim 160 \mu\text{W}$, incident $200 \mu\text{W}$), we show the obtained experimental graphs of transmission power in arbitrary units as function of the distance X from the position of the center of the resonance.

The IW structure used is with $e = 8 \mu\text{m}$ thickness, mirror reflectivity $R = 0.85$ and index of the refraction $n = 1.5$. The experimentally measured transmitted power (consideration for $\lambda = 0.6328 \mu\text{m}$) as a function of distance X is given in Fig.4. As it is

expected, the continuous decreasing of the transmitted power moving away from the center is registered.

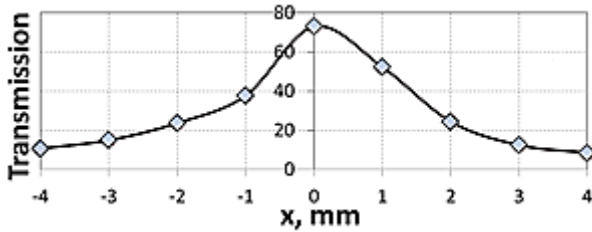


Fig.4. Experimentally measured transmission through the wedged interference structure (arbitrary units) as function of the distance X of the beam from the exact position of the resonance (for details - see the text).

It is of interest the intensity distribution in the transmitted spot. The typical experimentally observed distribution for the considered above case is shown in Fig.5 - with the solid line; the dropped line is the Gaussian distribution function. The investigation shows that in the case of the IWS, discussed above, the distribution in the intensity in transmitted beam repeats well the incident beam.

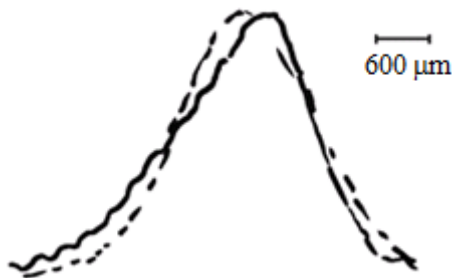


Fig.5. Intensity distribution in the transmitted (solid line) and in the incident beam (dashed line).

In general, for incident angles to 10-30 angular degrees (depending on the IW thickness), the distribution is with acceptable conservation. As example, the distortion is essential for the high wedge thickness of order of 200 μm and incidence angle higher than 15 angular degrees (the case for glass or quartz wedge separator), the change of the distribution is not acceptable for many applications. In this case the Interference Wedged Structure can be also applied as simplest controlled transmission filter.

4. Conclusion

In the work we have developed, as a pioneering study, new possibility for employment

of Interference Wedged Structure (IWS) as competitive for many cases beam splitting and power separation element. The essential moment is the simplicity of the control. Also, if the structure is prepared by materials with high illuminating power density resistance, the application is possible for beams with very high powers ($\sim \text{MW}/\text{cm}^2$ and GW/cm^2). The given theoretical estimations compared with the experimental study confirm the expected advantages - precisely and variably controlled ratio of the reflected and transmitted power, obtained by simple translation of the list-like wedged structure in its plan; separation is practically without energy losses; there are not changes of the geometry of the beam propagation - in transmission and in reflection during the ratio variation. The last is essential difference in comparison with application of other techniques - based on FP interferometer, due to the variation of the transmitted and reflected beams during the control by variation of the incident angle.

The presented in the work development and results concern the one angle IWS with thickness of order of 2 - 20 μm and incident angles for the illuminating beam less than 40 angular degrees. As limitation, can be noted that in general, good reproducibility of the shape of incident beam can be guaranteed for the illuminating beam incident angles to the 10-40 angular degrees (depending on the IW thickness). The distortion is essential for the high wedge thickness of order of 200 μm and incident angle higher than 15 angular degrees (the case for glass or quartz wedge separator), the change of the distribution is not acceptable for many applications. Also in these cases, for separation of needed power part of the incident beam, the Interference Wedged Structure can be applied as simplest controlled transmission filter by its sliding in its plane along the length of the angle arms.

ACKNOWLEDGEMENTS

The work is supported by National Scientific Found - Bulgaria, contr. DN 17/7 (2017).

REFERENCES

1. Born M. and E. Wolf (1999). *Principles of Optics*. Cambridge University Press.
2. Stoykova E. and Nenchev M. (2010). Gaussian Beam Interaction with Air-gap Fizeau Wedge. *J. Opt. Soc. of America*, 27, 58-68.

3. Nenchev M. and Stoykova E. (2001). Fizeau wedge with unequal mirrors for spectral control and coupling in a linear laser oscillator-amplifier system. *Appl. Optics*, 40 (27), 5402-5411 and the literature therein
4. Nenchev M., Meyer Y.H. (1981). Two-wavelength dye-laser operation using a reflecting Fizeau interferometer. *Applied Physics*, 24, 7-9.
5. Deneva M., Uzunova P., Nenchev M. (2007). Tunable subnanosecond laser pulse generation using an active mirror concept. *Optical and Quantum Electronics* 39, 193-212.
6. Nenchev M., Deneva M., Stoykova E. (2017). Development of composite wavelength tunable interference wedged structures for laser technology, spectroscopy and optical communications. *Photonica-2017, Intern. Confer. Book*, p.141, Belgrade, Bulg. Patent, BG Patent reg. No 110967/ed.2017/, Wavelength division /multiplexing devices
7. Kajava, T. Lauranto T., and Salomaa, R. (1993). Fizeau interferometer in spectral measurements. *J. Opt. Soc. of America*, 27, 58-68.
8. Nenchev M., Deneva M, Stoykova E., Project DN 17/7, NSF – Bulgaria (2017) ; Patent BG, Reg. No 110967 / ed.2017/, Wavelength division /multiplexing devices

For communications:

Marin Nenchev

E-mail: marnenchev@yahoo.com

Margarita Deneva

E-mail: mar.deneva@abv.bg

CONTACTS

Margarita Deneva, Valko Kazakov, Vania Plachkova, Kamen Ivanov, Marin Nenchev

*Technical University of Sofia, R&D Department, "Quantum and Optoelectronics" Scientific Laboratory (QOEL) and Plovdiv Branch, Department of Optoelectronics and Laser Engineering, FEA
Plovdiv 4000, 25 "Tsanko Diustabanov" Str.
Bulgaria*

Pepa Uzunova

*Medical University of Sofia,
Sofia 1431, 15 "Acad. Ivan Geshov" Blvd.
Bulgaria*

Elena Stoykova

*Institute of Optical Materials and Technologies, Bulgarian Academy of Sciences,
Sofia 1113, 109, "Acad. G. Bontchev" Str.
Bulgaria*

CALCULATION OF DIELECTRIC DISSIPATION FACTOR AT VARIABLE FREQUENCIES OF MODEL TRANSFORMER

ONUR KAYA, T. CETIN AKINCI, EMEL ONAL

Abstract: *The condition of the insulation is essential for secure and reliable operation of transformer. Measuring capacitance and dissipation/power factor helps us to determine insulation condition in bushings or between windings. Changes in capacitance can indicate mechanical displacements of windings or partial breakdown in bushings. Aging and degradation of the insulation, coupled with the ingress of water, increase the amount of energy that is converted to heat in the insulation. The rate of these losses is measured as dissipation factor. With testing systems, you can even determine the capacitance and dissipation/power factor at variable frequency. Therefore, aging phenomena can be detected earlier, and corresponding action such as repair, oil treatment or drying can be initiated. At this study, the dissipation factors of the model transformer is analyzed at variable different frequency range. Especially the analyses at low frequencies range of transformers are very important nowadays in terms of energy saving.*

1. Introduction

The power transformer is one of the most important parts of any electrical transmission and distribution service. Although the transformer is a stationary device and its design is quite simple, the maintenance of the power transformer is rather troublesome [1]. Despite major advances in power equipment design in recent years, the weak link in the chain is still an insulation system. A greasy paper insulation system in a power transformer can deteriorate under electrical, thermal and environmental stress even under normal operating conditions [2]. Unexpected failures result in major degradations in operating systems, resulting in unplanned outages and power distribution problems [3]. Many power transformers are old nowadays. It is a very expensive solution to replace transformers with new ones just because they are old. Because most of these transformers can work for many years. Since the life of a transformer is directly related to the insulation quality, monitoring of the insulation condition of the transformers is an important issue [4]. Accurately assessing and monitoring the condition of the oil-paper insulation system is necessary to understand how long the life of the transformer remains and not only increases the reliability of the power source, but also reduces maintenance costs [5]. At this study, the model transformer is analyzed at different frequencies range.

2. Measurement and Analysis Principles

Several recent techniques have been developed for monitoring transformers in recent years [6]. Some of these techniques are Recovery Voltage Measurement (RVM) or Polarization Depolarization Current Measure (PDC). Dielectric analysis can also be done in the frequency domain as in our work, as in Frequency Domain Spectroscopy (FDS) technique [7]. Frequency Domain Spectroscopy (FDS) diagnostic techniques have recently become more popular than other techniques. One of the reasons is that the measurement of the loss factor ($\tan \delta$) is independent of the transformer geometry (shape). Another advantage of FDS is that it is less noisy than other insulation measurement methods (PDC, RVM) [8]. Numerous studies on the application of FDS measurements have shown that measurements between high and low voltage are least influenced by the shape, weather, and external factors [6]. Frequency Domain Spectroscopy (FDS) can measure loss factors and capacity at low voltages and at all frequencies [9]. That is, it shows the general aging condition and moisture content of oil-paper insulation of the FDS transformer [10]. It is necessary to know the equivalent circuit of the capacitance to understand these analyzes well. Because an insulator can be represented a capacitance as mathematical model. Figure 1 shows a capacitance with parallel resistance as an equivalent circuit of insulator.

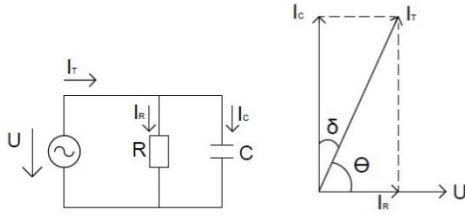


Fig. 1. . Equivalent capacitance circuit and vector diagram

As seen in Figure 1, the total current I_T , flows from the insulator. This current has two components, I_R (ohmic) and I_C (capacitive). The angle between the capacitive current I_C and the total I_T current is δ and the angle between the ohmic current I_R and the total I_T current is θ . These angle values give the information about the level of insulation. Here the loss factor is $\tan \delta = I_R / I_C$. The power factor is $\cos \theta = I_R / I_T$ [11]. Insulation of the transformer cannot be detected when measured with a constant frequency. A complete diagnosis can be made thanks to the measured loss factor (FDS) [11]. Since the aging of the transformer insulation is a slow-moving process, the analysis of the transformer can be applied using the loss factor ($\tan \delta$). The dielectric loss factor $\tan \delta$ is one of the key parameters for evaluating the high-voltage insulation condition. Because this method can remove the effects caused by random errors. However, the dielectric loss factor can be affected by various factors such as strong electromagnetic disturbance due to normal operation of devices and environmental temperature and humidity [12].

In this study, the insulation condition is observed by using (Insulation Diagnostics Systems) IDA 200 as shown in figure 2. IDA 200 is designed for diagnostic measurements of electrical insulation.

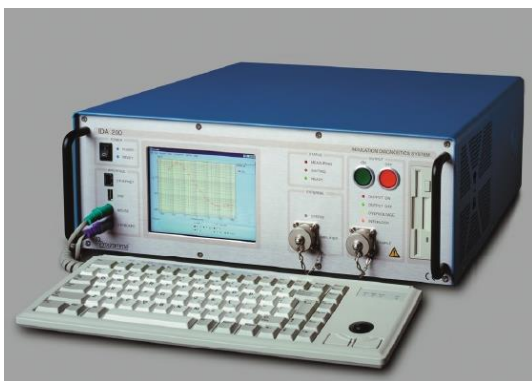


Fig. 2. Insulation diagnosis system equipment

By using IDA 200 it is possible to identify insulation materials in most high voltage installations (e.g power transformers, measuring transformers, bushings, paper- insulated cables).The diagnostic measurement is made by applying a relatively low voltage up to 140 V. The IDA 200 measures capacitance and dielectric

losses at different frequencies both above and below the frequency. By using filtering does not allow harmonics. Electrical insulation occurs in all three forms: solid (such as cellulosic paper and porcelain), liquid (such as mineral oil), and gas. Insulation systems with solid and liquid insulation are suitable for IDA 200 measurements. However, IDA 200 is not suitable for pure gas insulated systems. Most power transformer's insulation system consists of oil and cellulose. Both materials change their dielectric properties throughout the lifetime of the transformer. When only a constant frequency is measured, the property changes in different materials are indistinguishable. Analysis of the measured dissipation factor frequency characteristic allows the inspected insulation to be diagnosed more accurately. The system measures the impedance of a sample with variable voltage and frequency. A Digital Signal Processing (DSP) unit produces a test signal of a desired frequency. This signal is amplified by an internal amplifier and then applied to the sample. The voltage and current are measured using a voltage divider and an electrometer. For measurement input, the IDA 200 uses the DSP unit which multiplies the input (measurement) signals by the reference sine voltages and then combines the results on a series of loops. By using this method, noise are almost completely filtered. Thus, the IDA 200 operates at low voltage levels with high efficiency [13].

At this study, the model transformer is analysed by using IDA 200. The model transformer is called "Pancake Model Transformer". The model consists of eight shaped coils with ducts between them. The ratio of barriers and spacers to oil ranged from 15 to 100 % as described in the following table 1. This simulates the main insulation of different transformers.

Table 1. The ratio of oil and spacers of pancake transformer

Connection	Oil / Barriers	Oil / Spacers
CH-B	83 / 17	85 / 15
DG-CH	72 / 28	72 / 28
E-DG	50 / 50	45 / 55
F-E	0 / 100	0 / 100

3. The Results of Measurements

When measuring the dielectric loss factor $\tan \delta$, the transformer tank and its windings are treated as a natural capacity. In our two-winding transformer, as in our example, there are three capacitors, between the windings and the tank and also between the windings [14]. The two windings transformer with 25 MVA three phase and 110/20 kV used in this study is shown at figure 3.

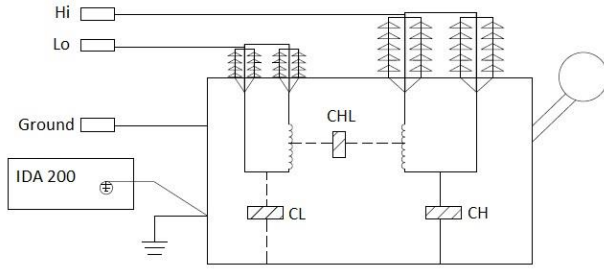


Fig. 3. Simple display of capacitances of two-winding transformers

Capacitance values are the values between the low voltage windings and the transformer tank and the high voltage windings and the transformer tank. Here;

C_L : Capacity value between low voltage windings and transformer tank

C_H : Capacity value between high voltage windings and transformer tank

C_{HL} : Capacity value between transformer high and low voltage windings

In the following graphs, the power factor, loss factor $\tan \delta$ and capacity curves of some of the CH-B connections of the pancake transformer for C_{HL} , C_H and C_L are shown in figure 4, 5, 6 respectively.

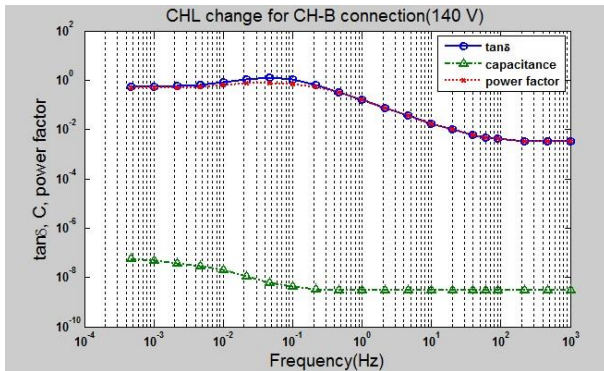


Fig. 4. C_{HL} change for CH-B connection (140 V)

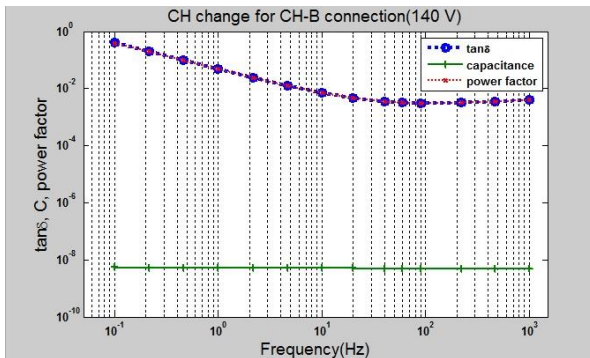


Fig. 5. C_H change for CH-B connection (140 V)

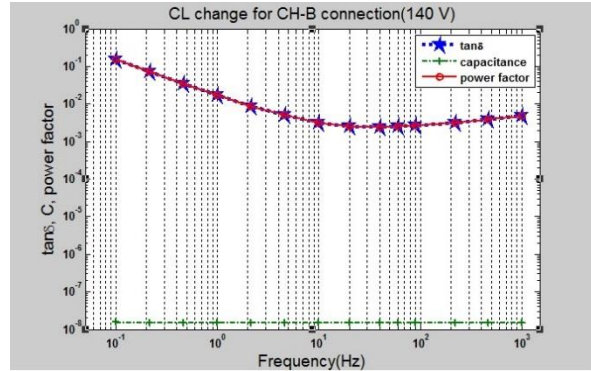


Fig. 6. C_L change for CH-B connection (140 V)

The high voltage bushings are critical components of the power transformer and particularly, capacitive high voltage bushings need care and regular tests to avoid sudden failures. These bushings have a measurement tap-point at their base and both the capacitance between this tap and the inner conductor (normally called C1) and the capacitance between the tap and ground (normally called C2) are measured. An increase of C1 indicates partial breakdowns of the internal layers. To determine bushing losses, dissipation factor tests are performed. H1, H2 and H3 are default labels for bushings of different phases.

H1C1: measurement of bushing H1 main insulation, C1, H2C1: measurement of bushing H2 main insulation, C1, H3C1: measurement of bushing H3 main insulation, C1, H1C2: measurement of bushing H1 insulation between test tap and ground sleeve, C2, H2C2: measurement of bushing H2 insulation between test tap and ground sleeve, C2, H3C2: measurement of bushing H3 insulation between test tap and ground sleeve, C2.

According to measurement results, dissipation factor, and power factor and capacitance C_1 variation are obtained for CH-B connection of pancake transformer can be shown in figure 7. C_1 capacitance value is evaluated as the main insulation of transformer. As seen from figures of 4, 5 and 6, dissipation factor and power factor values are high at low frequencies. As the frequency values increase, dissipation factor, and power factor also decrease. In figure 5 and 6, the capacitance values are not changed, but in figure 4, it is slightly higher than the values at low frequencies. The dissipation factor and capacitance of C_H , C_L , C_{HL} , C_1 and C_2 are shown for the frequency of 60 Hz in table 2. The dissipation factor and power factor of C_{HL} are higher than that of C_L and C_H for the frequencies of low range.

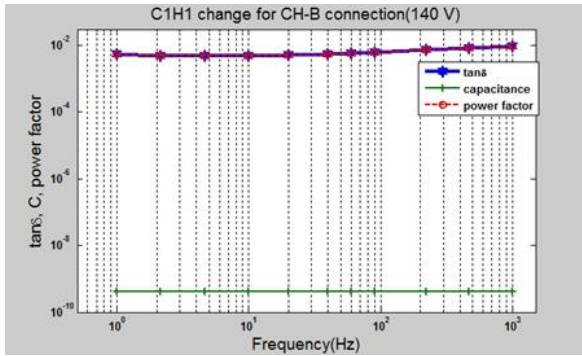


Fig. 7. C_{1H1} change for CH-B connection (140 V)

Table 2. $\tan \delta$ and capacitance values of CH-B connections of model transformer for the frequency of 60 Hz

Sweep	Tan δ	Capacitance(pF)
CH	0.3428	4981
CHL	0.2589	6984
CL	0.248	1.512E4
H1C1	0.5637	411.6
H2C1	0.5555	418.7
H3C1	0.5623	411.9
H1C2	0.3604	5388
H2C2	0.3604	5388
H3C2	0.3716	5441

The pancake transformers have different oil spacers ratio as shown in table 1. CH-B connections has the largest oil ratio while F-E connections has 0 amount of oil. The additional analysis for F-E connections with the aim of comparison is seen in figure 8.

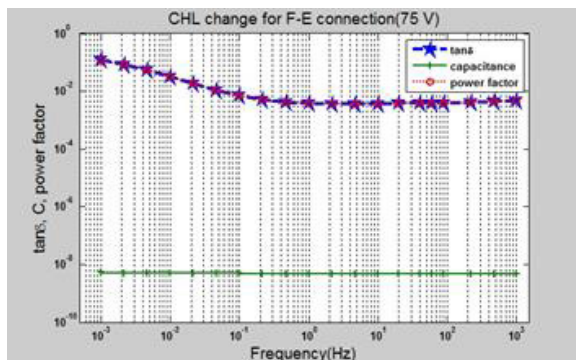


Fig. 8. C_{HL} change for F-E connection (75 V)

The variations of the dissipation factor and power factor of F-E connections are smaller than that of CH-B connections. This case can be explained by the fact that there are fewer losses due to lack of oil.

4. Conclusion

In this study, the dissipation, loss factor and capacitance values of C_{HL} , C_H and C_L , C_{1H1} are given for CH-B type connection of the pancake model transformer. In the loss factor and power factor curves of C_H , C_L and C_{1H1} are found to be very close to each other at the high frequency range. However, these values at low frequencies are very different each other. So the dissipation factor, power factor and capacitance values are more determinative for us. The dissipation factor of C_{HL} for a frequency of 0.1 Hz in CH-B mode is approximately 2.5 times greater than that of C_H and approximately 7 times greater than that of C_L . On the frequency of 60 Hz, the dissipation factor of C_{HL} is approximately 1.5 times higher than that of C_H and approximately 2 times higher than that of C_L . At the frequency of 1000 Hz, the rates are quite variable. At the low frequencies, the dissipation factor value is the highest for C_{HL} , while at the lower frequency, the dissipation factor is the lowest. While the dissipation factor of C_L and C_H are close to each other, the dissipation factor of C_{HL} are about 1.5 times smaller than that of C_L and C_H . Today, transmission and production at low frequencies is preferred due to the fact that losses are low. In particular, since diagnostic tests are performed at low frequencies such as 0.1 Hz, a more detailed analysis of low frequency experiments is required. The capacity and dissipation factor changes give us information about the type of error. For example, short circuit turns and changing of magnetizing are seen at the low frequency response while the medium frequency response is sensitive to hoop buckling and axial movement faults in the transformers. The high frequency response is sensitive to the part of the properties of windings and grounding problems. For this reason, the capacitance and dissipation analysis according to frequency are very important for detecting faults in transformers.

REFERENCES

- [1] Y. Hadjadj, I. Fofana, F. Meghnefi, H. Ezzaidi, "Assessing Oil Paper Insulation Conditions by Poles computed from Frequency Domain Spectroscopy" IEEE International Conference on Dielectric Liquids (ICDL), 2011 Norway.
- [2] Tang, M.; Lei, M.; Xu, H.; Lu, Y.; Wei, X.; Wu, P.; Zhang, G. Frequency domain characteristics and insulation condition evaluation of power transformer oil-paper insulation. In Proceedings of the 2012 International Conference on High Voltage Engineering and Application (ICHVE), Shanghai, China, 17–20 September 2012; pp. 443–446.
- [3] Betie, A.; Meghnefi, F.; Fofana, I.; Yéo, Z. On the Impacts of Ageing and Moisture

- on Dielectric Response of Oil Impregnated Paper Insulation Systems. In Proceedings of the IEEE Annual Conference on Electrical Insulation and Dielectric Phenomena (CEIDP), Montreal, Canada, 14–17 October 2012; pp. 219–222.
4. [4] I. Fofana, H. Hemmatjou, F. Meghnefi, M. Farzaneh, A. Setayeshmehr, H. Borsi and E. Gockenbach] “Low Temperature and Moisture Effects on Oil-Paper Insulation Dielectric Response in Frequency Domain” IEEE Electrical Insulation Conference EIC, 2009. 2009.
 5. [5] De Nigris, M.; Passaglia, R.; Berti, R.; Bergonzi, L.; Maggi, R. Application of Modern Techniques for the Condition Assessment of Power Transformers; International Council on Large Electric Systems: Paris, France, 2004.
 6. [6] C. Ekanayake ; T. K. Saha ; H. Ma ; D. Allan “Application of Polarization Based Measurement Techniques for Diagnosis of Field Transformers” IEEE Power and Energy Society General Meeting, 2010 USA.
 7. [7] [Impact of Temperature on the Frequency Domain Dielectric Spectroscopy for the Diagnosis of Power Transformer Insulation J H Yew, Member, T K Saha, Senior Member, A J Thomas, Student Member IEEE]
 8. [8] [Determination of Water Content in Transformer Solid Insulation by Frequency Domain Spectroscopy BELÉN GARCÍA, BAUDILIO VALECILLOS, JUAN CARLOS BURGOS Electrical Engineering Department University Carlos III, Madrid C/ Butarque 15, Leganés Madrid SPAIN]
 9. [9][Application of Low Frequency Dielectric Spectroscopy to Estimate Condition of Mineral Oil-A.A. Shayegani, H. Borsi, E. Gockenbach and H. Mohseni]
 10. [10] [On the Frequency Domain Dielectric Response of Oil-paper Insulation at Low Temperatures -I. Fofana, H. Hemmatjou, F. Meghnefi, M. Farzaneh, A. Setayeshmehr, H. Borsi and E. Gockenbach]
 11. [11] Insulation Diagnostics Spectrometer IDA 200: User’s Manual; Programma Electric AB: Täby, Sweden, 2000.
 12. [12] [Analytical Processing of On-line Monitored Dissipation Factor Based on Morphological Filter-N. Wang, F. C. Lu and H. M. Li]
 13. [13] [Güç Trafolarında Arıza Tespitine Yönelik Gelişmiş Tanı Testleri- Nihat PAMUK- TEİAŞ 5. İletim Tesis ve İşletme Grup Müdürlüğü Test Grup Başmühendisliği, 54100, Sakarya]
 14. [14][Güç Transformatörlerinde Güç Faktörü Testleri Ve Enerji Kalitesine Etkileri-Hakan Çuhadaroğlu, Yılmaz Uyaroğlu, Gökşin Sönmez, Haluk Yılmaz]
 15. [15][A Comparative Test and Consequent Improvements on Dielectric Response Methods M. Kochl*, S. Tenbohlen, M. Kruger, A. Kraetge University of Stuttgart, IEH, Pfaffenwaldring 47, 70569 Stuttgart, Germany 2 Omicron electronics GmbH, Oberes Ried 1, 6833 Klaus, Austria]

Authors’ contacts

Istanbul Technical University
Cognitive Systems Lab.
Istanbul Turkey
kayao16@itu.edu.tr
akincitc@itu.edu.tr
eonal@itu.edu.tr

TWO NEW TYPES OF COMPACT ULTRA-WIDE BANDWIDTH ANTENNAS FOR 3-AXIAL RF FIELD STRENGTH MEASUREMENT.

SARANG M PATIL

Abstract: *In later years, embedded systems and smart sensors have been of increasing interest for a machine to machine (M2M) communication and automation of services. Embedded systems also improve the standard of living by offering a higher level of functionality to the home, workplace and to the environment we live in. To manage the situation with an increasing density of mobile electronic devices where many are equipped with wireless communication technology, EMC issues must be considered at the planning stage. The EMC regulations are an essential step since they are helping people to realize that it is necessary to address electromagnetic compatibility problems. The problems have to be solved before products involving electronic circuits can be put on the market. Using electromagnetic compatibility test probe to identify situations where EMC problems may occur, the EMI threats can be reduced by EMC barriers in the design of electronic devices. Conventional field strength sensors use different detection Methods, each having advantages and disadvantages. Electromagnetic fields are conventionally measured using diode detectors or thermocouple detectors. The diode is limited in dynamic range. We present two new types of ultra-wideband antennas for EMC measurements: the calculable Rod-dipole antennas and the Hexagonal shape fractal dipole antenna. Three orthogonal dipole antennas connected with the hybrid coupler; result in a fast, broadband, and high dynamic range field strength probe. Both antennas have a compact size and Broad bandwidth for flexibility in EMC measurement tests. The antennas have a low manufacture cost and lightweight, easy for installation, the concept described, and simulation results are shown.*

Keywords: *Ansys HFSS, Broad Bandwidth Antenna, Electromagnetic Compatibility, Fractal Antenna, pre-compliance test, 3-axis field probe.*

1. Introduction

The necessity of electromagnetic compatibility (EMC) is not that obvious for people with no experience from electronic design. It is, however, necessary for modern electronic devices to be designed with electromagnetic compatibility issues in mind to fulfill all safety and protection requirements and offer the functionality and quality expected from the device. An apparatus both emits and receives electromagnetic energy conducted on attached cables or radiated from enclosures and cables. That is just a side effect when an electronic device is offering some functionality and cannot be avoided entirely.

Because of this undesired side effect there are legislated demands on the device regarding its electromagnetic properties. Which means that a manufacturer of an apparatus covered by the directive must design the equipment and specify how it should be used so that it does not cause

electromagnetic disturbance harming other devices that do comply with the directive? Also, the apparatus must be designed to be immune to an average level of electromagnetic energy in the environment it is intended to be used.

Most electronic design engineers well understand that it is essential to meet requirements on EMC. However, the lack of easy to use design tools to address EMC problems often makes EMC fixing necessary at a late stage in the design process. Better knowledge of EMC barriers as components and design blocks can simplify EMC considerations at an early design stage.

Many EMC antennas are available in the market, such as ETS-Lindgren EMC antennas [1], active LogPer directional antennas [2], etc.

Electric field strength is conventionally measured using thermocouple elements or diode detectors [3], [4]. The thermocouple is used for determining the real root mean square (RMS) of a signal, but it is too slow to measure the envelope

and peak of fast-changing fields, such as the ones described earlier. The diode is limited in dynamic range because the diode-based detection is square-law with a small signal but becomes linear with a significant signal until it is overloaded and damaged.

To measure all three components of electric field vector, a tailor-made antenna type called "Tripole" is most beneficial over conventional antenna; different EM wave field is measuring antenna with the comparison is discussed in previous work [5].

In this paper, we present two ultra-wideband antennas as alternatives to EMC antennas in the market. Both antennas have a compact size, excellent reflection coefficient, stable gain therefore linearly increased antenna factor with frequencies. Besides, these two antennas can have different port setups, offering flexibility to varying tests in EMC measurements. The manufacturing cost of the antennas is also low.

2. Three-Axial Rod Dipole Structure

Three-axis antenna consists of three orthogonal dipole antenna oriented as shown in figure 1. The response of each dipole is more strongly to linear polarized signal that matches the antenna elements. Each element output of this array is proportional to the vector field components E_x , E_y , and E_z .

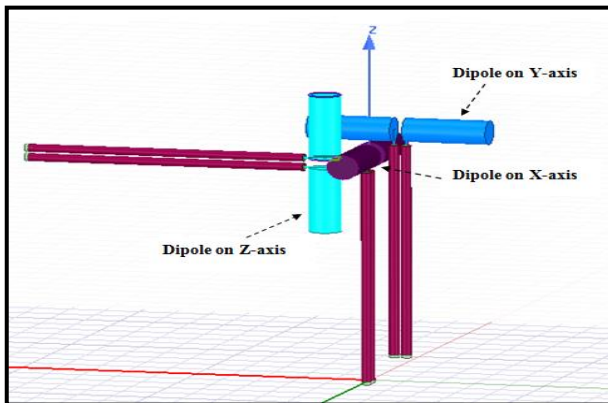


Fig. 1. Simulated Return loss for a dipole antenna.

The antenna consists of an array of three orthogonally arranged dipole elements, each with its Feed line along with hybrid balun. Wide-band, calculable dipole antenna along with balun design and dimensions, has explained in previous work [6], the single dipole antenna is designed for reference frequency 1.3GHz to cover the range of frequencies from 900MHz to 3.2GHz. Unique antenna geometry along with coaxial cable with specified length placed orthogonally to each other, like one dipole on X-axis, second on y-Axis and third on Z-axis

respectively. The result of this geometry explained in the following section.

2.1 Return Loss S11 of Single dipole

For the wide-band coverage from 700MHz to 3GHz, the reference frequency is set to 1.3GHz and accordingly a dipole dimension is calculated using a standard formula of the dipole. So dipole total length L is 10.3cm including the gap between two is $GL = 7\text{mm}$. For this sizes, the antenna is simulated using hfss v18.0.

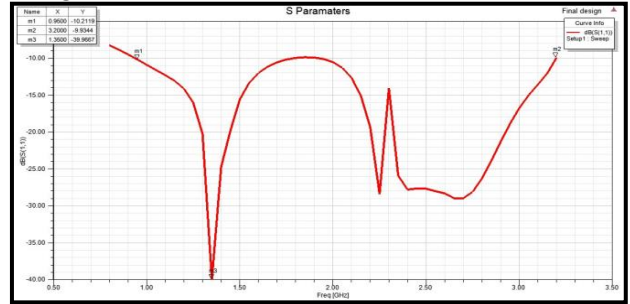


Fig.2. Simulated Return loss for a dipole antenna.

See Figure 2 is returned loss of calculable reference antenna with the balun, As per above graph, marker one start from 900MHz and end at a 3.2GHz frequency and for this 0.9-3.2GHz frequency range return loss S11 value less than -10dB. Maximum peak is detected at frequency 1.3GHz, and its value is -39.96dB.

2.2 Return Loss S11 of three orthogonal dipole

In electric field strength measuring case, for high accuracy, three crosses orthogonal dipole are arranged. Return loss s11 of XY plane is shown in figure 3; it shows that return loss is below -10dB for the range of frequency approximately from 900MHz to 3GHz, similarly return loss S11 of xz and yz plane are shown in figure 4 and 5 respectively.

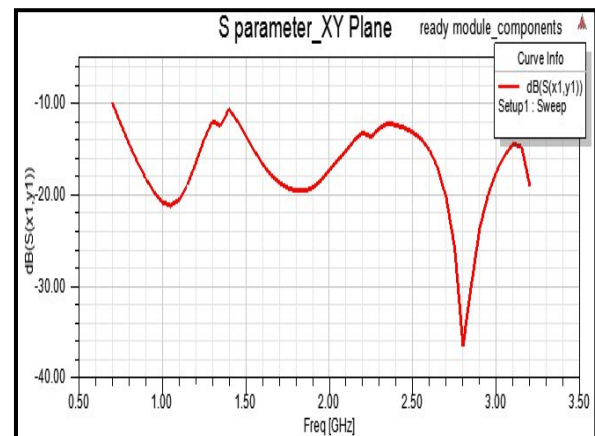


Fig.3. Simulated Return loss S11 in XY Plane

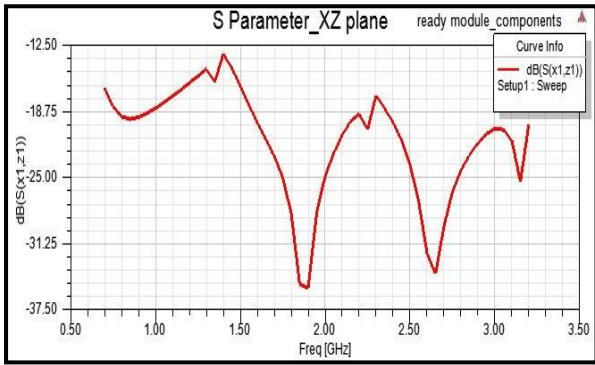


Fig.4. Simulated Return loss S11 in XZ plane

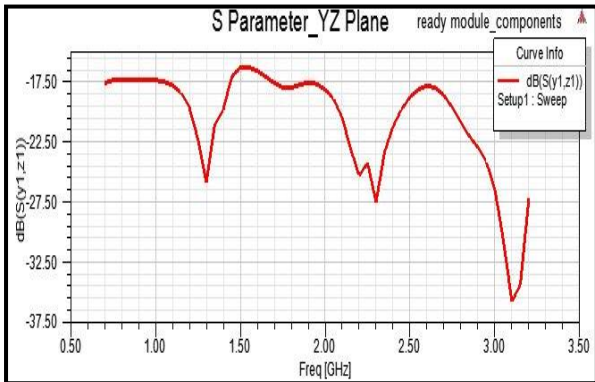


Fig.5. Simulated Return loss S11 in YZ plane

2.3 Polar radiation pattern

Polar radiation pattern at frequency 2GHz shown in figure 6; antenna provides a perfect pattern with maximum gain is 2.3 dB,

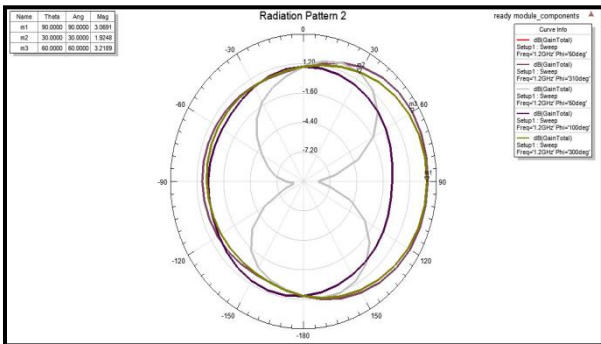


Fig.6. Polar radiation pattern of antenna at 2GHz.

2.4 3D Radiation and Directivity

Seen in figure 7, the Smith chart of the antenna, characteristics impedance is most crucial aspect in antenna designing, theoretically power and cable impedance was considered as 50Ω, a per the simulation results with Smith chart characteristics impedance is calculate, impedance value is approximately equal to theoretical consideration, as per the result characteristics impedance is 73.53 at 1.35 GHz frequency.

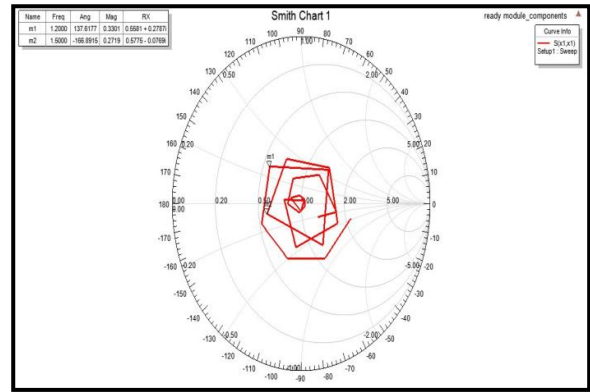


Fig.7. Smith chart of dipole antenna

2.5 3D Radiation and Directivity

3D radiation pattern and 3D directivity is shown in figure 8 and 9 respectively

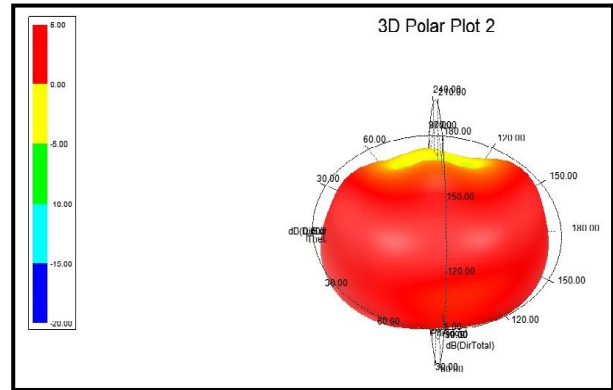


Fig.8. 3D Radiation pattern of an antenna.

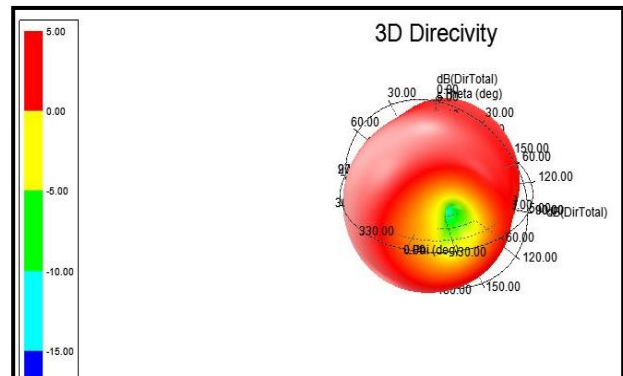


Fig.9. 3D Directivity plot.

3. Three-Axis Fractal Dipole Structure

Three-axis antenna consists of three orthogonal cross-coupled dipole-arms placed in X, Y and Z axis. The antenna orientation as shown in figure 10.

The antenna array consists of the three orthogonally placed dipole elements, each with its Feed line. Design of hexagonal fractal dipole antenna along with feeding techniques has explained in [7] previous work. A single dipoleantenna is designed for reference frequency

3GHz to cover the range of frequencies from 2 GHz to 12GHz. Unique antenna geometry along with coaxial connector with above-specified length placed orthogonally to each other, like one dipole on X-axis, second on a y-axis and third on Z-axis respectively. The result of this geometry explained in the following section.

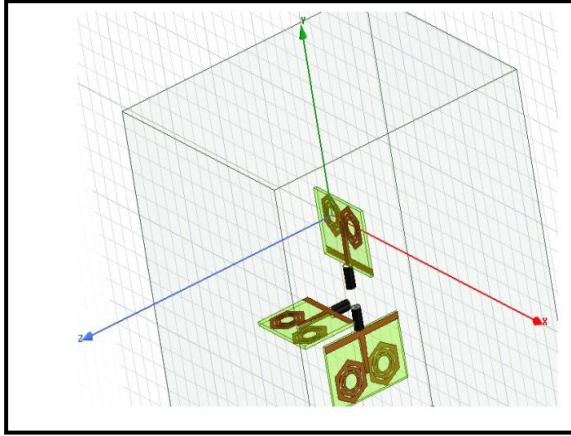


Fig. 10. Proposed HFSS module of 3-Axis Orthogonal fractal dipole.

3.1 Return Loss S11 of Single dipole

For the Ultra-wideband coverage from 2GHz to 12GHz, the reference frequency is set to 3GHz, and accordingly, dipole dimensions are designed using a hexagonal shaped fractal dipole. So dipole whole substrate dimension is 44mm*33mm including a gap between two arms.

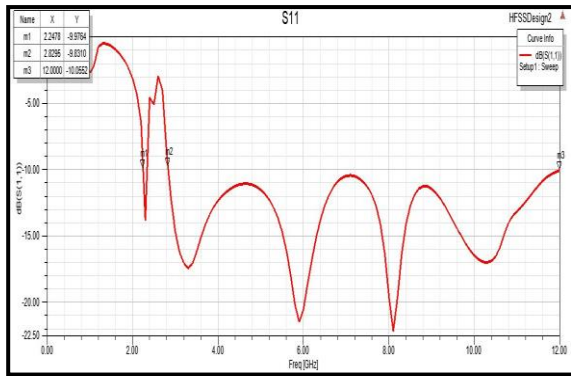


Fig.11. Simulated Return loss for dipole antenna

Figure 11 shows the return loss (S11) of the hexagonal fractal dipole antenna, As per above fig. 11, marker m1 start from 2.8GHz and end at the 12GHz frequency and for this frequency range return loss S11 value is less than -10dB. Maximum peak is detected at frequency 8.2 GHz, and its value is -22.dB.This design provides excellent agreement between mathematical model and simulator design.

3.2 Return Loss S11 of three orthogonal Dipole.

In electric field strength measuring case, for high accuracy measurement of field strength, three crosses orthogonal dipole are arranged. Return loss S11 shown in figure 12; it shows that return loss is below -10dB for the range of frequency approximately from 0.5 GHz to 12GHz with minimum value is -29.19. The graph of S21 and S22 is shown in figure 13 and 14

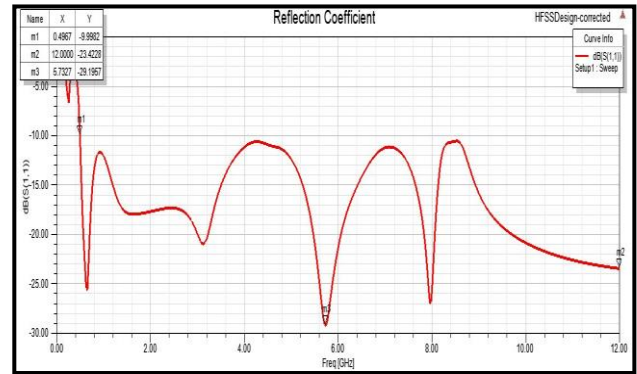


Fig.12. Simulated Return loss S11.

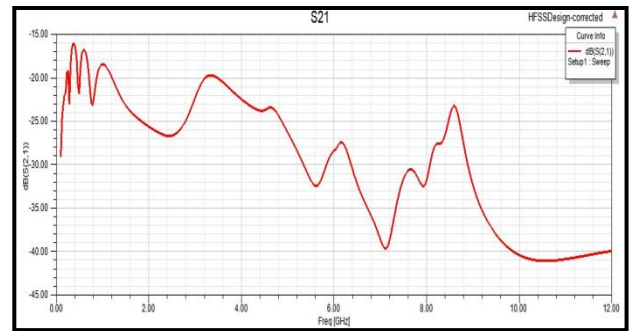


Fig.13. Simulated Return loss S21.

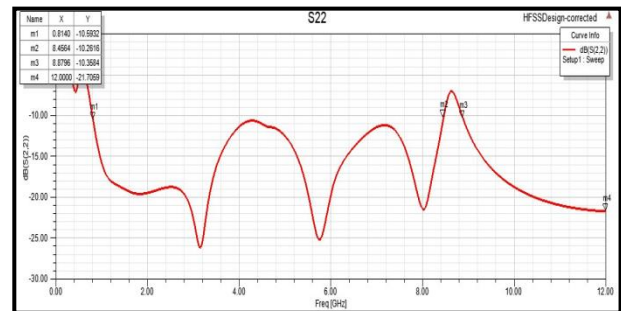


Fig.14. Simulated Return loss S22.

3.3.VSWR plot

VSWR plot shown in figure 15; it shows that antenna is a good radiator for the range of frequency approximately from 2.8 GHz to 8.83GHz with minimum value is 1 and second band is observed from 9.5-11.17GHz with a minimum amount less than 1.

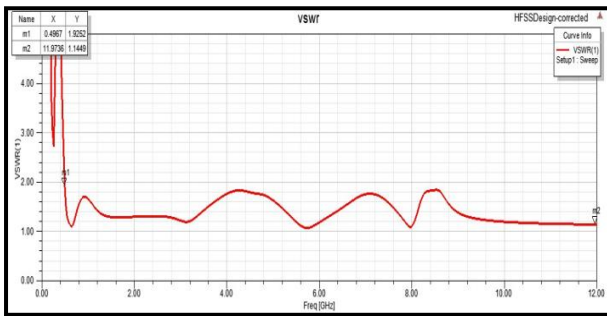


Fig.15. Simulated VSWR plot.

3.4.Polar radiation pattern of antenna at 3GHz.

Polar radiation pattern at frequency 3GHz shown in figure 16; the antenna provides the perfect design with maximum gain is 19.0 dB,

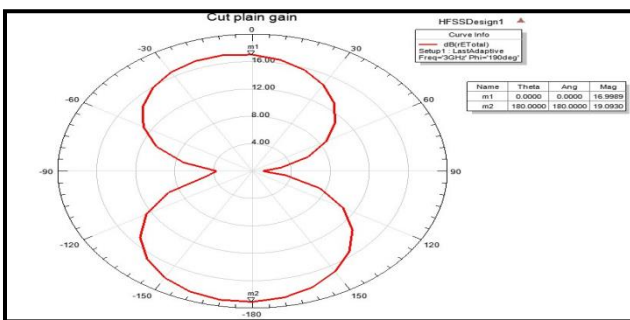


Fig. 16 The polar radiation pattern of an antenna at 3GHz.

3.5. 3D Radiation pattern and Directivity

3D radiation pattern and 3D directivity are shown in figure 17 and 18 respectively.

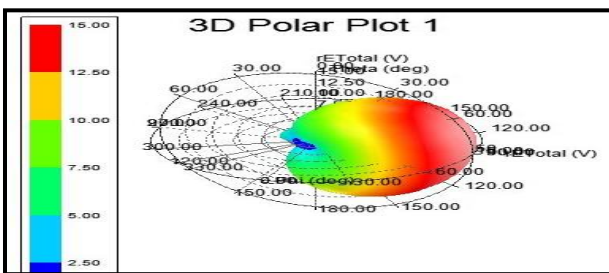


Fig.17. 3D Radiation Pattern of an Antenna.

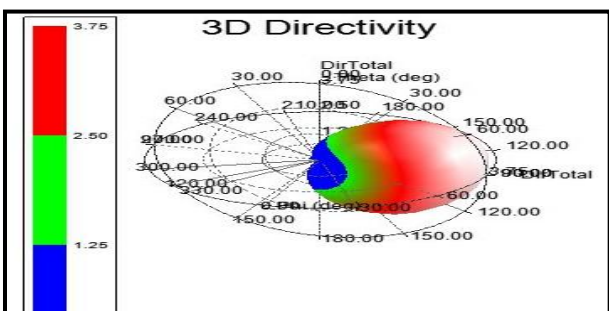


Fig.18. 3D Directivity Plot.

4. Conclusion

We present two new ultra-wideband antennas as alternatives to EMC antennas in the market. Both antennas have some characteristics, such as stable antenna gain and fixed phase center location, full bandwidth which may be interpreting for some EMC measurements.

Acknowledgment

The author would like to thank Prof Peter Petkov for helpful guidance to design the antenna module and important discussion, comments during experimentation, author also like to acknowledge the financial support given by Technical University of Sofia, Bulgaria. "This material is based upon work supported by the Bulgarian Science Fund under Grant No DN07/19/15.12.2016"

REFERENCES

1. ETS-Lindgren antennas: available at http://www.ets-lindgren.com/emc_antennas.
2. LogPer antennas: available at <http://www.aaronia.com/products/antennas/>.
3. H. I. Bassen and G. S. Smith, "Electric field probes—A review," *IEEE Trans. Antennas Propag.*, vol. AP-31, no. 5, pp. 710–718, Sep. 1983.
4. M. Kanda, "Standard probes for electromagnetic field measurements," *IEEE Trans. Antennas Propag.*, vol. AP-41, no. 10, pp. 1349–1364, Oct.1993.
5. Mr.sarang M. Patil, Prof.Peter Z Petkov, and Prof.boncho G Bonev, "A review on recent antenna designing techniques For electromagnetic compatibility(EMC) test, "International Scientific Conference on Engineering, Technologies and Systems TECHSYS 2017, Technical University – Sofia, Plovdiv branch, May 2017, ISSN Online: 2535-0048
6. Mr.Sarang M. Patil, Prof.Peter Z. Petkov, Prof.Boncho G. Bonev and Hitesh Singh," The Design and Modeling of Wideband Calculable Standard Rod-Dipole reference antenna for Broad Bandwidth Applications ,"IEEE International Conference on Infocom Technologies and Unmanned Systems (ICTUS'2017)",Amity Directorate of Engineering & Technology (ADET),Dubai during 18th-20th December 2017.(Accepted and Presented for Publication)
7. Mr.Sarang M. Patil, Prof.Peter Z. Petkov, Prof.Boncho G. Bonev and Hitesh Singh, "A Novel Design of Enhanced Ultra-wide

Bandwidth Fractal Dipole Antenna for Wireless Applications”,IEEE 5th International Conference on "Computing for Sustainable Global Development," Bharati Vidyapeeth's

Institute of Computer Applications and Management (BVICAM), New Delhi (INDIA)
14th - 16th March 2018(Accepted)

Authors' contacts

Organization:

Research Fellow,

Department of RCVT, Faculty of Telecommunication.

Technical University of Sofia

Sofia, Bulgaria

sarang.p86@gmail.com

AN EXPERIMENTAL SETUP FOR DETERMINATION OF THE RESONANT FREQUENCIES OF A MECHANICAL FRAME STRUCTURE

SVETLIN STOYANOV

Abstract: *An experimental setup is developed to measure, process, and analyze the free vibrations of a plane frame structure. The vibration signal processing is realized in the software systems LabView and Matlab. The vibration spectrograms are obtained. The resonance frequencies are determined and compared with theoretical results from the corresponding eigenvalue problem solution.*

Key words: *experimental setup, resonant frequencies, natural frequencies, eigenvalue problem, the Fourier transform, frame structure, finite elements method, Abaqus, LabView, Matlab*

1. Introduction

The resonant frequencies and the corresponding natural mode shapes of a mechanical structure can be obtained by modal or by operational analysis. In the case of modal analysis, the excitation applied to the mechanical structure investigated is fully known. Usually, one uses an impulse hammer to realize excitation force with a known value. However, when the vibration analysis is realized in working conditions, the input excitations are unknown or partially known. This analysis is named operational analysis.

In [4] is presented an investigation of the natural and resonant frequencies of a cantilever beam. The frame structures are, however, more commonly used in the areas of mechanical and civil engineering [2]. The determination of the frequency response function of a mechanical structure requires the value of the excitation force to be known [3]. However, the resonant frequencies can be obtained without knowing of the excitation force value, based only to the structure response [4]. Taking into account that the resonant frequencies are the most important information to generate or to avoid resonance, the aim of this work is to develop an experimental setup for determination of the resonant frequencies of a frame structure.

2. Theoretical background

By applying of Newton's second law of motion, one obtains the government equation of the mechanical structure free vibration:

$$[M]\{\ddot{X}\} + [C]\{\dot{X}\} + [K]\{X\} = \{0\}, \quad (1)$$

where $[M]$, $[C]$, and $[K]$ are the mass, damping, and stiffness matrix of the mechanical structure investigated; $\{X\}$, $\{\dot{X}\}$, and $\{\ddot{X}\}$ are the vectors, which contains the displacement, velocity, and acceleration of each structure degree of freedom. The displacement $\{X\}$ can be described as a harmonic function of a vibration amplitude $\{\phi\}$ and a circular frequency ω as follows:

$$\{X\} = \{\phi\} \sin(\omega t). \quad (2)$$

Substituting (2) in (1), and neglecting of the damping leads to the eigenvalue problem [1]:

$$(-\omega^2[M] + [K])\{\phi\} = 0 \quad (3)$$

with the characteristic equation

$$|[K] - \omega^2[M]| = 0. \quad (4)$$

The frequency ω represents the natural circular frequency, its square represents the eigenvalue with the corresponding eigenvector $\{\phi\}$, which describes the mode shape of the system.

When the mass matrix is positive definite, all eigenvalues are positive. Rigid body modes and instabilities cause the mass matrix to be indefinite. Rigid body modes produce zero eigenvalues and instabilities produce negative eigenvalues [1].

3. The experimental setup description

The experimental setup created includes a plane steel frame structure, a piezoelectric accelerometer, a signal conditioner, a data acquisition device, and a computer system. For the measured signal processing, a LabView virtual instrument is developed. A photo of the setup is

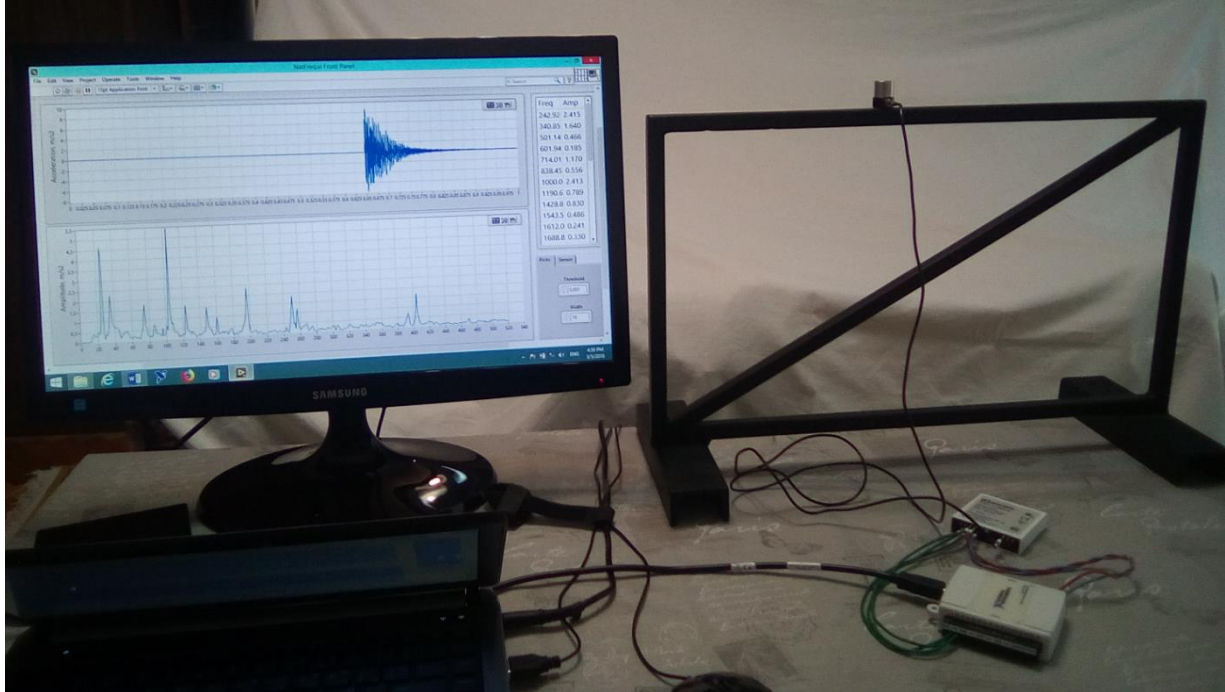


Fig. 1. A photo of the setup created

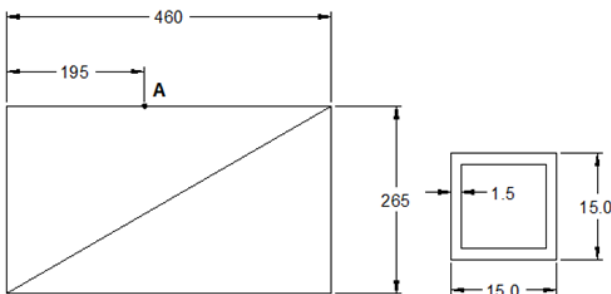


Fig. 2. A drawing of the frame used in the experimental setup

The voltage generated from the accelerometer is amplified by a signal conditioner SCM5B-40. This signal conditioner provides a channel of analog input, which is amplified, isolated, and converted to a high-level analog voltage output. The amplified analog voltage is converted then to a digital voltage through a data acquisition device NI USB-6009. After that, the vibration signal is processed by LabView virtual instrument developed. The virtual instrument is able to realize the Fourier transform and obtain the vibration spectrograms. Also, the virtual instrument developed finds the peak points on the spectrograms

shown on Fig. 1. A drawing of the plane frame used in the setup is shown on Fig. 2. The accelerometer is attached at point A. The type of the piezoelectric accelerometer used is KD 35. This accelerometer has voltage sensitivity of 5.01 mV/ms^{-2} and resonance frequency above 10 kHz .

and determine some of the resonant frequencies of the frame structure investigated.

4. Theoretical results

A CAD model of the frame structure is created in the integrated working environment of the software system Abaqus. The corresponding eigenproblem is solved according to (3) and the natural frequencies and mode shapes are determined. The first thirteen of the natural frequencies are shown in Table 1, and the first six natural mode shapes are presented on Fig. 3. Some of the nodes and antinodes are also shown.

Table 1. Theoretically obtained values for the natural frequencies

N	Natural frequency, Hz	N	Natural frequency, Hz
1	254	8	1112
2	315	9	1236
3	414	10	1540
4	706	11	1795
5	839	12	2111
6	901	13	2410
7	1039		

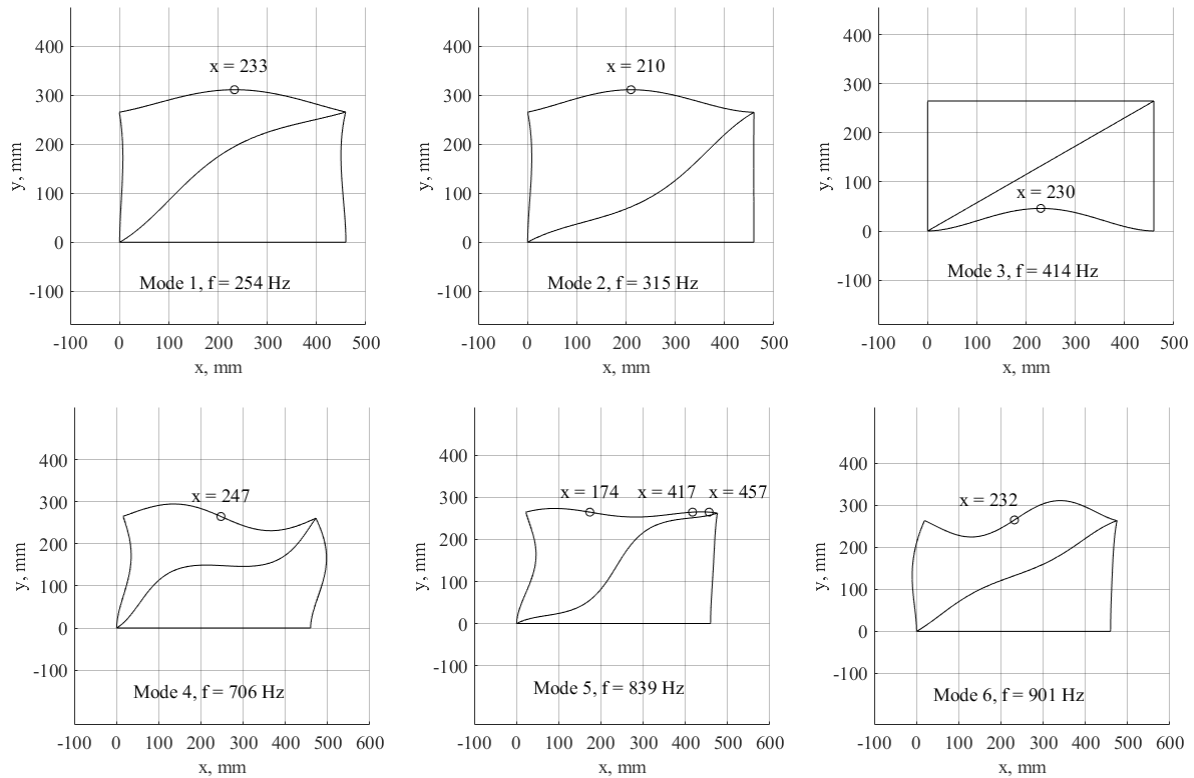


Fig. 3. Theoretically obtained natural mode shapes

5. Experimental results and comparison

The experimental setup created is used to determine some of the resonant frequencies of the frame structure used. For this purpose, an impulse initial excitation is applied at point A of the frame. This cause the frame to start running free damped vibrations. The vibration acceleration time-diagram is shown on Fig. 4. The vibration time-domain signal of a point of the system is polyharmonic, i.e.

it is a sum of harmonics each with its own amplitude and frequency. The frequencies of these harmonics are free damped vibration frequencies, i.e. they are some of the resonant frequencies of the structure investigated. Therefore, these resonant frequencies can be determined by decomposing the polyharmonic signal through the Fourier transform. The acceleration spectrograms obtained and the resonant peaks are shown on Fig. 5.

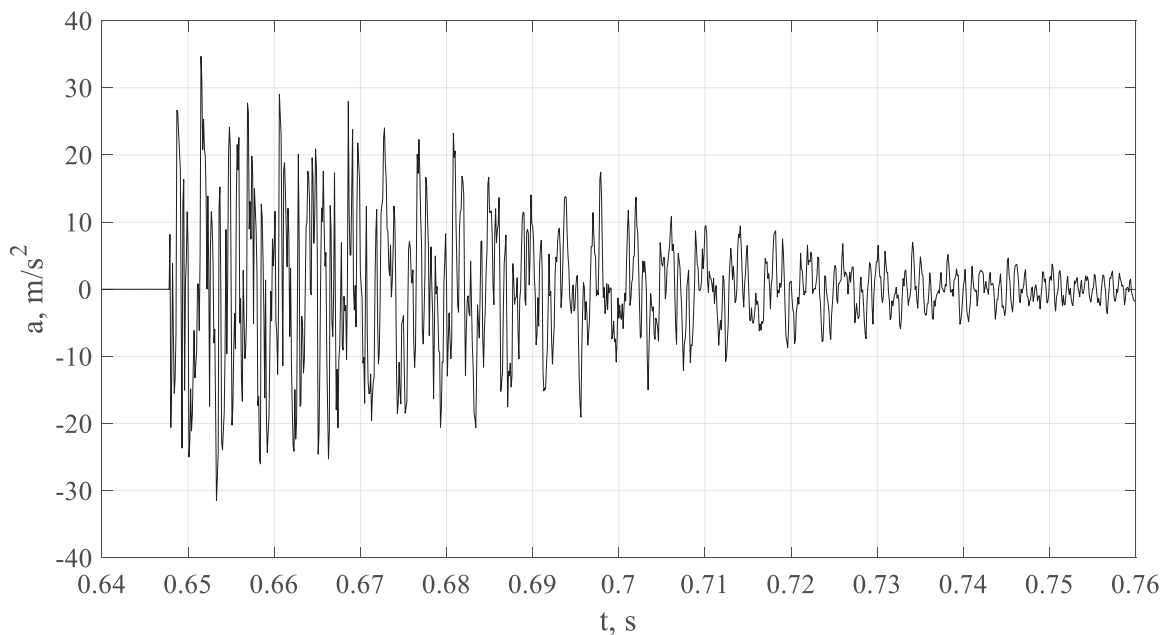


Fig. 4. A time-diagram of the measured vibration acceleration

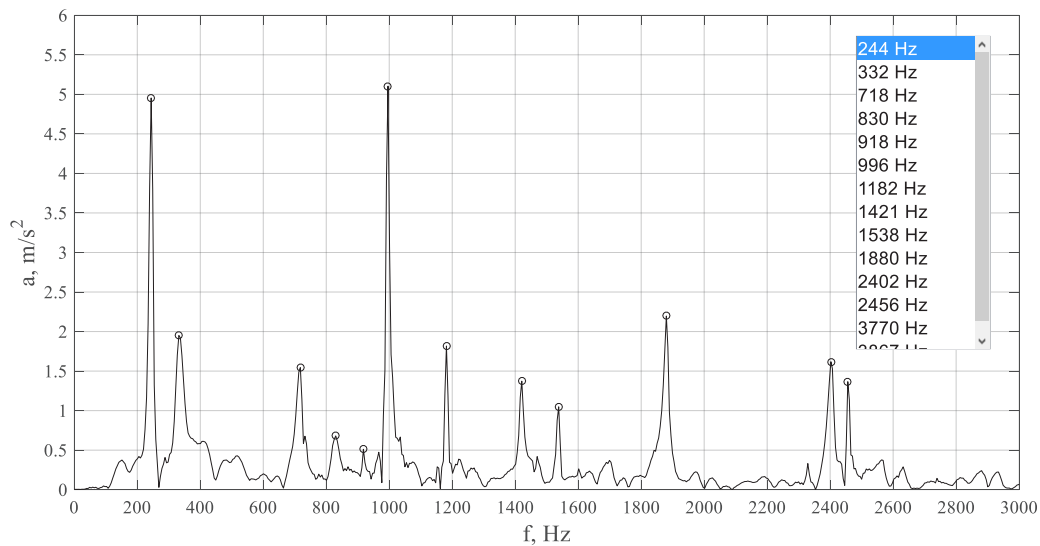


Fig. 5. A spectrogram of the measured vibration acceleration

From (3), one can see that the natural frequencies are derived neglecting the system damping. Therefore, they are higher than the corresponding resonance frequencies. The difference, however, can be usually neglected [4]. Thus, this comparative study is done by comparison the experimentally obtained resonance frequencies

to the theoretically obtained natural frequencies. The results are shown in Table 2. For the 3rd and 8th natural mode shape, the top frame member has no deformations (Fig. 3) and the corresponding resonant frequencies cannot be determined with an accelerometer mouthed at point A (Fig. 2).

Table 2. The results comparison

Mode N	Theoretically obtained natural frequency, Hz	Experimentally obtained natural frequency, Hz	Absolute difference, Hz	Relative difference, %
1	254	244	10	4
2	315	332	15	5
3	414	-	-	-
4	706	718	10	1
5	839	830	10	1
6	901	918	15	2
7	1039	996	44	4
8	1112	-	-	-
9	1236	1182	54	4

6. Conclusion

An experimental setup for determination of the resonant frequencies of a plane frame is created. The results obtained with the setup are compared to theoretical results and the difference is under 6%.

REFERENCES

1. Dessault Systemes Simulia Corp. (2011). *Abaqus 6.14 Theory Manual*. Providence, RI, USA
2. Madhu, Ps. and Venugopal T. (2014). Static Analysis, Design Modification and Modal Analysis of Structural Chassis Frame. *International Journal of Engineering Research and Applications (IJERA)*, volume 4.
3. Peng, Y., Li, B., and Mao, X. (2018). A method to obtain the in-process FRF of a

machine tool based on operational modal analysis and experiment modal analysis. *The International Journal of Advanced Manufacturing Technology*, volume 95.

4. Stoyanov, S. (2017). Sensors mass influence on the natural frequency of a cantilever beam. *Journal of the Technical University - Sofia Plovdiv branch, Bulgaria "Fundamental Sciences and Applications"*, volume 23.

Authors' contacts

Organization: "Angel Kanchev" University of Ruse

Address: 36 Zahari Stoyanov str., POB 7005, Ruse, Bulgaria

E-mail: sstoyanov@uni-ruse.bg

STUDYING AN AXIAL GENERATOR WITH ROTATING MAGNETS IN ITS STATOR WINDINGS

NIKOLA GEORGIEV

Abstract: *An axial generator with modified construction has been considered here, with rotating magnets in its stator windings. Simulation by the finite elements method Femm 4.2 has been carried out and both the distribution of the magnetic induction and the induced e.m.f. in one stator winding have been obtained. A model in OrCAD has been created, by means of which the phase voltage on an active load has been simulated and then compared to the experimentally measured phase voltage. A scheme of a voltage doubler has also been simulated in OrCAD and the obtained phase voltage has been compared to the measured rectified voltage. The obtained models simulate with good precision the operation of the studied generator.*

Key words: *axial, generator, rotating, magnets, permanent*

1. Introduction

Low power axial generators with permanent magnets in their rotors find application in practice most frequently as wind or hydro-generators due to their simple, reliable and easy to produce construction. There are no excitation windings or current in these generators, which leads to high efficiency in operation [1], while the considerable air gap reduces the magnetic attraction between the rotor and the stator, as well as the resistive moment [2].

Single-phase axial generators are easier to produce than the three-phase ones and the voltage in them is higher. Papers [3] and [4] consider similar single-phase generators with two rotors and a stator, for which the flux linkage has been calculated by the finite elements method, from where the r.m.s. value of the phase e.m.f. has been calculated.

The present paper presents a model of the magnetic field in one stator winding in a low power axial two-rotor generator with rare-earth magnets and one stator with rotating permanent magnets in it, developed by the finite elements method Femm 4.2.

2. Exposition

The axial generator, considered here, consists of two steel rotors with dimensions D200xH4 mm with sixteen rare-earth magnets each, measuring 20x20x10 mm. The stator, made of turbonit, measures 240x 240x4 mm and has eight windings

with 600 turns of a conductor with cross-sectional area $S=0,385 \text{ mm}^2$ in each. The air gap (the distance between two opposite magnets in the rotors) is $l_g=40 \text{ mm}$.

The two-dimensional method of finite elements Femm 4.2 [5] is used in the process of modeling, since this is the way of defining the magnetic induction along the vertical for a stator winding of the generator. After that the magnetic flux is defined and the induced e.m.f. in one of the stator windings of the generator with rotating magnets is calculated at active load of 10 Ω .

Fig. 1 shows the constructive scheme of the studied two-rotor generator with rotating magnets in its stator windings and the notations in the figure are as it follows: 1 – the turbonit stator; 2 – the stator windings; 3 – the steel rotors and 4 – the rotor magnets; by 5 the rotating rare-earth magnets in the stator windings are denoted.

