# Central European Researchers Journal

CERES

# A Word of Welcome from the Editors

Dear Colleagues, Readers and Authors,

WThe members of publishing team of the journal CERES are glad to present to you next number. We would like to recall that the Central European Researchers Journal has been created under the CERES project titled as "Centers of Excellence for young RESearchers" (ref.no.: 544137-TEMPUS-1-2013-1-SK-TEMPUS-JPHES). This fact causes the journal politic and publications specifics and context. The journal is established as the opportunity facility to support of investigations of young researchers and PhD student. It is an essential journal's mission. The project is finished a year ago and the journal was in new conditions. Of course, the mission of the journal will be saved and continued. But new environment and conditions for the journal development cause reorganization of the publishing team. The Editorial Board will be reorganized at the first step. The journal scope will be modified too. We hope that this reorganization will improve the journal popularisation and its status.

We thank all authors and readers who were with the journal CERES. Your help is a journal sustainability. We thank all reviewers. Your work was appreciated and useful for the journal formation. The list of reviewers is published for every year in the second number of journal (you can find it in the firs separated file of number under title "Front Cover/ Editorial Board/ Table of Contents"). Our special thanks to the technical editor Dr. Jozef Kostolny from University of Zilina (Slovakia). Dr. Jozef Kostolny created the journal design. He been elaborated and maintained the technical platform, prepared and uploaded all journal's numbers.

We believe that new publishing team will improve the journal and all authors and readers will stay with the journal.

With best wishes

*Prof. Vyacheslav Kharchenko*
*Prof. Elena Zaitseva*

# CONTENTS

# Trends of the development of internationalization and integration in the World educational space

U.A. Beizerau

**Abstract** – The article deals with the study of the  processes of integration and internationalization in the systems of higher education of the countries.  The main trends in the World educational space are discussed. The problems of creating unified educational space in Europe are analyzed. The research speculates on the topic of creating common educational environment in EU and the problems of facing  competition from other developed countries. Much attention is paid to the analyses of the development of the World market of higher education. Qualitative and quantitative analyses is produced.

**Keywords** –  integration; the system of education; educational space; globalization; internationalization

## I. INTRODUCTION

The interest shown by  the World and European community to the problems of integration in education is explained by the fact that in the conditions of intensive development of science and transition to information society strengthening of attention to the trends causing processes of unification of components of educational systems is required. In reference books the trend is defined as the direction of development [9,  p. 793]. Therefore, trends of development of integration processes in education are the directions of development, a way of integration of national education systems to uniform educational space.

The expansion of democracy happening almost in all countries and strengthening of the constitutional state raise the role in education of youth and adults in the spirit of democratic civic consciousness. Higher education broadens  the idea of  academic freedom and equality with growth of  higher education in society, universalizing  political system in  different countries. Formation of the new sociocultural values  shared  by the majority of the countries entering into the world community such as civilized free market and humanization of the public relations not only changes structure of education, expanding training of economists, managers, humanists (lawyers, sociologists, political  scientists), but also change the paradigm. These are the calls of  modern era, feature of  international life, generating main trends in education.

The major trend connected with the  increase of the role of science in production and society is the massification  of education. Education becomes massive on a large scale. In different countries the level of arrival of graduates of schools in higher educational institutions averages nearly 60%, and in North America – 84%. There is a prompt proliferation of students of higher education institutions. If in 1960 the number of students in the world according to UNESCO was 13 million, then in 1997 it  increased almost by 7 times and was 88,2 million. In 2003 the number of students in the World  exceeded 122 million people, according to forecasts in 2025 the number of students will be 150 million [14; 15].

Other important trend which develops especially dynamically from the second half of  the 20th century is the diversification of education in institutional forms, levels and contents. In the conditions of growth of variety of the educational institutions giving knowledge and skills in the sphere of brainwork, the role of classical universities, however, not only decreased, but increased. They try to play role  of the centers for formation of the sociocultural environment in the region. From year to year the number of students, researchers and teachers who work, live and communicate in the international environment grows.

U.A. Beizerau F. Scorina Gomel State University,  Gomel, Belarus (e-mail: bejzerov@gsu.by)

The financial analysis of the international education market shows that his total gross revenues from his activity exceed 100 bln. dollars, the number of the foreign students who are annually coming to other countries according to educational programs and training exceeds 5 million people. Statistical data of the Organization for Economic Cooperation and Development (OECD) show the long-term trend characterizing increase in number of foreign students. For the last 30 years the number of foreign students has increased almost by 5 times, from 800 thousand people in 1975 up to 4,1 million people in 2015. [15].

## II. INTERNATIONAL COOPERATION IN HIGHER EDUCATION

International cooperation changes the forms and kinds of activity, accumulating the potential for the solution of a triune task: achievements of such education level which would correspond to the needs of modern international society; alignments of level of National educational systems; training of qualified personnel for national economy. The role and value of the international organizations, funds and programs in the field of education and science increases in these conditions.

In the European Union within several decades carries out  policy, first of all in the field of higher education, supranational institutes of coordination and management are formed. 7 conventions on mutual recognition of documents on the termination of  educational institutions, training courses and the periods of training, diplomas of higher education, academic degrees creating a standard basis of integration process in the sphere of the higher education of the EU [331-333] were  prepared and adopted. Broad development was gained by bilateral and multilateral scientific and pedagogical cooperation of the universities, exchange of teachers and students,  with the  assistance of the created target supranational programs of the EU (COMETT, ERASMUS, LINGUA, SOCRATES, etc.).

In the last decades   international cooperation in the sphere of  higher education also Republic of Belarus has intensified
dramatically. In 1998 - 2005 Belarus ratified the fundamental conventions in the field of education adopted under the auspices of UNESCO and the Council of Europe. A significant amount of bilateral intergovernmental contracts on cooperation with the universities   of foreign countries were  signed, the number of direct contracts with the foreign universities has significantly increased. Preservation of national experience, traditions, consolidation and development of  its  undoubted advantages has to be an important condition of integration of the education system of Belarus into the World educational space. It is necessary to find optimal variants of consecutive integration of  the educational system of Belarus into the World education system, to preserve everything   valuable that    Belarusian education possesses, at the same time to carry out, taking into account the international experience, the changes demanded by time which will allow to provide the prospects of development of the country in new century.

The realization of the process of integration in interaction of education systems of various countries is enabled in the presence of  integrity of teaching and educational process, the uniform language environment of teaching and educational process, participation of subjects of teaching and educational process in all forms of joint activity, combination of local (national and local) and the World culture in the content of education, wide use of national educational traditions in educational process [2].

Formation on the continent of the most competitive,  dynamic, focused on science economy capable to constant, steady growth and creating the best conditions for any citizen of the European Union [5] became  new strategic objective of the European Union defined during the meeting of the Council of Europe in Lisbon on March 23 - 24, 2000. On the basis of offers of the Commissions of the European Union and offers of member countries of the Union, the Council has accepted "The report on the concrete prospects of development of

education systems" on February 12, 2001. This document became the first, which defined comprehensive and consecutive approach to national politicians in the field of education in the context of the European Union. In this document the main objectives were formulated: improvement of quality and efficiency of education systems of EU Member States; ensuring access to all education levels for all citizens; ensuring open character of the educational system [5]. The European Commission and the Council of Europe have developed some common goals and also defined the role of education in achieving of a strategic goal of integration of educational systems into uniform educational space.

1. Improvement of the quality standards of education. Education is the best method of social and cultural interaction and also a considerable social and economic asset which promotes formation of the European Union as the most competitive in economic and the most dynamically developing society in the social plan. It is necessary to improve quality of training of teachers and tutors at all levels, it is necessary to pay special attention to the development of professional skills and also adaptation of experts in new conditions. Improvement of quality of education at schools, universities and other educational institutions by the fullest, complex and rational use of the available resources, is the main prospect in future, also as well as the most widespread introduction of knowledge from technical and naturalistic areas, i.e. to pay bigger attention to study mathematics, physics, chemistry, biology and other natural sciences to provide competitiveness of Europe in future. Improvement of quality of standards of education designates also the best ratio of resources and requirements, granting an opportunity to educational institutions to develop partnership, to support their new, more significant role in the society of new type.

2. Ensuring broader and simple access to education. The European model of social integration has to be capable to give everyone an opportunity for receiving formal and informal education in any educational institutions and also to simplify the procedure of transition from one level on another, for example, to simplify Regulations of Admission in higher educational institutions, to provide the continuity of education throughout all life since the early childhood up to elderly age.

3. Ensuring open character of education system for the World community. This purpose assumes creation of the European educational space through increase in mobility, learning foreign languages at all levels of educational system, hardening of the available communications and establishment of new communications with the World systems (educations, sciences, etc.)

The strategic objective of the European Union has demanded the use of absolutely new innovative method which has the name of the Open method of coordination. This method, unlike earlier used methods, provides new framework of cooperation of member states in the line with convergence of policy of the governments and inclusion in the common, uniform big goal shared by all nations. The method is based on the following targets:

1) joint determination of results which have to be reached;

2) application of uniform techniques when maintaining statistics and uniform indicators, for more objective assessment of the situation and correction of actions;

3) implementation of joint programs and projects for stimulation of innovations and improvement of quality of education.

One more key element of the strategy is the concept of training at an extent of all life. It is important not only for increase in competitiveness at revenues to work, but also has social importance in respect of personal development.

Despite good results which have been achieved (now more than one million students take part in the ERASMUS program), one of the main obstacles for the people wishing to get education, or to work in other EU country is the fact that their qualification won't be recognized. In the EU recognition of professional qualification or education is carried out

according to the collection of directives of the European Commission which will be replaced with the only directive in which all existing professional qualifications will be reflected.

Study of new technologies becomes priority in constantly changing world. The European Commission has defined the action program under the name "Electronic Education". This program assumes acceleration of technologization and computerization of infrastructure of education at reasonable expenses for the purpose of training computer literacy as strengthening of ties between the states and people, schools and the universities at all levels – local, regional, national and European for a possible bigger number of people.

The plan approved by the European Commission in March, 2001 has provided a basis for the European cooperation. The wide range of resources of the EU from education, youth and research programs to the European Regional Fund of Development, European Social Fund and the European Investment Bank is included in this process of cooperation. The partial embodiment of this plan is the SOCRATES program. Along with the intra-European programs of cooperation in education, also programs with other countries, in particular, with Canada are carried out [5].

Every year hundreds of thousands of Europeans use an opportunity to get education or job abroad, participating in the European project, supported by the SOCRATES program. The TEMPUS program contributes to the development of cooperation with the purpose of modernization of the higher education from North Africa to Mongolia. Moreover, it allows to make European education equally open both for citizens of EU countries, and for the citizens of foreign states.

In 1976 Ministers of Education of the European countries for the first time made the decision to create information network as a basis for the best understanding of policy in the fields of education and structures of educational systems, in 9 countries of the European Community (the predecessor of the EU). This decision showed the principle of mutual respect and simultaneous cooperation in improvement of education systems, social security, etc. The European information network which has received the name "Eurydice" of the beginning to work in 1980.

In 1986 the form of exchange of information has passed into a form of exchange of students according to the ERASMUS program (a part of the SOCRATES program). This program is considered one of the most effective in the field of integration of educational systems and also the most successful initiative of the European Community. The program has been successfully distributed in 12 EU countries [12]. The got experience has been generalized and finished in the SOCRATES program which covers all fields of education at all levels. The SOCRATES program allows to finance training and researches. It is open for a large number of the public and private organizations connected with education. Management of the program is performed of the National agencies which open in all participating countries of this program.

The European space of continuous education (or educations throughout all life) will allow citizens of the EU to move freely within the continent for the purpose of employment or training in other region. The term "continuous education" covers both preschool education, and postdegree education, includes all forms of education (formal and informal, etc.)

EU member states came to the agreement on need of development of the concept and the strategy of continuous education. Blocks of this strategy which will help the states to carry out rapprochement of educational systems are developed. Transformation of traditional systems is the first step on the way to general availability of continuous education for all citizens. Besides it is necessary:

To develop partnership at all levels of public management (national, regional and local) and also between institutions of education (schools, the universities, etc.), civil society in a broad sense (business partnership, social partnership);

To define  needs of students and labor markets in the context of application new, first of all information technologies;

To search for resources by attraction of private and public investments into financing of educational projects;

To ensure  bigger availability of education, mainly by expansion of network of the local educational centers. Special efforts should be made for providing  equal educational opportunities for disabled people and representatives of ethnic and racial minorities;

To create  special culture of education for motivation of potential students;

To ensure recognition of the documents on education issued in the territory of the EU in the participating countries;

To recognize  all types of documents on education, to create corresponding mechanisms.

In the joint working program which defined the purposes for the  educational system, the European Commission suggested to accept  realistic plan directed to achievement of goals. According to this plan it was supposed: increase in investments into education; decrease in quantity of drop-outs from secondary school; increase in graduates of the universities on mathematics,; increase in population, having senior secondary education; education throughout all life.

Now the number of the young people who  left secondary school without certificate in Europe averages 15% of the total number of people of this age group. In different EU Member States this figure isn't identical. So, for example, in Sweden, Finland and Austria it makes 10,3%, and in Portugal – 45%, in Spain – 29%, Italy – 26% [5].

All EU countries had to balance the number  of graduates on mathematical and technical specialties and also increase total number of university graduates on these specialties. Now higher education institutions of the EU annually leave  550000 undergraduates  in mathematics, technical and exact sciences. (In comparison in the USA – 370000, in Japan – 240000). It is necessary to pay special attention to motivation of girls to obtaining specialties in the sphere of mathematics and other exact sciences.

### III. INTERNATIONALIZATION AND COMPETITION IN HIGHER EDUCATION

Higher education is international, much more, today, than in the Middle Ages. Only two decades ago the number of  students studying  abroad was scanty. Now, according to data of the International Finance Corporation, more than 4  million students (about 3% of total of students in the world)  studied abroad. Since the end of the 1990th years the size of the market of educational services in the field of the higher education grew for 7% a year. Annual income from payment for training made 30 billion dollars. This situation promoted strengthening of the competition between higher educational institutions for the right to have the most talented and diligent students and also subsidies. [7].

Two trends  in higher education – internationalization and the  competition – are interdependent.  More universities create conditions for convenience of potential students from abroad, and  they become more attractive  for foreigners. For example, young people from Germany seek to receive the first academic degree at the universities of Great Britain providing better training and which are less staffed with students than the universities of Germany.

The idea that the student is a consumer of educational services is new to the whole World. In Europe and other countries of the World the governments of these countries were  the main consumers of educational services within the last century. They considered that the most capable students of the country got educations on those specialties which the country needed. About 110000 students from the countries of continental Europe get higher education in Great Britain. However the largest numerical growth is observed among the Chinese students.

According to OECD, in 2015 a half of all foreign students preferred to get an education in five countries. Most foreign students in the world study in the USA - 18%, in Great Britain –

10%, Australia – 7%, in Germany – 7%, in France of 7% too. Further follow Canada - 5%, Russia – 4%, Spain – 2%. See table 1.

Table 1. The leading countries (in attracting foreign students)

| Country | Number of students in the country | Number of foreign students in the country |
|---|---|---|
| United States | 11748263 | 740482 |
| United Kingdom | 3582166 | 427686 |
| Australia | 1088366 | 249588 |
| France | 2616643 | 239344 |
| Germany | 2645504 | 206986 |
| Russia | 3070235 | 173627 |
| Japan | 3670435 | 150617 |
| Canada | 1015000 | 120960 |
| China | 52829775 | 88979 |
| Italy | 2780343 | 77732 |

As well as in any competitive area in education it is observed during competition appearance of new participants, and loss of leadership by other players too. So, for the last 10 years the share of the foreign students wishing to study in the USA  reduced from 23% to 18%. It was  caused by toughening of entrance rules for foreign students, after the terrorist events on September 11, 2001. Reduction of number of foreign students for 2% and for 1% in Great Britain is observed. At the same time Australia, New Zealand and Russia enjoys growth of number of foreign students approximately for 2%.

At the end of the XX  century-beginning of the XXI century the biggest  number of potential students went to study to the USA. About 740000 foreign students  arrive to country annually. They were attracted by excellent quality of education and grants, partially or completely covering expenses on training. However, applications for 2004/2005 for training in the USA from the Chinese students  reduced by 45% and for 30% from  Indian. It is connected mainly with toughening of the procedure of issuing of the American visa. At the same time the number of  Chinese students in Australia, for example, s grew  in the mid 200-th for 47%, and Indian for 52%. Total number of foreign students, nevertheless, reduced for 10%. France and the Netherlands  became serious participants in the market of  higher education [6; 15].

In the sphere of  the World educational space there  are also  radical changes. The British and Australian universities open the branches in China, Malaysia and the  United Arab Emirates. One of the reasons is reduction of  the cost of training. The university of Texas established relations with the  London university college, getting  local resources and base from  it. Nevertheless, there  are doubts about real profitability of these actions. Similar practice continues to extend.

Thus, the universities repeat experiment of large industrial corporations today, concentrating attention on income, advancing the services on foreign, cheaper markets. Despite all difficulties, internalization and the competition in educations – the objective process bringing positive changes.

Mass higher education forces  universities to become more and more versatile, more global and flexible. It is explained, first, by the fact that the higher education is democratized, becomes massificated, especially in the developed countries.

The proportion of adult population with the higher education in the developed countries almost  doubled since 1975  from 22% to 45%. Secondly, national economies and the  World economy become more and more knowledge-intensive, i.e. demands highly qualified personnel. The World is in power of "Smart revolution" where knowledge, replacing physical resources, becomes the chief conductor of economic growth. Statistically, during the period from 1985 to 2000 the contribution of high-tech and knowledge-intensive industries to

economy has increased from 51 to 59% in Germany, from 45 to 51% in Great Britain. The largest companies invested up to 1/3 means in science. The universities, thus, became "engines" of new, hi-tech economy. Thirdly, emergence of processes of globalization abolished barriers between the countries of the World, allowed to transform higher educational institutions also considerably, as well as business. The number of people, the citizens of the developed countries, getting education abroad doubled for the last 20 years and has reached more than 2 million. The universities open branches worldwide, and many countries do higher education one of the income items of the economy, exporting it [10]. Fourthly, increase in the competition has led to the fact that the traditional universities are forced to compete for attracting new students and grants on scientific research. According to the World Bank, annual expenses for education in the World make 300 billion dollars a year that is equal to 1% of the world GDP [15].

American higher education is the absolute leader in the World. 17 of 20 best universities of the world are in the USA. American universities take in this list the majority of places – 35 that makes 70% of the total number of higher education institutions. Now 70% of all Nobel laureates work at the universities of the country. In 2001 they published 30% of articles in natural-sciences and technical disciplines and 44% of articles in humanitarian disciplines [4].

European countries spend for education only 1,1% of their GDP while in the USA this figure exceeds 2,7%. American universities spend for each student 2,5 times more than European. It affects smaller load of audience, presence of the best professors and high quality of the conducted researches at the American universities. According to the European Commission about 400000 scientists – citizens of EU countries teach now and conduct researches in the USA.

In Europe four times less inventions are patented, than in the USA. Thus, in Europe there is no question how to catch up with the USA in the field of quality and efficiency of the higher education, but the question is how not to lag behind rates of development of education in China and other Asian regions. The main problems of the universities are identical in all European Union – complete control of the state, lack of independence in selection of students and teachers, impossibility to pay the work of professors in worthy volume that makes most of the European universities noncompetitive in the World market of the higher education. See figure 1.
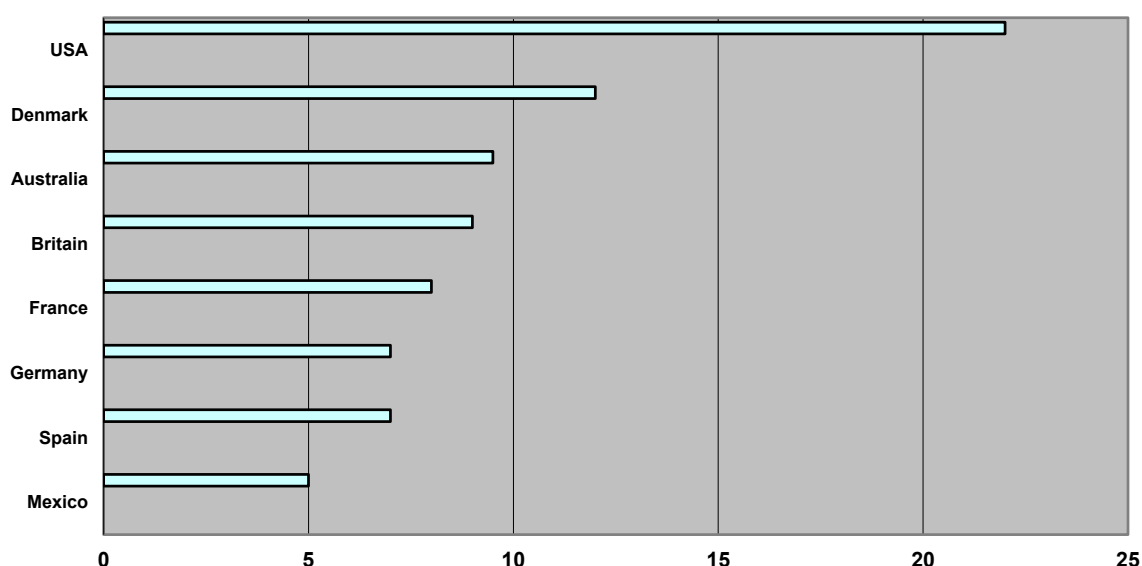


Fig. 1 Expenses for one student in thousands dollars per year

Russian and Belarusian higher education institutions lag behind the foreign competitors so far. For example, the most prestigious higher education institution of Russia – Moscow State University was only in the 77th place in the World ranking of higher education institutions, published by the university of Shanghai.