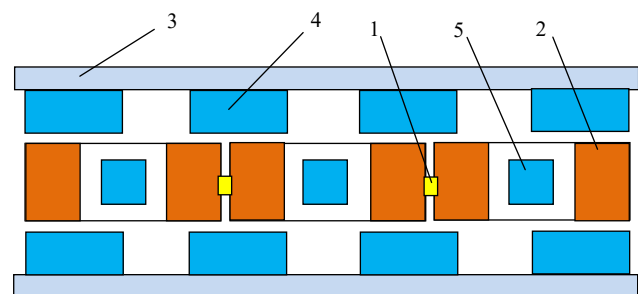


Fig.1. Generator with rotating magnets in its stator windings

By means of the method of finite elements Femm 4.2 the distribution of the magnetic field in the two-rotor single phase generator is obtained at

different angles of rotation of the permanent magnets in the stator windings – Fig. 2.

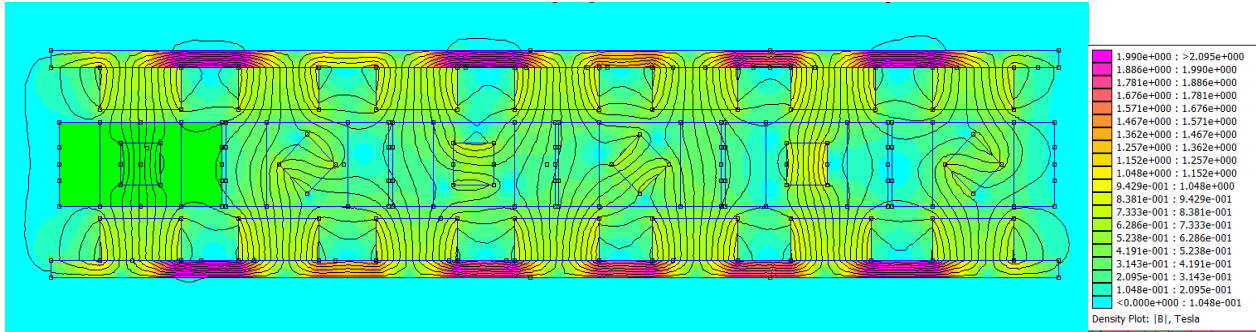


Fig.2. Distribution of the magnetic field in the generator with rotating magnets in its stator windings

θ - the angle measured in degrees.

Fig. 3 shows the integral magnetic induction for 1 m^3 along the vertical B_{cy} , as well as the volume of the stator winding V_c .

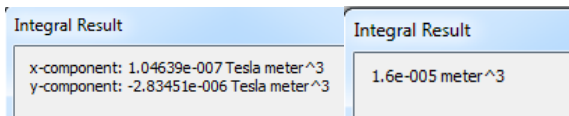


Fig.3.

With their help the magnetic induction for the volume of the stator winding B'_c is calculated and it is equal to

$$B'_c = \frac{B_{cy}}{V_c}, \quad T \quad (1)$$

By means of the defined magnetic induction it is possible to calculate the induced e.m.f. at different angles of rotation of the permanent magnet in one stator winding [6]

$$e(t) = N\omega AB'_c \cos \theta, \quad (2)$$

where: N is the number of turns in the winding;
 ω - the circular frequency;
 A - the cross-sectional area of the winding;

The circular frequency can be expressed by the number of revolutions per minute n and the number of pole pairs of the generator p

$$\omega = \frac{2\pi \cdot p}{60} n \quad (3)$$

From expressions (2) and (3) for the induced e.m.f. in one stator winding it is obtained

$$e(t) = N \frac{\pi \cdot p}{30} n AB'_c \cos \theta \quad (4)$$

With the help of the expression (1) and Fig. 3 the change of the integral magnetic induction for the volume of the stator winding is found – Fig. 4, while the induced e.m.f. at different angles of rotation of the permanent magnet in one stator winding at active load is obtained from expression (4) and Fig. 4 – Fig. 5.

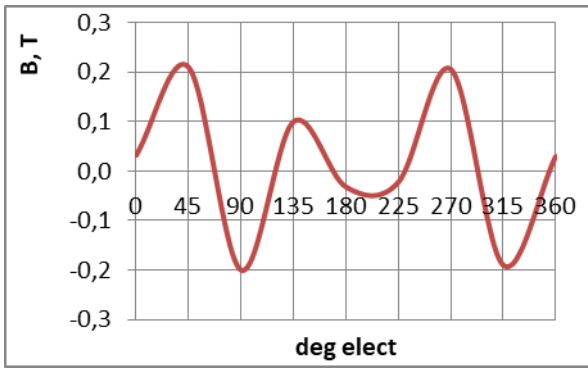


Fig.4.

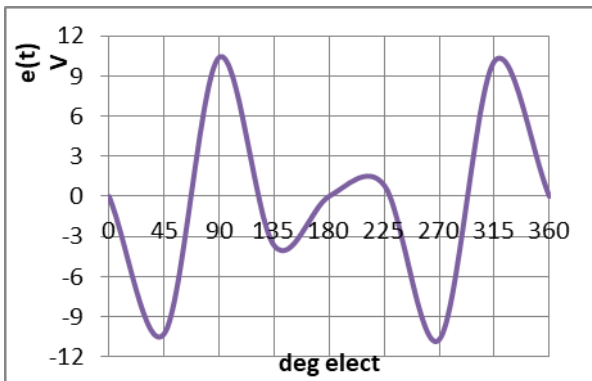


Fig.5.

By means of OrCAD – Fig. 6. – the instant form of the e.m.f. for one stator winding of the generator with rotating magnets in its stator windings is modeled. By V1, V2 and V3 here the electromotive voltages in instant form are denoted for the first, second and third harmonics respectively, while the active resistance and the inductivity of the stator winding are denoted by R1 and L1 correspondingly. R2 denotes the active load, connected to the winding.

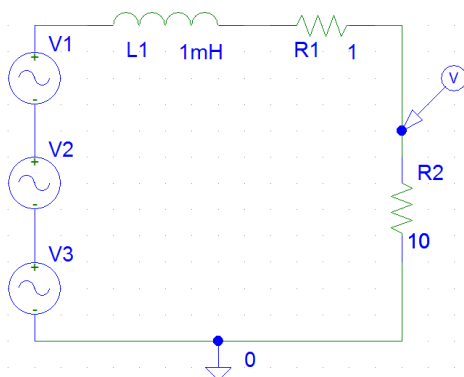


Fig.6.

The total e.m.f. in instant form for one stator winding is the sum of the electromotive voltages of the first, second and third harmonics

$$e(t) = e_1(t) + e_2(t) + e_3(t) \quad (5)$$

Fig. 7 presents the simulated change of the total phase voltage on the active load.

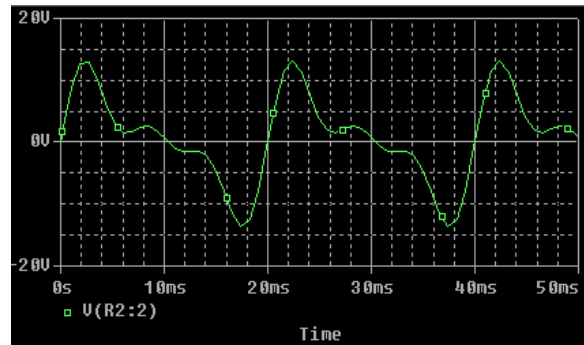


Fig.7.

Fig. 8 shows the voltage on the load in instant form for a stator winding with a rotating magnet, measured by an oscilloscope.



Fig.8.

When comparing the figures 7 and 8 it can be seen that the instant form of the phase voltage for one winding of the generator with rotating magnets in its stator windings in OrCAD well simulates the real phase voltage.

In order to test the simulations of the generator with rotating magnets in the stator windings, their phase voltage is rectified by means of the voltage doubler as in Fig. 9. The rectified phase voltage from the model in OrCAD is $U_{mod}=4,8 \text{ V}$ at $n=500 \text{ min}^{-1}$, Fig.10.

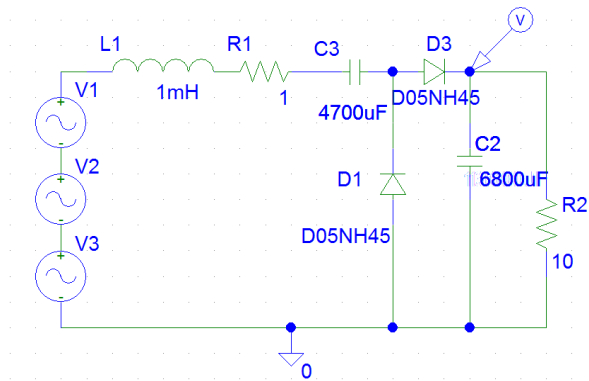


Fig.9.

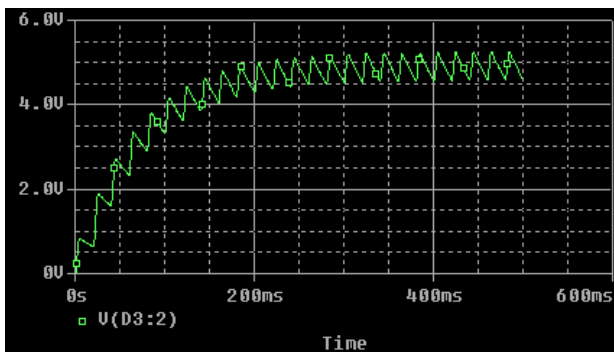


Fig.10.

Table 1 presents the phase voltages, both measured and calculated by the model, at $n=100$, 300 and 500 revolutions per minute for one stator winding, as well as the relative error of the model in OrCAD.

Table 1.

n, min^{-1}	100	300	500
U_{mod}, V	0,96	2,88	4,8
$U_{\text{meas}}, \text{V}$	0,91	2,74	4,62
$\delta, \%$	5,49	5,11	3,89

3.Conclusion

An axial generator with modified construction with rotating magnets in the stator windings has been considered. From a simulation by the method of finite elements Femm 4.2, the distribution of the magnetic induction at different angles of rotation of the permanent magnet in one stator winding has been obtained.

A model in OrCAD has been obtained, by means of which the phase voltage on an active load has been simulated and compared to the experimentally measured phase voltage. A scheme of the stator winding, as well as of a voltage doubler, have also been simulated in OrCAD and the obtained voltage has been compared to the measured rectified phase voltage.

The maximum relative error is comparatively low $\delta=5,49\%$, which confirms that the models, obtained by means of Femm 4.2 and OrCAD, simulate with good precision the operation of the studied generator.

References:

- [1] Gieras, J and Wang, R. (2004). *Axial flux permanent magnet brushless machines*, Dordrecht Netherlands: Kluwer Academic Publishers.
- [2] Mohammad, M. and Widyan, S. (2006). *Optimization, Construction and Test of Rare-Earth Permanent-Magnet Electrical Machines with New Topology for Wind Energy Applications*, Fakultät IV–Elektrotechnik und Informatik der Technischen Universität, Berlin, pp. 17-22.
- [3] Wang, R. and Kamper, M. (2005). *Optimal Design of a Coreless Stator Axial Flux Permanent-Magnet Generator*, IEEE Transactions on magnetics, vol. 41.
- [4] Wannakarn, P. and Kinnares, V. (2011). *Microcontroller based Grid Connected Inverter for Axial Flux Permanent Magnet Generator*, IEEE PEDS 2011, Singapore.
- [5] Meeker, D. (2006). *Finite Element Method Magnetics – Version 4.0*, User's Manual.
- [6] Gupta, B. R. and Singhal, V. (2007) *Energy conversion*, Paperback, ISBN 978-81-224-2061-6, pp.6-7.

Department of Electrical Engineering
 Technical University–Sofia, Branch Plovdiv
 25 Tsanko Dystabanov St.
 4000 Plovdiv
 BULGARIA
 E-mail: nikola.georgiev @tu-plovdiv.bg

IMPLEMENTATION OF A NOVEL FORCE COMPUTATION METHOD IN THE FEMM SOFTWARE

VASIL SPASOV, IVAN KOSTOV, IVAN HADZHIEV

Abstract: A novel force computation method is implemented in the Finite Element Method Magnetics software - the nodal force method. For this purpose supplementary code is developed in C++ and added to the FEMM source code. In addition to force computation, the extended version of FEMM enables also force visualization that was not possible until now. To demonstrate the capabilities of the extended version of FEMM, the forces of three models are computed and visualized – two current-carrying copper busbars, an AlNiCo permanent magnet with steel core and a benchmark non-linear dc electromagnet. A comparison is made between the newly implemented nodal force method and the available Maxwell's stress tensor method from the viewpoint of accuracy and visualization capabilities.

Key words: Finite element method, electromagnetic force, nodal force method, FEMM

1. Introduction

Finite Element Method Magnetics (FEMM) is a finite element software for solving 2D problems in low frequency magnetics and electrostatics [1]. The program addresses linear and nonlinear magnetostatic problems, time harmonic magnetic problems and others. FEMM has been extensively used in science, engineering, industry and for teaching electromagnetics in higher education [2]. It is a free, open source, accurate and low computational cost product. There is no limit on the problem size – the maximum number of finite elements and nodes is limited only by the amount of available memory. This enables to solve problems resulting in more than a million elements on a personal computer.

The aim of this paper is to extend the capabilities of the FEMM software when computing and visualizing electromagnetic forces of electrical devices. For that purpose the mathematical model of a novel electromagnetic force computation method - the nodal force method, is developed [3, 4]. This mathematical model is implemented in the FEMM postprocessor by developing a C++ code.

To verify the extended version of FEMM, the forces of three models are computed and visualized. The models are two current-carrying copper buses, an AlNiCo permanent magnet with steel core and a benchmark non-linear dc electromagnet.

The present paper is organized as follows. The derivation and analysis of the nodal force method are presented in Section 2. The implementation of the nodal force method and the force visualization enhancements to FEMM are discussed in Section 3. The accuracy of the extended version of FEMM is validated numerically in Section 4. Finally, conclusions are drawn in Section 5.

2. Derivation of the nodal force method

The nodal force method (NFM) is derived for the two-dimensional case using first-order nodal triangular finite elements. In NFM the work δW performed by electromagnetic force for displacement δu is [5]:

$$\delta W = - \iint \tau_{ik} \frac{\partial(\delta u_i)}{\partial k} d\Omega; \quad (i, k = x, y). \quad (1)$$

Here τ_{ik} are the Maxwell stress tensor components and Ω is the analyzed region.

The Maxwell stress tensor components are defined as:

$$[T] = \frac{1}{\mu_0} \begin{bmatrix} B_x^2 - 0.5B^2 & B_x B_y \\ B_y B_x & B_y^2 - 0.5B^2 \end{bmatrix} = \begin{bmatrix} \tau_{xx} & \tau_{xy} \\ \tau_{yx} & \tau_{yy} \end{bmatrix}, \quad (2)$$

where μ_0 is the permeability of air, B is the magnetic flux density magnitude and B_j ($j = x, y$) is the magnetic flux density component along the two axes.

The displacement is interpolated by the well-known continuous and piecewise-differentiable shape functions N_i of nodal triangular elements [6]:

$$\delta u_i = \sum_n N_i \delta u_{ni}, \quad (3)$$

where n is the number of element nodes and i is the shape function number.

Replacing (3) in (1) yields:

$$\delta W = \sum_n \left(- \iint \tau_{ik} \frac{\partial N_i}{\partial k} d\Omega \right) \delta u_{ni}; \quad (i, k = x, y). \quad (4)$$

On the other hand, the completed work is equal to:

$$\delta W = \sum_n f_{ni} \delta u_{ni}, \quad (5)$$

where f_{ni} is the i^{th} component ($i = x, y$) of the nodal force acting on node n .

After equating (4) to (5) it is obtained for the nodal force:

$$f_{ni} = - \iint \tau_{ik} \frac{\partial N_i}{\partial k} d\Omega. \quad (6)$$

Based on (6), the x component of the force of node n of one triangular finite element can be computed by the 2D finite element method as follows:

$$f_{nx} = -(\tau_{xx} b_k + \tau_{xy} c_k) \cdot S_e. \quad (7)$$

Here b_k and c_k are the shape functions coefficients of the nodal first-order triangle and S_e is its area.

The integration in (6) for a node is performed for all elements to which the node belongs. The total force acting on a part is obtained by summing up the nodal forces of all nodes included in the part.

Next the computer implementation of NFM will be analyzed. Formula (7) shows that the nodal force method uses only quantities that have already been computed during the finite element analysis. In other words, due to the absence of additional arithmetic operations, the NFM needs less CPU time as compared to the Maxwell's stress tensor method (MSTM).

To perform (7), no integration contour should be defined, as required by the MSTM in FEMM. Thus two more advantages are to be expected: the NFM can be implemented fully automatically in models of arbitrary shape and its accuracy is not affected by the choice of the integration contour needed by the MSTM [1].

Another important advantage of the NFM is that it directly computes the local electromagnetic force, i.e. the force acting on every finite element node. Therefore, to create vector plots of force, no additional post-processing is needed. In contrast, the Maxwell's stress tensor method calculates only

global force. Local force is often needed for the analysis and design of electrical devices and for post-processing purposes, as shown below in Section 4.

The above-mentioned advantages make the NFM very attractive to use. Due to these advantages, the NFM has been the method of choice to extend the capabilities of the FEMM software.

3. Implementation of the nodal force method in the FEMM software

The nodal force method is implemented in the Finite Element Method Magnetics software using object-oriented programming in C++ [7]. To visualize the computed force, supplementary code is developed and added to the FEMM source code.

First a base ForceAlgorithm class is created which defines the common attributes and member functions of all descendent algorithms [7]. The NodalForce class inherits the ForceAlgorithm characteristics and specifies additional methods related to the implementation of the nodal force algorithm. This enables to easily add other classes such as the VirtualWorkForce in the future.

The ForceAlgorithm class defines a virtual function solve(). This means that all derived algorithm classes from this class must provide their own implementation of the function, thus defining a common pattern of usage. As a child of ForceAlgorithm, the NodalForce class contains the code needed to actually perform the calculation in solve(). The implementation details of the above algorithm are hidden by abstraction in the ForceAlgorithm class.

The class ForceAlgorithm stands at the base of the class hierarchy. It contains methods and attributes which are common to the different algorithms. They include parameters of one finite element such as the magnetic flux density components and magnetic permeabilities along the two axes. A MeshNode array is created containing the element node coordinates. Memory is allocated dynamically in the constructor of ForceAlgorithm for the result returned from solve(). The destructor of ForceAlgorithm cares for memory deallocation.

The NodalForce class is a subclass of ForceAlgorithm and inherits all attributes from the parent class. The forces acting on the three nodes of one finite element are calculated in solve(). For this purpose the Maxwell stress tensor components in (2), the shape function coefficients and the force components in (7) are computed and assigned to variables. Then the method solve() returns the force vectors in the three nodes of the finite element being analyzed.

The invocation of the nodal force algorithm is done in the Block Integrals section of the code.

The Block Integral is applied over the selected integration areas from the model. After the integration areas are selected, invocation of the nodal force algorithm can be done.

The method `add_vector_result` performs the integration in (6). This integration is carried out by summing up the force components along the two axes on the nodes of the selected finite elements. The reference passing avoids duplication of objects in memory while calculating the nodal force.

Due to the extension of FEMM by the nodal force method, the original Block Integrals dialog for invoking the Nodal Force algorithm is modified by adding one more entry (Nodal Force) to the drop-down list. The modified dialog of the extended version of FEMM is shown in Fig. 1.

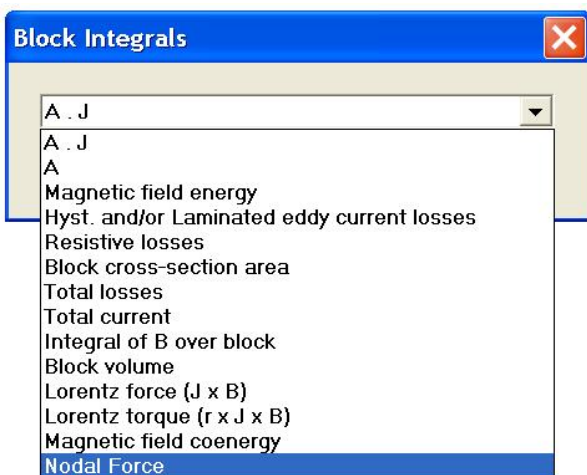


Fig. 1. Modified dialog for the Nodal Force algorithm

As shown in Section 2, an important advantage of the nodal force method is that it computes local force. To utilize this advantage, the visualization capabilities of FEMM are enhanced by developing supplementary code in C++ and adding it to the source code. The original View context menu is modified by adding one more entry (Force vectors) to the drop-down list in Fig. 2.

Fig. 3 shows the entirely new dialog designed for the purposes of force visualization.

The local forces acting on the finite elements nodes are displayed as vectors whose direction coincides with the direction of force. The length of these vectors is obtained by scaling the magnitude of forces using the Choose scale slider at the top of the dialog in Fig. 3. The value in the Maximum length edit box shows the longest vector length when the slider is set to the rightmost position.

The vectors of computed forces are drawn as arrows. As seen in Fig. 3, several options for the head angle, length and colour of arrows are

provided. Applications of the developed visualization enhancements are given in the next section.

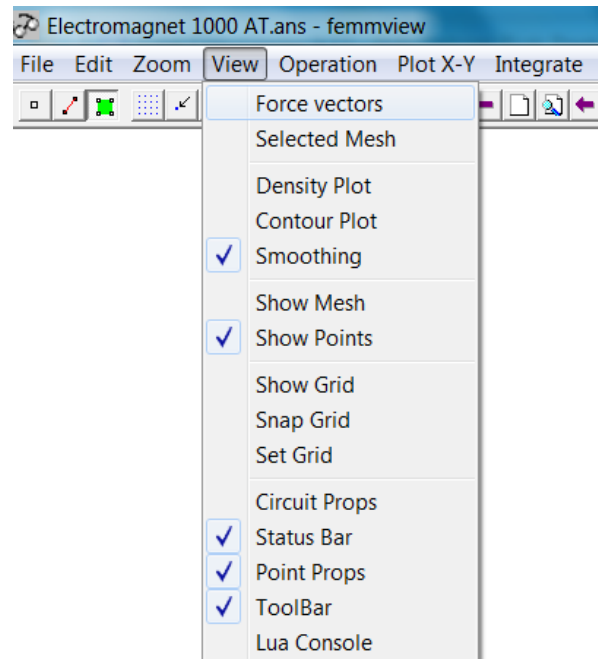


Fig. 2. Modified context menu for the Nodal Force visualization

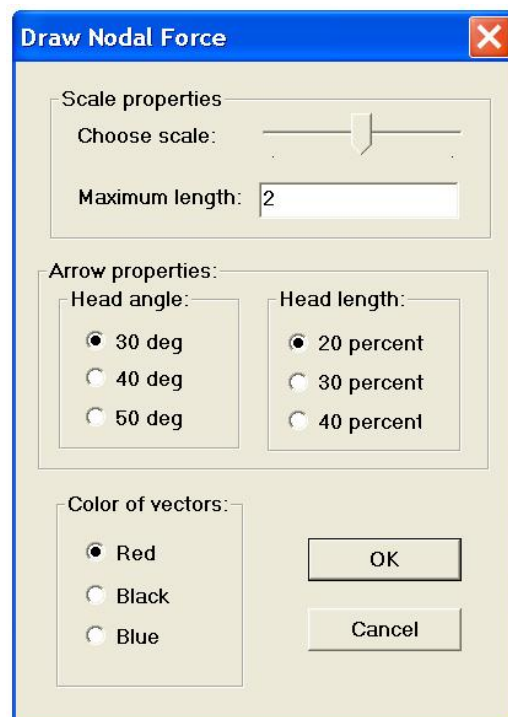


Fig. 3. New dialog for the Nodal Force visualization

The above described extensions to the Finite Element Method Magnetics source code are made according to the Document/View Architecture (MVC) [8].

4. Accuracy validation and visualization by the extended version of FEMM

In this section the accuracy of the extended version of FEMM is validated by comparing the implemented nodal force method with the Maxwell's stress tensor method, available in the conventional FEMM. For that purpose three models are analyzed and their electromagnetic forces are visualized by the new FEMM capabilities.

The first model consists of two copper busbars carrying currents of 100 kA in the same direction [9]. The geometry is shown in Fig. 4. The dimensions are in centimeters. The finite element mesh has 25427 nodes and 50295 triangles.

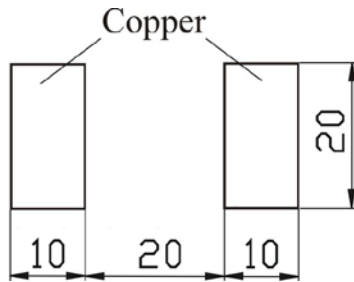


Fig. 4. Geometry of the copper busbars

Table 1 shows the total x-axis electromagnetic forces between the busbars when the currents flow in the same direction. The forces are computed analytically, by the Maxwell stress tensor method and by the nodal force method.

Table 1. Forces between the busbars

Analytical [N]	MSTM [N]	NFM [N]
6333	6328	6336

The vector plot of the local forces acting on the finite element nodes is shown in Fig. 5. The plot is generated using the visualization enhancements to FEMM described in the previous section.

The plot in Fig. 5 confirms the theory that conductors carrying currents of the same directions attract each other. As expected, the nodal forces act only in the current-carrying regions. Due to the nature of the nodal force method, the local force vectors originate from the finite elements nodes.

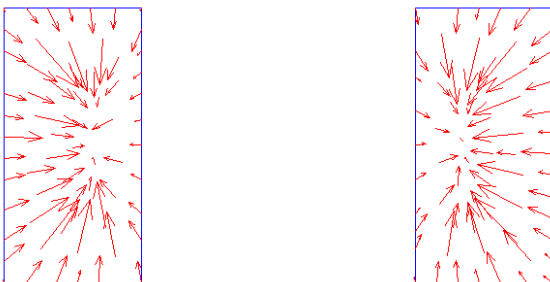


Fig. 5. Vector plot of the forces on the busbars

The second model is a non-linear AlNiCo permanent magnet with steel core [9]. The dimensions in centimeters are shown in Fig. 6. The finite element mesh has 71394 nodes and 141915 triangles.

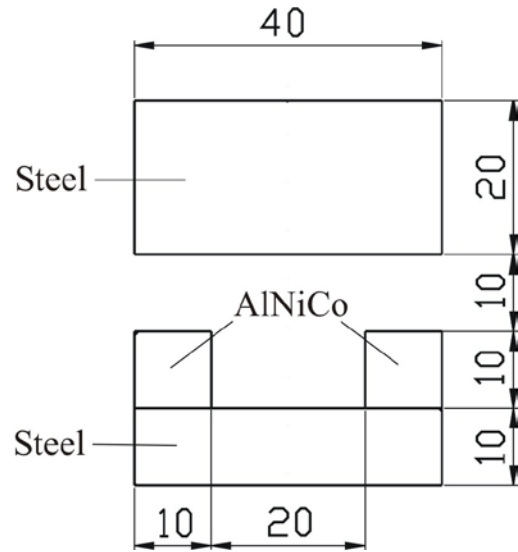


Fig. 6. Geometry of the AlNiCo magnet with steel core

Table 2 shows the total y-axis force acting on the steel armature, computed by the MSTM and NFM. The computed forces by the two methods are very close which confirms the accuracy of the built-in nodal force method.

Table 2. Forces on the steel armature

MSTM [N]	NFM [N]
1290	1293 N

The vector plot of the forces acting on the steel and on the permanent magnet is shown in Fig. 7. The plot is generated using the new visualization capabilities of FEMM.

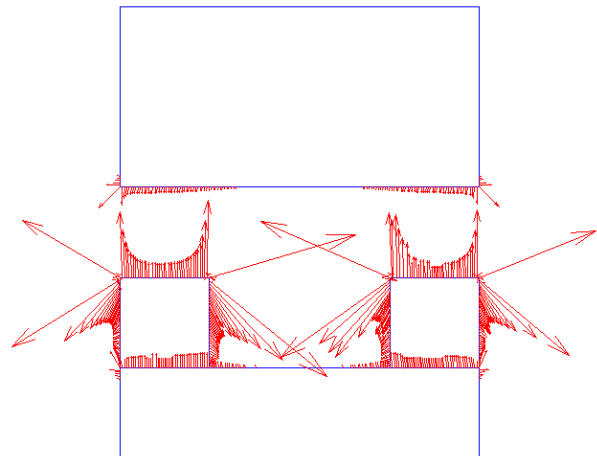


Fig. 7. Vector plot of the forces on the magnet

The third model is a non-linear dc electromagnet. This is a benchmark model used for evaluating the accuracy of the methods for force computation as well as for validation of computer programs [10].

The electromagnet has complex geometry, very small air gaps and uneven magnetic flux distribution. To saturate the steel, the excitation current is varied within a wide range.

The geometry of the model is shown in Fig. 8. The electromagnet is comprised of a steel yoke 1, coil 2 and a central pole 3. The yoke and the central pole are made of steel. The coil has 381 turns and is fed by dc current. The reluctivity curve of steel is given in [10]. To analyze the steel saturation effect, the total current in the coil has values 1000, 2000, 3000, 4000 and 5000 ampere turns. The finite element mesh has 81225 nodes and 162145 first order triangles.

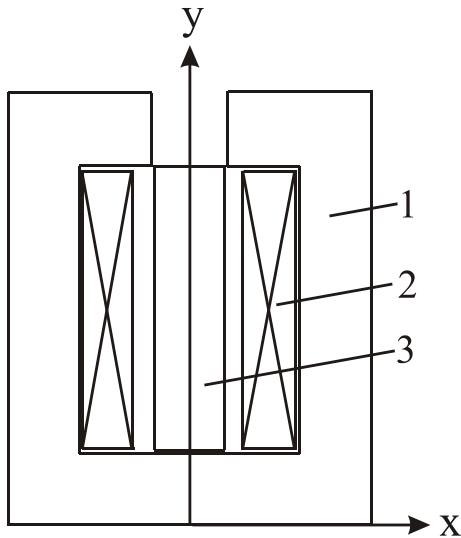


Fig. 8. Geometry of the benchmark electromagnet

Table 3 shows the y-axis forces acting on the central pole. They are computed by the newly implemented NFM and the standard MSTM in FEMM. The relative error in force by the NFM and the MSTM is determined by the formula:

$$\varepsilon_r = (F_{\text{NFM}} - F_{\text{MSTM}}) / F_{\text{MSTM}}, \quad (8)$$

where F_{NFM} and F_{MSTM} are the forces by the NFM and MSTM, respectively.

The force by the MSTM is used as reference in (8) instead of the measured values in [10], since the benchmark model analyzed in this paper is two-dimensional.

The absolute values of the relative errors in force are given in Table. 3. They show that the NFM, implemented in FEMM, has excellent accuracy, the maximum relative error in force being less than 1%.

Table 3. Y-axis forces acting on the central pole and relative error

current [A]	NFM [N]	MSTM [N]	ε_r [%]
1000	324.5	326.9	0.73
2000	1363.6	1369.5	0.43
3000	3170.4	3186.3	0.50
4000	5727.8	5741.5	0.24
5000	8693.8	8721.7	0.32

The vector plot of the nodal forces on the pole and yoke at current 1000 A is shown in Fig. 9. The plot is generated using the visualization enhancements to FEMM.

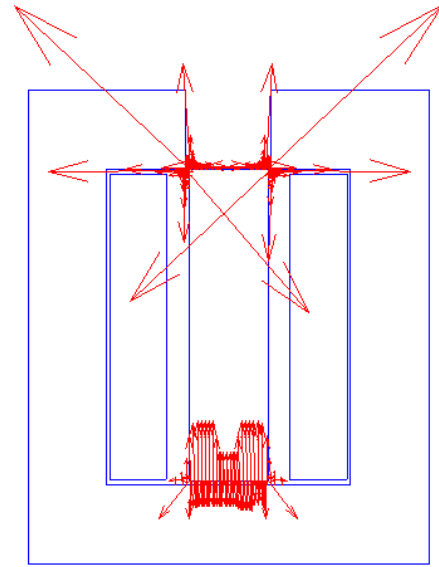


Fig. 9. Forces on the pole and yoke

A zoomed-in view of the force distribution in the upper left corner of the pole is given in Fig. 10.

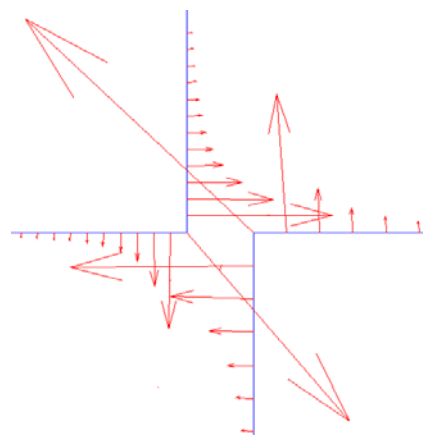


Fig. 10. Forces on the upper left corner of pole

Fig. 11 shows the forces on the lower left half of the pole and yoke.

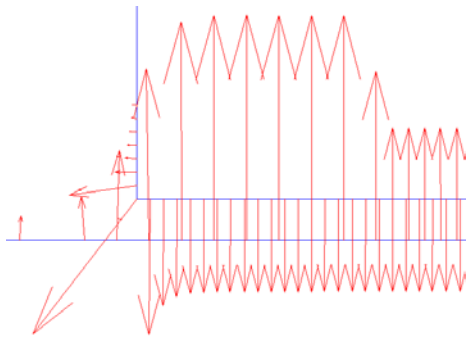


Fig. 11. Forces on the lower left half of the pole and yoke

The results from Figures 9, 10 and 11 show that the force has complex distribution and attracts the pole to the yoke. As expected, in the upper part of the pole the force is directed both along the x and y axes, while in the lower part it is mostly along the y axis. Due to the different mesh sizes in the yoke and central pole, the magnitudes of the force vectors on both sides of the air gap in Fig. 11 are different.

The results from Tables 1, 2 and 3 confirm the accuracy of the extended version of the Finite Element Method Magnetics software. There is a very good agreement between the forces computed by the NFM and the MSTM.

5. Conclusion

The capabilities of the FEMM software are enhanced by adding a novel method for force computation and visualization - the nodal force method. The NFM is implemented using object-oriented programming in C++.

The extended version of FEMM is validated by comparing the added-on nodal force method with the Maxwell's stress tensor method. The forces of three models are computed and visualized. The results show that the NFM has excellent accuracy.

The visualization capabilities of FEMM are also enhanced by developing supplementary code in C++ and adding it to the source code. Vector plots of the electromagnetic forces of the three models are created. The plots yield reasonable results.

Based on these plots it can be concluded, that nodal forces are localized only on the surface nodes of the steel and on all nodes of the current-carrying coil. This coincides well with the real physical situation where force acts only on the surface of magnetic materials, while the force in electric conductors manifests itself as a volumetric force (Lorentz force).

The extended version of the Finite Element Method Magnetics software can be used for research and design purposes, as well as for teaching numerical methods in electromagnetism and CAD systems at higher schools.

REFERENCES

1. Meeker, D. (2006). *FEMM 4.2 Magnetostatic tutorial*.
2. Baltzis, K. (2008). The FEMM package: a simple, fast and accurate open source electromagnetic tool in science and engineering. *Journal of Engineering Science and Technology Review*, Vol. 1, pp. 83-89.
3. Spasov, V., Noguchi, S., and Yamashita, H. (2001). Comparison of the methods for electromagnetic force computation by edge elements. *Proceedings of the Third Asian Symposium on Applied Electromagnetics, Hangzhou, China, May 28-30*, pp. 111-114.
4. Spasov, V., Noguchi, S., and Yamashita, H. (2001). Comparative analysis of the force computation methods in the 3D FEM with edge and nodal elements. *Electrical Engineering Research Conference, Kita Kyushu, Japan, SA-01-22, RM-01-90*, pp. 9-13.
5. Spasov, V. (2005). Computation of electromagnetic force by the nodal force method. *XIV-th International Symposium on Electrical Apparatus and Technologies SIELA 2005, Plovdiv, Vol. II*, pp. 139-144.
6. Salon, S. (1995). *Finite element analysis of electrical machines*, 247 p. Kluwer.
7. Preiss, B. (1998). *Data structures and algorithms with object-oriented design patterns in C++*, 688 p. Wiley.
8. Prosiše, J. (2003). *Programming Windows with MFC*, 1376 p. Microsoft Press.
9. Brandisky, K., and Yatcheva, I. (2002). *CAD systems in electromagnetism*, CIELA, 244 p. Sofia.
10. Takahashi, N., Nakata, T., et al. (1994). Investigation of a model to verify software for 3-D static force calculation. *IEEE Transactions on Magnetics*, Vol. 30, No. 5, pp. 3483-3487.

Assoc. Prof. Vasil Spasov, Ph.D.
Department of Electrical Engineering
E-mail: vasilspasov@yahoo.com

Assoc. Prof. Ivan Kostov, Ph.D.
Control Systems Department
E-mail: ijk@tu-plovdiv.bg

Assistant Prof. Ivan Hadzhiev, Ph.D.
Department of Electrical Engineering
E-mail: hadzhiev_tu@abv.bg

Technical University - Sofia
Branch Plovdiv
25 Tsanko Dyustabanov Str.

STUDYING THE ELECTRICAL AND THERMAL FIELDS, PRODUCED BY THE CURRENT IN THE INSULATION OF A MIDDLE VOLTAGE CABLE

IVAN HADZHIEV, DIAN MALAMOV, VASIL SPASOV, DIMITAR NEDYALKOV

Abstract: The paper presents numerical and experimental studies of the characteristics of a power supply cable for middle voltage of 20kV. Experimental studies of the current, flowing through the insulation of the cable at direct testing voltage, have been conducted. A computer model of the cable has been synthesized in the software program Comsol. Based on the model, the distribution of both the electrical and the magnetic field of the cable have been studied. Finally, a comparison between the obtained numerical and experimental results has been drawn.

Key words: electric field, finite element method, power cable, thermal field

1. Introduction

Power supply cables are widely used in electricity distribution networks. Underground cable lines are frequently used with middle voltage electricity grids. Most part of these cables are exploited for many years. During the exploitation their electrical insulation is ageing. The process of insulation ageing is speeded up by thermal overloads [1], [2], [3]. One of the main factors for failures in the power supply cables is the damaged electrical insulation, which is due to partial internal discharges inside its interior.

This paper describes experimental and numerical studies of the current, flowing through the insulation of a power supply cable AOSB type (aluminum conductor, individual coating of the strands, power cable, lead-clad) for middle voltage. The experimental tests have been conducted at direct testing voltage. The numerical studies have been carried out by means of a synthesized computer model of the cable in the programming package Comsol. Based on of the developed computer model, the distribution of both the electrical and the magnetic field of the cable have been obtained. Finally, a comparison between the obtained numerical and experimental results has been drawn.

2. Failures in the middle voltage electrical cables

Statistical analysis of the failures in the middle voltage electrical cables has been carried out. The results are shown in Fig. 1 ÷ Fig. 4.

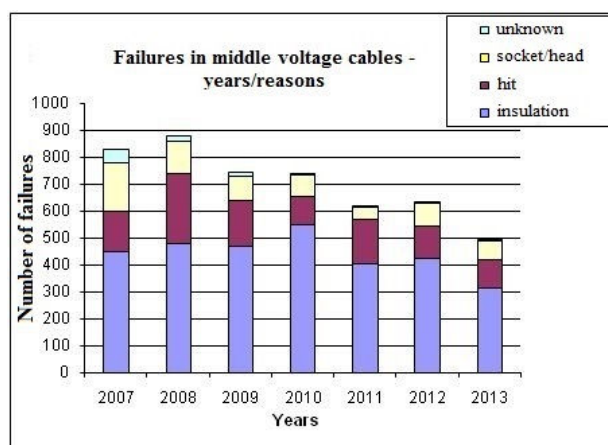


Fig. 1. Middle voltage cable failures – years/reasons

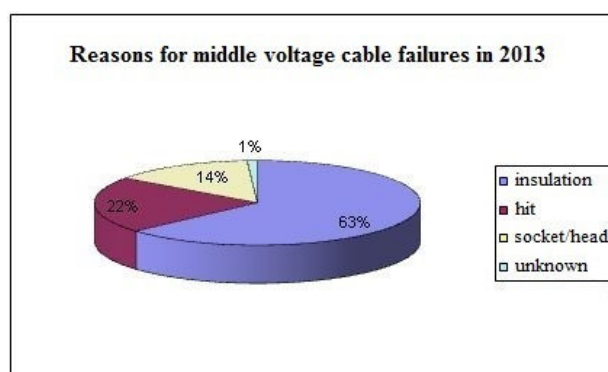


Fig. 2. Reasons for middle voltage cable failures in 2013

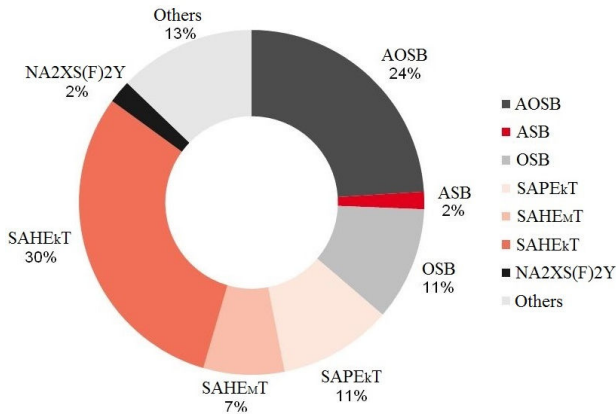


Fig. 3. Failures in the various types of middle voltage cables in 2013

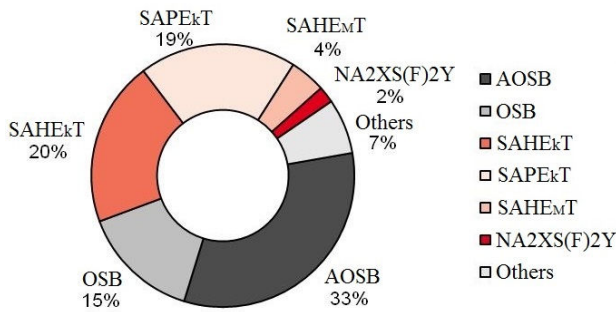


Fig. 4. Failures in the insulation of the various types of middle voltage cables around 2013

3. Experimental study of a middle voltage cable

3.1. Construction of the studied cable

Object of study is a middle voltage cable AOSB type (aluminum conductor, individual coating of the strands, power cable, lead-clad), whose construction is shown in Fig. 5. The data of the cable is as follows: diameter of the aluminum core - 7,7mm; thickness of the semiconductor shield - 0,3mm; thickness of the paper-oil insulation - 5,5mm; thickness of the lead mantle (coating) - 1,5mm; cable length - 0,8m.

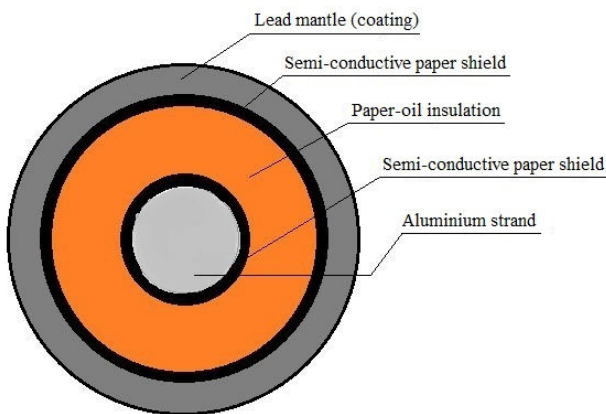


Fig. 5. Construction of a middle voltage cable AOSB type

3.2. Experimental results

Experimental tests were conducted with the described cable at direct testing voltage in accordance with [4]. The tests were carried out by a testing system BPS 5000-d, produced by the company SebaKMT. The values of the current, flowing through the insulation, are given in Table 1.

Table 1. Values of the current, flowing through the insulation in case of direct testing voltage for a AOSB type cable

№	Direct voltage [kV]	Current through the insulation [mA]
1	5	0
2	10	0,0007
3	15	0,002
4	20	0,0045
5	25	0,007
6	30	0,01
7	35	0,013
8	40	0,017
9	45	0,022
10	50	0,027

Fig. 6 shows the obtained graphical dependence of the current, flowing through the insulation, on the testing voltage.

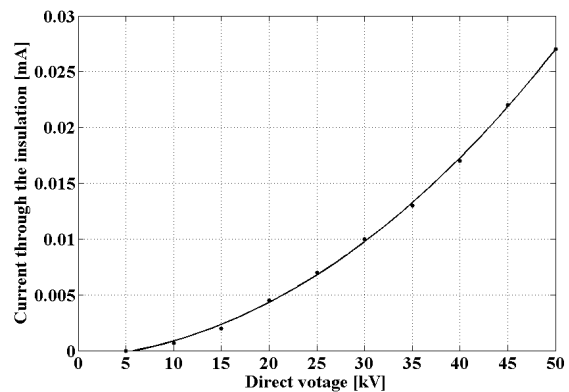


Fig. 6. Dependence of the current through the insulation on the testing voltage of a middle voltage cable AOSB type

4. Numerical study of a middle voltage cable

4.1. Mathematical model

The mathematical model consists of two components (electrical and thermal), related by the electrical conductance, which depends on the temperature. The electric field is described by the following set of equations [5]:

$$\nabla \cdot \mathbf{J} = 0; \quad (1)$$

$$\nabla \cdot \mathbf{D} = \rho; \quad (2)$$

$$\mathbf{J} = \sigma(T) \cdot \mathbf{E} + j\omega \mathbf{D} + \mathbf{J}_e; \quad (3)$$

$$\mathbf{E} = -\nabla V, \quad (4)$$

where: \mathbf{J} is the current density vector; \mathbf{J}_e is the current density of the external sources; \mathbf{D} is the electric flux density; ρ is the volume density of the electric charges; \mathbf{E} is the electric field intensity; $\sigma(T)$ is the electrical conductivity taking into account its dependence on the temperature; V is the electric scalar potential.

The electrical problem is solved with the following boundary conditions:

- electrical potential, corresponding to the testing voltage of the cable, is predetermined along the surface of the current-carrying core;
- electrical potential, equal to 0, is predetermined along the surface of the lead mantle (coating).

The mathematical model for defining the distribution of the thermal field is described by the thermal conductivity equation:

$$-\nabla \cdot (\lambda \nabla T) = q, \quad (5)$$

where: λ is the coefficient of thermal conductivity; q are heat sources in the power cable; T is the temperature.

The thermal problem is solved with the following boundary condition:

- ambient temperature of 20 [°C] is predetermined for the surface of the lead mantle (coating).

4.2. Computer model

Based on the cable construction and on the mathematical model, described above, an axial symmetrical computer model was synthesized in the software program Comsol [6]. The finite element method is used for the analysis of the model and its mesh is shown in Fig. 7. The electrical characteristics of the insulation (dielectric permeability, specific conductivity and $\text{tg}\delta$), and the thermal physical characteristics of the insulation

(the coefficient of thermal conductivity, specific heat capacity and volume density) are in accordance with the technical data of a new cable. The coupled electrical-thermal problem is solved both at direct and at alternating electric field.

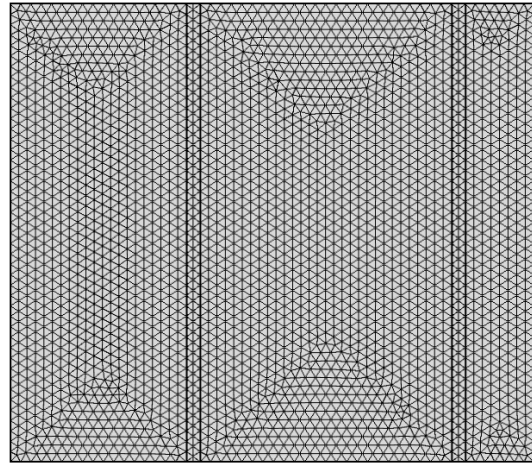


Fig. 7. Finite element mesh of the power cable

4.3. Numerical results at direct electric field

The studies were conducted without current in the current-carrying core at direct voltage of 25 kV. Figs. 8, 9 and 10 show the obtained results for the distribution of the scalar electrical potential, the intensity of the electric field and the specific losses, correspondingly, in a radial cross-section of the cable. Fig. 11 shows the vector distribution of the current density in the insulation of the cable.

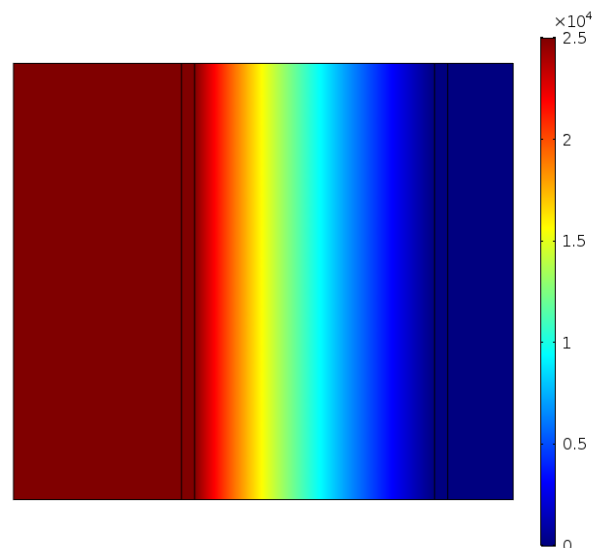


Fig. 8. Distribution of the scalar electric potential [V] in the cable at direct voltage

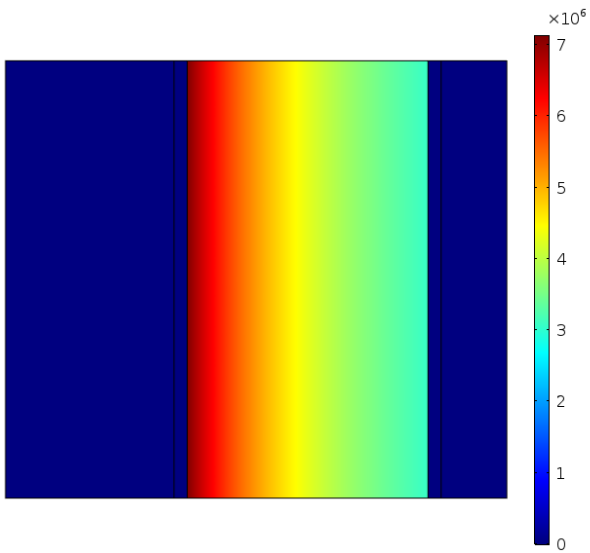


Fig. 9. Distribution of the electric field intensity [kV/m] in the cable at direct voltage

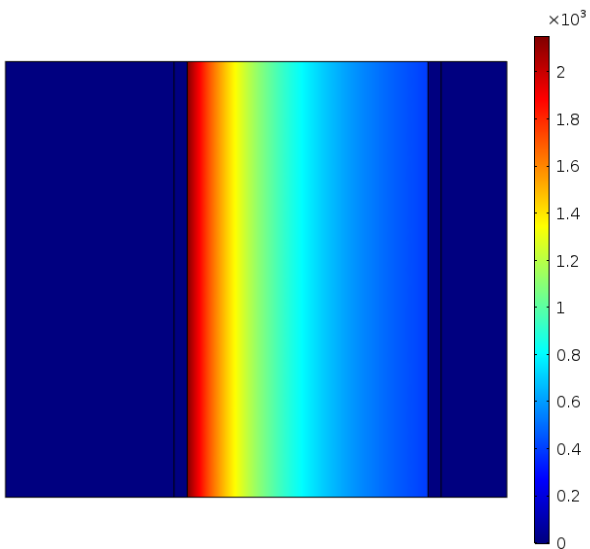


Fig. 10. Distribution of the specific losses [W/m³] in the cable at direct voltage

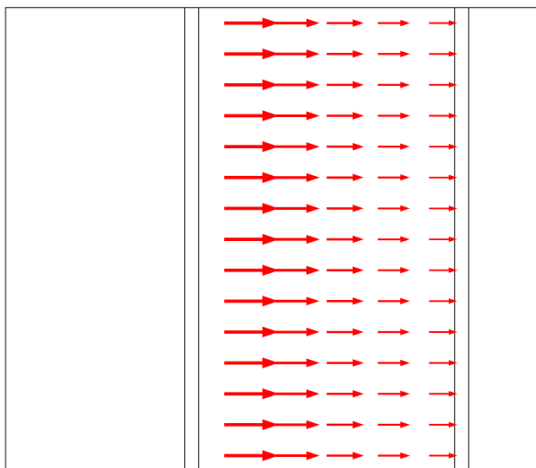
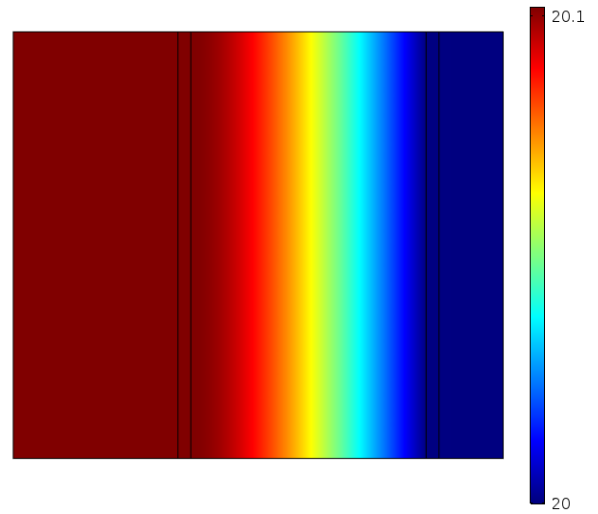
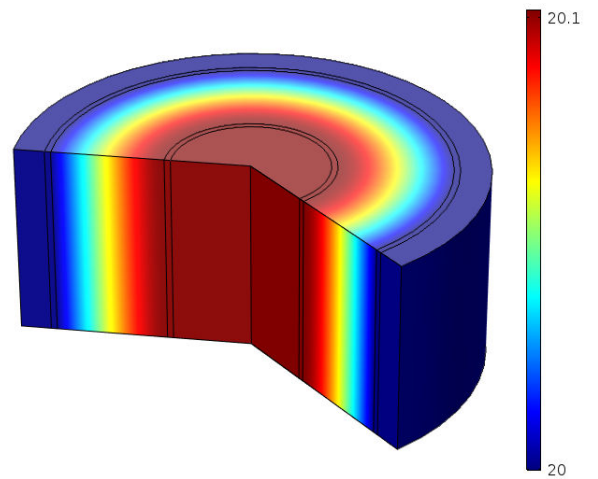


Fig. 11. Direction of the current density in the insulation of the cable at direct voltage

Figs. 12 and 13 illustrate the change of the thermal field in the cable.



a)



b)

Fig. 12. Distribution of the thermal field [°C] in: a) – a radial cross-section of the cable; b) – the volume of the cable

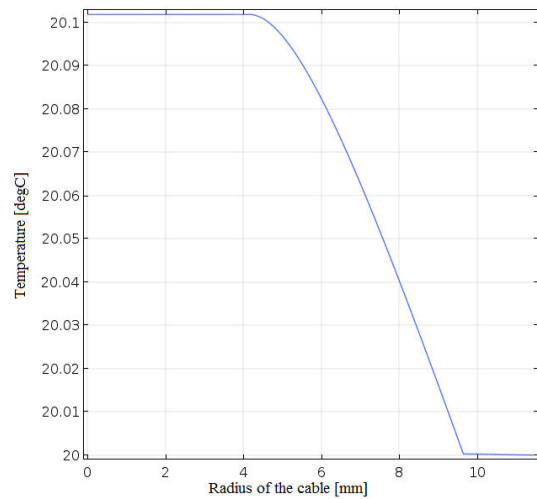


Fig. 13. Change of the temperature along the radius of the cable

4.4. Numerical results at alternating electric field

The studies were conducted without current flow in the current-carrying core at alternating voltage with r.m.s. value of 25 kV and frequency 50 Hz. Fig. 14 and 15 show the obtained results of the distribution of the scalar electric potential and the intensity of the electric field in a radial cross-section of the cable, correspondingly.

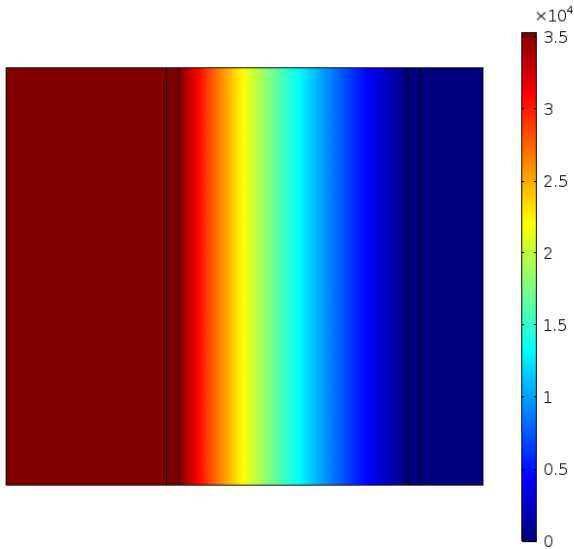


Fig. 14. Distribution of the scalar electric potential [V] in the cable at alternating voltage

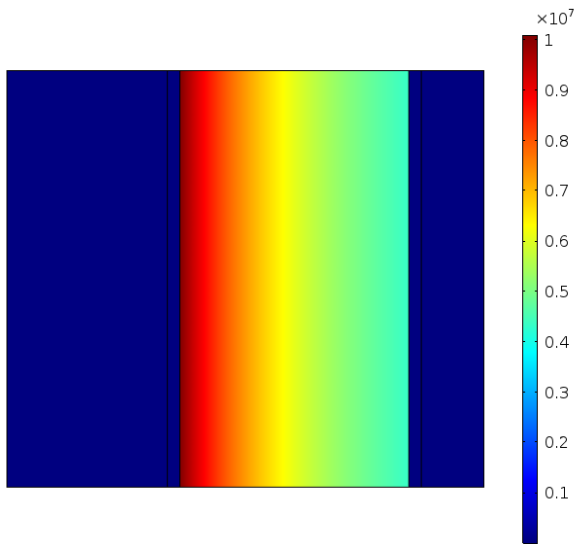


Fig. 15. Distribution of the electric field intensity [kV/m] in the cable at alternating voltage

Fig. 16 and 17 show the specific losses and the thermal field in the cable, respectively.

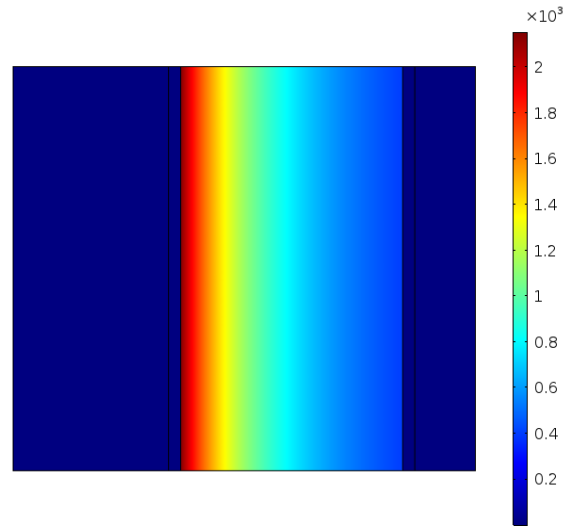
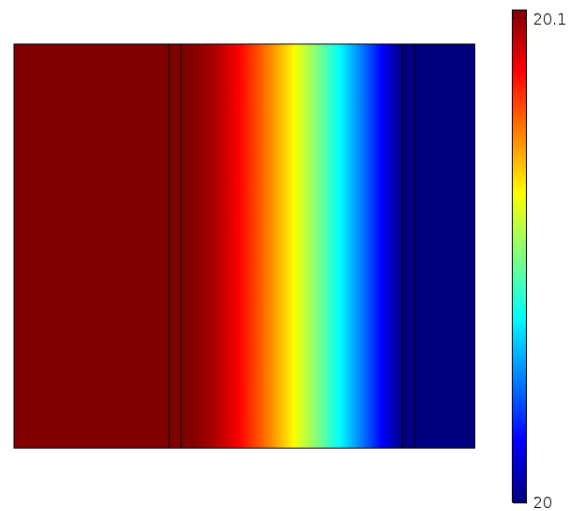
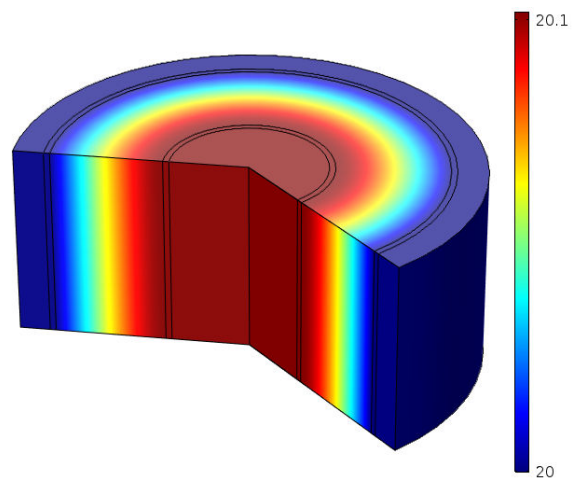


Fig. 16. Distribution of the specific losses [W/m³] in the cable at alternating voltage



a)



b)

Fig. 17. Distribution of the thermal field [°C] in: a) – a radial cross-section of the cable; b) – the volume of the cable

5. Comparison between the obtained numerical and experimental results

Fig. 18 shows the graphical dependencies of the current, flowing through the insulation, on the testing voltage, obtained both from the experiments (curve 1) and from the numerical studies (curve 2).

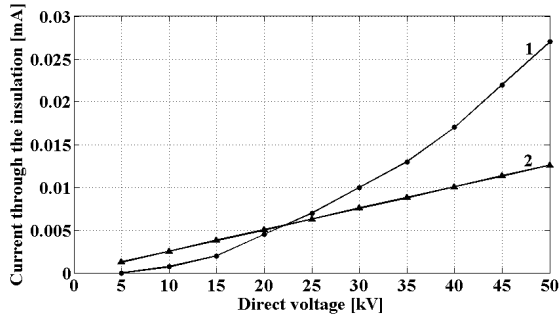


Fig. 18. Dependence of the current through the insulation on the voltage of the AOSB type cable at direct testing voltage

Fig. 19 shows the graphical dependencies of the current, flowing through the insulation, at direct voltage (curve 2) and alternating voltage (curve 1) with frequency 50Hz, obtained from the computer model.

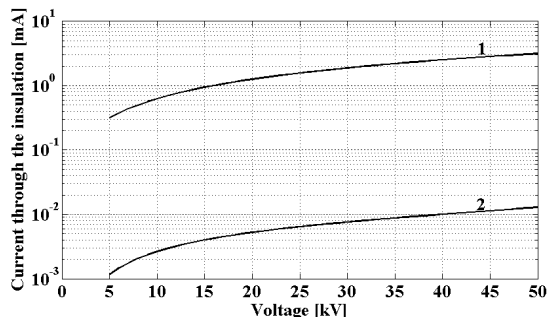


Fig. 19. Dependence of the current through the insulation on the testing voltage for a AOSB type cable

6. Conclusion

Based on the obtained experimental data, the following conclusions can be drawn:

- the highest incident record of failures in the insulation for middle voltage cables belongs to the following types of cables: AOSB, SAHEkT and SAPEkT;
- the main reason for failures in middle voltage cables is their insulation;
- the nonlinear relationship between the current through the insulation and the applied voltage is caused by the ageing of the insulation and deterioration of its parameters.

The current, flowing through the insulation at alternating voltage with frequency 50Hz, is much higher than the current at direct voltage. This requires the use of a powerful source when testing the insulation by means of alternating voltage.

ЛИТЕРАТУРА

1. Popa I., Dolan A. (2016). Numerical modeling of power cables, *SIELA – 19th International Symposium on Electrical Apparatus and Technologies*, pp. 262-265.
2. Li Y., Yongchun Y., Li Y., Yuan P., Li J. (2009). Coupled electromagnetic-thermal modeling the temperature distribution of XLPE cable, *Power and Energy Engineering Conference, APPEEC, Asia-Pacific*, pp. 1-4.
3. Garrido C., Otero F., Cidras J. (2003). Theoretical model to calculate steady-state and transient ampacity and temperature in buried cables, *IEEE Transactions on Power Delivery*, vol. 18, pp. 667-677.
4. BDS 3156:1977 – *Power cables for fixed installation with impregnated paper insulation* (in Bulgarian).
5. Yatchev I., Marinova, I. (2011). *Numerical methods and modeling of circuits and fields, Part one*, Sofia. (in Bulgarian).
6. *COMSOL Version 4.2 User's guide*. (2011).

Assistant Prof. Ivan Hadzhiev, Ph.D.
Department of Electrical Engineering
E-mail: hadzhiev_tu@abv.bg

Assoc. Prof. Dian Malamov, Ph.D.
Department of Electrical Engineering
E-mail: deanmalamov@abv.bg

Assoc. Prof. Vasil Spasov, Ph.D.
Department of Electrical Engineering
E-mail: vasilspasov@yahoo.com

Technical University - Sofia,
Branch Plovdiv
25 Tsanko Dyustabanov Str.
4000 Plovdiv, Bulgaria
Telephone number: +359 32 659686

Eng. Dimitar Nedyalkov, M.Sc.
Elektrorazpredelenie Yug EAD (EP Yug)
37 Hristo G. Danov Str.
4000 Plovdiv, Bulgaria
E-mail: dimitur.nedqtkov@abv.bg

COMPARISON OF PYTHON LIBRARIES USED FOR WEB DATA EXTRACTION

ERDİNÇ UZUN, TARIK YERLİKAYA, OĞUZ KIRAT

Abstract: *There are several libraries for extracting useful data from web pages in Python. In this study, we compare three different well-known extraction libraries including BeautifulSoup, lxml and regex. The experimental results indicate that regex achieves the best results with an average of 0.071 ms. However, it is difficult to generate correct extraction rules for regex when the number of inner elements is not known. In experiments, only %43.5 of the extraction rules are suitable for this task. In this case, BeautifulSoup and lxml, which are the DOM-based libraries, are used for extraction process. In experiments, lxml library yields the best results with an average of 9.074 ms.*

Key words: *HTML, DOM, Web Data Extraction, Python*

1. Introduction

Over the years, due to increased content on the Internet, it has become harder and harder to differentiate between meaningful and unnecessary contents on the web pages. Web data extraction [1] (also known as web content extraction, web scraping, web harvesting, etc.) is the process of extracting user required information from web pages. A large amount of data is continuously created and shared online. Web data extraction systems allow to easily collect these data with limited expert user effort [2]. In this study, we will explain how to access this data from web pages. Moreover, time results of libraries in Python that can be used extraction process are compared.

Web data extraction methods can be classified into three different categories: Wrapper based methods [3], DOM (Document Object Model) based methods [4] and machine learning based methods [5]. Wrapper in this task is a program that extracts data of a particular information from web pages. DOM-based methods utilize structure, tags and attributes of HTML. Machine learning-based methods are on state-of-the-art machine learning algorithm and these methods require labeled data obtained from the DOM. This paper focuses on libraries that can be used for this task. Overall performance in terms of accuracy, time efficiency, memory and processor efficiency varies dramatically depending on the library and algorithms used, even while using libraries classified in the same category.

In this paper, we will discuss time efficiencies of web extraction methods including

regular expressions and two different data extraction libraries in Python programming language called BeautifulSoup [6] and lxml [7], both can be categorized as a DOM based library.

2. Extraction from a Web Page

DOM [8] is an interface that categorizes each element in an HTML or any other XML-based document into nodes of a tree structure. Each node represents a part of the document and can contain other nodes. DOM standard is handled by World Wide Web Consortium, like HTML. The contents of DOM Tree can be changed by programming languages such as JavaScript, C#, Java, Python and etc.

A web page is downloaded after being requested by a web user via a browser. After downloading, the downloaded document is processed and the DOM elements are produced. Each DOM element is interpreted by the web browser to construct the DOM tree while displaying a web page. DOM elements can also be accessed and modified during the display in a web browser of a web page using client-side JavaScript.

There is a relationship between the DOM and the web page as shown in Fig. 1. This web is from www.collinsdictionary.com that presents dictionaries for English or bilingual word reference and plus thesauruses for expanding your word power. In this web domain, `<h2>` is an HTML element that defines the most important heading. Moreover, this element have an attribute (class) and value of attribute (`h2_entry`) that provide additional information for the element.

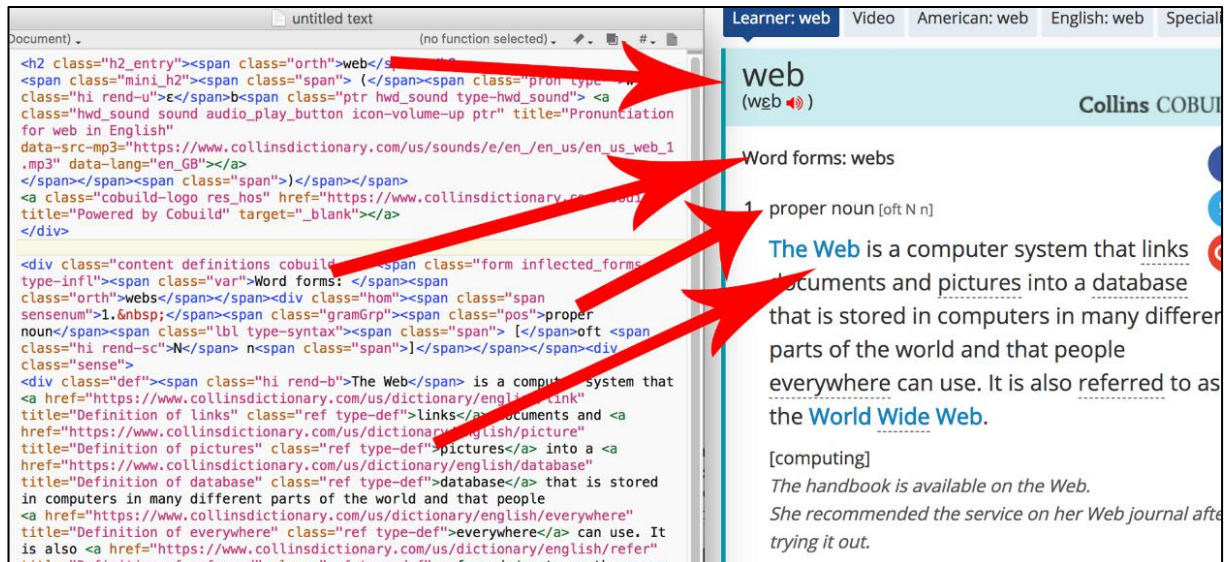


Fig. 1. Relationship between the DOM and the web page

There are a lot of HTML elements such as ``, `<div>`, `<a>`, and etc. in a web page. Fig. 1 shows only a crucial part of a web page that can be used in extraction process. Additionally, these elements have attributes such as `id`, `class`, `title`, `href`, `data-scr-mp3` and etc. `id` and `class` are widely used attributes useful for applying styles and manipulating an element with DOM and JavaScript. These information can be used for extraction process. For example, the following rules can be used in this task.

- `<h2 class="h2_entry">` : Title of a web page
- `<div class="content definitions cobuild br">` : All definitions
- `` : Word forms
- `<div class="hom">` : sub definitions

Web pages in a web domain have similar elements. When the appropriate element is selected from a web page, it can be used to extract other web pages of web domain. For example, these rules can be utilized in the other web pages of this domain. In experiments, we have downloaded 30 different web pages for 10 different web domain. Moreover, we have prepared 1-4 rules for every web domain.

3. Regular Expressions and Data Extraction Libraries for Python

Using regular expressions is the well-known technique that can be used in the extraction process. However, it can cause problems when the number of inner tags is ambiguous. In this situation, the DOM-based libraries can be utilized as a solution for these problems.

3.1. Regular Expressions

Regular expressions (sometimes called as `regex` or `regexp`) is a sequence of characters to define a search pattern. Regular expressions can be handled using the “`re`” module in Python. For this module, we firstly prepared extraction pattern shown in Code 1. For example, this code returns the pattern `<h2. class.h2_entry.>(.*?)</h1>` for an extraction rule of `<h2 class="h2_entry">`.

Code 1. Preparing a pattern for a given rule

```
def prepare_regex(pattern):
    return pattern.replace(" ", ".").replace("\\",
    ".").replace("'", ".").replace("=", ".") + "(.*?)"
    + "</" + parse_TagName(pattern) + ">"

def parse_TagName(element):
    return element[(element.find('<')
    + 1):(element.find(' ')]
```

For an extraction rule, there are two different extraction techniques: the whole document can be searched or the extraction process can be finalized by finding the first record. For some extraction rules, terminating the extraction process after the first extraction can improve the extraction processing time. “`re`” module supports two different extraction techniques. Code 2 finalizes the extraction process after finding the relevant content.

`re`.`search` method scans through string looking for the first location and return a corresponding match object. If no position in the string matches the pattern, Code 2 returns empty string. However, if more than one record exists, the entire document should be looked at. Code 3 extracts all content from a web page.

Code 2. Extracting the first record with re

```
import re
```



```
def extract(html, pattern):
    res = re.search(prepare_regex(pattern), html,
re.DOTALL)
    if res:
        return res.group(1)
    else:
        return ''
```

Code 3. Extracting all records with re

```
def extract_all(html, pattern):
    return re.findall(prepare_regex(pattern), html,
re.DOTALL)
```

re.findall method returns all non-overlapping matches of pattern in a web page, as a list of strings. In experiments, the effect of these extraction techniques will be investigated.

3.2. Data Extraction Libraries in Python

Two well-known extraction libraries, BeautifulSoup and lxml, are used for this task.

3.2.1. BeautifulSoup

BeautifulSoup is a Python data extraction library developed by Leonard Richardson and other open source developers. It is licensed under the Simplified BSD License and works on both Python 2.7+ and Python 3. It can parse HTML and XML documents and provides simple methods to interact with the DOM model.

Code 4. BeautifulSoup extraction methods

```
from bs4 import BeautifulSoup
def extract(html, pattern, parser):
    soup = BeautifulSoup(html, parser)
    return soup.find(parse_TagName(pattern),
    attrs=parse_Attributes(pattern,
    parser)).decode_contents(formatter="html")
def extract_all(html, pattern, parser):
    soup = BeautifulSoup(html, parser)
    temp_res =
    soup.find_all(parse_TagName(pattern),
    attrs=parse_Attributes(pattern, parser))
    html_res = []
    for res in temp_res:
        html_res.append(res.decode_contents(formatt
er="html"))
    return html_res
def parse_TagName(element):
    return element[(element.find('<') +
1):(element.find(' '))]
def parse_Attributes(element, parser):
    soup = BeautifulSoup(element, parser)
    e = soup.find(parse_TagName(element))
    return e.attrs;
```

BeautifulSoup supports two different extraction techniques like “re” module. Code 4 extracts only the first record in web page.

A BeautifulSoup object has two arguments. The first argument is the source of web page, and the second argument is the parser. The different parsers including html.parser, lxml, and html5lib can be adapted to the object of BeautifulSoup. html.parser comes with Python’s standard installation and provides a class named HTMLParser which can be used as a basic HTML and XHTML parser. lxml is feature-rich and easy-to-use library for processing XML and HTML. It provides Python bindings for the C libraries libxml2 and libxslt and mostly compatible with ElementTree API. It is also open source and licensed under the BSD license. html5lib conforms WHATWF HTML specification which allows the module to parse HTML content the same way a web browser does. It is mainly developed by James Graham and open source under the MIT license.

The find method uses when you only want to find one result. The find_all method scans the entire document looking for results. If the number of extraction result is one content, you can use the find() method for improving time efficiency of the extraction process. In Code 4, the parse_TagName method finds the element name of the pattern and parse_Attributes method returns all attributes and their values in list format. The decode_contents method renders the contents of this element in html format.

3.2.2. Lxml

Some web sites introduced BeautifulSoup [9] recommend to install and use lxml for speed. But also, it can be used stand-alone. lxml supports XPath [10] for extracting the content of a tree. XPath allows you to extract the content chunks into a list. XPath uses "path like" syntax to identify and navigate nodes in an html and xml document. Code 5 returns the XPath expression for a given extraction rule.

Code 5. Preparing of XPath expression for a given element

```
def prepare_XPath(pattern):
    root = etree.fromstring(pattern + '<' +
    parse_TagName(pattern) + '>')
    temp=""
    for att in root.keys():
        temp += '['@'+ att + '=' + root.get(att) +
        "]"
    return ".//" + root.tag + temp
```

For example, Code 5 returns the XPath expression `./h1[@class="h2_entry"]` for an

extraction rule of `<h2 class="h2_entry">`. Code 6 has two function for extracting all results and the first result.

Code 6. Extraction with lxml

```

from lxml import etree
from io import StringIO
def extract_all(html, pattern):
    parser = etree.HTMLParser()
    tree = etree.parse(StringIO(html), parser)
    result = etree.tostring(tree.getroot())
    root = tree.getroot()
    my_list = []
    yol = prepare_XPath(pattern)
    for elem in root.findall(yol):
        my_list.append(etree.tostring(elem,
            encoding='unicode'))
    return my_list
def extract(html, pattern):
    parser = etree.HTMLParser()
    tree = etree.parse(StringIO(html), parser)
    result = etree.tostring(tree.getroot())
    root = tree.getroot()
    elem = root.find(prepare_XPath(pattern))
    return etree.tostring(elem,
        encoding='unicode')

```

In Code 6, find method efficiently returns only the first match. findall method returns a list of matching Elements.

Table 1. Information about dataset

Domain	Category	Avg. (KB)
Aliexpress	Shopping	93.66
Bild	Newspaper	67.49
Booking	Trip	689.11
Collinsdictionary	Dictionary	38.93
Ebay	Shopping	297.35
Imkb	Movie	214.23
Sciencedirect	Articles	48.65
Tchibo	Shopping	27.61
Tutorialspoint	Articles	28.56
W3schools	Articles	58.01
		153.36

4. Experiments

We prepare a dataset which contains web pages on many different content types, including scientific articles, dictionary, movies, newspaper articles, shopping, and trip/hotel information. For constructing this dataset, we have designed a simple crawler to download web pages. Then, this crawler downloads 30 web pages for every domain. Table 1 gives the average file size of web domains for this

dataset. Moreover, we prepare extraction rules (like in Section 2) for every domain. All experiments are carried out on a computer using Intel Core i5-3.2Ghz processor and 8 GB RAM running Windows 10 operating system.

For measuring extraction time of these methods, we use time.clock method. This method returns wall-clock seconds elapsed, as a floating point number. This method is based on the Win32 function QueryPerformanceCounter.

4.1. Time results of Regex

There are two techniques in regular expressions. Table 2 indicates the extraction time results and whether the result of the extraction is correct or not.

Table 2. Time results and accuracy of regex

Method	Avg. (ms)
extract	0.071
extract_all	0.307
	0.189
Accuracy: 390 / 600 = %43.5	

As expected, focusing only on one result rather than looking at the entire document has yielded better results. Extract method is 4.324 times faster than Extract_All method. However, 43.5% of the extraction rules give the correct result. In this case, DOM-based methods can be considered as a solution.

4.2. Time results of BeautifulSoup

BeautifulSoup supports two different extraction techniques and three different parsers for this task. Table 3 indicates the extraction time results.

Table 3. Time results of BeautifulSoup

Method	Parser	Avg. (ms)
extract	html.parser	820.075
extract_all	html.parser	1192.317
extract	lxml	591.583
extract_all	lxml	1025.892
extract	html5lib	2191.472
extract_all	html5lib	2626.747
		1408.014

lxml parser for BeautifulSoup is the best parser for this task. html5lib and html.parser are just not very good results. As expected, the time results of extract methods are better than the time results of extract_all methods in all parsers. Finally, we examine lxml parser with its methods.

4.3. Time results of lxml

The lxml library is a binding for the C libraries libxml2 and libxslt. It is unique in that it combines the speed and XML feature completeness of these libraries with the simplicity of a native Python API. It can be used with BeautifulSoup, but it can be also utilized stand-alone. In this section, we examine this library stand-alone.

Table 4. Time results of lxml

Method	Avg. (ms)
extract	9.047
extract_all	9.480
	9.277

Time results of lxml is better than time results of BeautifulSoup. lxml parser in BeautifulSoup makes the results better, but the use of lxml stand-

alone provides a much better improvement. (See Fig. 2 and 3 for comparing all libraries)

5. Conclusion

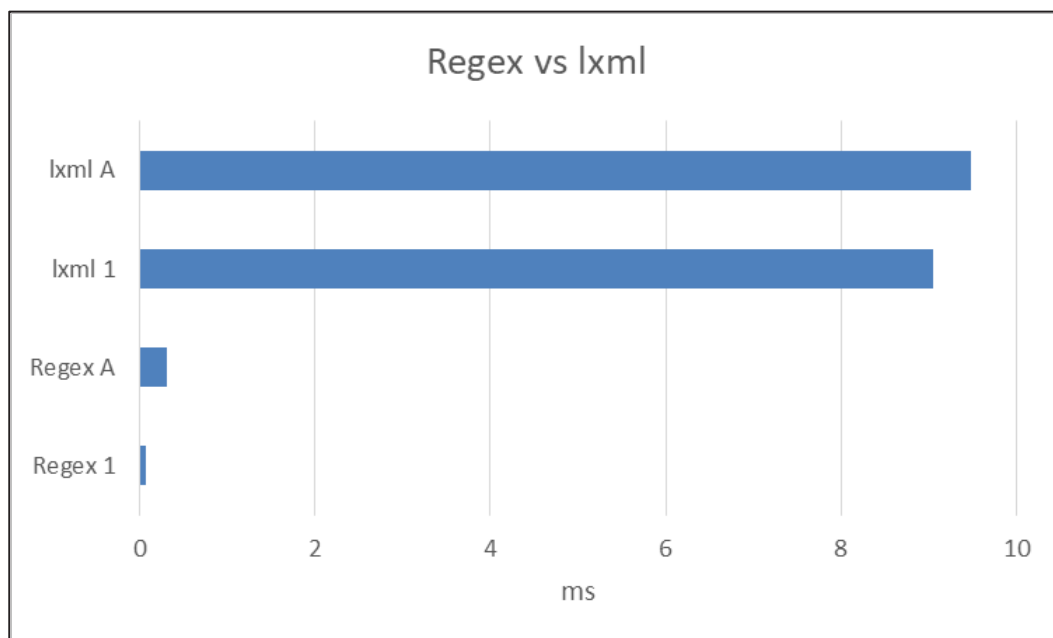
In this study, libraries of Python for extracting data from web pages are compared in order to understand their time durations. As expected, the experimental results show that regex gives better time duration with 0.071 ms. However, 43.5% of the extraction rules give the correct result. In this case, DOM-based libraries including BeautifulSoup and lxml can be considered as a solution. Lxml (stand-alone) provides better time results in DOM-based libraries. BeautifulSoup gave worse results because of extra processes for creating DOM even when using lxml parser.

In future work, we will need to develop more effective and effective methods for this task. Moreover, time results of other languages [11] on this task are compared.

Additional Information

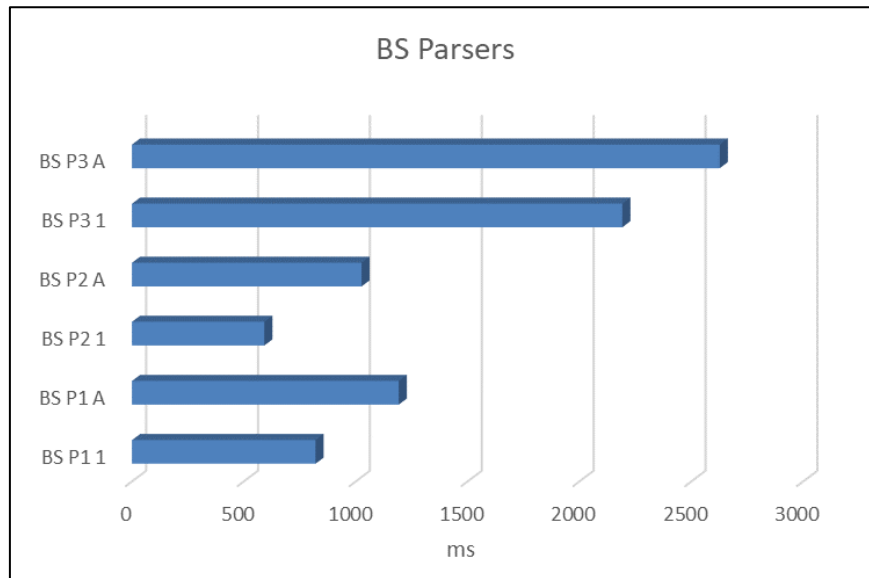
All codes are open-source. Source codes and dataset are as follows.

- <https://github.com/erdincuzun/WebDataExtractionInPython>



A: All extraction results, 1: only first result

Fig. 2. Average time duration results of Regex and lxml



A: All extraction results, 1: only first result, BS: BeautifulSoup, P1: html.parser, P2: lxml, P3: html5lib

Fig. 3. Average time duration results BeautifulSoup Parsers

REFERENCES

- Rahman, A.F.R., Alam, H. and Hartono, R. (2001). Content extraction from HTML documents, *International workshop on Web document Analysis*, pp. 7-10, 2001.
- Ferrara, E., De Meo, P., Fiumara, G., Baumgartner, R. (2014). Web data extraction, applications and techniques: A survey, *Knowledge-Based Systems*, Volume 70, 2014, pp. 301-323.
- Flesca, S., Manco, G., Masciari, E., Rende, E., Tagarelli, A. (2004). Web wrapper induction: a brief survey”, *In: AI Communications*, vol. 17, pp. 57–61. IOS Press, Amsterdam.
- Álvarez-Sabucedo, L. M., Anido-Rifón, L. E. and Santos-Gago, J. M. (2009). Reusing web contents: a DOM approach, *Softw. Pract. Exper.*, 39: 299–314. doi:10.1002/spe.901.
- Fu, L., Meng, Y., Xia Y. and Yu, H. (2010). Web Content Extraction based on Webpage Layout Analysis”, *Second International Conference on Information Technology and Computer Science*, Kiev, 2010, pp. 40-43.
- BeautifulSoup, <https://www.crummy.com/software/BeautifulSoup/>, (12.04.2018)
- lxml, <http://lxml.de/>, (12.04.2018)
- DOM, https://developer.mozilla.org/en-US/docs/Web/API/Document_Object_Model/Introduction, (12.04.2018)
- Crummy, <https://www.crummy.com/software/BeautifulSoup/bs4/doc/>, (12.04.2018)
- XPath, <https://www.w3.org/TR/xpath/>, (12.04.2018)
- Uzun, E., Buluş, H. N., Doruk, A., Özhan, E. (2017). Evaluation of Hap, AngleSharp and HtmlDocument in web content extraction. *International Scientific Conference'2017 (UNITECH'17)*, Gabrovo, Bulgaria, November 17-18, Vol. II – pp. 275-278.

Authors' contacts

Erdinç Uzun,
 Organization: Namık Kemal University,
 Çorlu Faculty of Engineering, Computer
 Engineering Department
 Address: NKÜ Çorlu Mühendislik Fakültesi
 Dekanlığı, Silahtarğa Mahallesi
 Üniversite 1.Sokak, No:13, 59860 Çorlu /
 Tekirdağ / TURKEY
 Phone (optional): +90 (282) 250 2325
 E-mail: erdincuzun@nku.edu.tr

Tarık Yerlikaya (Corresponding author)
 Oğuz Kırat
 Organization: Trakya University, Faculty of
 Engineering, Computer Engineering
 Department
 Address: Trakya Üniversitesi Ahmet
 Karadeniz Yerleşkesi Mühendislik
 Fakültesi 22020 Merkez / Edirne
 /TURKEY
 Phone (optional): +90 (284) 226 1217 /
 2215
 E-mail: tarikyer@trakya.edu.tr,
 ogzkirat@gmail.com

OBJECT-BASED ENTITY RELATIONSHIP DIAGRAM DRAWING LIBRARY: ENTREL.JS

ERDİNÇ UZUN, TARIK YERLİKAYA, OĞUZ KIRAT

Abstract: *An entity relationship diagram (ERD) is a visual helper for database design. ERD gives information about the relations of entity sets and the logical structure of databases. In this paper, we introduce an open source JavaScript Library named EntRel.JS in order to design sophisticated ERDs by writing simple codes. This library, which we have developed to facilitate this task, is based on our obfc.js library. It generates a SVG output on the client side for modern browsers. The SVG output provides storage efficiency when compared to existing ERD drawings created with existing drawing applications. Moreover, we present our animation library to gain action for elements in your web page.*

Key words: *Entity relationship diagrams, SVG, JavaScript*

1. Introduction

The internet has been evolving for over 25 years with some standards including HTML (HyperText Markup Language), CSS (Cascading Style Sheets), JavaScript, XML (eXtensible Markup Language) and so on. These standards are shared by the W3C throughout the world, and browser developers release updates to support these standards. The W3C [1] continues to set standards on a number of different issues. In this study, one of these standards is SVG (Scalable Vector Graphics) [2] that is used for drawing entity relationship diagrams (ERDs) and JavaScript language is used for creating these ERDs with simple codes.

ER modeling was developed for increasing the clarity of database design by Peter Chen[3]. This model is the result for systematic data analysis of system. It is usually drawn in a graphical form named ER diagrams including entities, their attributes and relationships between entities. This model is typically implemented as a database. Software developers and database designers can easily discuss the design of database and software over the ER diagram. Therefore, this subject is the basis of the database courses. ER design can be created very quickly with our library introduced in this study.

The ERD drawing process can be done on the server or client side by using the drawing applications including Microsoft Visio, OpenOffice / LibreOffice Draw, Dia, Diagramly and so on. The drawing file stored on the server side is displayed in the browser in the element. On the other hand, it can still be prepared on the server side and

displayed in the CANVAS [4] or SVG elements used for drawing graphics in HTML 5. In this study, the SVG is suitable for our study instead of CANVAS. CANVAS is typically used in web-based games, while it allows the drawing process with JavaScript codes. Because SVG is a vector-based system, the browser is not affected by the magnification. It was chosen for this reason. It can also be drawn on the client side by using JavaScript's advantage and holding fewer code on the server side. This makes the server load lighter and the drawing works are made on the client side. However, with the JavaScript file loaded once, it is only necessary for utilizing the AJAX (asynchronous JavaScript and XML) [5] in order to upload the required code. AJAX improves the bandwidth performance of web-based applications. In this study, the object creation codes are stored in the server, and drawing with the EntRel.JS library is performed on the client side.

The EntRel.JS library (Entity Relationship) is based on the obfc.JS library (Object-Based Flow Charts) [6] that we develop is an object-based library for drawing SVG flow charts across modern web browsers. It makes easily to construct objects, links and connections. Moreover, it dispatches a click event when an object or a line is clicked and descriptions can be added for all clicks. The EntRel.JS allows you to design complex ERDs by writing short codes. SVG output is parsed from these texts and the SVG output is produced on the client side.

Technologies such as SVG, JavaScript and AJAX are used in many different subjects such as numerical graphics, networking, geography,

medicine and electricity. For example, Saito and Ouyang [7] indicate how to draw data on the client side with ChartML. They use SVG and JavaScript in the client side and they produce graphs. Some studies [8-11] focus on how network topologies can be displayed on the browser by using SVG, JavaScript and AJAX. Yin and Zang [12] describe SVG and AJAX technologies in the Web geographic Information system. Fang and et al. [13] explain the use of these technologies in local thermal power plant management. Birr and et al. [14] introduces how three-dimensional medical data can be demonstrated in the Web environment. Alhirabi and Butler [15] perform the gene notation using SVG.

This study presents a JavaScript Library to create ERDs with objects. Moreover, you can also easily link objects to each other. Draw function of these objects produces a SVG output in browsers. Moreover, we introduce an animation library (animation.JS) that we developed.

2. obfc.JS library

Before explaining the EntRel.JS library, we give information about obfc.JS library that we developed earlier. obfc.JS has 24 different SVG shapes for drawing flow charts. In this section, we will describe the most basic features of this library.

2.1. Creating an object

Before creating an object, you add obfc.js file and SVG element to body of a web page. Then, you connect library to the id of the SVG element. Code 1 indicates these lines.

Code 1. Preparing library and creating an object

```
<script src="obfc.min.js"></script>
<svg id="demo" width="600" height="700">
</svg>
<script>
prepare_SVG("demo");
var object2 = add_theObject(new Process(300,
150, 1, ["Line 1", "Line 2"], 10));
</script>
```

To draw an object, add_theObject function can be used for the given SVG element. add_theObject is a function that adds the object into a given SVG element and returns this object. This returned object is used for drawing lines between two objects. There are 24 different SVG objects in obfc.JS. "Process" is the one of objects from 24 different objects. There are 9 parameters for creating an object. First two parameters are required, others are optional.

```
Object_Name(_middle_x, _middle_y, _size,
_text, _text_size, _description, _fill_color,
_stroke_color, _text_color);
```

- `_middle_x` and `middle_y`: centre of the object. (Required)
- `_size`: For example, default size of a process object width is 125 and height is 50. `_size` value is multiplied by these values.
- `_text` and `_text_size`: Text is written in the center of an object. This parameter can be defined string value or array["", ""...]. If your text is too long, you can use array for creating lines.
- `_description`: These value can be coded in HTML format. These value is displayed in a HTML element that contains "desc" id after clicking an object or a line.
- `_fill_color`: Default value is white. But, the object color can be determined with this parameter.
- `_stroke_color`: Default value is black.
- `_text_color`: Default value is black.

2.1. Creating lines between objects

After creating all objects, objects can be linked by using draw_theLine function.

Code 2. Creating lines between two objects

```
<script>
....
var o_line1 = draw_theLine(new Line(object1,
object2));
</script>
```

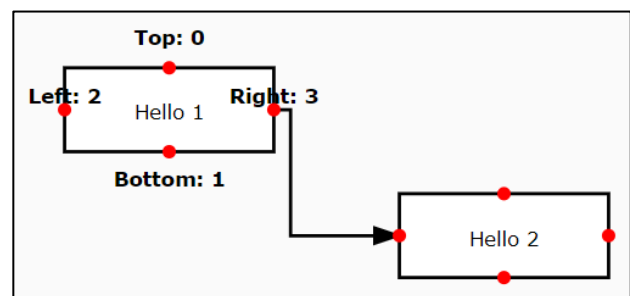


Fig. 1. Output of Code 2

Line is a function that determines the path for given two objects and their positions. This function has 9 parameters. First two parameter is required and others are optional.

```
Line(object1, object2, position1, position2,
_text, _text_size, _description, _stroke_color,
_text_color);
```

- `object1` ve `object2`: are variables that is defined in the previous section.

- `position1` and `position2`: are position information of objects. There are four positions for all shapes. Top=0, Bottom=1, Left=2 and Right=3. But, when these values are not entered or entered “-1”, this function automatically determines these position by calculating differences between all unused positions. (Unused position means that this position is used for creating lines)
- `_text` and `_text_size`: text in line. This function selects the longest sub-line for writing text.
- `description`, `_stroke_color`, `_text_color`: (same with previous section)

Moreover, you can determine connection points manually as shown in Code 3.

Code 3.

```
<script>
var object1 = add_theObject(new Process(100,
75, 1, "Hello 1", 12));
var object2 = add_theObject(new Process(300,
150, 1, "Hello 2", 12));
var o_line1 = draw_theLine(new Line(object1,
object2, 0, 1));
var o_line2 = draw_theLine(new Line(object1,
object2, 2, 2));
var o_line3 = draw_theLine(new Line(object1,
object2, 1, 3));
</script>
```

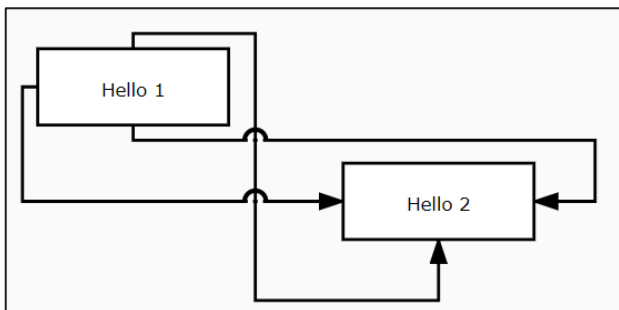


Fig. 2. Output of Code 3

In this example, position informations are added to link two objects. For example, 0 means top of object1 and 1 means bottom of object2. There are three lines in this example. Moreover, obfc.js has

jumping mechanism in collision of lines. For more information, you can visit the following web page:

<https://www.e-adys.com/obfc-js/object-based-flow-charts-obfc-js/>

3. EntRelJS

EntRelJS is JavaScript Library for creating ERDs with simple JavaScript methods. Figure 3 gives information about the location of the libraries in the Web. Libraries of obfc.JS, EntRel.JS and Animation.JS takes the drawing code from a server and parse these codes for drawing into an SVG element. The file sizes of obfc.JS, EntRel.JS and Animation.JS are 78 KB, 3.85 KB and 3.62 KB. The first installation of these libraries is loaded into the browser cache, and these libraries are not updated on other requests.

There are two main classes for creating diagrams as FlatLine and Entity. Entity class derived from Process class of obfc.js. For better understanding this object, you can examine Section 2.1.

```
Entity(_middle_x, _middle_y, _size, _text,
_attributes, _text_size, _weak, _description,
_fill_color, _stroke_color, _text_color)
```

The fifth parameter contain extra information about attributes of an entity. You can define 8 attribute in array format for an entity. For example:

```
[null, "AuthorID(PK)", "AuthorName",
"AuthorSurName"]
```

The first attribute contain a null value so that it can be drawn. The first four attributes are in top and others are in bottom. Weak parameter is used for a weak entity that cannot be uniquely identified by its attributes alone.

FlatLine connects the objects. FlatLine class derived from Line class of obfc.js. FlatLine don't contain an arrow in drawing. Code 4 is an example for Entity and FlatLine. Moreover, Fig 4 is output of Code 4.

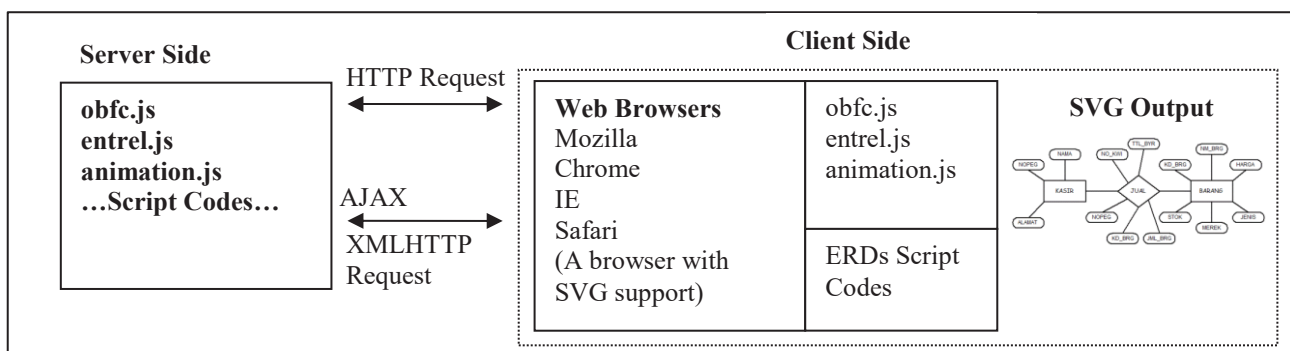


Fig. 3. Location of the *EntRel.JS* and *Animation.JS* in the Web**Code 4.** An example ERD

```

<script>
prepare_SVG("demo");
var object1 = add_theObject(new Entity(300,
350, 0.75, "Books", [null, null, null, null, null,
"BookID(PK)", "Title"], 16));
var object2 = add_theObject(new Entity(65, 150,
0.75, "Authors", [null, "AID(PK)", "AName",
"ASurname"], 16));
var object3 = add_theObject(new Entity(300,
150, 0.75, "Types", [null, "TID(PK)", "TName"],
16));
var object4 = add_theObject(new Entity(500,
150, 0.75, "Publishers", [null, "PID(PK)",
"PName", "Location"], 16));
var object5 = add_theObject(new Decision(65,
250, 0.75, ["R1"], 12));
var object6 = add_theObject(new Decision(300,
250, 0.75, ["R2"], 12));
var object7 = add_theObject(new Decision(500,
250, 0.75, ["R3"], 12));

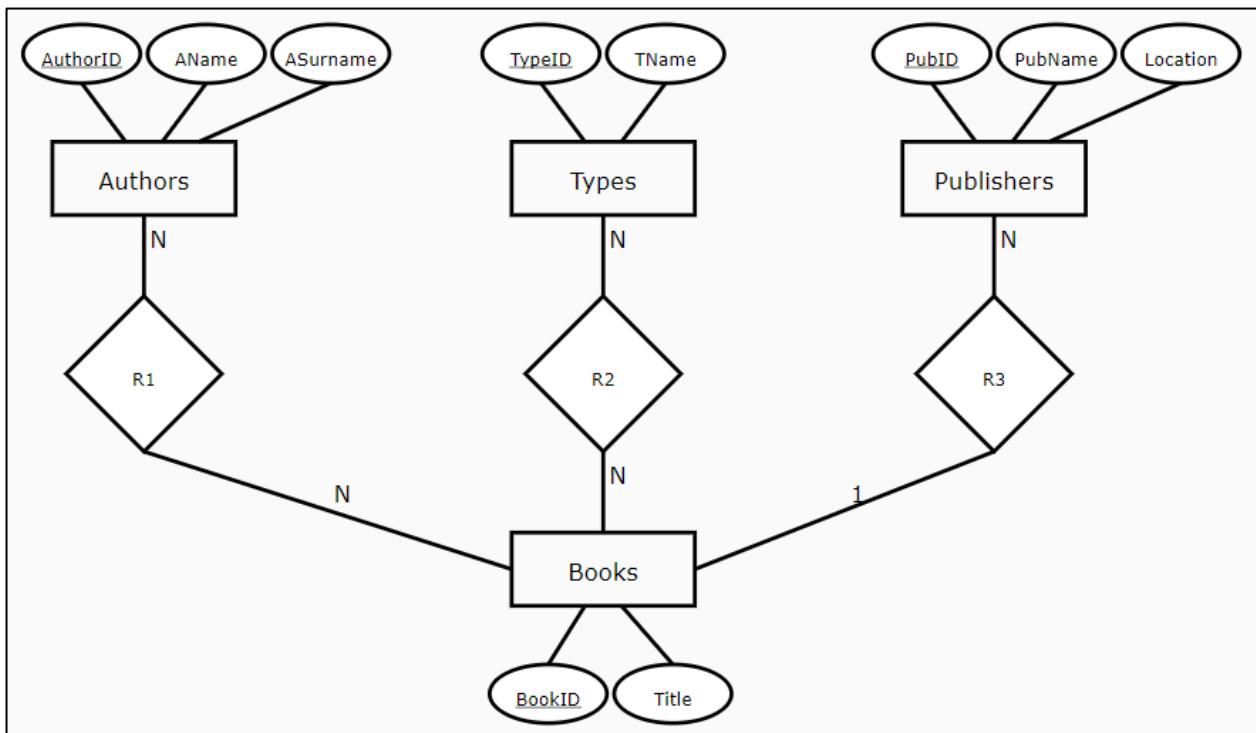
var o_line1 = draw_theLine(new
FlatLine(object2, object5, null, null, "N"));
var o_line2 = draw_theLine(new
FlatLine(object5, object1, 1, 2, "N"));
var o_line3 = draw_theLine(new

```

```

FlatLine(object3, object6, null, null, "N"));
var o_line4 = draw_theLine(new
FlatLine(object6, object1, null, null, "N"));
var o_line5 = draw_theLine(new
FlatLine(object4, object7, null, null, "N"));
var o_line6 = draw_theLine(new
FlatLine(object7, object1, 1, 3, "1"));
</script>

```

**Fig. 4.** Output of Code 3

In Code 4, (PK) keyword defines a primary key of an entity. Moreover, (DE) keyword defines

derived attributes and (MV) keyword is for multivalued attributes. There are four entity objects

and three Decision objects created from obfc.JS library. And these objects link with six FlatLine. Now let's add animation to this drawing.

4. animation.JS

With animation.js, you can easily animate your a desired element of child elements in your web page. Navigation bar of this animation contains links: Previous, Next, Show All, Hide All, Start Animation, and Stop Animation. Code 5 is a simple animation code.

Code 5. Creating an animation a web page

```
<script src="animation.min.js"></script>
<script>
initializeAnimation(null, ".animation", "div,p",
"#masthead,#secondary");
</script>
```

- Add animation.min.js your web page, then call initializeAnimation method for configuration of the animation.
- The first parameter of this method is the number of elements. But this parameter is used by other libraries: obfc.js and entrel.js.
- The second parameter is the base element of animation. You can use a selector for setting this parameter.
- Third parameter is inner elements in the base element. You can use more than one selector for selecting the desired elements.
- Fourth parameters is used to set opacity of the selected elements.

Moreover, you can set time interval for animation and opacity for elements. For example: (Append the following code to Code 5)

```
opacity = 0.1; //opacity of elements
timeInterval = 3000; //3 seconds
```

For creating links for navigating animation in a web page, you can append Code 6 for your web page.

Code 6. Navigation bar

```
<div id="anavbar" class="anavbar"
style="display:none">
<ul class="horizontal">
<li><a id="pre">Previous</a></li>
<li><a id="next">Next</a> </li>
<li><a id="show">Show All</a> </li>
<li><a id="hide">Hide All</a> </li>
<li><a id="start">Start Animation</a> </li>
<li><a id="stop">Stop Animation</a> </li>
</ul>
</div>
```

5. The use of obfc.JS and EntReel.JS with animation.JS

For creating animations, you can group the objects created by obfc.JS and EntReel.JS using an array. If two objects are together, these objects can be grouped together. (For example ([object5, object2])). You can append the following codes to Code 4 for creating animation in Fig 4.

```
var groups = [object1, [object5, object2],
[o_line1, o_line2], [object6, object3], [o_line3,
o_line4], [object7, object4], [o_line5, o_line6]];
prepareClassforAnimation(groups);
initializeAnimation(groups.length - 1);
```

prepareClassforAnimation method prepares objects for animation. initializeAnimation method starts animation. For testing animation, you can visit the following web page:

https://www.e-adys.com/web_tasarimi_programlama/entrel-js-creating-entity-relationship-diagrams-with-javascript-and-svg/

6. Conclusion

In this study, the open source EntRel.JS and animation.JS Libraries that we developed and the functions of these libraries to draw ERDs in the Web environment are introduced. Thanks to this library, drawings are made quickly with very little code. Unlike other drawing applications, the links between shapes are automatically feasible. SVG output occupies very little space according to the other image formats.

In future studies, we plan to develop a design library that enables drawing and drag-and-drop functionality through the Web page without writing codes. We are also aiming to group some designs and share them in a web environment. Finally, we intend to develop new libraries for different subjects including logic circuits, data structures, database management systems, software engineering and system analysis in order to draw different shapes.

Additional Information

All codes are open-source. Web addresses and help documents are as follows.

- <https://github.com/erdincuzun/obfc.js>
- <https://github.com/erdincuzun/entrel.js>
- <https://github.com/erdincuzun/Animation.js>
- <https://www.e-adys.com/>

REFERENCES

1. World Wide Web Consortium, <https://www.w3.org/>, (14.06.2018)
2. Scalable Vector Graphics, <https://www.w3.org/Graphics/SVG/>, (14.06.2018)
3. Chen, P. (1976). The Entity-Relationship Model - Toward a Unified View of Data. *ACM Transactions on Database Systems*. 1 (1): 9–36.
4. HTML 5 Canvas, <https://www.w3.org/TR/2dcontext/>, (14.06.2018)
5. AJAX, <https://developer.mozilla.org/en-US/docs/Web/Guide/AJAX>, (14.06.2018)
6. Uzun, E., Buluş, H. N. (2017). Object-based flowchart drawing library. International Conference on Computer Science and Engineering (UBMK 2017), Antalya, Turkey, 5-8 Oct. 2017, pp. 110-115.
7. Saito, T. and Ouyang, J. "Client-side data visualization," 2009 IEEE International Conference on Information Reuse & Integration, Las Vegas, NV, 2009, pp. 194-199. doi: 10.1109/IRI.2009.5211550.
8. Valle, R. D. T., Passos, D., Albuquerque, C. and Muchaluat Saade, D. C. (2008). Mesh Topology Viewer (MTV): an SVG-based interactive mesh network topology visualization tool. *2008 IEEE Symposium on Computers and Communications*, Marrakech, pp. 292-297.
9. Lin, T., Zou, F., Kienle, H. M. and Muller, H. A. (2008). A domain-customizable SVG-based graph editor for software visualizations. *2008 IEEE International Conference on Software Maintenance*, Beijing, 2008, pp. 466-467.
10. Fan, C., Wu, Y. and Wang, F. (2009). SVG based on Ajax and its application in graphical network topology management. *2009 IEEE International Conference on Communications Technology and Applications*, Beijing.
11. Kehe, W., Tingting, W., Yanwen, A. and Wenjing, Z. (2015). Study on the Drawing Method of Project Network Diagram. *2015 7th International Conference on Intelligent Human-Machine Systems and Cybernetics*, Hangzhou, 2015, pp. 95-98.
12. Yin, F. and Zhang, L. (2010). Research of WebGIS based on SVG and Ajax technology. *2nd IEEE International Conference on Information and Financial Engineering*, Chongqing, pp. 629-632.
13. Fang, W., Zhang, J., Hu, B., Zhang, Q. and Ha, X. (2011). Graphics and data web publishing for local thermal power plant management information system. *2011 International Conference on Multimedia Technology*, Hangzhou, 2011, pp. 337-340.
14. Birr, S., Mönch, J., Sommerfeld, D., Preim, U. and Preim, B. (2013). The LiverAnatomyExplorer: A WebGL-Based Surgical Teaching Tool. in *IEEE Computer Graphics and Applications*, vol. 33, no. 5, pp. 48-58.
15. Alhirabi, N. and Butler, G. (2015). A visual spreadsheet using HTML5 for whole genome display. *2015 IEEE Conference on Computational Intelligence in Bioinformatics and Computational Biology (CIBCB)*, Niagara Falls, ON, 2015, pp. 1-7.

Authors' contacts

Erdoğan Uzun,
 Organization: Namık Kemal University,
 Çorlu Faculty of Engineering, Computer
 Engineering Department
 Address: NKÜ Çorlu Mühendislik Fakültesi
 Dekanlığı, Silahtarağa Mahallesi
 Üniversite 1.Sokak, No:13, 59860 Çorlu /
 Tekirdağ / TURKEY
 Phone (optional): +90 (282) 250 2325
 E-mail: erdincuzun@nku.edu.tr

Tarık Yerlikaya (Corresponding author)
 Oğuz Kırat
 Organization: Trakya University, Faculty of
 Engineering, Computer Engineering
 Department
 Address: Trakya Üniversitesi Ahmet
 Karadeniz Yerleşkesi Mühendislik
 Fakültesi 22020 Merkez / Edirne
 /TURKEY
 Phone (optional): +90 (284) 226 1217 /
 2215
 E-mail: tarikyer@trakya.edu.tr,
 ogzkirat@gmail.com

AN ALTERNATIVE EXECUTION MODEL FOR OPTIMUM BIG DATA HANDLING IN IOT-WSN CLOUD SYSTEMS

GÜNGÖR YILDIRIM, YETKİN TATAR

Abstract: *Wireless sensor networks (WSNs) are one of the fundamental IoT subsystems, which can produce big data. In line with this, the quality of data collected by sub-WSN systems is an important parameter. Eliminating redundant information and making the collected data understandable/interpretable will increase the efficiency and enable an easy integration with other IoT systems. In the paper, for the cloud system involving sub-WSN systems, an intermediate execution model which improves the quality of the collected data is proposed. The model includes semantic and data fusion/aggregation components in order to make the data more understandable and optimal. The model also proposes an approach that enables an interactive data analysis to be done between the data clients and the provider system. With this interactive model, it is aimed that the clients can execute their own data fusion and aggregation algorithms on the understandable big data set.*

Key words: *WSN, IoT, Sensor Cloud, Big Data*

1. Introduction

Nowadays, many different technologies may come together under IoT systems and play role in the solution of more complicated problems. IoT systems, which inherently have a heterogeneous structure, usually obtain the measurement data needed from sensor networks. Therefore, WSNs are one the fundamental subsystems of IoT projects. Traditional WSNs are generally closed systems designed for goal-oriented applications. The general processes such as data collection, storage and analysis are performed within the same system. In traditional WSNs, the data gathered are stored in three different ways. These are "internal storage", "Centralized storage" and "hybrid storage" [1]. In the internal storage, the data gathered are stored in the resources of the system in a distributed fashion. This type of storage is often preferred in the WSN systems with a large number of nodes. The centralized storage is a method usually used in small scale WSNs. In the centralized storage, the data is collected and handled in a central unit which is often external. The hybrid method is the combination of the other methods. It is not efficient to use the three methods in big scale IoT WSN systems due to some factors, such as the limited resources of WSNs and the possibility that the data gathered may reach big data sizes. IoT systems often deploy WSNs as subsystems. The WSN systems that integrate with IoT systems may have a

traditional structure or advanced features which use the IoT technologies such as 6LowPAN, RPL. Furthermore, these different WSNs may exist under the roof of a single IoT system (e.g., sensor cloud systems). One of the basic problems that arise in IoT systems which involve many different WSNs or a big scale WSN is the big and heterogeneous data management. The storage and handling of big data in IoT systems are usually carried out by a separated subunit involving big data management technologies [12, 13]. On the other hand, although today's big data technologies have advanced features and enough flexibility, the quality of the gathered data may need to be improved. When considered the adverse situations such as redundancy and repeating, which decrease the quality of both data and communication, the importance of the quality of the big data can be better understood in terms of the performance of an IoT system. For this purpose, data aggregation and data fusion operations are used in traditional WSNs [2]. However, these are not sufficient. Data aggregation operations are performed generally on the WSN nodes and usually goal-orientation. On the other hand, it is too difficult to execute the same data aggregation operations in an IoT system involving many heterogeneous WSNs. In addition, the content of the data which are demanded by the clients of an IoT system is also an important issue. The fact that the clients can access the interpreted

meaningful data instead of raw sensor data increases the service quality and enables an efficient big data management.

This paper proposes an effective intermediate execution model which can increase the quality and usability of the big data for IoT WSN systems like sensor clouds. The model focuses mostly on creating client-defined meaningful data, and efficient big data collection. In addition to these, an analytic approach which expresses the relation between WSNs and big data creation and shows its effect on the energy consumption in a WSN system, is also presented in the paper.

2. Related Works

In 2020 if we consider that the measured data from climate is going to reach 100 PB and 150 million sensors will be used in the experiment LHC in Europe, hence, it can be better understood the importance of the relation between WSN and big data [3] The literature studies about the relation of WSNs and big data are mostly related to the optimization processes taking place before the big data is stored. In a study [4], an algorithm that creates small data sets reflecting the character of the real big data is introduced. In the study, it is stated that the small data sets provide a successful big data management. The work in [5] focuses on the creation of big data in cluster based WSN systems which have a large number of sensor nodes. For these types of WSNs, algorithms towards the selection of optimum cluster head and the obtainment of optimum big data by mobile nodes are proposed. In the study, MULE and SENMA models are introduced for the mobile data collection. The study states that the design of big scale WSN systems can be quickly accomplished and the optimum big data quality can be achieved. In the other study [6] adaptive data collection algorithms towards the optimum big data management for periodical WSN systems are proposed. In the system operations, the algorithms take into account the association between node sampling rate and the current state of the physical phenomena. The study expresses that the big data collection can be achieved more efficiently by the algorithms proposed. Different approaches towards the integration of WSN and big data technologies can be also found in the literature [7-11].

3. Theoretical Relation Between WSNs and Big Data

The theoretical expression of WSN-big data relation could be useful in terms of showing the process of creation and communication of big data and the effect on the energy consumption in the

nodes. Accordingly, assume that there are “n” sensor nodes in a WSN running on a periodical operating mode.

$$W = \{d_1, d_2, \dots, d_n\}, \quad 1 \leq i \leq n, \quad (1)$$

$$T = \{t_1, t_2, t_3, \dots, t_j\}, \quad j \geq 1, \quad (2)$$

There are “j” resources types in the WSN and they are the member of the set “T”. For instance, “t1”, “t2”, “t3” denote analog LM35 temperature sensor, digital DS18B20 temperature sensor and ADXL345 3-axis accelerometer sensor, respectively. “ \ddot{E} ” is the resource entity matrix with “j x n” dimension, showing which node has what resources. A sensor node may theoretically have “j” resources. In this case, “ \ddot{E} ” can be expressed as;

$$\ddot{E} = \begin{bmatrix} k_1^{d_1} & k_2^{d_1} & \dots & k_j^{d_1} \\ k_1^{d_2} & k_2^{d_2} & \dots & k_j^{d_2} \\ \dots & \dots & \dots & \dots \\ k_1^{d_n} & k_2^{d_n} & \dots & k_j^{d_n} \end{bmatrix} \quad (3)$$

The entity value of each member of the matrix can be defined as a Boolean type.

$$k_j^{d_i} = \begin{cases} 1 & , \text{ if } t_j \in d_i \\ 0 & , \text{ otherwise} \end{cases} \quad (4)$$

On the other hand, the pre-defined matrix “B” denotes the data size produced by each resource in the system. The matrix is given as;

$$B = [b_1 \quad b_2 \quad \dots \quad b_j] \quad (5)$$

Besides, the node state matrix given in Eq.6 shows whether the nodes in the WSN are in active state or not;

$$N_s = [s_1 \quad s_2 \quad \dots \quad s_n], \quad \begin{matrix} s_i = 0, & d_i \text{ sleep/off} \\ s_i = 1 & d_i \text{ active} \end{matrix} \quad (6)$$

If it is assumed that all the nodes in the system run with a frequency, “f” in a duration of “P”, and their clocks are synchronous, the data size “D (byte)” gathered in the WSN will be;

$$D = f \times [N_s \times [\ddot{E} \times B^T]] \quad (7)$$

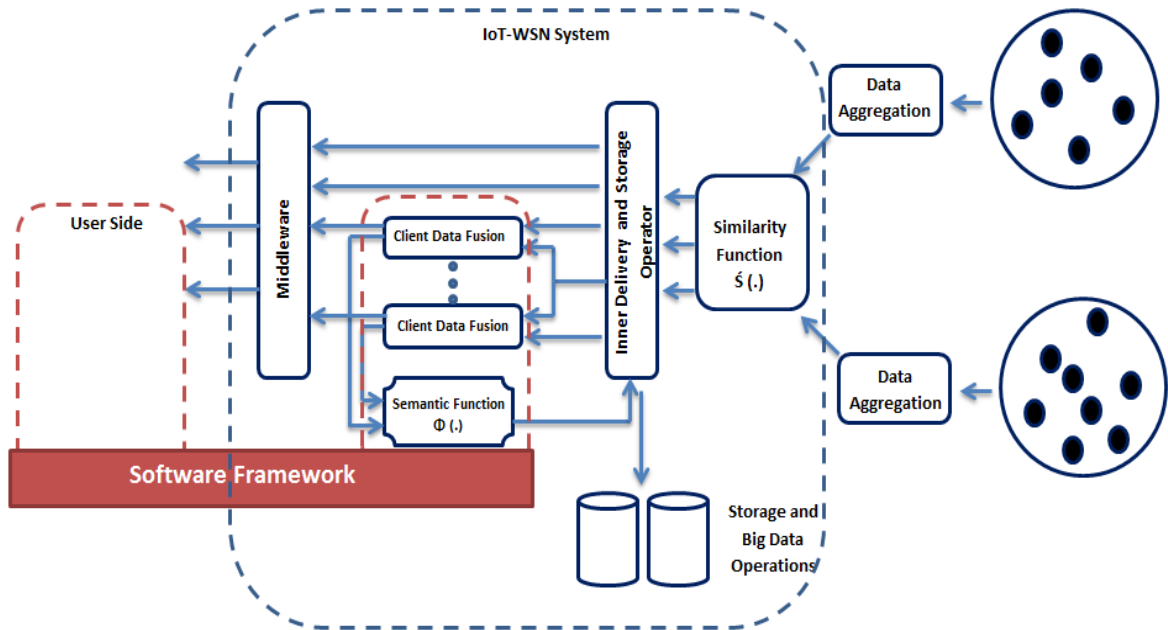


Fig.I. The Proposed Pre-Execution Model for Big Data Management in IoT-WSN Systems

The radio operations of a node are the phase in which energy consumption is the highest. According to the energy consumption model given in [14,15], the power consumptions in transmitting (E^{Tx}) and receiving modes (E^{Rx}) of a node in the WSN;

$$E^{Tx} = b_i \times (N + \epsilon \times d^2) \quad (8)$$

$$E^{Rx} = b_i \times N \quad (9)$$

where; N and ϵ are the radio unit and amplifier constants; " b_i " denotes the data size to be send/received and " d " is the transmission distance

As can be seen, in the case where n , j and d are too big, the data size and energy consumption will be too big. If the limited resources of a WSN are taken under consideration, the storage techniques used in traditional WSN systems will not be suitable for big scale IoT WSN systems.

4. Basic System Model

The type and data operations in IoT-WSN can be more different than those of traditional WSNs. For example, in the traditional WSNs, while the "data aggregation" is used in order to eliminate the redundancy and repeating situations, data fusion is used for inference and correlation operations. However, these methods alone are not sufficient for an IoT system since it can have different types of clients and resources. For this reason, the handling of the data to be collected or stored according to a specific working model is essential for efficient big

data operations. Besides, the type of the service in which the data gathered will be used is quite important. Namely, the clients getting services from IoT-WSN systems usually take raw or filtered data. In addition to this, some clients may demand the system to store the data sent from their own systems. On the other hand, clients are more interested on meaningful data rather than the raw data. The fact that the interpretation/semantic operations carried out in the provider side will bring a significant service quality and variation. Furthermore, these meaningful data can be either shared with other clients or can be used in other data mining solutions. What is important here is that a software infrastructure between the clients and the IoT-WSN system is deployed.

In Fig. I, an intermediate system infrastructure, which can perform efficient big data operations in an IoT-WSN system, is presented [15]. The proposed system focuses on three important operations. The first one is the regulatory operations which eliminate redundancy/repeating situations. The regulatory operations consist of two sub-processes. The data aggregation process is a central operation and usually performed in WSNs. This process is optional because of the fact that all WSNs cannot offer it. However, the similarity detection process is applied to all data from sub-WSNs. The similarity detection process is necessary to obtain the optimum data. In the process, the similarity ratio on different data sets is found and then the found ratio is compared with a threshold

value. Cosine and Dice are some of the similarity functions preferred for this purpose in literature [6]. After these processes, the obtained data can be shared with the clients or stored in the system. The intermediate execution system proposes a working model that the clients can interpret and customize the data at which they demand on the provider side. This consists of two sub-processes. The first is the client interpretation process in which the clients perform their own data fusion operations. The second is the fact that output of the data fusion is sent to the semantic function unit. The data obtained at the end of these processes can be shared with other clients demanding it, or can be stored on the big data file system. The most important advantage of the system is that it enables the clients to quickly interpret the data from different resources existing in different WSNs. The client-defined meaningful data can help other clients to select the analysis data for data mining operations. The model is based on a software framework which can be used on both the client side and provider side. With the help of the framework, all operations prepared by clients on the client side can be easily executed on the provider side. The created meaningful data will also increase the quality and usability of the big data on the system.

5. Conclusion

In this work, it has been discussed how WSN systems can be big data sources, and most efficiently how the created big data can be managed. In the study, it is also shown that the data storage techniques, used in traditional WSNs, are not suitable for IoT-WSN systems like sensor clouds. Nowadays, the big data management in IoT-WSN systems is among the active study areas. Especially, improvements in the quality of big data are an important subtitle in this area. Direct storage of raw sensor data in IoT-WSN systems leads to reducing the quality of big data. In addition, some adverse factors such as redundancy and repeating make the analysis of the stored big data harder. As a solution, an intermediate execution working model, which both eliminates the redundancy in the data coming and enables the data to be interpreted and shared, has been proposed. The originality of the proposed model is that the system allows the clients to carry out their data fusion operations on the provider side. This can be provided with a software framework infrastructure. This type of framework, which could be developed by today's object oriented languages, can provide quite flexibility for the big data management of IoT-WSN systems. The model used in both simulation and practical environment is open to development.

Acknowledgment

This study has been supported by the project of dept. of FUBAP of F.U, "Mo-bile Communication Technologies and Wireless Sensor Networks Laboratory", project no: MF1420..

REFERENCES

1. Hung C.C., Hsieh C.-C.,(2017), *Big Data Management On Wireless Sensor Networks*, Big Data Analytics for Sensor-Network Collected Intelligence, doi: <http://dx.doi.org/10.1016/B978-0-12-809393-1.00005>
2. Chhabra S., Singh D., (2015), *Data Fusion and Data Aggregation/Summarization Techniques in WSNs: A Review*, International Journal of Computer Applications, Vol. 121 – No.19.
3. Overpeck J. T. Meehl, G. A. Bony, S., and Easterling D. R., (2011), *Climate data challenges in the 21st century*, Science(Washington), , vol. 331, no. 6018, pp. 700–702.
4. Brumfiel G.,(2011), *Down the petabyte highway*. *Nature* ,vol. 469, no. 20, pp. 282–283
5. Cheng S., Cai Z., Li J., Fang X., (2015), *Drawing Dominant Dataset From Big Sensory Data in Wireless Sensor Networks*, IEEE Conference on Computer Communications INFOCOM,, doi: 10.1109/INFOCOM.2015.7218420
6. Medlej M., (2014), *Big data management for periodic wireless sensor Networks*. *Doctoral Thesis*, The University of Franche Comté
7. Rios L.G., Diguez J.A.I., (2014) *Big Data Infrastructure for analyzing data generated by Wireless Sensor Networks*, IEEE International Congress on Big Data proc, doi:10.1109/BigData.Congress.2014.142
8. Ang K.L.M., Seng J.K. P., Zungeru A. M.,(2017), *Optimizing Energy Consumption for Big Data Collection in Large-Scale Wireless Sensor Networks With Mobile Collectors*, IEEE Systems Journal, doi: 10.1109/JSYST.2016.2630691
9. Zeng J., Wang T., Lai Y., (2017), *Data Delivery from WSNs to Cloud based on a Fog Structure*, International Conference on Advanced Cloud and Big Data proc., pp:104-109
10. Christos A., Anadiotis G., Morabito G.o, Palazzo S. (2016),*An SDN-Assisted Framework for Optimal Deployment of MapReduce Functions in WSNs*, IEEE

- Transactions on Mobile Computing, Vol. 15, No. 9
11. Chung W. Y., Yu P.S, Huang C.J, (2013) *Cloud Computing System Based on Wireless Sen-sor Network*, Proceedings of the 2013 Federated Conference on Computer Science and Information Systems, pp. 877–880
 12. Yildirim G., Hallac I.R., Aydin G., Tatar Y, (2015), *Running genetic algorithms on Hadoop for solving high dimensional optimization problems*, Application of Information and Communication Technologies (AICT), 9th International Conference on, 2015, doi: 10.1109/ICAICT.2015.7338506
 13. Yildirim G., Aydin G., Alli H., Tatar Y., (2014), *An Analysis of Chaos-Based the FCW Op-timization Algorithm by Hadoop*. Eleco 2014 Elektrik – Elektronik – Bilgisayar ve Biyomedikal Mühendisligi Sempozyumu
 14. Heinzelman W., Chandrakasan A., and Balakrishnan H.,(2000), *Energy-Efficient Commu-nication Protocol for Wireless Microsensor Networks*. In Proceedings of the Hawaii Conference on System Sciences
 15. Yildirim G., (2018), *Design Of A Distributed-Parallel Cyber Physical System Based On Virtual Wireless Network*, Doctoral Thesis, Firat University Institute of Science and Technology, Turkey

Gungor YILDIRIM
Firat Univ. Computer Engineering
Elazig/Turkey
gyildirim23@gmail.com

Yetkin TATAR
Firat Univ. Computer Engineering
Elazig/Turkey
ytatar@firat.edu.tr

SECURITY PATTERNS FOR MICROSERVICES LOCATED ON DIFFERENT VENDORS

TIHOMIR TENEV, DIMITAR BIROV

Abstract: *The security holds a significant part of designing a cloud application, since it operates, in many cases, with sensitive information. There is an agile architecture style, which is mainly used in constructing a distributed application with advantage of independent partition scalability - Microservice architecture style. For mitigating security threats as disclosing or tampering of sensitive information, and not only, we suggest using Security patterns. In that paper, we estimate the issues that stands while designing a secure application, where the microservices are located in different Platform as a Service vendors. The first contribution is in making a vulnerable analysis and discovering the threats in this regard by using STRIDE. The second contribution is in providing recommendations of the proposed security patterns. That approach facilitates developers in better finding of the appropriate security pattern depend on their scenario.*

Key words: *Software architecture, Microservices, Security Patterns, Vulnerable analysis*

1. Introduction

The security is one of the main concern in designing distributed application, since it operates, in many cases, with sensitive information. Working toward mitigating of the security gaps, at earlier phase of application constructing, prevents sequential problems. There is an agile architecture style [1], which can be used for building a cloud application with advantage of independent partition scalability - Microservice architecture style [2].

Cloud system is functional part of distributed system [3] and has ability to grow with contemporary activities, which rapidly increases its complexity. The National Institute of Standards and Technology (NIST) [4] categorizes cloud context in meaning of three service models: Software as a Service (SaaS), Platform as a Service (PaaS) and Infrastructure as a Service (IaaS). The significant model, that we consider here, is Platform as a Service. It provides platform for deploying customer application, and exclude supporting of any hardware assets. In our case example, a Microservice based application is split and located on two PaaS vendors (Fig. 1).

Microservice architecture style is a new trend for an agile and decoupled partitioning of business logic [5], where each part is adapted to a bounded context. Those constraints lead to splitting an application into small services (*microservices*) accessible with lightweight mechanism. One approach for enhancing the security with microservices is by using Security patterns [6].

The roots of Security patterns are Design patterns [7]. Both types describe a solution for a problem, which occurs frequently in a certain case. Furthermore, each pattern should consider certain software constraints to implement the resolution accordingly. Taking into account the advantages of using Design patterns, we trust Security pattern approach for giving a resolution for a specific domain, as is security for microservices.

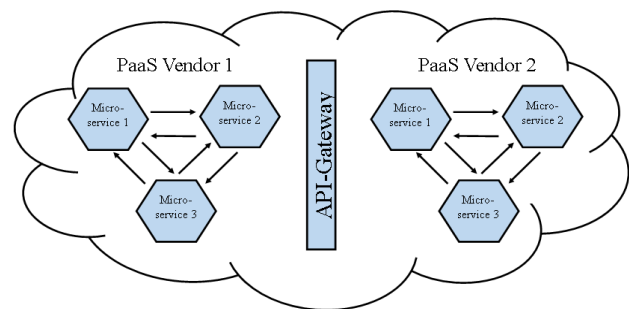


Fig. 1. *Cloud ecosystem with two PaaS vendors*

The problem that we consider here is how to enhance the security within a microservice based cloud application, if the same is split and spread on two PaaS vendors. One of the advantages of using cloud based application is the possibility of splitting the application and distributing on different PaaS vendors as Amazon AWS [8], Microsoft Azure [9], Google App Engine [10] etc. In such a case, the lightweight communication among the parties

happens either directly with REST [11] or by using API Gateway [12] as mediator.

In that paper, we contribute with vulnerable analysis and recommendation of using security patterns. We make the vulnerable analysis by using threat modelling approach named STRIDE [13]. Here STRIDE stands for *Spoofing, Tampering, Repudiation, Information Disclosure, Denial of Service* and *Elevation of Privilege*. Considering this, we were able to connect each of the security patterns to at least one STRIDE category and then to shape the recommendations. That methodology increases the comprehensibility and make our recommendations about the area of using security patterns more trustworthy.

The paper is structured as follows: Section 2 explains what each of the STRIDE categories means and provide vulnerable analysis in regard to distributing microservices on different PaaS vendors. Section 4 shows a set of security patterns and gives recommendation for each of them. Section 5 shows other papers in that matter. Section 6 derives conclusion and place the next steps for future work.

2. Vulnerable Analysis

The possibility of spreading microservices on different PaaS vendors makes microservice architecture style one of the most preferable approaches in nowadays. However, this benefit raises a concern in design perspective – how to enhance the security. In our paper we made a vulnerable analysis to shatter the distributing concept into subdivisions. Such an approach gives more clarity about the points that a certain security pattern protects.

For making the vulnerable analysis we decided to use threat model approach entitled STRIDE [13]. Here STRIDE stands for *Spoofing, Tampering, Repudiation, Information Disclosure, Denial of Service* and *Elevation of Privilege*. Each of those categories is explained toward distributing microservices on different PaaS vendors for identifying the types of attacks that a microservice may encounter:

Spoofing is a type of fraud where a violator tries to gain access to a microservice system or information by pretending to be the user. For

example, when a certain user waits for a microservice to trigger some event, and the event is triggered by an unauthorized user, usually this leads to irrelevant proceeding.

Tampering means making illegal changes of a data flow when a user should not. In many cases, tampering leads to adjusting the content of a file, memory unit or data, transferred over network. In some cases, the microservice data perhaps is persisted in a filesystem and keeping the integrity is a must.

Repudiation is in meaning of rejecting to accept something. An example is when a user did something, but claims he didn't touch it. This in many cases lead to less responsibilities. Such threat can be caught with logging of each activity performed by a user and/or a microservice.

Information Disclosure is in situation when an unauthorized user/microservice can operate with an information, which is forbidden for disclosing. The data source perhaps come from a running process, filesystem or data flow. The best way to mitigate that risk is to enhance the confidentiality.

Denial of Service mostly stands for exhausting of a hardware asset - memory, CPU, data store etc. Example is filling up the network bandwidth, which inflict high degree of response time among microservices. That category requires improving the microservice availability.

Elevation of Privilege observes cases, where a user has the possibility to do something without required access rights. For example, only administrators should operate with major services within an Operating System. Looking for authorization methods will mitigate those threats.

3. Security Recommendations

Standing by the principles of those three things – distributing microservices to different PaaS vendors, the needs of security enhancement and making a vulnerable analysis; we prepared a list of security patterns (Table 1) along with recommendations that can assist developers in building a secure application based on microservice architecture style.

Table 1. Security Patterns for microservices deployed on different vendors

Patterns	Spoofing	Tampering	Repudiation	Information Disclosure	Denial of Service	Elevation of Privilege
3rd Party Communication		x		x		
AGENCY GUARD		x		x		
AGENT AUTHENTICATOR	x					
Application Firewall	x				x	
Cloud Access Security Broker	x					x
Integration Reverse Proxy					x	
Known Partners	x			x		x

Each security pattern has its own structure. In many cases, it consists of several sections - *Intent, Context, Problem, Solution, Structure, Implementation, Consequences* and *Known uses*. However, we decided to emphasize on two of them - *Context* and *Solution*. They helped us to shape the recommendations. The *Context* describes the nature of a situation, which includes domain assumption and expectation of a system environment. The *Solution* guide us how to solve a problem by providing a decision.

In next several paragraphs we provide recommendations for each pattern from Table 1:

3rd Party Communication [14] gives several requirements in the process of negotiation of a business relationship between PaaS vendors. It has as a main intention to restrict all the information they persist and convey. Following all the requirements will prevent further data *Tampering* and/or *Information Disclosure*.

AGENCY GUARD [15] gives more clarity about implementing a guard mediator for accessing microservices. Direct access to a microservice is forbidden due to the risk of misuse. Therefore, redirecting all the traffic to pass through the API Gateway prevents *Tampering* and *Information Disclosure*.

AGENT AUTHENTICATOR [15] provides solution on how to secure the calls between two vendors by using authenticator for access granting and session keeping. The pattern can be enforced in API Gateway to prevent *Spoofing*.

Application Firewall [16] [17] encourage developers to use application firewall, based on specific policies, for filtering incoming requests. There is a case where a malicious user tries to gain access from Vendor 1 to Vendor 2. Applied in API Gateway it works against *Spoofing* and *Denial of Service*.

Cloud Access Security Broker [18] provides complete solution and can be used as third-party application for conveying authenticated users among microservices deployed on different places. It works towards *Spoofing* and *Elevation of Privilege*.

Integration Reverse Proxy [6] advises developers to store and frequently update all the metadata about each distributed software component in one data source. In that context, a microservice metadata, as addresses, access method etc., should be persist in one place, because sometimes a vendor may adjust its network configuration and this perhaps lead to microservice unavailability. The patter works against *Daniel of Service*.

Known Partners [6] advises developers to seek for vendors, which are in sufficient relationship. At least they are negotiated already and authenticate themselves in a secure manner. This pattern helps in preventing *Spoofing, Information Disclosure* and *Elevation of Privilege*.

4. Related Works

In [19], Roman Malisetti reviews a secure method, which relies on REST communication. The patterns he enforces are: Transport level security (TLS/SSL), which provides secure peer-to-peer authentication; OAuth, which enables consumers to access services through UI API, without using service credentials; Token-based authentication, which can be applied with OAuth together and can be used for exposing services over REST or SOAP. Here we enrich our list with more security patterns for different aspects.

Similar as [19], Anivella Sudhakar [20] represents techniques for REST securing. He describes differences between securing of REST and HTTP. Such mechanisms are: HTTP Authentication Schemes, which consist of Basic Authentication Scheme and Digest Authentication Scheme; Token-Based Authentication, which supports in authenticating of REST services; Transport Layer Security (TLS) and Secure Socket Layer (SSL); OAuth and OpenID.

Hafiz et al [21] show similar method, however, their categorization present an approach, which follows “one pattern per one STRIDE point”. Here is not the same, because a pattern may participate in more than one STRIDE point.

The focus at [14] is mainly against securing of web applications. They distinguish security patterns in two directions: procedural and structural. Structural patterns can be applied in already completed product, while procedures are aimed in phase of planning and writing software. Many of listed patterns don't match our paper. However, “Application Firewall” is borrowed to enrich our classification.

Fern et al [22] consider only three security patterns: Authorization, Role-Based Access Control, and Multilevel Security. They argue that these are the only three basic patterns that can be applied at each level of entire system. However, there are many scenarios, which require specific patterns. For instance, Known Partners [6] advises developers to seek for vendors, which are in sufficient relationship.

5. Conclusion

Our paper aims to show how security patterns can assist in designing a secure cloud application based on Microservice architecture

style. In particular we consider the possibility of spreading microservices on different PaaS vendors. In that context, we did a vulnerable analysis for estimating the threaten points, which need taking into account. The approach that we followed is based on the STRIDE. It guides us to shatter the subject into six categories - *Spoofing, Tampering, Repudiation, Information Disclosure, Denial of Service* and *Elevation of Privilege*. After we complete with the estimation, we start looking for the most applicable security patterns. Initially, the patterns provide rather common solutions and hence we prepared recommendations for each of them. The benefit here is that the reader can recognize and use them toward designing a secure application.

Further work is considered in researching of other security patterns, which match the rest aspects of an application align with Microservice architecture style – the data that a microservice persist, the accounts and identity etc. After covering the aspects, the follow step is to adopt a running case example to illustrate and discuss the patterns in more detail.

REFERENCES

1. Newman, S. (2015). *Building Microservices*, O'Reilly Media Inc.
2. Fowler, M. and Lewis, J. (2014). Microservices a definition of this new architectural term. from: <http://martinfowler.com/articles/microservices.html>.
3. Richard, M. (2015), *Software Architecture Patterns*, Page 27 – 35, O'Reilly Media Inc.
4. Mell, P. and Grance, T. (2009). The NIST Definition of Cloud Computing, Pages 800-145
5. Oberhauser, R. and Stigler, S. (2017). Microflows: Enabling Agile Business Process Modeling to Orchestrate Semantically-Annotated Microservices, *Proceeding of the Seventh International Symposium on Business Modeling and Software Design (BMSD 2017)*, pp 19-28.
6. Schumacher, M., Fernandez-Buglioni, E., Hybertson, D., Bushmann, F. and Sommerlad, P. (2006). *Security Patterns Integrating Security and Systems Engineering*
7. Gamma, E., Helm, R., Johnson, R. and Vlissides, J. (1995). Design Patterns – Elements of Reusable Object-Oriented Software, *Addison-Wesley Professional*.
8. Amazon Web Services Home Page, <https://aws.amazon.com/>
9. Microsoft Azure Home Page, <https://azure.microsoft.com/en-us>
10. Google App Engine Home Page, <https://cloud.google.com/appengine/>
11. Fielding, R. (2000). Representational State Transfer (REST), Chapter 5.
12. Richardson, C. and Smith, F. (2016). *MICROSERVICES From Design to Deployment*, pp 15-20, NGiNX Inc 2016.
13. Shostack, A. (2014). *THREAT MODELING: Designing for Security 1st Edition*, WILEY.
14. Romanosky, S. (2001). Security design patterns part 1.
15. Mouratidis, H., Giorgini, P., and Schumacher, M. (2003). Security patterns for agent systems, *in Proceedings of the European Conference on Pattern Languages of Programs. UVK - Universitaetsverlag Konstanz*, pp. 399–416.
16. Nelly Delessy-Gassant, S. R., Fernandez E. B. and Larrondo-Petrie, M. M. (2004). Patterns for application firewalls, *in Proceedings of the Conference on Pattern Languages of Programs*, pp. 1–19.
17. Fernandez, E. B., Larrondo-petrie, M. M., Seliya, N., Delessy, N., and Herzberg, A. (2003). A pattern language for firewalls, *in Proceedings of the Conference on Pattern Languages of Programs*, pp. 1–13.
18. Fernandez, E. B., Yoshioka, N. and Washizaki, H. (2015). Patterns for Security and Privacy in Cloud Ecosystems, *Evolving Security and Privacy Requirements Engineering (ESPREE), 2015 IEEE 2nd Workshop*.
19. Maliseti, R. (2011). Securing RESTful Services with Token-Based Authentication, *CA Technology Exchange, vol.1, 2011*, Page 43-48.
20. Sudhakar, A. (2011), Techniques for Securing REST, *CA Technology Exchange, vol.1, 2011*, Page 32-40.
21. Hafiz, M., Adamczyk, P. and Johnson, R. (2012). Growing a pattern language (for security), *Onward!*, Pages 139-158.
22. Fern, E. and Pan, R. (2001). A pattern language for security models, *PLoP 2001 Conference*.

Authors Contact:

Name: Tihomir Tenev

Organization: Faculty of Mathematics and Informatics, Sofia University “St. Kliment Ohridski”

Address: 5 James Bouchier blvd., Sofia, Bulgaria

Email: tenevtih@gmail.com

Name: Dimitar Birov

Organization: Faculty of Mathematics and Informatics, Sofia University “St. Kliment Ohridski”

Address: 5 James Bouchier blvd., Sofia, Bulgaria

Email: birov@fmi.uni-sofia.bg

EMBEDDED AUDIO CODING USING LAPLACE TRANSFORM FOR TURKISH LETTERS

MEHMET VURAL,¹ MUHARREM TUNCAY GENÇOĞLU^{2*}

Abstract: *In this paper a different cryptographic method is introduced by using Power series transform. A new algorithm for cryptography is produced. The extended Laplace transform of the exponential function is used to encode an explicit text. The key is generated by applying the modular arithmetic rules to the coefficients obtained in the transformation. Here, ASCII codes used to hide the mathematically generated keys strengthen the encryption. Text steganography is used to make it difficult to break the password. The made encryption is reinforced by audio steganography. To hide the presence of the cipher text, it is embedded in another open text with a stenographic method. Later, this text is buried in an audio file. For decryption it is seen that the inverse of the Power series transform can be used for decryption easily. Experimental results are obtained by making a simulation of the proposed method. As a result, embedded text is increased security by hiding inside an audio file.*

Keywords: *Cryptography, Power Series Transform, Data Encryption, Embedded Image 2000 AMS Classification: 94A60, 11T71, 14G50, 68P25*

1. Introduction

The confidential communication, with the technological progress it has varied in terms of form and methods have maintained continuous its importance. To be very important of privacy in applications; protected information before hand of third parties were aimed to sending related destination and studies in this direction were made [2,3,5,6]. Network security problem has become very important in recent years. E-banking, e-commerce, e-government, e-mail, SMS services, security of ATMs, and the existence of financial information has become indispensable in our lives. Protecting the information that is processed and transferred in these environments or to ensure safety is of great importance. There are many threats such as unauthorized access, damage, etc. while data communication is being performed in the digital environment. In order to eliminate these threats many encryption techniques are developed [5-7, 9]. Cryptography is the all of mathematical technical studies related to information security. Cryptology is a cipher science and ensures security of information.