In the Republic of Belarus an undisputed leader and the most prestigious higher education institution is Belarussian State University. Training cost in higher education institutions of the country varies from 800 to 1100 US dollars on a preparatory course, from 1500 to 2500 dollars a year on basic courses of liberal arts, technical, economic colleges, 2400 – 3600 dollars on basic courses of medical schools, 2500 – 4000 dollars for a year of training in a postgraduate study [8].

In the developing countries universities become available for a wide range of the population. The governments of these countries expand network of higher educational institutions. The main reason of spasmodic development of higher education in Third World countries – increase in demand for qualified specialists in the developing economy of regions. Also the number of the countries where secondary education becomes obligatory grows.

In the field of the higher education, also as well as in many other areas, China and India have outstripped other developing states. If in the early eighties only 2–3% of graduates of schools continued training at the universities, then in 2003 this figure has made 17%. In the sphere of post-degree education in the period of 1999-2003 12 times more theses were defended, than during the period from 1982 to 1989. The number of graduate students almost tripled.

In China communication between universities and industry is stronger, than in any of the countries. During the period from 1992 to 2003 the majority of dissertations was defended in applied disciplines, 38% – in technical disciplines, 22% – in natural sciences and 15% – in medicine.

In the mid 2000-th foreigners made 30% of doctoral candidates in the USA and 38% of doctoral candidates in Great Britain. The number of foreign students studying for bachelor degree (Undergraduate) in the USA – increases by 8%, in Great Britain for 10% annually. Today in the USA 20% of employed professors are foreigners.

Now thanks to cooperation in the field of education between the countries and integration of national education systems into world educational space the number of the students who study abroad grows annually. This growth was promoted by opportunities and desire of students to study at the best universities of the World, 50% from them are citizens of developing countries. The championship in in export of students among developing countries belongs to China (10% of number of all foreign students) and India (4%) [1].

The majority of the developed countries try to attract talented young people with the help of different methods. The University of Oxford increased a set of foreign students twice that now makes 15% of all students. The mobility of students is promoted also by the measures of its encouragement taken by the European Union [3].

Since 1950th the USA dominates on the international market of higher education. Now about 740000 foreign students study in the country. It gives sure and indirect gains: first, payment for training by foreign students brings to America nearly 15 billion dollars annually, secondly, the American economy receives highly qualified specialists [15].

## IV. CONCLUSION

Growth of a number of foreign students is connected with the following factors: first, the strategy led by the states on preservation and development of political and social connections between the countries. Especially in connection with the creation of European educational space. Secondly, significant increase in number of the students. Thirdly, democratization of the prices in the conditions of globalization in the sphere of transport connection, etc.

Thus, educational migration today is one of the most dynamically developing types of the territorial movement of the population. Education in modern conditions is the most difficult phenomenon getting into all spheres of public life. The transformations which are carried out in the field of higher education define a lot of things in each country: standard of living of the nation, welfare of the population, national security of the state. Intellectual potential of the state is the main factor and resource of prosperity of each country. Training of professionals always was and will be an important strategic problem of the state. Higher educational institutions turn into the centers of gravity of youth and create strong conditions for creation of steady educational migration.

## REFERENCES

[1] Beizerau, U. A. Trends of the development of the higher education within integration of educational systems into uniform educational space, Current problems of modern humanitarian education: materials of the 3rd republican scientific conference of young scientists and graduate students, Minsk, November 29, Minsk: RIVSh, 2006, pp . 278 – 285.
[2] Beizerau, U.A. Models of integration of the Belarusian system of education, Gomel, GSU, 2007, 108 p.
[3] Clever stuff. Brittania redux, The Economist, (8514), 2007, pp. 9 –11.
[4] China chipped, The Economist, (8409), 2005, pp. 38-40.
[5] EU cooperation in higher education, European Commission, Brussels, 2003.
[6] Free degrees to fly, The Economist, (8415), 2005, pp. 63 – 65.
[7] Higher Ed Inc, The Economist, (8443), 2005, pp. 19 – 20.
[8] How much does the higher education cost. Interfax, 2017.
[9] Ozegov, S.I. Dictionary of Russian language, Moscow, ITI of Technology, 2005, 944 p.
[10] Scholars for dollars, The Economist, (8405), 2004, pp. 55-56.
[11] Structure of the Educational and Initial Training Systems in the European Union, Luxembourg, 1995.
[12] The Amsterdam Treaty, 1999, Art. 149, Pt.1.
[13] The Amsterdam Treaty, 1999, Art. 150, Pt.4.
[14] World education report, UNESKO publishing, 2000.
[15] World education report. To knowledge society, UNESKO publishing, 2016.

# Tool for topological reliability analysis of reversible logic circuits

Peter Sedlacek

***Abstract***—It is assumed that progress with classical logic gates development can have its boundaries in the near future. Therefore, it is needed to make research with other approaches. One of them is application of reversible logic. In this paper, we present a method for reliability analysis of logic circuits composed of reversible logic gates. The method is based on structure function, and we implemented it in a software tool that also allows us to design reversible logic gates and circuits.

***Keywords*** – reversible logic, logic circuit, reliability analysis

## I. INTRODUCTION

Reversible logic circuits are getting to the foreground in presence, as it is assumed that progress with present irreversible logic gates can have its limits in the near future [1].

The computing devices loss part of information in the process of computing. Irreversible logic functions are always connected with physical non-reimbursement, and so, there is a minimal amount of heat generated by the device for every irreversible function (Fig. 1). This disappearing serves for signals to become independent of their history. But it is possible for computing to run without this loss of information. This can be ensured that way, calculations will be reversible, what means, it is possible to return from any state to any previous. Landauer pointed out, that majority of physical laws are reversible, so if you have the full information about state of closed system in some time, it is possible, at least in principle, to perform these laws in reverse order and get exact state of system in any previous time [2].



Fig. 1 Landauer's barrier and trend in energy consumption [3]

In 1973 Charles Bennett in [4] showed that it is possible to build full reversible computer capable to perform any calculations without needing of overwhelm memory with temporary data. He reverted operations, that produce temporary result. Reversible computer designed this way would be in principle capable to cross Landauer's barrier, but design of a more complex reversible computer is demanding.

An important aspect when designing a reversible logic circuits is its reliability, which is the main subject of this paper.

P. Sedlacek, Faculty of Management science and Informatics, University of Zilina, Slovakia (e-mail: Peter.Sedlacek@st.fri.uniza.sk)

## II. RELIABILITY ANALYSIS

### A. Structure Function

The structure function is used for mathematic description of system and reflects system operation depending on its components and theirs states [5]. This function is defined as a map of the following form:

$$\phi(x_1, x_2, \ldots, x_n) = \phi(\boldsymbol{x}) \colon \{0, 1\}^n \to \{0, 1\}, \qquad (1)$$

where $n$ is a number of components in the system, $x_i$ is a state of component $i$, for $i = 1, 2, \ldots, n$, and $\boldsymbol{x} = (x_1, x_2, \ldots, x_n)$ is a vector of components states (state vector).

The structure function can be used to calculate system availability or perform importance analysis, which deals with evaluation of influence of system components on the system [5]. However, the key issue is how to find or construct the structure function of a system.

### B. Structure Function of Logic Circuit

Reliability analysis of logic circuits based on structure function has been investigated in [6]. The authors of the paper proposed a method for finding the structure function of a logic circuit. They assumed that the relevant components of a circuit are logic gates that the circuit is composed of. For a specific combination of input signals (e.g., (1,1) in case of a gate with two inputs), a gate (e.g., AND with two inputs) is working if its output agrees with the expected output (1 in this specific case). If the output does not agree with the expected one, then the gate is considered to be non-working for the specific combination of the input signals. Similarly, the logic circuit is functioning for a specific combination of input signals if its output agrees with the expected one. Otherwise, it is failed. This indicates that the functionality (and so the structure function) of a logic circuit depends not only on states of its components (logic gates) but also on values of the input signals of the circuit.

To show how the structure function of a logic circuit looks like, let us consider a logic circuit composed of $n$ logic gates, and let us assume that the circuit has $k$ inputs and $m$ outputs. The expected output of the circuit is defined by a vector-valued function of the following form:

$$\boldsymbol{F}(y_1, y_2, \ldots, y_k) = \boldsymbol{F}(\boldsymbol{y}) = \big(F_1(\boldsymbol{y}), F_2(\boldsymbol{y}), \ldots, F_m(\boldsymbol{y})\big) \colon \{0,1\}^k \to \{0,1\}^m, \qquad (2)$$

where $y_l$ defines value of the $l$-th input signal, for $l = 1, 2, \ldots, k$, $\boldsymbol{y} = (y_1, y_2, \ldots, y_k)$ is a vector of input signals (input vector), $F_t(\boldsymbol{y})$ agrees with output $t$, for $t = 1, 2, \ldots, m$, and $\boldsymbol{F}(\boldsymbol{y}) = \big(F_1(\boldsymbol{y}), F_2(\boldsymbol{y}), \ldots, F_m(\boldsymbol{y})\big)$ represents a vector of functions defining values of the output signals. This vector-valued function is created based on the logic gates and their expected behavior. However, if gates are unreliable, i.e., they can fail and generate an output that does not agree with the expected one, the output of the circuit depends also on states of the gates. In this case, the output of the circuit agrees with the following map:

$$\boldsymbol{F}_o(\boldsymbol{y}; \boldsymbol{x}) = \boldsymbol{F}(\boldsymbol{y}) = \big(F_{1,o}(\boldsymbol{y}; \boldsymbol{x}), F_{2,o}(\boldsymbol{y}; \boldsymbol{x}), \ldots, F_{m,o}(\boldsymbol{y}; \boldsymbol{x})\big) \colon \{0,1\}^{k+n} \to \{0,1\}^m, \qquad (3)$$

where $\boldsymbol{x} = (x_1, x_2, \ldots, x_n)$ is a state vector defining states of the logic gates, $F_{t,o}(\boldsymbol{y}; \boldsymbol{x})$ agrees with real output $t$, for $t = 1, 2, \ldots, m$, and $\boldsymbol{F}_o(\boldsymbol{y}; \boldsymbol{x}) = \big(F_{1,o}(\boldsymbol{y}; \boldsymbol{x}), F_{2,o}(\boldsymbol{y}; \boldsymbol{x}), \ldots, F_{m,o}(\boldsymbol{y}; \boldsymbol{x})\big)$ represents a vector of functions defining real values of the output signals.

Based on (2) and (3), the structure function of a logic circuit has been defined in [6] as the equivalence of the real outputs of the circuit (which assumes that the gates are unreliable) and the expected ones (which assumes that the gates are perfectly reliable), i.e.:

$$\phi(\boldsymbol{x}; \boldsymbol{y}) = \boldsymbol{F}_o(\boldsymbol{y}; \boldsymbol{x}) \leftrightarrow \boldsymbol{F}(\boldsymbol{y})$$
$$= \{F_{1,o}(\boldsymbol{y}; \boldsymbol{x}) \leftrightarrow F_1(\boldsymbol{y})\} \wedge \{F_{2,o}(\boldsymbol{y}; \boldsymbol{x}) \leftrightarrow F_2(\boldsymbol{y})\} \wedge \ldots \wedge \{F_{m,o}(\boldsymbol{y}; \boldsymbol{x}) \leftrightarrow F_m(\boldsymbol{y})\}, \qquad (4)$$

where ↔ is logical biconditional and ∧ logical conjunction. So, the structure function of a classic logic circuit can be obtained as the conjunction of the Boolean functions defining the situations in which the output signals of the circuit agrees with the expected ones. Please note that this function is quite different from (1) since its output depends not only on states of the components (state vector $x$) but also on the environment, which is included in input vector $y$.

### C. Structure Function of Reversible Logic Circuit

Reliability analysis of reversible logic circuit can be done in similar way as in the case of irreversible logic circuit presented above. In general, an irreversible logic gate performs one logic operation. Output of logic gate can be interpreted as value 0 or 1. On the other hand, reversible logic gate performs multiple logic functions. Output of such gate is now vector of values. As failure state we consider when output vector is different from the expected output vector at least in one element.

The other difference is that reversible logic circuits contains garbage bits and constant inputs. Garbage bit is an output of logic circuit that is not connected, i.e., we are not interested whether this output agrees with the expected one or not. On the other hand, a constant input is an input set to a specific value, so we do not have to verify cases when this input is set to other value. So, the set of output vectors we have to explore is smaller. This implies that the structure function of a reversible logic circuit can be defined as follows:

$$\phi(x; y_r) = F_o(y; x) \leftrightarrow F(y) = \wedge_{t \in \mathbf{R}}\{F_{t,o}(y; x) \leftrightarrow F_t(y)\}, \tag{5}$$

where $y_r$ is a vector of input signals from which the constant signals are excluded and $\mathbf{R}$ is a set of the output signals from which the garbage bits are excluded. This formula implies if a failure of a logic gate of the circuit causes that just garbage bits are different from the expected values, then this situation will not be considered as a failure of the reversible circuit.

### D. Example

For illustration of obtaining the structure function of a reversible logic circuit, let us consider an example of a full adder from [7], which is composed of two Modified Islam Gates (MIGs) and one Controlled Operation Gate (COG) introduced in [8]. This circuit is depicted in Fig. 2. The block diagrams defining operation of MIG and COG are shown in Fig. 3, where symbol ⊕ denotes logic operation XOR and symbol ′ agrees with logic complement, i.e., NOT. As we mentioned before, an unreliable gate performs different logic function as working one. For our example, let us assume, the unreliable gate performs function defined by following formula: $F(y) = 0$. That means output of a broken gate is always $0$. We can notice that the circuit in Fig. 2 has 2 constant inputs (denoted by symbol 0) and 4 garbage outputs (denoted as $g_1$, $g_2$, $g_3$, and $g_4$). They are not interesting for reliability analysis.



Fig. 2 Full adder using MIG and COG gates (according to [7])

(a) Modified Islam gate          (b) Controlled operation gate

Fig. 3 Block diagrams of reversible logic gates (according to [7])

To find structure function of the circuit, we build two truth tables. The first one contains all combinations of input values and their respective output values (Table 1). This table describes situation when all components are functioning, i.e., it defines the expected outputs of the circuit for individual combinations of the input signals. The second table is extended by all combinations of logic gate states and respective output values. This one describes real output values of logic circuit when a specific gate is in function or failure state (Table 2). We can notice that none of these tables contain columns for constant inputs and garbage outputs, as we mentioned before.

Table 1 Selected part of truth table for full adder in Fig. 2

| $A$ | $B$ | $C_{in-}B_{in}|$ | $Ctrl$ | $S/D$ | $C/B$ |
|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 0 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 1 | 1 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 1 | 1 | 1 |
| 0 | 1 | 1 | 0 | 0 | 1 |
| 0 | 1 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 1 | 0 | 0 | 1 | 1 | 1 |
| 1 | 0 | 1 | 0 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 | 0 |

Table 2 Selected part of truth table including failure states for full adder in Fig. 2

| $A$ | $B$ | $C_{in-}B_{in}|$ | $Ctrl$ | $x_1$ | $x_2$ | $x_3$ | $S/D$ | $C/B$ | $\phi$ |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 |
| 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 |
| 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 |
| 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 |
| 0 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 0 |

Based on the tables introduced above, we can find the structure function of the circuit. This can be done by comparing the expected values of outputs (from Table 1) with their real values (from Table 2). If these values match, the structure function will acquire value 1, otherwise, it will take value 0. For illustration, we placed values of the structure function of the circuit into the last column of Table 2.

Based on the structure function, we can compute several reliability measures. One of them is topological availability, which was introduced in [9] as the relative frequency of system state 1. It can be computed simply as a proportion of states at which the structure function takes value 1. In our case, it is computed as value $56/128$ because the structure function of the circuit is defined with respect to 7 variables ($A$, $B$, $C_{in} - B_{in}|$, $Ctrl$, $x_1$, $x_2$, and $x_3$), i.e., it is defined at $2^7 = 128$ different points, and it takes value 1 in 56 out of these points.

### III. IMPLEMENTATION OF TOOL FOR RELIABILITY ANALYSIS OF REVERSIBLE LOGIC CIRCUITS

We implemented the method for obtaining the structure function of a reversible logic circuit that was described above in a form of a software tool. We set several goals that our tool should do. Firstly, we should be able to create and modify reversible logic gates. Then we should be able to create reversible logic circuits from these gates and, finally, we should be able to find the structure function of the circuit and, based on it, compute its availability. We decided to implement it in C# language.

#### A. Model-View-View-Model (MVVM) Architecture

We decided to implement our tool using MVVM architecture [10] depicted in Fig. 4 as it allows relatively high independence of individual layers, e.g., a complete redesign of user interface will not require change in model and vice versa. As we can see in Fig. 4, this architecture is composed of the following three layers:

1) *Model*, which is a representation of business objects. It is optimized for logic relationships and operations between individual entities without its representation in user interface. The Model communicates with View-Model layer via events.

2) *View*, which is a representation of user interface. It displays information for users and responses to user inputs sends to View-Model layer.

3) *View-Model*, which is a bridge between Model layer and View layer. Every class of the View has its corresponding class in View-Model. This layer obtains information from the Model and processes it into a form suitable for the View. Simultaneously, it responds on requests from the View by updating the Model.



Fig. 4 Diagram of MVVM architecture [10]

*B. User Interface*

Based on the requirements specified above, our tool has three main components that are:

1) *editor for creating logic circuits* (Fig. 5) – this component allows us to build a logic circuit from the gates, which are located in the left part of the editor. We are able to add inputs in logic circuits and connect individual components to each other. We are also able to turn on/off inputs and see values of the output signals.



Fig. 5 Editor for creating logic circuits

2) *gate editor* (Fig. 6) – this component is used for creating and editing logic gates. Gates are identified by their names. A function that the gate performed is specified by truth table. We can also specify function for failure state of the gate. A logic gate created by this editor is stored into logic gate library, which is responsible for reading/saving gates into file and accessing gates upon request.



Fig. 6 Gate editor

3) *reliability analysis calculator* (Fig. 7) – it works with two truth tables as we described above. Based on these tables, it is able to find the structure function of the circuit that is created using the editor mentioned above. Based on the structure function, it allows us to calculate topological availability of the circuit.



Fig. 7 Reliability analysis calculator – table including failure states

## IV. CONCLUSION

In this paper, we proposed a method for topological reliability analysis of reversible logic circuits based on the method described in [6]. It uses the same approach, i.e., a logic circuit is considered to be working correctly if and only if all real outputs match with the expected ones for a given combination of inputs. The main difference in reliability analysis of reversible circuits is by using garbage bits and constant inputs that are not used in reliability analysis as we do not consider as failure state when real value of a garbage bit is different from the expected one and do not take into account a possibility of a change of a value of constant inputs.

We implemented the method for analysis of reversible logic circuits in our own tool that was also presented in this paper. The tool allows us to create and modify reversible logic gates, design logic circuits from created gates, find the structure function of the circuit and investigate its topological availability based on the structure function.

### ACKNOWLEDGMENT

### REFERENCES

[1] M. P. Frank: *Throwing Computing into Reverse,* IEEE Spectrum (Volume: 54, Issue: 9, September 2017)

[2] R. Landauer: *Irreversibility and Heat Generation in the Computing Process,* IBM Journal of Research and Development (Volume: 5, Issue: 3, July 1961), pp 183-191.

[3] R. Hanson: *Slowing Computer Gains,* http://www.overcomingbias.com/2013/03/slowing-computer-gains.html [accessed 2018-04-28].

[4] C. Bennett: *Logical Reversibility of Computation,* IBM, Journal of Research and Development (Volume: 17, Issue: 6, 1973), pp. 525-532.

[5] W. Kuo and X. Zhu: *Importance Measures in Reliability, Risk, and Optimization: Principles and Applications,* Wiley, Chichester, UK, 2012.

[6] M. Kvassay, E. Zaitseva, V. Levashenko, and J. Kostolny: *Reliability analysis of multiple-outputs logic circuits based on structure function approach,* IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (Volume: 36, Issue: 3, 2017) pp. 398-411.