The main goal of cryptography is to allow communication of two people through non-secure channels. Encryption is the process of blocking information to make it unreadable without special knowledge. These operations are expressed using an algorithm. In general this is called the symmetric algorithms. For encryption and decryption must be

used the same secret key in the symmetric algorithms [2]. The converse is also true. The security of these algorithms is related with the secret key [5]. The original information is known as plain text and cryptic text is encrypted format of this text. Encrypted text message contains all of the information in plain text message but it is not a readable format by a human or a computer without a suitable mechanism to decryption. The cipher is often expressed with parameters called key which is as part of the external information. Decryption is almost impossible without an appropriate key. Advanced Encryption Standard (AES) method is the most used. Encryption converts data into an incomprehensible format and makes difficult to access the actual data, however, cannot ensure the confidentiality of communications. Steganography as word meaning means hidden text or covered text. The objective of the Steganography is hide the presence of a message and is create a channel to the implicit [8]. It is art of storing information which cannot be detected the presence [9]. The aim of this study; the steganography and cryptography are used together to increase the security for confidential data. Power series are used for cryptography. Later this process is supported by an 8-bit ASCII code and is held in high security applications for confidential data combined with steganography. In the second section of the study; definitions and some standard results are given for the proposed

method respectively. In the third section, flow diagrams are given together with recommended method and practice. In the fourth section, the evaluations of the results from the study are situated.

1.1. Our Contribution

We wrote a Power series transformation algorithm which can be used in the encryption methods existed in the literature. Next, we buried a hidden text into any audio by cryptology and steganography. We suggested a hybrid method for crypto machines.

In this article, a new information security model based on the Taylor series for steganography is proposed. The proposed model is used for both cryptography and steganography. The

2. Preliminaries

Definition 2.1.

Let $f(t)$ be defined for $t > 0$. We say f is of exponential order if there exist numbers $\alpha, M > 0$ so that

$$|f(t)| \leq M e^{\alpha t} \quad (2.1)$$

If $f(t)$ is exponential function, then we have $f(t) = \infty$ for $t \rightarrow \infty$ [1].

Definition 2.2.

Let $f(t)$ be given for $t \geq 0$ and assume the function satisfy the property of α exponential order and $t, s \in \mathbb{R}$. The Laplace transform of $f(t)$ is defined by [1]

$$F(s) = \int_0^{\infty} e^{-st} f(t) dt \quad (2.2)$$

Let's define a new transformation function by expanding the Laplace transformation using Definitions 1 and 2.

Definition 2.3.

Transformation of $f(t)$ for every $t \geq 0$ is defined as:

$$F(h) = T[f(t)] = \int_0^{\infty} \frac{1}{h} e^{-\frac{t}{h}} f(t) dt. \quad (2.3)$$

(Extended Power Series Transformation) We present $f(t) = T^{-1}[F(h)]$ to define the inverse transformation of $f(t)$. Obtained extended power series transformation has the following standard results [3]

$$1. T\{t^n\} = \frac{n!}{s^{n+1}} \Rightarrow T^{-1}\left\{\frac{1}{s^{n+1}}\right\} = \frac{t^n}{n!}$$

characterization of the proposed model is given as follows.

- A new method, in which steganography and cryptography are used together, is proposed.
- The proposed steganography method is a media independent steganography method. This method is used for both in audio steganography and text steganography.
- A Taylor series based coding method is defined mathematically and practically.
- The method is simulated and the simulation results are shown in the experimental results clearly.

$$2. T\{t^n e^{st}\} = \frac{n! \cdot h^n}{(1-sh)^{n+1}} \Rightarrow T^{-1}\left\{\frac{h^n}{(1-sh)^{n+1}}\right\} = \frac{t^n \cdot e^{st}}{n!} \quad (t \geq 0) \quad (2.4)$$

Definition 2.4.

In order to keep the text information in the computer memory computer system assigns a numerical value to each letter or symbol. This process depends on the encoding system. By setting the numerical value of symbols, in order to represent non-numeric or alphabetic type of information on the computer the most commonly used as the coding system is used in ASCII coding system.

Definition 2.5.

The process of fitting a data or message into another object is called steganography. The goal is to conceal the existence of the message [6].

3. The Proposed Method

Combining cryptography and steganography methods, the application stages that increase the data security and privacy of this proposed hybrid model are as follows.

3.1. Encryption

In this method, a new encryption method based on the Taylor series is proposed. The steps of the proposed method are as follows.

Step 1: The Taylor series is expanded with e^t . Then this value is multiplied by t^3 to generalize the mathematical relation to be used in the encryption algorithm.

Step 2: The number values corresponding to the letters in the alphabet are applied to the text to be encrypted.

Step 3: The numbers found are replaced in the generalized encryption algorithm.

Step 4: The Power series transform is applied to the function obtained from here.

Step 5: The obtained coefficients are found mod 28 values.

Step 6: Instead of these numbers, the encryption keys are found by taking the quotients in the mode operation.

Step 7: The text is encrypted by writing the letters corresponding to these keys.

Step 8: Encrypted text is converted to ASCII code and the corresponding numbers are found. Then these numbers are converted into a binary system.

Step 9: These numbers are hidden into any text by a method that the user will specify.

Step 10: Create Stego Object and sender sends this embedded audio file.

Example

Assume that we want to send the message "FIRAT". Firstly we consider extended Taylor series with e^t :

$$\begin{aligned} f(x) &= f(a) + \frac{f'(a)}{1!}(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \\ &\dots + \frac{f^n(a)}{n!}(x-a)^n + \dots \\ &= \sum_{n=0}^{\infty} \frac{f^n(a)}{n!}(x-a)^n. \end{aligned} \quad (3.1)$$

Then, if we expand:

$$e^t = 1 + \frac{t}{1!} + \frac{t^2}{2!} + \frac{t^3}{3!} + \dots = \sum_{n=0}^{\infty} \frac{t^n}{n!} \quad (3.2)$$

With t^3 , then we get:

$$t^3 e^t = t^3 + \frac{t^4}{1!} + \frac{t^5}{2!} + \frac{t^6}{3!} + \dots = \sum_{n=0}^{\infty} \frac{t^{n+3}}{n!} \quad (3.3)$$

Therefore, we obtain:

$$f(t) = \sum_{n=0}^{\infty} K_n \frac{t^{n+3}}{n!}. \quad (3.4)$$

If we enumerate letters of the alphabet from scratch "FIRAT" plain text is equal 6,9,19,0,22. If we write $K_0=6, K_1=9, K_2=19, K_3=0, K_4=22$ in to (3.4), we get

$$\begin{aligned} f(t) &= \sum_{n=0}^{\infty} K_n \frac{t^{n+3}}{n!} \\ &= K_0 \frac{t^3}{0!} + K_1 \frac{t^4}{1!} + K_2 \frac{t^5}{2!} + K_3 \frac{t^6}{3!} + K_4 \frac{t^7}{4!} \end{aligned} \quad (3.5)$$

If we apply extended power series transformation to both sides of (3.5), we get

$$\begin{aligned} T[f(t)](h) &= T\left[\sum_{n=0}^{\infty} K_n \frac{t^{n+3}}{n!}\right](h) \\ &= T\left[K_0 \frac{t^3}{0!} + K_1 \frac{t^4}{1!} + K_2 \frac{t^5}{2!} + K_3 \frac{t^6}{3!} + K_4 \frac{t^7}{4!}\right](h) \\ &= 6.3! h^3 + 9.4! h^4 + 19.5! \frac{h^5}{2!} + 0.6! \frac{h^6}{3!} \\ &\quad + 22.7! \frac{h^7}{4!} \\ \sum_{n=0}^{\infty} K_n (n+3)! \frac{h^{n+3}}{n!} &= 36h^3 + 216h^4 + \\ 1140 \frac{h^5}{2!} + 0 \frac{h^6}{3!} + 4620 \frac{h^7}{4!}. \end{aligned} \quad (3.6)$$

The provisions of 36,216,1140,0,4620 in the modes (28) are (K_n) 8,20,20,0,0. If we write quotient in mode operation instead of these numbers, we obtain the key (K'_n) 1, 7,40,0,165. "FIRAT" plain text converts "HSSAA" by (3.3).

If we convert "HSSAA" encrypted text to 8-bit characters in the ASCII code we obtain 72, 83,83,65,65. If these codes are written in binary system we get the keys $(1001000)_2, (1010001)_2, (1010001)_2, (1000001)_2, (1000001)_2$. ASCII 8 bit keys are in the text as follows: A provision giving the binary number system in the space between each word of the text namely if the number between two words 1 then we get $(1)_2$ and we define 2 with $(0)_2$.

Table 1. Embedded Text

<p>Günümüzde gelişen teknoloji ile birlikte gizli tutulması gereken verilerin dışarıdan ulaşılması ve bazı sistemler yardımıyla, bilginin kaynağından başka bir yere aktarılması veya bilginin değiştirilmesi bilginin sahipleri için çok ciddi sorunlar oluşturmaktadır. Bilginin gizli kalmaması bireyleri, toplumları ve devletleri çok kötü şekillerde etkileyeceği için bu konu önemli ve tehlikeli bir tehdit oluşturmaktadır</p>	<p>Today, with the development of technology, accessing the data that need to be kept secret from outside and transferring the information from other sources through the help of some systems or changing the information constitute very serious problems for the owners of the information. This is an important and dangerous threat because the inability of information to be kept secret will affect individuals, societies and states in very bad ways.</p>
---	---

Sender also send this text clearly with (1, 7, 40, 0,165) secret key.

Hence theorem can be following

Theorem 3.1.

The given plain text in terms of (K_n) , under Laplace transform of $K_n \frac{t^{n+3}}{n!}(h)$, can be converted to cipher text,

$$(K'_n) = (K_n) - 28q_n \quad (n=0, 1, 2, \dots) \quad (3.7)$$

Where a key

$$q_n = \frac{K_n - K'_n}{28} \quad (n=0, 1, 2, \dots) \quad (3.8)$$

3.2. Decryption

Steps of the proposed decryption method are given below.

Step 1: The receiver opens embedded audio file and writes hidden ASCII codes into the text.

Step 2: He finds the letters corresponding to these codes.

Step 3: He obtains the numbers corresponding to the letters.

Step 4: These numbers and the secret key are written in place of the inverse power series.

Step 5: The letters corresponding to the coefficients obtained here are written.

Step 6: The first clear text is obtained.

Example

The recipient receives a text message and by reading the spaces between words with software that will get the data buried create the necessary numerical equivalents. If these numbers are divided into 8-bits groups then ASCII provision of the data buried has been obtained. We can see the hidden data buried "HSSAA" in the Table 2.

Table 2. Embedded text and solution

Günümüzde gelişen teknoloji ile birlikte gizli tutulması gereken verilerin dışarıdan ulaşılması ve bazı sistemler yardımıyla, bilginin kaynağından başka bir yere aktarılması veya bilginin değiştirilmesi bilginin sahipleri için çok ciddi sorunlar oluşturmaktadır. Bilginin gizli kalmaması bireyleri, toplumları ve devletleri çok kötü şekillerde etkileyeceği için bu konu önemli ve tehlikeli bir tehdit oluşturmaktadır					
Text Bits	100 1000	101 0001	101 0001	100 0001	100 0001
Secret Data	72	83	83	65	65
Şifreli Message	H	S	S	A	A

If we write H,S,S,A,A→8,20,20,0,0 and secret key values (1,7,40,0,165) into

$$A_n = \frac{K_n - K'_n}{28}$$

$$36 = 28x1 + 8$$

$$216 = 28x7 + 20$$

$$1140 = 28x40 + 20$$

$$0 = 28x0 + 0$$

$$4620 = 28x165 + 0 \text{ are obtained.}$$

If we apply these values 36,216,1140,0,4620 to the

$$\sum_{n=0}^{\infty} K_n (n+3)! \frac{h^{n+3}}{n!}$$

then, we get

$$\sum_{n=0}^{\infty} K_n (n+3)! \frac{h^{n+3}}{n!} = 36h^3 + 216h^4 + 1140 \frac{h^5}{2!} + 0 \frac{h^6}{3!} + 4620 \frac{h^7}{4!}$$

$$= 6.3! h^3 + 9.4! h^4 + 19.5! \frac{h^5}{2!} + 0.6! \frac{h^6}{3!} + 22.7! \frac{h^7}{4!}$$

(3.7)

If we apply inverse Extended Power Series Transformation to both sides of the (3.7), then we get

$$T^{-1} \left[\sum_{n=0}^{\infty} K_n (n+3)! \frac{h^{n+3}}{n!} \right] =$$

$$T^{-1} \left[6.3! h^3 + 9.4! h^4 + 19.5! \frac{h^5}{2!} + 0.6! \frac{h^6}{3!} + 22.7! \frac{h^7}{4!} \right]$$

$$\sum_{n=0}^{\infty} K_n \frac{t^{n+3}}{n!} = 6. t^3 + 9. t^4 + 19. \frac{t^5}{2!} + 0. \frac{t^6}{3!} + 22. \frac{t^7}{4!}$$

If we convert the K_n coefficients we will get the first plain text 6,9,19,0,22→F,I,R,A,T.

Hence theorem can be following

Theorem 3.2.

The given cipher text in terms of (K'_n) , with a given key q_n , can be converted to plain text (K_n) under the inverse Laplace transform of

$$T^{-1} \left[\sum_{n=0}^{\infty} K_n (n+3)! \frac{h^{n+3}}{n!} \right] = \sum_{n=0}^{\infty} K_n \frac{t^{n+3}}{n!},$$

Where

$$K_n = 28q_n + K'_n \quad (n=0, 1, 2, \dots).$$

Operations performed in this section are shown in Figure 1 and Figure 2.

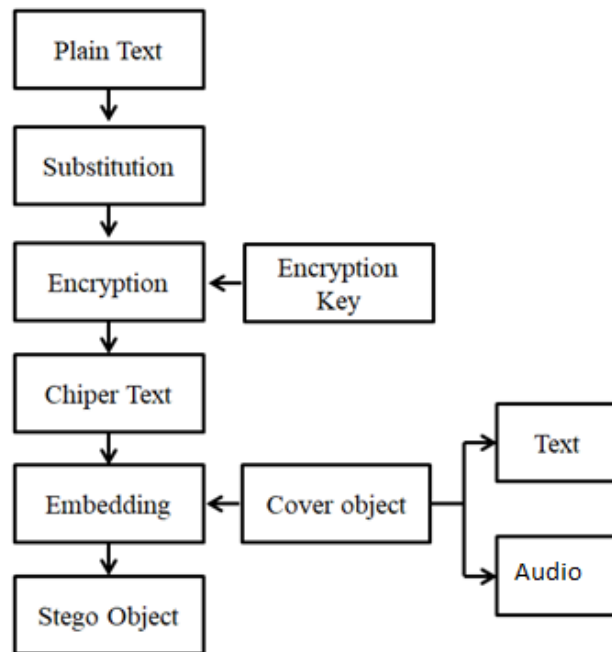


Fig. 1. Flow Diagram of Encryption System

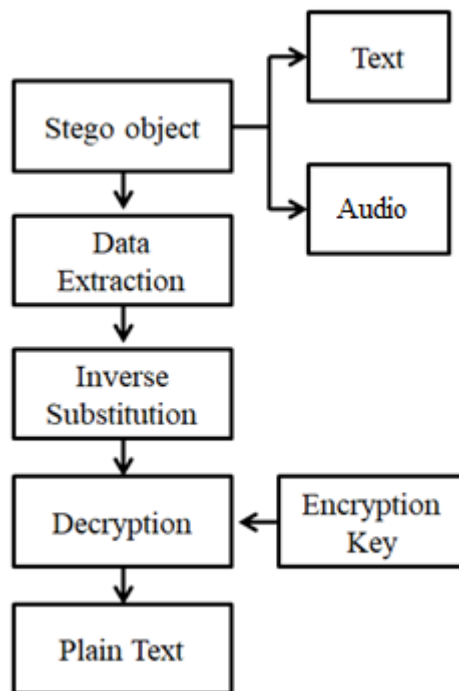


Fig. 2. Flow Diagram of Decryption System

4. Experimental Results

In this article, recommended the method for both text and image media is applied. A simulation program is carried out to obtain experimental results. Text steganography and image steganography results are obtained in the performed

simulation. Simulation of text steganography windows Fig. 3 and 4 are shown. Fig. 3 shows encryption method and data embedding window and Fig.4 shows the data extraction and decryption window, respectively.

sifreleme

Çalıştır

Tüm verileri sil

Word dosyası oluştur

Dosya Ekle

Gönder

Açık Metin: FIRAT

Sayı Dönüşümü: 6 9 19 0 22

Algoritma: 36 216 1140 0 4620

Mod: 8 20 20 0 0

Şifreli Metin: H S S A A

Ascii: 7283836565

Binary: 001001000001010011001010011001000001001000001

Gömülecek Metin: Günümüzde gelişen teknoloji ile birlikte gizli tutulması gereken verilerin dışardan ul...

Gömülü Metin: Günümüzde gelişen teknoloji ile birlikte gizli tutulması gereken verilerin dışarda...

Anahtar: 1 7 40 0 165

Alıcının e-mail:

Fig. 3. Encryption window of the simulation

SAyı Dönüşümü

Gömülü Metni Giriniz: Günümüzde gelişen teknoloji ile birlikte gizli tutulması gereken verile...

Ascii Binary: 001001000001010011001010011001000001001000001

Ascii: 72 83 83 65 65

Şifreli Metin: HSSAA

Mod: 8 20 20 0 0

Anahtarı Giriniz: 1 7 40 0 165

Sayı Dönüşümü: 6 9 19 0 22

Şifremiz: FIRAT

Şifreyi Cöz

Tüm verileri sil

Fig. 4. Decryption window of the simulation

In this part of the simulation, audios are used as data hiding media. Original audio window is shown in Fig. 5.

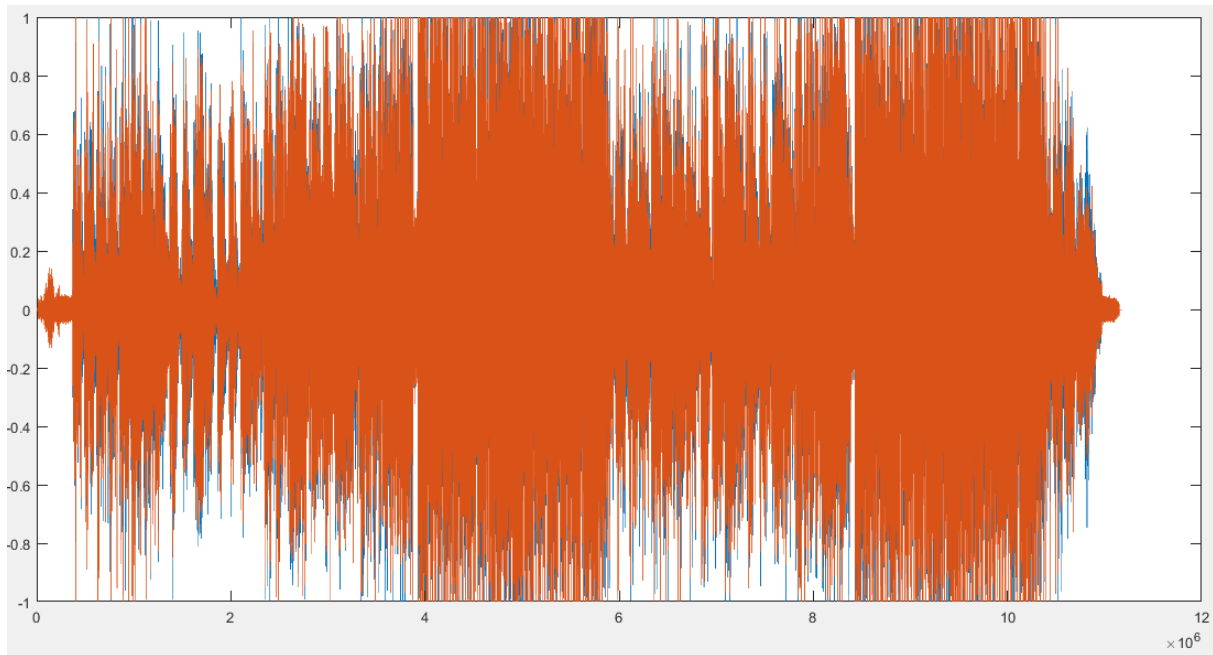


Fig. 5. Original window for audio.

Data encryption and embedding window is shown in Fig. 6

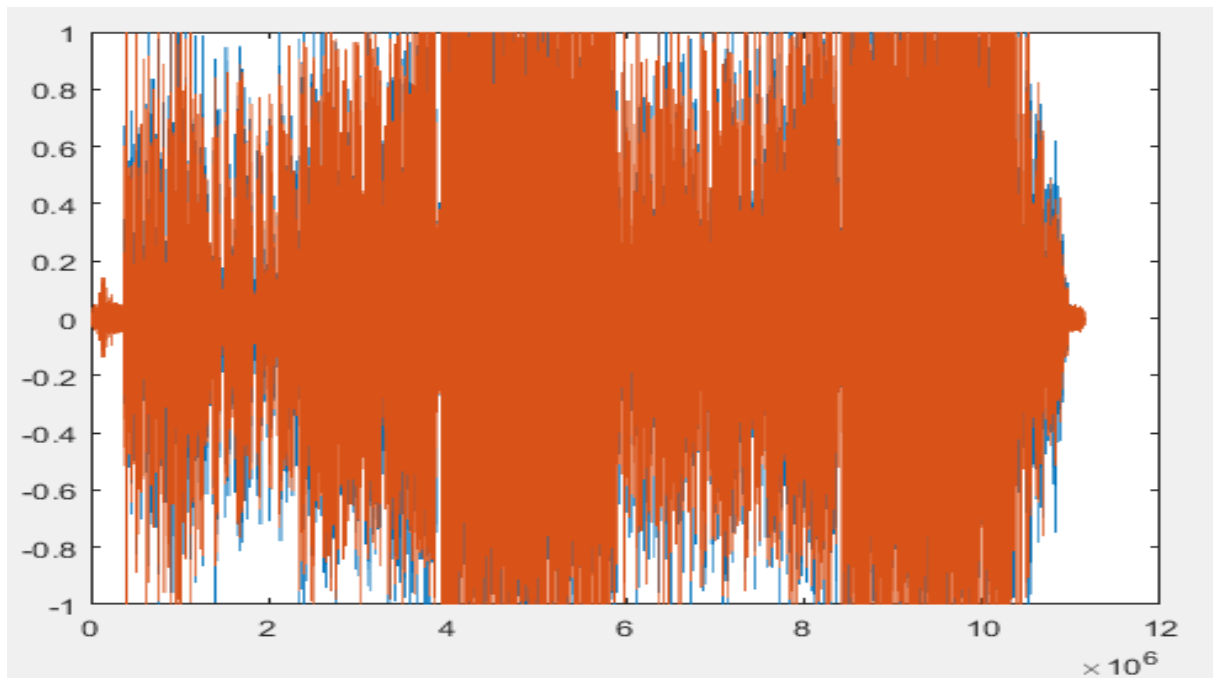


Fig. 6. Data encryption and embedding window for audio.

Matlab window for audio is shown in Fig. 7.

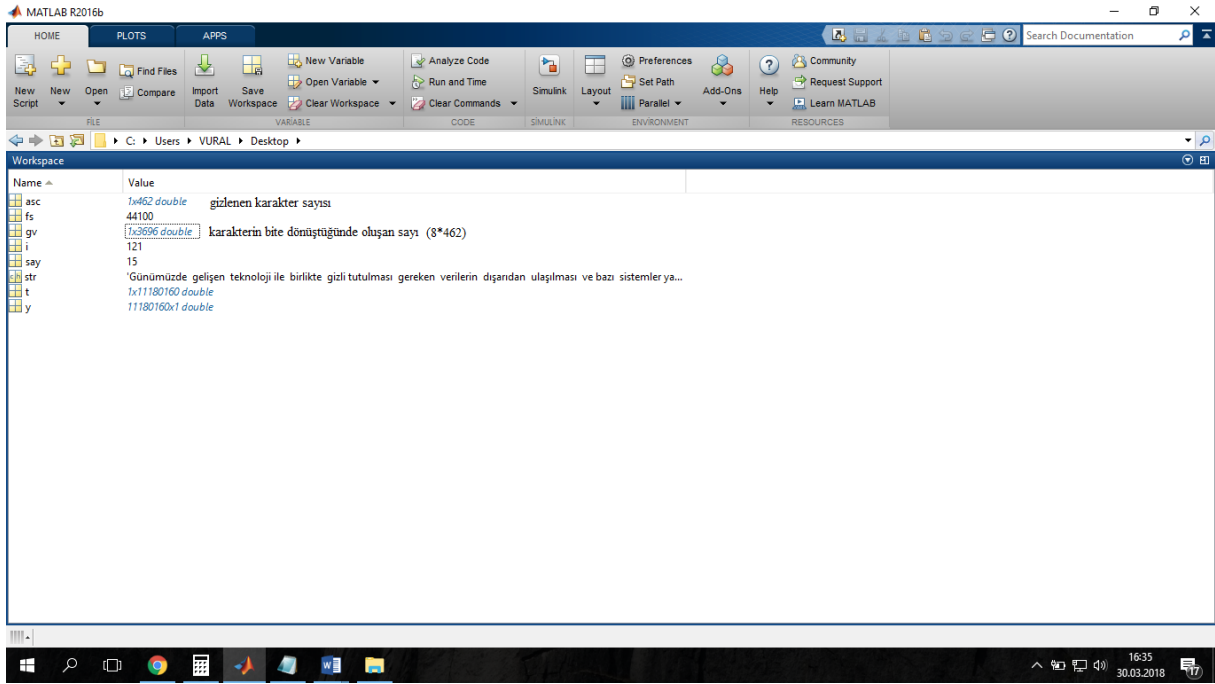


Fig. 7. Matlab window for audio.

5. Conclusion and Recommendations

The proposed algorithm is created by using the power series transformation. Keys generated using this algorithm is applied to the method known as substitution method in the literature. In our practice, the keys obtained by the proposed method that emerged as an end result of digitization are used. A hybrid model is developed by using steganography to provide high security and explained in detail. The user can hide this message which is obtained by taking the coefficient q_n instead of the coefficient (K'_n) and also presence of this message is hidden with ASCII code. Then, using another password, an encrypted text can be hidden into a text by the proposed method. Then the encrypted text is embedded in an audio file. In this way the security level of the data can be increased. As a result, Embedded text is increased security by hiding inside an audio file.

References

- [1] Aydın, M., Gökmen, G., Kuryel, B., Gündüz, G. Diferansiyel Denklemler ve Uygulamaları, Barış Yayınları. SS 332-349, 1990.
- [2] Delfs, H. and Knebl, H. Introduction to Cryptography Principles and Applications, Springer 2007.
- [3] Gençoğlu, M.T. Use of Integral Transform in Cryptology. Science and Eng. J of Fırat Univ., 28 (2), 217-220, 2016.
- [4] Johnson, N.F. and Jajodia, S. Exploring steganography: Seeing the unseen, Computer, 31 (2), 26-34, 1998.

- [5] Koç, Ç.K. Cryptographic Engineering, Springer. PP 125-128, 2009.
- [6] Martin, K.M. Everyday Cryptography Fundamental Principles and Applications, Oxford University Press 2012.
- [7] Paar, C. and Pelzl, J. Understanding Cryptography, Springer 2010.
- [8] Usha, S., Sathish Kumal G. A, Boopathybagan., K. A. Secure Triple Level Encryption Method Using Cryptography and Steganography, International Conference on Computer Science and Network Technology, IEEE 2011.
- [9] Yalman, Y. and Ertürk, İ. Kişisel Bilgi Güvenliğinin Sağlanmasında Steganografi Biliminin Kullanımı. ÜNAK 2009 Bilgi Çağında Varoluş "Fırsatlar Ve Tehditler" Sempozyumu 2009.

Authors' contacts

² Fırat University, Vocational School of Technical Sciences, 23119 Elazığ, Turkey,

¹ Fırat University, graduate Student, Faculty of Technology, 23119 Elazığ, Turkey,

*Correspondence: mt.gencoglu@firat.edu.tr

MANAGED ACTIVE DIRECTORY IN DIRECTORY-AS-A-SERVICE

VENETA ALEKSIEVA, SVETOSLAV SLAVOV

Abstract: *Active Directory (AD) is the one of the last categories to make the transition to the cloud. AD as SaaS will give solution to the IT administrators can take advantage of the synchronized directories between on-premises to the cloud, and by doing that the same identity will be used on both environments. In this paper is proposed an web-based management system for AD, which provided a seamless and simple experience to the IT administrators and synchronize directories between on-premises to the cloud, and by doing that the same identity is used on both environments.*

Key words: *Active Directory, Directory as a Service, SaaS*

1. Introduction

Nowadays with the advent of SaaS and managed services, a number of IT management tool categories are making the leap to be delivered as outsourced services. Active Directory (AD), which is a Windows OS directory service since 1999, is the fundament on which network security is built. However, not only have IT resources dramatically changed in the last two decades since the beginning of AD, but IT administrators don't want to have to deal with the costs and time it takes to manage Active Directory hardware. With all of the benefits of moving to the cloud, including lower upfront costs and maintenance, AD is the one of the last categories to make the transition to the cloud. The Directory-as-a-Service allows IT administrators to create a secure and centralized environment that offers users access to all of their IT resources [1,2].

In the cloud era, any organization can have zillions of applications available to end users, so having usernames and passwords for every single application like we used to have in the past is just not practical anymore. AD as SaaS will give solution to the IT administrators with one set of credentials to authenticate to systems (Mac, Linux and Windows), legacy and web-based applications, on-prem and virtual files, and wired and wireless networks. Furthermore, IT administrators don't have to worry about upkeep, maintenance, or the complex management that often comes with an on-prem identity management solution. Moreover, an organization can take advantage of the synchronized directories between on-premises to the cloud, and by doing that the same identity will be used on both

environments, providing a seamless and simple experience to the end users.

A comprehensive web based directory service platform is highly sought after in modern IT organizations. This approach present only web-based system for management of ADs.

2. AD and Directory-as-a-Service

Active Directory worked as an organization would have AD located on-prem, and an internal IT team was responsible for managing the identity provider. AD worked on a direct connect model, so any IT resources needed to be close to the AD server. If hosting AD at a third party location will require the VPNs and will increase networking and security work.

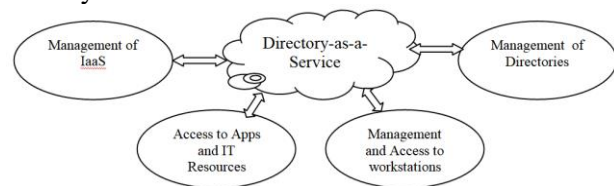


Fig. 1. *Directory-as-a-Service*

Directory-as-a-Service (fig.1) simplifies task execution on devices including globally updating policy settings, modifying registry settings, applying patches, and changing system configurations. It ensures consistency across company environment, by allowing IT administrators to group like objects and apply the same policies and configurations across them. Requests to authenticate users are sent to the Cloud via LDAP protocol. The Cloud agent can also be deployed on the Windows, Mac, and Linux devices

for task and policy management, survivability, and security auditing.

Transition of AD to Directory-as-a-Service need the services to be treated as abstract objects, decomposed set of more basic components or service objects. The creation of new services and new service instances is simplified, and will allow the IT administrator to validate and analyze the performance of these higher-level services in a fashion that parallels the construction of those services and service instances.

Together, these basic AD services and support functions can be orchestrated to produce more sophisticated service constructs that can be easily replicated and customized, reserved, activated, and then operationally monitored and verified throughout the lifecycle of the entire AD service.

3. Related works

Since twenty years the EU project GÉANT[3] operates the pan-European GÉANT network (more than 73 national networks-NRENS), delivering advanced multi-domain services and facilitates joint-research activity that drives innovation and providing the best possible infrastructure to ensure that Europe remains in the forefront of research. Software Defined Networking (SDN) and technologies such as Network Function Virtualization (NFV) and Network as a Service (NaaS) offer national research and education networking (NREN) organizations the ability to overcome the limitations imposed by more traditional service provider technologies. But its version of AD isn't a cloud version of the on-prem Active Directory solution.

Some cloud organizations such as Amazon Web Services (AWS) have introduced a managed Active Directory solution for their own IaaS platform [4]. AWS exposes a series of domain join settings during the virtual machine instance creation process in the Virtual Private Cloud (VPC). AWS gives two options to base your Active Directory environment - on the Standard or Enterprise edition of Windows Server. [5]

Azure has its version of AD, called Azure Active Directory. But it isn't a cloud version of the on-prem Active Directory solution. Azure AD can be integrated with an existing Windows Server Active Directory, giving organizations the ability to leverage their existing on-premises identity investments to manage access to cloud based SaaS applications. [6] Active Directory Domain Services is hierarchical Database with forest and domains. Whereas Azure AD is flat system without any forest, domains and trusts.

In the case of both AWS and Azure, the end result is that their managed Active Directory approach is really for use within their cloud infrastructure.

In [7] is presented the advantages of JumpCloud Directory-as-a-Service® against the traditional services as Microsoft® Active Directory® (AD) and OpenLDAP™. While both of these solutions are great for homogenous, on-prem IT networks, the issue with AD and OpenLDAP is that modern IT organizations seek to eliminate the majority of their on-prem infrastructure in favor of cloud solutions.

Up to this moment Active Directory as a Service solution which is a replacement to the on-prem Active Directory doesn't really exist.

4. Web-based system for management of ADs, named MATEX

The developed Web-based system (MATEX) manages ADs. The server side is based on PowerShell skripts and cmdlets. The system has developed with Visual Studio 2010, ASP.NET, C#, and Windows Powershell. The architecture of the MATEX is presented on the Figure 2.

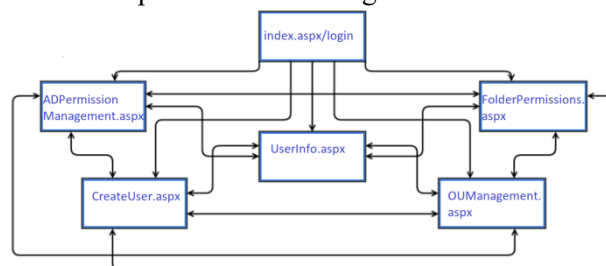


Fig. 2. The Architecture of the MATEX

The MATEX has worked on three virtual machines (VM) - two servers (BGVARNADC01, JPTOKYODC01) and one client USTESTPC01. On the first server BGVARNADC01 (Windows Server 2012R2 Std) / VM have started AD Primary Domain Controller (PDC) and all FSMO roles holder, DHCP Primary server, ISS Web server, File server, SMTP server, Backup server. On the second server JPTOKYODC01 (Windows Server 2012R2 Std) / VM have started AD Domain Controller server, DHCP backup server, WSUS server, File Server, Backup Server. The third VM USTESTPC01 (Windows 7 Pro)/VM has role of client workstation and there have configured Dynamic DNS record to joined into the domain and IP settings from DHCP server.

The Matex.com domain is made up of two domain controllers (DCs) - BGVARNADC01 and JPTOKYODC01, which a bidirectionally replicate all AD containers and DNS containers. Thus, a change or action done on one DC is replicated and

is visible to the other. The default replication time is one hour.

The Organization Unit (OU) structure is divided into 3 levels for easier management and classification of users and resources (Region Level - Africa, Asia, Europe ...; Country Level - China, India, Japan ...; Site Level - Fuji, Kobe, Tokyo). It is presented on Figure 3.

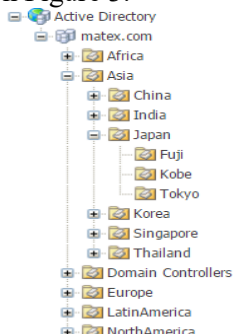


Fig. 3. The OU structure

The DNS functionality in MATEX is integrated into the Active Directory as dynamic sign-up allowed only for domain customers. This fact does not allow personal computers from outside the domain to be registered. Aging / Scavenging functionality is set, and both options are 7 days. This means, that dynamic DNS record will be automatically deleted after 14 days if the client machine is inactive during this time. (See fig.4)

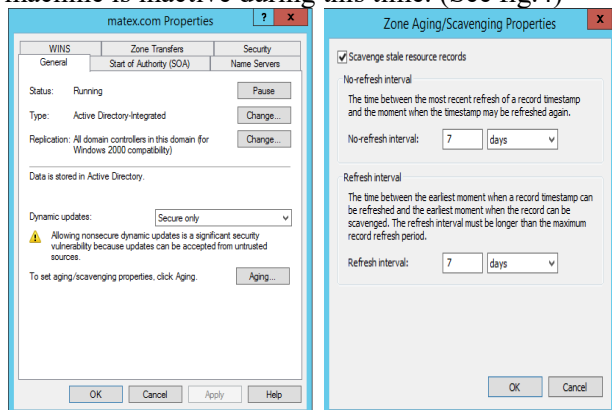


Fig. 4. The Zone Aging / Scavenging properties

Dynamic zone matex.com allows the IT administrators to automatically register and create a record of all Windows machines in the domain. The DNS structure has the additional functionality to create static entries in the local DNS zones. (See fig.5)

BGVARNADC01 has installed a DHCP feature and a DHCP scope is set to deliver IP addresses and network settings to customers in the MATEX network. The first 20 addresses are not distributed, but they are reserved for devices that

require a static IP address. The rest of the IP addresses (from 192.168.159.20 to 192.168.159.254) can be taken by customers. (See fig.6)

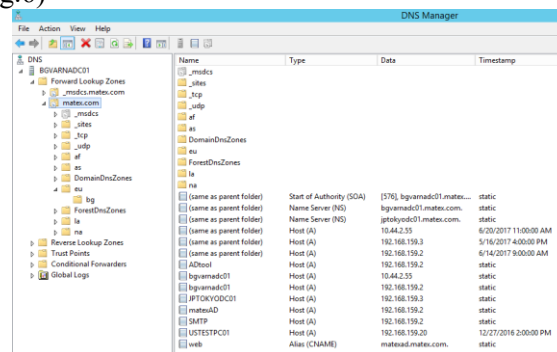


Fig. 5. The DNS Manager

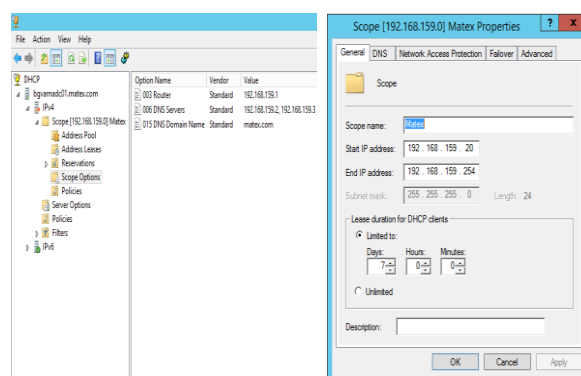


Fig. 6. The DHCP service

Since the server operating system is W2012R2, a DHCP failover is set up between them. This means that if the leading DHCP server (in case on fig.7 BGVARNADC01) is inaccessible for 5 minutes, JPTOKYODC01 will activate MATEX scope on himself and will start distributing addresses. The idea is not to disrupt the user's ability to work. When BGVARNADC01 becomes available, roles on both servers will be re-exchanged.

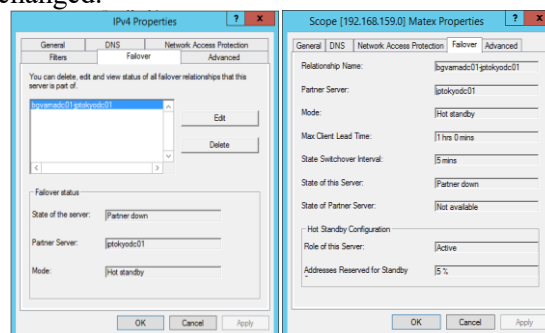


Fig. 7. The DHCP failover

Both servers are also used as file servers where a group folder structure has been created for each department. (fig.8)

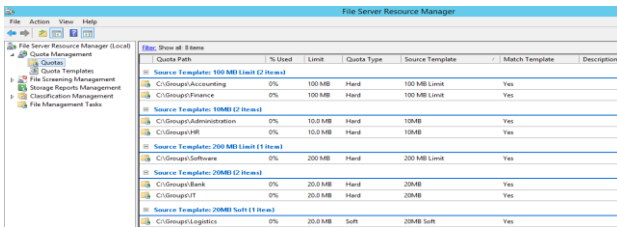


Fig. 8. The File Server Resource Manager

The feature File Server Resource Manager is installed to manage directory contents and data size. Quotas are set for each folder. Hard quotas do not allow the limit to be exceeded, and soft quotas allow only by informing them that the limit is exceeded. Each quota is set to send a notification by mail if the folder size exceeds 90%. (fig.9).

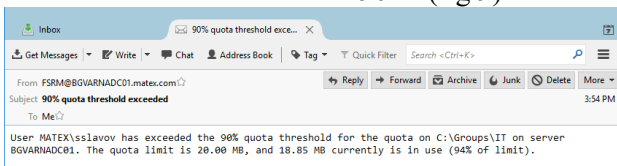


Fig. 9. The Exceeded Quotas in the File Server Resource Manager

File screening is also set up, which prohibits the uploading of audio and video files to any of the servers. If a user attempts to write such a file, the system will not allow and will send an email to the server administrator.

The group directories on both servers are backup by the built-in Windows Server Backup feature. It is set every Friday to make a backup of the group structure to the destination - backup folder on the server. After the completion of each backup, an email is sent with the result - successful or not.

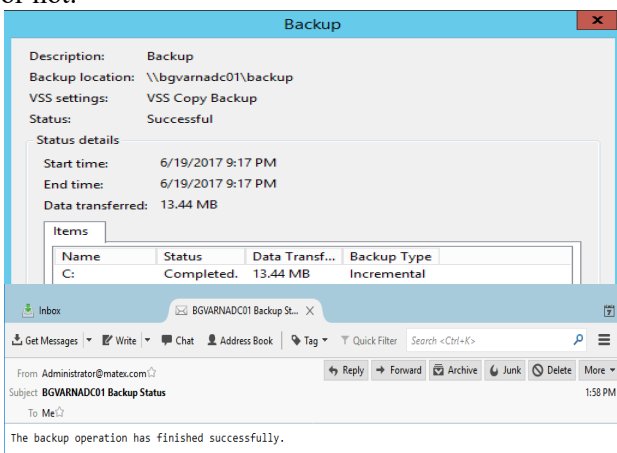


Fig. 10. The Windows Server Backup

To install all the updates from Microsoft to the operating systems in the domain, Windows Server Update Services (WSUS) feature of JPTOKYODC01 is install. (Fig.11) Approved for installation and download are all

patches for the Windows 10 client operating system and for Windows Server 2012 R2 servers.

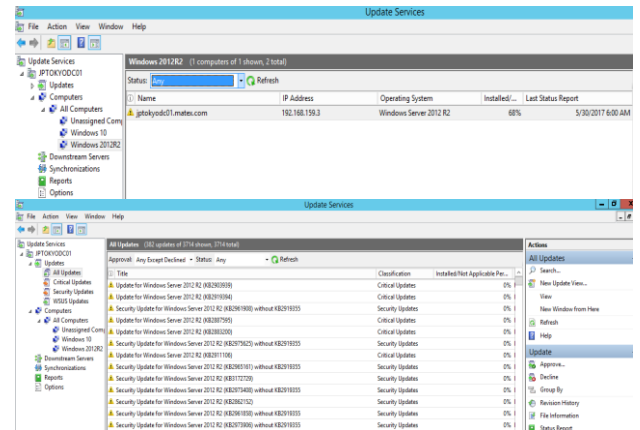


Fig. 11. The Windows Server Update Services

The MATEX features offers:

- Manage (create / delete) user accounts in the domain. When creating, fill in the full name, the user name account, password, OU where is account, e-mail, and phone number are created. When deleting, only the username is filled in.(fig.12-left)
- AD user information, which includes when the last password was changed and the account was locked. It can help administrators in troubleshooting of an user account issue. This information can be send to e-mail to the administrator.(fig.12-right)

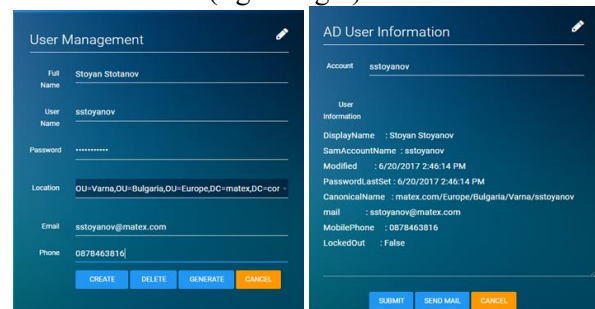


Fig. 12. Manage of users

- Manage (add / remove) rights to manage the various objects in OUs. There is a possibility to grant user rights to manage only computer objects, only groups, only users in the corresponding OU or manage all objects in the OU (OU Management).(fig.13-left)
- The MATEX has the functionality to create County (2nd) level OUs and Site (3rd) level OUs, while creating and empowering the four delegation groups - OU Management, Computer Account Management, User Account Management and Group Account Management. (fig.13-right)

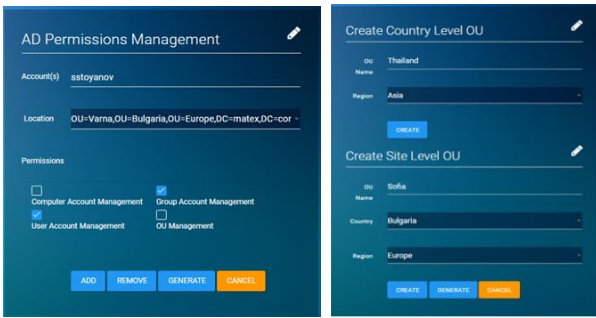


Fig. 13. AD Permissions Management

- Manage access to the group directories on file servers and generate all group folders on a particular file server, select the required folder(s), select the account(s), which should be given / removed rights and the type of rights (FULL or READ).

5. Experimental Results

To determine the performance of the MATEX are applied various tests, by simulating employability for a certain number of users, who use the site at the same time. The numbers of requests per second are presented on figures 14,15,16,17.

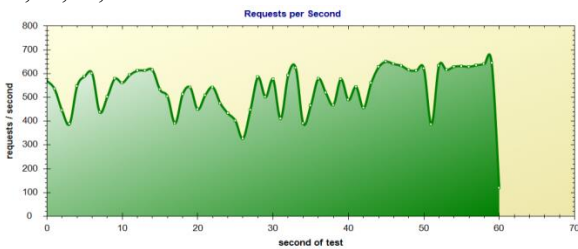


Fig. 14. Time: 60 seconds / Users: 5

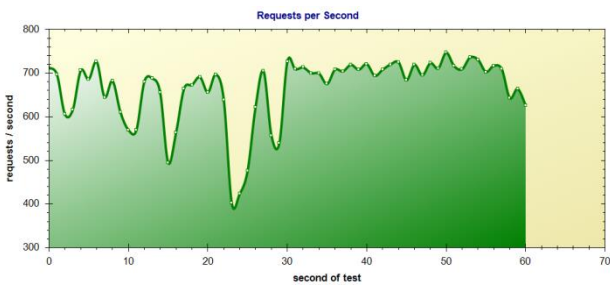


Fig. 15. Time: 60 seconds / Users: 10

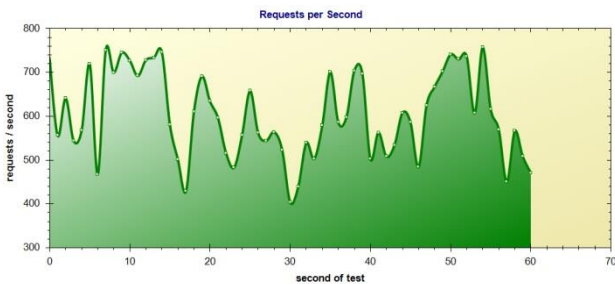


Fig. 16. Time: 60 seconds / Users: 15

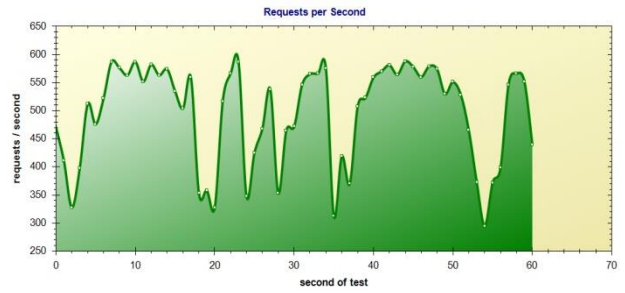


Fig. 17. Time: 60 seconds / Users: 20

The MATEX manages simultaneous queries of 5, 10 and 15 users simultaneously at the same time. In case of 20 simultaneously active users, there are already observed moments where productivity falls significantly at certain times.

All web pages and scripts from the MATEX load for approximately the same amount of time for multiple user queries. The time taken per requests are presented on figures 18,19,20,21.

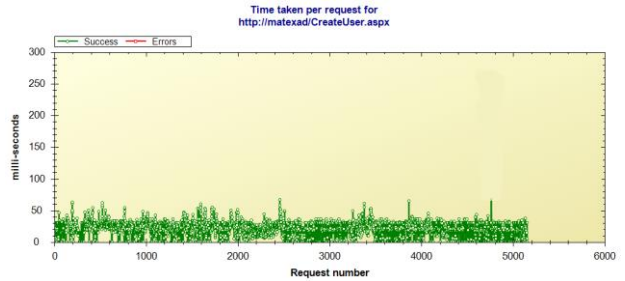


Fig. 18. Time taken per request for <http://matexad/CreateUser.aspx>

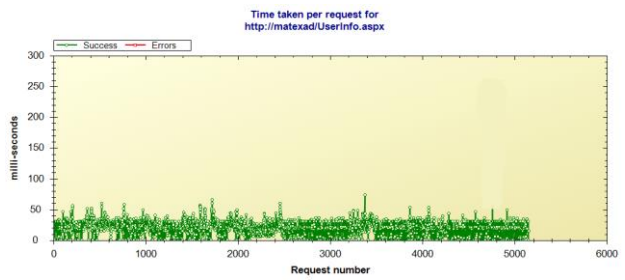


Fig. 19. Time taken per request for <http://matexad/UserInfo.aspx>

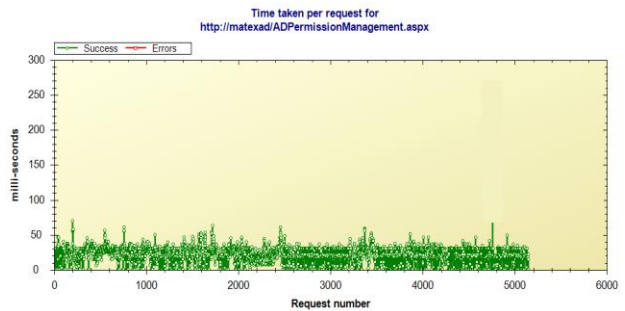


Fig. 20. Time taken per request for <http://matexad/ADPermissionManagement.aspx>

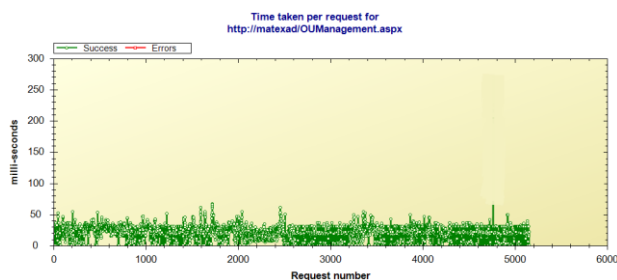


Fig. 21. Time taken per request for <http://matexad/OUManagement.aspx>

The average time of multiple "clicks" was measured for a different number of users (5,10,15) made at the same time. The results are presented in the Table 1. For the number of users 5 (fig.22), 10 (fig.23), and 15(fig.24) the average measured time of multiple "clicks" is very small.

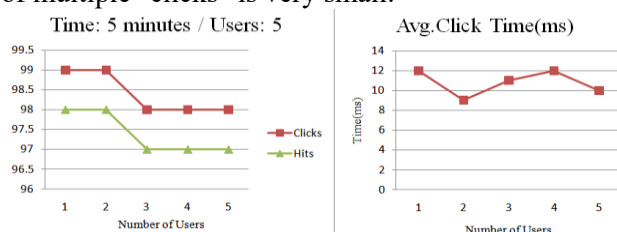


Fig. 22. The average time of multiple "clicks" for 5 users

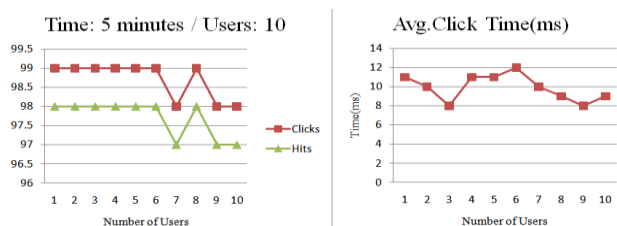


Fig. 23. The average time of multiple "clicks" for 10 users

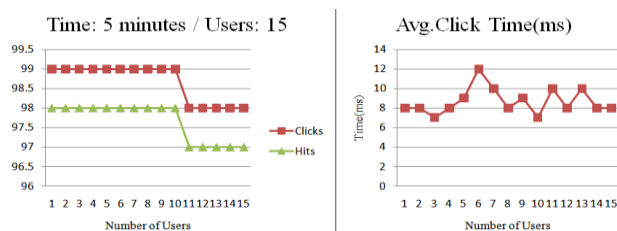


Fig. 24. The average time of multiple "clicks" for 15 users

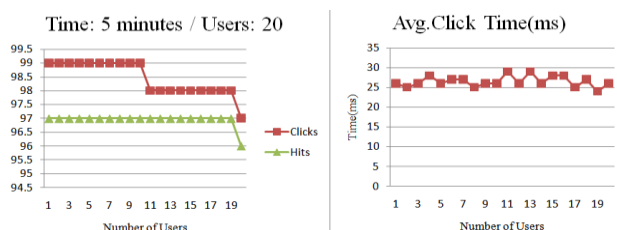


Fig. 25. The average time of multiple "clicks" for 20 users

For 20 users (fig.25), the average time doubled. The conclusion is that 20 is the limit of both active users who use the MATEX, because over this number the productivity and performance slow down.

6. Conclusion

In this paper is proposed an web-based management system for AD, which provided a seamless and simple experience to the IT administrators. To determine the performance of the MATEX are applied various tests, by simulating employability for a certain number of users, who use the web-based system in the same time. It works well with up to 20 active users simultaneously.

REFERENCES

1. Bhargava R., What is DaaS – A Directory-as-a-Service®, <https://jumpcloud.com/blog/what-is-daas-directory-as-a-service/>, September 2014.
2. Bluhm N., Managed Active Directory®(AD), <https://securityboulevard.com/2018/01/manage-d-active-directory-ad/>, January 2018.
3. GEANT, <https://www.geant.org>
4. Start Building on AWS Today, <https://aws.amazon.com/>, [last visited on March 2018]
5. Posey B., How To Create an AWS-Based Active Directory Forest, <https://virtualizationreview.com/articles/2018/01/08/create-an-aws-based-ad-forest.aspx>, January 2018.
6. Gilbert J., M. Tillman, What is Azure Active Directory?, <https://docs.microsoft.com/en-us/azure/active-directory/active-directory-what-is>, July 2017.
7. Lujan V., Web Based Directory Service Platform, <https://jumpcloud.com/blog/web-based-directory-service-platform/>, January 2018.

Veneta Aleksieva, Svetoslav Slavov
 Organization: Technical University of Varna
 Address: Str. Studentska 1, Varna , 9010
 Phone (optional): +35953383439
 E-mail: valeksieva@tu-varna.bg,
 sslavov@lirex.bg

OVERCOMING THE SECURITY ISSUES OF NOSQL DATABASES

TONY KARAVASILEV, ELENA SOMOVA

Abstract: *With the current escalating popularity and use of NoSQL databases, the amount of sensitive data stored in these types of systems is increasing significantly, which exposes a lot of security vulnerabilities, threats and risks. This paper presents effective ways to mitigate or even completely overcome them. The purpose of the developed practical tests using MongoDB is to evaluate how applying those security measures can affect the overall system performance. The results of this experimental research are presented in this article.*

Key words: *security, encryption, databases, NoSQL, MongoDB, RESTful API, firewall*

1. Introduction

Due to the spreading use of modern cloud computing solutions and the increasingly larger data volumes for storage, the adoption of a new class of non-relational databases has arisen, also known as NoSQL or "Not Only SQL". Such kind of databases have existed even before relational and object-oriented database management systems but was resurrected and developed in the recent years by information technology companies for providing private problem solutions for their growing distributed web applications with millions of users. The main advantage of these databases is that they cope up with processing and storing unstructured data way better than standard relational SQL solutions. [1]

The simplicity of their design, schema-less models and primitive query languages allows them to perform, scale and distribute much better in certain situations. These databases provide an alternative for storing some types of data more efficiently than other ones and support advanced clustering solutions, including load balancing and transparent backup features.

Using a NoSQL database depends entirely on the problem it must solve and the data type it is going to store. Depending on the type of unstructured data, there are four main subclasses of non-relation databases to choose from:

- Document-oriented – document string formats;
- Column-oriented – column or tabular nesting;
- Graph-based – graph structures with nodes;
- Key-value store – unique associative arrays.

Some solution can also provide more than one model available and are called multi-model

systems, this may or may not include relational functionalities. Other traditional SQL solutions have merged a lot of the NoSQL capabilities (clustering, sharding, XML/JSON document objects and linear structures like the array type) and are forming a new database class called NewSQL, unfortunately inheriting a lot of the non-relational security and data integrity problems. Clearly, there are no set standards and these databases allow a bit more flexibility. This the main reasons why they are typically used for caching purposes, search engine implementations, file storage and activity logging.

Unlike relational database management systems, these types of software products do not have complex mechanisms to guarantee data consistency and have almost no support for security features at the database level. Having that said, this makes them vulnerable to both security threats and irreversible data loss, which is a serious problem building up over the years. [2]

The main purpose of this article is pointing out the security risks introduced by the use of modern non-relational databases and effective ways to mitigate or completely overcome them. There are multiple issues pointed out and a detailed explanation of how they can be eliminated with the use of end-to-end encryption and a security-oriented configuration of all services. This paper includes a practical implementation of how to integrate the explained security precautions and extended tests showing how this affects the overall system performance and the data storage volume size.

This research is also a part of ongoing work in developing an advanced PHP object-oriented software framework for cryptographic services.

2. NoSQL security issues and their remedies

When using NoSQL as a solution for dealing with sensitive and top secret data in the real world a few problems may occur. Even when using paid non-relational database systems, paid vendor support and hiring well-paid professional database administrators, you can introduce big security holes in your system and compromise the overall data privacy, putting your company at imminent risk.

The next sections show the most frequent issues that come with these type of systems and effective ways of overcoming them completely.

2.1. Lack of authorization features

In general, most professional NoSQL solutions have either only basic access control mechanism or do not support any at all. This proposes a huge problem in the overall application security and leaves out an open possibility of hostile access without any credentials or restrictions set. [1]

A common mistake that big data fans do is to use the default product installation credentials and configuration. It is a must to enable and use strong authentication credential. In the case where there is a huge lack of authorization functionalities, access control restrictions, user role capabilities and auditing features in the chosen NoSQL system, the only thing we can do to mitigate the problem is to build a RESTful API (Application Programming Interface) around our database solution. A lot of in-memory key-value solutions and search engines suffer from this problem by design. Building our own application layer of access restriction on top of an unprotected system is a common professional security approach. Also, implementing access tokens as credentials for the Web API clients is a great way of boosting security.

Overcoming this security issue is obligatory and must not be underestimated. It is important to note that message queue broker systems also suffer from the same authentication security problems, although they are not technically database solutions.

2.2. Transport encryption and client drivers

Sadly, a lot of modern NoSQL products do not provide network transport layer encryption over TLS/SSL. The lack of this security trait is both on server and client side. The support of different programming language consumption drivers is not good enough and may introduce the risk of data corruption, theft or even loss. Faulty drivers may injure the system stability and performance a lot. [2]

If your desired solution supports transport encryption, always use TLS/SSL for client-server communications and try to avoid the use of self-signed generated certificates. When your NoSQL has no such capabilities implemented, your best

option is to restrict non-encrypted connections via the use of firewalls and develop a RESTful API on top of your database which accepts only communications via HTTPS (HTTP over TLS/SSL) encrypted traffic, using the most supported client programming language and up to date drivers.

Mitigating this security threat actually upgrades the proposed solution in the last authorization section by forcing the use of encrypted communication protocols and highly supported client library connection drivers.

2.3. Missing database encryption features

Most of the relational database systems have built-in encryption storage engines, encryption aggregate functions and data integrity features. In the other corner, modern NoSQL solutions lack such at-rest encryption functionalities and store data as plain text which imposes too many security risks.

There are only a few paid cloud or enterprise NoSQL products that support native storage encryption engines and embedded recovery functionalities. Mitigating this problem can be done by developing a transparent encryption application layer. This kind of software layer encrypts data before sending it to the database and decrypts it before returning it to your software. The transparent middleware can either be implemented directly in your application's client connection internal library or on another server, acting as a communication proxy between the database and the application. The only limitation created by using this approach is that you cannot search directly inside the encrypted fields via the database query language. [3]

A good approach is to include this feature when building a RESTful API. It is also encouraging to combine it with native storage encryption functions when there are such available.

2.4. NoSQL Injections and CSRF attacks

With the emergence of new query formats and languages, most of the old SQL injection techniques are pointless but this does not make NoSQL immune to query injections. Every non-relation product can be attacked or exploited, depending on the type of your database language query, used message formats and its available native application programming interfaces. [4]

The main cause for the rise of these NoSQL injection attacks is because a lot of non-relational database systems provide either an embedded RESTful API or load one via third-party extensions, that use JSON or XML formats. This can enable an attacker to execute valid malicious code, such as native JavaScript code, via the request's payload data. As spoken before, building your own custom RESTful API on top of the database is a common

trait but if done incorrectly may easily be exploited via NoSQL injections. Also, it is possible to accomplish some cross-site request forgery attacks (CSRF) and user session hijacking.

The best way of avoiding this kind of attacks is via input sanitization, especially when creating your own RESTful API. The cleansing of the received data includes validations, filtering, whitelisting, blacklisting, regular expressions and escaping of special characters. Another good practice is adding a per client pseudo-randomly generated token with time-to-live (TTL) expiration or regeneration period for avoiding forged user requests. Every request must contain a valid token or the request will not be handled and the user may become suspended after a given number of retries. The disabling or firewall blocking of unused native APIs and extensions is always recommended.

Mitigating this security risk upgrades the previously proposed solution by adding extra input validations and further request integrity verification enchantments.

2.5. Cluster desynchronization issues

Most NoSQL products advertise their decentralization features like sharding and clustering as a data-friendly way to scale out, containing no single point of failure. It is important to note that almost all SQL solutions support these features but have more reliable and time-tested realizations than non-relation systems. [2]

When using sharding, each data partition or shard is held on a separate database server instance and distributing the data on more than one machine. Some features can include storing duplicates of the shards on other servers for backup reasons. The main problem is that even if one server crashes you may end up with losing a certain amount of shards without any real backup and compromise the overall stored information. Knowing this, a malicious attacker may exploit a misconfigured cluster or penetrate known vendor vulnerabilities to destroy, hold for ransom or modify your data.

As previously noted, another traditional solution for scaling out, load balancing or creating active backups is called clustering. Supported types may include master-slave replication, master-master replication, shared-nothing clustering, auto-sharding, hybrid storage and other high availability distributed approaches. Most of those clustering techniques provide a certain amount of guarantee for data integrity, backup and availability on hardware or software failures. There are a lot of hidden problems like cluster desynchronization and some failover dead-locking situations where your information gets corrupted or even irreversible lost. This may lead to security leaks where sensitive data

is being returned to users or gets permanently defective without any real backup.

As with any other product, there is a huge gap between what is being advertised or sold and what you get in reality. A large portion of the NoSQL community follows certain solutions for philosophical reasons rather than practically proven production use cases. In theory, both SQL and NoSQL clustering solutions can fully eliminate all failure cases but in reality, when misused, misconfigured or not implemented correctly they can create even bigger security and integrity issues.

The contra measures you can take involve hourly backups to at least two separate physical devices and creating simulations of all know crisis situations before using the clustering configuration in a production environment. Also, keep your database solution up to date, always encrypt your backups and never store them on the same production machine. Only then you can fully harvest their true performance and backup gains.

2.6. Virtualization leaks and disk theft risks

Even when you have encrypted your database and have taken security measures you are still not immune to physical disk theft or virtualization snapshot leaks. [5]

The security threats of stealing the physical disk or the virtual machine backup clone files involve gaining database credentials or sensitive data via log files analysis, raw cache files, unencrypted database diagnostic tables or persistent in-memory data structures. Other risks include gaining access to guest virtual machine clones or virtualization snapshots which contain memory dumps of a passed machine state and are full of unencrypted database structures, active in-memory key-value collections or even loaded application credentials.

To mitigate these security issues you must apply transparent disk encryption on all physical disks, virtualization host environments and virtual machine guests with a strong encryption key. This operating system feature can be used without causing any decrease in the overall system performance. Note that if you lose or forget your secret password, you will not be able to start up your operating system or restore any usable data from your drive.

3. Analyzing and improving the MongoDB database security

With the increasing use of MongoDB in both startups and enterprise solutions worldwide and also being one of the most feature-rich NoSQL databases the need of implementing security protection has become huge. The next sections

provide a detailed practical analysis of security hotspots and how to tighten up the overall database protection using previously discussed approaches.

3.1. Overview of MongoDB features

MongoDB is the most popular document-oriented store that uses JSON (JavaScript Object Notation) format documents providing flexible and easily changeable schemas. It also provides server-side scripting with JavaScript and binary-encoded serialization of JSON-like documents. This open-source NoSQL software is provided for free and also has an enterprise paid version with extended features and live support. [6]

MongoDB provides a huge variety of high availability clustering features like load balancing, replication and sharding. In some cases it also can be used efficiently as a file storage server or powerful caching system. The query language supports aggregation, range queries, regular expressions and different field indexing types.

The main problem is that the default security configuration of MongoDB has been exploited many times in production setups during the years and even held for ransom. The next sections will discuss how to avoid a security breach by taking certain precautions.

3.2. Enforcing authorization, auditing and input data sanitization

After the product installation, the database access is publicly exposed without any credentials or verification configured, making it easy for anyone to connect and take full control of the database. A lot of system architects neglect MongoDB's initial configuration and are commonly hacked for it.

To avoid this, you must enable all native authentication features. First of all, you have to add a user administrator for the MongoDB instance and define different limited access roles for every other client account that can connect to the database. Next, you must enable the native system auditing facility for keeping track of all configuration changes and log access history. For even further hardening, you can set running of the database processes with a dedicated operating system user account that has limited permissions. In some cases, disabling server-side scripting on database level can remove the possibility of some types of NoSQL. Also, if you use a cluster setup then never forget to define proper authentication between cluster members and always use long complex credentials.

Finally, to reduce the access, even more, it is a great idea to develop an internal RESTful API that connects to the database by using only a limited user account. Also, when developing such adapter

software, you can completely sanitize your data input, use advanced session forgery protection techniques and create complicated authentication features. Remember to allow only direct connections from the internal API and block all native MongoDB client communications via network or system firewall.

This way you ensure that malicious code execution or unauthorized access to the database is not possible and will not disable the use of native MongoDB high availability cluster configurations. Also, you can reliably scale out and increase performance by deploying multiple instances of your API and using a network traffic load balancer for distributing the incoming request between them.

It is important to note that unlike MongoDB a huge amount of the NoSQL solutions do not provide even basic authentication features. Either way, you would have no real choice but to develop your own internal database adapter software.

3.3. Using encrypted communications and limiting network exposure

When installing the product, a lot of people leave out the default access port and connection protocol publicly exposed. You must always change the default port and switch to encrypted TLS/SSL communications for both all your cluster servers and client machines. This way your data is protected in-transfer and cannot be altered by man-in-the-middle attacks (MITM).

You must never leave a MongoDB server instance visible over the Internet or accessible in non-management computer networks. You can even disable the MongoDB networking service and switch to using UNIX sockets instead, especially if you are going to use only one server instance and hide it behind an isolated RESTful API.

However, enabling the network is a must when using clustering configurations and the most professional way of protecting any type of server instances is to combine the use of network firewalls with Virtual Local Area Networks (VLAN). This way you can partition and isolate different networks with limited access to other devices or computers.

Always limit the network exposure on every service you use and switch to the use of encrypted protocol connections only. Remember, a service that is not visible or accessible cannot be easily exploited, flooded or hacked.

3.4. Applying data storage encryption

The data protection at-rest is truly important but most databases do not provide native encryption functions or secure storage engines. For example, MongoDB provides native encryption only for its enterprise paid version since a few years back. It is

recommended to use it when available and also develop a transparent encryption middleware. [6]

When creating such encryption application layer, you must always encrypt sensitive fields before inserting into the database and decrypt after fetching them back. This can be included in your RESTful API logic. The main limitation of using such middleware is that searching inside encrypted document values is not possible without having to first decrypt them all.

Another good habit is to always enable the operating system transparent disk encryption to ensure data and system logs safety even if a physical disk theft occurs. The use of this feature will not harm performance and significantly increase the overall server security. You can also encrypt service log files over time, implement encrypted application file logging adapter classes and minimize the amount of service related details.

4. Testing environment specification

For the results to be adequate we have chosen to run the tests in a virtual machine environment created with Oracle VM VirtualBox version 5.2.8 hypervisor. The setup consists of two virtual machines. The first one is running Apache 2.4.18 with PHP 7.2.3-FPM (FastCGI Process Manager implementation) for connecting to the database and executing the experiments. The second one has a single MongoDB 3.6.3 server instance for the data storage purposes.

The specification of the allocated resources for each of the machines is shown in Table 1.

Table 1. Virtual machine specification

	Detail
CPU	Intel i7-6700HQ, 2 cores, 2.59GHz, 6 MB L3
RAM	DDR4, SODIMM, 4096 MB, 2.40 MHz
GPU	Intel HD Graphics 530, 16 MB, 2.40 MHz
HDD	42GB, 7200 RPM, 32 MB cache, 2GB swap
LAN	VirtualBox Intel PRO/1000 MT 82540EM
OS	Ubuntu Server 16.04.3 LTS x64, Kernel 4.4.0

The virtual machines have been installed with all available updates, kernel drivers and virtualization-specific packages. The connection between them will be over the host-only networking mode embedded in VirtualBox to avoid any network slowness and ensure the executed tests accuracy. On the first virtual machine, all settings are set by default, with the exception of boosting the values for maximum random-access memory (RAM) usage for PHP. The second machine will be configured twice for every experiment. All tests will first be executed with the default insecure configuration. Next, the tests will be repeated with MongoDB authorization enabled, using a self-

signed TSL/SSL certificate for communication and applying an encryption middleware via PHP.

The created encryption middleware uses the AES-256 CTR algorithm via the OpenSSL PHP [7] native extension functions and the Base64 encoding core functions for converting to a storage-friendly format. The main purpose of this setup is to simulate both plain and encrypted pseudo-API to MongoDB communications and to compare results.

5. Costs of implementing security measures

Although applying end-to-end encryption is a must, are there any consequences of using it in your production environment? To answer this question, we have created several practical experiments using MongoDB to evaluate the performance and storage costs. The next sections describe the executed tests and show their results.

5.1. Experiments suite overview

The experiments include the insertion of 10000000 (ten million) records containing pseudo-random cryptographically generated strings with a fixed length of 1000 printable ASCII characters and the fetching of 20 records from the middle of the collection via an extra inserted 10-digit integer field. The integer field will act as a unique creation identifier for easier lookup and will be created with an ascending index. Since we will be retrieving records from the middle of the collection, it would not matter if we use an ascending or descending index lookup. We will also leave the default MongoDB “_id” field creation but use projection when querying the database to not get it with other results. Also, all encrypted data will be converted to Base64 strings for storage in MongoDB documents.

The time spent executing a test shown is just for the section of the program that does iteration, encryption, decryption, data insertion, collection lookup and records retrieval. The time needed for generating cryptographically the pseudo-random strings is explicitly excluded. Every single experiment is executed 10 times and the average result of those runs is taken as final. Execution time results will be shown in seconds with 6-digit precision after the decimal point and storage size results will be displayed in bytes.

5.2. Record insertion results

This experiment will test the situations when you need to store big string data and compare the average insertion time from PHP and the record storage size on disk. The results for both plain record creation and using the transparent encryption middleware are shown in Table 2.

Table 2. *Ten million records insertion*

	Plain	With Encryption
Total Time	1493.585353	1620.437595
Average Time	0.000149	0.000162
Average Object	1052	1388
Index Storage	212541440	212492588
Collection Size	10629316608	14346072064
Database Size	11156746240	14977986560

As we can see from the results, when applying encryption, the overall database storage size has significantly increased by 34.25%. Also, the total record creation time has become with 8.49% slower but is still rapid.

Having in mind that we would probably encrypt only sensitive data fields like passwords and credit card numbers, the cost of applying encryption is relatively tolerable.

5.3. Record retrieval results

The second test will apply to the scenarios where we need to execute a complex query lookup in huge data collections like the created ones in the previous experiment. To simulate this, we will query the database to fetch the first twenty records after the fifth million record, using our creation identifier. After that, drop the created extra index for our identifier and run the query again to see a more precise contrast between plain and application decryption fetching. The results for record retrieval experiments are shown in Table 3.

Table 3. *Twenty records retrieval*

	Plain	With Decryption
With index	0.000644	0.000855
Without index	11.763420	16.326245

The time increase caused by the use of application decryption with index fetching is 32.76% and with non-index retrieval is 38.79%. It is important to note that the experiment also showed the huge performance boost of using field indexing.

The results show that the application decryption will not slow us down significantly when the correct schema approach is being applied.

6. Conclusion

This paper has created a practical analysis of NoSQL database solutions and evaluated the

performance and storage costs of applying end-to-end encryption. It summarizes the best approaches to overcoming common NoSQL security problems.

The most interesting results from the experiments are:

- Building an API around the NoSQL solution in a security-driven matter and setting up an isolated network is the best safeguard approach;
- Using a transparent encryption middleware can increase the disk storage size significantly but does not hurt the overall system performance;
- Searching directly inside encrypted fields via the database query language is only available when using native database encryption engines;
- Querying the database by index scan is more than 18000 times faster than using full scan.

REFERENCES

1. J. Sadalage, P., and Fowler, M. (2012). NoSQL Distilled: A Brief Guide to the Emerging World of Polyglot Persistence. ISBN-13: 978-0321826626.
2. Okman, L., Gal-Oz, N., Gonen, Y., Gudes, Eh., and Abramov, J. (2011). Security Issues in NoSQL Databases. DOI: 10.1109/TrustCom.2011.70.
3. Tian, X., Huang, B., and Wu, M. (2014). A transparent middleware for encrypting data in MongoDB. DOI: 10.1109/IWECA.2014.6845768.
4. Ron, Av., Shulman-Peleg, Al., and Bronshtein, Em. (2015). No SQL, No Injection? Examining NoSQL Security. arXiv:1506.04082 [cs.CR].
5. Grubbs, P., Ristenpart, Th., and Shmatikov, V. (2017). Why Your Encrypted Database Is Not Secure. DOI: 10.1145/3102980.3103007.
6. <https://docs.mongodb.com/manual/index.html>
7. <http://php.net/manual/en/ref.openssl.php>

Contacts:

UNIVERSITY OF PLOVDIV PAISI
HILENDARSKI
24 TZAR ASEN
PLOVDIV

E-mail: tony.karavasilev@gmail.com

E-mail: eledel@uni-plovdiv.bg

VIRTUAIZATION AND CONTAINERIZATION SYSTEMS FOR BIG DATA

DANIEL TFRIFONOV, HRISTO VALCHANOV

Abstract: *In the case of processing unstructured, semi-structured and structured data as well as research data of historical and statistical nature, the traditional relational database management systems are not a rational choice and have been replaced by other solutions such as large data storage systems - Big data. Distributed data processing systems such as Apache Hadoop are better solution. Building such a system requires the availability of multiple machines, which is a serious investment even for the large companies. The use of virtualization platforms can solve a number of existing problems. Increasingly, however, an attractive technology emerges as an alternative to virtualization - the containerization technology. Containers solve some of the problems typical of hypervisors and virtual machines. The performance comparison of large data processing system implemented on virtual machines and containers is presented in this paper.*

Key words: *Big data, Apache Hadoop, Virtualization, Containerization*

1. Introduction

The amount of generated data daily grows faster. All this amount of information is required not only to be stored but also processed. Data becomes more and more diverse in nature and less structured. In the case of processing unstructured, semi-structured and structured data as well as research data of historical and statistical nature, the traditional relational database management systems are not a rational choice and have been replaced by other solutions such as large data storage systems - Big data. Big data is used not only by all software giants like Google, Microsoft, AWS, Facebook, Ebay, Booking, New York Stock Exchange, but also by many smaller companies and research centers [1].

The desktop analytics tools and relational databases often have problems dealing with large amounts of heterogeneous data, so distributed data processing systems such as Apache Hadoop [2] are better solution. Hadoop is a cluster system of a small number of controlling nodes and a large number of computing and data storage nodes. This robust distribution allows calculations on huge volumes of data to be done for a very short time, with each node processing only its local data.

Developing of such a system is a serious investment even for large companies, because dozen machines are needed for achieving benefit from the distributed processing. The use of virtualization platforms (VPs) can solve a number of existing problems. This can be done in several

directions. These platforms provide economical calculations, reducing the need for investment in new equipment. At the same time, VPs reduce the cost of support, including reinstalling operating systems and software. Increasingly, however, an attractive technology emerges as an alternative to VP - Container Technology. Containers solve some of the problems typical for hypervisors and virtual machines. Because of their simple architecture they provide better performance than virtual machines. At the same time, they allow fastest and more flexible providing of resources and availability of new applications.

This paper presents a study of the performance of large volume data processing systems based on virtual machines and containers.

2. Virtualization

Virtualization is defined as "an abstract presentation of the computer's physical resources on a virtual computer using specialized software" [3]. The virtualization platform virtualizes the hardware resources of the computers. It creates a fully functional virtual machine on which a standalone operating system and applications can run, just like on a real computer.

Virtualization is implemented with dedicated hypervisor software that distributes the computer hardware resources dynamically and transparently between virtual machines. Several operating systems can work simultaneously on a physical computer and share hardware resources

with each other. By encapsulating the entire machine, including a processor, memory, operating system and network devices, the virtual machine is fully compatible with all standard x86 operating systems, applications and drivers.

There are two main types of virtualization environments: hosted and bare-metal (Fig. 1).

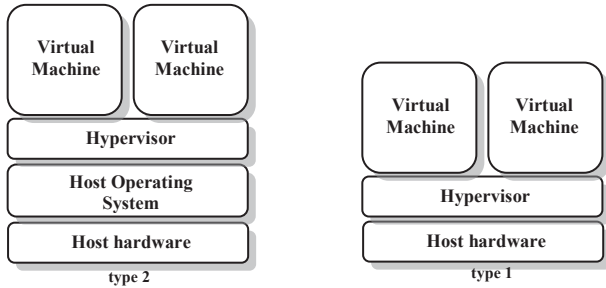


Fig. 1. Hypervisor types

In hosted environments (type 2) hypervisors are software applications running under the operating system. The hypervisor controls the lower-level resources that are assigned by the operating system. This type of hypervisors is mainly used in systems where different I/O devices are needed and they can be supported by the host operating system. Another usage is in low performance client systems. Examples of such type of hypervisors are: Microsoft Virtual Server [4], VMware Server, and VMware Workstation [5].

In bare-metal environments (type 1) the hypervisors are software systems that work directly on host hardware. This results in higher performance and performance. This type of hypervisors is a preferred virtualization approach. Examples of such hypervisors are: Citrix XenServer [6], VMware ESXi [7], Microsoft Hyper-V [8].

3. Containerization

The difference between virtualization and containerization is mainly at the place of the virtualization layer and the way the system resources are used. Containerization, also called “container-based virtualization”, “para-virtualization” or “application virtualization”, is a virtualization method for deploying and running of distributed applications at the operating system level without the need to run an entire virtual machine for each application. Instead, multiple isolated systems, called containers, are run on a single host and have access to the kernel of the operating system (Fig. 2).

The container is a lightweight, stand-alone, executable software package that includes everything it needed to do: code, libraries, system applications and settings. The software inside the container is executed in the same way, no matter the

environment. Containers isolate the software from the environment, for example between the development environment and the work environment. They help to reduce conflicts between teams that use different software on the same infrastructure.

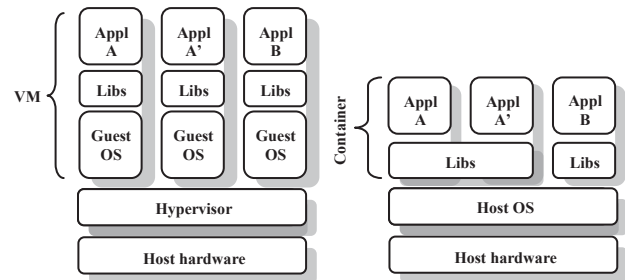


Fig. 2. System with virtual machines and with containers

Because the containers share the same operating system kernel with the host machine, they can be more effective than virtual machines, which require separate operating systems. The host operating system restricts the access of the container to the physical resources such as processor and memory, so a container cannot consume all system resources. The containers start faster and use less computing time and memory than virtual machines. Containers share common files that reduce disk space usage. There are a number of container-based solutions such as Linux-VServer [9], OpenVZ [10], Docker [11].

4. Big data processing

One of the widely used systems for processing large volumes of data is Hadoop. Hadoop is an open-source software developed in Java. Hadoop has a number of code execution tools and scripts in different programming languages. It consists of two main parts - a storage part and a data processing part. The storage component is the distributed file system - Hadoop Distributed File System (HDFS). It has a block organization, each block having a size of 256MB. The blocks that make up a file are distributed across all nodes of the distributed system. In order to achieve a redundancy, Hadoop maintains block replications. The data store is Hive. Hive offers easy data aggregation, temporary queries and other large data analyzes. Queries use a language similar to SQL, known as HiveQL.

The data processing part is MapReduce. This is a programming model for parallel processing of multiple tasks. One of the main aspects of MapReduce programming is that MapReduce splits the tasks in such a way that they allow their parallel execution on a distributed system of computing nodes. Contrary to traditional

relational database management systems that cannot grow to handle large amounts of data, the programming in the Hadoop MapReduce environment allows users to run applications on a huge number of machines, which also includes the processing of thousands of terabytes data.

5. Experimental study and results

The VMware ESXi 6.5 hypervisor is chosen as a platform for virtualization, and the Docker as a container. The choice of these two platforms is dictated by their wide use, high productivity and capabilities.

The test environment is built on the DELL CS24-TY server system, which has the following hardware:

- 2 CPU 4 cores Xeon L5520 with Hyper Threading technology;
- 72Gb RAM;
- 1,8 TB storage.

On the machine, the two platforms are installed sequentially and the corresponding performance tests on the Big Data system are performed on them. There are 5 separate Ubuntu 16.4 LTS (64-bit) virtual machines installed on VMware. In virtual machines, a Hadoop system is installed, which is divided into 3 computational nodes (data nodes), 1 master (name node) and 1 secondary (secondary name node). The Docker is installed on the Ubuntu 16.4 LTS (64-bit) operating system, and 5 separate containers are created. The same Hadoop system is placed in the containers.

The tests are selected to meet the requirements of different hardware devices: processor, memory, and storage devices.

The first test tests the performance of the system when processing data in semi-structured form. Such processing is consistently imposed on Big Data systems, as input data quite often come from heterogeneous sources and are in different formats. The test is a Java application (WordCount) executed directly by Hadoop and calls MapReduce. MapReduce crawls the file and divides the text into separate words by removing the punctuation – the “Map” phase of the task. The next phase is “Reduce”. It reduces the result to a file containing name-value pairs, indicating for each encountered word how many times they are detected. The processing data are one of the latest Wikipedia (EN) archives with only English content used. It is provided as a bz2 archive in which there is xml file of approximately 62GB. Figure 3 shows the results of the test.

The results show an advantage of ~ 4.5% for containerization and Docker. This is expected, given that the maximum amount of memory is used during the test, and each virtual machine uses

approximately 1.3GB of memory that is otherwise used by Hadoop. The access to the disk storage for virtualization is potentially slightly slower, which also slightly increases the lead of the Docker. Another important point is the long time to complete the test. This is due to the low performance of the system drives, and it is possible for higher speed devices to differ from the results obtained in this test setup.

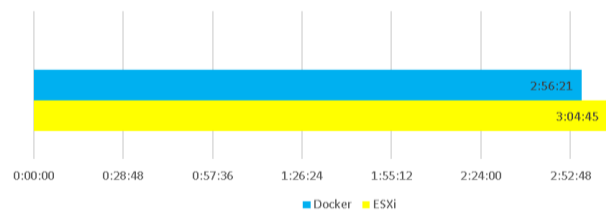


Fig. 3. WordCount test results

The second test is complex in terms of processor work and I/O operations - Hive dataset import. This test uses MapReduce again. It reads the contents of the text file *enwiki-20170701-pages-articles-multistream.xml* and imports it into a non-relation table MongoDB. The number of columns of the table is 2, the first is the title of the article, and the second is the entire text. The text processing requires processor time and memory, while reading from the file and writing to the database loads up the discs. The writing requires 3 times more I/O operations due to the replication factor. The results are shown in Fig.4.

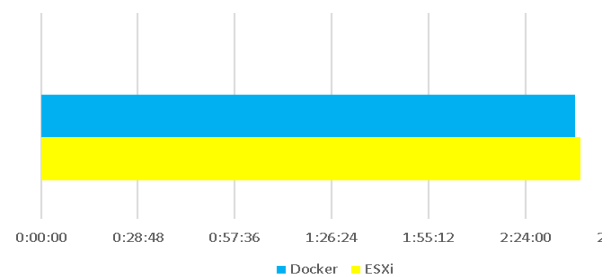


Fig. 4. Hive dataset test results

It should be noted that this test does not give a clear result for the advantage of virtualization or containerization. It runs faster than the first test because the text is stored in a different format in the database, and data processing is limited to retrieving from the text and writing to the database. The difference in WordCount and Hive database import execution time is not great again due to the specificity of disk storage. This test should end much faster if there is a separate hard drive or even SSD for each data node. In such a study of test systems with real and virtual machines where two virtual machines are running on a real one, the

virtual ones perform better because they utilize the processor's use. In contrast, the real ones are not fully loaded at any point during the test time. In this study, no such behavior is observed since the hardware used is a single computer system and the load in both cases is fully.

The third test evaluates the storage system. It runs in two parts: TestDFSIO-write records 100GB of data in the Hadoop-HDFS file system, and TestDFSIO-read reads them back. Because of the replication factor, the write test triggers 3 times more I/O operations than the reading test and generates significant network traffic. The results are presented in Fig.5.

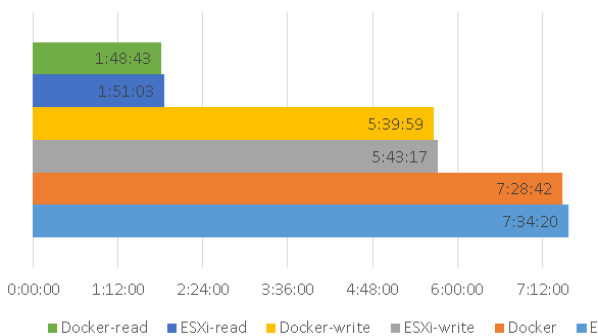


Fig. 5. DFSIO test results

In this test the results are much more interesting. It is clear that the tests with writing are about 3 times slower than those with reading. This is normal and it is precisely because of the replication factor and the specificity of the test cluster. The replication factor is 3, and the same number are the data nodes, i.e. each computing machine takes up the entire volume of data on its disk, while reading is performed simultaneously from the all three machines. If the replication factor remains 3 and the computational nodes are more, then the difference in read and write times will be proportionally smaller. In large systems with hundreds of computing nodes, there is virtually no difference between read and write speeds, although the replication factor there may even reach 5. A slight advantage is again for containerization due to the fact that the disk access through the virtual disk controller is bit more slowly than through the interface provided by containerization.

6. Conclusions

On the basis of the tests carried out, it was found that containerization produced somewhat better results in most cases. If there already have an installed server with a hypervisor, it does not need to be reinstalled with a single operating system and containers because the productivity gap is not that

big. Additionally, other virtual machines can be used in the hypervisor while the cluster is not in use. It should be noted that tests generally go much slower than expected. The reason is in shared two discs for all 5 virtual machines in one case, and containers in the other. As a recommendation for implementation of a system of this type, it is better to use at least one disk and more precisely a SSD drive for each virtual machine or container. Then the performance of the Hadoop system would be more real. The main difference in performance comes from having a larger amount of memory used in the container system because otherwise it is occupied by the operating systems of virtual machines (5 x 1.3GB) and the hypervisor (2.26GB).

As recommendation when building new systems for big data processing, it is profitable to implement them with containers because the performance of hardware is slightly better, and it also makes it easier to administer the system - just one operating system for maintenance.

REFERENCES

1. Want To Use Big Data? <https://www.forbes.com/sites/bernardmarr/2017/08/14/want-to-use-big-data-why-not-start-via-google-facebook-amazon-etc/#5dd460173d5d>.
2. T. White. Hadoop: The Definitive Guide. O'Reilly, 2015.
3. J. Drews. Going Virtual, Network Computing, Vol. 17, No. 9, p. ES 5, 2006.
4. Microsoft Virtual Server. <https://www.microsoft.com/windowsserversystem/virtualserver/>.
5. VMware. <https://www.vmware.com>.
6. Xen Server. <https://xenserver.org/>
7. ESXi. <https://www.vmware.com/products/esxi-and-esx.html>.
8. Server Virtualization. <http://www.microsoft.com>.
9. Linux-VServer. <http://linux-vserver.org/>.
10. OpenVZ Linux Containers Wiki. <http://openvz.org/>.
11. Docker – Build, Ship, and Run Any App, Anywhere. <http://www.docker.com/>.

Contacts

Hristo Valchanov, Daniel Trifonov
 Technical University of Varna
 9010, Varna, 1 Studentska Str.
 phone: +359 52 383 278
 E-mail: hristo@tu-varna.bg

PERFORMANCE ESTIMATION OF PARALLEL APPLICATION FOR SOLAR IMAGES PROCESSING

DIMITAR GARNEVSKI, PETYA PAVLOVA

Abstract: Many of the implemented parallel versions of applications are based on sequential ones. In the process of creating the parallel version or improving it, it is necessary to estimate the performance of the application using different metrics. The problem gains more weight in cases where different techniques are used to create parallelism in the studied application. In the current work will be examined a performance estimation of a created parallel version of an application for processing a series of images of the solar corona.

Key words: parallel computing, performance, image processing, algorithms, OpenCL, OpenMP, MPI

1. Introduction

The choice to a parallel processing environment is based on the capability to achieve maximum performance on given devices, as well as the applicability of a paradigm for parallel programming to the problem solved. Regardless of whether a sequential or parallel algorithm is estimated, the primary parameter used for measurement is time, as well as parameters that are associated with it. By calculating the time required to execute an algorithm or its individual parts, the appropriate conclusions about the performance of an algorithm can be made. During estimation of software and measuring their parameters, the following parameters can be calculated and evaluated:

Speedup - the times of successive and parallel execution, where the time of consecutive execution is the sum of the total calculation time required for each task, and the time in the parallel execution is a planned time frame for a limited number of processors.

$$S_p = \left(\sum_{i=1}^n T_i \right) / T_p \quad (1)$$

Efficiency - parallel program performance is a measure of CPU usage where S_p = Speedup, N_p = Number of processors

$$EFF = S_p / N_p \quad (2)$$

Overheads - measured the extra time required by the parallel program to perform the calculations

$$O_H = T_p - (T_s / N_p) \quad (3)$$

where O_H = Overheads, T_p = Parallel time, T_s = Serial time, N_p = Number of processors

In terms of distributed computing can be defined also:

- Fork time - time needed to distribute data between a number of processors
- Join Time - the time needed to collect the results of a number of processors

2. Parallel application for solar images processing

The method for motion tracking and mapping, based on solar corona images [1] proposes an ability for modeling the dynamical changes in solar prominences during its evolution. Parallel implementation of it was presented in [2].

The process of image processing is divided into the following basic steps:

- Loading images into the application
- Pre-processing of input images with filtering, which includes conversion to the internal format, clearing noise, adding artificial sunshine [6][7]
- Implemented algorithm creates a motion map by correlating two frames and generates visualization of the motion
- Write output motion map to the disk

Sample input image and generated output from the application are shown on Fig.1 and Fig.2



Fig. 1. Input sample image. Wavelength 304\AA , captured with SDO/AIA on 22.06.2014



Fig. 2. Motion map generated from two input gray-scale images of the solar corona with wavelength 304\AA .

Estimated parallel software uses following paradigms for parallel programming:

- Massively parallel processing with OpenCL [3]. Used as the primary platform for the implementation of image processing filters
- Shared memory parallel processing with OpenMP [5]. Used as an ancillary tool for partitioning code that perform the conversion and processing operations involving different data structures, as well as parts of filters that perform operations such as reduction and can therefore be relatively efficiently implemented using OpenMP and to run on a CPU.
- Distributed computing with MPI [4]. Library aims to provide communication when performing parallel processing applications on distributed systems. In current work, MPI is used in a console version of the application that can be run on a cluster. Using MPI, tasks are distributed across application instances, as well as data processing during the processing process.

2.1. Single process application

Estimated parallel software can be divided into number of common blocks shown on Fig.3. Single process parallel application uses OpenCL and OpenMP at pre-processing stage. Pre-processing stage includes sequential execution of three filter: noise normalization, Mexican hat filter and median filter. Processing stage which involves motion map generation and its conversion to result image (via HSV color space) uses OpenCL kernel.

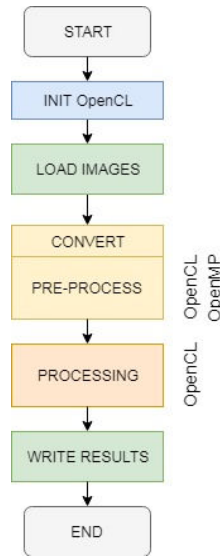


Fig. 3. Parallel computing on single node

2.2. Distributed processing

In the implemented application, MPI is used to synchronize the execution over multiple nodes and to perform data exchange between nodes. Lack of dependencies on the image pre-processing stage, each image being processed using OpenCL filters within a single node, the filters being executed on a GPGPU or CPU computing device. At the pre-processing stage, MPI's primary task is to allocate the number of tasks (inputs) for each of the nodes involved in the processing.

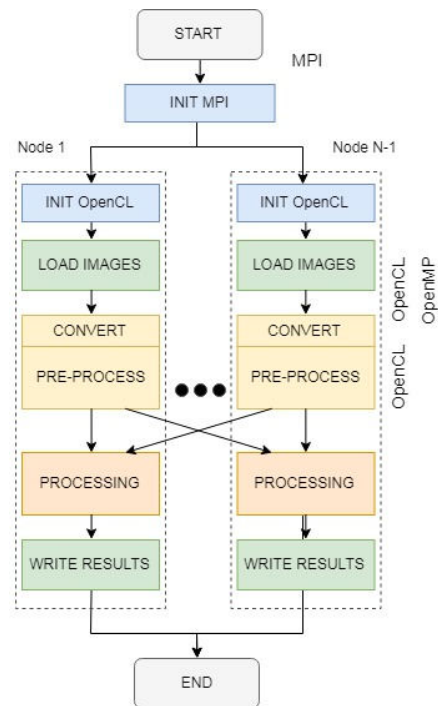


Fig. 4. Distributed computing in the estimated application

3. Test data

The application presented in the second part uses a series of input images. As input data are used images of solar disk captured using the SDO (Solar Dynamics Observatory) space probe, the AIA (Atmospheric Imaging Assembly) tool. Images were captured on 22.06.2014 with wavelength 304Å. Sample input image (top 50% part) is shown on Fig.1. Original images from SDO/AIA are with size 4096 x 4096 pixels and file format JPEG 2000. Due some limitations of serial version of the algorithm we use scaled down to 1024 x 1024 pixels images in JPEG file format. Available count of input images is 100.

4. Performance estimation

Before performance estimation we must specify our test environment in which we perform software execution and time measurements. All measurements were performed on machine with Intel Core i5 2520M CPU (with support of SSE4.2, EM64T, AVX), 8 GB of RAM, hard drive is a SSD and Linux OS. OpenCL platform is version OpenCL 1.2 with AMD provided driver (due lack of Intel drivers for older CPUs). In case of distributed computing with MPI we use local area network and second machine with identical parameters.

4.1. Single process application

Based on Fig.4 which represents processing on a multi node we can define time measuring points for the application. As the serial and parallel versions of the application uses same input and output of the framework for i/o operations (OpenCV) we skip time measuring in stages "load images" and "write results" and concentrate over pre-processing and processing stages.

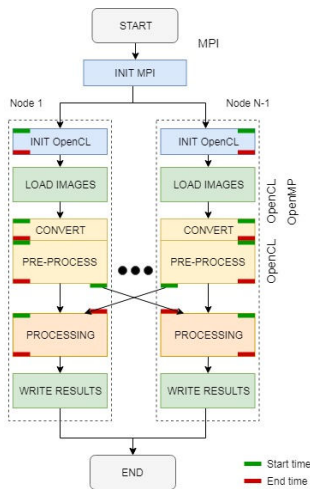


Fig. 5. Time measurement moments

Pre-processing stage on Fig.5 was divided into convert (prepares image structure) and filtering (three filters mentioned in 2.1) operations. Summarized results (10 series of input 10 images) of sequential and parallel execution times of this operations is shown in Table 1, Fig. 6 and Fig. 7.

Table 1. Execution times (seconds) of preprocessing (for single image)

Series	Serial		Parallel	
	Convert	Filtering	Convert	Filtering
1	0.235	0.560	0.008	0.049
2	0.220	0.548	0.006	0.046
3	0.220	0.548	0.008	0.047
4	0.220	0.548	0.007	0.045
5	0.220	0.559	0.006	0.045
6	0.220	0.546	0.007	0.046
7	0.219	0.546	0.007	0.045
8	0.222	0.551	0.007	0.046
9	0.228	0.554	0.007	0.046
10	0.224	0.549	0.006	0.045
AVG	0.223	0.551	0.010	0.046

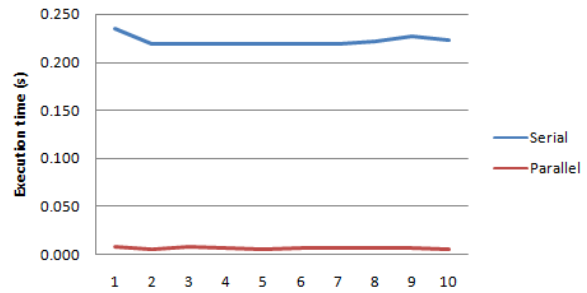


Fig. 6. Comparison of convert operation

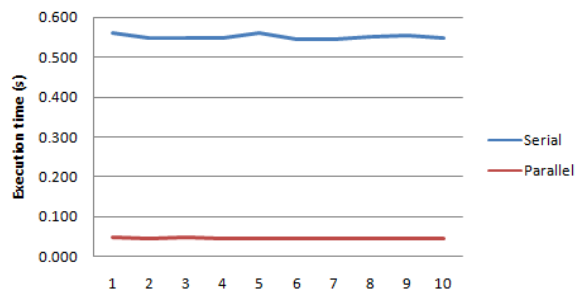


Fig. 7. Comparison of pre-processing stage

Execution times of common processing algorithm of motion map construction by correlation of two images is shown in Table 2 and Fig. 8.

Table 2. Execution times (seconds) of correlation algorithm (per iteration)

Series	Serial	Parallel	S_p
1	1.239	0.085	15
2	1.225	0.093	13
3	1.225	0.083	15
4	1.225	0.082	15
5	1.233	0.083	15
6	1.224	0.083	15
7	1.231	0.088	14
8	1.283	0.087	15
9	1.233	0.083	15
10	1.234	0.083	15
AVG	1.235	0.085	15

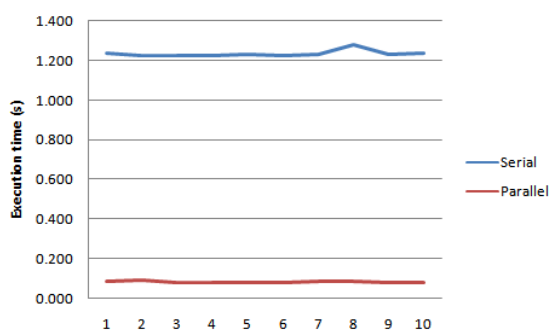


Fig. 8. Comparison of execution times of correlation algorithm

In terms of OpenCL environment we also can specify overhead, generated from initialization of OpenCL context, queues and compilation of kernels source code. Average time need for this initialization after 10 starts of the application is around 0.859 s.

4.2. Distributed processing

Version of the parallel application with MPI distribute calculations between multiple nodes in defined network. Application splits available work between nodes in sequential batches with size M or $M+1$ images (count of input images / N_p). Exchange of the data is minimized, so each node will send/receive maximum two images after the end of pre-process stage. First node only sends 1 image, and the last one receive one image. Average time need for image exchange in local area network after 10 starts of the application is around 0.02 s for each exchange and performed synchronization. Total execution time of each node include overhead from MPI initialization and communication (images exchange) and overhead from OpenCL initialization.

5. Conclusions

At the end we can make some conclusions about performance of the estimated applications. Individual stages of the parallel version on single node achieved about 15 times speedup compared to the serial version. Overhead from OpenCL initialization (~ 0.859 s) can be earned back after processing of single pair of images. Distributed version of the application preserves gained speedup due significantly small footprint of exchange operations that involves MPI. For two identical nodes in MPI network final speedup of the parallel version is 7. Depending of number of nodes and number of input pictures performance of the distributed version may vary, so the basic objective in this case is balance between size of batch on each node and number of data exchange operation between nodes.

REFERENCES

1. Petya Pavlova, Dimitar Garnevski, Kostadinka Koleva (2016). Optimization of a motion tracking and mapping method based on images of the solar corona. *Bulgarian Astronomical Journal, Volume 24*
2. Dimitar Garnevski. Implementation of software architecture for complex analysis of solar corona images, TECHSYS 2017, "Fundamental Sciences and Applications" vol.23, ISSN 1310-8271
3. Benedict Gaster, Lee Howes. Heterogeneous Computing with OpenCL
4. Marc Snir, Jack Dongarra, Janusz S. Kowalik, Steven Huss-Lederman, Steve W. Otto, David W. Walker. MPI: The Complete Reference
5. Barbara Chapman, Gabriele Jost, Ruud van der Pas. Using OpenMP. Portable Shared Memory Parallel Programming
6. Al Bovik. (2004). *Handbook of image and video processing*, Academic Press. Canada.
7. Rafael C. Gonzalez Richard E. Woods (2002) *Digital Image Processing*, 2nd Edition, Prentice Hall

Dimitar Garnevski
 Department of Electrical Engineering
 Technical University–Sofia, Branch Plovdiv
 25 Tsanko Diustabanov St.
 4000 Plovdiv
 E-mail: garnevsky_dm@abv.bg

PROGRAMMING TOOLS FOR RELIABLE USER AUTHENTICATION ACROSS MULTIPLE SERVERS

VLADIMIR DIMITROV, DENITSA TSONINA

Abstract: *This article discusses the design and the implementation of a software product that allows a user to make a single sign-on instance into a multi-server application and thus to gain access of its all remote resources that are referenced to him regardless of their location. The implementation of the authentication server is relatively simple, which implies many possibilities for future expanding of its functionality and also providing additional reliability.*

Key words: *server for authentication, OAuth, OAuth 2.0, token, REST architecture*

1. Introduction

All of the big software companies are developing their own authentication server. This is due to the fact that almost every software implies the storage of large amounts of data. Often for some reason, these data are stored on multiple servers. In this way, a higher security is achieved, a better system recovery capability for a sudden problem and modularity. However, the user does not understand that the information that he uses is in different locations, as his access to the data, wherever they may be, takes place instantly without transition from one system to another. On the other hand, the use of a separate authentication server in modular systems helps for minimizing the cyberattacks because the server is as secure as possible and it does not have to be applied the same security rules to the other servers, because the information on them cannot be identified for a particular user. The authentication servers are also applicable when it comes to functionality of the type SSO (Single Sign-On). In this case, the signing with username and password in one system allows the automatic signing of a particular user into a number of predefined other systems [1]. Such functionality is extremely popular on social networks and all add-on applications and games that expand them as well as most major email servers. The ability to use the same sign-up in multiple different systems saves not only the constant need for inputting of the username and password, but also the need for creating a new account in every application.

As part of a multifunctional system the purpose of the authentication server is to systemize

all user data into one single location and to separate this data from the other data in the system while ensuring the security of all data. However, the users of the system are still continuing to have access to all data that relates to them and it is considered to be accessed by them. It is assumed that on the servers containing user's information there is also information that they do not have access to – for example – when a client of a bank accesses his account for Internet banking, he only has access to his own accounts, not to those of the other bank customers. At the same time, users enter their username and password only once. This raises the question of how to find out which client can access certain information without knowing his username and password when part of that information is kept on another server. It is therefore necessary to invent a way in which the users to authenticate themselves at the other servers without entering their passwords again in order to provide a higher security for their own data.

The proposed solution is based on Spring OAuth2 authentication and consists of three main parts:

1. Authentication server.
2. Resource server – 1.
3. Resource server – 2.

The user communicates with the three servers via REST queries. The authentication server and the resource servers do not communicate directly with each other.

The system is implemented with the tools offered by the programming language Java. The user login the system with an email and a password that he has specified. There cannot be two users

with the same password. In case of a wrong password or a non-existing user in the database, the server returns an error message and the user can correct their data, then to try login again to the system. Before a system login for the user to be possible, he first must be registered to it as a user. This happens via a dedicated button “register”. The required fields upon a new user registration are three – First Name, Last Name and Email Address. Each of the three fields is validated as such as possible – the names cannot contain special characters and numbers, the email field cannot contain any special characters other than dot, dash and underscore. The capital letters for the email are also forbidden, because in the context of email addresses they are meaningless. The email field must necessarily follow the form xxxx@xx.xx.

After filling in all the required fields the user presses a button and this operation sends an email with an activation code to the email address used for the registration. Then a user activation page is loaded in the UI (User Interface). The time for receiving the activation email depends on the workload of the SMTP server, but generally takes no more than a few minutes. The application also offers restoring of a forgotten password. Because the passwords are stored in hashed format in the database, administrators do not have an access to their real values and cannot be given to third parties or to their owners in case of a forgotten password. For this reason, in case of a forgotten password to the user is sent an email from where he can change his password with a new one. The email is sent to the email address with which the registration was made. Password cannot match the username and must necessarily contain at least one lower case letter, an upper case letter, a special character, and a number. The minimum password length is 8 characters. In case of a mismatch, the application alerts the user for an error and refuses to send a request to the server. Both on the login screen and the account activation screen, the fields in which the user enters his password are encrypted and the password text itself is not displayed anywhere in the UI. Once the user has entered a matching passwords and the correct activation code, his account is activated. Prior to this activation, it is impossible for an inactive user to enter the system in any way, either through the UI or through the REST API.

Figure 1 shows how an authentication server works. Generally speaking, the authentication server works the following way: The user wants to access information that is in the main application. To do so, however, he must verify his identity. Because the user has already entered his password and username when he signed in the application, they are kept in the authentication

server. In order not to re-enter again this credentials when accessing the other servers in the system, the authentication server entrusts to him the so-called token. At the same time, the server with the information which the user is trying to access is set to maintain and require token authentication. By obtaining a token from the user, the server checks its validity and whether it is issued by its authentication server. This happens through a password that is kept on both servers and with which the token is signed. If this password does not match, the user access to the information he is trying to access is denied. If the password matches and the token is valid, the user gets access to the information. Each token has a pre-defined term of validity. Once it expires, it cannot be used anymore [2].

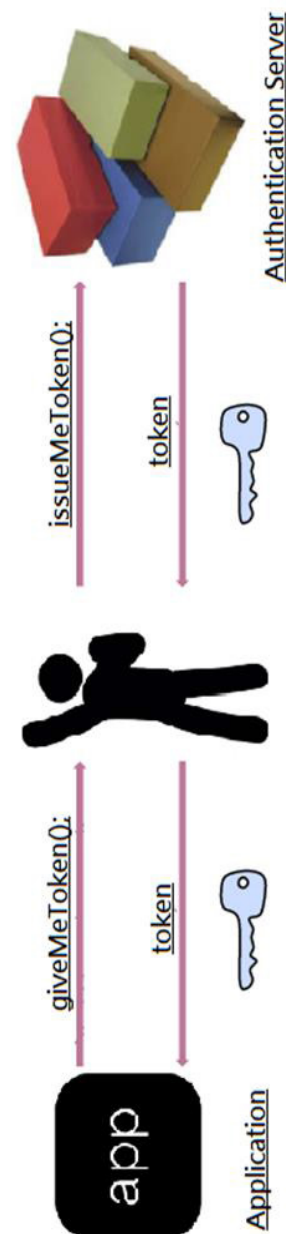


Fig. 1. Function of an authentication server

Figure 2 shows the application structure that uses an authentication server.

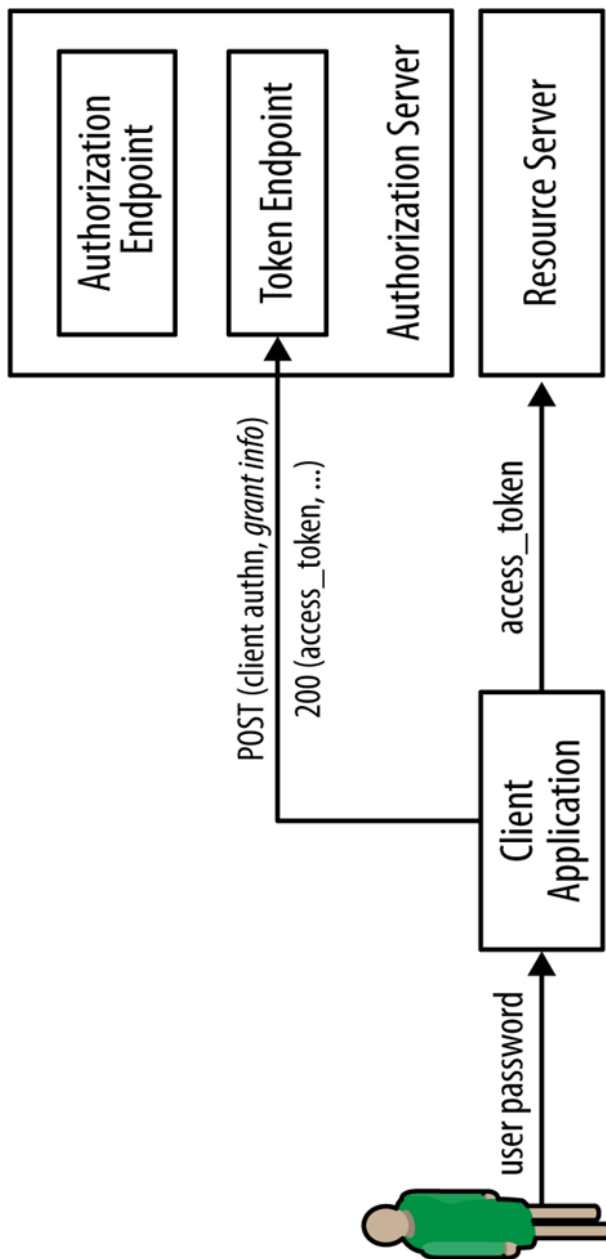


Fig. 2. An application architecture using an authentication server

Several programming solutions [3] [4] [5], which provide extremely reliable authentication across multiple servers, have been addressed in the design and implementation of the application. The provided solutions are accompanied by numerous additional functionalities, directly and indirectly related to the user authentication. It is important to note that the cost of these solutions is not negligible, and often companies prefer to implement their own authentication server, which, although limited, is much cheaper solution.

2. Application architecture

The application consists of 3 parts – an authentication server, and 2 servers that stores user's information. In order to improve the structured presentation of the information, the application imitates a system for tax declaration of property. One of the resource servers provides access to the user's real estate and the other to the vehicle property.

In the section for the real estate property there are several fields: "estate counter", "estate type", "address" and "tax assessment". The address of the property is also used as its identifier.

Similar to the real estate the section for the vehicle property also has a field for the number of properties to be declared – "vehicle counter". The other fields in this section are "model" and "power". They are also essentials because the car tax is calculated according to their power. The model field servers more as a vehicle's identifier for the user himself.

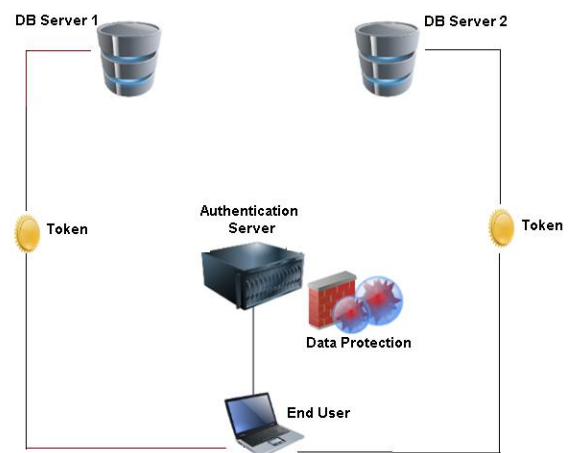


Fig. 3. Application architecture

The purpose of the authentication server is to provide unified access to the three servers without providing the user's password to the other two servers. By this way, the user data that may be subject to a hacking attacks is kept on a single location that can be further assured and secured. On the other hand, the information on the other two resource servers cannot be identifiable with an actual user, because the only user information that is kept apart from the data about its estate or vehicle property, is a randomly generated identification number.

The access to the authentication server is gain through a username and password, and to the resource servers - via a token-encrypted string of characters signed with a password that assures that the user is really the one to whom he is claiming to be.

2.1. Use-case UML diagram

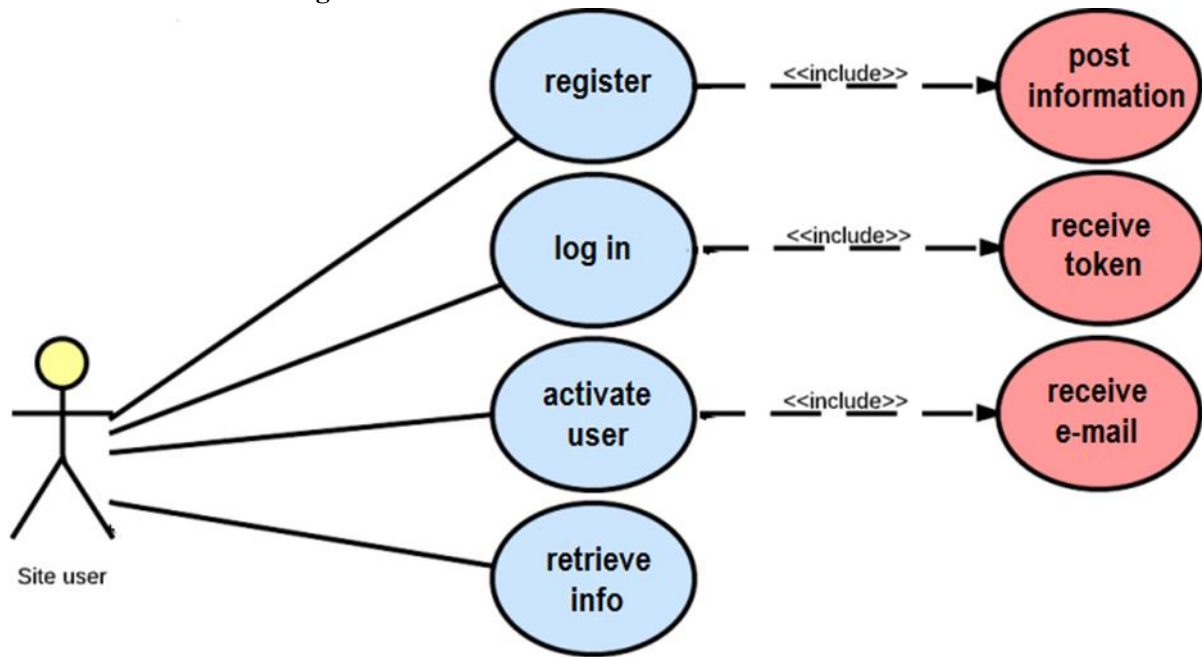


Fig. 4. Use-case diagram

2.2. Deployment diagram

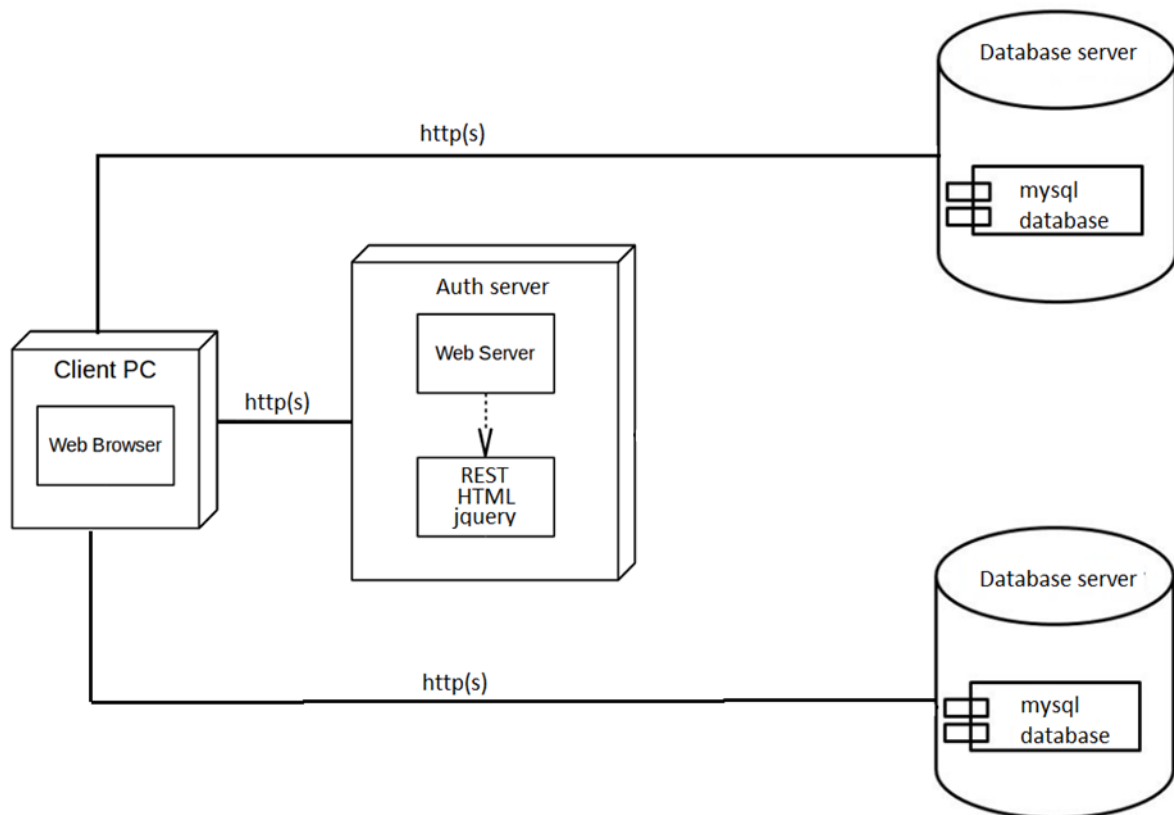


Fig. 5. Deployment diagram

2.3. Interaction diagram

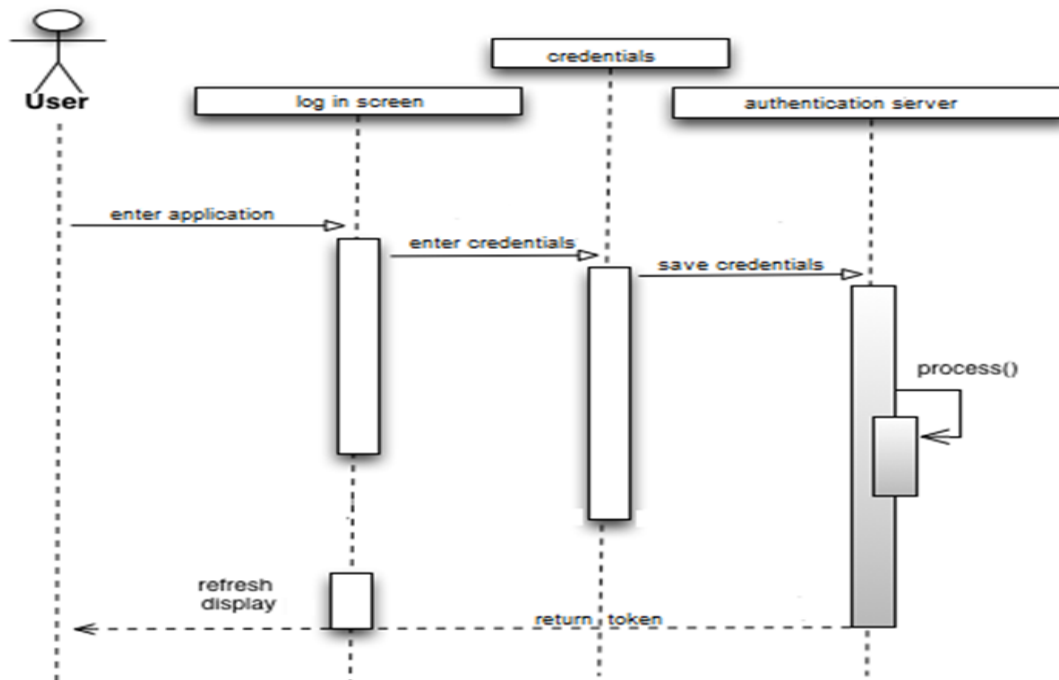


Fig. 6. Interaction diagram

3. Authentication server implementation

The authentication server is based on Spring OAuth 2.0 framework and it follows the REST architecture. After running, the server remains in ON state and expects REST queries for executing.

In order to create a new user in the database the required fields are "first_name", "last_name" and "email". The identity of the registered users is provided by the email column, i.e. there cannot be two or more entries with the same e-mail address. For each user is automatically generated UUID (Universally Unique Identifier), which is also unique and it help for providing the connection between the authentication server and the resource servers. The passwords are stored in hash format using SHA-256. The password is known only by the user who has create it and is not visible in the database or from any other part of the application. For added security, a randomly generated string called salt is concatenated before encrypting the password. For every password is randomly generated a new salt. The hashing allows for later authentication without storing the password as a clear text. Since the salt strings should not be remembered by the humans, they do not burden the users and at the same time they provide an additional protection for their personal data. Because the salt strings are different, they mostly protect frequently used passwords or the users who use the same password in multiple systems. Despite

the identity of the password, their hash appears to be completely different thanks to the different salt.

Upon registration, the user receives an email with an activation code (invitation token). In case that he has activated his account with this code the field in the database for this user is marked as "active" or "1". Otherwise, the user is inactive and the field remains "0". The inactive user cannot sign in because the email that they have provided is believed to be wrong or somebody else's.

The activation request sends as a payload the password and the confirmed password after checking that their values match completely. The pre-received by email activation code is sent as part of the Request URL in the main part of the request. In case of password mismatch, the request to server is not sent at all, since the page validation takes care of data consistency. The send activation code in the request is compared to the one in the database and if no inactive user with such code is present, the server returns an error message: "Request processing failed; nested exception is java.lang.IllegalArgumentException: Activation token doesn't exist".

The activation code to be sent by email is stored in the column "invitation token". It is a standard randomly generated string and is not further secured by encryption algorithms, because it does not contain any important information. The column "invitation token expiry date" marks the

exact date and time that the code expires. By default, the validity is set to 2 days from the date of generating of the code but it can be changed. After activating a user, the field “invitation_token” automatically becomes null.

An integration with a Gmail account is implemented that is also the sender of the emails on behalf of the application.

After the activation of the user, he can now access the three servers. By logging in to the authentication server, it automatically gains access to the other two resources servers via a token, which is also called JWT. The token is a long string that at first seems randomly generated. The tokens are signed by the server’s key so that the server is able to verify that the signature is legitimate. The tokens are designed to be compact, URL-safe and usable especially in a web browser for SSO. The tokens can be verified and encrypted. They are base-64 encoded and usually consist of three parts: header, payload and signature.

The O.Auth.SECRET_KEY constant is the binding link between the authentication server and the resource servers. If the token is signed with a code other than the one on the resource server, the user's access will be denied. The same would happen if the token has expired. The identity of the user is authenticated by his personal identification number (UUID). This ensures that each user has access only to their own information, but not to the others user’s information.

The two resource servers on which the information is stored are relatively simple by design. Their only purpose is to store and provide data after appropriate authentication. They do not store passwords in their database, the only identifier they kept is the UUID of each user, which is used for binding the link between the user and the data that the server can provide to him.

A demonstration of the system when a new user registration is made can be viewed via this link:

<https://my.pcloud.com/publink/show?code=XZMQAm7ZWJey85DuM7J4J7k9ADEsI8UBowak>

4. Conclusion

The result of this work is a solution to a common problem. Internet security and convenience are essential in the age of the IT – a time when more and more of our personal data are exposed to hacker attacks and a number of other Internet threats. Separating the user data from the resources that users access is critical to protecting passwords and other sensitive information [6]. For more convenient for the user is to enter only once his credentials, instead of several times, for a system that uses several remote storage servers. Of

course, the server can also be used for transition between two or more system and SSO-type functionality. At the same time, the separate authentication server greatly improves the structure of the programming code, because the functionality associated with the user data is stored on a single server and it is not implemented to each of the resource servers.

The realization of the authentication server is relatively simple, which implies many possibilities for extending its functionality and providing addition reliability and convenience. The simplified realization does not result in a loss of performance because the key algorithms for signing and authenticating the origin of the data are implemented.

The program implementation of the authentication server gives unlimited possibilities for future development of the idea. It is possible to use a more reliable algorithm for signing the tokens, as well as encrypting of the entire token before sending it and decrypting it after it is received. However, the added security of this type would increase the time for issuing a token, which requires additional work to optimize the process in conditions of a slow keys and one more step with encryption and decryption. Such optimization could be accomplished by splitting the application of a multiple threads. However, they must be perfectly synchronized, in order to avoid deadlocks that would lead to the opposite of the desired effect.

REFERENCES

1. Wang, G., Yu, J. and Xie, Q. (2013). Security Analysis of a Single Sign-On Mechanism for Distributed Computer Networks. *IEEE Transactions on Industrial Informatics*, Volume 9, Issue 1, 294-302.
2. <https://www.rfc-editor.org/pdf/rfc7523.txt.pdf>
3. <https://www.vasco.com/products/management-platforms/identikey-authentication-server.html>
4. <https://www.rcdevs.com/products/openotp/>
5. <https://www.hidglobal.com/products/software/activid/activid-authentication-server>
6. Beltran, V. (2016). Characterization of web single sign-in protocols. *IEEE Communications Magazine*, Volume 54, Issue 7, 24-30. Department of Computer Systems

Technical University of Sofia
8 Kliment Ohridski Blvd.
1000 Sofia
BULGARIA
E-mail: vldimitrov@tu-sofia.bg
E-mail: denitsa.yordanova@abv.bg

TESTING AND DIAGNOSTICS OF COMPUTER SYSTEMS

ATANAS KOSTADINOV

Abstract: *This paper presents some important aspects of the educational process in testing and diagnostics of computer systems subject in Technical University – Sofia, Plovdiv branch. Many years we teach to the students how to deal with computer systems problems and especially with personal computers ones. We hope that received, refreshed and obtained knowledge in above-mentioned subject could help the computer specialists to diagnose and eventually to solve some simple problems which could be appeared in the computer systems.*

Key words: *testing, diagnostics, computer systems, personal computers, syllabus, software and hardware tools*

1. Introduction

Testing and diagnostics of computer systems (TDCS) is a compulsory elective subject in Technical University – Sofia, Plovdiv branch. Similar courses could be found in different universities, as well as professional education institutions in abroad [1, 2, 3, 4, 5, 6, 7, 8, 9]. This specific knowledge could be used in order to be performed diagnostics and eventually to be solved some simple problems which appeared during computer systems work [10, 11]. Small number of students can continue increasing your knowledge and experience in case of working as personal computer (PC) support engineers.

Some years ago, a laboratory manual was published [12]. Because the content of this publication is a quite large, it could be used in lectures, too. Our experience shows that this subject is interesting for the students. Most of them pass the written exam in the form of test from the first time. Small group of students must appear for the second and third time in order to pass the exam.

2. Diagnostics of computer systems

In this section of the paper a short information about the lectures and laboratory exercises topics of TDCS will be given. In the lectures as well in the laboratory exercises the testing and diagnostics of computer systems main components are presented. In every computer system they are the processor (Central Processing Unit – CPU), the memories, the video subsystem, etc. In the next lines will be given information about these components, their diagnostics and used test

programs. All main PC components are tested during the POST (Power-on self-test). Unfortunately, this initial test does not guarantee that if the component passes the test everything is OK. There are additional comprehensive tests which could find more problems than the above-mentioned one.

2.1. Processor diagnostics and test programs

As processor tests are used different tasks and corresponding programs. They could be as example calculating complex matrix, calculating Pi, fast Fourier transformation, different sorting algorithms, etc [13]. In order to be performed above-mentioned tests, we should have at least a partially working PC.

There are additional CPU tests as burn-in and stability ones. They very heavily load the processor and due to that the more power and heat are generated. During work of these test programs, it should be observed the temperature of the CPU. The results of testing the microprocessor are presented in Fig. 1 [13].



Fig. 1. CPU testing

2.2. RAM (Random Access Memory) diagnostics and test programs

RAM is another very important PC component. RAM test includes write, read and compare operations [12, 14] realized with help of previously checked CPU. In every of the memory cell is written a digital value called sample. Then is performed memory read operation. These two numbers (written and read ones) are compared each other. If there is no difference between them, we could expect that the checked memory cell is working properly otherwise it will be a problem with this cell.

In this way, it have to be checked every of the RAM cell. It is possible after finishing the test to be started new one again (next pass) or the tests to be started one by one for a given period of time. A special case is using the DOS (Disk Operation System). This operating system is smaller than others used in the PCs today. The test program can check a part of RAM, and then DOS could be loaded in the verified cells. Using this approach, it is able to be tested all available memory cells.

A part of RAM testing using Windows operating system is shown in Fig. 2 [14].

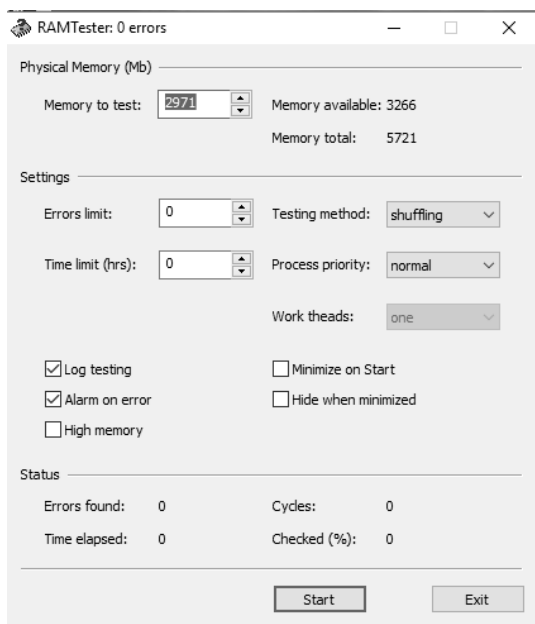


Fig. 2. RAM testing

2.3. HDD (Hard Disk Drive) diagnostics and test programs

Hard disk drive is another type of nonvolatile memory. The very large capacity of contemporary HDDs has able to store an extremely big amount of digital information (terabytes). Hard disk drive contains the operating system as

well different types of software as example software drivers, application programs, files with the results of application programs, Internet files, etc.

There are two major producers of HDDs – the American firms Seagate and Western Digital. Based on this, the different test programs are able to be used in PCs. In Fig. 3, the results from Western Digital Data LifeGuard Diagnostic [15] test program for Windows operating system are presented.

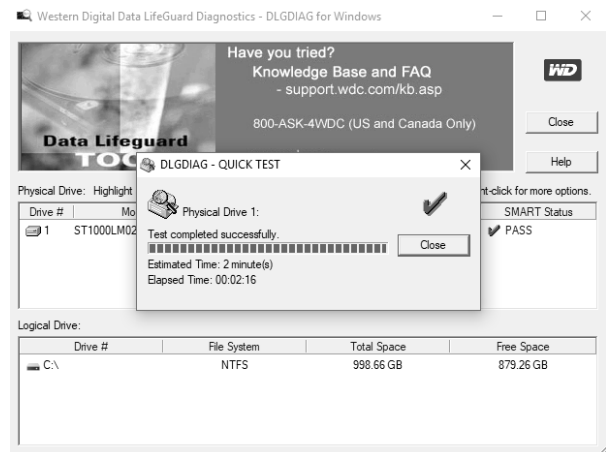


Fig. 3. HDD testing

2.4. Video subsystem diagnostics and test programs

Main components of the PC video subsystem are the video card and the monitor. The video processor located on the video card could be tested using different programs. Video RAM memory which is another major component of the video card is able to be tested as it was explained in previous paragraph especially using write, read and compare operations. This type of testing is presented in Fig. 4 [16].

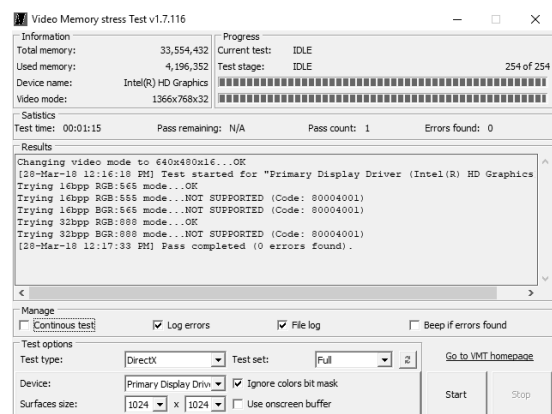


Fig. 4. Video RAM testing

The PC's video monitor can be tested using specialized and vendor-specific software. These utilities are able to check the colors, geometry of the screen, convergence of the pixels, color gradients, color spectrum, etc. In Fig. 5 [17] are shown the test results of the concrete video monitor.

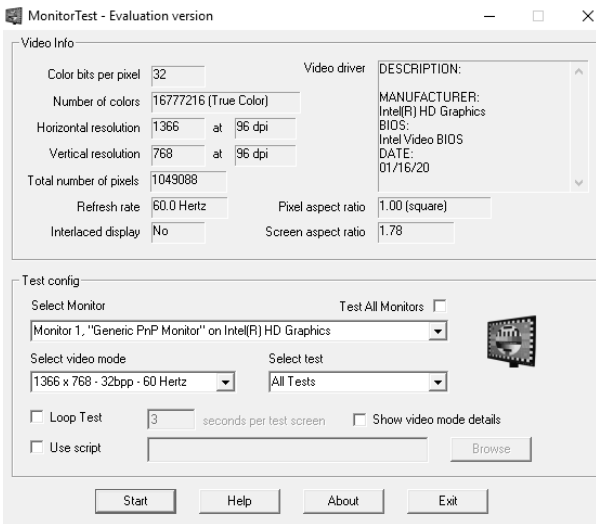


Fig. 5. PC's video monitor test results

3. Benchmarking computer systems

The benchmark tests are very often used in order to be measured the performance of the whole computer system or the performance of the PC building components [18]. The benchmarking tools are consisted usually of one program or a set of programs used in the process of performance evaluation of computer systems [19]. The different metrics could be used [20] as final results. They are able to be MIPS (Million instructions per second), FLOPS (Floating-point operations per second), MB/s and so on. The results of the benchmarking can be used in comparison between different computer systems and also between different PC components. In order this evaluation to be correct should be applied the same test programs and the same version of the benchmarking software.

In Fig. 6, 7, 8 and 9 are presented the test results received from benchmark tool [21] which evaluate the performance of the whole computer system, disk and video subsystems.



Fig. 6. Whole computer system benchmark results



Fig. 7. Disk subsystem benchmark results



Fig. 8. 2D Video subsystem benchmark results



Fig. 9. 3D Video subsystem benchmark results

4. Conclusions

In this paper, some important aspects of the educational process in testing and diagnostics subject are presented.

This bachelor's compulsory elective course is very suitable for our students, because they have to learn or have to refresh their knowledge in above-mentioned topic.

The computer systems are very complex structure consisting of many different hardware and software components.

The obtained and refreshed skills will help the students to diagnose and eventually to solve some problems which could appear during computer systems work.

5. Acknowledgments

I would like to thanks to my colleague Assistant Professor PhD Mollov for his cooperation during presenting this subject to the bachelor's students and for his help in publishing a manual of the laboratory exercises.

REFERENCES

1. <http://www.thestug.org/tutorials/Basic%20PC%20Maintenance%20and%20Backup.pdf>
2. http://www.navttc.org/downloads/curricula/VOC/Computer_Hardware_Technician.pdf
3. http://www.lssc.edu/faculty/betti_mcturk/.../Syllabus/HARDWARE%20SYLLABUS.pdf
4. <https://www.lewisu.edu/academics/comsci/pdf/Lewis-AAS-CIT.pdf>
5. http://math.mercyhurst.edu/~platte/syllabi/com_p_oper_term2_11-12.php
6. <https://www.scitraining.ca/pc-repair>
7. <http://mrkalifm.weebly.com/computer-diagnostic--repair.html>
8. <https://www.waynesville.k12.mo.us/Page/4144>
9. <http://dec.nure.ua/en/en-info-discipline-ceme/>
10. Schneider, W. (1981). Basic computer troubleshooting and preventive computer maintenance operation. *Behavior research methods & Instrumentation*, volume (13), pp. 153-162.
11. Lemeš, S. (2014). Održavanje racunarskih sistema (Maintenance of computer systems). Brdarević, S. and Jašarević S. (ed.), *Proceedings of the 3rd conference Održavanje (Maintenance) 2014*, pp. 53-60. University of Zenica, Zenica, Bosnia and Herzegovina.
12. Mollov, V. and Kostadinov, A. (2014). *Testing and diagnostics of computer systems (in Bulgarian)*. Technical University - Sofia, Sofia, Bulgaria.
13. <http://www.7byte.com/index.php?page=hotcpu>
14. <http://cpu.rightmark.org/>
15. <https://support.wdc.com/downloads.aspx?p=3>
16. https://mikelab.kiev.ua/index_en.php?page=PROGRAMS/vmt_en
17. <https://www.passmark.com/products/monitortest.htm>
18. Menascé, D., and Almeida, V. (2001), *Capacity planning for Web services: metrics, models, and methods*, Prentice Hall, Upper Saddle River, USA.
19. Bouckaert, S., Gerwen, J., Moerman, I., Phillips, S. and Wilander, J. (2010), Benchmarking computers and computer networks, *EU FIRE*, White Paper.
20. Eeckhout, L. (2010), *Computer Architecture Performance Evaluation Methods*, Morgan & Claypool Publishers, San Rafael, USA.
21. <https://www.passmark.com/products/pt.htm>

Authors' contacts

Organization: Technical University – Sofia,
Plovdiv branch

Address: 25 Tsanko Diustabanov Str.

Phone (optional): +359 32 659 726

E-mail: kostadat@tu-plovdiv.bg

USING GRAPHIC PROCESSING UNITS FOR IMPROVING PERFORMANCE OF DEEP LEARNING PROCESSES

TEODORA HRISTEVA, SPAS TOMOV, MARIA MARINOVA

Abstract: GPUs have proven to be significantly important for deep learning because they can process lots of little bits of data in parallel and deep learning networks are designed to analyze massive amounts of data at speed. GPU-accelerated computing has now grown into a mainstream movement supported by the latest operating systems. We made an experiment running a neural network on CPU and GPU alternatively to prove the performance advantages of using GPUs for deep learning tasks.

Key words: deep learning, CUDA, GPU, TensorFlow

1. INTRODUCTION

The GPU's advanced capabilities were originally used primarily for 3D game rendering. But now those capabilities are being harnessed more broadly to accelerate computational workloads in areas such as deep learning, financial modeling, cutting-edge scientific research and much more. GPUs are optimized for taking huge batches of data and performing the same operation over and over very quickly, unlike PC microprocessors, which tend to skip all over the place.

Architecturally, the CPU is composed of just few cores with lots of cache memory that can handle a few software threads at a time. In contrast, a GPU is composed of hundreds of cores that can handle thousands of threads simultaneously. The ability of a GPU with 100+ cores to process thousands of threads can accelerate some software by 100x over a CPU alone. What is more, the GPU achieves this acceleration while being more power- and cost-efficient than a CPU.

In the PC of today the GPU can now take on many multimedia tasks, such as accelerating Adobe Flash video, transcoding (translating) video between different formats, image recognition, virus pattern matching and others.

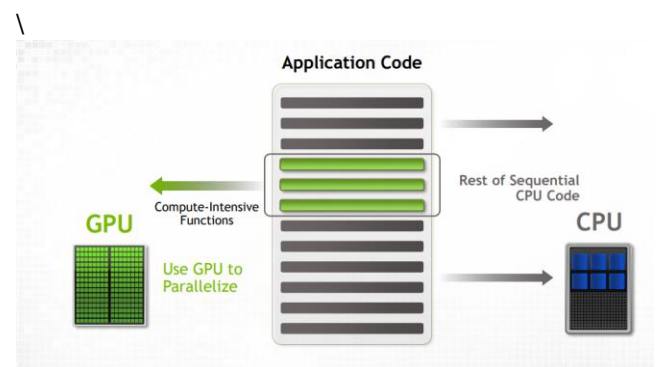


Fig. 1. GPU computations

2. PROCESS DESCRIPTION

2.1 CUDA architecture

CUDA is a parallel computing platform and application programming interface (API) model created by Nvidia. It allows software developers and software engineers to use a CUDA-enabled graphics processing unit (GPU) for general purpose. The CUDA platform is a software layer that gives direct access to the GPU's virtual instruction set and parallel computational elements, for the execution of compute kernels.

The CUDA platform is designed to work with programming languages such as C, C++, and Fortran. This accessibility makes it easier for specialists in

parallel programming to use GPU resources, in contrast to prior APIs like Direct3D and OpenGL, which required advanced skills in graphics programming.

The graphics processing unit (GPU), as a specialized computer processor, addresses the demands of real-time high-resolution 3D graphics compute-intensive tasks. Recently GPUs had evolved into highly parallel multi-core systems allowing very efficient manipulation of large blocks of data. This design is more effective than general-purpose central processing unit (CPUs) for algorithms in situations where processing large blocks of data is done in parallel.

Here is an example of CUDA processing flow:

- Copy data from main memory to GPU memory
- CPU initiates the GPU compute kernel
- GPU's CUDA cores execute the kernel in parallel
- Copy the resulting data from GPU memory to main memory

The CUDA platform is accessible to software developers through CUDA-accelerated libraries, compiler directives such as OpenACC, and extensions to industry-standard programming languages including C, C++ and Fortran.

In addition to libraries, compiler directives, CUDA C/C++ and CUDA Fortran, the CUDA platform supports other computational interfaces, including the Khronos Group's OpenCL, Microsoft's DirectCompute, OpenGL Compute Shaders and C++ AMP. Third party wrappers are also available for Python, Perl, Fortran, Java, Ruby, Lua, Common Lisp, Haskell, R, MATLAB, IDL, and native support in Mathematica.