[7] J. B. Chacko and P. Whing: *Low Delay based Full Adder / Substractor by MIG and COG Reversible Logic Gate*, Computation Intelligence and Communication Networks (October 2017).

[8] S. Mamataj, B. Das, A. Rahaman: *An Ease Implementation of 4-Bit Arithmetic Circuit for 8 Operation by using a New Reversible Cog Gate*, International Journal of Advanced Research in Electrical, Electronics and Instrumentation Engineering (Volume: 3, Issue: 1, January 2014).

[9] M. Kvassay and E. Zaitseva: *Topological Analysis of Multi-state Systems Based on Direct Partial Logic Derivatives*, in Recent Advances in Multi-state Systems Reliability, A. Lisnianski, I. Frenkel, and A. Karagrigoriou, Eds. Cham, CH: Springer International Publishing, 2018, pp. 265–281.

[10] *Implementing the Model-View-ViewModel Pattern*, https://msdn.microsoft.com/en-us/library/ff798384.aspx [accessed 2018-04-28].

# Methods of Human Reliability Analysis

Andrej Forgac

*Abstract*— The article considers the methods of Human Reliability Analysis, which can be divided into qualitative and quantitative. First, the goal and the process, that consists of a series of steps of Human Reliability Analysis is described. Then some specific qualitative methods are described and it is pointed out that qualitative methods are not sufficiently investigated.

*Keywords*— human factor, Human Reliability Analysis, Reliability Engineering

## I. INTRODUCTION

People and systems are not error-proof, and that improved reliability requires an understanding of error problems, leading to improved mitigation strategies and therefore it is necessary to follow the concept of Human Reliability Analysis (HRA).

A number of HRA techniques have been developed for use in a variety of industries. Quantitative techniques refer to human tasks and associated error rates to calculate an average error probability for a particular task. Qualitative techniques guide a group of experts through a structured discussion to develop an estimate of failure probability, given specific information and assumptions about tasks and conditions [12].

Most techniques in HRA are qualitative. It is caused by the specific of initial data for human factor analysis. This data are often incompletely specified and ambiguous. There are complications to represent this data according to mathematical models that are used in Reliability Engineering, for example as structure function, reliability block diagram, Markov model, Universal Generation Function, Petri network, etc. Therefore typical method for quantitative reliability evaluation cannot be used without special adaption for the specific of human factor analysis.

## II. HUMAN RELIABILITY ANALYSIS (HRA)

### A. Definition of HRA

Human error is an important factor to be considered in the design and risk assessment of large complex systems, especially when the human is a crucial part of the system, such as nuclear power plant operations, air traffic control, and grounding of oil tankers [1]. There is a special area in the Reliability Engineering named *Human Reliability Analysis* that investigates the influence of Human into different system. *HRA is a comprehensive and structured methodology that applies qualitative and often also quantitative methods to assess the human contribution to risk* [2]. In HRA there is already developed many methods for estimating and analyzing the behavior and impact of a human factor on system performance or its error rate. The HRA aims to identify a potential system failure resulting from human error, to analyze the causes and identify appropriate countermeasures to prevent as much as possible and reduce the associated risk. Human errors affect between 60% and 90% of all industrial and transport accidents [3].

### B. Basic conception of HRA

The HRA process consists of a series of steps that includes problem definition, task analysis, human error identification, human error representation, and human error

A. Forgac, University of Žilina, Žilina, Slovak Republic (e-mail: andrej.forgac@.fri.uniza.sk).

quantification. How each step of the HRA is conducted depends on the HRA method used and the purpose of the analysis. The results of the HRA may reveal the need for error management to reduce errors or mitigate their effects [2].
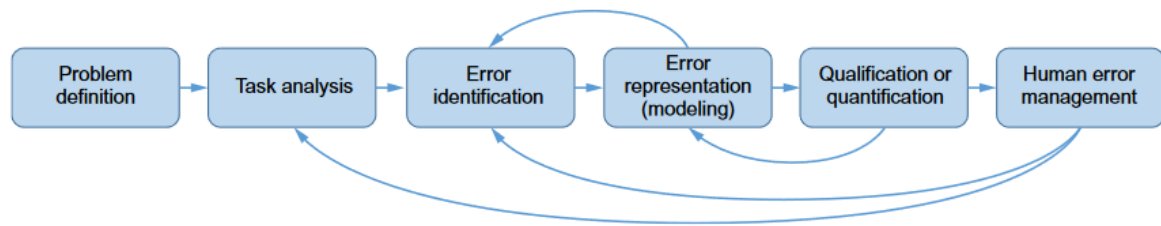


Fig. 1 Basic steps in the HRA process

*1)Problem definition*

The first step in the HRA process is used to determine the scope and type of qualitative or quantitative analysis, the tasks that will be evaluated and the human actions that will be assessed.

*Scope of the HRA* - Two factors need to be considered to determine the scope of the analysis: the purpose of the analysis and the system's vulnerability to human error. An optimal system design is less complex (less prone to error), is error tolerant (provides the capability to detect and correct errors), and allows the user flexibility to recover from failures.

*Type of HRA* – analysis of HRA can be qualitative or quantitative. In analyses as accident investigations, problem report evaluations, and general process improvements, the severity of the consequences and likelihood of occurrence are expressed for example through words like high, medium, or low, in other words qualitatively. In quantitative assessments, the consequences are expressed numerically for example the number of people potentially hurt or killed and their likelihoods of occurrence are expressed as probabilities or frequencies.

*2)Task analysis*

Second step in the HRA process is task analysis - a method that consists of systematically identifying and breaking down each task into the steps and substeps that constitute the human activities necessary to achieve a system's goal. Human activities can include both physical actions, such as installing a device; and cognitive processes, such as diagnosis, calculations, and decision making. The goal of task analysis is to decompose the functions into tasks, tasks into subtasks, and subtasks into human actions. There are over 25 variations of task analysis, each designed to accomplish different goals including task data collection, task description, simulation, behavior assessment, and task requirement evaluation.

*3)Human error identification*

The identification of human errors is third and the most important step of the HRA because failing to identify a critical human error will result in the omission of that error's contribution to risk in the HRA and to the underestimation of the overall system risk. This part of the analysis should include all of the actions that could adversely affect the system's reliability. This is done through the evaluation of the basic human actions to determine what errors can occur, and which can potentially contribute to undesired outcomes. The analyst must determine not only the types of human error that can occur, but also the factors that could contribute to the errors' occurrence.

*4)Human error representation*

Human error representation is often described as modeling because it helps illustrate the data, relationships, and conclusions that cannot be as easily described with words. Human error representation allows the analyst to look deeper and develop a better understanding of the causes, vulnerabilities, recoveries, and possible risk mitigation approaches that could be used to address accident scenarios. Tools available for human error representation include

master logic diagrams (MLD), event sequence diagrams (ESD), event trees, fault trees or generic error models, and influence diagrams.

*5)Human error qualification and quantification*

Quantification is the process used to assign probabilities to human errors. The steps in quantification depend on the method used, and the method used depends mostly on the resources available (usually time and money), the experience level of the analyst, and the available relevant data. The data must be sufficient to allow the analyst to estimate the frequency with which the errors may occur and the number of opportunities for these events. Once the human errors have been modeled and quantified, and the HRA has been completed, risk calculations can be performed to evaluate the overall system risk.

*6)Human error management*

The HRA analysts rank the errors that significantly contribute to risk in decreasing order of importance and make decisions on whether and how to manage human error appropriately. The analysts can decide to include barriers to prevent errors, provide a means to detect and correct the errors or mitigate the negative effects of the error. The human error management philosophy is founded on two basic principles: that humans, no matter how well trained and how experienced, will make mistakes; and that potential human errors can be identified and be prevented, corrected, or their effects can be mitigated. As human error is effectively managed, human reliability improves and, consequently, system reliability improves. Conversely, when human error is not adequately managed, human reliability is lower and the overall system reliability suffers [2].

## III. HRA TECHNIQUES

*A. The range and scope of HRA techniques*

There are two approaches in Reliability Analysis:

➢ Qualitative evaluation

− aims to identify, classify and rank the failure modes, or event combinations that would lead to system failures

➢ Quantitative evaluation

− aims to evaluate in terms of probabilities the attributes of dependability (reliability, availability, safety)

In broad terms, HRA consists of a qualitative phase followed by, if necessary, a quantitative phase. Different HRA methods have different approaches to completing the qualitative and quantitative phases of the analysis. Some methods, such as those associated with root cause analysis, are primarily qualitative; however, even those methods described as quantitative begin with a qualitative analysis. In most HRA methods, the qualitative phase consists of identifying potential human errors and analyzing them in terms of those factors that might contribute to a human making the error [2].

The techniques can be grouped into five categories spanning the principal types and purpose of HRA analysis. Some techniques are primarily descriptive or concern basic data gathering (Table 1). These are often used as a prelude to more sophisticated approaches involving simulation, human error analysis and human error quantification. Techniques may be used separately, but more often in combination [4].

TABLE I

The range and scope of HRA techniques

| Type of technique | Description |
|---|---|
| Data Collection | Collection of information on incidents, goals, tasks, etc. |
| Task Description | Taking the data collected and portraying this in a useful form |
| Task Simulation | Simulating the task as described and changing aspects of it to identify problems |
| Human Error Identification and Analysis | Uses task description, simulation and/or contextual factors to identify the potential errors |
| Human Error Quantification | Estimated the probability of the errors identified |

### B. HAZOP (Hazard and Operability Study)

The HAZOP method is mainly used in the chemical industry. It is a very flexible method that is used for large technological units but can also be used for small devices. This is a method suitable for large and small organizations [5].

The HAZOP study is used to identify hazard scenarios that impact receptors such as people, the environment and property, as well as operability scenarios where the concern is with the capacity of the process to function properly. The HAZOP method is based on an assessment of the likelihood of threats and the resulting risks. Its main goal is to identify possible risk scenarios - thus enabling it to identify the dangerous conditions that may occur on the investigated device. The method looks for so-called critical locations, and then evaluates potential risks and hazardous situations. This is a team-based multidisciplinary method where team members look for scenarios in a joint discussion, for example, using the brainstorming method. The results are formulated in the final recommendation that aims to improve processes or systems [6].

Steps of the HAZOP method:
1. Identify causes
2. Estimation of possible consequences and risks
3. Proposals for risk elimination measures
4. Valuation

### C. FTA (Fault Tree Analysis)

The FTA method was first used in 1962 by Bell Telephone Laboratories and was perfected by Boeing. The method has been used wherever complex systems have to be solved, and to find or reduce failures or improve quality, especially in sectors such as energy, space research, aviation, nuclear power and others [7].

FTA is the analytical technique used to evaluate the probability of failure or the reliability of complex systems. Because of its versatility, it is well-known in many areas, particularly in the areas of risk management and quality management, and safety management. It is applicable as a preventive method as well as a method of analyzing an already existing problem (for example, a crash). The FTA usually follows an FMEA analysis and is designed for complex systems.

The FTA method is based on an analysis of a peak event or a problem (a generally negative phenomenon such as crash, failure, poor quality, high costs) and helps to systematically identify the factors causing the problem or negatively affect the functionality of the system. Its aim is to analyze in detail - to find the causes of the negative phenomenon and further reduce the likelihood of its occurrence. For a simple system, it is preferable to use FMEA or HAZOP methods. FTA is a systematic method for analyzing the cause of risks by adopting a deductive method, in which a specific risk that is only qualitatively recognized from a relevant primary system is placed as the top event in the tree for deductive reasoning [8].

## IV. CONCLUSION

Most techniques in HRA are qualitative. It is caused by the specific of initial data for human factor analysis. This data, as a rule, incompletely specified and ambiguous. Typical approach for the data collection is expert evaluation [9], [10]. There are complications to represent this data according to mathematical models that are used in Reliability Engineering, for example as structure function, reliability block diagram, Markov model, Universal Generation Function, Petri network, etc. [11]. Therefore typical method for quantitative reliability evaluation cannot be used without special adaption for the specific of human factor analysis. For example, techniques of FTA used for analysis of technical system have specifics in HRA application.

## ACKNOWLEDGMENT

## REFERENCES

[1] Su, X., Mahadevan, S., Xu, P., & Deng, Y. (2015). Dependence Assessment in Human Reliability Analysis Using Evidence Theory and AHP. Risk Analysis, 35(7), 1296–1316.W.-K. Chen, *Linear Networks and Systems* (Book style).  Belmont, CA: Wadsworth, 1993, pp. 123–135.

[2] M. Philippart, Human reliability analysis methods and tools. In *Space Safety and Human Performance* (s. 501-568). Oxford ; Cambridge, MA: Butterworth-Heinemann, 2018.

[3] M. Rausand, Risk Assessment: Theory, Methods, and Applications, Wiley, Hoboken, NJ, 2011.

[4] M. Lyons, S. Adams, M. Woloshynowych and C. Vincent, (2004) Human reliability analysis in healthcare: A review of techniques, *Int. J. of Risk and Safety in Medicine*, Vol. 16 (4), pp. 223–237.

[5] P. Baybutt, (2015). A critique of the Hazard and Operability (HAZOP) study. *Journal of Loss Prevention in the Process Industries*, 52-58.

[6] E. Zio, (2007). *Introduction to the Basics of Reliability and Risk Analysis (Series on quality, reliability and engineering statistics ; v. 13)*. Singapore: World Scientific, p.222.

[7] K.-C. Hyun, S. Min, H. Choi, J. Park & I.-M. Lee , (2015). Risk analysis using fault-tree analysis (FTA) and analytic hierarchy process (AHP) applicable to shield TBM tunnels. *Tunnelling and Underground Space Technology*, 121-129.

[8] D. M. Shalev & J. Tiran, (2007). Condition-based fault tree analysis (CBFTA): A new method for improved fault tree analysis (FTA), reliability and safety calculations. *Reliability Engineering and System Safety* , 1231-1241.

[9] B. S. Dhillon (2003), *Human Reliability and Error in Medicine*, World Scientific.

[10] T. Aven, B. Heide, (2009), Reliability and validity of risk analysis, *Reliability Engineering and System Safety*, Vol. 94 (11), pp.1862–1868.

[11] A. Lisnianski, G. Levitin, (2003), Multi-state System Reliability. Assessment, Optimization and Applications. Singapore, SG: World Scientific.

[12] "Navigation," *Human error » NOPSEMA*. [Online]. Available: https://www.nopsema.gov.au/resources/human-factors/human-reliability-analysis/. [Accessed: 11-Dec-2018]

# Remote sensing satellite virtual constellation optimizing with target recognition probability

Alexsandr M. Kondratov, Oleg V. Maslenko

***Abstract***—The improved method of satellite systems selection for the electro-optical observation of the specified objects and regions of the Earth is described. The method provides the probability of targets correct detection (recognition). Other indicators of system selected efficiency, which reflect the controlled area, the timeliness of space imagery and the cost should not be less than a predetermined threshold.

***Keywords***—Earth observation, remote sensing satellite system, virtual constellation, target recognition probability.

## I. INTRODUCTION

Currently, the market of information services in the field of remote sensing of the Earth is developing rapidly. Applications for satellite imagery of the necessary regions of the Earth are realized by the electro-optic imaging satellite systems. Satellite images are distributed by dealers that provide their services in Internet. They offer both new-ordered and archive images for anywhere Earth territory. When ordering the satellite images (both new and archive), a great number of parameters that affect the efficiency of the mission accomplishment should be taken into account.

The main elements of any satellite system for electro-optical surveillance are spacecraft and imaging equipment mounted on them. The satellite system selection depends on the geographical, geometrical and physical characteristics of areas and objects of interest; imaging conditions (time of year, time of day, state of the atmosphere, cloudiness level); technical specifications of the imaging equipment; spacecraft position in space and time in relation to the areas of interest and data reception points; state affiliation of the satellite system owner; information requirements (accuracy, reliability, timeliness of acquisition, cost). Therefore, special algorithms are required for optimal solving the problem of a suitable satellite system selection.

The most common way to forecast the efficiency of the remote sensing task solution is the service simulation of satellite system application [1]. During the simulation the swath of Earth's surface and their positioning on the area of interest are computed. The time intervals when imaging is possible are calculated. The orbits by which the spacecraft passes over the area of interest at a specified time are selected. The illuminance conditions of the area of interest, the viewing angle inclination, imaging repeat cycle, the cost of a satellite imagery scene and other parameters are determined [2, 3]. Thereafter, the problem of the most suitable satellite systems selection is solved. By comparison, the following indicators of the efficiency of satellite observation systems are analyzed: the time of data obtaining and the information renewal rate; productivity of a satellite system that is specified as the total area imaged per day; spatial resolution of the satellite image; positional accuracy of the captured images in geodetical coordinate system, the geometric distortions of image, etc. [4]. Most of the available techniques for satellite observation planning rely on sophisticated and accurate models of the spacecraft orbital motion, however, as a rule, without consideration the probabilistic characteristics of the target detection and recognition [5].

A. M. Kondratov, CASRE, National Academy of Sciences of Ukraine,  (e-mail: dro.huston@gmail.com)
Oleg V. Maslenko, CASRE, National Academy of Sciences of Ukraine,

## II. METHOD

To determine the expected time of satellite imaging, it is necessary to set the geographical coordinates of the area of interest and to define the moments when the path of viewing axis crossing the area boundaries [6, 7]. Initial data are the Kepler elements of solar synchronous orbits. The inclination of the viewing axis for the roll $\eta$ and the current moment of the spacecraft flying time $t_j$, $j = 0, 1, 2, \ldots$ are considered to be known. Required values are the geographic coordinates of the observation point: geographic latitude $\varphi_j = \varphi(t_j)$ and longitude $\lambda_j = \lambda(t_j)$ at specific time interval $t_j$ within the boundaries of the orbit selected. At first, the coordinates of the viewing point are calculated in accordance with a spacecraft coordinate system. Then, the coordinate system is converted to the other one, in which the geographic coordinates of the viewing point are calculated. Further, the path of the viewing axis is designed as a set of the viewing points with allowance for the flying time. Current coordinates of the viewing points in a geocentric spherical coordinate system at a time point $t_j$ [6, 8] are following:

$$\varphi_j = \arcsin\left(\sin u_j \sin i \cos \psi_\eta + \cos i \sin \psi_\eta\right) \tag{1}$$

$$\lambda_j = \pm \arccos\left(\frac{(a_{11} \cos \psi_\eta + a_{31} \sin \psi_\eta) \cos \theta_\varsigma + (a_{12} \cos \psi_\eta + a_{32} \sin \psi_\eta) \sin \theta_\varsigma}{\sqrt{(a_{11} \cos \psi_\eta + a_{31} \sin \psi_\eta)^2 + (a_{12} \cos \psi_\eta + a_{32} \sin \psi_\eta)^2}}\right) \tag{2}$$

where $\psi_\eta$ is the geocentric angle between the radius-vector of a spacecraft and the radius-vector of a viewing point, $a_{11}$, $a_{12}$, $a_{31}$, $a_{32}$ are the elements of transition matrix from the inertial to the Greenwich geocentric coordinate system.

The calculation of a latitude argument at a time point $t_j$ is carried out according to the formula:

$$u(t_j) = \sqrt{\frac{\mu_0}{a^3}} \Delta t_j \tag{3}$$

where $\Delta t_j = t_j - t_\Omega$ is the time interval since the moment the spacecraft was located in the ascending orbit.

An optical axis path makes it possible to set the time when satellite survey begins as a time point when the path coordinates match the geographical coordinates of the area of interest at a certain orbit.

Inasmuch as the area of interest is defined as a trapezoid mostly, the start of satellite survey can be set as [9]:

$$t_m^{sa}(n) = \begin{cases} t_m^S(n), & \text{if } (\varphi_j = \Phi_m^S) \wedge (\Lambda_m^W \leq \lambda_j \leq \Lambda_m^E) \wedge (S \to N) = 1 \\ t_m^W(n), & \text{if } (\lambda_j = \Lambda_m^W) \wedge (\Phi_m^S \leq \varphi_j \leq \Phi_m^N) \wedge (S \to N) = 1 \\ t_m^N(n), & \text{if } (\varphi_j = \Phi_m^N) \wedge (\Lambda_m^W \leq \lambda_j \leq \Lambda_m^E) \wedge (N \to S) = 1 \\ t_m^E(n), & \text{if } (\lambda_j = \Lambda_m^E) \wedge (\Phi_m^S \leq \varphi_j \leq \Phi_m^N) \wedge (N \to S) = 1 \end{cases} \tag{4}$$

where $\Phi_m^S$ and $\Phi_m^N$ are the southern and northern latitudes of the spherical trapezium sides of the $m$-th region, $\Lambda_m^W$ and $\Lambda_m^E$ are the western and eastern longitude of the spherical trapezium sides of the $m$-th region, $t_m^S(n)$, $t_m^N(n)$ are the time points when the sighting axis path crosses the southern and northern boundaries of the $m$-th region at the $n$-th orbit of the spacecraft, $t_m^W(n)$, $t_m^E(n)$ are the time points when the sighting axis path crosses the western and eastern boundaries of the $m$-th region at the $n$-th orbit of the spacecraft, $S \to N$ is the spacecraft orbital motion from south to north, and $N \to S$ from north to south.

The region area $S$ defined as the spherical trapezoid can be calculated [9] by formula:

$$S(\Phi,\Lambda) = R^2 \, (\Lambda_m^E - \Lambda_m^W) \cdot (\sin\Phi_m^N - \sin\Phi_m^S) \qquad (5)$$

where $R$ is the Earth's radius.

The instantaneous projection area of the imaging coverage in case of the viewing axis inclination from the nadir in relation to the flat Earth's surface can be calculated using a formula:

$$S(\eta) = H^2 \, \mathrm{tg}\,\beta \, [tg(\alpha+\eta) + tg(\alpha-\eta)] \cdot [\sec(\alpha+\eta) + \sec(\alpha-\eta)] \qquad (6)$$

The presence or absence of visibility conditions of the $\mu$-th spacecraft for the $m$-th region during the satellite observation can be determined [6] using a logic function

$$\Phi_m^F(n_\mu) = \begin{cases} 1, \text{if } (K_m^S(\mu) \geq \overline{K}_m^S) \wedge (K_m^T(\mu) \geq \overline{K}_m^T) \wedge (\beta_m^c(\mu) \geq \overline{\beta}_m^c) \wedge (Q_m^\xi(\mu) \leq \overline{Q}_m^\xi) = 1 \\ 0, \text{if } (K_m^S(\mu) \geq \overline{K}_m^S) \wedge (K_m^T(\mu) \geq \overline{K}_m^T) \wedge (\beta_m^c(\mu) \geq \overline{\beta}_m^c) \wedge (Q_m^\xi(\eta) \leq \overline{Q}_m^\xi) = 0 \end{cases} \qquad (7)$$

where $K_m^S(\mu)$ is the spatial coefficient of the $m$-th region coverage by the swath of the $\mu$-th spacecraft, which should not be less than the pre-defined valid value $K_m^S$. It can be found as a ratio of the expected area of the imaging of the $m$-th region by $\mu$-th the spacecraft to the total area of the region:

$$K_m^S(\mu) = S_m(\mu)/S_m \,, \quad K_m^S(\mu) \to \max \qquad (8)$$

Secondly, $K_m^T(\mu)$ is the time coefficient of the $m$-th region coverage by the swath of the $\mu$-th spacecraft, which should not be less than the pre-defined valid value $\overline{K}_m^T$.

The conditions of the temporal cover of the $m$-th region by the swath of the $\mu$-th spacecraft can be submitted [6] as

$$[T_m(\mu) \in \overline{T}_m] \wedge [T_m(\mu) \geq \overline{T}_m] = 1 \qquad (9)$$

where $T_m(\mu) = t_m^{end}(\mu) - t_m^{start}(\mu)$ is the duration of an expected imaging interval, $\overline{T}_m = \overline{t}_m^{end} - \overline{t}_m^{start}$ is the duration of a specified imaging interval, $t_m^{end}(\mu)$, $\overline{t}_m^{end}$ are the expected and specified imaging termination time, $t_m^{start}(\mu)$, $\overline{t}_o^{start}$ are the expected and specified imaging start time.

The temporal coefficient of coverage can be found as the ratio of the expected duration of the $m$-th spacecraft over the $\mu$-th area $T_m(\mu)$ to the specified time of visibility of the area:

$$K_m^T(\mu) = T_m(\mu)/\overline{T}_m \,, \quad K_m^T(\mu) \to \max \qquad (10)$$

Thirdly, $\beta_m(\mu)$ is the current of the Sun's elevation angle for the satellite survey of the $m$-th region by the $\mu$-th spacecraft, which needs to meet the requirement $\beta_m(\mu) \geq \overline{\beta}_m$, where $\overline{\beta}_m$ are the permissible minimum angle of the Sun's elevation [6]:

$$\beta_m(\mu) = \begin{cases} \dfrac{(t_{local} - t_{sunset}) \beta_{\max}}{12 - t_{sunset}} & \text{if } t_{local} < 12^h \\[3mm] \dfrac{(24 - t_{local} - t_{sunset}) \beta_{\max}}{12 - t_{sunset}} & \text{if } t_{local} \geq 12^h \end{cases} \qquad (11)$$

where $t_{local}$ is the current local time, $t_{sunset}$ is local time of sunset; $\beta_{\max}$ is the maximum angle of the Sun's elevation within the area of interest.

In the fourth place, $Q_m^\xi(\mu)$ is the forecasted cloud-covered area of the $m$-th region for the $\mu$-th spacecraft, taking into account the coefficient of transparency of the atmosphere $\xi$, which should be no less than the valid value $\overline{Q}_m^\xi$. Information about clouds over certain areas of the Earth can be obtained from the world meteorological services or other relevant institutions.

The expected cost of the image of the observed part of the $m$-th region, acquired by the $\mu$-th spacecraft $C_m(\mu)$, can be estimated using the following data: the area of the $m$-th region observed by the $\mu$-th spacecraft $S_m(\mu)$; commercial offers of the Earth observing systems operators concerning the minimum scene size in order $S_{\min}$; cost of an image $C_1$ per 1 sq. km; threshold area $S_n$, which exceeds the operator's ability to reduce the cost $C_1$; the cost of image acquisition $C_2$ per 1 sq. km with a discount ($C_2 \leq C_1$).

Then the total image cost can be calculated by the following formula:

$$C_m(\mu) = \begin{cases} C_1 S_{\min} & \text{if } S_m(\mu) < S_{\min} \\ C_1 S_m(\mu) & \text{if } S_n > S_m(\mu) \geq S_{\min} \\ C_2 S_m(\mu) & \text{if } S_m(\mu) \geq S_n \end{cases} \tag{12}$$

Then the normalized cost factor of the satellite image should be used in the form

$$K_m(\mu) = \begin{cases} C_m(\mu)/C_m^{\max} & \text{if } C_m(\mu) < C_m^{\max} \\ 1 & \text{if } C_m(\mu) \geq C_m^{\max} \end{cases} \tag{13}$$

where $C_m^{\max}$ is the maximum cost of the image of the $m$-th area acquired by available spacecraft.

The coefficient (13) can take a range of values (14)

$$0 \leq K_m(\mu) \leq 1 \tag{14}$$

At the same time, the following term should be compiled with

$$K_m(\mu) \to \min \tag{15}$$

and the selection of suitable satellite system should be carried in accordance with smallest values of the coefficient (15) from the ordered set

$$K_m(1) < K_m(2) < ... < K_m(\mu) \quad , \quad \mu = \overline{1,\mathrm{M}} \tag{16}$$

The probability of correct detection (recognition) of the target $P(x,\theta)$ can be described by a generalized relation

$$P(x,\theta) \cong 1 - \prod_i \varepsilon(x_i, \theta_i) \tag{17}$$

where $x$ is the input vector of optical signal, $\theta$ is the set of parameters, $\varepsilon(x,\theta)$ is the probability of error. The probability of error can be written as [10]

$$\varepsilon(x,\theta) \cong 1 - \Phi\left(\frac{\Delta x \sqrt{n}}{2\sigma}\right) \tag{18}$$

where $\Phi(x) = \dfrac{2}{\sqrt{2\pi}} \int_0^x e^{-\frac{u^2}{2}} du$ is the probability integral [11], $\Delta x$ is the difference between target and background optical signals, $\sigma$ is the standard deviation of the optical signals, $n$ is the number of resolution elements within the target image.

Equation (17) describes the process of object detection. In passing to recognition, the dependence $\varepsilon(x,\theta)$ becomes more complicated. This fact can be taken into account using the Johnson criteria [12] or any other model that describes image recognition to the required information level.

### III. RESULT

In this research, the model of satellite systems selection for the electro-optical imaging has been applied. For this purpose, the satellite survey of a specified area has been simulated. The area of interest was specified as shown in Fig. 1.
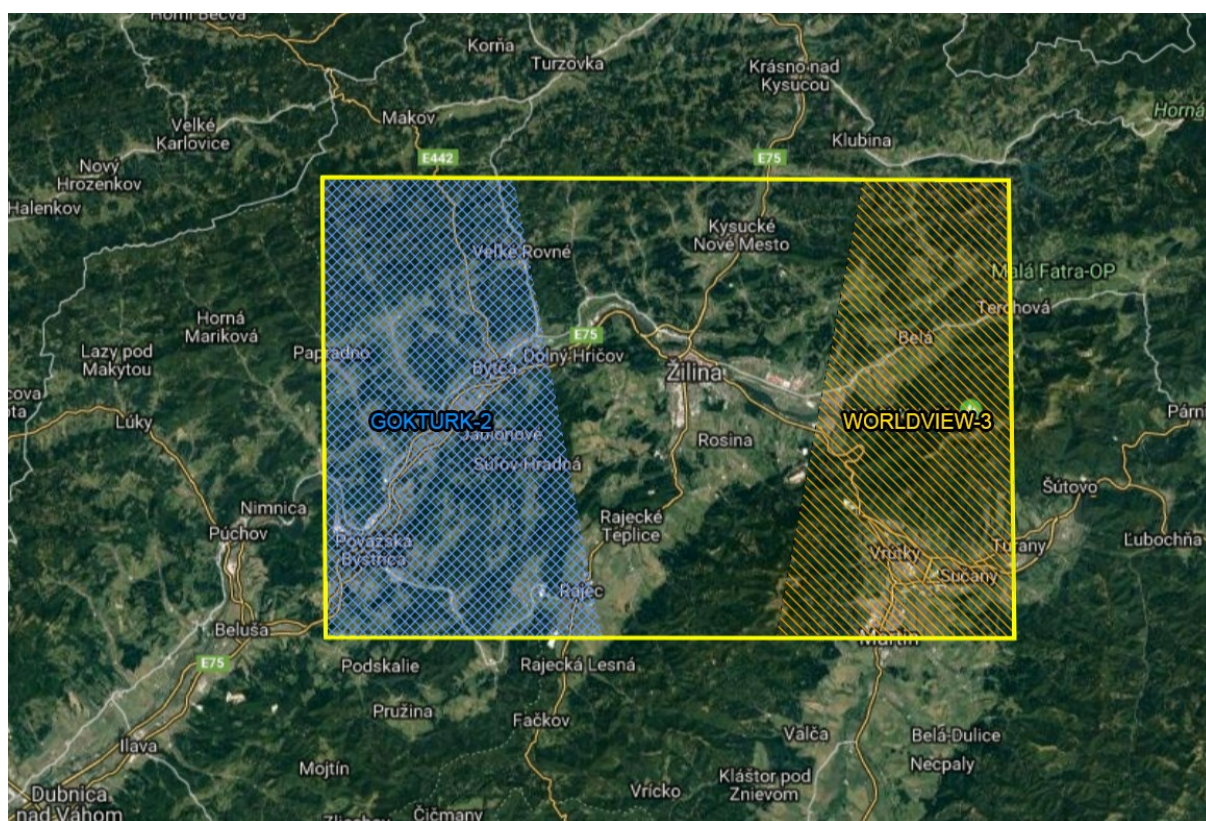
Fig. 1. The area of interest for satellite system selection

Also the some satellite system's swaths are plotted within the area of interest.
The simulation's results are presented in the Table I.

TABLE I
THE RESULTS OF THE SATELLITE SYSTEMS EVALUATION

| Satellite system | Time frame: form 6:00 a.m., 3 January 2019 till 6:00 a.m., 13 January, 2019 | | Time required to cover the region completely by one spacecraft swath | |
| --- | --- | --- | --- | --- |
| | Detection probability | Recognition probability | Hours | Days |
| WorldVie-1 | 0,9978 | 0,9912 | 122,4 | 5,1 |
| Geoeye 1 | 0,9985 | 0,9941 | 123,912 | 5,163 |
| Worldview-2 | 0,9981 | 0,9925 | 123,888 | 5,162 |
| SPOT 6 | 0,9651 | 0,8681 | 3,24 | 0,135 |
| GokTurk 2 | 0,9461 | 0,8017 | 74,208 | 3,092 |
| KompSat 2 | 0,9912 | 0,9653 | 122,88 | 5,12 |
| KompSat 3 | 0,9957 | 0,9828 | 77,592 | 3,233 |
| Pleiades 1A | 0,9978 | 0,9912 | 52,296 | 2,179 |
| WorldView-3 | 0,9991 | 0,9966 | 148,344 | 6,181 |
| KazEOSat 1 | 0,9912 | 0,9653 | 51,816 | 2,159 |
| KazEOSat 2 | 0,6874 | 0,2244 | 27 | 1,125 |
| Pleiades 1B | 0,9978 | 0,9912 | 76,176 | 3,174 |
| SPOT 7 | 0,9651 | 0,8681 | 27,144 | 1,131 |
| NigeriaSat 2 | 0,9461 | 0,8017 | 98,832 | 4,118 |
| Cartosat 2A | 0,9963 | 0,9852 | 26,928 | 1,122 |
| Cartosat 2B | 0,9963 | 0,9852 | 74,904 | 3,121 |
| DubaiSat 2 | 0,9912 | 0,9653 | 51,072 | 2,128 |
| Deimos 2 | 0,9912 | 0,9653 | 123,192 | 5,133 |
| WorldView-4 | 0,9991 | 0,9966 | 76,32 | 3,18 |
| Sentinel 2A | 0,4118 | 0,0291 | 28,272 | 1,178 |
| Sentinel 2B | 0,4118 | 0,0291 | 3,936 | 0,164 |

As it follows from the analysis of Table 1 data, the optimal satellite system for target detection is SPOT 6, and for target recognition is Cartosat 2A. Other satellite systems either do not meet the probability of detection (recognition), or consume more time to capture the entire area of interest.

## IV. CONCLUSIONS

The improved method of satellite systems selection for the electro-optical observation of the specified targets and regions of the Earth is presented. The method provides the selection of electro-optical observation satellite systems with allowance for area of interest characteristics, conditions and timeliness of satellite survey, and also the cost of satellite imagery. The fundamental advantage of the method is the providing the probability required for targets correct recognition by the satellite imagery.

### REFERENCES

[1]  I.A. Glazkova, V.V. Malyshev, and V.V. Darnopykh, "Estimation of perspective micro-satellite Earth observation system efficiency on the base of imitative modeling (in Russian)", *Computer Science and Control*, vol. 16, no. 6, pp. 125-134, June 2009.

[2]  *STK User's Guide*. Exton, PA: Analytical Graphics, Inc., 2004, 536 p.

[3]  V.M. Vishnyakov, "Optimization of orbital constellation parameters of the satellite system for emergency monitoring (in Russian)", *Current Problems in Remote Sensing of the Earth from Space*, vol. 1, no. 2, pp. 222-237, June 2005.

[4]  O.D. Fedorovskyi, M.V. Artiushenko, and Z.V. Kozlov, "Parametric synthesis of space systems for remote sensing of the Earth on the basis of the genetic method (in Russian)", *Space Science & Technology*, vol. 10, no. 1, pp. 54-60, March 2004.

[5]  V. Malyshev, and V. Bobronnikov, "Mission planning for remote sensing satellite constellation", in: *Mission Design & Implementation of Satellite Constellations*, ed. by Jozef C. van der Ha, Dordrecht: Kluwer Academic Publishers, 1998, pp. 431-437.

[6]  P.V. Friz, *The Spacecraft Orbital Movement Fundamentals* (in Ukrainian), Zhytomyr: Korolev Zhytomyr Military Institute, 2012, 348 p.

[7]  *Spacecraft Flight Theory Fundamentals* (in Russian), ed. by G.S. Narimanov, Moscow: Machine Building, 1972, 608 p.

[8]  B.S. Skrebushevsky, *Spacecraft Orbits Formation* (in Russian), Moscow: Machine Building, 1990, 256 p.

[9]  P.V. Friz, "Improved mathematical apparatus for determining the observed area of the landed earth region in space monitoring problems" (in Ukrainian), *The Journal of Zhytomyr State Technological University. Series: Engineering*, vol. 1, no. 2, pp. 126-134, Feb. 2017.

[10] S.A. Stankevich, "Estimating the linear resolution of digital aerospace images" (in Russian), *Space Science & Technology*, vol. 8, no. 2-3, pp. 103-106, May 2002.

[11] T.T. Soong, *Fundamentals of Probability and Statistics for Engineers*, Hoboken, NJ: Wiley, 2005, 498 p.

[12] T.A. Sjaardema, C.S. Smith, and G.C. Birch, *History and Evolution of the Johnson Criteria*, Oak Ridge, TN: Sandia National Laboratories, 2015, 40 p.

# Reliability Prediction of Electronic Devices, Considering the Gradual Failures

Sergei M. Borovikov, Evgeni N. Shneiderov

*Abstract*—The authors suggest a method of the reliability prediction for electronic devices, considering possible gradual failures. Reliability prediction value of the new samples can be given using the degradation (ageing) parameters' model obtained by preliminary research of a sample of products.

*Keywords*—Reliability; prediction; electronic devices; physical-statistical models; gradual failures.

## I. INTRODUCTION

During the operation time of electronic devices (ED) in electronic circuits their functional parameter (denote by $y$) changes and can be considered as a time-function $t - y(t)$. The gradual change of parameter $y(t)$ and its output outside the established norms defines such concept as a gradual failure. Reliability to gradual failure of ED characterizes ability to save the level of functional parameter $y(t)$ within the norms, specified in the technical documentation or by the customer, for a given time $t_G$ at selected modes and conditions.

This reliability was named "parametric". A quantitative measure of parametric reliability level is the probability $P(t_G)$. It can be defined as [1–3]

$$P(t_G) = P\{a \leq y(t) \leq b, t \leq t_G\}, \qquad (1)$$

where $P\{\ldots\}$ – probability of execution of the condition specified in the brackets.

Parametric reliability can be considered as a component of the general reliability of ED.

Unexpected failures causes can be largely eliminated as the result of development of electronic devices production technology [1–3]. But it is impossible to eliminate gradual failures that reflect the inherent material properties of ED, in particular aging. This is the reason of rising interest to gradual (wear-out) failures of ED. It is known [1–4] that gradual failures and therefore parametric reliability of ED can be predicted. It is topical to develop a method of obtaining the predicted value of parametric reliability.

## II. PHYSICAL MODELING OF THE PARAMETERS' DEGRADATION

To obtain reliable prediction to gradual failures and, therefore, parametric reliability of ED we must have a quantitative model of reliability as a function of degradation function parameter $y(t)$ of time, temperature, electrical load, and other operational factors [1, 2, 5–9]. This model is based on a study of the behavior of ED, not only at the time of failure, but also in gradual change of the functional parameter $y(t)$, that is, the study of the kinetics of failure, and can be obtained from the probabilistic and statistical methods. Degradation model of a functional parameter of ED, constructed in such way, is called physical-statistical [1, 2, 5, 6]. As soon as physical-statistical model of the degradation of the functional parameter $y(t)$ will be obtained, finding the probability defined by (1) is possible from a mathematical point of view. Obtaining physical-statistical model of degradation is facilitated by physical experiment, during the modeling of the most common conditions of failure mechanisms and processes of physical and chemical degradation of the functional parameter $y(t)$. Let us explain some of the mathematical aspects of physical modeling

S.M. Borovikov, Belarusian State University of Informatics & Radioelectronics, Minsk, Belarus (e-mail: bsm@bsuir.by).
E.N. Shneiderov, Belarusian State University of Informatics & Radioelectronics, Minsk, Belarus (e-mail: shneiderov@bsuir.by).

of degradation of functional parameters and parametric reliability prediction of electronic devices on these models.
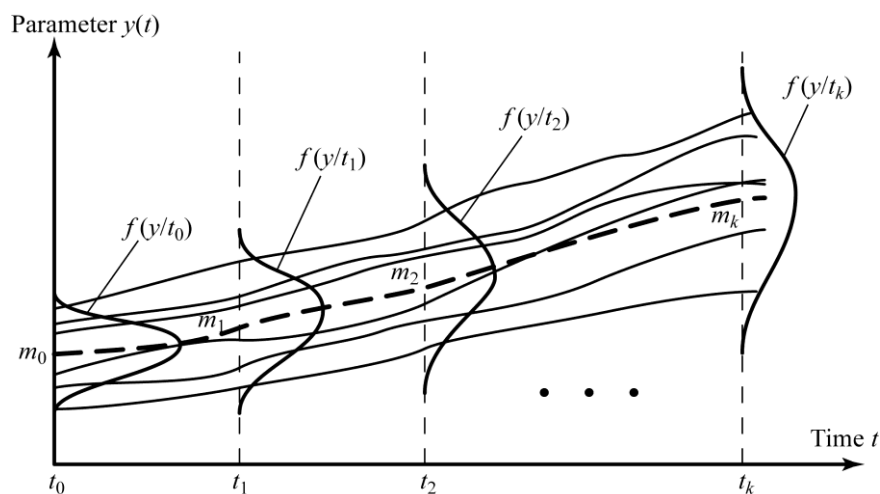
In many cases, the quantitative characteristic of parametric reliability $P(t_G)$ defined by (1) may be obtained on the basis of knowledge of the distribution law of the functional parameter $y(t)$ at the initial time, such as conventional (for the time $t = 0$) density $f(y \mid t = 0)$, and change of function $y(t)$ at time

$$y(t) = \varphi(y_0, t),\qquad (2)$$

where $\varphi$ – symbol of the functional connection; $y_0$ – value $y(t)$ at time $t = 0$.