In the computer game industry, GPUs are used for graphics rendering, and for game physics calculations (physical effects such as debris, smoke, fire, fluids). CUDA has also been used to accelerate non-graphical applications in computational biology, cryptography, machine learning and other fields.

CUDA provides both a low level API and a higher level API. The initial CUDA SDK was made public for Microsoft Windows and Linux. Mac OS X support was later added in version 2.0.[13] CUDA works with all Nvidia GPUs from the G8x series onwards, including GeForce, Quadro and the Tesla line. CUDA is compatible with most standard

operating systems. Nvidia states that programs developed for the G8x series will also work without modification on all future Nvidia video cards, due to binary compatibility.

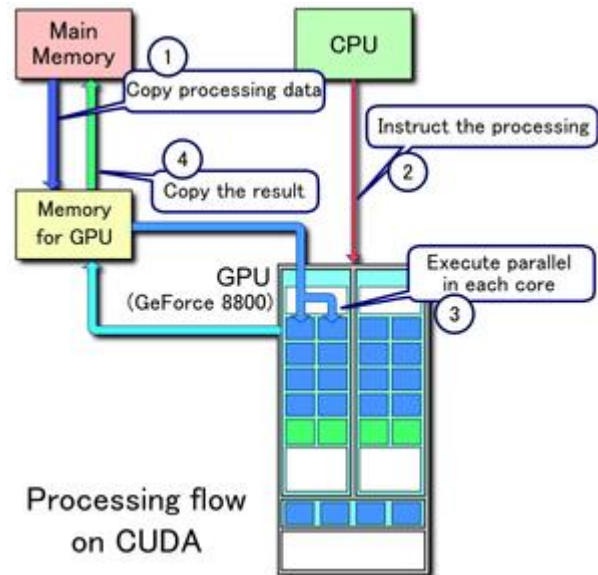


Fig. 2. Processing flow on CUDA

2.2 Pascal architecture

Pascal is the codename for a GPU microarchitecture developed by Nvidia, as the successor to the Maxwell architecture.

Architectural improvements of the GP100 architecture include the following:

- In Pascal, an SM (streaming multiprocessor) consists of 64 CUDA cores. Maxwell packed 128, Kepler 192, Fermi 32 and Tesla only 8 CUDA cores into an SM; the GP100 SM is partitioned into two processing blocks, each having 32 single-precision CUDA Cores, an instruction buffer, a warp scheduler, 2 texture mapping units and 2 dispatch units.

- CUDA Compute Capability 6.0.
- High Bandwidth Memory 2 — some cards feature 16 GiB HBM2 in four stacks with a total of 4096-bit bus with a memory bandwidth of 720 GB/s.

- Unified memory — a memory architecture, where the CPU and GPU can access both main system memory and memory on the graphics card with the help of a technology called "Page Migration Engine".

- NVLink — a high-bandwidth bus between the CPU and GPU, and between multiple

GPUs. Allows much higher transfer speeds than those achievable by using PCI Express; estimated to provide between 80 and 200 GB/s.

- 16-bit (FP16) floating-point operations (colloquially "half precision") can be executed at twice the rate of 32-bit floating-point operations ("single precision") and 64-bit floating-point operations (colloquially "double precision") executed at half the rate of 32-bit floating point operations.

- More registers — twice the amount of registers per CUDA core compared to Maxwell.

- More shared memory.

- Dynamic load balancing scheduling system. This allows the scheduler to dynamically adjust the amount of the GPU assigned to multiple tasks, ensuring that the GPU remains saturated with work except when there is no more work that can safely be distributed to distribute. Nvidia therefore has safely enabled asynchronous compute in Pascal's driver.

- Instruction-level and thread-level preemption..

2.3 TensorFlow

TensorFlow is an open-source software library for dataflow programming across a range of tasks. It is a symbolic math library, and is also used for machine learning applications such as neural networks.

TensorFlow was developed by the Google Brain team for internal Google use. It was released under the Apache 2.0 open source license on November 9, 2015.

TensorFlow is Google Brain's second generation system. While the reference implementation runs on single devices, TensorFlow can run on multiple CPUs and GPUs (with optional CUDA and SYCL extensions for general-purpose computing on graphics processing units). TensorFlow is available on 64-bit Linux, macOS, Windows, and mobile computing platforms including Android and iOS.

TensorFlow computations are expressed as stateful dataflow graphs. The name TensorFlow derives from the operations that such neural networks perform on multidimensional data arrays. These arrays are referred to as "tensors".

EXPERIMENTAL TEST

For our experiment we used:

- 64-bit desktop
- Windows 7
- GPU with Pascal architecture
- CUDA® Toolkit 9.0.
- Python 3
- TensorFlow
- Iris flower data set – a popular data set, introduced by Ronald Fisher

We created a program that makes the following steps:

- Import and parse the data sets.
- Create feature columns to describe the data.
- Select the type of model
- Train the model.
- Evaluate the model's effectiveness.
- Let the trained model make predictions.

We executed the program using the computer's CPU and then run it again, this time using the GPU and measured the time of execution in both cases.

RESULTS

The result are as follows:

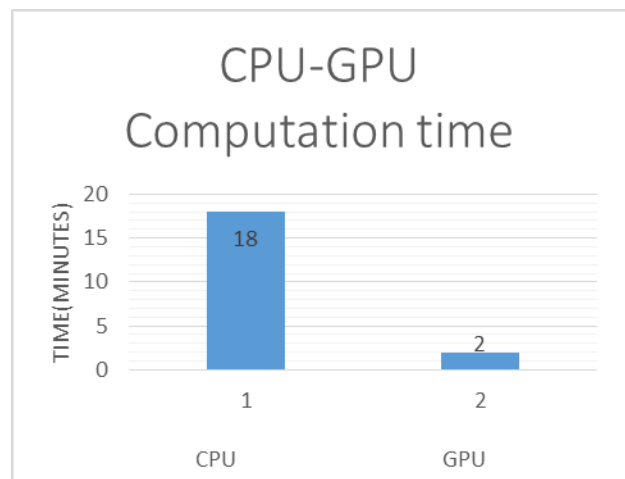


Fig. 3. Results from the experiment

We received 18 minutes when running on CPU and only 2 minutes on GPU.

3. CONCLUSION

Using GPUs is for deep learning computational task is a good way to gain speed and greater performance. GPUs have proven to be significantly important for deep learning because they can process a lot of data in parallel and neural networks are designed to analyze massive amounts of data at speed. GPU-accelerated computing has now grown into a mainstream movement supported by the latest operating systems. The reason for the wide and mainstream acceptance is that the GPU is a computational powerhouse, and its capabilities are growing faster than those of the x86 CPU. Our experiment shows that GPUs can be significantly helpful in the complex computations in the neural networks.

BIBLIOGRAPHY

1. <https://developer.nvidia.com>
2. <https://www.notebookcheck.net/Nvidia-Pascal-Architecture-Overview>
3. Nick McClure, "TensorFlow Machine Learning Cookbook", Packt Publishing, 2017
4. A. Géron, Hands-On Machine Learning with Scikit-Learn and TensorFlow, O'Reilly,
5. Ian Goodfellow, Yoshua Bengio, Aaron Courville-Deep Learning, MIT Press
6. Li Deng, Dong Yu Deep Learning: Methods and Applications, Now, 2014.
7. Dean, Jeff; Monga, Rajat; et al., "TensorFlow: Large-scale machine learning on heterogeneous systems" Google Research, November 10, 2015

Teodora Hristeva
 Technical University - Sofia, Branch Plovdiv
 Telephone: +359 899 862 656
 Email: thisteva@gmail.com

MOBILE CROWDSOURCING FOR VIDEO STREAMING

GEORGI ILIEV

Abstract: *This paper focuses on the concept of collective internet activities, known as the crowdsourcing paradigm. A general architecture of a mobile crowdsourcing system for video data streaming is proposed. The implementation of the idea in Footlikers platform is presented with extended capabilities for real-time broadcasting and receiving amateur football match video. A key software part of the system is a mobile application with iOS and Android support, which permits the users of the Footlikers platform to live stream picture and sound from football games to other users using the application. The proposed crowdsourcing architecture is realized using main client-server (Footlikers.com) and WOWZA video streaming engine, deployed on Amazon EC2 machine. Results of the program realization of the developed system prototype are presented, intended for amateur football competitions based and organized in France, Belgium and Luxembourg.*

Key Words: *mobile application, crowdsourcing architecture, Amazon Compute Cloud (EC2), Representational State Transfer (REST)*

1. Introduction

Modern ICT and software applications provide great opportunities for communication and information sharing. When such activities are performed by organized groups of users to achieve a common goal (with or without payment), they refer to the crowdsourcing category. Crowdsourcing finds application in a variety of areas of practice, from the use of collective work of the crowd through social network platforms such as Wikipedia (from 2001), Freelancer (from 2004), Facebook and Twitter to areas like urban surveillance [1], environmental sensing [2], disaster management [3], crowd intelligence [4], and more. The crowdsourcing paradigm finds growing interest and development in the scientific literature. New concepts, principles and classifications are being developed to build and operate crowdsourcing platforms, systems, tools and services [5, 6].

In the field of mobile services and crowdsourcing systems there is a huge amount of publications (see [7, 8, 6] and the literature cited there). Specifically, they concern the development of standards, concepts and tools for creating and sharing movies and videos in the context of mobile crowdsourcing and the Big Data paradigm, in which this type of information is included. Overviews of the current state of development and future challenges associated with the processes of video streaming, main architectures and services adopted over the

years for streaming live and pre-recorded videos, including the internet of things paradigm are presented in [8]. More aspects, including necessary RESTful web-services and cloud computing could be found in [10, 11].

The aim of this study is to propose a concept and architecture for the creation of a mobile crowdsourcing system which provides the extended features for live broadcasting and receiving amateur video from the football matches. An implementation is developed as a mobile crowdsourcing application in the frame of the social network platform footlikers.com [12]. An application in the same field but without crowdsourcing is described in [13].

2. General architecture for mobile video crowdsourcing

Designing a crowdsourcing system for video streaming from mobile devices to other mobile devices involves the development of a specific architecture and a mobile application for quality transmission of information. To achieve greater flexibility and parallel asynchronous video processing, the system is expected to work in cloud environments.

We propose the concept of general crowdsourcing architecture in the case of mobile video streaming, which is presented schematically in Fig. 1. In a simplified case this comprises three main components: active crowd of broadcasters, i.e., users

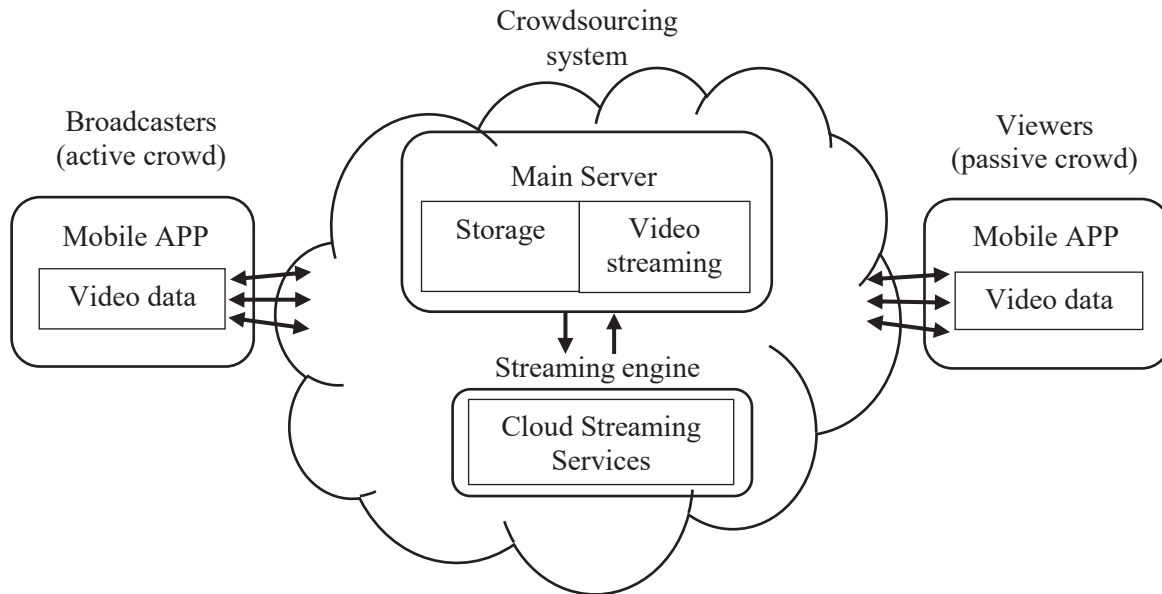


Fig. 1. General mobile crowdsourcing architecture for video streaming.

who generate and upload the video films from the football games using their mobile phones or tablets, platform for processing the video data and a passive crowd of viewers (i.e., users who watch). For the users, a specialized mobile application for video streaming is needed.

The crowdsourcing platform is assuming to have at least two subsystems or components – a main server and a streaming engine. The last one is expected to provide cloud services for multi scalable video streaming with adjustable video quality depending on the internet connection speed and type of the current mobile device. The described simpler approach to create video streaming from a mobile device, as shown in Fig. 1, is the more standard and thus typical for any live stream application.

3. The Footlikers crowdsourcing application

We further describe the author's system Footlikers for live video streaming among registered users of the system. We will note that user participation is voluntary and the user can choose their own type of crowdsourcing activity. For the broadcaster this also depends on the permission of the football club authorities.

3.1. System functionalities

For the creation of the system the following functional requirements were specified:

- To ensure secure access control the existing users of the platform Footlikers must be able to login in the mobile application by using their current account registration in Footlikers. New user account registrations can only be handled in the website of Footlikers.
- The application must support only two types of user sessions:
 - **Broadcaster** – can stream picture and sound from camera and microphone of his mobile device by using the mobile application
 - **Viewer** – can watch the live video feed on his mobile device inside the application.
- The viewer must be able to see and consult the week match schedule of his favorite football teams. When a specific match is being broadcasted, there is a green icon indicator which informs the Viewer that he can tap on the selected line and initiate the reception of the live feed from the Broadcaster.
- There is no limit regarding the number of viewers that want to access and use the application to view matches. Every user of the mobile application has access to every football match which is being broadcasted, not only his own favorite teams. There is a built-in search engine which helps the viewer to do that. No payment is required for the broadcasting and viewing of live video feed at this moment.
- The Broadcaster must be able to choose a desired football match from the match schedule, for which he has permissions to broadcast. These permissions are set via the administrative club panel from the website by the club owners or official club administrators. When he is ready with the selection, the live stream is launched. The Broadcaster must also receive picture and sound from his own broadcast on his mobile device. The Broadcaster must be able to fully

control the live feed, by start, stop or pause functionalities.

- Every type of user of the mobile application has the possibility to use the search engine to find a specific football match from the week's match schedule. The search engine is working with the name of the desired team in a football match. The personal team favorites and preferences are by default proposed to the users on top of the match schedule for the day.
- The match schedule is created in the website Footlikers [12]. The changes of the match program are on a weekly basis.
- The live broadcast is only available to users that are using the mobile application of Footlikers.com

3.2. Football match scenarios

A football match has a multitude of different statuses and scenarios, which could be taken into account:

1. Match start
2. End of first-half
3. Start of second-half
4. Match end

These are normally the standard four stages of any football match. After which, depending on the match result and the type of the competition it can continue with additional stages:

5. Start of first additional time
6. End of first additional time
7. Start of second additional time
8. End of second additional time

If the winner is still not decided then a last set of stages is available:

9. Penalty session start
10. Penalty session end
11. Match end

For processing the video transmissions, these basic time stages require additional processing of user requests to the server.

3.3. Architecture

The Footlikers system architecture to deliver the above-mentioned functionalities is shown in Fig. 2. It follows the general approach from Fig. 1. The first of the main components is the Footlikers mobile application, designed for the users. The mobile application was always sending requests to the main server footlikers.com [12]. This is realized via the RESTFUL API (Representational State Transfer), an application program interface (API) that uses different HTTP requests. This is followed by the initialization and usage of the WOWZA streaming server engine [14], operating within AMAZON EC2 Machine, to deliver the video stream to all Viewer mobile devices. So, in fact, any kind of info throughout the football game was being handled by RESTFUL API. The communication between the users and the WOWZA engine is processed with the Real-Time Messaging Protocol (RTMP) [15].

4. Software implementation

The main software development of the mobile application is based on the utilization of ADOBE AIR platform for developing hybrid mobile applications [16]. The programming language which is used to code the business logic is ActionScript. It is then pre-compiled and re-compiled to build two versions: one for iOS and one for Android.

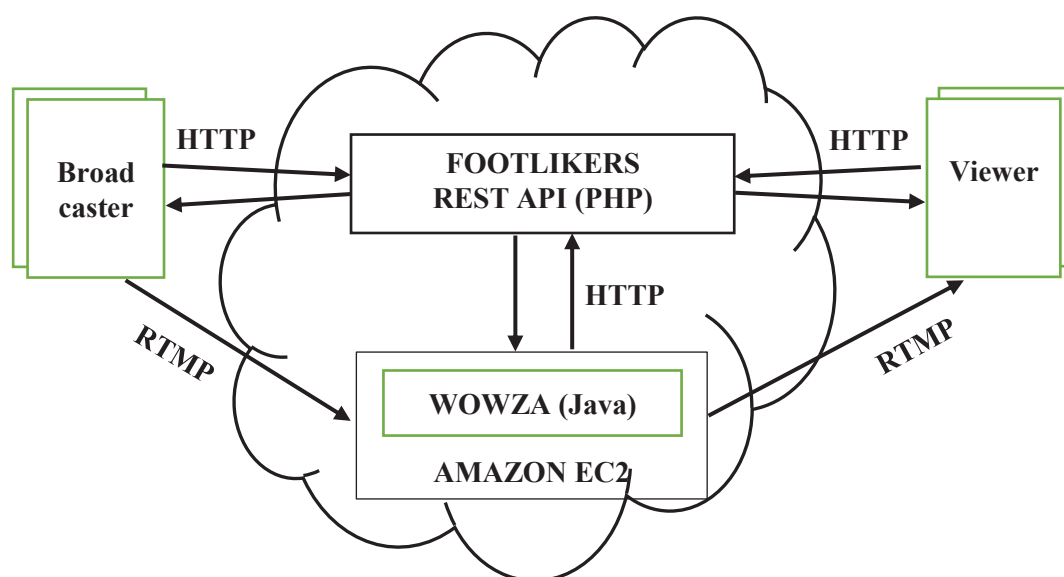


Fig. 2. Footlikers mobile crowdsourcing architecture for video streaming.

When we arrive at the final screen of the application, depending on our session (Broadcaster/Viewer), the mobile application uses NetConnection, which is a native class from the ADOBE AIR Framework, which initiates an open connection to the WOWZA server. This is followed by NetStream via RTMP. In general we use the standard build-in framework classes of ADOBE AIR like: APIManager, VideoStream Manager and others.

The user interface of the mobile application is in French. The welcome screen with login fields is shown in Fig. 3.

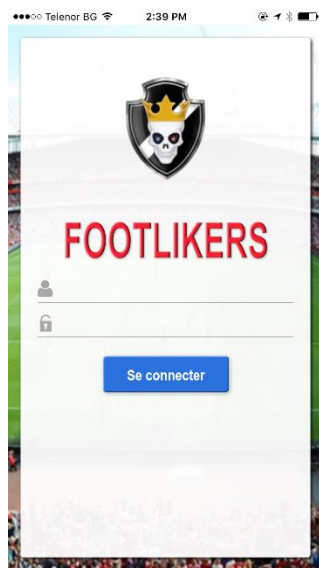


Fig. 3. Footlikers Login screen.

The login requests towards the RESTFUL API are being handled via a specific URL depending on the type of event and access token which has to pass to guarantee credentials to the API. After a successful login, we have the choice between the two modes of the application: Broadcaster or Viewer as shown in Fig. 4.

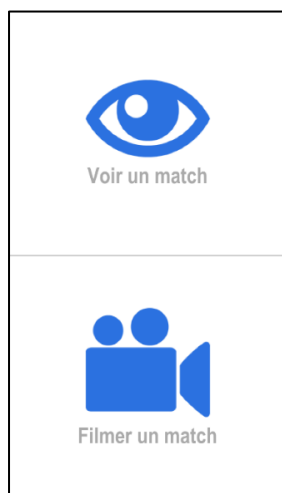


Fig. 4. Mode selection screen.

Depending on the choice, the UI adapts accordingly and shows the next screens, which are of course different for Broadcaster and Viewer. The following screen (Fig. 5) is with the match schedule for Broadcaster, and analogous is for the Viewer (Fig. 6). It is the event which happens after you tap on the football match which differs.

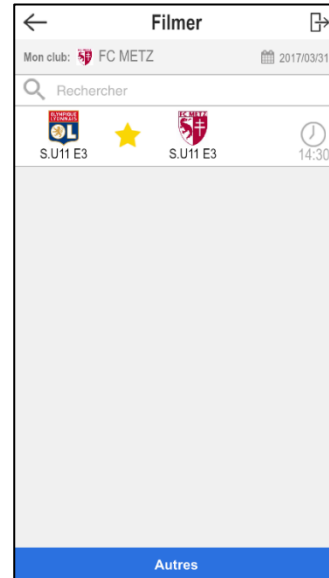


Fig. 5. Broadcaster match schedule screen.

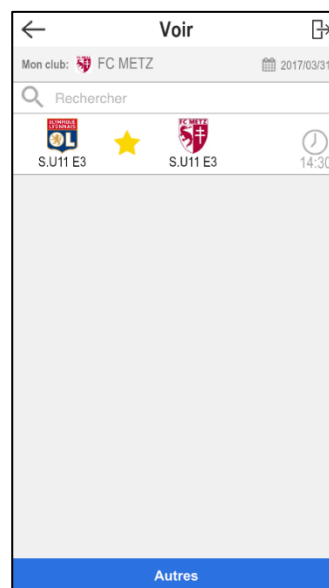


Fig. 6. Viewer match schedule screen.

After selecting a possible football match, NetConnection starts and the application will demand for microphone and camera permission from the device, depending on iOS/Android settings of the user. From this point on we trigger the video stream protocol RTMP and receive/send video and audio data with our mobile device.

Fig. 7 shows an example of Broadcaster screen and Fig. 8 shows the Viewer screen.

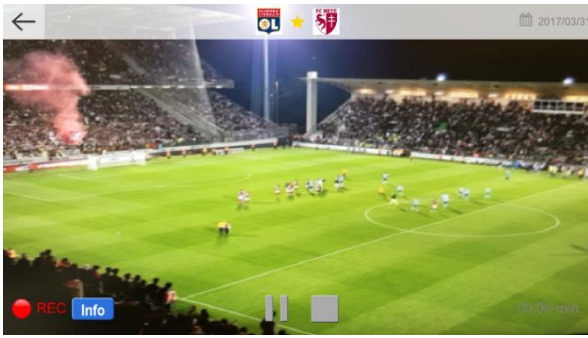


Fig. 7. Broadcaster filming screen.



Fig. 8. Viewer watch screen.

As for the built-in search engine of the application it is working as it follows: we search for “Paris Saint Germain”, the engine reacts on each keystroke of the touch screen keyboard and triggers an event which makes a query that locally searches in the application. To make the search engine much more economic in terms of requests, it does not send or receive any data. The whole match schedule of the week is downloaded as a JSON object when the user authenticates on screen. So the search is actually a basic search on a prefetched dataset. To access the search screen, there is a blue button “Other” (“Autres”) at the bottom as shown on Fig. 5 and Fig. 6.

5. Conclusion

The proposed general architecture and the developed crowdsourcing system for video streaming is exclusively designed for a specific social group – football fans, coaches, staff, journalists and other people interested in football competitions of amateur leagues in France, Belgium and Luxembourg.

In general, for this kind of sport events, there always has been a lack of multimedia tools so that football fans watch their favorite local teams without having to go to the stadium. But the users that are registered in Footlikers will have this ability to do so and to be always with their team when they can't attend the match.

In terms of software implementation, the developed crowdsourcing application Footlikers

delivers a vast range of services which are strictly planned and organized.

The proposed architecture can be easily scaled and altered to include many other features like posting comments, sharing moments from the match, transmitting pictures etc.

Another feature which is being planned is the possibility to stream not only to mobile devices, but also to users logged in the Footlikers website, that will benefit from larger and more convenient screen.

ACKNOWLEDGEMENT

This research was partially supported by the Scientific and Research Department of University of Plovdiv Paisii Hilendarski, grant MU17-FMI-003.

REFERENCES

1. Monahan, T. and Mokos, J.T. (2013). Crowdsourcing urban surveillance: The development of homeland security markets for environmental sensor networks. *Geoforum*, volume 49, pp. 279-288.
2. Ganti, R.K., Ye, F., and Lei, H. (2011). Mobile crowdsensing: current state and future challenges, *IEEE Communications Magazine*, volume 49, 11, pp. 32-39.
3. Poblet, M., García-Cuesta, E., and Casanovas, P. (2014). Crowdsourcing tools for disaster management: A review of platforms and methods. In Casanovas, P. et al. (eds.) *AI Approaches to the Complexity of Legal Systems*, pp. 261-274. LNCS 8929. Springer, Berlin.
4. Peng, X., Gu, J., Tan, T.H., Sun, J., Yu, Y., Nuseibeh, B., and Zhau, W. (2018). CrowdService: Optimizing Mobile Crowdsourcing and Service Composition. *ACM Transactions on Internet Technology (TOIT)*, volume 18, Issue 2, Issue-in-Progress, Article No. 19, ACM New York, NY, USA.
5. Brabham, D.C. (2013). *Crowdsourcing*. The MIT Press, Cambridge, Massachusetts.
6. Prpić, J., and Kietzmann, J. (2018). Crowd Science 2018: The Promise of IT-Mediated Crowds. In *Proceedings of the 51st Hawaii International Conference on System Sciences, (HICSS)*, pp. 4085-4093. <http://scholarspace.manoa.hawaii.edu/bitstream/10125/50402/1/paper0515.pdf>
7. Hetmank, L. (2013). Components and functions of crowdsourcing systems – a systematic literature review. *Wirtschaftsinformatik Proceedings 2013*, volume 4, pp. 55-69.
8. Fuchs-Kittowski, F. and Faust, D. (2014). Architecture of mobile crowdsourcing systems. In: Baloian, N., Burstein, F., Ogata, H., Santoro, F., and Zurita, G. (eds.) *Collaboration and Technology*. CRIWG, 2014. Lecture Notes in

- Computer Science, volume 8658, pp. 121-136. Springer, Cham.
9. Pereira, R. and Pereira, E.G. (2016). Video streaming: Overview and challenges in the internet of things. In Dobre, C. and Xhafa, F. (eds.), *Pervasive Computing, Next Generation Platforms for Intelligent Data Collection*, pp. 417-444.
 10. Richardson, L., Ruby, S., and Hansson, D.H. (2007). *RESTful Web Services*. O'Reilly Media, Sebastopol.
 11. Christensen, J.H. (2009). Using RESTful web-services and cloud computing to create next generation mobile applications. In *Proceedings of OOPSLA '09, 24th ACM SIGPLAN conference companion on Object oriented programming systems languages and applications*, ACM, pp. 627-634, Orlando, Florida, USA.
 12. *Footlikers*. www.footlikers.com
 13. Mickulicz, N.D. Narasimhan, P. and Gandhi, R. (2013). YinzCam: Experiences with in-venue mobile video and replays. In *Proceedings of LISA '13: 27th Large Installation System Administration Conference*, November 3–8, 2013, Washington, D.C., pp. 133-143.
 14. *WOWZA Streaming Engine*. <https://www.wowza.com/products/streaming-engine>
 15. Parmar, H. and Thornburgh M. (eds.) (2012). *Adobe's Real Time Messaging Protocol*. Adobe. http://www.images.adobe.com/content/dam/Adobe/en/devnet/rtmp/pdf/rtmp_specification_1.0.pdf
 16. *Adobe AIR technology*. <http://www.adobe.com/products/air.html>

Contacts:

Georgi Iliev

University of Plovdiv Paisii Hilendarski

E-mail: iliev86@gmail.com

FAST GAUSSIAN FILTERING FOR SPEED FOCUSED FPGA BASED CANNY EDGE DETECTION COMPUTATIONS

DIMITRE KROMICHEV

Abstract: *The paper is focused on a novel computational approach guaranteeing fastest execution of Gaussian smoothing to be implemented in a speed orientated FPGA based Canny. Analyzed are the two capital dimensions of FPGA speed accomplishment - maximum clock frequency and minimum clock cycles required for obtaining a correct result, in view of integer arithmetic capabilities and advanced organization of computations within the weighted average filter. Presented is a new integer division replacing algorithm ensuring optimal calculation options in terms of speed. Exact speed gain is set forth in comparison to the available Gaussian smoothing realizations*

Key words: *FPGA, Canny, Gaussian filtering, speed, organization of computations, integer arithmetic, clock frequency, clock cycle, speed gain*

1. Introduction

For speed orientated FPGA based Canny Gaussian smoothing represents a focal point due to the following specific aspects: 1) It commences the Canny computations, and thus generally determines the efficiency of pipelining in terms of speed; 2) The weighted average filter's size can vary vastly, the sole restriction being for the side to be an odd number; 3) The weighted average filter's coefficients' magnitudes can vary widely, and can be on the order of several thousands; 4) Mathematically, the function calculating the weighted average pixel value includes integer arithmetic operations that are indispensable for all of the Canny algorithm; 5) Actually, it is one of these operations that defines the highest clock frequency which is to be employed as a timing analysis constraint for testing the whole bulk of Canny computations's being tangible/intangible with respect to accomplishing the optimal speed goal; 6) Being a square neighbourhood operation, it encompasses all the capital parameters impacting Canny computations in terms of speed; 7) It brings to the fore the demand for a reliable mechanism which guarantees the invulnerability of Canny's organization of computations to the variability of input image and filter's matrix in terms of size.

All these facts having been taken into account, the objective of this paper is to propose an advanced technology of Gaussian filtering focused entirely on speed. The task is to define the optimal clock frequency and most advantageous number of

clock cycles for the smoothing function to execute accurately and reliably, and on that basis, a novel organization of computations will be set forth. Comparison with the available realizations of Gaussian smoothing will be drawn to exhibit the achieved speed gain. Relevant to the conducted experiments and conclusions arrived at are only gray-scale images. The targeted hardware is Intel (formerly Altera) FPGAs (hereafter referred to only as Altera FPGAs). Quartus II TimeQuest Analyzer is used for setting timing analysis constraints and testing the feasibility of the proposed computational approach.

2. Literature survey

All the implementations of Gaussian smoothing for FPGA based Canny described in the literature boil down to several computational approaches: 1) Sequential. It generally features traversing successively all of the square neighbourhood pixel by pixel to convolve with the required coefficients in the filter matrix with each multiplication result being added to the previous one stored in a register. Then comes the division by the normalization factor [8]. Main flaw in view of speed: privation of pipelining and parallel calculations; total dependence upon filter's size; repetitive use of pixels for filtering a single image row/column 2) Separability orientated. Gaussian filter being separable, for a coefficient matrix of size $Z \times Z$ two consecutive sets of multiplications are applied: $1 \times Z$ followed by $Z \times 1$. The results of

multiplication having been summed up, division by the normalization factor is executed. This approach is implemented in two variants: with [5] or without distributivity [7]. Capital flaws in terms of speed: extensive use of sequential logic for intermediate storage and successive multiplications including the same square neighbourhood pixels; repetitive use of pixels for filtering a single image row/column; 3) Symmetry orientated. The total number of different coefficient magnitudes in a Gaussian filter is

$$N = [(Z+1)*(Z+3)]/8 \quad (1)$$

where

- N is the number of different coefficient magnitudes;
- Z is an odd number; it is the side of a Gaussian filter of particular size.

Following this, all image pixels positioned under filter matrix coefficients of identical magnitude are summed up, and then the multiplications by the targeted coefficients are realized. Then all the multiplication results are summed up. Finally, division by the normalization factor is performed. Here, there are several variants: all pixels pertaining to a coefficient magnitude are summed up, and then all multiplications are realized simultaneously [6]; different coefficient magnitudes being distributed unevenly across the filter matrix, and the image pixels being set at the Gaussian filter's input in sequential order, when all the pixels pertaining to a certain coefficient's magnitude are summed up the multiplication by the targeted coefficient is realized immediately [4]; column/row-wise multiplications and additions with intermediate storage.[6]. With speed being in focus, the flaws of this approach are: layers of addition with their number being directly proportional to the size of filter matrix; repetitive use of pixels for filtering a single image row/column.

1. The proposed technology of Gaussian smoothing targeting optimal speed

3.1. Integer arithmetic performance

Speed performance of Canny on FPGA has two standard peremptory dimensions - maximum clock rate and minimum clock cycles to calculate a plausible result. The former being a priority, the optimal operating frequency of the Gaussian smoothing indispensable integer arithmetic should be tackled in the first place. The filtering function's mathematics demands multiplication, addition and division, the latter being totally incompatible with accomplishing the speed goal in terms of its conventional execution on Altera FPGAs [1][2]. Compared to addition, it is way too slow under the same test conditions. Therefore, we have

implemented an integer division replacing algorithm which guarantees less than half the clock cycle at the highest clock frequencies the Altera FPGA families can sustain. It is based on bit slicing and weighted average filter coefficients' modification by employing a filter specific constant calculated in advance. This algorithm has no division procedure and the only mathematics it includes is a single addition of the rounding bit to the integer number represented by the sliced set of consecutive bits. Therefore, the utilized integer division replacing algorithm is not impacted by the magnitude of the dividend, and can be entirely relied upon in terms of being the fastest among all of the Gaussian smoothing arithmetic operations, whatever the size of the filtering matrix and the coefficients' values. Total mathematical accuracy of results is ensured.

The optimal clock rate contest between the remaining two integer arithmetic operations - multiplication and addition, is in favour of the latter. On Altera FPGAs the "+" operator and the LPM function are realized by default as ripple carry adders, which, the dedicated carry chains being taken into account [1][2], provide the fastest results for the purposes of Canny implementation. Thus, of all the integer arithmetic Gaussian smoothing is based upon, it is the multiplication executed employing the embedded hard multipliers that proves to be the slowest, and consequently serves as a reference point for determining the highest clock frequency to be used as constraint in the timing analysis.

To recapitulate, there are two practical limits for the highest clock frequency Altera FPGA based Canny, and Gaussian smoothing as its module, can achieve - the performance of on-chip memory and the performance of hard multipliers. A prerequisite to weighted average filtering's being definitive for the entire Canny algorithm is reducing all successive computational modules' integer arithmetic operations to a very limited set comprising only the fastest available, which are rated in descending order as follows: division (using the integer division replacing algorithm), addition, subtraction (it is as fast as addition), and hard multiplication. With respect to the fact that a hard multiplier performance never exceeds that of the on-chip memory [1][2], there is only one conclusion to be drawn: the optimal clock rate of entire Canny algorithm is solely determined by the maximum clock frequency of the embedded multipliers for the highest speed grade on a particular Altera FPGA family.

Scrutinizing the entire range of Altera's FPGAs, considerable differences in performance are to be encountered with most of the families when it

comes to the different operational modes of hard multipliers. The latter depend on input data width. The fastest - that of 9x9 configuration, is not applicable taking into account the possible Gaussian filters' coefficients' magnitudes. Thus, although there are options for larger input data widths available for hard multipliers in the upper class Altera FPGAs, it is actually the 18x18 hard multiplier operational mode in its highest speed grade that sets the standard for optimal clock frequency for Gaussian smoothing, and consequently, for all of the Altera FPGA orientated Canny.

Dealing with the other optimal speed requirement formulated as the minimum amount of clock cycles taken to achieve accurate result at the highest clock frequency, Gaussian filtering integer arithmetic presents the following picture. Targeting a single clock cycle execution, it is of crucial importance not to use input and output registers for the appropriate multiplication LPM_MULT function [3]. Otherwise, whatever the clock frequency, this integer arithmetic operation will require two additional clock cycles.

The hard multiplier's maximum frequency is the actual measure for the minimum clock period as a reciprocal to this frequency. It has to be long enough to satisfy the generalized indispensable timing analysis inequality:

$$T_{clk} \geq D_{clk-to-Q(1)} + D_{logic} + T_{set-up(2)} \pm T_{skew} \quad (2)$$

where

T_{clk}	is the clock period
$D_{clk-to-Q(1)}$	is the clock-to-output propagation delay of the prepositioned register;
T_{logic}	is the propagation delay of the combinational logic (the hard multiplier);
$T_{set-up(2)}$	is the set up time of the postpositioned register.
T_{skew}	is the delay of the signal between the clock inputs of the two registers.

In view of the fact that in Gaussian filtering it is only T_{logic} that is practically manageable in terms of selecting the appropriate hard multiplier input widths, two main consequences are to follow: 1) the total propagation delay of any combinational logic in the Gaussian module as well as the entire Canny cannot exceed the propagation delay of the employed hard multiplier; 2) the targeted shortest clock cycle can only accommodate a single multiplication, i.e. accumulating the numerator of the Gaussian filtering fraction requires more than one clock cycle to add up two consecutive

convolution results. The exact number of these additional clock cycles depends upon the ripple carry adder's propagation delay, which is directly proportional to the input data width. Therefore, a very serious situation needs to be tackled here.

For equal input data widths the adder has a smaller propagation delay than the multiplier. The difficulty is that both the size and coefficient magnitudes of Gaussian filter can vary enormously. Thus, aiming at proposing an advanced Gaussian smoothing organization of computations, a comfortable positive slack should be guaranteed for arbitrary filter matrix sizes and coefficient values. Therefore, it is quite computationally unreliable to calculate the numerator of the averaging fraction by sequentially adding up all the multiplication results for the input image pixels pertaining to a square neighbourhood directly, as they are set at the multiplier's output - this can lead to values at the adder inputs which will violate the constraints following from (2), and force the system into metastability and failure. In order to ensure that adding a convolution result to the sum of previous convolutions results will be executed within a single clock cycle under all conditions, a substantially different approach is needed. Its feasibility is due to the fact that the delay of the proposed integer division replacing algorithm is much smaller than the hard multiplier's delay. And this is to be presented in the paragraphs to follow.

3.2. Organization of computations

3.2.1. Basic assumptions

The optimal speed focused Gaussian smoothing computations require the following basic assumptions: 1) Image pixels are set at the Gaussian filter's input in sequential order at the rate of a pixel per clock cycle; 2) The clock frequency of setting the image pixels at the input is equal to the clock rate Gaussian filter operates at; 3) The image is processed in a row-wise fashion; 4) All image pixels needed for the filtering operation are accessible in a specifically predefined mode.

3.2.2. Integer arithmetic operations' execution order

For an input image pixel to be filtered all the pixel values belonging to the square neighbourhood defined by the Gaussian smoothing matrix is multiplied by the appropriate coefficients. The convolutions are executed sequentially, the way the image pixels are set at the filter's input - one per clock cycle. No distributivity is employed.

Targeting the optimal clock frequency and obeying the restrictions imposed by (2), within a clock period a single multiplication of an image

pixel by a coefficient is executed. The resulting value is stored into a register. The next clock cycle accommodates two consecutive integer arithmetic operations: the convolution value stored in the register is divided by the Gaussian filter's normalization factor, and then the division result is added to the previous division results, the achieved value being stored into another register functioning as accumulator. All the image pixels belonging to the square neighbourhood defined by the Gaussian smoothing matrix of size $Z \times Z$ having been traversed by means of applying the proposed integer arithmetic operations' execution order, the value stored into the accumulator register after the $(Z \times Z)$ th integer division is the final value representing the required Gaussian smoothed image pixel. The flow of computations being entirely pipelined, starting from clock cycle # 2 each clock cycle a multiplication, a division, and an addition are realized.

In terms of speed, there are two facts of crucial importance. Whatever the size of Gaussian filter and its coefficients' magnitudes, the combined propagation delays of integer division and addition is always smaller than the propagation delay of the hard multiplier. This is due to the following:

- 1) The employed integer division replacing algorithm based on bit slicing requires the adding of a single bit to a set of consecutive bits whose count is always significantly smaller than the multiplier's inputs. With the fast dedicated carry chains on Altera FPGAs, as well as the circumstance that addition is faster than multiplication under the same test conditions, taken into account, it is a very small portion of the hard multiplier's delay that the division replacing algorithm's delay is equal to.
- 2) The maximum value resulting from the integer division can be 255. Thus, in the whole computational process of accumulating the division results for the purpose of calculating the final smoothed pixel value, in terms of propagation delay, the worst possible case scenario is an 8-bit operand at the adder's input. Actually, the larger the filter matrix, the smaller the values to be accumulated through addition. Therefore, the difference between multiplier's delay and integer division replacing algorithm's delay is always large enough to accommodate the combined transport delay and adder's delay.

Consequently, the proposed execution order of arithmetic operations guarantees that the optimal clock frequency defined by the hard multiplier's topmost performance cannot be impacted by the size of the Gaussian filter and the magnitude of its coefficients. So, this technology ensures that the optimal clock rate of the entire FPGA orientated Canny is independent from the

smoothing matrix numerical characteristics. For any number of input image pixels pertaining to a square neighbourhood which are set consecutively at the Gaussian module's input the required count of clock cycles to be processed at the highest clock frequency is expressed as

$$C_{\text{cycle}} = N_{\text{pixel}} + 1 \quad (3)$$

where

- C_{cycle} is the total of clock cycles taken for a set of consecutive image pixels from a square neighbourhood to be processed;
- N_{pixel} is the number of image pixels from square neighbourhood which are set consecutively at the filter's input.

3.2.3. Parallel filtering of consecutive input image pixels

The smoothing of pixel # 1 pertaining to an input image row using a Gaussian filter of size $Z \times Z$ is realized in the following fashion : 1) in clock cycle # 1, pixel # 1 from column # $Z - (Z-1)$ of the square neighbourhood is multiplied by the appropriate coefficient from column # $Z - (Z-1)$ of the filter matrix; 2) in clock cycle # 2, the convolution # 1 result is divided by the normalization factor, and the resulting value is stored in the accumulator register; simultaneously, pixel # 2 from column # $Z - (Z-1)$ is multiplied by the appropriate coefficient from column # $Z - (Z-1)$; 3) in clock cycle # 3, pixel # 3 from column # $Z - (Z-1)$ is multiplied by the appropriate coefficient from column # $Z - (Z-1)$; simultaneously, the convolution # 2 result is divided by the normalization factor, the resulting value is added to the previous division result, and the value is stored in the accumulator register. The pipelined calculations are repeated until the entire column is traversed.

Then the process continues with the pixels and coefficients from column # $Z - (Z-2)$. And here is the point where the parallel smoothing of consecutive input image pixels starts. This is based on the fact that the pixels from column # $Z - (Z-2)$ are at the same time the pixels pertaining to column # $Z - (Z-1)$ in the square neighbourhood of the second input image pixel to be filtered. Thus, in terms of speed, the essential efficiency of the proposed organization of computations is focused on the parallel multiplication of one and the same image pixels by the appropriate coefficients pertaining to different columns of the Gaussian filter.

For a filter matrix of size $Z \times Z$, the count of consecutive input image pixels that are simultaneously being subjected to smoothing is

equal to Z . The latter being an odd number, Gaussian filter features a central column with relation to which the right hand columns represent an outward bound projection of the left hand columns in terms of coefficients' magnitudes. Therefore, for each clock cycle an image pixel will be multiplied by Z number of coefficients of which $(Z-1)/2$ number of coefficients have one and the same magnitude. Thus, for each clock cycle actually $(Z-1)/2$ number of multiplications prove to be definitely redundant. As an immediate consequence, the total count of hard multipliers employed in the parallel smoothing is considerably reduced to $(Z+1)/2$. The value resulting from a single multiplication is directly used in the filtering computations involving all the input image pixels for which, in a particular clock cycle, the smoothing matrix coefficients coincide in terms of magnitude.

In accomplishing the goal of Gaussian filtering optimal speed, the decisive advantage of the proposed organization of computations is defined by the fact that an input image pixel value is used only once during the entire bulk of calculations demanded for the smoothing of a whole input image row. In this way, applying this computational approach, smoothing an image pixel by utilizing a Gaussian matrix of size $Z \times Z$ requires exactly

$$Z+1 \quad (4)$$

clock cycles at the maximum clock frequency that is tangibly guaranteed by the hard multiplier performance for a particular FPGA family.

4. Evaluating the proposed approach to Gaussian smoothing in terms of speed

4.1. Research methodology

The feasibility of the presented advanced technology of Gaussian filtering should most plausibly be proved and reliably assessed on the platform of its being compared with other methods that have been used and described in the literature so far. The Gaussian matrix utilized to conduct the analyses and arrive at the conclusions has the following parameters: size - 5×5 , σ - 1.4, normalization factor - 159. Altera's Quartus Version 12.0 and TimeQuest Analyzer are employed in scrutinizing the proposed technology's computational accuracy and setting the constraints required to reveal the veracity of its optimal speed capabilities. The targeted Altera FPGAs this research deals with are the 130 nm, 90 nm, 65 nm, 40 nm, and 28 nm Cyclone and Stratix families.

The experimentally secured results focus exhaustively on the basic optimal speed indicators – maximum clock frequency of execution and minimum clock cycles demanded for the Gaussian

smoothing of a single input image pixel. The reference value for maximum clock rate is determined on the basis of the 18×18 hard multiplier performance in its highest speed grade for a particular FPGA family. Taking into account this numerical limitation, the comparative evaluation benefits from studying and analyzing the impact of several major integer division techniques on the two basic speed parameters.

4.2. Results and analyses

The achieved results for the highest clock frequencies of Gaussian filtering computations for the targeted Altera FPGA families are shown in Table 1.

Table 1. Maximum clock frequencies of Gaussian smoothing

FPGA family	18x18 hard multiplier Fmax (MHz)	Gaussian smoothing Fmax for various integr division techniques (MHz)		
		Altera LPM_DIVIDE (Divider) IP Core	Multiplication by the reciprocal of the divisor	Our algorithm
Cyclone I	187	28	88	187
Cyclone II	246	32	121	246
Cyclone III	279	39	136	279
Cyclone IV	282	41	139	282
Cyclone V	286	44	277	286
Stratix I	217	34	107	217
Stratix II	367	57	181	367
Stratix III	458	70	442	458
Stratix IV	534	81	521	534
Stratix V	595	97	588	595

The total amount of clock cycles required for filtering a single image pixel with Gaussian matrix of size 5×5 through employing various computational approaches are shown in Table 2.

Table 2. Number of clock cycles required for smoothing a single image pixel with 5×5 Gaussian filter

Approach	Hard multiplier Fmax applicable or not applicable	Total number of clock cycles employing various integr division techniques		
		Altera Divider IP Core	Multiplication by the reciprocal of the divisor	Our algorithm
Sequential	No	32	28	26
Separability without distributivity	No	34	30	28
Separability with distributivity	No	30	28	26
Symmetry with simultaneous multiplication	No	32	28	26
Symmetry with sequential multiplication	No	34	39	28
Symmetry with intermediate storage	No	22	18	16
Our organization of computations	Yes	-	-	6

These results demonstrate that the proposed advanced organization of computations has no match on a comparative basis in terms of both maximum clock frequency as defined by the hard multiplier performance and minimum amount of clock cycles demanded for smoothing an image pixel. Taking into account the combined impact of the two basic speed parameters on the overall speed enhancement, the comparative evaluation shows that the presented technology of Gaussian smoothing secures a speed increase on the order of 6.1 - 12.7 times across the entire bulk of targeted Altera FPGA families. That fact renders the proposed computational approach optimal as well as most reliable for the purposes of a speed focused FPGA based implementation of Canny.

5. Conclusion

Presented in this paper is an advanced technology of Gaussian filtering computations aimed at securing optimal speed. Integer arithmetic performance is scrutinized and the upper speed limit of Gaussian calculations, and therefore the entire FPGA based Canny, is defined. Employed is a new purposely designed integer division replacing algorithm based on bit slicing which allows for an appropriate ordering of arithmetic operations so that the maximum clock frequency is always guaranteed.

A new organization of computations is set forth. It is capable of ensuring the parallel filtering of consecutive image pixels whose number is equal to the side of the smoothing matrix. The proposed technology is evaluated on a comparative and experimental basis to prove its being reliable as well as tangible for the purposes of securing optimal computational results in terms of the two capital speed parameters – maximum clock rate and minimum clock cycles required for a single image pixel to be filtered.

REFERENCES

1. Altera Corporation. Cyclone I - V Device Handbook. 2014, Volume 1
2. Altera Corporation. Stratix I - V Device Handbook. 2014, Volume 1
3. Altera Corporation. Integer Arithmetic IP Cores User Guide, 2014
4. Chandrashekar N.S., K. R. Nataraj. Design and Implementation of a Modified Canny Edge Detector based on FPGA, *International Journal of Advanced Electrical and Electronics Engineering*, (IJAEET). 2013, Vol.2, (1), pp. 17-21
5. Chandrashekar N.S., K. R. Nataraj. NMS and Thresholding Architecture used for FPGA based Canny Edge Detector for Area Optimization, *Proceeding of International Conference on Control, Communication and Power Engineering*. 2013, pp. 80-84
6. Divya. D., P. S.. Sushma. FPGA Implementation of a Distributed Canny Edge Detector, *International Journal of Advanced Computational Engineering and Networking*. 2013, Vol. 1, (5), pp. 46-51
7. Shamlee V., Jeyamani.. A Split Canny Edge Detection: Algorithm and its FPGA Implementation. *International Journal of Science and Research (IJSR)*. 2014, Vol. 3, pp. 1198 -1205
8. Veeranagoudapatil, Chitra Prabhu (2015). Distributed Canny Edge Detector: Algorithm & FPGA Implementation, *International Journal for Research in Applied Science & Engineering Technology (IJRASET)*, Vol. 3 (5), pp. 586-588

Plovdiv University Paisii Hilendarski
Plovdiv 4000
24 Tzar Asen Street
e-mail: dkromichev@yahoo.com

EFFICIENT COMPUTATION OF ORTHOGONAL GRADIENTS TO BE USED IN SPEED FOCUSED FPGA BASED CANNY EDGE DETECTION

DIMITRE KROMICHEV

Abstract: *In a speed focused FPGA orientated Canny implementation, the orthogonal gradients module ought to secure both reliable results for the next stage to work with and efficient organization of computations capable of guaranteeing the execution of integer arithmetic for minimum amount of clock cycles at the highest clock rate. The goal is to boost pipelining efficiency of the entire Canny algorithm in terms of ensuring that the square neighbourhood based Sobel filtering calculations are up to the task of materializing a non-intermittent mathematically accurate flow of computations once the processing has started.*

Key words: *orthogonal gradients, FPGA, Canny, algorithm, pipelining, clock frequency, clock cycle*

1. Introduction

Canny is a gradient based contour detector. The two most important requirements for its FPGA based hardware implementation are accuracy and speed, the former being a function of the mathematical exactness of computational mechanisms within each of the modules and framed by the utmost utilization of FPGA capabilities. Speed depends upon several factors, among them of tremendous significance being the organization of computation. Sobel filtering is Canny's second stage and a square neighbourhood operation, and as such definitely accounts for both the performance of its own calculations and the speed capabilities of the entire contour detecting algorithm in terms of boosting pipelining efficiency. FPGA functionalities and characteristics should be taken into account by the computational approaches in terms of utilizing those favourable to fast executions whereas avoiding operations and methods serving as bottlenecks. Thus, in FPGA-based Canny every single algorithm requires a multifaceted approach targeting the results' being reliable as well as feasible. Speed is worthy of being set as an accomplishable goal only on the peremptory platform of total mathematical accuracy of calculations thus ensuring the detected contours' veracity.

The objective of this paper is to present an efficient organization of computations for Canny's

orthogonal gradients module aimed at speeding up the FPGA-based image processing calculations.

The task is to describe in detail the sequence of steps, to thoroughly analyze the computational reliability, to expose the characteristics, and to point out the applicability of the proposed computational approach in view of the FPGA functionalities. The targeted hardware is Intel (formerly Altera) FPGAs (hereafter referred to only as Altera FPGAs). Quartus II TimeQuest Analyzer is used for setting timing analysis constraints and testing the feasibility of the proposed computational approach. Relevant to the analyses and conclusions arrived at in this paper are only gray-scale images.

2. Literature survey

Described in the literature are two main approaches addressing Sobel filtering on FPGA. Most commonly, the proposed computational techniques rely only on sequentially adding up the results of multiplying the positive and negative coefficients in the Sobel masks with the Gaussian filtered pixels in the appropriate positions in the square neighbourhood to calculate the x- and y-gradients[1][4][5][8][10][11][12][13][14][15]. This approach is aimed at speeding up execution at the expense of economizing on calculations. The flaws here are: lack of precision; difficulties with tackling the negative values and additional steps to be taken in the next Canny module for the purpose of

avoiding division by zero; taking into account that the computed results can exceed the maximum value of a gray-scale image pixel, an appropriate scaling is demanded. The other approach is applied comparatively more rarely [2][3][6][7][9][16]. Here there is some averaging employed following the multiplication and addition steps - in these cases the divisor is 2. Flaws: redundand operations to work with the negative values and avoid division by zero; need for scaling.

3. Orthogonal gradients' mathematics

Gaussian filtering having been applied, the total number of pixels in the input image is reduced, and the total number of filtered image pixels is :

$$M*N - \{[M - (Z-1)] * (Z-1) + N * (Z-1)\} \quad (1)$$

where

MxN is the input image size
ZxZ is the filter matrix size.

The exact number of available Gaussian smoothed pixels required for this module to commence the parallel computations of the vertical and horizontal gradients applying the two Sobel filters should be:

$$2*(N-(Z-1)) + 4 . \quad (2)$$

The filtered fixels are stored in revolving Dual-port RAMs for the purposes of guaranteeing the access to all the pixels under the Sobel masks within a single clock cycle, thus employing pipelining of computations, and saving memory.

As gradient filters, the two Sobel matrices have positive and negative coefficients. If in a 3x3 neighborhood C(x,y) is the central pixel, and the numbers of neighboring pixels are as shown (Fig. 1),

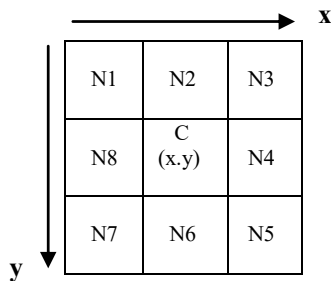


Fig. 1. 3x3 neighborhood pixels

the equations used to calculate the x gradient and the y gradient for C(x,y) are of two types:

1) Exact mathematics - used to accurately calculate the gradients:

$$G_x = \frac{(N_3 + 2N_4 + N_5) - (N_1 + 2N_8 + N_7)}{4}$$

(3)

$$G_y = \frac{(N_7 + 2N_6 + N_5) - (N_1 + 2N_2 + N_3)}{4} , \quad (4)$$

where

N1..N8 are the neighboring pixels in the 3x3 neighborhood as defined by the two Sobel filters (Fig. 2.)

-1	0	1
-2	0	2
-1	0	1

G_x

-1	-2	-1
0	0	0
1	2	1

G_y

Fig. 2. Sobel masks for x and y

The values calculated here are within the interval [-255,255]. There are two important facts here: 1) Division by the sum of positive coefficients in the Sobel masks avoids the need for scaling; 2) The difference in signs of the values calculated for both gradients is applied in the speed optimization of gradient direction computations.

The multipliers and the divisors being a power of 2, integer arithmetic guarantees a maximum speed of computations in this module.

2) Approximation.

$$G_x = [(N_3 + 2N_4 + N_5) - (N_1 + 2N_8 + N_7)] \quad (5)$$

$$G_y = [(N_7 + 2N_6 + N_5) - (N_1 + 2N_2 + N_3)] , \quad (6)$$

where

N1..N8 are the neighboring pixels in the 3x3 neighborhood.

With division by the power of 2 being dropped out, the values calculated here are within the interval [-1020,1020]. Consequently, proper scaling is required for all the results to fit within the maximum gray scale image pixel range, and that can possibly introduce disproportionalities impacting the qualities of the detected contours.

4. Efficient computation of orthogonal gradients

Being a square neighbourhood operation, Sobel filtering can start executing as soon as the condition determined by (2) has been satisfied. From that clock cycle on the organization of computations

should guarantee the non-intermittent securing of x- and y-gradient values at the input of the next Canny module for the pipelining in all of the contour detector to be able to function with optimal efficiency. Therefore, the smallest Gaussian filter's size, which is 3x3, defines the plausible upper speed limit measured in clock cycles for the Sobel filtering to successfully accomplish the goal of ensuring a continuous data flow beneficial to boosting pipelining. Taking into account that the feasible maximum speed of Gaussian smoothing with a coefficient matrix of size $Z \times Z$ is

$$Z+1 \quad (7)$$

clock cycles at the highest clock frequency determined by the optimal hard multiplier's speed in a particular Altera FPGA family, Sobel filtering has the task of securing mathematically accurate x- and y-gradient values at the inputs of the next Canny's stage every

$$(Z+2)th \quad (8)$$

clock cycle. Thus, (8) represents the upper speed limit of pipelining speed for in orthogonal gradients calculation.

Mathematical exactness of results requiring the use of (3) and (4), the proposed efficient organization of computations encompasses the following sequence of steps.

- 1) Being stored in Dual-port RAMs, all the pixels from the square neighbourhoods defined by the two Sobel filter matrices (Fig. 2.) are accessible within a single clock cycle. This allows for parallel execution of addition and multiplication by 2, the latter being achieved through emulating shift left operation.
- 2) In the next clock cycle, addition and multiplication results from the previous cycle are added.
- 3) As a result of the calculations in the previous clock cycle there are four values in stock. Two of them can be only positive or zero, and the other two can only be negative or zero. These values' being positive or negative is determined exclusively on the basis of their being calculated through utilizing coefficients pertaining to particular positions in the Sobel filters – the columns of negative coefficients as opposed to the columns of positive coefficients. Therefore, with their signs being positionally bound, these four values can be processed just as positive integers. The essence of this approach is to calculate which value is larger, and depending on its position information on the signs of the x- and y-gradients is sent directly to the next Canny's module. This is accomplished by means of

executing twelve parallel subtractions and checking with relation to zero.

- 4) The two resulting positive values computed in the previous cycle are divided parallelly by 4 in terms of emulating right shift operation. This secures the exact values of both gradients.

Thus, with this organization of computations, the accurate calculation of orthogonal gradients is completed in 4 clock cycles. Therefore, the condition defined by (8) is satisfied.

5. Conclusion

In this paper, presented and analyzed is an efficient technology for computing the orthogonal gradients to be implemented in speed focused FPGA based Canny. A thorough description of the approach's peculiarities is set forth. The algorithm is scrutinized in terms of its applicability, computational reliability and speed characteristics. The technological capabilities of efficiently avoiding speed eroding computational approaches and satisfying the peremptory demand for mathematical accuracy and plausibility of results pinpoints the feasibility of the proposed technology as a tangible tool for enhancing Canny's performance on FPGA.

REFERENCES

1. Anup Singh Ramprakash Singh Rajput and Samina Jafar (2013). Improved Distributed Canny Edge Detector in VHDL, *International Journal of Electrical, Electronics & Communication Engineering*, Vol. 3 (6), pp. 291-295
2. Aravindh G. and Manikandababu C. S. (2015), Algorithm and Implementation of Distributed Canny Edge Detector on FPGA, *ARPN Journal of Engineering and Applied Sciences*, Vol. 10 (7), pp. 3208-3216
3. Chaithra N.M. and Ramana Reddy K. V. (2013). Implementation of Canny Edge Detection Algorithm on FPGA and displaying Image through VGA Interface, *International Journal of Engineering and Advanced Technology (IJEAT)*, Volume 2 (6), pp. 243-247
4. Chandrashekar N.S. and Nataraj K.R. (2013). Design and Implementation of a Modified Canny Edge Detector Based on FPGA, *International Journal of Advanced Electrical and Electronics Engineering*, (IJAEET), Vol. 2 (1), pp. 16-21
5. Parminder Kaur and Ravi Kant (2014), A Review on: Comparison and Analysis of Edge Detection Techniques, *International Journal of Engineering Research and General Science*, Volume 2 (3) pp. 102-109

- 6 . Ping Zhou, Wenjun Ye, Yaojie Xia and Qi Wang (2011), An Improved Canny Algorithm for Edge Detection, *Journal of Computational Information Systems*, Vol.7 (5), pp. 1516-1523
7. Pooja Ameta and Mahesh Kumar Porwal (2015). A Review on Edge Detection Technique, *International Journal of Advanced Engineering Research and Science (IJAERS)*, Vol 2 (4), pp. 48-53
8. Poonam Dhankhar and Neha Sahu (2013). A Review and Research of Edge Detection Techniques for Image Segmentation, *International Journal of Computer Science and Mobile Computing*, IJCSMC, Vol. 2 (7), pp.86 – 92
9. Qian Xu, Chakrabarti C., Karam L. J. (2011), A distributed Canny edge detector and its implementation on FPGA, *Digital Signal Processing Workshop and IEEE Signal Processing Education Workshop (DSP/SPE)*, pp. 500 – 505
10. Rajwinder Kaur, Monika Verma, Kalpna and Harish Kundra (2014), Classification of Various Edge Detectors, *International Journal of Computer Science and Information Technology (IJCSIT)* Vol. 2 (5), pp.16-23
11. Raman Maini, Himanshu Aggarwal (2009). Study and Comparison of Various Image Edge Detection Techniques, *International Journal of Image Processing (IJIP)*, Vol. 3 (1), pp. 1-12
12. Ramgundewar, Pallavi, Hingway, S.P. and Mankar, K. (2015). Design of modified Canny Edge Detector based on FPGA for Portable Device, *Journal of The International Association of Advanced Technology and Science*, Volume 16 (2), pp. 210-214
13. Rashmi, Mukesh Kumar and Rohini Saxena (2013). Algorithm and Technique on Various Edge Detection: A Survey, *Signal & Image Processing : An International Journal (SIPIJ)* Vol.4 (3), pp. 65-75
14. William McIlhagga (2011). The Canny Edge Detector Revisited, *International Journal of Computer Vision*, Vol. 91 (3), pp. 251-261
15. Worthington, P.L (2002). Enhanced Canny edge detection using curvature consistency, *Proceedings. 16th International Conference on Pattern Recognition*, vol.1, pp. 596–599
16. Veeranagoudapatil, Chitra Prabhu (2015). Distributed Canny Edge Detector: Algorithm & FPGA Implementation, *International Journal for Research in Applied Science & Engineering Technology (IJRASET)*, Vol. 3 (5), pp. 586-588

Plovdiv University Paisii Hilendarski
Plovdiv 4000
24 Tzar Asen Street
e-mail: dkromichev@yahoo.com

A SURVEY OF METHODS AND TECHNOLOGIES FOR BUILDING OF SMART HOMES

GEORGI PAZHEV

Abstract: *This paper explores the main concepts for development of smart home and presents basic architectures employed. From the wide variety of developments in the recent years, a survey of some of the most significant ones has been made. They use technologies such as cloud computing, IoT (Internet of Things), and interfaces such as ZigBee, Bluetooth, Ethernet, WiFi.*

Key words: *smart home, Internet of Things (IoT), Advanced Metering Infrastructure (AMI), Cloud Computing*

1. Introduction

Smart home has been subject of research over the past fifty years, and the interest continues to present days. There are several projects which are created about this topic – from smart grid infrastructure to development of smart home projects based on industrial automation built by commercial products and open source hardware. There are three main categories of smart home (SH) technologies [1]: in-home; home-to-grid; home/grid-to-enterprise.

In-home technologies include local monitoring and control capabilities. They address the intelligent management of devices available in the SH, extracting and utilizing both internal and external information. If present, they provide the optimum usage of the locally produced energy, supplying local loads and injecting the surplus on the grid. Alternatively, as a result of a demand-side management strategy, it reduces local consumption for satisfying external power peak demand, thus improving customers’ profits, due to favorable tariffs.

Home-to-grid technologies include measurement capabilities, remote control and monitoring. These are mostly used to interconnect houses and to connect them with grid operators and utilities, thus enabling reciprocal real-time information exchanges.

Home/grid-to-enterprise technologies are mainly used to link the information generated within the smart home with enterprise services. They support the management of the infrastructure

via decision-support functionality that can be used to apply control strategies.

The basic node of the energy management system, which could be a whole or a part of the smart home is the metering unit. The metering unit represents the interface between the grid and the end user – it is main element of full integration of smart home with the smart grid (Fig. 1) [1]. In the 1990s the utilities began introducing of automated meter reading (smart meter) with the ability of unidirectional communications of the energy consumption to a central unit by means of power-lines or wireless communications, thus yielding significant reductions in billing costs and inaccuracy. A smart meter is a digital, advanced device with high accuracy, control, and configuration functionality with better theft-detection ability. Meter communications can be from the meter to the devices in the buildings and from the meter to the energy utility.

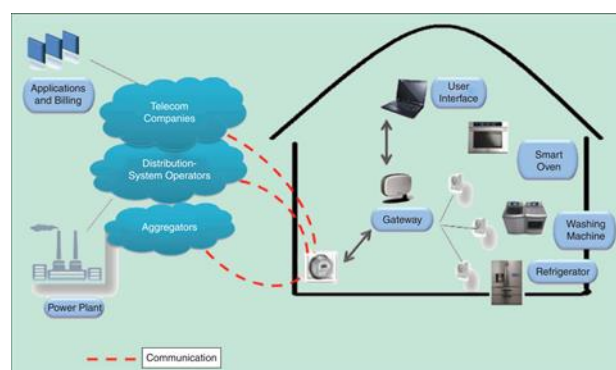


Fig.1. Smart Home with smart metering

AMI (Advanced Metering Infrastructure) represents a full exploitation of meters' capabilities, allowing full automation of the billing process and adding flexibility to time-of-use billing. AMI meters maintain continuous bidirectional communications with the utility and automatically read either on schedule or on demand by the enterprise billing system. This mechanism provides an opportunity for the utility to perform real-time system analysis and gather feedback on power utilization as well as to upload information on the smart meter, allowing local policy aiming for energy management and consumption reduction.

2. Conceptual basics for building of smart homes

Smart homes nowadays are equipped with many devices such as smart meter, in-home displays, renewable energy sources and storage, and smart appliances such as washing machine, refrigerators, TV, oven, thermostat, HVAC, lights, and plugs for electrical cars [2]. A home area network is considered the backbone communication network that connects these devices. It is a two-way communication that is utilized in the demand response, advance metering infrastructure, distributed energy generations and storage. The U.S Department of Energy presents two types of home area network architectures namely: utility managed and utility and consumer managed.

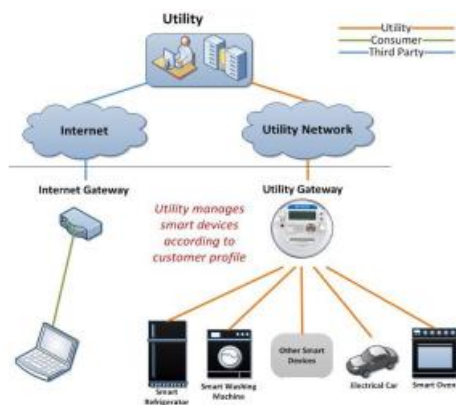


Fig. 2. Utility Managed architecture

2.1. Utility managed architecture

In this architecture, the utility monitors controls and manages smart home appliances via its private network. This architecture is shown in Fig. 2 [2]. The smart appliances are connected to the utility control server via utility gateway and network. The smart appliances send information to a smart meter, which collects and stores information data and then send it to utility control server through utility gateway. The utility is capable of monitoring the power consumption for billing purposes and it can control the appliances. The user

is connected to the utility via internet. He has only access to information that is provided by the utility—he has no control and can only monitor the performance of appliances.

2.2 Utility and consumer managed architecture

This architecture is shown in Fig. 3 [2]. The internet gateway and the utility gateway are connected in home with intermediate hub. The users and the appliances can exchange data and control commands through this common gateway/hub. Home owners can access their home appliances' control system directly through Internet gateway as well as through the utility server. Furthermore, any relevant information about the appliances can also be delegated to a third part through the Internet.

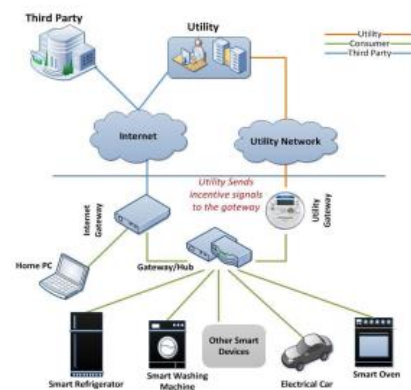


Fig. 3. Utility and Consumer Managed architecture

3. Related work

The smart home is created because of the following main purposes – comfortability, entertainment, health care, energy management and surveillance and security. Chandra Sukanya Nandyala and Haeng-Kon Kim proposed healthcare system architecture, which is based on both Cloud and Fog computing as shown in Fig. 4 [3].

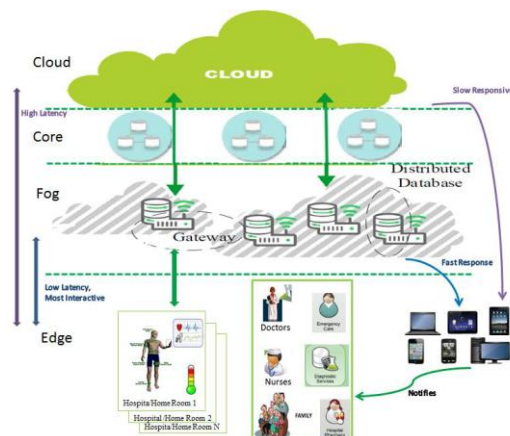


Fig. 4. Fog and IoT-based real-time u-healthcare monitoring

This architecture is based on four tiers: Health Sensor tier, Fog tier, Core tier and Cloud tier. The Health Sensor tier is designed for M2M (Machine to Machine) interactions and collects, process the data, and issues control commands to the actuators. It also filters the data to be consumed locally, and sends the rest to the higher tiers. The Fog (Edge) tier works like sensing, control and correlation. It supports wired and wireless connectivity. This tier must support many different protocols, such as Zigbee, IEEE 802.11, 3G and 4G to accommodate a variety of endpoints. Additionally, this tier must be modular to scale to meet the growth requirements. The components and services offered within one module should be similar so that additional modules can be added in a short span of time. The health sensors from health sensor tier collect data and forward that information to the controllers. The controller can forward any information gathered from the sensors to other devices in the Fog. The health controller is able to process this data locally, analyze and determine optimal health patterns to take action on it. Using this information the controller will send signals to actuators in the system to transmit data or notifies to medical staff and family members via mobile devices.

The Core tier tasks are to provide paths to carry as well as transfer data and network information between numerous sub-networks. The traffic profile is the critical variation between IoT and traditional core network layers.

The Cloud tier tasks are to host applications and to manage the IoT architecture. This tier contains data centers for network management and applications.

For reducing the energy consumption several decisions for building of home energy management system are proposed. Ihsan Ullah suggests an architecture of home energy management system, which is based on load priority for high power appliances like air cooling system, water heating, electric vehicle charger and air conditioning and execution of demand response (DR) events [4]. A DR event is defined as an action taken to alter the electricity demand in response to the changes in the electricity prices over time. A customer that takes part in the DR program can be informed of a DR event by an external signal from the retailer through his smart meter. During a DR event the home energy management algorithm allows the residents to run their appliances as long as the total household consumption remains below the specified demand limit. At the same time home energy management algorithm takes into account

the load priority and customer comfort preferences. Jongbae Kim et al. suggest an IoT (Internet of Things) home energy management system for dynamic home area networks (Fig. 5) [5]. They propose an energy management system based on user-centric service domains where the user-centric services are provided autonomously, based on the contextual information related to users and environments. The user uses a nomadic agent (smart phone or tablet) to enter the service domain. The networked devices (e.g. networked appliances, networked lighting systems, smart meters, and networked PV systems) in home area networks exchange data via the nomadic agent, by constructing peer-to-peer (P2P) connections.