The value of probability $P(t_G)$ is the result of changes in the statistical distribution $f(y \mid t = t_G)$ of parameter $y(t)$ during the work time $t_G$, $t_G = t_1, t_2, \ldots, t_k$ (Fig. 1).

If the function $y(t)$ is monotonic, then parameters' distribution of samples ED is stored in any time points [10–13]. In this case we can talk about saving for a long time not only the form of the distribution, such as the conditional (for time $t$) distribution density of the functional parameter $f(y \mid t)$, but also a strong correlation of the parameter $y(t)$ for different time points (see Fig. 1). The close correlation is confirmed by experimental studies on high power bipolar transistors (BT) of many types for functional parameters such as collector-emitter saturation voltage ($U_{CE(sat)}$) and static value of the forward current transfer ratio ($h_{21E}$) [13]. Correlation matrix for the parameter $U_{CE(sat)}$ for high power BT of KT872A type, as an example, is given in Table I. As time points we considered values of 0, 3 840, 8 320, 12 800, and 17 280 hours.



$t_0, t_1, t_2, \ldots t_k$ – measuring time points;
$m_0, m_1, m_2, \ldots m_k$ – expectation values of parameter $y$ at appropriate time points (dashed line)

Fig. 1 Changing in density of functional parameter $y(t)$ at work ED

TABLE I
The Correlation of Parameter $U_{CE(sat)}$ for KT872A in Time Points

| Time Section, h | 0 | 3 840 | 8 320 | 12 800 | 17 280 |
|---|---|---|---|---|---|
| 0 | 1,0000 | – | – | – | – |
| 3 840 | 0,9588 | 1,0000 | – | – | – |
| 8 320 | 0,9240 | 0,9899 | 1,0000 | – | – |
| 12 800 | 0,9170 | 0,9868 | 0,9955 | 1,0000 | – |
| 17 280 | 0,8931 | 0,9755 | 0,9881 | 0,9969 | 1,0000 |

The close correlation between the functional parameters $U_{CE(sat)}$ (or $h_{21E}$) in different time points can be considered as a basis for prediction the gradual failures, therefore parametric reliability of electronic devices according to statistical data of the parameter at the initial time ($t = 0$). Further parameters $U_{CE(sat)}$ and $h_{21E}$ generally will be regarded as a parameter $y$.

### III. USING OF PHYSICAL-STATISTICAL MODELS FOR RELIABILITY PREDICTION

Approximate analytic expression of conditional distribution density $f(y \mid t)$ of parameter $y$ for every time $t = t_i$ can be obtained through mathematical transformation of initial distribution $f(y \mid t = 0)$:

$$f(y \mid t = t_i) = \psi[w(y \mid t = 0), t_i],\qquad(3)$$

where $\psi$ – symbol of the functional dependence.

Physical-chemical characteristics of the degradation of the functional parameter $y(t)$, obtained by averaging over the studied ED samples, will be included in the form of the coefficients in the right side of equation (3).

Finding the exact analytical expressions for the function $f(y|t = t_i)$ involves considerable mathematical difficulties. So idealization of parameter $y(t)$ and simplifications allowable in determining $f(y|t = t_i)$ justify themselves, because they provide an opportunity to determine, at least approximately quantitative characterization of parametric reliability $P(t_i)$ with the adopted in probability theory rules of finding the value

$$P(t_i) = P\{a \le y(t) \le b, \quad t \le t_i\},$$

using distribution law of random variables [14]:

$$P(t_i) = \int_a^b f\left(y|t = t_i\right) dy = F\left(b|t_i\right) - F\left(a|t_i\right),\qquad(4)$$

where $F(a|t_i)$, $F(b/t_i)$ – values of the conditional (for time $t_i$) distribution function $F(y|t)$ of the functional parameter $y$, calculated for values $y = a$ and $y = b$.

With streamlined technological process of manufacturing of ED is often observed a normal distribution of the product parameters. We take normal distribution of the functional parameter $y$ as a basis for obtaining the model of degradation. The conditional distribution density $y$ for the given time point $t$ in this case is written as

$$f\left(y|t\right) = \frac{1}{\sqrt{2\pi} \cdot \sigma\left(y|t\right)} \exp\left\{ -\frac{\left[y - m\left(y|t\right)\right]^2}{\left[\sigma\left(y|t\right)\right]^2} \right\},\qquad(5)$$

where $m(y|t)$, $\sigma(y|t)$ – characteristics (parameters) of the normal distribution.

The values $m(y|t)$, $\sigma(y|t)$ represent the average value and standard deviation of the functional parameter $y$ in time point $t$ and in the implicit form include the physical and chemical characteristics of its degradation for interested time $t$. According to the expression (3) values $m(y|t)$, $\sigma(y|t)$ defined as a function of time $t$ and values $m(y|t = 0)$ and $\sigma(y|t = 0)$, which are the parameters of the normal distribution law at the initial time ($t = 0$):

$$m(y|t) = \varphi_1[t, m(y|t = 0), \sigma(y|t = 0)];\qquad(6)$$

$$\sigma(y|t) = \varphi_2[t, m(y|t = 0), \sigma(y|t = 0)], \tag{7}$$

where $\varphi_1$, $\varphi_2$ – symbols of functional dependencies that must be determined.

The conditional distribution density (5), obtained from the expressions (6) and (7), can be considered as physical-statistical model of degradation of the functional parameter $y(t)$. To obtain this model prior studies of a samples of interest ED type are need. Such samples will be called the training. Its size $n$ must be at least 60 ... 100 specimens. The construction of the model using training samples, includes a number of steps that can be found in the works [1, 2, 5–9].

The obtained physical-statistical model of the degradation of the functional parameter $y(t)$ in the form of the conditional distribution density (5) can further be used in practice for the multiple parametric reliability prediction of new samples of studied type ED. Prediction is obtained as the probability determined by the expression (1). Predictive value of this probability $P(t)_{pr}$, according to (4) and the hypothesis of normal distribution of the functional parameter $y$ in time sections $t = t_i$, determined by the expression

$$P(t_i)_{pr} = \Phi\left[\frac{b - m(y|t_i)}{\sigma(y|t_i)}\right] - \Phi\left[\frac{a - m(y|t_i)}{\sigma(y|t_i)}\right], \tag{8}$$

where $i = 1, 2, \ldots, k$; $\Phi[\ldots]$ – tabulated normal distribution function [3, 14], found for the argument given in brackets; $m(y|t_i)$ and $\sigma(y|t_i)$ – parameters the normal distribution, calculated from the expressions (6) and (7) for time $t = t_i$.

According to data obtained during the test of operating time (physical modeling) for specimens of the other (control) samples can be found experimental values of the level of parameter reliability $P(t_i)_{ex}$, appropriate to time points $t_i$. For these goals can be used the expression

$$P(t_i)_{ex} = \frac{r(a \leq y \leq b)}{r}, \quad i = 1, 2, \ldots, k, \tag{9}$$

where $r(a \leq y \leq b)$ – number of specimens of the control sample for which the function parameter $y(t)$ at time $t_i$ is within the norms from $a$ to $b$; $r$ – the total number of specimens in the control sample (control sample volume).

Possibility of using the building of physical and statistical model for degradation of parameter $y(t)$ for prediction the parameter reliability of new samples of ED for the time points in the range of operating time $(t_1 \ldots t_k)$ can be seen from the average prediction error of parametric reliability. For determining this error there where proposed an expression [1, 2, 13]

$$\Delta_{av} = \sqrt{\frac{1}{k}\sum_{i=1}^{k}\left(\frac{P(t_i)_{pr} - P(t_i)_{ex}}{P(t_i)_{ex}}\right)^2} \cdot 100\%, \tag{10}$$

where $k$ – number of time sections, for which predicted and experimental values of the level of parametric reliability were found; $P(t_i)_{pr}$ – predictive value of the parametric reliability level for ED of control sample obtained from (8) for the $i$-th time point; $P(t_i)_{ex}$ – experimental value of the level of parametric reliability for ED of control sample, calculated by the expression (9) for the $i$-th time point; $i = 1, 2, \ldots, k$.

The proposed method of predicting the parametric reliability of ED was tested on the powerful BT of KT872A type. Parameters $h_{21E}$ and $U_{CE(sat)}$ were studied as functional. Electrical measurement modes of parameters conformed to the requirements of BT technical documentation. Below as an illustration present the results for the parameter $U_{CE(sat)}$.

For physical modeling of parameter $U_{CE(sat)}$ degradation two samples were formed: training samples in the amount of $n = 200$ specimens and control sample in the amount of $r = 300$ specimens.

Training samples are used for obtaining physical-statistical degradation model $U_{CE(sat)}$. Control samples were intended to evaluate the error of group prediction. In relation to the control samples at the initial time ($t = 0$) the problem of predicting the group parametric reliability for time points $t_i$ (3 840, 8 320, 12 800, and 17 280 h) was solved, and then a physical simulation of operating time with control of parameter $U_{CE(sat)}$ value on time points was performed.

Physical modeling for the degradation of parameter $U_{CE(sat)}$ was to conduct for BT on standard procedures [15–19] of accelerated forced testing equivalent 17 280 hours in terms of the functioning BT in normal operating conditions. To obtain expressions of the form (6), (7) the application Microsoft Excel, the package "Data Analysis" tool "Regression" was used.

Expressions of the form (6) and (7), obtained using the training sample of specimens, for a parameter $U_{CE(sat)}$:

$$m(U_{CE(sat)}|t) = 0{,}6417m_0 + 1{,}2745\sigma_0 + 1{,}8088(t)^{0,5}, \qquad (11)$$

$$\sigma(U_{CE(sat)}|t) = 1{,}1292m_0 - 280{,}08[m \cdot (\sigma_0)^{-1}]^{0,5} + 2{,}1304(t)^{0,5}, \qquad (12)$$

where $m_0$, $\sigma_0$ – average value and standard deviation of $U_{CE(sat)}$ at the initial time ($t = 0$).

Values $m(U_{CE(sat)}|t)$ and $\sigma(U_{CE(sat)}|t)$, determined by the expressions (11) and (12) are the physical characteristics of the statistical model for degradation of parameter $U_{CE(sat)}$.

For the control sample in the amount $r = 90$ specimens by the expression (8) were obtained predicted parametric reliability values $P(t_i)_{pr}$ of BT for operating time $t_i$. Table II shows the values of the level of parameter reliability for BT, in line with expectations and the experimental observations for several values of norms for parameter $U_{CE(sat)}$, established by the consumer. Condition specified in the figured brackets of expression (1) for time sections $t_i$ was chosen in the form $U_{CE(sat)} \le U_{norm}$, where $U_{norm}$ – norm for the parameter $U_{CE(sat)}$, specified by the BT consumer.

TABLE II
THE RESULTS OF THE PARAMETRIC RELIABILITY PREDICTION ON THE PARAMETER $U_{CE(SAT)}$ IN BT OF CONTROL SAMPLE

| The value of $U_{norm}$, mV | Probability $P(t_i)$ at time $t_i$, h: | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | 3 840 | | 8 320 | | 12 800 | | 17 280 | |
| | $P(t_1)_{pr}$ | $P(t_1)_{ex}$ | $P(t_2)_{pr}$ | $P(t_2)_{ex}$ | $P(t_3)_{pr}$ | $P(t_3)_{ex}$ | $P(t_4)_{pr}$ | $P(t_4)_{ex}$ |
| 900 | 0,781 | 0,853 | 0,679 | 0,735 | 0,618 | 0,706 | 0,576 | 0,677 |
| 1000 | 0,876 | 0,883 | 0,780 | 0,794 | 0,715 | 0,824 | 0,668 | 0,735 |
| 1200 | 0,972 | 0,941 | 0,917 | 0,941 | 0,865 | 0,912 | 0,821 | 0,883 |
| 1400 | 0,996 | 0,971 | 0,977 | 0,941 | 0,95 | 0,941 | 0,920 | 0,912 |
| 1600 | 0,9997 | 1,000 | 0,995 | 0,971 | 0,985 | 0,941 | 0,971 | 0,912 |

Experimental assessment of parametric reliability level of control sample of BT on the parameter $U_{CE(sat)}$ obtained by expression (9). The results of this assessment in the form of

probabilities $P(t_i)_{ex}$ made to the Table II. Using data from the Table II, by expression (10) average prediction error $\Delta_{av}$ of parametrical reliability for different norm values of parameter $U_{CE(sat)}$ (Table III) was calculated.

TABLE III
VALUES OF THE AVERAGE PREDICTION ERROR $\Delta_{av}$

| Value of $U_{CE(sat)}$, specified by the consumer, mV | 900 | 1 000 | 1 200 | 1 400 | 1 600 |
|---|---|---|---|---|---|
| Average prediction error $\Delta_{av,}$ % | 11,26 | 8,09 | 4,83 | 2,40 | 4,18 |

Table III shows that the prediction error are acceptable to the practice in time points from $t = 0$ to $t = t_k = 17\,280$ h.

## IV. CONCLUSION

Obtaining physical and statistical degradation models and group prediction on the example of functional parameters for considered BT types (КТ872А, КТ8272В и КТ8271В) confirmed, that the proposed method enables experimentally, using physical modeling of the ED functional parameter degradation and the statistical analysis of simulation results, to obtain physical-statistical degradation model of this parameter. It is found once for considered type of ED by examining of training samples. The resulting physical-statistical degradation model enables to solve the problem of group prediction at the initial time ($t = 0$) for other samples of the same type of ED.

The solution is to determine the probability of the fact that the ED function parameter will be in the range of specified norms during the time of interest.

## REFERENCES

[1] Borovikov S. M. Statistical forecasting for rejection of potentially unreliable electronic products. – M.: "New Knowledge" publishing, 2013. – 343 p.
[2] Borovikov S. M. "Reliability prediction of electronic equipment"/ it was recommended for publication by the university (BSUIR). – Minsk: MSHRC, 2010. – 308 p.
[3] Borovikov S. M. Theoretical foundations of design, technology and reliability: a book for students of engineering specialties. – Minsk: DesignPRO, 1998. – 336 p.
[4] European Organization of the Quality Control Glassary. – Bern: EOQC. 1988. – 24 p.
[5] Physical-statistical models of functional parameters of the degradation electronics products / S. M. Borovikov, [and other] // NAS Reports of Belarus. – Vol. 51, No. 6, 2007. – pp. 105–109.
[6] The reliability prediction of electronic devices based on a mathematical model of functional parameter degradation / S. M. Borovikov, E. N. Shneiderov, [and other] // BSUIR Reports: Electronics, Materials, Technology, Computer Science. – No. 6 (36), 2008. – pp. 32–39.
[7] Borovikov S. M., Shneiderov E. N. The parametric reliability prediction of electronic devices using Weibull distribution // Reports BSUIR. – No. 7 (61), 2011. – pp. 31–37.
[8] Borovikov S. M., Shneiderov E. N. Forecasting method of parametric reliability of electronic devices by model of functional parameter's degradation / Reports BSUIR. – No. 6 (84), 2014. – pp. 5–11.
[9] Borovikov S., Shneiderov E., Burak I. Models Based on the Weibull-Gnedenko Distribution for the Description of the Degradation of Functional Parameters of Electronic Devices / Computational Problems of Electrical Engineering. – Vol. 6, No.1, 2016. – pp. 1–8.
[10] Borovikov, S. M., Beresnevich A. I., Shalak A. V. Prediction of functional parameters of semiconductor devices using degradation models / Modern Radioelectronics: scientific research and personnel training: collection of materials on the results of international scientific and practical conf., Minsk, April 10–11, 2007; ed. by prof. N. A.Tsyrelchuk. – Minsk : MRC, 2007. – Part 1. – pp.105–107.

[11] Beresnevich, A. I. Prediction of the reliability of electronic products using degradation models of functional parameters / Technical means of information protection: thesis of the report VIII of the Belarusian-Russian scientific and technical conf., Braslav, Belarus, May 24–28, 2010. – Minsk: BSUIR, 2010. – p. 70.

[12] Physical and statistical bases of predicting the reliability of electronic products on models of parameter degradation / S. M. Borovikov, [and other] // Modern problems of radio engineering and telecommunications: materials of the 7th international youth scientific and technical conf. RT-2011, 11–15 April 2011, Sevastopol, Ukraine. – Sevastopol: SevNTU, 2011. – p. 429.

[13] Borovikov S. M., Shneiderov E. N. Correlation of products' functional parameters of the electronic equipment in the time points as a basis for predicting parametric reliability / Modern Means of Communication: papers XVI International scientific and technical conf., Minsk, 27–29 september 2011. – Minsk: HSCC, 2011. – p. 81.

[14] Ventcel E. S., Probability theory: a book for higher. tech. universities; 10th edition. – M.: "High school" publishing, 2006. – 575 p.

[15] Quick Logic Reliability Report / pASIC, Vialink and Quick Logic Corp. – Orleans, 1998. – 21 p.

[16] Robinson L. E. Life expectancy in electronic components and the 10th rule / Testing. – No. 1, 1998. – p. 16.

[17] Acceleration Factors SSB −1.003. – Arlington : EIA Government Electronics and Information Technology Association Engineering Department, 1998. − 14 p.

[18] Bipolar Power Transistors Data Book 1998 / TEMIC Semiconductor GmbH. DGT-005-1297, 1997.

[19] Borovikov S. M., Shneiderov E. N., Plebanovich V. I., Berasnevich A. I., and Burak I. A. Experimental research of electronic products degradation / Reports BSUIR. – No. 2 (104), 2017. – pp. 45-52.

# The Impact of Text Pre-processing and Term Weighting on Al-Hadith Al-Shareef Classification

Ahmed S. J. Abu Hammad

*Abstract*—Preprocessing is one of the key components in a typical text classification framework. The preprocessing step usually consists of tasks such as tokenization, filtering, lemmatization and stemming. This paper studies the impact of text pre-processing and totally different term weighting schemes on Al-Hadith Al-Shareef Classification. Additionally, thereto, presents and compares the effectiveness of three distinct automatics learning algorithms for classifying Al-Hadith Al-Shareef into eight selective books depending on Sahih Muslim. To the best of our knowledge, there is still no published study on this data set. The automatic learning algorithms are Naïve Bayes (NB), Support Vector Machines (SVM), and Complement Naïve Bayes (CNB) with 10-fold cross-validation. We used Term Frequency-Inverse Document Frequency (TF-IDF), Term Frequency (TF), Term Occurrences (TO), and Binary Term Occurrences (BTO) techniques to compute the relative frequency for every word in a very specific document. The results indicate that term stemming and pruning, document normalization, and term weighting dramatically reduce reductional, enhance text representation and directly impact text mining performance. What is more, classification results show that the CNB achieved promising results compared with other supervised methods in classifying A-Hadith. CNB obtains 91.22% accuracy and 91.86% F-measure.

*Keywords*— Arabic Text Classification, Arabic Text Mining, Arabic Morphological Analysis, Term weighting.

## I. INTRODUCTION

Islam based on two fundamental laws: Al-Qur'an as the set of words of Allah and Al-Hadith that documenting words, deeds, provisions, and approvals of Mohammad as the prophet of Allah. Hadith was compiled and classified by many Imams such as Imam Bukhari, Imam Muslim, and Imam Tirmidzi, etc. All of them based on one source prophet Mohammad (peace and blessings of Allah be upon him). Imam Muslim is one amongst the known Imam that according to Ulama. Imam Muslim spent nearly fifteen years to compile over 3000 Hadiths without repetition [25]. Referring to [28], Figure 1 is the component of Hadith.
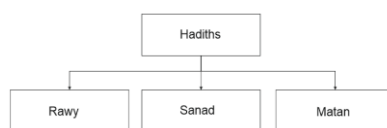


Fig. 1 Hadith components.

Sanad is that the chain of the conveyor of every Hadith, this part present at the beginning of Hadith. Matan is that the content of Hadith, present after the Sanad, and at last Rawy, this is the person or Imam that compile Hadith such as Imam Muslim.

By the exponential growth of digitalized document, emerge the necessity of a system that ready to extract high-quality information, that's why automatic Text Classification (TC) become widespread.

TC task goes through three main steps: text pre-processing, text classification and evaluation. Text pre-processing phase is to make the text documents appropriate to train the classifier. Then, the classifier is built and tuned employing a learning technique against the training

Ahmed Abu Hammad, University College of Science and Technology, Khan Younis, Palestine (e-mail: asj_hammad@hotmail.com).

dataset. Finally, the classifier gets evaluated by some evaluation measures, i.e. recall, precision; etc. The careful description of those steps is often found in [29, 30, 31].

Several existing classification algorithms are used to classify English text corpora such as: SVM [6, 33], NB [6, 7, 33], NB [6, 7, 33], Decision Trees (DTs) [6, 7], k-Nearest Neighbor (KNN) [33], Artificial Neural Networks (ANNs) [33] et al. However, little research works are conducted on Arabic corpora, chiefly since the Arabic language is very wealthy and needs special treatments like order verbs, morphological analysis, etc. Notably, in Arabic morphology, words have affluent meanings and contain a good deal of grammatical and lexical information [32]. Additionally, in grammar structure, Arabic sentence formation differs from English. During this regard, the Arabic text documents are required, significant processing to build an accurate classification model. Therefore, few scholars have applied a variety of classification approaches to the matter of Arabic text classification, i.e. NB [3, 10] [13], SVM [2, 15, 22], KNN [22] and DTs [2, 16]. Even so, researchers conclude that the Arabic text classification may be a terribly difficult task because of language complexity.

This paper studies the impact of text pre-processing techniques and different term weighting schemes on Arabic corpus collected manually from Islam's lawsuit and indicative website. Additionally, presents and compares varied classification rules mining methods associated with the matter of Arabic text classification. Primaries, NB, SVM, and CNB learning methods are applied to classify Sahih Muslim Arabic corpus into one of the predefined categories (books) to measure their performance and effectiveness with reference to different text evaluation metrics like accuracy, precision, recall, and F-measure measures. Experiments are going to be conducted on a specific set of AL-Hadith from Muslim book, wherever eight selective books were chosen as categories so as to run these experiments.

The sub-sequence sections are organized as follows: section 2 contains related works. Section 3 introduces the corpus; we used to test our learning methods and the pre-processing done to the text. Finally, experimental results and evaluation, and conclusions are presented in Section 4 and Section 5 respectively.

## II. RELATED WORKS

The Arabic language is the mother tongue of more than 300 million people; it is considered for religious reasons the language of Islam, and it is ranked as the fifth most spoken language around the world [26]. Unlike Latin-based alphabets, the orientation of writing in Arabic is from right to left; the Arabic alphabet consists of 28 letters. Arabic, in general, is a challenging language because it has a very complex morphology as compared to English. This is due to the unique nature of Arabic morphological principle, which is highly inflectional and derivational [9, 11, 14].

El-kourdi et. al [10] used an NB classifier to classify an in-house collection of Arabic documents. The collections include five classes and three hundred web documents for every class and have used many partitions of the data set. They have concluded that there is some indication that the performance of the NB algorithm in classifying Arabic documents is not sensitive to the Arabic root extraction algorithm, additionally to their own root extraction algorithm; they used other root extraction algorithms. The average accuracy reported was about 68.78%.

Duwairi [8] compared the performance of NB, KNN, and distance-based classifiers for Arabic text categorization. The collected corpus contains a thousand documents that vary in length and writing styles and comprise ten classes every class consists of a hundred documents. The author used stemming to reduce the number of features extracted from documents. The recall, precision, error rate and fallout measures were used to compare the accuracy of classifiers. The results showed that the performance of NB classifier outperformed the other two classifiers.

Al-harbi et. al [2] evaluated the performance of two popular classification algorithms C5.0 decision tree and SVM on classifying Arabic text using the seven different Arabic corpora such as (Saudi News Papers, WEB Sites, Arabic Poems). They have implemented a tool for Arabic text classification to accomplish feature extraction and selection. They have concluded that the C5.0 decision tree algorithm outperformed SVM in terms of accuracy whereas the SVM, average accuracy was 68.65%, while the average accuracy for the C5.0 was 78.42%.

Hattab et. al [17] applied the SVM model in classifying Arabic text documents. The results compared with the other traditional classifiers NB classifier, KNN classifier, and Rocchio classifier. Their experimental results performed on a set of 1132 documents, showing that Rocchio classifier gave better results when the size of the feature set is small while SVM outperformed the other classifiers when the size of the feature set was large enough. The classification rate exceeds 90% when using more than 4000 features.

Al-khatib [4] compared the effectiveness of four different learning algorithms for classifying Al-Hadith Al-Shareef into eight selective books depending on Sahih Bukhari. The testing corpus has 1500 Hadiths that vary in length distributed eight books. The learning algorithms are the Rocchio algorithm, KNN, NB, and SVM. He used the Term TF-IDF technique to compute the relative frequency for each word in a particular document. His results showed that the best accuracy was reported for the SVM algorithm in AL-Hadith Classifications since the precision value is the smallest one for all results. KNN and NB algorithms had a good accuracy in Al-Hadith classifications, and the worst accuracy is reported for the Rocchio algorithm in AL-Hadith classifications since the precision value is the largest one.

Jbara [19] examined the knowledge discovery from AL-Hadith through a classification algorithm in order to classify AL-Hadith to one of the thirteen predefined classes (books) from Sahih AL-Bukhari. The testing corpus has 1321 Hadiths that vary in length distributed over thirteen books. The author used a supervised method called Stem Expansion (SEC) to discover knowledge from AL-Hadith by assigning each Hadith to one book (class) of predefined classes. His results showed that SEC performed better in classifying AL-Hadith against existing classification methods (WBC and AL-Kabi) according to the most reliable measurements (recall, precision, and F-Measure) in the text classification field.

We found that there's a significant shortcoming of the Arabic classification studies during this field. Each study is restricted to a limited range of classification algorithms. This research studies the impact of text pre-processing and different term weighting schemes on Arabic text classification. Additionally, presents and compares distinct classification methods that may use the same corpus in order to evaluate such algorithms and choose the one most suited to the considered case study. This guarantees that the various algorithms had the same conditions and also the same setting in all the experiments.

## III. The Corpus and The Text Pre-processing

### A. The Corpus

In this work, we tend to build an in-house corpus of Arabic texts collected from [18], that referred to as MHAC to perform our experimentation; the corpus includes 1,306 text documents and classified in eight classes that chosen from Sahih Muslim. The corpus contains concerning 24,127 district features after stop words removal. We generate all text representations for MHAC corpus to assess the obtained classification results. The generated text representations for MHAC corpus are: (Light stemming, Stemming) and percentual term pruning (min threshold = 3%, max threshold = 30%) with (TF-IDF, TF, TO, and BTO). Table 1 shows statistical information concerning the books included within the experiments along with its name in English and Arabic as it was used by Sahih Muslim.

TABLE I
UNITS FOR MAGNETIC PROPERTIES

| Book (Class) Name | اسم الكتاب | Number of text documents | Number of distinct features after stop words removal |
|---|---|---|---|
| The Book of Prayers | كتاب الصلاة | 238 | 4062 |
| The Book of Zakat | كتاب الزكاة | 168 | 4758 |
| The Book of Fasting | كتاب الصيام | 200 | 3050 |
| The Book of Marriage | كتاب النكاح | 124 | 2602 |
| The Book of Transactions | كتاب البيوع | 115 | 1021 |
| The Book of Musaqah | كتاب المساقاة | 131 | 2266 |
| The Book of Drinks | كتاب الأشربة | 185 | 3454 |
| The Book of Greetings | كتاب السلام | 145 | 2914 |
| | Total | 1306 | 24127 |

## B.  The Text Pre-processing

One of the widely utilized methods for text mining presentations is viewing the text as a Bag of Tokens (BOT) (words, n-grams). Under that model, we can already classify text [5].

Before applying any algorithm, for both training and testing data, some pre-processing will be conducted on each Hadith. It includes removing Sanad, tokenizing string to words, removing punctuation and diacritic marks, applying stop words removal, applying the proper term stemming and pruning methods as feature reduction techniques, normalizing the tokenized words and finally applying the appropriate term weighting scheme to enhance text document representation as feature vectors. We utilize the open-source machine learning tool Rapid Miner for text pre-processing. Table 2 shows all steps of pre-processing for AL-Hadith.

TABLE II
RESULTS OF PRE-PROCESSING PHASE STEPS FOR AL-HADITH

| Step | Result of the step |
|---|---|
| Removing Sanad | أن رسول الله صلى الله عليه وسلم قال حق المسلم على المسلم ست. قيل ما هن يا رسول الله؟ قال: إذا لقيته فسلم عليه، وإذا دعاك فأجبه، وإذا استنصحك فانصح له، وإذا عطس فحمد الله فشمته، وإذا مرض فعده، وإذا مات فاتبعه. |
| Tokenization | {"أن","رسول","الله","صلى","الله","عليه","وسلم","قال","حق","المسلم","على","المسلم","ست",".",".","قيل","ما","هن","يا","رسول","الله","؟",":","قال","إذا","لقيته","فسلم","عليه","،","وإذا","دعاك","فأجبه","،","وإذا","استنصحك","فانصح","له","،","وإذا","عطس","فحمد","الله","فشمته","،","وإذا","مرض","فعده","،","وإذا","مات","فاتبعه","."} |
| Removing Punctuation and Diacritic Marks | {"أن","رسول","الله","صلى","الله","عليه","وسلم","قال","حق","المسلم","على","المسلم","ست","قيل","ما","هن","يا","رسول","الله","قال","إذا","لقيته","فسلم","عليه","وإذا","دعاك","فأجبه","وإذا","استنصحك","فانصح","له","وإذا","عطس","فحمد","الله","فشمته","وإذا","مرض","فعده","وإذا","مات","فاتبعه"} |
| Removing Stop Words | {"رسول","الله","صلى","الله","وسلم","قال","حق","المسلم","المسلم","قيل","يا","رسول","الله","قال","لقيته","فسلم","وإذا","دعاك","فأجبه","وإذا","استنصحك","فانصح","وإذا","عطس","فحمد","الله","فشمته","وإذا","مرض","فعده","وإذا","مات","فاتبعه"} |
| Light Stemming | {"رسول","له","صل","سلم","قال","حق","مسلم","مسلم","قيل","يا","رسول","له","قال","لقيت","فسلم","اذا","دعاك","فاج","فانصح","اذا","فحمد","عطس","اذا","فشمت","اذا","مرض","فعد","اذا","مات","فاتبع"} |
| Filter Tokens: | {"رسول","له","صل","سلم","قال","حق","مسلم","مسلم","قيل","يا","رسول","له","قال","لقيته","فسلم","اذا","دعاك","فاج ب","اذا","استنصحك","فانصح","اذا","عطس","فحمد","له","فشمت","اذا","مرض","فعد","اذا","مات","فاتبع"} |
| Generate 2-Grams | {"رسول","رسول_له","له","له_صل","صل","صل_له","له","له_سلم","سلم","سلم_قال","قال","قال_حق","حق","حق_مسلم","م سلم","مسلم_مسلم","مسلم","مسلم_قيل","قيل","قيل_يا","يا","يا_رسول","رسول","رسول_له","له","له_قال","قال","قال_لقيت","ل قيت","لقيت_فسلم","فسلم","فسلم_اذا","اذا","اذا_دعاك","دعاك","دعاك_فاجب","فاجب","فاجب_اذا","اذا","اذا_استنصحك","استن صحك","استنصحك_فانصح","فانصح","فانصح_اذا","اذا","اذا_عطس","عطس","عطس_فحمد","فحمد","فحمد_له","له","له_فشمت","فشمت","فشمت_اذا","اذا","اذا_مرض","مرض","مرض_فعد","فعد","فعد_اذا","اذا","اذا_مات","مات","مات_فاتبع","فاتبع"} |

In linguistics, morphology is the identification, analysis, and description of the structure of morphemes and other units of meaning in a language like words, affixes, and parts of speech. For the Arabic language, there are two different morphological analysis techniques; stemming

and light stemming. Stemming is the process for reducing inflected (or sometimes derived) words to their stem, base or root form – generally a written word form. Stemming algorithm by Khoja [21] is one of the well-known Arabic stemmers. Light stemming, in contrast, removes common affixes from words without reducing them to their stems and keeps the words' meanings unaffected [1, 12, 24]. A light stemmer [23] is a standard Arabic light stemmer.

The aim of term weighting is to enhance text document representation as feature vectors. Popular term weighting schemes are TF-IDF, TF, TO, and BTO. BTO indicates the absence or presence of a word with Boolean 0 or 1 respectively. TF(t,d) is the number that the term t occurred in document d. TO be the number of occurrences of term t in document d. TF-IDF is a weight often used in retrieval and text mining. This weight is a statistical measure used to assess how important a word is to a document in a collection or corpus. Term frequency tf(t, d) is the number that the term t occurred in document d. Document frequency df(t) is the number of documents in which the term t occurred at least once. The inverse document frequency can be calculated from document frequency using the formula: log(num of Docs/num of Docs with word i). A reasonable measure of term importance may then be obtained by using the product of the term frequency and the inverse document frequency (TF*IDF) [12, 20, 24, 27].

## IV. EXPERIMENTAL RESULTS AND EVALUATION

We perform experiments on Arabic MHAC corpus collected manually from Islam's lawsuit and indicative website [18]. The corpus includes 1,306 text documents belonging to one of the eight categories (the book of Prayers, the book of Zakat, the book of Fasting, the book of Marriage, the book of Transactions, the book of Musaqah, the book Drinks, and the book of Greetings) that chosen from Sahih Muslim. For text classification, we use NB, SVM, and CNB with 10-fold cross-validation. We split the corpus into two parts (90% of the corpus for training and the remaining 10% to test) using stratified sampling, which keeps class distributions remain the same after splitting. We split the corpus in this way to achieve higher classification results.

For assessing the classification results, we use confusion matrices that are the primary source of performance measurement for the classification problem. We have assessed the obtained classification results utilizing the most common classification measures such as accuracy, precision, recall, and F-measure.

The average classification results are depicted in Figure 2. The morphological analysis (stemming, light stemming), term pruning and term weighting schemes (TF-IDF, TF, TO, BTO) have an obvious impact on the classifier performance as shown in Figure 2. The Figure emphasizes that light stemming, and TO representation for CNB classifier has the best classification results (the accuracy is 91.22%, and the F-measure is 91.86%).
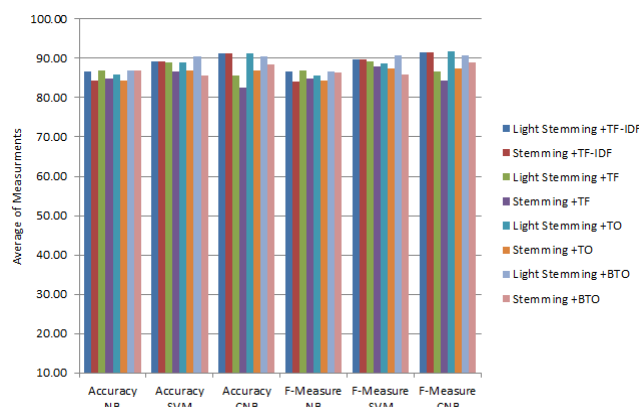


Fig. 2 The classification results for MHAC text representations.

Several observations can be made by analyzing the results in Figure 2. First, using pre-processing techniques like Arabic stop word remover and Arabic stemmer will enhance the accuracy and the F-measure of the classifiers. Second, light stemming has the best classification results this is because lighting stemming is more proper than stemming from linguistics and semantic viewpoint and keeps the word meanings unaffected. Furthermore, classifiers are very sensitive to term weighting schemes because they depend on the distance function to determine the nearest neighbors. For example, the BTO weighting scheme has the worst classification results because the text representation is 0 or 1.
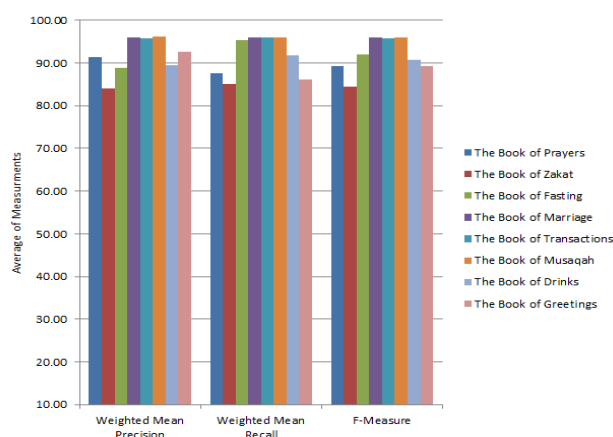


Fig. 3 The classification results for CNB (light stemming + TO)

Figure 3 shows the classification results for the optimal text representation of MHAC corpus (light stemming + TO for CNB) in each of the domain categories. From Figure 2, we can see that the best F-measure is recorded in the book of Musaqah that because the book of Musaqah has limited space of words that are limited and cleared compared with other books. Moreover, it shows that the book of Zakat has the lowest F-measure may be that also because the book of Zakat has a large space domain.

## V. CONCLUSION AND FUTURE WORKS

This paper studies the impact of text pre-processing and different term weighting schemes on Arabic text classification. In addition, presents and compares the effectiveness of three distinct automatic learning algorithms for classifying Al-Hadith Al-Shareef into eight selective books depending on Sahih Muslim. To the best of our knowledge, there is still no published study on this data set. The classifiers have been tested using Arabic text corpus collected manually by us from the Sahih Muslim, which cover eight books: the book of Prayers, the book of Zakat, the book of Fasting, the book of Marriage, the book of Transactions, the book of Musaqah, the book of Drinks, and the book of Greetings. The learning algorithms are NB, SVM and CNB with 10-fold cross-validation are applied to classify Sahih Muslim Arabic corpus. Moreover, we used TF-IDF, TF, TO, BTO and techniques to compute the relative frequency for each word in a particular document. The results indicate that term stemming and pruning, document normalization, and term weighting dramatically reduce dimensionality, enhance text representation and directly impact text mining performance. Furthermore, classification results show that the CNB achieved promising results compared with other supervised methods in classifying A-Hadith. CNB obtains 91.22% accuracy and 91.86% F-measure.

Possible directions for future work include conducting additional experiments using further text collections to make sure the results that we got. Additionally, we tend to decide to use the other feature choice and weighting methods and compare them with the methods already used. Additionally, enhancing the accuracy of the system, more than one classification method can

be merged with each other to increase the accuracy. Finally, it's possible to build a system which can accept as input an archive of texts like Islamic books archive and some category (subject), and as a result, it will give all the texts, which are related to this category.

## REFERENCES

[1] ABABNEH M., ALNOBANI A., AlSHALABI R., and KANAAN G., "Building an Effective Rule-Based Light Stemmer for Arabic Language to Improve Search Effectiveness", *in Proceedings of the International Arab Journal of Information Technology*, vol. 9, no. 4, pp. 368-372, 2012.

[2] AL-HARBI S., ALMUHAREB A., AL-THUBAITY A., KHORSHEED M. and AL-RAJEH A., "Automatic Arabic Text Classification", *in Proceedings of the 9th International Conference on the Statistical Analysis of Textual Data*, Lyon-, France, 2008.

[3] AL-KABI M. N. and AL-SINJILAWI S., "A Comparative Study of the Efficiency of Different Measures to Classify Arabic Text", *University of Sharjah Journal of Pure and Applied Sciences*, 2007.

[4] AL-KHATIB M., "Classification of Al-Hadith Al-Shareef Using Data Mining Algorithm", *Proceedings of European, Mediterranean & Middle Eastern Conference on Information Systems*, Abu Dhabi, UAE, 2010.

[5] AL-SHALABI R., KANAAN G. and GHARAIBEH M., "Arabic Text Categorization Using KNN Algorithm", *in Proceedings of the 4th International Multiconference on Computer Science and Information Technology*, pp. 5-7, 2006.

[6] BERGER H. and MERKL D., A Comparison of Text-Categorization Methods Applied to N-gram Frequency Statistics, *AI 2004: Advances in Artificial Intelligence, Springer*, pp. 998-1003, 2005.

[7] DUMAIS S., PLATT J., HECKERMAN D. and SAHAMI M., "Inductive Learning Algorithms and Representations for Text Categorization", *in Proceedings of the 7th International Conference on Information and Knowledge Management, ACM,* pp. 148-155, 1998.

[8] DUWAIRI R., "Arabic Text Categorization", *in Proceedings of the International Arab Journal of Information Technology*, vol. 4, no. 2, pp. 125-132, 2007.

[9] EL-HALEES A., "Arabic Opinion Mining Using Combined Classification Approach", *in Proceedings of the International Arab Conference on Information Technology (ACIT'2011)*, Riyadh, Saudi Arabia, 2011.

[10] EL KOURDI M., BENSAID A. and E-RACHIDI T., "Automatic Arabic Document Categorization Based on the Naive Bayes Algorithm", *in Proceedings of the Workshop on Computational Approaches to Arabic Script-based Languages,* Association for Computational Linguistics, pp. 51-58, 2004.