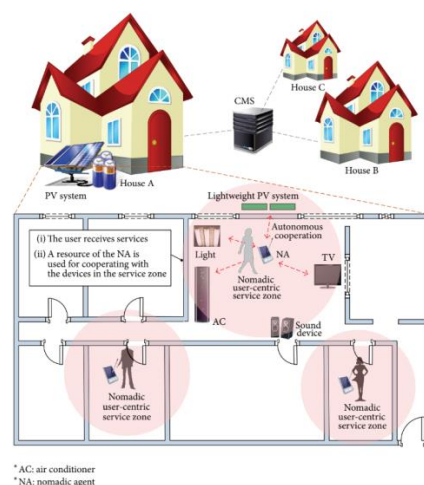


Fig. 5. Home energy management system for dynamic home area networks

The implementation of smart homes is based on creation of heterogeneous sensor networks. There are variety of network interfaces and protocols for building automation. The sensor network is using gateways to provide transparent work of the applications. The sensors and actuators provide variety wired and wireless interfaces like I²C, 1-wire, Bluetooth, ZigBee, WiFi and so on. The sensors measure the parameters of environment and send data to a control gateway. This gateway collects data, stores and analyze them and could forward them to cloud network for the further analysis. Variety of proposals exist for building home automation. Zubir Nabi at al. suggest a cloud based smart home environment named Clome. The Clome architecture is shown in Fig. 6 [6]. All entities inside the house are connected to the outside world through a programmable network switch. All applications and smart appliances have a CPU-heavy end staged in the cloud represented by S and a thin-client end represented by C. Applications in the cloud can make use of both simple storage and a transactional database. Users interact with and control the system and all applications through a

natural user interface. Traditional devices such as PCs, smart phones, and tablets are also connected to the same network and can also be used to access applications staged on the cloud. The implementation of this architecture is based on high speed OpenFlow gateway switch which connects the home appliance to the cloud. For implementation of natural user interface is used a Microsoft Xbox Kinect device, which contains built in RGB camera, depth sensor and microphone for detection of user gestures and voice commands.

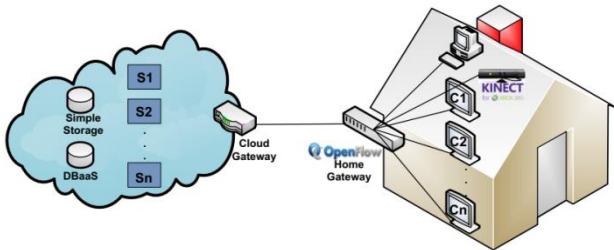


Fig. 6. Clome Architecture

Anindya Maiti suggests another cloud based home automation model called Home Automation as a Service (HAaaS) [7]. This cloud service architecture is based on PaaS, where computer hardware, operating system, data storage and network bandwidth are outsourced, while application and data are managed by the HAaaS (Fig. 7) [7].

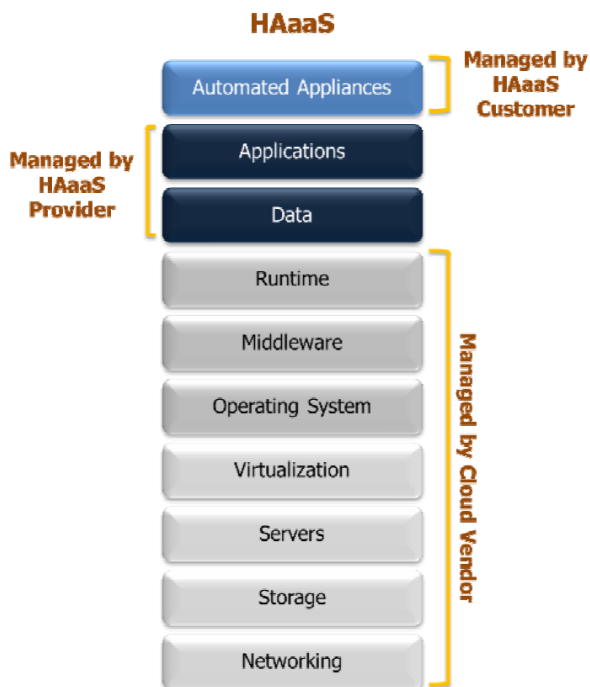


Fig. 7. Home Automation as a Service Architecture

For realizing of amalgamation of cloud and home automation the bridge link is the internet and there are considered two approaches to realize them

– using an internet gateway or using Internet of Things (IoT).

Chih-Yung Chang et al. suggest a smart home architecture based on created IoT Access Point [8]. Fig. 8 [8] depicts the application scenario of the IoT AP for smart life where three network types are considered. The first one is Ethernet, which allows an Access Point connecting to Internet; the second one is Wi-Fi, which provides Internet connection for handheld devices via an Access Point; the third one is ZigBee, which is characterized by low-power and commonly embedded in sensors for environmental monitoring or event detection.

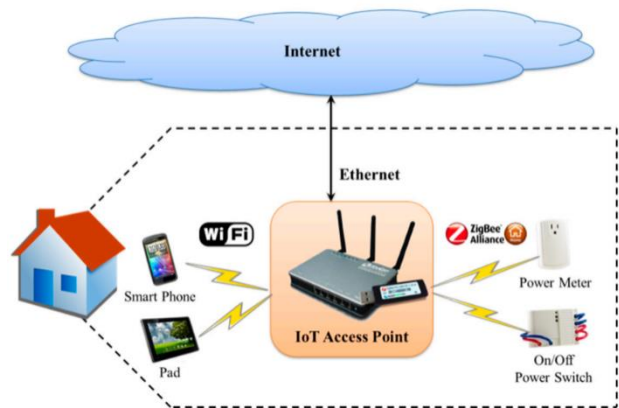


Fig. 8. Smart Home Architecture based on IoT Access Point

Moataz Soliman et al. suggest a smart home architecture, which integrates Internet of Things with Web services and Cloud Computing as shown in Fig. 9 [9].

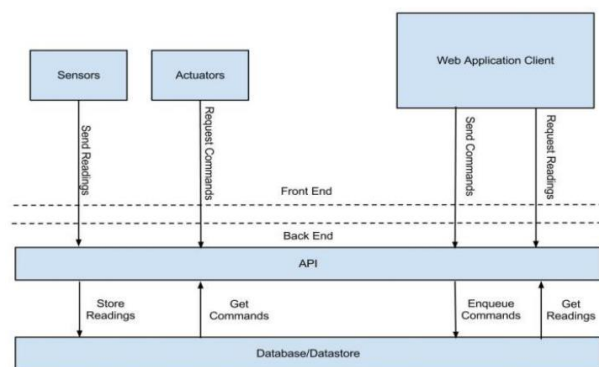


Fig. 9. Smart home architecture, based on integration of IoT with Web services

This architecture contains the following major components – microcontroller-enabled sensors; microcontroller-enabled actuators; database/data store; Server API layer and web application. The microcontroller-enabled sensors

measure the home conditions. The microcontroller-enabled actuators receive commands transferred by the microcontroller for performing certain actions. The commands are issued based on the interaction between the microcontroller and Cloud services. The database/data store component stores data from microcontroller-enabled sensors and Cloud services for data analysis and visualization, and serves as command queue being sent to actuators as well. The Server API layer, which is between front-end and back-end, facilitates processing the data received from the sensors and storing the data in database. It also receives commands from the web application client to control the actuators and stores the commands in database. The actuators make requests to consume the commands in the database through the server. The Web application, which is serving as Cloud services, enables to measure and visualize sensor data, and control devices using a mobile device.

Another smart home automation work is the voice assistant Amazon Echo (Amazon Alexa). This assistant is made by Amazon and suggested in vary variants – echo dot, echo, echo plus, echo spot and echo show. Echo device is a smart speaker, which can play music and execute voice commands given by the user. This smart speaker provides features like hands free calling and messaging, voice recognition and control of appliances. The implementation of the features is based on skills, which Echo looks for them at the Amazon Alexa Service Platform as shown in Fig. 10.

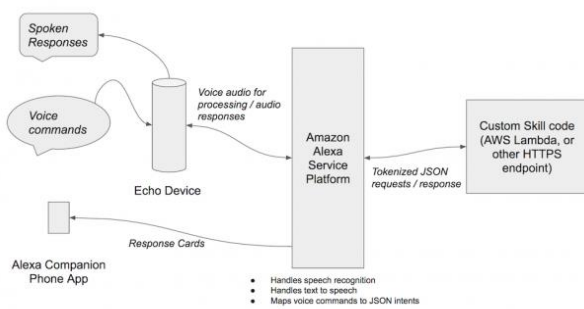


Fig. 10. Amazon Echo skill architecture

As shown in Fig. 10 [10] the user speaks to Echo, using trigger words so that Echo knows that it is being addressed, and identifies the skill that the user wishes to interact with. The Echo device sends a request to the Alexa Service Platform, which handles speech recognition, turning the user's speech into tokens identifying the "intent" and any associated contextual parameters. The intent and parameters for the user's request are then sent as a JSON encoded text document to the server side skill implementation for processing. The server side skill

receives the JSON via a HTTP request. The skill code parses the JSON, reading the intent and context, and then performs suitable processing to retrieve data appropriate to those. A response JSON document is then sent back to the Alexa Service Platform containing both the text that Alexa should speak to the user and markup containing text and an optional image URL for a "card" that appears in the Alexa companion smartphone app. The Alexa Service Platform receives the response, and uses text to speech to speak the response to the user whilst also pushing a card to the Alexa companion app.

4. The ideal smart home

Which are the requirements for the ideal smart home? That is the basic question, which needs to be answered. The smart home has to respond and satisfy the needs of very different groups of domestics. For this it is necessary to know the requirements and needs of the household, interested in having a Smart Home.

A survey conducted in 2012 in Poland explored the behavior of households to predict the usage of the most frequently used home appliances. Krzysztof Gajowniczek and Tomasz Zabkowski explored the usage of electricity by using smart meters in households [11]. The households are situated in a flat of about 140 m² floor area and were equipped with various home appliances including a washing machine, refrigerator, dishwasher, iron, electric oven, two TV sets, audio set, coffee maker, desk lamps, computer, and a couple of light bulbs. The data were gathered during 60 days, starting from 29 August until 27 October 2012. For the analysis they extracted 44 days for which they gathered a set of user behavioral information such as devices' operational characteristics at the household. With their analysis they are found the probabilities for activation of each home appliance at the whole day and they are grouping activation of each appliance at the whole day and whole week. Such surveys give basic information for further developments so they could be based on the consumer requirements.

The smart home proposals, which are mentioned in section 3 could be compared by following parameters – latency, scalability, hardware requirements, complexity, bandwidth size-requirements and their dependencies. This comparison is represented in table 1.

Table 1. Comparison of smart home proposals

Smart Home Proposal	Latency	Scalability	Hardware requirements	Complexity	Bandwidth size requirements	Dependencies
IoT and Fog based monitoring u-healthcare system	low	high	high	high	middle	None
Clome	high	low	high	high	high	depends of cloud supported appliances
IoT Access Point	low	high	low	low	low	None
IoT Energy Management System for dynamic home network based on nomadic agent	low	high	low	low	low	depends of location of the nomadic agent
HAaaS	high	high	high	high	high	None
Amazon Alexa	high	high	high	high	high	depends of cloud supported appliances
Smart home based on integration of IoT with Web Services	middle	high	low	low	low	None

As represented in table 1 the proposals are evaluated in scale of following grades: low, middle and high. In the first proposal, which is based on Fog computing, the hardware requirements are set to high because of the requirements of computer hardware to store local information need for the instantly decisions and multiple gateways about creation of the Fog network (Fog tier). The complexity of the structure of this proposal is high because it uses three paradigms simultaneously – Cloud, Fog and IoT. These paradigms require provision of variety of networks, network protocols and gateways for their integration. In this proposal the usage of these paradigms makes the latency to the low grade, which is a very good advantage. Decisions that need to be made immediately are made by the Fog, but the decisions that take a long time to do so are made by the Cloud tier.

The second proposal “Clome” is evaluated to high bandwidth requirements, high hardware requirements, low scalability, high complexity and high latency. The complexity is high because the usage of different hardware platforms for each appliance and for the need of high bandwidth OpenFlow gateway switch, which connects the home automation to the Cloud. The scalability is low because of the dependency by the application suite provided by the Cloud. From one side there are variety of home appliances, which does not provide Cloud integration and from other side is the requirement of each appliance to be compatible with the Cloud. The hardware requirements are high because the usage of Microsoft Xbox Kinect and OpenFlow switch. The latency is high because the system must reference to the Cloud platform for each decision.

The IoT Access point is evaluated to low complexity because it uses only three types of networks – ZigBee, WiFi and Ethernet. The scalability is set to high because both interfaces ZigBee and WiFi are highly scalable and most of the appliances provide WiFi and ZigBee. In this proposal there is an achievement of low latency with low hardware and low bandwidth requirements.

The IoT Energy Management system for dynamic home network with the nomadic agent is characterized with low latency, low complexity, high scalability and low bandwidth and hardware requirements. In this proposal there are the same achievement as in IoT Access Point – low latency based on low complexity with low hardware and bandwidth requirements. The only disadvantage is the dependency of the location with the nomadic agent. This is a disadvantage because the following reason - only the nomadic agent could have connection with the central server. The consequence of this generates the dependency of location of the nomadic agent. When the nomadic agent is out of the room or house the user cannot control the appliances remotely.

In both proposals HAaaS and Amazon Alexa all parameters are evaluated to high grade. They are open platforms, which are based on IoT and Cloud computing paradigms. Their latency is high because both systems are controlled by the Cloud. Especially Amazon Alexa sends the audio data obtained by the user speech to the Cloud. Then Cloud process tokenized JSON requests to find the skill, after that executes the skill with parameters given by the user’s speech and responses with the result to the Amazon Echo sound device. All transactions to the Cloud machine for each command by home appliance or by user speech leads the latency to high grade. Therefore it requires not only high bandwidth requirements but also high hardware requirements need for the processing of the data.

The last proposal is characterized with low complexity – it uses ZigBee network for connection with all home appliances and uses Ethernet for interaction with remote server, which is used to collect sensor data and controlling of the appliances. This network is highly scalable – it provides variety of connected devices at the same time. The latency is set to middle grade because there are needed requests to the server’s database for loading the commands, which must be executed by the sensor or actuator. If the control commands, which must be executed by each actuator or sensor, are loaded by them only one time (during initial configuration of the home automation) by usage of local permanent memory with multiple cycles (flash ROM or EEPROM), then the latency will get low. The bandwidth and hardware requirements are also evaluated to low grade, because the parts of this home automation proposal does not send big data, which could load the communication environment and hardware platforms.

5. Conclusion

In the modern world of rapid technological progress, modern technologies have entered not only in industry and industrial production, but also in the everyday life of all households. This striving to facilitate and improve everyday life and lifestyle predetermine the growing interest and research into smart homes. The variety of networks and technologies available so far are a good prerequisite for building a smart home - from M2M and IoT to Cloud and Fog Computing.

REFERENCES

1. Gungor V., Sahin D. et al., (2012). Smart Grid and smart homes-key players and pilot projects, *IEEE Industrial Electronics Magazine*, December, p.18-34
2. Hafeez A., Kandil N., et al., (2014). Smart Home Area Networks Protocols within Smart Grid Context. *Engineering and Technology Publishing*, Vol. 9, № 9, p.665-671
3. Nandyala, C., Kim, H. (2016). From Cloud to Fog and IoT-Based Real-Time U-Healthcare Monitoring for Smart Homes and Hospitals, *International Journal of Smart Home*, Vol. 10, No. 2, pp. 187-196
4. Ullah Ihsan (2015). A survey of home energy management system for residential customers. *IEEE 29th International Conference on Advanced Information Networking and Applications*.
5. Kim J., Byun J., Jeong D. et al. (2015). An IoT – Based Home Energy Management System over Dynamic Home Area Networks, *International Journal of Distributed Sensor Networks*, Hindawi Publishing Corporation, Volume 2015, 15 pages
6. Nabi Z., Alvi A., (2014). Clome: The practical implications of Cloud-based smart home, Cornell University, arXiv:1405.0047 [cs.CY]
7. Maiti, A. (2012). Home Automation as a Service, *International Journal of Computer Networks and Wireless Communications*, Volume 2, № 3, p. 421-427
8. Chang C., Kuo C., Chen J et al. (2015). Design and Implementation of an IoT Access Point for Smart Home, *Applied Sciences*, 5, p. 1882-1903
9. Soliman, M. et al. (2013). Smart Home: Integrating Internet of Things with Web Services and Cloud Computing. *IEEE International Conference on Cloud Computing Technology and Science*. pp. 317-320
10. Modus. (2016). Build an Alexa Skill with Python and AWS Lambda, <https://moduscreate.com/blog/build-an-alexa-skill-with-python-and-aws-lambda/>
11. Gajowniczek K., Zabkowski T. (2015). Data Mining Techniques for Detecting Household Characteristics Based on Smart Meter Data, *Energies*, Volume 8, pp. 7407-7427

Contacts:

Georgi Pazhev
 Technical University of Sofia – Branch
 Plovdiv
 Plovdiv, bul Sankt Peterburg 63
 Phone: +359897454436
 E-mail: georgpajev@gmail.com

NATURAL BARRIER WALL DESIGN FOR RADIATION AND NOISE PROBLEM ON THE ECOLOGICAL AIRPORTS

HATICE POLAT¹, KUBRA CELIK², HALUK EREN^{3*}

Abstract: *The airports located nearby to cities are affecting human life negatively due to emitted radiation and noise. Ecological airports are aiming to increase the life quality. The ecological barriers built around the airport could decrease the emitted radiation and noise. With this study, the airports will be designed more environmentally friendly places for humans. This study will investigate the ecological airports' factors in terms of plant cover, the geometric shape of the plants and density of the cover that can be placed around the airports. Thus, the installation cost will be calculated with an optimization function with different parameters. With this function, a proposal has been developed to ensure that such an ecological discipline is compatible with time, labor, monetary cost, and climate or terrestrial conditions. The structure of airports also affects climatic events. This will also include factors such as greenhouse effect reduction, fresh air supply, and oxygen production.*

Key words: *ecological airports, airport design, natural acoustic barriers*

1. Introduction

The airports emits radiation [1][2] and noise to [3] nearby urban areas. Airplanes create noise, especially on departure and landings [4]. Emitted radiation created from both airport station and devices. These two factors are quite harmful for human health and ecological balance. Noise affects people's hearing health and perception in the negative direction, disrupting the physiological and psychological balance [5][6]. Radiation, on the other hand, affects the human body and cell genetics [7].

90% of aviation activities are done in ground place, and the remaining 10% are done in the air. The people living in the near-urban sides are being influenced from airports. Accordingly, social and environmental responsibility is important for airports. Many airports have begun to develop a sustainable environmental design for future [8].

Sustainability policies for airports affect the quality of life of future generations. Airports should also be a model and leader for the environmental transformation of the zone in which they are involved. The environmental impacts resulting from the construction and operation of airports can be solved by the planned international "biopolitics" to be developed in the context of "environmental bioethics" [9].

Urban settlements should not be allowed by governments, especially near airports [10]. The structure of settlements that will evolve over time should be considered. Noise barrier zone areas could be established according to the intensity of the noise in the residential areas near the airport [11] [12] [13].

When sound barrier walls are provided by artificial materials [14], the cost of these materials is high. In addition, because of their psychological and aesthetic effects, plant materials are preferred on noise barriers [15][16]. The General Directorate of Forestry of Turkey uses plants to prevent the noise and exhaust gases on highways. Trees are selected according to their height. For example, the bush type plants are grown in the first line and then long trees are planted for noise barrier walls [17].

The noise barrier should be planted from the noise side, beginning with the bush type trees and short trees, and at the inside with the trees with thick leaves and conifer trees. The dense plant cover is very important in reducing noise [18]. Planted areas and walls reduce also radiation, balance heat and cold and reduce the noise level [19][20]. For airport design; the barrier walls should be considered.

1.1. Problem Statements

Factors affecting the life quality of people living in the surrounding urban areas of airports need to be reduced. These negative factors are electromagnetic radiation emission and noise problem.

1.2. Proposed Approaches

In this study, it is aimed to produce a solution that will prevent the noise problem and radiation around the airports and protect the urban areas from harmful effects. In this work, an optimization study of barrier zone consist of plants is done between urban area and airports.

1.3. Contributions

The airports have become indispensable parts of our lives. For this reason, airports might be designed more environmentally. As an answer to this question; the airports should be located far away from the urban areas. And a forestry area should be planted between airport zone and urban areas. In Fig. 1 sketched forestry area nearby the airport is given. Forestry area is designed as acoustic and radiation barriers for urban areas.



Fig. 1: Forestry area between airport zone and urban areas

The area between the airport and the urban area should not be empty, a forestry area could be grown between them. Green plants could be grown in the airport. The livestock should not be done near airports. The farm animals could be affected by harmful emissions from airports like radiation and noise. Trees absorb the negative factors that affect the living organisms. Thus, biodiversity is ensured. There may be greenery park zones for people in this kind of barrier areas near airports. People can be encouraged to voluntarily take care of the forest area. Thus, peoples could spend time in nature. Lastly, an estimate of cost must be made.

For instance, Fig. 2 shows a similar example of Frankfurt Airport. There is a green area between Frankfurt Airport and the nearest urban area.



Fig. 2: Satellite image of Frankfurt Airport

1.4. Outline (Layout)

In the second session, the system theory of barrier zone design is given for selected airport. In this part, the selected barrier plant cover properties are discussed and sketched. In the subsections, the installation and operation costs are estimated with the developed equation.

2. System Theory (Preliminaries)

The ecological airport design needs system diagram. In Fig 3. the generated system diagram in this study is given.

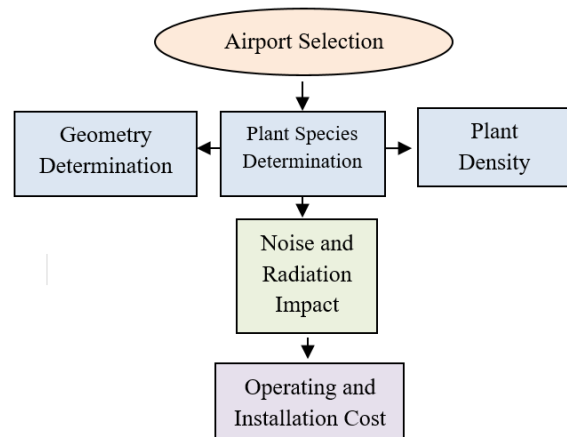


Fig. 3: System diagram

The barrier plant cover could be located between airport and urban area zone. When determining the plant variety, the following must be considered in Table 1. After planting the proper barrier plant cover, the different type trees could be grown on the urban side of the forestry area.

Table 1: Selected barrier plant cover properties

Selected plants should be large and hard-leafed
Trees with a common leaf texture should be preferred
Trees that do not fall in the winter should be selected
The leaf texture should have lied to the ground
The trees should be arrayed densely

Bushes should be grown on the airport side of the barrier plant cover. Trees such as cypress, hornbeam, oak, beech, spruce, juniper can be planted behind the bush line. On the system diagram of the barrier plant cover, the trees could be shown as symbols. In this study, the triangle symbol indicates trees such as cypress, pine, poplar, beech, spruce. The round symbol symbolizes trees such as hornbeam, oak, plane, ash. Trees that are not planted for barrier walls can be expressed as rectangles. Fig. 4 shows the plants that are schematized in this study.

In the design of the plant density, 10 plant species were examined. The distance between the plants must be minimum for the barrier wall. The trees should be positioned closest to each other considering the root structure of the trees. The spacing between the leaves of both trees may be 60 cm. Thus, precautions against a possible forest fire have been taken. The plants must be as deep and at least 5 m high as possible in order to be effective in noise suppression [13].

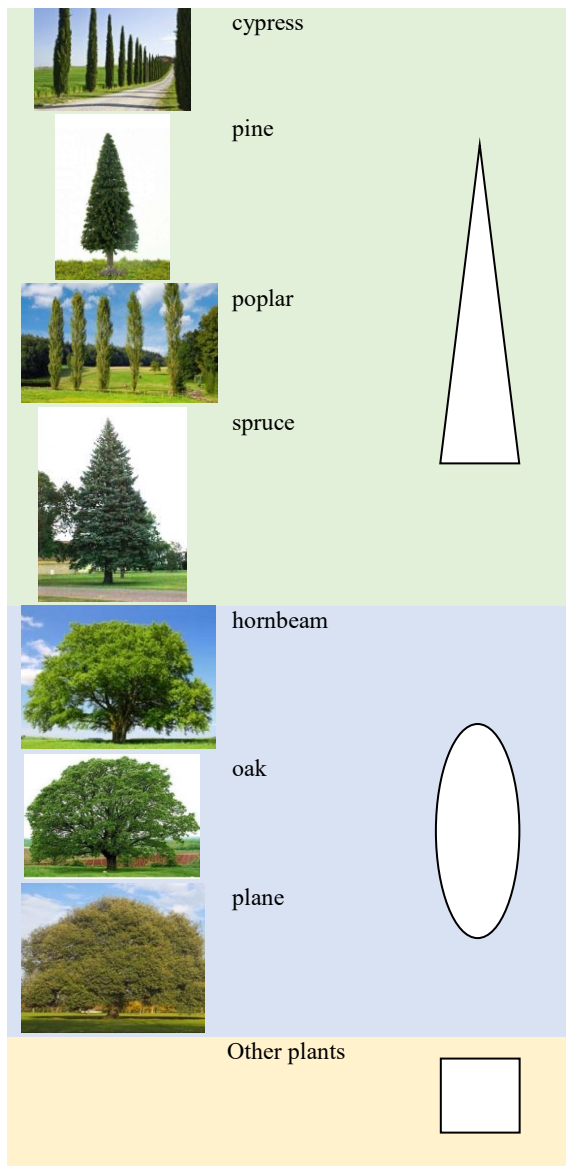


Fig. 4: The plants schematized in this study

Fig. 5 shows the schematic of the designed barrier plant cover in this study for radiation and acoustic barriers



Fig. 5: Schematic of the designed barrier plant cover

According to the recent studies, the barrier wall must be at least 16-20 m. The altitude should be at least 14 m. The direction, place and severity of the voice that the airport has emitted around should be

measured. Thus, it is determined in which area the acoustic waves are concentrated. The barrier wall should be designed thicker in these intense areas.

2.1. Installation and operation cost

The cost of installation of the designed forest area can be calculated. The cost formula (X) can be calculated as follows;

$$X = nLC + mPC + MC \tag{1}$$

where LC indicates labor cost, n shows the number of workman, PC defines plant cost, m shows the tree number, MC remarks material cost.

3. Discussions

The calculated results about ecological airport installation costs are given in Table 2.

Table 2: Total installation cost

X	46.000 Euro, Total Installation cost
n	100 workman
LC	400 euro
m	500 tree
PC	10 euro for each tree
MC	1000 euro total

The total installation cost X is calculated as;
 $X=40.000+5.000+1.000=46.000$ Euro

The operation cost of designed forest area per month can be calculated. For daily operations 10 workmen are needed. And artesian wells can be installed for irrigation. In Table 3 total operation cost is given.

Table 3: Total operation cost

Y	4500 Euro, Total operation cost per month
n	10 workman
LC	400 euro per month
MC for irrigation	500 euro

The total operation cost per month Y could be calculated as;
 $Y:10 \times 400 + 500 = 4500$ Euro

3.1. Change of barrier plant cover over time

After planting the barrier cover, the trees will be change during time. The height and density will be the key parameters in this variation period. In Fig. 6 the estimated plant cover change is given for each century.

The trees in the forest can die after a certain time and thus the density of plant cover can be reduced. However, if the forest continues to be maintained, new trees will be planted instead of dead trees. Thus, continuity of plant cover is ensured.

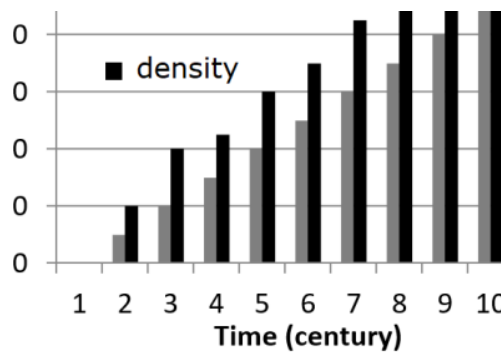


Fig. 6: Barrier plant cover changing depending on time

4. Conclusions

In this study, noise and radiation problems which may occur in residential areas around the airport were sought. Between the urban area and the airport, a barrier zone consisting of a forested area is designed. In the study the trees that can be used in the barrier zone are modeled. Installation and operation costs are calculated in the forest area. The plant cover changing during time is examined. This study could help the governments and airport designers for architecting greener ecological transportation places. In the future works, the planted barrier wall theory could be developed with the estimations of radiation or light effects and new generation material solutions.

REFERENCES

- [1] Bergot T., Escobar J., and Masson V. (2015). Effect of small-scale surface heterogeneities and buildings on radiation fog: Large-eddy simulation study at Paris-Charles de Gaulle airport. *Q. J. R. Meteorol. Soc.*, vol. 141, no. 686, pp. 285–298.
- [2] Stolaki S., Pytharoulis I., and Karacostas T. (2012). A Study of Fog Characteristics Using a Coupled WRF–COBEL Model Over Thessaloniki Airport, Greece. *Pure Appl. Geophys.*, vol. 169, no. 5–6, pp. 961–981.
- [3] Brueckner J. (2008). Airport noise regulation, airline service quality, and social welfare. *Transportation Research Part B: Methodological*, Volume 42, Issue 1, Pages 19-37.
- [4] Lu, C., & Morrell, P. (2006). Determination and applications of environmental costs at different sized airports—aircraft noise and engine emissions. *Transportation*, 33(1), 45-61.
- [5] Stansfeld S. and Matheson M. (2003). Noise pollution: non-auditory effects on health. *Br. Med. Bull.*, vol. 68, no. 1, pp. 243–257.
- [6] Goines, L., & Hagler, L. (2007). Noise pollution: a modern plague. *Southern Medical Journal-Birmingham Alabama*, 100(3), 287.
- [7] Zamanian, A., & Hardiman, C. (2005). Electromagnetic radiation and human health: A review of sources and effects. *High Frequency Electronics*, 4(3), 16-26.
- [8] Ruble, V. M. (2011). Status Report on Sustainable Airports in the United States: Case Study of Chicago O'Hare International Airport. Illinois: Department of Political Science in the Graduate School Southern Illinois University Carbondale.
- [9] Oto N., (2010). *Havaalanlarının Çevresel Etkileri, Çevre Dostu Havaalanı Planlama, Uygulama ve İşletme Esasları; Esenboğa Havalimanı Örneği*. Doctorate Dissertation Ankara University, Social Science Institute.
- [10] Carley M. and Christie I. (2017). Managing Sustainable Development. *Routledge*, 75-78.
- [11] Kang, J. (2006). *Urban sound environment*. CRC Press.
- [12] Zhang, J., Wu, X., Zheng, L., & He, L. (2017). On the specification of general sound environmental zones and planning. In *Inter-Noise and Noise-Con Congress and Conference Proceedings* (Vol. 255, No. 4, pp. 3423-3434). Institute of Noise Control Engineering.
- [13] Zaporozhets, O., Tokarev, V., & Attenborough, K. (2011). *Aircraft Noise: Assessment, prediction and control*. CRC Press.
- [14] Kotzen, B., & English, C. (2014). *Environmental noise barriers: a guide to their acoustic and visual design*.

- CRC Press.
- [15] Kotzen, B. (2002). Plants and environmental noise barriers. In *International Conference on Urban Horticulture 643*(pp. 265-275).
- [16] Peleszezak, P. (2007). *U.S. Patent Application No. 11/630,210*.
- [17] "ogm." [Online]. Available: <https://www.ogm.gov.tr/ekutuphane/Sayfalar/Projeler.aspx>. [Accessed: 04-Apr-2018].
- [18] Givoni B. (1991). Impact of planted areas on urban environmental quality: A review. *Atmos. Environ. Part B. Urban Atmos.*, vol. 25, no. 3, pp. 289–299.
- [19] Ottel  M., Perini K., Fraaij A., Haas E., and Raiteri R. (2011). Comparative life cycle analysis for green faades and living wall systems. *Energy Build.*, vol. 43, no. 12, pp. 3419–3429.
- [20] Saadatian O. *et al.* (2013). A review of energy aspects of green roofs. *Renew. Sustain. Energy Rev.*, vol. 23, pp. 155–168.

Contacts:

Organization : Firat University
Department :School of Aviation, Dept.
of Air Traffic
Management,

Address: Sivil Havacılık
Yüksekokulu Yazıkonak
Beldesi Konak Mahallesi
Ali Şengez Bulvarı 23180
Elazığ – TURKEY

Phone : +90 424 237 00 00

E-mail: he.edu.tr@gmail.com

POSSIBILITIES FOR A COMPARATIVE STUDY OF THE VIBRATIONS IN A COMPLEX PENDULUM JAW CRUSHER

ZHIVKO ILIEV, GEORGI DINEV

Abstract: *An algorithm is developed for the functional analysis of the movement of complex pendulum jaw crushers. This makes it possible to gather comprehensive information about the occurrence of defects in the course of the exploitation. The efficiency of applying the FEMA method in improving the quality of constituent details with pronounced defects is given solid grounds through this analysis. A particular case is discussed: analysis of the state through monitoring and diagnostics of antifriction bearings.*

Key words: *complex pendulum jaw crusher, vibrations, bearings*

1. Major considerations

To guarantee the quality of the products, the Failure Mode and Effects Analysis (FEMA) method is increasingly being applied as a means of reducing errors during systems analysis that is aimed at detecting the causes of these errors. The risk for defects in the mechanisms to occur is assessed through FEMA. The assessment is based on a previously performed Pareto-analysis for establishing the factors that dominate the efficiency of the constituent details [1, 4, 10].

2. Functional analysis of the mechanism

In order to carry out a functional analysis, a structural classification is offered of the separate details of a complex pendulum jaw crusher (see Table 1). The object of this study is a crusher operating within the “Balsha” Mining Plant. It is given in figures 1, 2, and 3 [2].

2.1. Analysis of the defects with the bearings through the FEMA method

The FEMA method can be employed in the analysis of possible errors that may arise during the production and exploitation of heavily loaded details of the crusher head. For this purpose, the state of the mechanism needs to be known through the Pareto method with the dominating factors that best reflect the occasional failures in the course of its exploitation [9, 11, 12]. Also, to diagnose bearing units and other elements in terms of

vibrations, a software package can be used for the processing of experimental data related to their normal state of efficiency. These methods can be employed in decision making when it comes to replacing machine elements. The conclusions are drawn by experts. The applicability of this approach is illustrated by the following example of the replacement of an antifriction bearing.

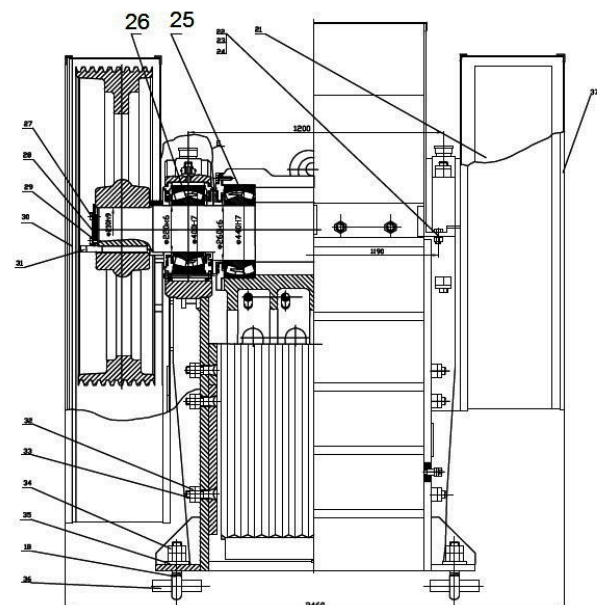


Fig.1 Complex pendulum jaw crusher

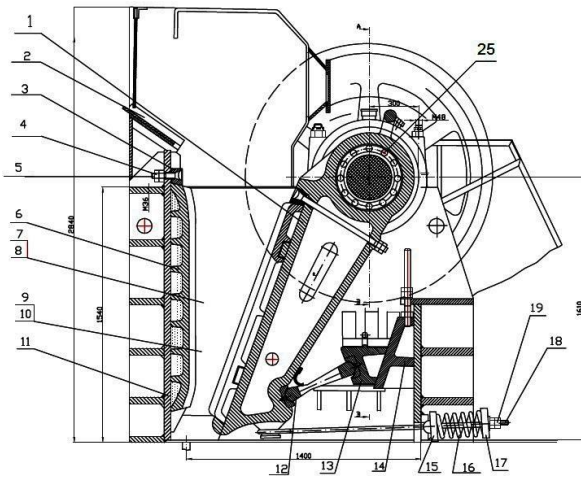


Fig.2 Cross section of a complex pendulum jaw crusher

Table 1. Structural classification of the details

№	Details	Major functions
1	Eccentric shaft	1.1. Transmission of motion to the jaw 1.2. Transmission of axial and radial forces to the bearings
2	Pendulum jaw eccentric shaft	2. Transmission of radial and axial forces to the jaw
3	Pendulum jaw	3. Regulation of the fraction size
4	Hull	4.1. Supporting the pendulum jaw shaft 4.2. Shaft bearing in the body
5	Bearing joints	5.1. Shaft bearing with the jaw 5.2. Transmission of forces to the body
6	Bearing caps	6.1. Maintaining the bearing in an axial direction 6.2. Protection against dirt

7	Cap bolts	7.1. Transmission of axial forces to the body
8	Crusher case	8.1. Holds the constituent details of the crusher. 8.2. Takes the load from the shaft and the pendulum jaw. 8.3. Protects against external influences in the course of the operation of the crusher.

2.2. Means of measurement.

The hardware part comprises the Arduino Mega development part and MMA7361 sensors for the measurement of acceleration along three directions. The parameters of the development system that are relevant to the device developed are:

- an ATmega 2560 microcontroller with a clock frequency of 16 MHz;
- a 10-bit analog-to-digital converter;
- 16 analog inputs and 54 input/output digital pins;
- the maximum data exchange rate is 115200 bps;
- SRAM - 8kB.

The MMA7361 sensor for measuring acceleration offers the following opportunities:

- Selectable range (+/- 1.5 or 6g);
- Maximum sensitivity 800mV/g;
- Analog output signal;

Power supply for Arduino Mega is provided via the USB port of the portable computer and the voltage to the MMA7361 sensors is provided by the development system itself.



Fig. 3 General view of the crusher

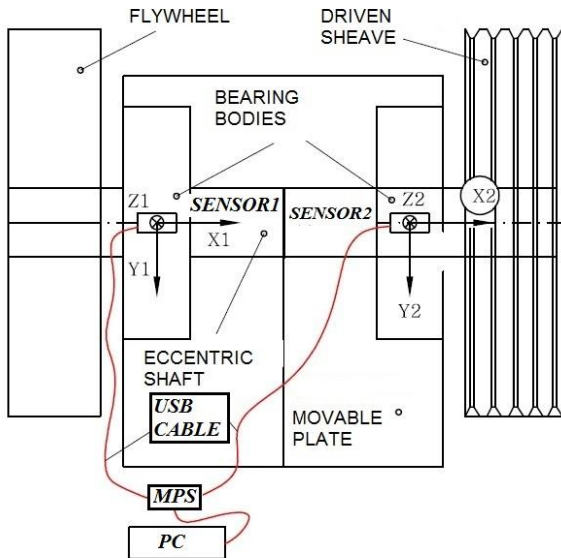


Fig. 4 Diagram of the sensor location on the crusher

The software allows unlimited recording of each of the 12 directions. The minimum sampling interval is 4.1615 milliseconds, the limiting condition being the number of measuring directions and the interface abilities to connect to the portable computer (HP 6730 s). The interval can be programmed to extend or be shortened by reducing the measuring directions. Numeric integration of the records has been made to obtain the dominant spectral density frequencies.

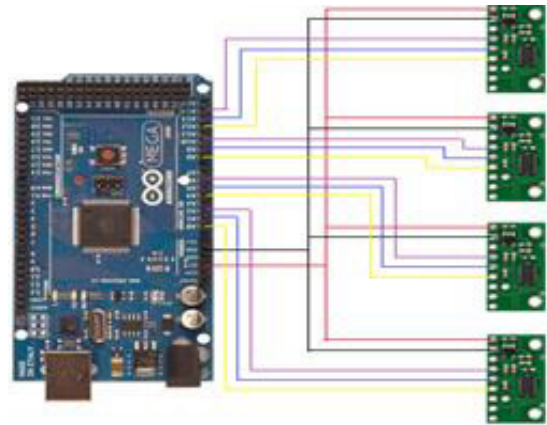


Fig. 5 Diagram of the measuring hardware

Figure 4 shows a diagram of the location of the measuring sensors, in positions 1 and 2 respectively. The microprocessor system (MPS) connects to a portable computer (PC) via a USB cable. Each sensor has its own measurement coordinate system. Thus, sensor 1 has the frame of axes X_1 , Y_1 , and Z_1 . Respectively, sensor 2 has the frame of axes X_2 , Y_2 , and Z_2 . The principle of the electrical circuit for the measuring hardware is given in figure 5.

3. Result analysis

The dominant frequencies obtained are part of the function of the spectral density of the vibrations during the operation of the crusher [12].

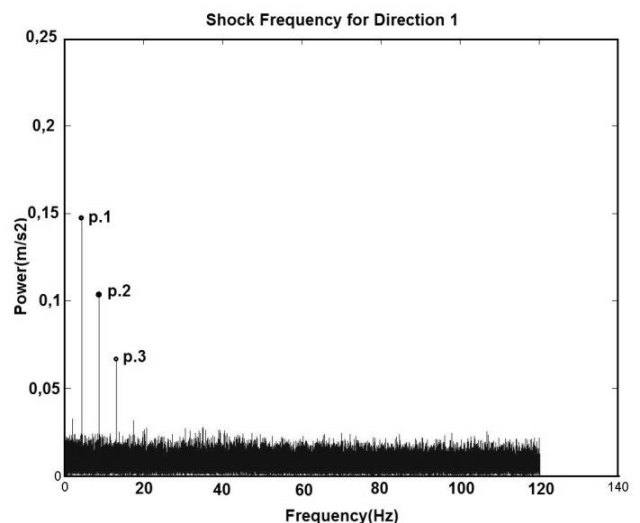


Fig. 6 Dominant frequencies for sensor 1 and axis X_1

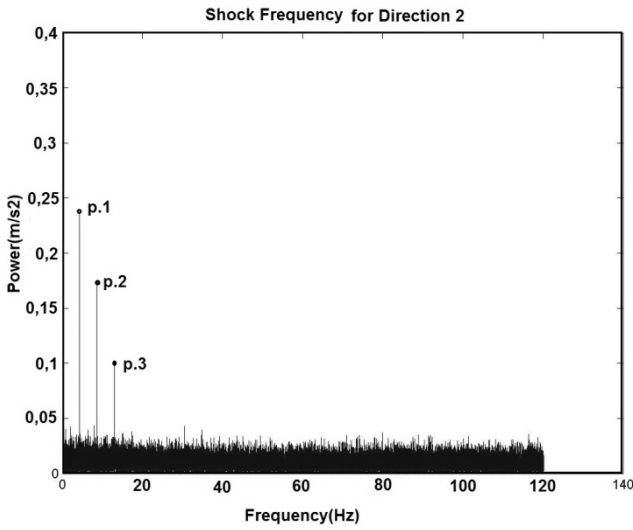


Fig. 7 Dominant frequencies for sensor 1 and axis Y_1

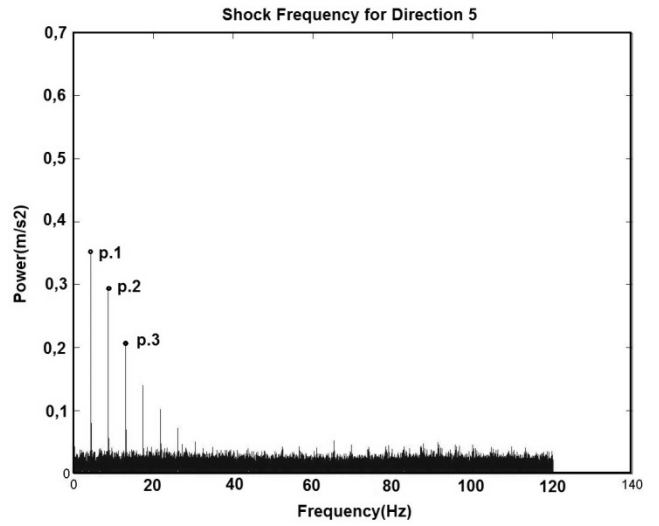


Fig. 10 Dominant frequencies for sensor 2 and axis Y_2

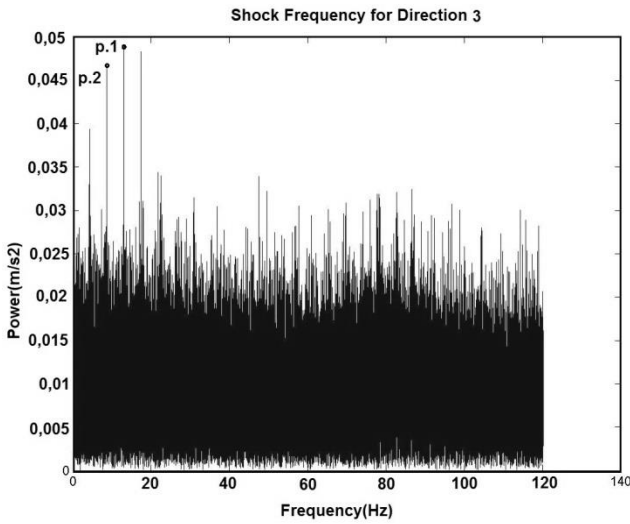


Fig. 8 Dominant frequencies for sensor 1 and axis Z_1

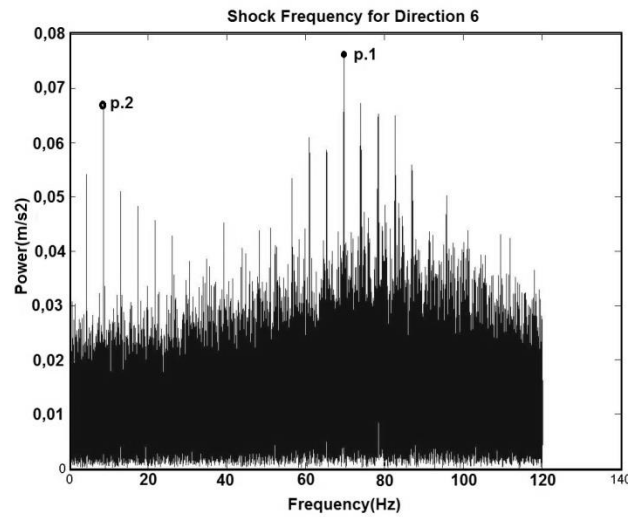


Fig. 11 Dominant frequencies for sensor 2 and axis Z_2

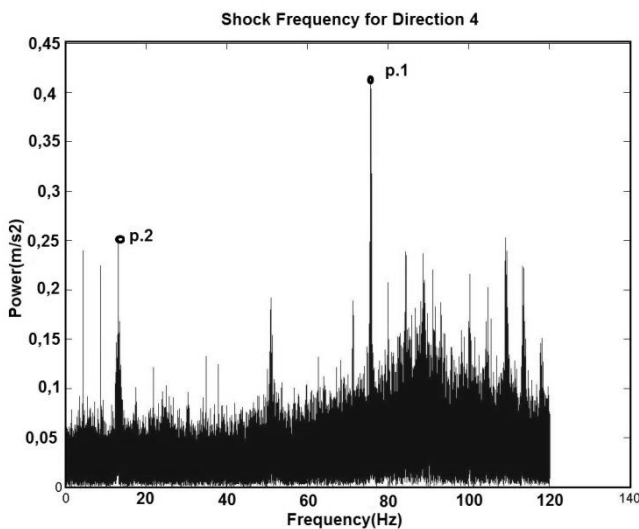


Fig. 9 Dominant frequencies for sensor 2 and axis X_2

If a particular element, like a bearing, a shaft, or an axle, is damaged, shock loads occur in the course of their operation. These loads depend on the state of the elements, as well as on the way they are transmitted through the other parts of the crusher. For this reason, the frequency characteristics are examined [6]. The frequencies generated by various defects are clearly distinguished and provide information during machine vibro-diagnostics.

Figure 6 gives the dominant frequencies for sensor 1 and axis X_1 . Peaks can be seen at three characteristic points on the graph, namely 0.15 m/s^2 at point 1; 0.1 m/s^2 at point 2; and 0.07 m/s^2 at point 3.

The dominant frequencies for sensor 2 and axis X_2 are given in figure 9. Particular points are marked that correspond to the frequency peaks. They are point 1 and point 2 with the respective

values of 0.4 m/s^2 and 0.25 m/s^2 . With the two sensors and the respective axes of X_1 and X_2 , the maximum values are within the frequency range of up to 20Hz, whereby their common origin is indicated. As for the second sensor, the peak of 0.4 m/s^2 is registered around 80Hz but is missing with the first sensor and obviously differs in nature.

Figure 7 shows the dominant frequencies for sensor 1 and axis Y_1 . Three peaks are clearly seen and are marked with three points: 0.24 m/s^2 at point 1; 0.17 m/s^2 at point 2; and 0.1 m/s^2 at point 3.

For axis Y_2 and sensor 2, the peaks are also within the frequency range of up to 20Hz. As can be seen from the graph in figure 10, the values for the three specific points are: 0.2 m/s^2 at point 3; 0.3 m/s^2 at point 2; and 0.35 m/s^2 at point 1.

The oscillographs for axes Z_1 and Z_2 from sensors 1 and 2 are given in figures 8 and 11 respectively. Two peaks are found in the first graph (fig. 8): 0.05 m/s^2 at point 1 and slightly over 0.045 m/s^2 at point 2.

Figure 11, too, shows peaks at two characteristic points: 0.078 m/s^2 at point 1 and 0.068 m/s^2 at point 2. Point 2 is within the range of up to 20Hz which coincides with the peaks in figures 5, 6, and 9. Point 1, however, coincides with the peak at point 1 in figure 8. These values are indicative of the common origin of the former (those of up to 20Hz) but also indicate as to another cause of the origin of the peaks at frequencies around 80Hz.

3.1. Study of the frequencies due to bearing defects [2, 5].

Input parameters:

$n = 250 \text{ min}^{-1}$ - revolutions per minute of the eccentric shaft;

d – diameter of the rolling balls, *mm*

For bearing SKF 23 244 ($\varnothing 220 \times \varnothing 400 \times 144$) as shown in fig.1 at position 26:

$D_{\text{external}}=400 \text{ mm}$; $D_{\text{internal}}=220 \text{ mm}$

For bearing 23 152 ($\varnothing 260 \times \varnothing 440 \times 144$) as shown in fig.1 at position 25;

$\beta = 15^\circ$ – contact angle (SKF 23 244, SKF 23 152);

$E = 2 \cdot 10^{11}$, N/m^2 - module of elasticity of steel (SKF 23 244, SKF 23 152);

$\rho = 7300$, kg/dm^3 - specific weight of the rolling balls; $D_{\text{external}}=440 \text{ mm}$;

$D_{\text{internal}}=260 \text{ mm}$

$D_{\text{internal}}=260 \text{ mm}$

$$D = \frac{D_{\text{external}} + D_{\text{internal}}}{2}, \text{ mm} \quad (1)$$

$$d = 0,25 \cdot (D_{\text{external}} - D_{\text{internal}}), \text{ mm} \quad (2)$$

- Rotational frequency of the eccentric shaft [5]:

$$f_{\text{shaft}} = \frac{n}{60}, \text{ Hz} \quad (3)$$

- Rotational frequency of the separator [5]:

$$f_{\text{separator}} = \frac{f_{\text{shaft}}}{2} \left(1 - \frac{d}{D} \cos \beta \right), \text{ Hz} \quad (4)$$

- Rotational frequency of the rolling balls:

$$f_{\text{rolling balls}} = \frac{f_{\text{shaft}}}{2} \cdot \frac{D}{d} \left[1 - \left(\frac{D}{d} \right)^2 \cdot \cos^2 \beta \right], \text{ Hz} \quad (5)$$

- Vibration due to a defect in the form of the rolling balls:

$$f_1 = \left(\frac{D+d}{d} \right) \cdot \left(\frac{D-d}{d} \right) \cdot \frac{n}{30}, \text{ Hz} \quad (6)$$

- Vibration due to a defect in the form of the inside track:

$$f_2 = \left(\frac{D+d}{D} \right) \cdot \frac{n \cdot z}{120}, \text{ Hz} \quad (7)$$

Z - number of rolling balls in a row

$$z = 5 \cdot \frac{D_{\text{external}} + D_{\text{internal}}}{D_{\text{external}} - D_{\text{internal}}} \quad (8)$$

- Vibration due to a defect in the form of the outside track:

$$f_3 = \left(\frac{D-d}{D} \right) \cdot \frac{n \cdot z}{120} \quad (9)$$

- Resonance frequency of the rolling balls:

$$f_{\text{resonance}} = \frac{0,848}{d} \cdot \frac{E}{2 \cdot \rho} \quad (10)$$

- Vibration due to a separator clearance:

$$f = \frac{1}{2} \cdot \left(1 - \frac{d}{D} \right) \cdot \frac{n}{60} \quad (11)$$

- Optimum pendulum frequency of the jaw

$\alpha = 30^\circ$ - angle of seizure of the material;

$K = 0.8$ – factor taking into account the friction force during unloading of the material

$S=26$, *mm* - Jaw run

- Jaw vibration:

$$z = 0,5 \cdot K \cdot \sqrt{\frac{g \cdot t \cdot g \alpha}{2.5}}, \text{ Hz} \quad (12)$$

4. Conclusion

The results obtained from the analytical study of the frequencies resulting from bearing defects are given in Table 2. The peaks of the dominant frequencies that appear around 80Hz (figures 8 and 11) may be attributed to defects in the outside form of the bearings. The frequencies for these defects have been analytically obtained and are within the range of 85.8-98.8 Hz.

In terms of frequency, the first harmonics that have their peaks around and up to 20 Hz (figures 6, 7, 8, and 10) coincide or are very close to those resulting from the separator rotation: 17.88–18.21 Hz (Table 2). Therefore, coincidences due to defects in the bearing bodies are out of the question.

The remaining frequencies in the analytical study are much greater than those measured in the study which is indicative to the lack of damaged parts and elements.

Table 2. Results of frequencies due to defects

	SKF 23 244	SKF 23 152
D,mm	310	350
d,mm	45	45
f shaft	41.6 Hz	41.6 Hz
f separ.	17.88Hz	18.21 Hz
f rolling	138Hz	157Hz
f1	387Hz	495Hz
f2	115Hz	128Hz
z	7	7
f3	85.8Hz	98.8Hz
freson.	258MHz	258MHz
f	1.78Hz	1.8Hz
z	4.17Hz	4.17Hz

4.1. Conclusions

An algorithm is developed for the functional analysis of the movement of complex pendulum jaw crushers. This makes it possible to gather comprehensive information about the occurrence of defects in the course of the exploitation. The efficiency of applying the FEMA method in improving the quality of constituent details with pronounced defects is given solid grounds through this analysis. In the case of machine element replacement, methods are offered for decision taking and for drawing conclusions based on frequency amplitudes [5, 7, 8, 12].

REFERENCES

1. Bragin V., F. Choban (1997), "Assessment of the risk and consequences of failures of integrated systems and process construction"(in russian), *Market and quality*, Yaroslavl.
2. Iliev Zh., Bogdanov I., Ivanov N., (2015) "Analysis of the vibration state of the eccentric shaft with the bearings of a complex pendulum jaw crusher"- *XVI Balkan mineral processing congress, Belgrad, Serbia*.
3. Barakov A., N. Barakov, A. Azovcev,(2000) "*Monitoring and diagnostics*"(in russian). St. Petersburg.

4. Dinev G., (2004) "Methods for improving and restoring of gear wheels from gearboxes" (in bulgarian), *SOFTTRADE*.
5. Pojidaeva V., (2013), "Vibrodiagnostics of machines and equipment in the mining and the processing industry" (in bulgarian) *ISBN 978-954-353-221-6*, Sofia.
6. Mitrev R., Janosevic D., Marinkovic D., (2017), "Dynamical modelling of hydraulic excavator considered as a multibody system" *Tehnicki vjjsnik 24(supplement 2).doi: 10.17559/TV-20151215150306*.
7. Hadjiiski V., S. Stefanov(2007.), "*Computer engineering analysis of machine elements*"(in bulgarian), COSMOSWORKS, Publishing house of XT,
8. Genadiev G., (2001), „Spending the resource of the machines and the quality of the system maintenance and repair“(in bulgarian), *Publishing virtual center, culture and research*, Sofia.
9. Koriikov C. (1998). "Quality management“(in bulgarian), *The printing base of University of Rousse*, Rousse.
10. Dinev G., (2009), "Bases of the design of the gears"(in bulgarian), *Soft trade*,
11. Marin D., N Predancea, Dan-Michail Marin(2010) „Impact test in working space of milling centers“, *Proceeding of the Fifth International Conference Optimization on of the Robots and Manipulators, OPTIROB, Published by Research Publishing Services*, pp.114-118.
12. Panov V. (2009), „Analysis of the vibration state of a ball mill“(in bulgarian), *XVIII ISTC with an international participation ADP*, Sofia.

Faculty of Mining Electromechanics
University of Mining and Geology"St.Ivan
Rilski", Sofia
Prof. Boyan Kamenov Str.
1700 Sofia
BULGARIA
E-mail: halkopirit@mail.bg

Faculty of Mechanical Engineering
Technical University of Sofia
8, Kliment Ochridski Blvd.
1000 Sofia
BULGARIA
E-mail: gdinev@tu-sofia.bg

CONCEPTUAL IMPACT MODEL OF PROCESS MANAGEMENT ON THE MEAT INDUSTRY ENTERPRISES IN BULGARIA

TONI MIHOVA, VALENTINA NIKOLOVA-ALEXIEVA, MINA ANGELOVA

Abstract: *The survey puts lots of emphasis on an approbation of a conceptual model for evaluation of influencing factors on the BPM activity of Bulgarian meat-processing enterprises, which allows to determine the direction and magnitude of the impact of BPM on their innovation, efficiency and competitiveness.*

Key words: *Conceptual Model, Business Process Management-BPM, meat-processing enterprises, reengineering, Six Sigma Lean*

1. Introduction

The main problem facing meat processing plants is related to the effective management of their processes. Another problem is their insufficient competitiveness [1], [2], [3] due to the fact that most of them are forced to solve their current problems, especially in the years of severe financial crisis and afterwards, rather than concentrating its efforts for its strategic development [2], [8], [9]. This problem is also based on the inadequate ability to document and manage the main and auxiliary processes in the meat industry [2], [3], [9], [11].

The study supports **the thesis** that using the concept of business process management - BPM in the meat industry enterprises is a process approach that makes businesses more mature in their processes. The approach implements the best management Systems, principles, tools and techniques for documenting and managing processes, building a process architecture of the enterprise, and implementing it is a key success factor.

The purpose of the study is to reveal the dependencies and relationships between various factors of the management of business processes in enterprises of meat industry and their impact on the innovation activity, efficiency and competitiveness.

In order to achieve this goal, a survey was carried out in 156 different sized, state-owned and owned enterprises in the meat industry on the territory of Bulgaria in the period between February 2017 and October 2017. 138 of them completed the survey cards in full and provided an adequate response to the BPM tools they used, the innovations that implemented and evaluated the factors. Moreover, there is sufficient information about their financial position in the Commercial Register for the period 2009-2016.

The classification of enterprises by various signs shows their diversity and wide coverage. The following **materials and methods** were used in this study:

Target population - small, medium and large enterprises in the meat processing industry.

Method for collecting empirical data and tools - a personal interview with a paper-based questionnaire developed according to a specialized methodology [3],[7],[8],[11], incorporating the latest concepts in BPM, [1], [6], [7] Reengineering [4], [8], Six Sigma Lean, Redesign, TPM, Kaizen, "20 keys", "5S" and Outsourcing of Business Processes.

Sample type - zoning (stratified) sample.

Sample volume - 156 meat industries responding positively to the invitation to participate in the interview. Data from AMB, BSAF and BIA [5] were used as baseline data for the general population.

Scope of respondents - executive directors, financial managers, project managers, IT managers, external and internal consultants, experts.

Method for empirical data processing and tooling - SPSS 14.0.

2. Methodical conditions

A description of a methodological approach aimed at defining the main steps and interrelationship of the results from an empirical survey of business subjects from the Meat processing sector in Bulgaria is proposed.

The main objective of the approach is to present a system for assessing the impact of an identified complex of factors influencing the efficiency, innovation and competitiveness of meat processing enterprises at sectoral and regional level.

The following stages are included:

Stage 1: Investigation of localization factors. The purpose of this phase of the study is to investigate the views of economic operators, the

status of localization factors in the region, and how many of these factors influence the impact of process management.

The main outcome of the Stage 1 analysis is the formulation of a Primary profile of enterprises applying process management as a basis for the type of applied management approaches (functional, resource, process) and business strategies (to produce a new product, to expand production capacity, improving quality, increasing market share, etc.), as it illustrates, albeit indirectly, the direction of the enterprise's innovation activity.

Stage 2: An exploration of maturity in terms of its processes. This stage complements the analysis and the results of the previous stage in the following aspects: while preserving the objective, the tasks and the general set of the survey introduce additional criteria to the subjects surveyed; modifying the consultation tool to fine-tune the thematic focus. On the basis of the findings from Stage 2, the main result is achieved: formulation of the intermediate profile of the enterprises with introduced process management. The account includes information about process management tools, innovation activity, financial performance, and competitors' business process advantages, and gives feedback on which processes are more mature, what advanced process management approaches they use, what management software they use, what is their innovative potential, what innovation strategies they use, what are the factors influencing efficient process management, which leads to increased innovation activity and increased competitiveness and which of these factors influence the localization and strategic direction of these enterprises.

Stage 3: Survey of Influential Factors on Business Process Management. At this stage of the methodological approach the analysis deepened on defining the extent and strength of the impact of influential factors. The outcome of the Stage 3 analysis is the definition of the Integrated Enterprise Profile, which shows the factors with their impact on the various meat processing companies, according to their process maturity, innovative activity and degree of competitiveness.

The main result of the overall application of the described methodological approach is the possibility to derive a conceptual model describing the impact of the factors influencing the efficiency of the process management in the surveyed enterprises. The conceptual model presents the logical framework of the steps for the proposed solution, namely:

Step 1: The starting point for research and evaluation is the definition of Business Process Management (BPM) as a dependent variable.

Step 2: Exploratory hypotheses are formulated.

Step 3: A correlation analysis follows, which determines the strength and direction of impact of the influencing factors.

Step 4: Using a regression analysis to measure the level (power) depending on the selected (dependent) variable by changes in the independent variables i.e. demonstrated in the presence of dependency.

Step 5: Check the significance of the hypothesis describing dependence and draw conclusions.

Approbation of the model in its part concerning the definition of dependent variables and compiling research hypotheses within this survey is the following:

Dependent variable - BPM is determined by the level of process maturity of enterprises; from the applied process management toolkit; of the innovative solutions used so far and the frequency of future application of the different types of innovation; by the type of technology being used; the development of the market environment and access to skilled human resources.

Independent variables - Possible factors influencing the thematic blocks of the questionnaire.

Formulation of hypotheses - Based on the preliminary statistical analysis, the following research hypotheses are to be examined:

H1: Business Process Management (BPM) is **directly related** to the business start-up period, its capital structure and process maturity level; the reasons for starting the enterprise and the availability and implementation of strategic planning documents in the enterprise; the adequacy of process management tools used; from the introduction of teamwork; from the qualification of managers;

H2: BPM is **proportional, depending on:** the expansion of business entities; positive attitudes towards the development of business units and the region in which they operate; innovation activity; experience in previous innovation projects; interaction between different types of innovative solutions; the possibility of joint projects with a business entity from another industry or competitor; inadequate provision of production and technical staff; of new technological solutions; from creating a competitive advantage;

H3: BPM is **heavily influenced by** access to qualified human resources; appropriate placement locations; developed transport; support programs; the development of the market environment and the growth of business users; the productivity and efficiency of the business entity; of organizational

culture; the introduction of environmental standards; from sustainable development initiatives;

H4: BPM is **less influenced** by the state of production, transport and social infrastructure; the assessment of local government activities.

3. Basic results of the empiric survey

Taking into account the peculiarities of the market and the characteristics of the production process of the enterprises, the survey and analysis were carried out separately for each subsector of the meat processing industry - 'poultry', 'slaughterhouses', 'meat processing' and 'minced meat and meat preparations'.

It is assumed that the sub-sectoral affiliation of enterprises is determined by the predominantly produced output (more than 60%). Through this analysis a complete picture of the status and development of BPM in the meat processing enterprises is being developed. As a result of the analysis of the state of the meat processing industry and the number of respondents who responded adequately and completely to all questions, the authors conclude that, from the point of view of BPM, it is more appropriate to group companies by market segments i.e. the study should continue to regroup the participating enterprises in the market survey. After sifting the enterprises that do not apply the BPM tools and do not complete the questionnaire, the number of enterprises is reduced to 72 and the classification of the enterprises in the sample is as follows (see Figure 1):

The main object of the three-step empirical study is to bring out summary analytical characteristics of the studied meat processing plants, with the focus of the analysis being to define their level of maturity in terms of processes and the influential factors that affect an effective BPM. In order to achieve this goal, it is proposed to produce a corresponding profile (*primary, intermediate, and integrated*) after each study step of the surveyed enterprises to present specific elements of their activity.

The primary profile of the enterprises differentiates them according to the degree of process maturity in the meat processing subsectors, according to their business profile, legal-economic form, their localization in the territory of Bulgaria and their innovation activity in five groups: "2nd level of maturity"; "2-3 level of maturity"; "3rd level of maturity"; "3-4 level of maturity" and "4th level of maturity".

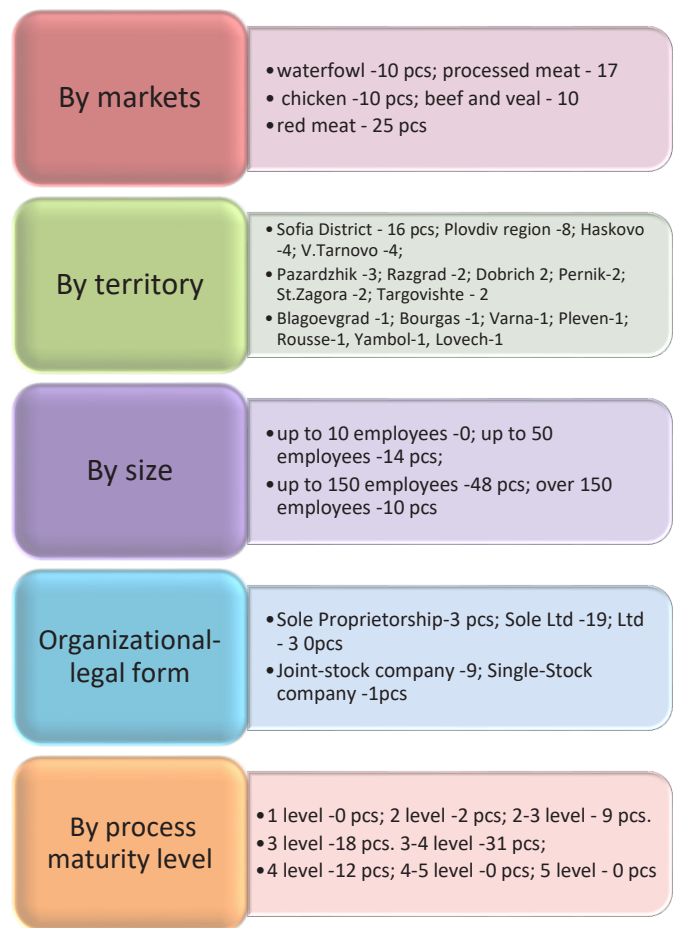


Fig. 1. Signs of classification of surveyed enterprises

The definition of the level of maturity in terms of business processes of the enterprises from the surveyed population is done by means of the answers received to questions № 8 and № 9 of the questionnaire (from the first stage of the survey), respectively on the knowledge and documentation of their processes, as well and the way the processes are organized. On this basis, the following types of enterprises are defined:

• **2 level of maturity**, according to the CMM scale - 3%: which means that they start to differentiate their main processes and are oriented towards applying some instruments of the process coordination but are still subject to the chaotic survival decisions and encounter huge difficulties to remain competitive on the market.

• **2-3 level of maturity**, according to the CMM scale - 12%: they have organized and documented processes and seek to understand how management and support processes have to be matched to support the main processes.

• **3 maturity level**, according to the CMM scale - 25%: they have organized and documented processes, monitoring both their basic and auxiliary processes.

• **3-4 level of maturity**, according to CMM scale - 43%: the processes are documented, for

engineering and for management. It is standardized and integrated into a methodology.

• **4 level of maturity**, according to the CMM scale - 17%: enterprises that understand how their processes work and adapt their strategy and that they have high innovation activity to remain flexible and competitive on the market but have not yet introduced building process architecture.

The definition of the type of innovation activity of the enterprises from the surveyed group is made with the help of the answers given to question № 12 concerning the prepared and implemented company strategies. On this basis, the following types of enterprises were identified:

• **Process innovators:** implement a strategy for creating new processes and expanding production capacity (answer „a“). Their share is 42%.

• **Product Innovators:** Applied a strategy for producing a new product (answer "b"). Their share is 31%.

• **Organizational innovators:** implemented a strategy for radically changing the organization of work and creating new relationships with other enterprises (answer "c"), their share being 2%.

• **Marketing innovators:** Subgroup A - applied a strategy for expanding the market share (12%) and subgroup B - applied strategies for entering new markets (11%).

• **Eco-innovators:** they have put in place procedures to reduce the harmful impact on the environment. Their share is 2%.

It can be **concluded** that the majority of the meat processing enterprises surveyed strive to monitor and document their processes. Regarding the degree of process maturity, the main conclusion that can be made is that there are no meat industry enterprises that have achieved a fifth degree of maturity, on the scale of CMM, namely those who continuously monitor their processes, build a process architecture, appoint teams for each process and continuously innovate to remain adequate to market requirements. Another important conclusion is that there are no enterprises that are at the first level of maturity on the scale of CMM, ie those who do not know and do not follow their processes at all, act operatically and their survival is due to chaos, rather than adequate strategic management, they do not use Business Process Management (BPM) tools to remain competitive on the market.

Another conclusion is that the surveyed business structures exhibit their innovation activity most often through process innovations to expand production capacity; followed by marketing innovations; followed by product innovations such as the introduction of new products for the organization and, last but not least, through

organizational innovation, in the form of introducing new management systems and new relationships with other companies.

Regarding the relationship between BPM tools used, demonstrated innovation strategies and subsector profile, the following **conclusions** are drawn:

• **"2nd level of maturity"** are rather product innovators that are relatively more influenced by the choice of legal form, year of establishment of the company, choice of settlement and choice in the formation and implementation of new product strategies, number of partners.

• **„2-3 level of maturity"** are both product innovators and marketing innovators of the first type, ie they apply strategies to expand the market share that are affected in relative terms by the settlement in which they operate, the location and the legal form.

• **„3rd level of maturity"** are both product innovators and process innovators that are influenced by the year of creation, the legal form and the number of partners.

• **„3-4 level of maturity"** are process and organizational innovators as well as marketing innovators of the second type, ie they enter new markets, which are partly influenced by the number of partners, the staff size and the choice of settlement.

• **"4th level of maturity"** eco-innovation, with process, organizational and marketing innovations tailored to environmental standards, food safety and sustainable development. They are influenced more strongly by the choice of the settlement, the year of creation and the legal form.

On the basis of the analyzed results for the different categories of enterprises, the following summaries can be drawn, regarding the main features of the **Intermediate profile** and its categories:

All five enterprise groups declare that they have experience in implementing BPM projects so far, with positive responses being predominantly high (over 70%). The highest value for this indicator is the "4 level of maturity" (84.6%), followed by "3-4 level of maturity" (81.3%), followed by "3 level of maturity" (75%) and "2-3 level of maturity" (74.2%), with the "2 level of maturity" (65%).