[11] ELKATEB S., BLACK W., VOSSEN P., FARWELL D., RODRÍGUEZ H., PEASE A. and ALKHALIFA M., "Arabic WordNet and the Challenges of Arabic", *in Proceedings of Arabic NLP/MT Conference*, London, UK, 2006.

[12] FELDMAN R. and SANGER J., The Text Mining Handbook: Advanced Approaches in Analyzing Unstructured Data, *Cambridge University Press*, 2007.

[13] HADI W., THABTAH F., ALHAWARI S. and ABABNEH J., "Naive Bayesian and K-nearest Neighbour to Categorize Arabic Text Data", *in Proceedings of the European Simulation and Modelling Conference.* Le Havre, France, pp. 196-200, 2008.

[14] HAMMAD A., An Approach for Detecting Spam in Arabic Opinion Reviews, 2013.

[15] HARRAG F. and EL-QAWASMAH E., "Neural Network for Arabic Text Classification", *in Proceedings of the 2nd International Conference on Applications of Digital Information and Web, IEEE*, pp. 778-783, 2009.

[16] HARRAG F., EL-QAWASMEH E. and PICHAPPAN P., "Improving Arabic Text Categorization Using Decision Trees", *in Proceedings of the 1st International Conference on Networked Digital Technologies*, *IEEE*, pp. 110-115, 2009.

[17] HATTAB A. and HUSSEIN A., "Arabic Content Classification System Using Statistical Bayes Classifier with Words Detection and Correction", *in Proceedings of World of Computer Science & Information Technology Journal*, vol. 2, pp. 193, 2012.

[18] Islam website Ministry of Islamic Affairs. April 2014. [Online]. Available: http://hadith.al-islam.com/.

[19] JBARA K., "Knowledge Discovery in Al-Hadith Using Text Classification Algorithm", *in Proceedings of American Science Journal*, vol. 6, no. 11, 2010.

[20] JING L., HUANG H. and SHI H., "Improved Feature Selection Approach TFIDF in Text Mining", *in Proceedings of the 1st International Conference* on *Machine Learning and Cybernetics, IEEE*, Beijing, pp. 944-946, 2002.

[21] KHOJA S. and GARSIDE R., "Stemming Arabic Text", *Computing Department, Lancaster University*, Lancaster, UK, 1999.

[22] KHREISAT L., "Arabic Text Classification Using N-gram Frequency Statistics a Comparative Study", *in Proceedings of the 2006 International Conference on Data Mining (DMIN'06)*, Las Vegas, USA, pp. 78-82, 2006.

[23]  LARKEY L., BALLESTEROS L. and CONNELL M., "Light Stemming for Arabic Information Retrieval", *Arabic Computational Morphology: Knowledge-based and Empirical Methods, Springer*, pp. 221-243, 2007.

[24]  LEWICKI P. and HILL T., Statistics: Methods and Applications, *Statsoft*, 2006.

[25]  RYDING K., A Reference Grammar of Modern Standard Arabic, *Cambridge University Press*, 2005.

[26]  SAID D., WANAS N., DARWISH N. and HEGAZY N., "A Study of Arabic Text Preprocessing Methods for Text Categorization", *in Proceedings of the 2nd International Conference on Arabic Language Resources and Tools*, Cairo, Egypt, 2009.

[27]  SAUBAN M. and PFAHRINGER B., "Text Categorization Using Document Profiling", *in Proceedings of 7th European Conference Principles and Practice of Knowledge Discovery in Databases (PKDD'03)*, pp. 411-422, 2003.

[28]  SEBASTIANI F., "Machine Learning in Automated Text Categorization", *in Proceedings of ACM Computing Surveys (CSUR) Journal*, vol. 34, pp. 1-47, 2002.

[29]  SEBASTIANI F., "Text Categorization", *in Proceedings of the Text Mining and its Applications to Intelligence, CRM and Knowledge*, UK, pp. 109-129, 2005.

[30]  SONG M. and WU Y., Handbook of Research on Text and Web Mining Techologies, Information Science Reference*,* USA, 2009.

[31]  YANG Y. and LIU X., "A Re-examination of Text Categorization Methods", *in Proceedings of the 22nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval,* pp. 42-49, 1999.

[32]  YANG Y., SLATTERY S. and GHANI R., "A Study of Approaches to Hypertext Categorization", *in Proceedings of the Journal of Intelligent Information Systems*, pp. 219-241, 2002.

# Mining Educational Data
# to Analyze Students' Performance
## (A Case with University College of Science and Technology Students)

Ahmed S. J. Abu Hammad

*Abstract*—Data mining is a new information process technology; it has gained some impressive results in higher education, such as detection of abnormal values in the result sheets of the students, a prediction about students' performance and so on. Predicting students' performance is critical for educational institutions because strategic programs can be planned for improving or maintaining students' performance during their period of studies. In this paper, data mining methods are applied at University College of Science and Technology (UCST) – Khan Younis on students enrolled. After data preparation and pre-processing, different techniques of data mining are applied to detect association, classification, clustering, and outlier detection rules. In every one of these four tasks, we offer the extracted knowledge and describe its significance in the educational domain. This paper can provide strong information support, decision support and work direction for administrators of UCST, thus, promote the comprehensive development of the educational system's and improve students' performance in UCST.

*Keywords*—Educational Data Mining, Association Rules, Classification, Clustering, Outlier Detection, Predicting Students' Performance.

## I. INTRODUCTION

Educational Data Mining (EDM), concerns with developing methods that discover knowledge from data originating from an educational context. The data can be gathered from historical and operational data reside in the databases of educational institutes. The student data can be personal or academic [8, 10].

The principal aim of institutions of higher education is to supply quality education for their students and to get better the quality of administrative decisions. One way to succeed in doing the highest level of quality in the higher education system is by discovering knowledge from educational data to study the main attributes that may impact the students' performance. The discovered knowledge can be used to offer helpful and constructive recommendations to the academic planners in institutions of higher education to promote their decision-making process, to improve students 'academic performance and reduce failure rate, to better understand students' behavior, to help instructors, to enhance teaching and many other advantages [1,9].

Educational data mining utilizes numerous techniques such as decision tree, rule induction, k-nearest neighbor, naive Bayesian and numerous others. By utilizing these techniques, numerous sorts of knowledge can be discovered such as association rules, classifications, and clustering [1,13,14].

This paper examines the educational domain of data mining utilizing a case study from the enrolled students' data collected from the University College of Science and Technology - Khan Younis. It showed what sort of data could be gathered, how could we pre-process the data, how to apply data mining techniques on the data, and finally how can we have profited from the discovered knowledge. There are numerous sorts of knowledge can be revealed from the data. In this work, we checked the most common ones which are association rules, classification, clustering and outlier detection. The Rapid Miner software is used for applying the methods to the enrolled student's data set.

The model will predict for students' performance, and relate them with other factors such as gender, address details, general secondary average, general secondary section, specialization of

Ahmed Abu Hammad, University College of Science and Technology, Khan Younis, Palestine (e-mail: asj_hammad@hotmail.com).

the student, and hours completed number. The discovered knowledge, supply a university college management with a helpful and constructive recommendation to conquer the issue of low grades of students and to progress students' academic performance.

The subsequence sections are organized as follows: section II contains related works in educational data mining. Section III illustrates the data set and the preparation and processing methods performed. Then the following section presents our experiments about applying data mining techniques to the educational data. Finally, conclusion and future work are presented.

## II. RELATED WORKS

Data mining in higher education is a recent research field and this area of research is gaining popularity because of its potentials to educational institutes.

El-Halees A. [7], gave a case study that used educational data mining to analyze students' learning behavior. The goal of his study is to show how useful data mining can be used in higher education to improve a student' performance. He used students' data from the database course and collected all available data including personal records and academic records of students, course records and data came from the e-learning system. Then, he applied data mining techniques to discover many kinds of knowledge such as association rules and classification rules using a decision tree. Also, he clustered the student into groups using EMclustering, and detected all outliers in the data using outlier analysis. Finally, he presented how can we benefited from the discovered knowledge to improve the performance of the student.

Nghe N. et al. [11] have compared the accuracy of the decision tree and Bayesian network to predict a student's Grade Point Average (GPA) at the end of the third year of undergraduate and at the end of the first year of postgraduate from two different institutes. All data set has 20,492 and 936 complete student records respectively. The results show that the decision tree algorithm was significantly more accurate than the Bayesian network algorithm for predicting student performance. The accuracy was further improved by using re-sampling technique especially for decision tree in all cases of classes. In the same time, can reduce misclassification especially on minority class of imbalanced datasets because decision tree algorithm tends to focus on a local optimum.

Pal B. [4], utilized the classification as a data mining technique to assess a student' performance, they applied the decision tree method for classification. The point of their study is to extract knowledge that characterizes students' performance in the end semester examination. They utilized students' data including Attendance, Class test, Seminar, and Assignment marks. This study helps earlier in deciding the dropouts and students who need particular attention and allow the teacher to give appropriate advising.

Bekele R. et al. [5] have examined the Bayesian network to predict students' performance of high school students containing 8 attributes (social and personal attributes). Dataset has 514 complete student records respectively. The paper demonstrates an application of the Bayesian approach in the field of education and shows that the Bayesian network classifier has the potential to be used as a tool for prediction of student performance.

Bidgoli M. et al. [6] use the data mining classification technique to predict students' final grades based on their web-use feature. By discovering the successful patterns of students in various categories, the university can predict the final grade of every student. Therefore, it helps to identify students at risk early and allow the instructor to provide appropriate advice in a timely manner. From this case study, it can be concluded that data mining is effective in predicting a student's performances in the educational domain. The result has an impact in improving the transition rate, and the process indicator of a higher learning to institute by improving the student assessment process.

Ayesha S. et al. [2], applied a k-means clustering algorithm as a data mining technique to anticipate students' learning activities in a students' database including class quizzes, mid and final exam and assignments. This correlated information will be reported to the class teacher before the conduction of the final exam. This study helps the teachers to minimize the failing ratio by making appropriate strides at the perfect time and improve the performance of students.

Paris I. et al. [3] have compared the accuracy of data mining methods to predict students' grade. Dataset has 2427 complete records for Bachelor of Computer Science students at University Putra Malaysia (UPM) admitted from 2000 to 2004. The results show that combining different classifiers improved the prediction accuracy compare to single classifiers. The results also show that the resampling technique has not improved the accuracy of prediction in all cases. The results also show that the hidden naive Bayes method consistently outperformed.

### III. STUDENT ENROLLMENTS DATASET AND PRE-PROCESSING

The admissions and registration department is currently gathering demographic, geographic, exam scores, financial information, so on., from applicants as part of the admissions and registration operation. There is too historical data available indicating the actual enrollment status of applicants along with all the other attributes that were gathered as a component of the admissions and registration operation [1].

In this study, data gathered from the UCST. The dataset contains the student enrolled data. The dataset made obtainable has 20 different attributes for each applicant including the decision result attribute. There are in all about 1173 records available. Table I shows the attributes, their types, and description that exist in the data set as taken from the source database.

TABLE I
THE ENROLLED STUDENT'S DATASET DESCRIPTION

| Attribute | DESCRIPTION | DATA TYPE | SELECTED |
|---|---|---|---|
| Seq. | A sequence for the record | Number | |
| Institution | A name for the institution | String | |
| ID | An identifier for the record | Number | |
| Name | Student's named | String | |
| DOB | Date of birth | String | |
| Gender | Student's gender | String | √ |
| Nationality | Student's nationality | String | |
| City | Student's address details | String | √ |
| GS Source | Source general secondary | String | |
| GS Year | Year general secondary | Number | |
| GS Avg | Average general secondary | Number | √ |
| GS Sec | Section general secondary | Number | √ |
| YI Join | Year institution join | Number | |
| YI Term | Year institution term | Number | |
| Std Level | Student's level | Number | |
| College | Student's college | String | |
| Specialization | Student's specialization | String | √ |
| HC Num | Hours Completed number | Number | √ |
| GPA | Student's a cumulative grade point average | Number | |
| Grade | Student's performance | String | √ |

An important step in the data mining process is data preprocessing. One of the challenges that face the knowledge discovery process in predicting students' performance is poor data quality. For this reason, we try to prepare our data carefully to obtain accurate and correct results. As part of the preparation and preprocessing of the data set, irrelevant and weakly relevant attributes should be removed. The attributes marked as selected as seen in Table 1 are processed via the Rapid Miner software to apply the data mining methods on them. The attributes, for example, ID are not chosen to be part of the mining process; this is because they do not provide any knowledge for the data set processing, likewise, they have vast differences which make them irrelevant for data mining [1,10].

The following steps are executed as part of the preparation and preprocessing of the data set:

- Selected attributes have little missing (no more than 27 value). Then we try to fill the missing with appropriate values. So, we used to replace the missing value method. This method enables the substitution of the missing values by the minimum, maximum or average statistics calculated on the basis of existing values for all or selected attributes. Moreover, we can also replace the missing values by some pre-defined values (e.g., zero or values that we consider that provide a better fit to data). Here we the substitute of the missing values by the average calculated on the basis of existing values for all selected attributes.
- GS Avg and HC Num attributes contain many values that cannot easily identify interesting patterns in the data from which to create a model. So, we use Discretize by User Specification method. This method allows numerical attributes to be placed in bins where the boundaries of the bins are defined by the user. This converts numerical attributes into nominal ones as required by some algorithms. Here we the substitute of GS Avg attribute by classes (Excellent – Very Good – Good – Acceptable – Fail) and HC Num attribute by classes (First – Second).

After applying the pre-processing and preparation methods, we try to analyze the data visually and figure out the grade distribution of the students which are in the pivot of the predicting students' performance, Figure 1 depicts the grade distribution of the students.
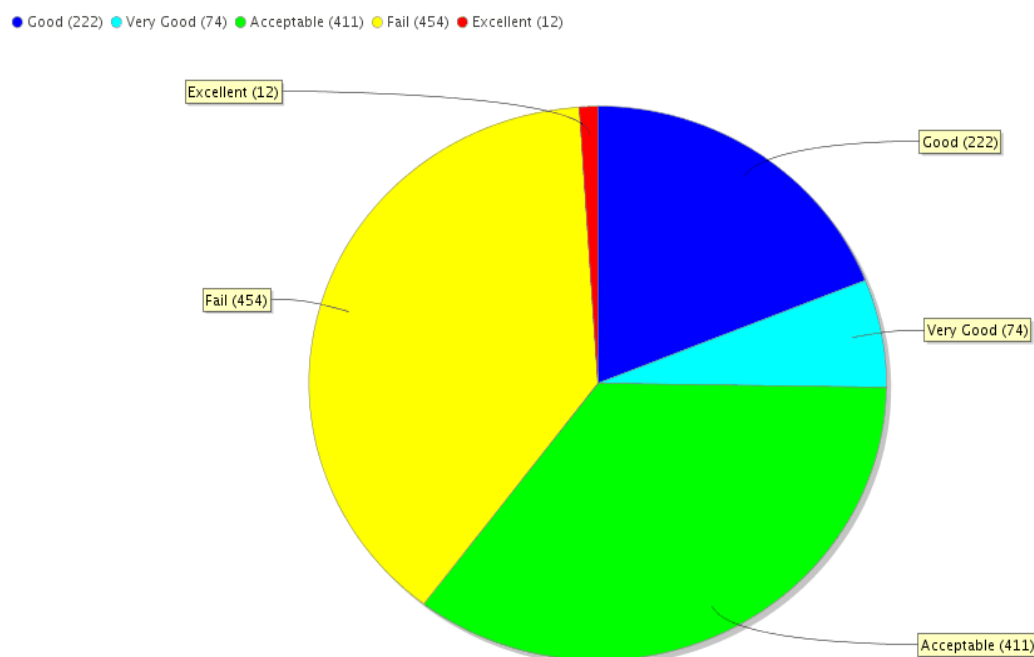


Fig. 1 The distribution of enrolled students according to their grades.

## IV. APPLICATION OF DATA MINING TECHNIQUES TO STUDENT ENROLMENTS DATASET: RESULTS AND DISCUSSION

Before applying the data mining techniques on the data set, there should be a methodology that governs our work. Figure 2 presents the work methodology used in this paper, which is based on the framework proposed in [8]. The methodology starts from the problem definition, then preprocessing which are debated in the introduction and the data set and preprocessing sections, then we come to the data mining methods which are an association, classification, clustering, and outlier detection, followed by the evaluation of results and patterns, finally the knowledge representation process [10, 12].
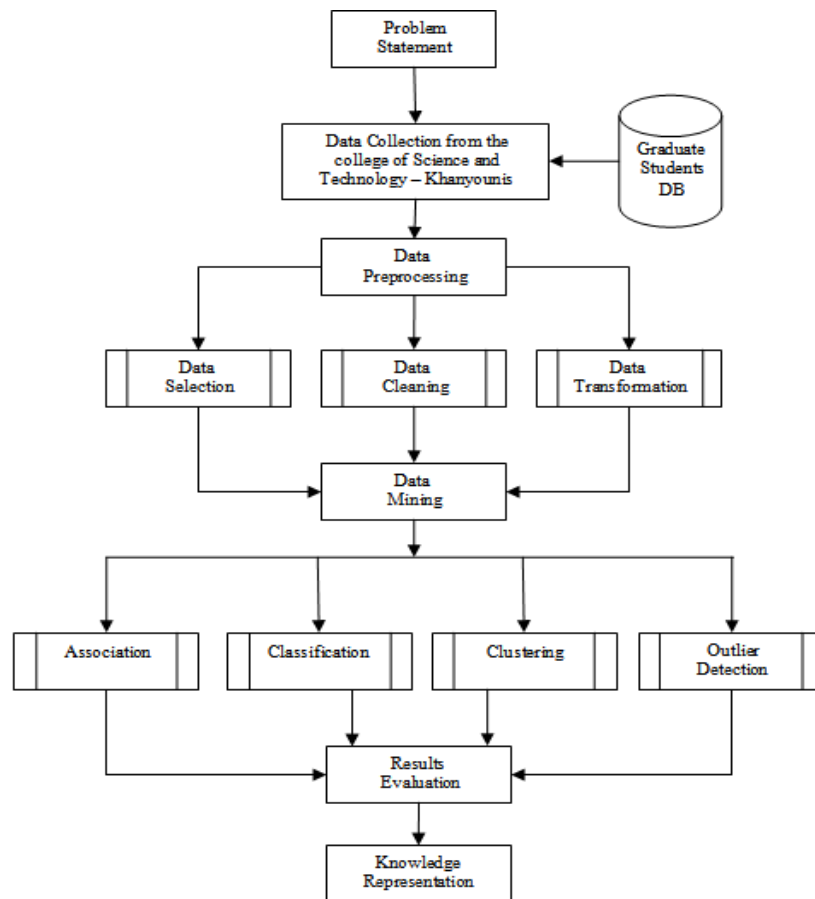
Fig. 2 Data mining work methodology [8].

In this section; we illustrate the results of applying the data mining techniques to the data of our case study, for every one of the four data mining tasks; association, classification, clustering and outlier detection, and how we can profit from the discovered knowledge.

*4.1. Association method*

Association method is the discovery of association rules showing attribute-value conditions that occur frequently together in a given set of data. Association analysis is widely used for data analysis. More formally, association rules are of the form X=>Y, i.e.,"A1^----^Am → B1^----^Bn", where Ai (for i to m) and Bj (j to n) are attribute-value pairs. The association rule X=>Y is interpreted as database tuples that satisfy the conditions in X are also likely to satisfy the conditions in Y " [10, 12].

As part of the association method, before applying the FP-Growth algorithm requires the change of nominal attributes into binomial attributes, the FP-Growth algorithm is applied to student enrollment data set with min confidence = 0.89, Table II illustrates some useful rules extracted from student enrolment data ordered by confidence.

TABLE II
ASSOCIATIONS RULES FOR STUDENT ENROLMENT DATA.

| # | RULE | CONFIDENCE |
|---|------|------------|
| 1 | [GS Avg = Weak, Specialization = Business] --> [GS Sec = Literary] | 0.965 |
| 2 | [Specialization = Business] --> [GS Sec = Literary] | 0.961 |
| 3 | [City = Khan Younis, Specialization = Business] --> [GS Sec = Literary] | 0.957 |
| 4 | [City = Khan Younis, GS Sec = Literary, Gender = Male] --> [GS Avg = Weak] | 0.937 |
| 5 | [GS Sec = Literary, Gender = Male, HC Num = Second] --> [GS Avg = Weak] | 0.929 |
| 6 | [GS Sec = Literary, Gender = Male] --> [GS Avg = Weak] | 0.927 |
| 7 | [City = Khan Younis, GS Sec = Literary, Specialization = Business] --> [GS Avg = Weak] | 0.915 |
| 8 | [City = Khan Younis, Specialization = Business] --> [GS Avg = Weak] | 0.909 |
| 9 | [GS Sec = Literary, Specialization = Business] --> [GS Avg = Weak] | 0.906 |
| 10 | [Specialization = Business] --> [GS Avg = Weak] | 0.902 |

Rules #1, #2, #3 and #4, can be used to predict the general secondary section of the student. For example, from rule #1, we understand that there is a general secondary section = Literary of the student if he is general secondary average = Weak and specialization = Business. Rules #5, #6, #7, #8, #9, #10, #11 and #12 provide with better understanding for general secondary average. For example, from rule #5, we understand that there is a general secondary average = Weak of the student if he is city = Khan Younis, general secondary section = Literary, and gender = male.