All five enterprise groups declare that they have used the various BAT tools, such as the **4th maturity level** - mainly using Process Reengineering, Balanced Cards - BSc; process frames - (SCOR, ITIL), Six SIGMA for defect reduction and quality improvement, LEAN technology, SIX SIGMA LEAN, quality assurance standards HACCP, ISO 22000, ISO 9001, OHSAS

18001; "20 keys" and KAIZEN have introduced process management systems - BPMS, ERP systems, BI, CRM, workflow; from **3-4 maturity level** - mainly based on process reengineering, redesign of poorly functioning processes, standardization and implementation of quality systems - HACCP, ISO 22000, ISO 9001; Just in time, LEAN technology, 20 keys, ERP systems; CRM, workflow; from **3rd maturity level** - BMP tools used are process reengineering, work organization change systems, SIX SIGMA, 5 "S", continuous improvement, CRM, workflow; from **2-3 maturity level** process documentation, quality standards - HACCP, ISO 9001, improvement of marriage and defects, through SIX SIGMA, improvement of working environment through 5 "S", process automation - workflow; from **2nd maturity level** - process documentation, process monitoring, 5 "S", implementation of quality standards and process automation.

All five groups of companies declare that they have experience in implementing innovative projects so far, with positive responses being predominantly high (73%). With the greatest value of this indicator are eco-innovators (89%), followed by product innovators (80%), followed by process innovators (79%), followed by the organizational innovators (75%), marketing innovators type 2 - pass in new markets (73%) followed by marketing innovators type 1 - expanding market share (72%).

The most preferred for future deployment are the following BPM tools: using process frameworks (SCOR, ITIL) and building process architectures are most preferred for **4 maturity levels** (78%), process reengineering is most preferred for **3-4 maturity level** (82%); SIX SIGMA and LEAN technology are the most preferred for **3 maturity level**; **2-3 level of maturity** prefer Balanced score cards - BSc and 20 "keys"; **2 maturity level** prefer the implementation of quality standards; "Just-in-time" and 5 "S".

The most preferred information systems for future BPM-related projects are: **2 maturity level** - workflow and CRM (72%); for **2-3 level of maturity** - workflow, CRM, ERP (87%); for **3 maturity level** - BP Modeling, CRM, ERP and BI (82%); for **3-4 maturity level** - BP Modeling and BI (79%); for **4 level** - BP Architecture; BI and BPMS (88%).

Most preferred for future innovate are **eco-related products** which have the highest index in the product innovators (81%), followed by marketing innovators type 2 - (73%). After them are preferably **eco-related technologies**, which have the highest index in the product innovators (68%). Next, **process innovations** are equally approved for product and organizational innovators (48%).

Services are relatively more preferable only to marketing innovators type 2 - entering new markets (56%).

A supplier, international company or consultant are the most preferred partners for joint BPM projects. A consultant is primarily for the **2nd maturity level** (56%) and **2-3 maturity level** (67%). International company is a very important partner for the **3rd maturity level** and **3-4 maturity level**, respectively (73% and 68%). Supplier is the most important partner for enterprises in the **4 maturity level** (78%).

Supplier, client, international company and consultant are the most preferred partners for joint innovation. Type 1 product and marketing innovators would first choose a supplier or an international company (46%) and then contact a customer (42-48%) or a consultant (33-42%). In organizational innovators and marketing innovators type 2, the consultant is ranked first in favor (36-47%) and supplier (36.5-32%). The least favored by all partner groups is local authorities, competitors and funding institutions.

Consumer preferences are the most significant factor for the five types of enterprises by maturity, with the highest value for the **4 maturity level** (92%). Human resources in the region are further influenced, with the highest value being for enterprises with a **3 maturity level** (62%). For enterprises with 3-4 level of maturity and 2-3 level of maturity, intellectual property protection (43%) and the legal base (38%) rank third. In contrast, **2 levels of maturity** rank third among the factors of availability of innovation partner (37%) and access to support programs (36-42,9%). The least significant for all categories of enterprises is the development of the local education system, the existence of clusters and the cooperation of branch organizations.

4. Conclusion

The proposed conceptual model of BPM's impacting factors can help businesses in the meat industry in Bulgaria to orientate where they are today and to serve as a navigator on their way to maturity in managing their business processes. Leading European and global companies focus their efforts on moving from the 4th to the 5th level of process maturity. They already have a process architecture that builds on their management system.

In the contemporary business environment, in order to become competitive and manage their processes efficiently, Bulgarian enterprises in the meat industry have to undertake a more synthetic and modern approach to process change that combines the best of process management, redesign, refinement and process automation.

REFERENCES

1. Алексиева В. (2013) „Управление на бизнес процесите в българските предприятия” ISBN-978-954-24-0189-6, Изд. „Интелексперт-96“
2. Mihova T., Nikolova-Alexieva V., (2016) ”Systematic comparison of existing and new approaches for monitoring compliance rules over business processes”, International virtual journal for science, techniques and innovations for the industry- MTM’15/2016, ISSN-1313-0226
3. Nikolova-Alexieva V., (2017) “Business Process Initiatives Applied to Bulgarian Meat Industry Enterprises“, Economics World Journal, ISSN: 2328-7144, ELSEVER
4. Nikolova- Alexieva V., (2012)”BPM model at the enterprise level”, World Conference of Business, Economics and Management-WBEM, Antalya, Turkey,
5. Data of AMB, BSAW, Agrostatistics at the Ministry of Agriculture and Food, data from FoodDrinkEurope, www.estat.com, www.nsi.bg
6. Angelova, M., Nikolova – Alexieva, V. (2017) “Opportunities for raising the entrepreneurial culture – a factor for competitiveness of the Bulgarian economy”; International scientific conference „Applied Modeling in Economics, Finance and Social Sciences“, Hisar, Bulgaria <http://dataconferences.org/page/publications>
7. Angelova, M., Pastarmadjieva, D. (2017) “Challenges and opportunities for flexible crediting of small and medium-sized enterprises in Bulgaria”, International Conference on Engineering, Technologies and Systems TECHSYS 2017, Technical University – Sofia, Plovdiv branch, Plovdiv, Bulgaria
8. Nikolova-Alexieva V., (2013) “Process maturity analysis of the Bulgarian enterprises” LUMEN, Iasi, Romania
9. Angelova M., (2017) Entrepreneurship in Bulgaria – possible or not for young people, International scientific journal "Machines. Technologies. Materials" (Print ISSN 1313-0226, Web ISSN 1314-507X), 217, pp.44
10. Nikolova – Alexieva V., Mihova T., Gigova T, (2014), How Bulgarian busines uses process modelling tools, II Міжнародної науково-практичної конференції „Особливості формування ефективної інноваційно-інвестиційної моделі розвитку підприємства в сучасних умовах господарювання“, м. Житомир, 20–21 Листопада 2014 року, ISBN 978-966-683-426-6 © ЖДТУ
11. Mihova T.; Alexieva – Nikolova V., (2016), “Principles and instruments of lean methodology” International Conference on Engineering, Technologies and Systems TECHSYS 2016, Technical University – Sofia, Plovdiv branch 26 – 28 May

Toni Mihova, Tehcnical University Sofia, branch Plovdiv, 00359893690655, expert9@abv.bg

Valentina Nikolova-Alexieva, University of food technologies, 4000 Plovdiv, 0035988569669 valentina_nikolova@abv.bg

Mina Angelova, University of Plovdiv Paisii Hilendarski, 24 Tzar Asen, 4000, Plovdiv, 00359887461272, mina.marinova@abv.bg

ANALYSIS OF THE MEAN FIELD APPROXIMATION FOR TRAINING THE DEEP BOLTZMANN MACHINE

TODOR TODOROV, GEORGI TSANEV

Abstract: *The developers of the deep machine learning are essentially inspired from the activities of the human brain. The main goal in this area is to get a finite software model of human recognition approach. The deep learning architecture is based on the mean field approximation. Many authors a priori assume that the mean field approximation problem has a solution for all random initial guesses. Some of them just make a few steps applied the fixed point iteration to establish whether neurons in hidden layers are active or not. The present paper describes the area of application of the mean field approximation for training multilayer Boltzmann machine. We have a strong proof on the fact that mean field approximations are not convergent for all random initial guesses. The convergence strongly depends on the norms of the weight matrices. The new results are supported by computer implemented examples.*

Key words: *Deep Boltzmann Machine, mean field approximation, gradient iterative methods.*

1. Introduction

The contemporary voice control of machines is related to deep belief learning. The Deep Boltzmann Machine became practically usable after R. Salakhutdinov and G. E. Hinton [1] had developed the mean field approximation algorithm to establish the state of neurons in hidden layers to be active. The mean field approximation requires the visible neurons to be fixed to the training data when a fixed-point iteration is performed. R. Salakhutdinov and G. E. Hinton extend their variational approach in the following publications [2,3,4]. After the pioneering paper of R. Salakhutdinov and G. E. Hinton [1] a lot of researchers [5,6,7,8] etc. have applied the mean field approximation when investigating Deep Boltzmann Machines. Having in mind that such procedures are executed a lot of times while training a neural network, the real application of a Deep Boltzmann Machine strongly depends on the convergence rate of the fixed point iteration. Most of the authors (see for instance [3] and [6]) have used the fixed-point iteration in order to obtain the probability of the neurons in a hidden layer to be active. Unfortunately this method has too low rate of convergence and very bad behavior when the weight matrices are non-positive definite or ill-conditioned. The two-point step size gradient method was obtained by J. Barzilai and J. Borwein

[9] way back in 1988. The method became popular [10] very fast because of the advantages like: no line search procedures, easy implementation and only slight dependence of the initial guesses etc.

The contributions of the present paper are described as follows. The paper is devoted to a numerical algorithm for solving the mean field approximation problem. The weak formulation of the original problem is transformed to an unconstrained minimization problem. Sufficient conditions for existence and uniqueness are rigorously proved. A strict proof that the weak solution can be used in the process of deep learning is made. The minimization problem is solved by the two-point step size gradient method. We emphasize on the fact that the choice of a steplength for the Barzilai-Borwein method is a crucial point with respect to the rate of convergence. Some steplengths assure convergence but with very slow rate. Inappropriate choice of the steplength causes divergence of the sequence of approximate solutions. The successful steplength for the two-point step size gradient method depends on the objective functional. The initial version of the Barzilai-Borwein method [9] was developed for quadratics in the two-dimensional case. Further, similar results were obtained by M. Raydan [11] in the n -dimensional case. T. D. Todorov [12] found a new steplength for quartics, which gives much

better results than the classical steplength proposed by J. Barzilai and J. Borwein [9]. The number of necessary iterations for satisfying the stop criterion strongly depends of the variable steplength. Two original steplengths for solving the mean field approximation problem are obtained in the present investigation. The new steplengths are compared with the classical one and the fixed-point iteration is compared with the Barzilai-Borwein method. The time for training Deep Boltzmann Machine strongly depends on the initial guesses of the weighted matrices. Lower norms of the weighted matrices assure higher rate of convergence solving the mean field approximation problem and hence decreasing of the necessary time for training.

The rest of the paper is organized as follows. The problem of interest is defined in Section 2. The weak formulation and the associate unconstrained minimization problem is considered in Section 3. Sufficient conditions for existence and uniqueness of the weak solution are found in the same section. An iterative method for solving the associate unconstrained minimization problem is investigated in Section 4. Here a detail proof for convergence of the Barzilai-Borwein method is presented.

2. Setting the problem

We begin with some basic definitions and denotations. The vector space \mathbb{R}^n with the standard basis $\langle e_1, e_2, \dots, e_n \rangle$ is provided with the Euclidean norm $\|\cdot\|$ and the corresponding scalar product. The vectors $\underline{e}_i, i = 1, 2, \dots, n$ are the columns of the corresponding identity matrix and $\hat{e} = \sum_{i=1}^n \underline{e}_i$. The norm in $C^k(\bar{\Omega})$, $k \in \mathbb{N}$ is denoted by $\|\cdot\|_{k,\Omega}$ and the norm in $L^k(\Omega)$ by $\|\cdot\|_{k,\Omega}$. Let F be a k -differentiable map from \mathbb{R}^n to \mathbb{R}^n . Then the norm of the k -th Fréchet derivative $D^k F(x)$ is given by

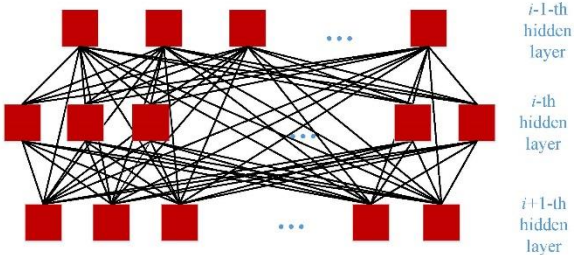


Fig. 1. Neural network where all the hidden units in all layers have the same length.

$$\|D^k F(x)\| = \sup_{\substack{\|\xi_i\| \leq 1 \\ 1 \leq i \leq k}} \|D^k F(x)(\xi_1, \xi_2, \dots, \xi_n)\|.$$

A deep Boltzmann machine with p hidden layers and no more than s hidden units in each layer is an object of interest in the present paper. We assume that all hidden layers are connected excluding horizontal connections within a fixed layer. In the first stage of our investigation we suppose that all hidden units in all layers have the same length m , Figure 1. The latter means that all matrices W_{ij} are square. Usually the weight matrices are rectangular, Figure 2. This case is considered further. We compile the following mean field approximation problem

$$\begin{pmatrix} v_1 \\ v_2 \\ \dots \\ v_r \end{pmatrix} = \text{Sigm} \left(\begin{pmatrix} a_1 \\ a_2 \\ \dots \\ a_r \end{pmatrix} + \begin{pmatrix} O & W_{12} & \dots & W_{1r} \\ W_{21} & O & \dots & W_{2r} \\ \dots & \dots & \dots & \dots \\ W_{r1} & W_{r2} & \dots & O \end{pmatrix} \begin{pmatrix} v_1 \\ v_2 \\ \dots \\ v_r \end{pmatrix} \right), \quad (1)$$

where r is the number of all hidden units in the neural network and $\text{Sigm}(\underline{v})$ is the multivariate sigmoid. Some blocks W_{ij} are zero matrices since there are not connections within layers and usually the number of hidden units is not the same in the different layers. Concatenating the vectors $\underline{v}_i, i = 1, 2, \dots, r$ we obtain an n -dimensional vector $\underline{v}(v_1, v_2, \dots, v_n)$. Then the problem (1) takes the following compact form

$$\begin{cases} \text{Find } \underline{u} \in \mathbb{V} \text{ such that} \\ \underline{u} = \text{Sigm}(\underline{a} + W\underline{u}), \underline{a} \in \mathbb{V}. \end{cases} \quad (2)$$

Here W is $n \times n$ block matrix and

$$\mathbb{V} = \{ \underline{v} \in \mathbb{R}^n \mid 0 < v_i < 1, i = 1, 2, \dots, n \}.$$

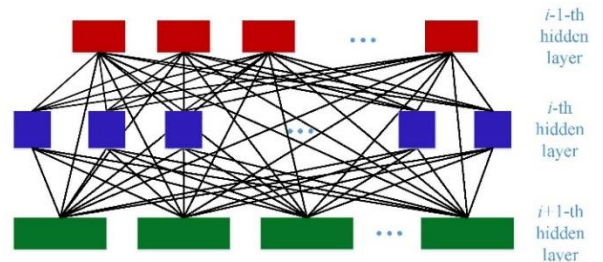


Fig. 2. Neural network with various lengths of the hidden units.

3. The weak formulation and the associate unconstrained minimization problem

Let τ_h be a uniform triangulation of the interval $\Omega=[0,1]$ with n linear finite elements. The associate finite element space \mathbb{V}_h is spanned by nodal basis functions φ_i , $i=1,2,\dots,n$. The space \mathbb{V}_h is provided with the inner scalar product

$$(u, v) = \int_{\Omega} uv dx.$$

A linear operator $T: \mathbb{V} \rightarrow \mathbb{V}_h$ defined by

$$T\underline{v} = \underline{v} \cdot \underline{\varphi}, \quad \underline{v} \in \mathbb{V}, \quad \underline{\varphi} = (\varphi_1, \varphi_2, \dots, \varphi_n)$$

generates a subset $\hat{\mathbb{V}}_h = \{v = T\underline{v} \mid \underline{v} \in \mathbb{V}\}$ of \mathbb{V}_h .

The weak problem is compiled as follows

$$\begin{cases} \text{Find } \underline{u} \in \hat{\mathbb{V}}_h \text{ such that} \\ (u, v) - (s(u), v) = 0, \quad \forall v \in \hat{\mathbb{V}}_h, \end{cases} \quad (3)$$

where $s(v) = T \text{Sigm}(\underline{a} + W\underline{v})$ and $\underline{v} = T^{-1}v$. The weak formulation is a very important part of our investigation. We look for a finite dimensional Hilbert space \mathbb{X} of continuous functions and a bijection between \mathbb{V} and \mathbb{X} . Note that the Gramian matrix of the basis functions is positive definite. As a particular example in our consideration we chose the finite element space $\hat{\mathbb{V}}_h$. In this case the Gramian matrix is actually the mass matrix. The interval $[0,1]$ is not essential for the proof of the main result. Any other finite closed interval can be successfully applied but we use the canonical one-dimensional simplex for the sake of simplicity. The choice of the piecewise linear trial functions is made for the same reasons. The operator T is the wanted bijection.

Lemma 1 All solutions of the weak problem are also solutions of the equation (2).

Proof. The zero vector does not belong to the set \mathbb{V} and the mass matrix M is symmetric and positive definite. Therefore the assertion of the Lemma results from the following equivalences:

$$\begin{aligned} (u, v) - (s(u), v) &= 0 \\ \Leftrightarrow \underline{u} M \underline{v} - \text{Sigm}(\underline{a} + W\underline{u}) M \underline{v} &= 0 \\ \Leftrightarrow (\underline{u} - \text{Sigm}(\underline{a} + W\underline{u})) M \underline{v} &= 0. \end{aligned}$$

Definition 1 The square matrix W is said to be \mathbb{V} -positive definite if $\underline{v}^T W \underline{v} > 0$, $\forall \underline{v} \in \mathbb{V}$.

The main goal of the present investigation is to present problem (2) as a minimization problem. Therefore we define an objective functional

$$J(v) = \frac{1}{2}(v, v) - \left(\int_0^v s(t) dt, 1 \right)$$

to associate weak form (3) with the following minimization problem

$$\arg \min_{v \in \hat{\mathbb{V}}_h} J(v). \quad (4)$$

Theorem 1 establishes existence and uniqueness of the solution of minimization problem (4). To prove this theorem we need the following denotations: $\text{sigm}(x)$ is the well-known logistic function; $\text{diag}(\lambda_i, i=1,2,\dots,n)$ is $n \times n$ a diagonal matrix with λ_i in the main diagonal; I is the identity $n \times n$ matrix;

$$\Lambda(\underline{v}) = \text{diag}(\text{Sigm}(\underline{a} + W\underline{v}), \text{diag}(\hat{e} - \text{Sigm}(\underline{a} + W\underline{u}))).$$

Theorem 1 If the weight matrix W is V -positive definite with spectral norm $\|W\| < 4$ then the problem (4) has a unique solution.

Proof. The functional $J(v)$ is continuous, i.e. we have to prove that $J(v)$ is bounded below, coercive and convex. We estimate the second term in the objective functional using Hölder inequality

$$\left(\int_0^v s(t) dt, 1 \right) \leq \|s\|_{0,\Omega} \|v\|_{1,\Omega} \leq |\Omega|^{\frac{1}{2}} \|s\|_{0,\Omega} \|v\|_{2,\Omega}$$

and the functional

$$\begin{aligned} J(v) &\geq \frac{1}{2} \|v\|_{2,\Omega}^2 - \|s\|_{0,\Omega} \|v\|_{2,\Omega} \\ &= \frac{1}{2} \left(\|v\|_{2,\Omega} - 2 \|s\|_{0,\Omega} \right) \|v\|_{2,\Omega}, \quad \forall v \in \hat{\mathbb{V}}_h. \end{aligned}$$

The latter inequality shows us that $J(v)$ is bounded below and coercive. It remains to prove that $J(v)$ is convex. For this purpose we calculate Fréchet derivatives of $J(v)$:

$$\begin{aligned} DJ(u)v &= (u, v) - (s(u), v), \\ D^2 J(u)(v, v) &= (v, v) - (DsD(\underline{a} + W\underline{u}), \varphi, v) \\ &= (v, v) - (\Lambda(\underline{u}) W \underline{v} \cdot \varphi, v) = \underline{v}^T M \underline{v} - (\Lambda(\underline{u}) W \underline{v})^T M \underline{v} \\ &= \left(\underline{v}^T - (\Lambda(\underline{u}) W \underline{v})^T \right) M \underline{v} = \underline{v}^T \left(I - (\Lambda(\underline{u}) W)^T \right) M \underline{v}. \end{aligned}$$

Since $\|\Lambda(\underline{u})\| \leq \frac{1}{4}$ and $\|W\| < 4$ the second

derivative $D^2 J(u)(v, v) > 0$ and the objective functional is strongly convex.

Corollary 1 If the conditions of Theorem 1 hold then the problem (3) has a unique solution.

Proof. The problem (3) is equivalent to $DJ(u)v = 0$, which means that the solution u of (3) is the unique stationary point of the functional $J(v)$.

4. An iterative method for solving the unconstrained minimization problem

The object of investigation in this section is a two-point gradient method provided with various steplengths. We define the Barzilai-Borwein method

$$(u_{k+1}, v) = (u_k, v) - \frac{1}{\chi_k} DJ(u_k)v, \quad k \geq 1 \quad (5)$$

with a step χ_k for the unconstrained minimization problem (4). The choice of the steplength is very important with respect to the rate of convergence and the overall necessary computational work. In the present investigation we propose the following original steplengths:

$$\beta_k = \frac{\|u_k\|_{2,\Omega}^2 + \|u_{k-1}\|_{2,\Omega}^2}{\|s(u_k)\|_{2,\Omega}^2 + \|s(u_{k-1})\|_{2,\Omega}^2}$$

$$\gamma_k = \frac{\sum_{i=k-1}^k \left((I - Ds(u_i)) DJ(u_i), DJ(u_i) \right)}{\sum_{i=k-1}^k \|DJ(u_i)\|_{2,\Omega}^2}$$

Additionally, we define the classical steplength

$$\alpha_k = \frac{\Delta DJ(u_k) \Delta u_k}{\|\Delta u_k\|_{2,\Omega}^2}, \quad \text{where } \Delta u_k = u_k - u_{k-1},$$

$$\Delta DJ(u_k) = DJ(u_k) - DJ(u_{k-1}).$$

The identity matrix and the identity operator are denoted by the same letter I . The steplength γ_k is a preconditioned version of the step length

$$\hat{\gamma}_k = \frac{\sum_{i=k-1}^k D^2 J(u_i) (DJ(u_i), DJ(u_i))}{\sum_{i=k-1}^k \|DJ(u_i)\|_{2,\Omega}^2}.$$

The steplength $\hat{\gamma}_k$ gives better results than γ_k but $\hat{\gamma}_k$ essentially increase computational complexity of the problem. That is why further we analyze the method with the steplength γ_k . Similar approach can be applied when β_k or $\hat{\gamma}_k$ are used.

Theorem 2 assures the convergence of the Barzilai-Borwein method independently of the initial guesses.

Theorem 2 The main result. Let the weight matrix W is V -positive definite and

$$\|W\| = \rho_0 < 4 \quad (6)$$

Then the sequence $\{u_k\}$ generated by the two-point step size gradient method (5) with the steplength γ_k converges Q -linearly to the weak solution u .

Proof. The theorem is proved under random initial guesses. We just suppose that $u_0, u_1 \in \hat{V}_h$. The

$$\text{equality } (u_{k+1}, v) = (u_k, v) - \frac{1}{\gamma_k} DJ(u_k)v, \quad k \geq 1$$

is transformed into the equality

$$\gamma_k (e_{k+1}, v) = \gamma_k (e_k, v) + DJ(u)v - DJ(u_k)v$$

having in mind that u is the solution of (3) and $e_k = u_k - u$ is the error in the approximate solution u_k . Applying the mean value theorem and replacing $v = e_{k+1}$ we obtain

$$(e_{k+1}, e_{k+1}) = (e_k, e_{k+1}) - \frac{1}{\alpha_k} D^2 J(\tilde{u}_k)(e_k, e_{k+1}), \quad (7)$$

where $\tilde{u}_k = u_k - \mathcal{G}_k e_k$ for some $\mathcal{G}_k \in (0, 1)$. Since

$$0 < \|\Lambda(v)\| \leq \frac{1}{4}, \quad \forall v \in \mathbb{V},$$

the steplength $\hat{\gamma}_k$ is estimated by $\frac{3}{4} < \gamma_k < \frac{4}{3}$

The inequality $1 - \frac{\rho_0}{4} \leq \|D^2 J(\tilde{u}_k)\|$ follows from (6). We continue estimating the left hand side of (7)

$$(e_{k+1}, e_{k+1}) < \left\| I - \frac{1}{\gamma_k} D^2 J(\tilde{u}_k) \right\| \|e_k\|_{2,\Omega} \|e_{k+1}\|_{2,\Omega} < \mu_k \|e_k\|_{2,\Omega} \|e_{k+1}\|_{2,\Omega}.$$

The contraction factor

$$\mu_k \text{ is bounded above by } \mu_k < \frac{1}{4} + \frac{3}{16} \rho_0 = \mu < 1.$$

Finally

$$\|e_{k+1}\|_{2,\Omega} < \mu_k \|e_k\|_{2,\Omega} < \mu^k \|e_1\|_{2,\Omega}, \quad \forall k \in \mathbb{N}.$$

Further we suppose that all hidden units in the i -th layer have the same length m_i and $m = \max_{i=1,2,\dots,p} m_i$.

Let h_{ij} be the j -th hidden unit in the i -th hidden layer and $W(h_{ij}, h_{jk})$ be the weight matrix generated by the connection between the neurons h_{ij} and h_{jk} , $i \neq k$. We only consider the particular case $0 < m_i < m_k < m$. In this case the weight matrix $W(h_{ij}, h_{jk}) = (A(m_i \times m_i) B(m_i \times m_k - m_i))$ is a block rectangular matrix. Similar approach can be used in the other cases. To obtain a problem

corresponding to (2) we extend all weight matrices and hidden units as follows:

$$\left\{ \begin{array}{l} \underline{h}_{ij} \mapsto \hat{h}_{ij}(\underline{h}_{ij}, \underline{o}_i), \underline{h}_{kl} \mapsto \hat{h}_{kl}(\underline{h}_{kl}, \underline{o}_k) \\ W(\underline{h}_{ij}, \underline{h}_{kl}) \mapsto \hat{W} = \begin{pmatrix} A & B & C \\ O & I_{22} & D \\ C^T & D^T & I_{33} \end{pmatrix}, \end{array} \right.$$

where $\underline{o}_i \in \mathbb{R}^{m-m_i}$ and $\underline{o}_k \in \mathbb{R}^{m-m_k}$ are zero vectors, $O(m_k - m_i \times m_i)$, $C(m_i \times m - m_k)$ and $D(m_k - m_i \times m - m_k)$ are zero matrices, and $I_{22}(m_k - m_i \times m_k - m_i)$ and $I_{33}(m - m_k \times m - m_k)$ are identical matrices. Then $\hat{h}_{ij}, \hat{h}_{kl} \in \mathbb{R}^m$ and \hat{W} is a $m \times m$ matrix. Thus we reduce the rectangular version of the mean field approximation problem to a square one.

Table 1. The number of necessary iterations for satisfying the stop criterion when $\|W\| = 1$.

$v_k \setminus n$	8	100	512	1024	2048	4096	6400
α_k	7	8	8	8	9	9	9
β_k	8	9	9	9	9	9	9
γ_k	11	13	14	15	15	15	16
FPI	7	8	8	9	9	9	9

Table 2. The number $\eta(\chi_k, n, 2)$ obtained by random initial guesses.

$v_k \setminus n$	8	100	512	1024	2048	4096	6400
α_k	9	10	11	11	11	11	13
β_k	11	11	12	12	12	13	13
γ_k	12	16	18	18	18	18	19
FPI	11	11	12	12	12	12	13

5. Experiments

From the computational point of view method (5) takes the following attractive form

$$\underline{u}_{k+1} = \frac{1}{\chi_k} (\text{Sigm}(\underline{a} + W\underline{u}_k) - (1 - v_k)\underline{u}_k), \quad k \geq 1$$

with the steplengths:

$$\alpha_k = \frac{\Delta s_k \cdot \Delta \underline{u}_k}{\|\Delta \underline{u}_k\|^2}, \quad s_k = \underline{u}_k - \text{Sigm}(\underline{u}_k),$$

$$\beta_k = \frac{\|\underline{u}_k\|^2 + \|\underline{u}_{k-1}\|^2}{\|\text{Sigm}(\underline{u}_k)\|^2 + \|\text{Sigm}(\underline{u}_{k-1})\|^2},$$

$$\gamma_k = \frac{\sum_{i=k-1}^k s_i (I - \Lambda(\underline{u}_i)) s_i}{\sum_{i=k-1}^k \|s_i\|^2}.$$

We consider different cases with respect to the features of the weight matrix. The stop criterion

$$\|DJ(\underline{u}_k)\| < \varepsilon. \quad (8)$$

with $\varepsilon = 10^{-13}$ is used throughout all our considerations. The number of necessary iterations $\eta(\chi_k, n, \|W\|)$ for satisfying (8) is an object of interest in this section. The fixed-point iteration is denoted by FPI in all tables.

Let's start with the case where W is V -positive definite and $\|W\| < 4$, i.e. the weight matrix satisfies the requirements of Theorem 2. In this case we have very fast convergence for all steplengths. The same rate of convergence is also established for the fixed-point iteration. Table 1 and 2 indicate that the rate of convergence does not depend neither on the number of the hidden neurons nor the initial guesses. Further we analyze the case when $\|W\| > 4$ and the requirement for V -positivity is broken. We establish low rate of convergence in the cases when $4 < \|W\| \leq 20$. Some examples are presented in Table 3. In this case the traditional steplength proposed by J. Barzilai and J. Borwein [9] assures nonmonotone convergence of the error norm. The sequence $\{\|\underline{e}_k\|\}$ converges monotonically to zero when the method is applied with steplengths β_k and γ_k . The best results for ill-conditioned matrix are obtained by the steplength γ_k . The number $\eta(\chi_k, n, \|W\| > 4)$ grows nonlinearly when the number of unknowns n tends to infinity. The convergence of the two point step size method is not influenced by whether the weight matrix is singular or not. The initial guesses does not affect on the rate of convergence as well. All computational tests with $\|W\| \gg 20.5$ and $n \geq 4096$ have failed since the Barzilai-Borwein method is divergent with all considered steplengths as well as the fixed-point iteration.

6. Conclusion

The Barzilai-Borwein method for solving the mean field approximation problem is studied in the present paper. Two original steplengths, which assure monotone decreasing of the norm of the error in approximate solutions are found. The mean filed approximation problem is reduced to a minimization one. The uniqueness and existence of

the weak solution are proved. A rigorous proof of the convergence theorem for the Barzilai-Borwein method with the steplength γ_k is made.

Table 3. The number of necessary iterations for satisfying the stop criterion when $\|W\| > 4$.

v_k	$n = 4096,$ $\ W\ = 20.1$	$n = 6400,$ $\ W\ = 17$
α_k	divergence	930
β_k	1296	237
γ_k	1037	204
FPI	1616	281

The same approach can be applied for the steplength β_k but it is totally inapplicable for the steplength α_k . The classical steplength α_k introduced by J. Barzilai and J. Borwein [9] for quadratics assures nonmonotone decreasing of the norm $\|e_k\|$. A comparison between the two point step size gradient method and the fixed-point iteration method is made. The computational tests indicates that

$$\eta(\beta_k, n, \|W\|) \approx \eta(\text{FPI}, n, \|W\|), \quad \|W\| < 4,$$

$$\eta(\gamma_k, n, \|W\|) \ll \eta(\text{FPI}, n, \|W\|), \quad \|W\| \gg 4.$$

Moreover, the steplength γ_k is superior in the case of ill-conditioned weight matrices. The features of the weight matrix affect the convergence of both considered methods. The mean field approximation problem has no solution in any case. The lack of solutions leads to improper working of the corresponding neural network. Any divergent mean field approximation procedure generates confusions in the process of deep machine learning. To avoid this difficulty, it is best to work with normalized weight matrices. Note that there is a higher rate of convergence with a smaller norm of the weight matrix. This is very important from the practical point of view.

Finally, a designing of artificial neural networks based on the mean field approximation without control on the features of the weight matrix can lead to confusions in the deep learning process.

REFERENCES

- Salakhutdinov R. and Hinton G. E., Deep Boltzmann machines, *In Proc. of the Int. Conf. on Artificial Intelligence and Statistics (AISTATS 2009)*, pp. 448-455, 2009.
- Salakhutdinov R., Learning Deep Boltzmann Machines using Adaptive MCMC, *Appearing in Proceedings of 27 th International Conference on Machine Learning*, Haifa, Israel, 2010.
- Salakhutdinov R., Larochelle H., Efficient Learning of Deep Boltzmann Machines, *Journal of Machine Learning Research*, vol. 9, pp. 693-700, 2010.
- Hinton G., Salakhutdinov R., An Efficient Learning Procedure for Deep Boltzmann Machines, *Neural Computation*, vol. 24, no. 8, pp. 1967-2006, 2012.
- Cho K., Raiko T., Ilin A., Karhunen J., A Two-Stage Pretraining Algorithm for Deep Boltzmann Machines, *Artificial Neural Networks and Machine Learning-ICANN 2013*, Vol. 8131 of the series Lecture Notes in Computer Science, pp. 106-113.
- Cho K., Raiko T., Ilin A., Gaussian-Bernoulli deep Boltzmann machine, *IEEE International Joint Conference on Neural Networks*, Dallas, Texas, USA, pp. 1-7, 2013.
- Dreameau A., Boltzmann Machine and Mean-Field Approximation for Structured Sparse Decompositions, *IEEE Transactions on Signal Processing*, Vol. 60, Issue 7, 2012, pp. 3425-3438.
- Srivastava N., Salakhutdinov R., Multimodal learning with Deep Boltzmann Machines, *Journal of Machine Learning Research*, vol. 15, pp. 2949-2980, 2014 .
- Barzilai J., Borwein J. M., Two-point step size gradient methods, *IMA Journal of Numerical Analysis*, Vol. 8, Issue 1, 1988, pp. 141-148.
- Birgin E. G., Martinez J. M., Raydan M., Spectral Projected Gradient methods: Review and Perspectives, *Journal of Statistical Software*, Vol. 60, Issue 3, 2014, pp. 1-21.
- Raydan M., On the Barzilai and Borwein choice of steplength for the gradient method, *IMA Journal of Numerical Analysis* 13, 1993, 321-326.
- Todorov T. D., Nonlocal problem for a general second-order elliptic operator, *Computers & Mathematics with Applications*, vol. 69, issue 5, 2015, pp. 411-422.

Authors' contacts.

Department of Mathematics and Informatics,
Technical University of Gabrovo,
5300 Gabrovo,
e-mail: t.todorov@yahoo.com

CORROSION PROTECTION WITH INHIBITORS QUATERNARY AMMONIUM BROMIDES

ANGELINA POPOVA

Abstract: *Two quaternary ammonium bromides are used as inhibitors of mild steel corrosion in 1 M HCl. Their behavior is studied with the application of gravimetric and potentiodynamic voltammetry methods. Additional gravimetric experiments are carried in 1 M H₂SO₄ aiming a comparison of the protective properties observed. The information summarized leads to the conclusion that the inhibitive properties of the bromides investigated depend on their concentration and molecular structure.*

Key words: *mild steel, corrosion, inhibitors.*

1. Introduction

The enormous material losses in the field of industry caused by corrosion make the latter a significant economic problem. That is why the elaboration of methods for corrosion protection is a research trend of great priority [1-3].

The use of corrosion inhibitors is an approach of great importance [4,5]. In an acid medium the inhibition is exercised by a layer of the inhibitor's adsorbed molecules. The elucidation of a connection between the protective properties and the molecular structure of the organic substances used as inhibitors in aqueous acid solutions is of a profound interest.

It is found that compounds of the group of the quaternary ammonium salts provide very good protective properties. The molecular structure of the substances studied is chosen to outline the effect of the adsorbed species area through comparing the behavior of molecules of a different size (Table 1). The compounds used in this investigation are specifically synthesized and studied for the first time as inhibitors..

2. Experimental

Two classical techniques were used to determine the corrosion inhibitor characteristics of the quaternary ammonium bromides – gravimetry and potentiodynamic voltammetry.

The gravimetric measurements were carried out at definite time intervals of 24h at a room temperature (20±2°C) using an analytical balance. The specimens of an area of 11.3 cm² were of a round shape to avoid edges effects attributed to high-speed corrosion proceeding. The preliminary treatment included pickling in a solution containing

concentrated HNO₃ and H₂SO₄, washing with distilled water and an ethanol-ether mixture. Three specimens were immersed simultaneously in every beaker containing 250 mL of the test solution.

The potentiodynamic polarization experiments were carried out in a conventional three-compartment electrochemical cell. A mild steel cylinder pressed into a Teflon holder served as a working electrode (WE). Its working area of 0.5 cm² remained precisely fixed. A saturated calomel electrode (SCE) connected through a salt bridge was used as a reference electrode, while platinum sheet acted as counter electrode. Prior to each experiment the WE was wet abraded with 600-grade emery paper, rinsed with distilled water and an ethanol-ether mixture. Then it was inserted immediately into the glass cell which contained 250 mL of the test solution.

EG&G Instruments PAR model 273 potentiostat monitored by an IBM personal computer via a GPIB-IIA interface and M342 software were used to run the tests as well as to collect and treat the experimental data.

The polarization curves were recorded from ca -250 mV to +130 mV vs. the measured corrosion potential E_{corr} with a scan rate of 0.2 mV/s starting one minute after the immersion of the WE in the test solution. The anodic (b_a) and cathodic (b_c) Tafel constants, the corrosion potential (E_{corr}) and the corrosion current density (j_{corr}) were determined using PARCalc342 data analysis technique.

All plots and calculated parameters were mean values of at least five independent experiments. Standard deviations were also reported.

The inhibitor efficiency IE (%) was calculated from the gravimetric and potentiodynamic measurements using the relations:

$$IE(\%) = \frac{w_0 - w_i}{w_0} \times 100 \quad (1)$$

$$IE(\%) = \frac{j_{corr,0} - j_{corr,i}}{j_{corr,0}} \times 100 \quad (2)$$

where w_0 and w_1 in g/m^2h are the average corrosion rates in absence and presence of an inhibitor, while $j_{corr,0}$ and $j_{corr,i}$ are the corresponding corrosion current densities.

3. Results and Discussion

3.1. Inhibiting efficiency

Gravimetric tests. The quaternary ammonium salts are studied in 1M HCl and 1M H_2SO_4 solutions of a wide concentration range. The highest possible concentration value is limited by the plateau reached in the corresponding concentration dependence of the inhibiting efficiency or the compound solubility (10^{-3} M for SS). The mean values of the corrosion rate, w (g/m^2h), and that the corresponding inhibiting efficiency, IE (%), evaluated in accordance with Eq. (1), are determined for each concentration studied. Fig. 1 illustrates the results obtained in 1 M HCl medium, while Fig. 2 – those in 1 M H_2SO_4 solution. It is worth noting that the value of the surface coverage degree, θ , is plotted on the ordinate axis assuming that $IE(\%) = \theta \times 100$. The same data is used in adsorption characteristics determination. The figures pointed above show that the efficiency of the compounds studied depends on their concentration. It is seen that the protective effect increases initially with the increase of the concentration, c_i . It changes slightly upon reaching a definite concentration value accepted as an optimal one. The latter has a different characteristic value for each inhibitor. The difference in the

Table 1. Investigated quaternary ammonium salts as inhibitors

Compounds	Structural formulae
3-Methylbenzo[d]thiazol-3-ium bromide (MTB)	
3,4-Dihydro-2H-benzo[4,5]thiazolo[2,3-b][1,3]thiazin-5-ium bromide (SS)	

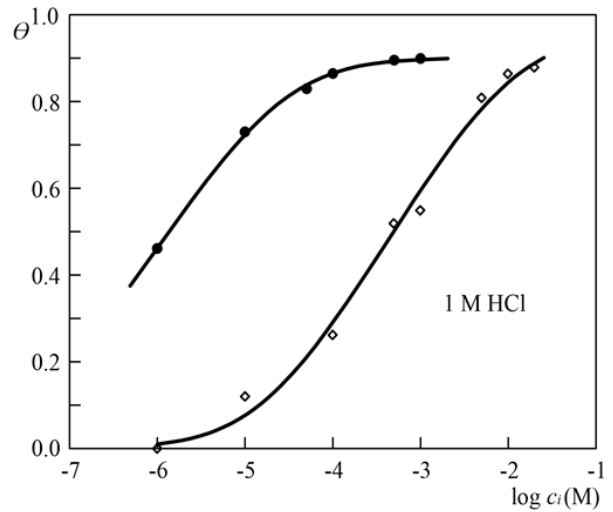


Fig. 1. Inhibition efficiency and adsorption isotherms of quaternary ammonium bromides in 1 M HCl - experimental gravimetric data referring to SS (●) and MTB (◇).

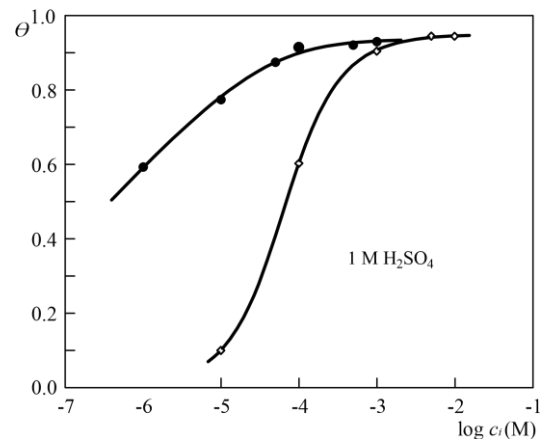


Fig. 2. Inhibition efficiency and adsorption isotherms of quaternary ammonium bromides in 1 M H_2SO_4 - experimental gravimetric data referring to SS (●) and MTB (◇).

inhibiting properties of the three compounds is obviously connected with the difference in their molecular structure. It is more vividly expressed at low and medium concentration values.

The comparison of the inhibitors in 1 M HCl in view of the maximal efficiency reached leads to the following line: SS (90.0%) > MTB (88.0%).

At a concentration of 1×10^{-4} M the line changes to: SS (86.5%) > MTB (26.0%).

The same line is obtained at lower inhibitor concentrations as well.

The inhibitor's sequence following their highest efficiency in 1 M H₂SO₄ is as follows: MTB (94.0%) ≈ SS (92.5%).

At a concentration of 1x10⁻⁴ M the line changes to: SS (91.0%) > MTB (60.0%)

The results of the gravimetric investigation show that all two compounds have inhibitive properties in 1 M HCl and 1 M H₂SO₄. The latter results are slightly better than those in HCl. The difference is better outlined in presence of MTB.

SS provides the best protective properties in 1 M HCl in the whole concentration range studied. It is also the best inhibitor among those investigated in 1 M H₂SO₄ at the concentration values used. In fact MTB is slightly better but at concentrations higher than SS maximal concentration studied. The latter is in fact insoluble at these MTB concentrations.

Potentiodynamic voltammetry tests. The potentiodynamic voltammetry investigations are carried out in 1 M HCl at various concentrations of the inhibitors. The polarization curves recorded provide the determination of the electrochemical parameters as the corrosion current density, j_{corr} , the corrosion potential, E_{corr} , the cathodic and anodic Tafel slopes, b_c and b_a , correspondingly, the polarization resistance, R_p . The variation of the electrochemical parameters provides to follow the effect of the inhibitors on the kinetics of the corrosion process. The values obtained for SS are listed in Table 2.

Table 2. Electrochemical parameters

SS				
c_i M	E_{corr} mV	$-b_c$ mV/dec	b_a mV/dec	j_{corr} μA/cm ²
1x10 ⁻⁶	-520 ± 2	136 ± 3	70 ± 4	458 ± 47
1x10 ⁻⁵	-520 ± 3	136 ± 4	75 ± 4	315 ± 16
1x10 ⁻⁴	-530 ± 2	140 ± 4	89 ± 2	195 ± 12
5x10 ⁻⁴	-532 ± 2	140 ± 5	110 ± 7	160 ± 4
1x10 ⁻³	-540 ± 2	143 ± 4	145 ± 11	146 ± 6

With concentration increase E_{corr} shifts in a positive direction in presence of MTB, while the presence of SS brings about a shift in the opposite direction. The values of the Tafel slopes, b_c and b_a , increase with the increase of all compounds concentration. The results pointed above lead to the conclusion that all compounds are in fact general mixed type inhibitors. MTB shows slightly better expressed anodic behavior, while that of SS tends to the cathodic one.

The corrosion current density, which is a measure of the corrosion rate decreases with increase of the inhibitor's concentration, c_i . This is valid for all compounds studied. The inhibiting efficiency, IE (%), is evaluated on the ground of the values of j_{corr} with the application of Eq. (2). The dependence of IE (%) on c_i is presented in Fig. 4. It is evident that the inhibiting effect increases with the concentration increase in case of all compounds studied. This result is in correspondence with the gravimetric findings. It is worth noting that the values of IE (%) found potentiodynamically are generally lower than those obtained gravimetrically. This is most probably due to the different exposure time in the solution – it is 24 h in the course of the gravimetric tests, while the polarization curves are potentiodynamically recorded immediately after the electrode immersion. Irrespective of this difference the sequence of the inhibitors in respect to their efficiency stays unchanged.

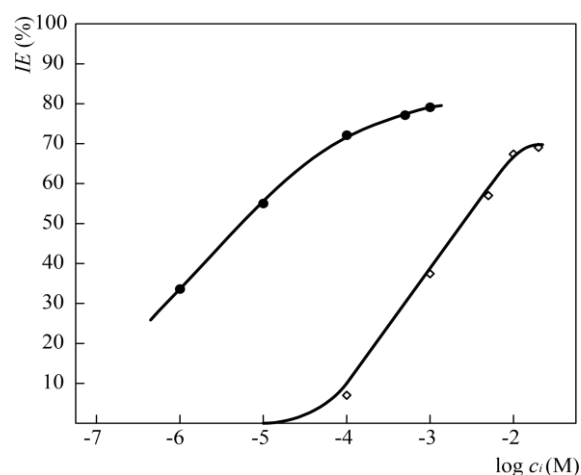


Fig. 3. Dependence of the inhibiting efficiency, IE (%), on the inhibitors concentration, c_i (M) obtained on the ground of the potentiodynamic voltammetric data obtained: SS (●) and MTB (◇).

Thus the inhibitors line obtained on the ground of the highest effect recorded at the highest available concentration is as follows: SS (79.19) > MTB (69.1%).

At a concentration of 1x10⁻⁴ M the line becomes: SS (72.1%) > MTB (10.0%). It has to be underlined that the lines just presented coincide with those obtained gravimetrically in 1 M HCl.

The results obtained with potentiodynamic voltammetry application show that SS has the best protective properties in the whole concentration range studied. It can be concluded that the

molecular structure affects also the inhibitive properties displayed.

3.2. On the inhibiting activity and molecular structure relation

The different inhibiting properties shown by the two compounds obtained in the course of study with the application of different methods are due to the difference in their molecular structure.

The quaternary ammonium bromides dissociate in acidic solutions giving ammonium cations and a bromine anion, Br⁻. We assume that the cation structure determines the difference of the inhibitive properties. The effect of Br⁻ presence in the solution cannot be excluded because of this anion disposition to specific adsorption on the metal surface [5, 6]. This can result to a change of the surface charge or a synergistic effect, i.e. Br⁻ contribute to the inhibiting effect observed. On the other hand, Br⁻ can favor the adsorption by decreasing the cations repulsion in the adsorption layer (the positive values of the interaction adsorption parameter indicate actual attraction in the adsorption layer).

We assume that the juxtaposition of MTB on one hand and SS on the other can outline the effect of the adsorbed species surface area (in fact of the cationic part in this case). The experimental results show that MTB is a weaker inhibitor when compared to SS at an identical concentration. This is most probably due to the smaller surface of its cationic part.

At this stage of our investigations we assume a probability of Br⁻ specific adsorption on the metal surface in both acidic solutions. This brings about a negative charge to the surface and explains the relatively close inhibiting efficiency in both acids. This effect can also favor the physical adsorption of the molecules cationic part.

4. Conclusions

The two compounds investigated show protective properties in case of mild steel corrosion in 1 M HCl and 1 M H₂SO₄. Their inhibiting effect increases with their concentration increase.

The results obtained show that the molecular structure affects essentially the inhibitive properties when compared at identical concentrations. We assume that physical adsorption of the molecules cationic parts takes place at the negatively charged metal surface.

REFERENCES

1. Shrier, L. L. (1981), *Corrosion*, Metallurgia, Moscow.
2. Zhuk, N. P. (1976), *Course of Corrosion and Metal Protection*, Metallurgia, Moscow.
3. Roberge, P. R. (2000), *Handbook of Corrosion Engineering*, McCrow-Hill, New York.
4. Rozenfeld, I.L. (1977), *Corrosion inhibitors*, Khimiya, Moscow.
5. Popova, A., Christov, M., Vasilev, A., and Zwetanova, A. (2011), Mono- and dicationic benzothiazolic quaternary ammonium bromides as mild steel corrosion inhibitors. Part I: Gravimetric and voltammetric results, *Corrosion Science*, 53, 679-686.

Department of Physical Chemistry
University of Chemical Technology and
Metallurgy
8, Kl. Ohridski Bld. ,1756 Sofia,
BULGARIA
E-mail: apopova@uctm.edu

STUDY OF CORROSION BEHAVIOUR OF ALUMINIUM ALLOYS EN AW-2011 and EN AW-2024

KALINA KAMARSKA

Abstract: *Aluminium alloys are construction materials widely used in automotive, aircraft and chemical industries. In their exploitation, in natural or technological environments, they are in contact with aggressive components and interacting with them are progressively destroyed. The corrosion study of these materials is of immense technological importance due to their growing industrial application. This article presents the results of a study of the corrosion behaviour of aluminium alloys EN AW-2011 and EN AW-2024 in nitric acid (HNO₃) and sodium chloride (NaCl) solutions with different pH values of the medium. The corrosion resistance of aluminium alloy samples is determined by a gravimetric method. The results obtained show that the corrosion rate of the aluminium alloys tested depends on the nature of the medium and the concentration of nitric acid.*

Keywords: *aluminium alloys EN AW-2011 and EN AW-2024, corrosion resistance*

1. Introduction

Aluminium alloys are applied in industry thanks to a number of valuable properties such as low weight, high strength, hardness and the ability to form a thick protective layer of Al₂O₃. The corrosion behaviour of aluminium and its alloys is largely determined by the chemical resistance of this layer [1]. Any impact that contributes to the removal of this layer or its disintegration enhances the corrosion of aluminium. The nature (pH), composition, concentration, temperature and other parameters of the corrosion environment are essential for its corrosion rate [2]. In a neutral environment, aluminium is moderately stable, and in acidic and alkaline environment it has a high corrosion rate [3]. The presence of aggressive ions, such as chloride, activate the corrosion process, impair the stability and integrity of the passive layer [4]. Nitric acid acts in a different way on aluminium - diluted nitric acid (10÷60%) intensely destroys it, while the concentrated passivates its surface and corrosion decreases.

The corrosion behaviour of aluminium alloys also depends on the characteristics of the parent metal, the chemical composition and the number of the alloying elements [5]. To improve the mechanical properties of aluminium, certain elements are added to it, but they give rise to electrochemical heterogeneity in the microstructure and are among the main causes of corrosion of its alloys [6].

The studied aluminium alloys EN AW-2011 and EN AW-2024 refer to the series of 2xxx alloys

in which the main alloying element is copper (up to 5%) and as additives, small amounts (total of 2-3%) of magnesium, manganese, iron, silicon etc.

These alloys are especially suitable for parts and structures requiring high strength and are often used to make parts that require good strength at temperatures up to 150°C. EN AW-2011 is used for making nuts, bolts, screws, studs, car parts and more. EN AW-2024 is used for making high-strength parts, bolts, screws as well as high-strength structural components, truck wheels, aircraft, car parts, etc. Compared to other aluminium alloys, alloys of this series have lower corrosion resistance [7]. One reason for this is the much higher amount of copper that weakens the protective properties of the oxide layer [8].

The purpose of this study is to determine the corrosion resistance of aluminium alloys EN AW-2011 and EN AW-2024 in nitric acid solutions and in sodium chloride solutions with different pH values of the medium, and thus to compare their corrosion behaviour in these environments.

2. Experimental

Aluminium alloy samples with a total surface area of 8.54 cm² were studied.

2.1. Study of the corrosion rate of aluminium alloys EN-AW 2011 and EN-AW 2024 in solutions with different concentrations of nitric acid

The corrosion behaviour of the aluminium samples is found in 10%, 30%, 50% and 60% nitric acid solutions at 50±5°C. Before testing, the

samples have been placed in ethyl alcohol for 5 minutes, washed with distilled water and dried. They are then placed in the said nitric acid solutions for 4 hours. The weight of the samples was measured before (m_1) and after (m_2) testing using an Acculab ATILON analytical scales accurate to $\pm 0,0001\text{g}$.

In order to assess the corrosion resistance of the alloys under test, in the said conditions, the gravimetric method was used, and by the change in the weight of the tested samples in the corrosive medium, the rate of corrosion (K_m) was determined:

$$K_m = (m_1 - m_2) / S \cdot t \text{ [g/m}^2 \cdot \text{h]} \quad (1),$$

where

m_1 – the weight of the sample, g;

m_2 – the weight of the sample after the corrosion test, g;

S – the area of the sample, m^2 ;

t – test time, h.

On the basis of the value of K_m , conclusions are drawn for the corrosion behaviour of aluminium alloy samples.

2.2. Study of corrosion rate of aluminium alloys EN-AW 2011 and EN-AW 2024 in solutions with different pH

The corrosion resistance of aluminium samples is determined in two solutions with a different concentration of NaCl (1% and 3%) at pH values of 1 to 13. To achieve the required pH of the medium, hydrochloric acid or sodium hydroxide is added to the solutions. The pH of the medium is measured with a laboratory pH meter MS2006. Prior to testing, the samples were placed in ethyl alcohol for 5 minutes, washed with distilled water and dried. They are then placed in 1% NaCl solution and in 3% NaCl solutions at a pH of 1 to 13 at room temperature for 4 hours. The weight of the samples was measured before (m_1) and after (m_2) the Acculab ATILON analytical scales test with accuracy to within $\pm 0,0001\text{g}$. The corrosion resistance of aluminium alloy samples is determined by a gravimetric method, using formula (1) for calculation.

3. Results and discussion

3.1. Corrosion test results of aluminium alloys EN AW-2011 and EN AW-2024 in solutions with different concentrations of nitric acid

Figure 1 shows the comparative graph showing the corrosion rate dependence of aluminium alloys EN AW-2011 and EN AW-2024 on the concentration of nitric acid at a temperature of $50 \pm 5^\circ\text{C}$. The data obtained for their corrosion rates (Table 1) are processed by the gravimetric

method and serve to assess their corrosion resistance.

The results show that the corrosion rate of the aluminium alloys studied depends on the concentration of nitric acid. In diluted nitric acid, aluminium alloys are actively dissolved [9, 10, 11].

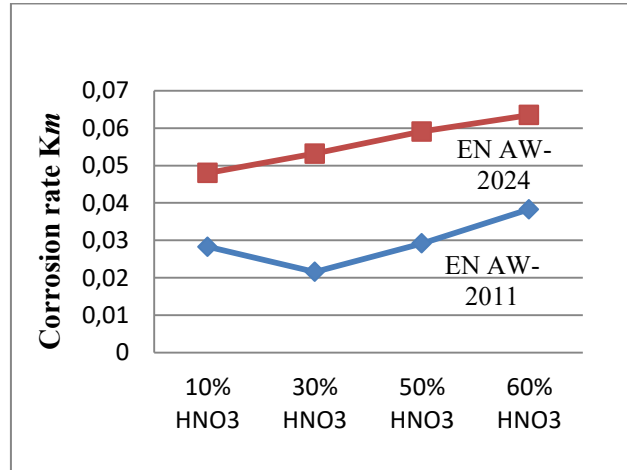


Fig. 1. Corrosion behaviour of aluminium alloys EN AW-2011 and EN AW-2024 in nitric acid solutions at temperature $50 \pm 5^\circ\text{C}$

With an increase of nitric acid concentration between 10%–60% at $50 \pm 5^\circ\text{C}$, the corrosion rate of said alloys increases (Fig.1). The corrosion rate value of EN AW-2011 is lower and is more suitable for use in a similar environment.

Table 1.

Corrosion rate of samples of aluminium alloys EN AW-2011 and EN AW-2024 at concentrations of nitric acid of 10% to 60% at temperature $50 \pm 5^\circ\text{C}$

Test no.	Concentration of nitric acid (C_{HNO_3}), %	Corrosion rate of EN AW-2011 (K_m), $\text{g/m}^2 \cdot \text{h}$	Corrosion rate of EN AW-2024 (K_m), $\text{g/m}^2 \cdot \text{h}$
1	10% HNO_3	0,0275	0,0480
2	30% HNO_3	0,0283	0,0532
3	50% HNO_3	0,0292	0,0591
4	60% HNO_3	0,0383	0,0635

3.2. Results of corrosion rate of aluminium alloys EN-AW 2011 and EN-AW 2024 in 1% and 3% NaCl in different pH

The data obtained are processed by the gravimetric method and presented in tables (Table 2 and Table 3).

Aluminium alloys tested demonstrate different behaviour in 1% NaCl solution in acidic medium (pH 1 - 2). With EN AW-2011 there was a higher weight loss and a significantly higher corrosion rate than EN AW-2024. The EN AW-

2024 alloy has a higher corrosion resistance and is suitable for use in a similar environment. In a 1% solution of NaCl in alkaline medium, the two alloys show similar behavior and are relatively stable.

As the NaCl concentration increases, it is noted that the rate of corrosion increases [12]. This is more noticeable in acidic (pH 1) and alkaline medium (pH 13).

In 3% NaCl in a strongly acid environment, the EN AW-2011 corrosion value is higher than that of EN AW-2024. In 3% NaCl in a strongly alkaline environment these alloys are also unstable. This is probably due to the presence of aggressive chloride ions which break the stability of the protective film on the surface of aluminium alloys and are the cause of their corrosion [4]. Upon monitoring the corrosion behaviour of EN AW-2011 and EN AW-2024 samples in a strongly alkaline medium (pH 13) it was noted that in a 3% NaCl solution there is a vigorous release of hydrogen gas and the biggest loss of weight of aluminium samples is observed. This indicates that these alloys show instability in the specified environment, they are very quickly destroyed, and they are not suitable for making parts that work in similar conditions.

Table 2.

Corrosion rate of aluminium alloys EN AW-2011 and EN AW-2024 in 1% NaCl in different pH

pH	NaCl , %	Km EN AW-2011	Km EN AW-2024
1	1%	0,0371	0,0009
2	1%	0,0351	0,0006
7	1%	0,0017	0,0006
12	1%	0,0020	0,0020
13	1%	0,0026	0,0020

Table 3.

Corrosion rate of aluminium alloys EN AW-2011 and EN AW-2024 in a 3% solution of NaCl in a different pH

pH	NaCl , %	Km EN AW-2011	Km EN AW-2024
1	3%	0,0462	0,0030
2	3%	0,0371	0,0017
7	3%	0,0023	0,0014
12	3%	0,0043	0,0020
13	3%	0,1984	0,2263

4. Conclusion

1. The results obtained show that in diluted nitric acid (10÷60%) at a temperature of 50±5°C, the studied aluminium alloys are intensively destroyed. Increasing the concentration of HNO₃ increases their corrosion rate. In the studied

corrosion environment EN AW-2011 has a higher corrosion resistance than EN AW-2024 and is more suitable for use in similar conditions.

2. In the studied highly acidic environments, in NaCl solutions, EN AW-2024 aluminium alloy has a higher corrosion resistance than EN AW-2011, and it is more suitable for making parts that work in environments of similar composition.

3. In a strongly alkaline environment (pH 13) with increasing NaCl concentration, the corrosion rate of the investigated alloys EN AW-2011 and EN AW-2024 is increased and they are corrosion non-resistant.

REFERENCES

1. Rachev, R. (2000). *Corrosion and protection of metals*, Novi znanija, Sofia.
2. Veleva, M., Kopchev, P., Obreshkov, K. (1987). *Chemistry*, Nauka i izkustvo, Sofia.
3. Veleva, M., Stoychev, D., Kopchev, P., Obreshkov, K. (1999). *Chemistry of construction and exploitation materials*. Multiprint, Sofia.
4. Huang, I.-Wen. (2016). Uniform Corrosion and General Dissolution of Aluminium Alloys 2024-T3, 6061-T6, and 7075-T6. Dissertation for the Degree Doctor of Philosophy in the Graduate School. Columbus, Ohio. Volume (79-01), 197
5. Esquivel, J. and R. K. Gupta. (2017). Corrosion Behaviour and Hardness of Al-M (M: Mo, Si, Ti, Cr) Alloys. *Acta Metallurgica Sinica(English letters)*, 30, 333-341.
6. Sukiman, N., Zhou, X., Birbilis, N., Hughes, A., Mol, J. Garcia, S., Thompson G. (2012). Durability and corrosion of aluminium and its alloys: overview, property space, techniques and developments. In *Aluminium Alloys – New Trends in Fabrication and Applications*, 224 – 262 Intech Publications, Rijeka.
7. Davis, J. R. (2001). *Alloying: Understanding the Basics*, ASM International, USA.
8. Vargel, C. (2004). *Corrosion of Aluminium*. Elsevier Ltd., New York.
9. Schütze, M., Wieser, D., Bender, R. (2010). *Corrosion Resistance of Aluminium and Aluminium Alloys*, 636. Wiley-VCH Verlag GmbH, Weinheim.
10. Yaroslavtseva, O., Ostanina, T. N., Rudoy, V.M., Murashova, I.B. (2015). *Corrosion*

- and metallic protection. Ural University publishing house, Yekaterinburg.
11. Ghali, E. (2000). Aluminum and Aluminum Alloys. In Revie, R. W. (ed.), *Uhlig's Corrosion Handbook, Second Edition*, 677-715. Wiley, Hoboken.
 12. Cicolin, D., Trueba, M., Trasattin S. (2013). Effect of chloride concentration, pH and dissolved oxygen, on the repassivation

of 6082-T6 Al alloy. *Electrochimica Acta*, 124,27–35.

Department of Mathematics, Physics,
Chemistry
Technical University–Sofia, Branch Plovdiv
25 Tsanko Dystabanov St.
4000 Plovdiv
BULGARIA
E-mail:kalina0506@gmail.com

VARIATIONAL PRINCIPLE FOR A CLASS OF NONLOCAL BOUNDARY VALUE PROBLEMS

GEORGI P. PASKALEV

Abstract. For considered nonlocal boundary value problem for hyperbolic-parabolic type PDE an implicit symmetrizing operator [4] is build. Equivalence of the problem to the problem of minimization of quadratic functional is proved. Existence and uniqueness of the generalized solution are obtained.

Key words: hyperbolic-parabolic type equation, nonlocal problem, variational principle, negative norm of Lax.

1. Introduction

Let $m \geq 1$, $x = (x_1, \dots, x_m)$ and $D \subset R^m$

be a bounded domain with a boundary ∂D and $G = D \times (0, T)$, $\Gamma = \partial D \times (0, T)$.

$M[u] = \sum_{|\alpha|, |\beta|=1} a_{\alpha\beta}(x) D_x^\alpha D_x^\beta u$ is a strong elliptic

operator, where $a_{\alpha\beta}(x) \in C^\infty(\bar{D})$

$$a_{\alpha\beta}(x) = a_{\beta\alpha}(x) \forall x \in \bar{G} \forall \alpha, \beta: |\alpha| = |\beta| = 1.$$

Let we have

$$k(t, x) \leq 0 \quad \forall (t, x) \in G,$$

$$k(T, x) = k(0, x) < 0 \quad \forall x \in \bar{D},$$

$$b(T, x) = b(0, x) \quad \forall x \in \bar{D}, \quad C = \text{const.} > 0.$$

Consider the equation

$$L[u] = f(t, x), \quad (1)$$

where

$$L[u] \equiv k(t, x) D_t^2 u + M[u] + b(t, x) D_t u - Cu.$$

We shall propose also that

$$2b(t, x) - D_t k(t, x) > 0 \quad \forall (t, x) \in \bar{G}.$$

The equation (1) is a of hyperbolic-parabolic type in D and on the bottoms of the cylindrical domain D it is of hyperbolic type.

Let $\lambda \neq 0$, $|\lambda| < 1$ is a real number. Consider the following boundary value problem. To find a solution to the equation (1) in G, satisfying the boundary conditions:

$$u|_\Gamma = 0, \quad D_t^i u(T, x) = \lambda D_t^i u(0, x), \quad i = 0, 1. \quad (2)$$

2. Function spaces and definition

Let $\tilde{C}^\infty(\bar{G})$ be the space of infinitely smooth in $C^\infty(\bar{G})$ functions, satisfying the boundary conditions (2) and $\tilde{C}_*^\infty(\bar{G})$ is the corresponding space of infinitely smooth in \bar{G} functions, satisfying the boundary conditions ad joint to (2):

$$v|_\Gamma = 0, \quad D_t^i v(0, x) = \lambda D_t^i v(T, x), \quad i = 0, 1. \quad (3)$$

Define $H^1(G)$ as the closure of $\tilde{C}^\infty(\bar{G})$ with respect to the norm

$$\|v\|_1^2 = \int_G \sum_{q+|\alpha|\leq 1} (D_t^q D_x^\alpha u)^2 dt dx. \quad (4)$$

Let $H^{1,*}(G)$ is the closure of the space $\tilde{C}_*^\infty(\bar{G})$ with respect to the norm (4). The scalar product in $L_2(G)$ we shall denote by $(\cdot, \cdot)_0$ and the corresponding norm - by $\|\cdot\|_0$.

Let the constant κ is defined by the following equality:

$$e^{\kappa T} = \frac{1}{\lambda^2}.$$

For any element $v \in H^{1,*}(G)$ and if Δ is the Laplace operator, denote by Kv the solution of the following nonlocal problem (1) with boundary conditions (3):

$$L^*(Kv) = e^{\kappa t} [v - \Delta v],$$

$$Kv|_{\Gamma} = 0, D_t^i Kv(0, x) = \lambda D_t^i Kv(T, x), i = 0, 1.$$

From the a priori estimates, near to the obtained in the paper [1]:

$$(Lu, e^{\kappa t} D_t u)_0 \geq C_1 \|u\|_0^2 \quad \forall u \in \tilde{C}^\infty(\bar{G}),$$

$$(L^* v, -e^{-\kappa t} D_t v)_0 \geq C_2 \|v\|_0^2 \quad \forall v \in \tilde{C}_*^\infty(\bar{G}),$$

where C_1, C_2 are positive constants, it follows that if the constant C is sufficiently large, then the solution Kv exists and is unique.

By integration by parts, now, for any functions $u, v \in \tilde{C}^\infty(\bar{G})$, it obtains that

$$(Lu, Kv)_0 = [u, v]_1^\# \quad \forall u, v \in \tilde{C}^\infty(\bar{G}), \quad (5)$$

where

$$\|u\|_1^{\#2} = \int_G e^{\kappa t} [(\gamma - \kappa^2) D_t^2 u + \sum_{|\alpha|\leq 1} (D_x^\alpha u)^2] dt dx,$$

is a norm, equivalent to (4) and such that

$$\gamma = \text{const.} > \kappa^2.$$

Changing the places of u and v , we have

$$(Lv, Ku)_0 = [u, v]_1^\# \quad \forall u, v \in \tilde{C}^\infty(\bar{G}). \quad (6)$$

From (5) and (6) follows that the operator L is

K -symmetric [2].

Now for $u \equiv v$ we have

$$(Lu, Ku)_0 = \|u\|_1^{\#2} \quad \forall u \in \tilde{C}^\infty(\bar{G}),$$

which means that L also is K -positive

operator [2].

Definition: The function $u \in \tilde{H}^1(G)$ is called a generalized solution to the problem (1),(2) if

$$(f, Kv)_0 = [u, v]_1^\# \quad \forall v \in \tilde{H}^1(G).$$

3.Main results

Theorem 1: The function $u \in \tilde{H}^1(G)$ is a generalized solution to the problem (1),(2) if and only if u realizes the minimum of the quadratic functional

$$D[u] = \|u\|_1^{\#2} - 2(f, Ku)_0 \quad (7)$$

in the space $\tilde{H}^1(G)$.

Theorem 2: For each function $f \in L(G)$ there exists a unique solution to the variational problem to minimize the quadratic functional (7) in the space $\tilde{H}^1(G)$. The problem (7) is correctly posed - to the small variance of f in $L^2(G)$ corresponds a small variance of the solution in $\tilde{H}^1(G)$.

4. Proofs

To prove the above theorems we need to obtain the following estimate.

Lemma: There exists a positive constant $C_4 > 0$, such that

$$\|Kv\|_0 \leq C_4 \|v\|_1 \quad \forall v \in \tilde{C}^\infty(\bar{G}).$$

Proof.

Since $Kv \in \tilde{H}^1(G)$, using the a priori

estimate from [1] we have

$$\|Kv\|_0 \leq C_3 \|L^*(Kv)\|_{-1} = C_3 \|e^{x^2} [v - \Delta v]\|_{-1} \leq$$

$$C_3 \|v - \Delta v\|_{-1} = C_4 \sup_u \frac{|(v - \Delta v, u)_0|}{\|u\|_1} =$$

$$C_4 \sup_u \frac{|(v, u)_1|}{\|u\|_1} \leq C_4 \|v\|_1,$$

where $C_3 = \text{const.} > 0$, $\|L^*[Kv]\|_{-1}$ is a negative norm of Lax and L^* is the corresponding formally ad joint operator to L .

The proofs follow the standard scheme, used in the papers [2-7].

As is well known, the presented approach is classical for elliptic partial differential equations. For some classes non elliptic partial differential equations with constant coefficients the above variational approach is presented in the papers [3-

7]. For such equations we can use the Fourier transform to build the symmetrizing operator and to obtain the corresponding estimate. From the other hand for equations with variable coefficients it is necessary to build a symmetrizing operator in implicit form – for example as a solution to the ad joint problem with a “special” right part.

REFERENCES

1. **Karatoprakliev, G.D.** (1987) Nonlocal boundary value problems for a class of mixed type equations, (Russ.), *Diff. Uravnenia*, Minsk, V.23, no.1, P.78-84.
2. **Shalov, V.M.** (1963) A solving of non self ad joint equations by the use the variational method. (Russ.), *Dokl. A. N. SSSR*, V.151, no.3, P.511-512.
3. **Shalov, V.M.** (1965) A principle of the minimum of quadratic functional. (Russ.), *Diff. Uravnenia*, Minsk, V.1, no.10, P.1338-1365.
4. **Paskalev, G.P.** (2017) Variational method for a class of higher order hyperbolic equations. *Journal of the Technical University at Plovdiv*, “Fundamental Sciences and Applications”, V.23, P.63-67.
5. **Paskalev, G.P.** (2013) Shalov’s variational method for the multidimensional wave equation. *Journal of the Technical University at Plovdiv*, “Fundamental Sciences and Applications”, V.19, P.141-145.
6. **Paskalev, G.P.** (1999) Variational method for the multidimensional heat equation. *Journal of the Technical University at Plovdiv* “Fundamental Sciences and Applications”, V.8, P.69-78.
7. **Paskalev, G.P.** (1992) On the investigation of boundary value problem for ultraparabolic equation with constant coefficients using variational method. (Russ.), *Diff. Uravnenia*, Minsk, V.28, no.9, P.1640-1641.

Department of Mathematics, Physics and Chemistry
 Technical University-Sofia, Plovdiv Branch
 25 Tsanko Dyustabanov Str
 4000 Plovdiv BULGARIA
 e-mail: g.p.paskalev@abv.bg