### 4.2. Classification method

Classification is the processing of finding a set of models (or functions) which describe and distinguish data classes or concepts, for the purposes of being able to use the model to predict the class of objects whose class label is unknown [10]. In this paper, the classification approaches are used to predict students' performance and present how other attributes affect them. The derived model may be represented in various forms, such as decision trees, rule-based classifier, k-nearest neighbors, or naive Bayesian classifiers. Here we apply five classification techniques above on our data set.

#### 4.2.1. Decision trees

Decision trees have become one of the most powerful and popular approaches in knowledge discovery and data mining [12], Figure 3 depicts the decision tree that resulted from applying the decision tree classification algorithm on the grade as a target class.

As it is seen from the Figure 3, the attributes that influence the category of the target class are GS Avg, City, GS Section, and HC Num, the model presented in Figure 3 has the accuracy of 41.88%. To interpret the rules in the decision tree, the most left branch of the decision tree says that, if the average general secondary = Acceptable and city = Gaza, then grade= Good.
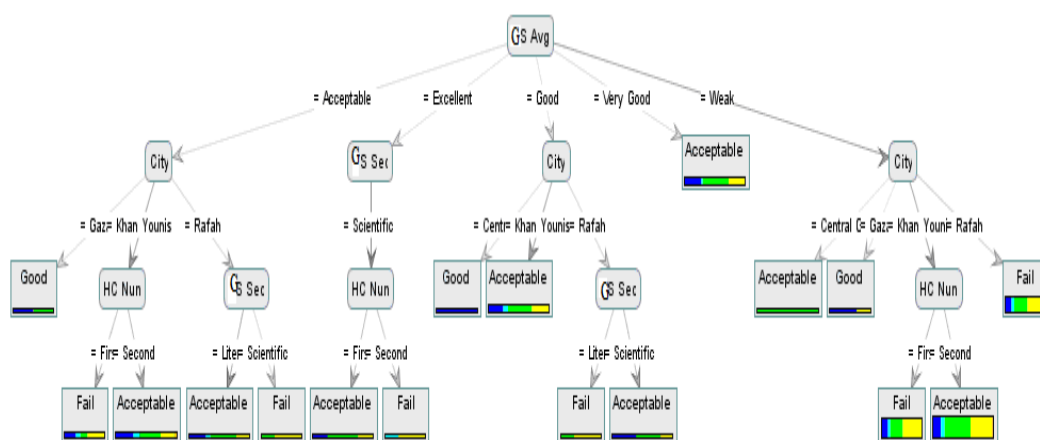


Fig. 3 Decision tree classification for Grade as a target class.

#### 4.2.2. Rule-based classifier

A rule-based classifier is a technique for extracting a set of rules that demonstrate the relationships between the attributes of a dataset and the class label using a collection of" if ... then ..." rules [10]. Table III depicts some of the rule-based classifiers that resulted from applying the rule-based classifier algorithm on the grade as a target class.

As it is observed from Table III, the attributes that influence the category of the target class (grade) are Specialization, GS Avg, City, GS Section, and HC Num, the model presented in table III has the accuracy of 61.88%. To interpret the rules in the rule-based classifier, for example from rule #1 we understand that If Specialization = Graphic Design and HC Num = First then Grade = Acceptable.

TABLE III
THE RULE-BASED CLASSIFIER FOR GRADE AS A TARGET CLASS.

| # | RULE |
|---|------|
| 1 | If Specialization = Graphic Design and HC Num = First then Grade = Acceptable |
| 2 | If Specialization = Graphic Design and GS Sec = Literary then Grade = Acceptable |
| 3 | If Specialization = Computer networks and the Internet Graphic and GS Sec = Literary then Grade = Acceptable |
| 4 | If City = Gaza and GS Avg = Weak then Grade = Acceptable |
| 5 | If Specialization = Medical laboratories and HC Num = Second then Grade = Acceptable |

### 4.2.3. k-Nearest Neighbors

K-nearest neighbor is a supervised learning algorithm where the result of new instance query is classified based on the majority of K-nearest neighbor category. The purpose of this algorithm is to classify a new object based on attributes and training samples. The classifiers do not use any model to fit and only based on memory. Given a query point, we find K number of objects or (training points) closest to the query point. The classification is using majority vote among the classification of the K objects [12]. The classification model that resulted from applying the K-NN algorithm; the model has five classes with the accuracy of 54.37%.

### 4.2.4. Naive Bayesian Classifiers

Naive Bayesian Classifiers or Naive Bayes is a technique for estimating probabilities of individual variable values, given a class, from training data and to then allow the use of these probabilities to classify new entities [10], Naive Bayes on our dataset presents the accuracy of 58.50%.

### 4.3. Clustering method

Unlike classification and prediction, which analyze class-labeled data objects, clustering analyzes data objects without consulting a known class label. In general, the class labels are not present in the training data simply because they are not known to begin with. Clustering can be used to generate such labels. The objects are clustered or grouped based on the principle of maximizing the intra-class similarity and minimizing the interclass similarity. The clustering method applied in this paper is the "k-means"; the objective of this k-means test is to choose the best cluster center to be the centroid [10,12]. The k-means algorithm requires the change of nominal attributes into numerical. The clustering method produced a model with four clusters. Figure 4 demonstrates the resulting "Centroid Table" where from the figure we can see the average value of each attribute in each cluster; for example, the cluster labeled "Cluster_0" has an average of grade: 1.931 and this cluster has 348 items which represent to about %29.67 of the records. The centroid of a cluster represents the most typical case.

| Attribute | cluster_0 | cluster_1 | cluster_2 | cluster_3 |
|-----------|-----------|-----------|-----------|-----------|
| GS Avg | 0.361 | 0.256 | 1.130 | 0.137 |
| HC Num | 0.491 | 0.470 | 0.519 | 0.601 |
| Gender | 0.451 | 0.441 | 0.652 | 0.668 |
| City | 0.379 | 0.265 | 0.370 | 0.208 |
| GS Sec | 0.896 | 0.613 | 0.556 | 0.815 |
| Specializatio | 6.555 | 1.879 | 10.848 | 14.626 |
| Grade | 1.931 | 1.907 | 1.952 | 2.109 |

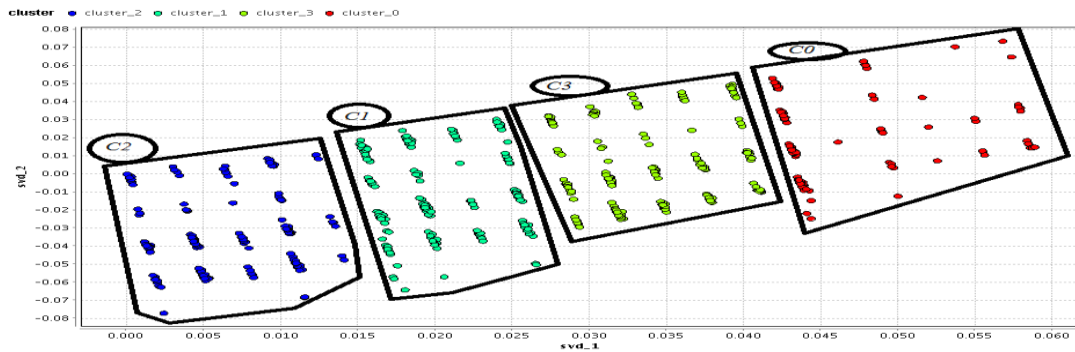Fig. 4 Resulting clusters after applying the k-means algorithm.

Fig. 5 Clusters distribution plot with SVD applied.

Figure 5 depicts a graphical representation of the clusters distribution after applying the Single Value Decomposition (SVD) method, which reduces the number of attributes to two in order to easily plot the clusters.

### 4.4. Outlier method

A database may contain data objects that do not comply with the general behavior of the data and are called outliers. The analysis of these outliers may help in fraud detection and predicting abnormal values [12], so in this paper, we apply two outlier methods, one of which is based on distance-based approach, the other is the density-based approach.

### 4.4.1. Distance-based approach

A popular method of identifying outliers is by examining the distance to an example's nearest neighbors, and the result of applying this method is to flag the records either to be the outlier or not, with true or false [10]. Figure 6 depicts a graphical representation of the region of detected outliers, after applying (SVD) method.
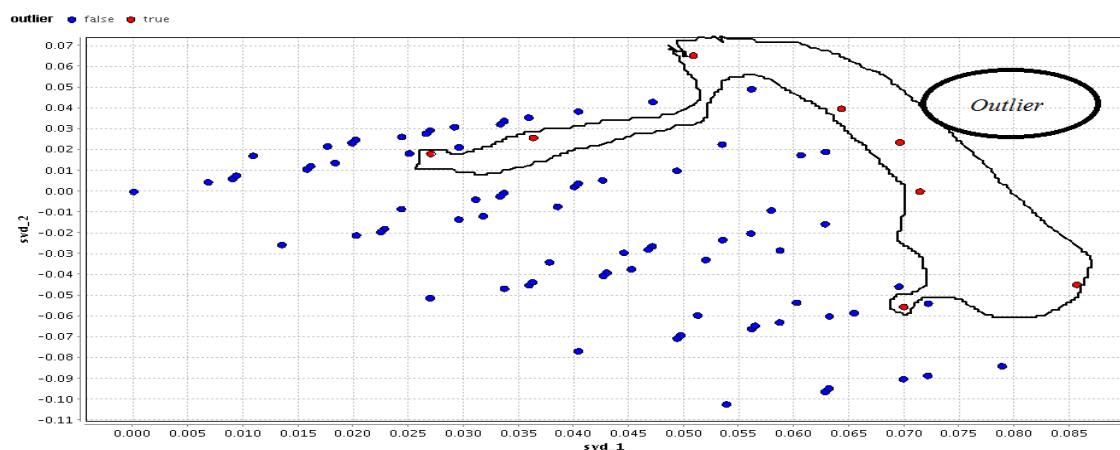


Fig. 6 Outlier (Distance-Based) plot with SVD applied.

### 4.4.2. Density-based approach

Compute local densities of particular regions and declare instances in low-density regions as potential outliers [12]. Figure 7 depicts a graphical representation detected outliers, by using Local Outlier Factor (LOF) approach, after applying (SVD) method, which represents too large percentage of the outlier.
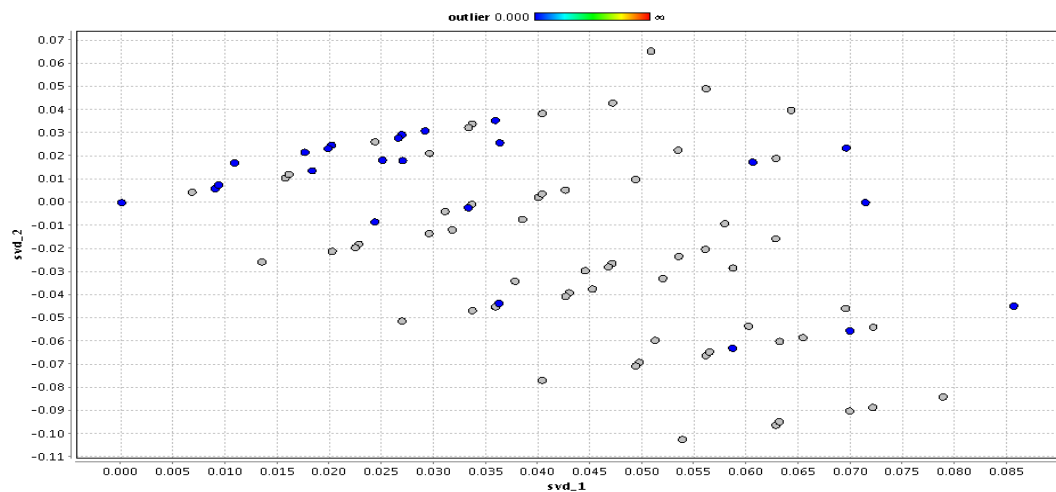
Fig. 7 Outlier (LOF) plot with SVD applied.

## V. RESULTS DISCUSSION AND ANALYSIS

The data set of 1173 students used in this study was obtained at University College of Science and Technology (UCST) – Khan Younis enrolled student.

When analyzing the results of applying the association and classification methods, it is apparent as discussed in sections 4 that there exists a direct relationship between the grade and the specification. The classification model can be used to predict the categories of the target class of grade. The accuracy of the models suggests using Rule-Based classifier for predicting the grade since it has the accuracy of 61.88%, while decision tree has the accuracy of 41.88%, the Naive Bayesian algorithm has the accuracy of 58.50%, and K-NN algorithm has the accuracy of 54.37%. The low accuracy values of the classification models that resulted would recommend using alternate classification algorithm or enhance the quality of data set during the preprocessing phase.

When analyzing clustering and outlier results, the same knowledge can be induced by both methods. Figure 5, 6 and 7 present graphs resulted from the clusters and outlier methods, it is apparent that the cluster labeled "cluster_0", which represents 29.67% of the total records. Also, the outliers based on distance presented in Figure 6, and the outliers based on density presented in Figure 7, which represents a too large percentage of the outlier.

## VI. CONCLUSION AND FUTURE WORK

In this paper data mining methods are applied at University College of Science and Technology (UCST) – Khan Younis enrolled students. The goal of this paper is to present how data mining can help solve low students' performance in UCST, by discovering patterns and building a classification model that relates grade as well as other affecting factors. Knowledge discovery processes were applied which included pre-processing, data mining, patterned evaluation, and knowledge representation. The classification model has resulted which is a basis for predicting students' performance, but the model's accuracy is not high. The paper addresses the issue of data quality as the main limitation that undermines producing an accurate classification model or generating useful patterns out of the data mining methods.

In future experiments, we want to measure the compressibility of each classification model and use data with more information about the students (i.e. profile and curriculum) and of higher quality (complete data about students that have done all the course activities). In this way, we could measure how the quantity and quality of the data can affect the performance of the algorithms.

## References

[1]   AbuTair M., and El-Halees A., "Mining Educational Data to Improve Students' Performance: A Case Study", *International Journal of Information and Communication Technology Research*, Volume2 No.2, February2012, ISSN2223-4985.

[2]   Ayesha S., Mustafa T., Sattar A., and Inayat M., "Data Mining Model for Higher Education System", *European Journal of Scientific Research*, vol. 43, no. 1, pp. 24-29, 2010.

[3]   Paris I., Affendey L., and Mustapha N., "Improving Academic Performance Prediction Using Voting Technique in Data Mining", *World Academy of Science, Engineering and Technology*, 2010.

[4]   Baradwaj B., and Pal S., "Mining Educational Data to Analyze Student s' Performance", *International Journal of Advanced Computer Science and Applications*, vol. 2, no. 6, pp. 63-69, 2011.

[5]   Bekele R., and Menzel W., "A Bayesian Approach to Predict Performance of a Student (BAPPS): A case with Ethiopian students", *in Proceedings of the International Conference on Artificial Intelligence and Applications (AIA-2005)*, Vienna, Austria, 2005.

[6]   Bidgoli M., Kashy D., Kortemeyer G., and PunchBehrouz W., " Predicting Student Performance: An Application of Data Mining Methods with the Educational Web-based System lon-capa", *33rd ASEE/IEEE Frontiers in Education Conference*, IEEE, 2003.

[7]   El-Halees A., "Mining Students Data to Analyze Learning Behavior: A Case Study", *The 2008 International Arab Conference of Information Technology (ACIT2008) – Conference Proceedings, University of Sfax*, Tunisia, Dec 15- 18, 2008.

[8]   Han J. and Kamber M., "Data Mining: Concepts and Techniques", *Morgan Kaufmann*, 2000.

[9]   Kumar V., and Chadha A., "An Empirical Study of the Applications of Data Mining Techniques in Higher Education", *International Journal of Advanced Computer Science and Applications*, vol. 2, no. 3, pp. 80-84, 2011.

[10]  Mannila H., "Data Mining: Machine Learning, Statistics, and Databases", *IEEE*, 1996.

[11]  Nghe N., Janecek P., and Haddawy P., "A Comparative Analysis of Techniques for Predicting Academic Performance", *ASEE/IEEE Frontiers in Education Conference*, pp. T2G7-T2G12, 2007.

[12]  Olson D., and Delen D., "Advanced Data Mining Techniques ", *ISBN 978-3-540-76917-0*, 2008.

[13]  Romero C., and Ventura S., "Educational Data Mining: A Survey from 1995 to 2005", *Expert Systems with Applications (33)*, pp. 135-146, 2007.

[14]  Romero C., Ventura S., and García E., "Data Mining in Course Management Systems: Moodle Case Study and Tutorial", *Computers & Education*, vol. 51, no. 1, pp. 368-384, 2008.

# Comparison of
# Software for Medical Segmentation

Zuzana Kresanova, Jozef Kostolny

*Abstract*—Modern medicine cannot do without the use of computer technology. For the processing of image data in medicine, multiple procedures and one of them are medical data segmentation. In this paper, we focus on an overview of individual medical image processing techniques, and the application of image processing algorithms in tools used in medicine.

*Keywords*—segmentation, medical data, software tool, 3d models

## I. INTRODUCTION

More and more industries see a better future with more computing. They are also fully aware of this fact in the field of medicine. Medical software is increasingly helping doctors work. It is an effective tool for teaching new doctors, but also for diagnosis. It also helps in training heavy interventions on the human body. Because the human body is very sensitive, and misconduct might have fatal consequences, emphasis is placed on displaying 3D models of the human body as closely as possible. At present, there are tools that can create 3D models with millimeter accuracy. However, getting this model is not easy.

Creating a 3D model is preceded by getting data. The most common apparatus for obtaining images are CT, MRI, X-ray and PET instruments. These tools provide 2D images in various file formats. Various algorithms for image data segmentation are used to obtain the 3D model. In the work, we present a few basic principles of segmentation.

The principles are further implemented and concrete implementations of segmentation algorithms are created. These are implemented in multiple successful medical tools that are freely available to the general user, so we will not focus on implementing them, but comparing the resulting 3D models.

To evaluate 3D models, we will set up an evaluation methodology that will focus on evaluating the tools we will work with and the visual results of segmentation algorithms.

## II. 3D SEGMENTATION OF MEDICAL DATA

Nowadays, with high-quality data processing technologies, various methods for 3D tissue reconstruction are widely used in the medical field. One of the essential parts of treatment is not only diagnosis but also the determination of targeted therapy. On this basis, increasing emphasis is now being placed on planning treatment using accurate navigation of the therapeutic tool, or by setting the patient's position directly based on reconstructed CT or MRI data. At the same time, care is taken to minimize and optimize damage to healthy tissues in the patient.

The process of creating a 3D model requires not only CT data but also depth information. Various tools such as depth cameras and laser are used to obtain this information. This information can then create a depth map. Thus, the depth map is a set that content the spatial position of individual points in a three-dimensional scene. The complexity of 3D tissue reconstruction is mainly due to the complexity of the human body, which consists of different

Z. Kresanova, Faculty of management science, University of Zilina, Zilina, Slovakia
J. Kostolny, Faculty of management science, University of Zilina, Zilina, Slovakia, (email: jozef.kostolny@fri.uniza.sk)

structures with irregular shapes and sizes. The problem also arises when the human body contains foreign bodies that appear in images. These are, for example, implants or metal plates. For 3D data reconstruction, one of the following methods can be used:

1. Surface reconstruction
2. Bulk reconstruction

### A.  Surface Reconstruction

The density threshold is selected by the physician itself (e.g., selects the density that corresponds to the bone density). Thus, we can create multiple 3D models in which we display bodies based on different thresholds, i.e. Each type of tissue displays different color because bone, muscle or even cartilage have different density. The 3D data then searches for edges and points of the surface. Of these, the surface is interpolated by two-dimensional patches.

### B.  Bulk Reconstruction

Compared to the surface, the volume reconstruction is not limited by the display of the surfaces themselves. With this kind of reconstruction, the transparency and colors are used to illustrate the better volume in the images. This allows e.g. the pelvic bone appears as semi-transparent or even at an oblique angle.

In the first step, 3D reconstructions of medical data are looking for points that are uniquely identifiable between multiple input data (multiple images). Subsequently, the key points are paired within each segment. This includes calculating 3D coordinates of points from found correspondences. The next step is to create a 3D model and edit it (clean holes). The creation process can be seen in Fig. 1 [1].
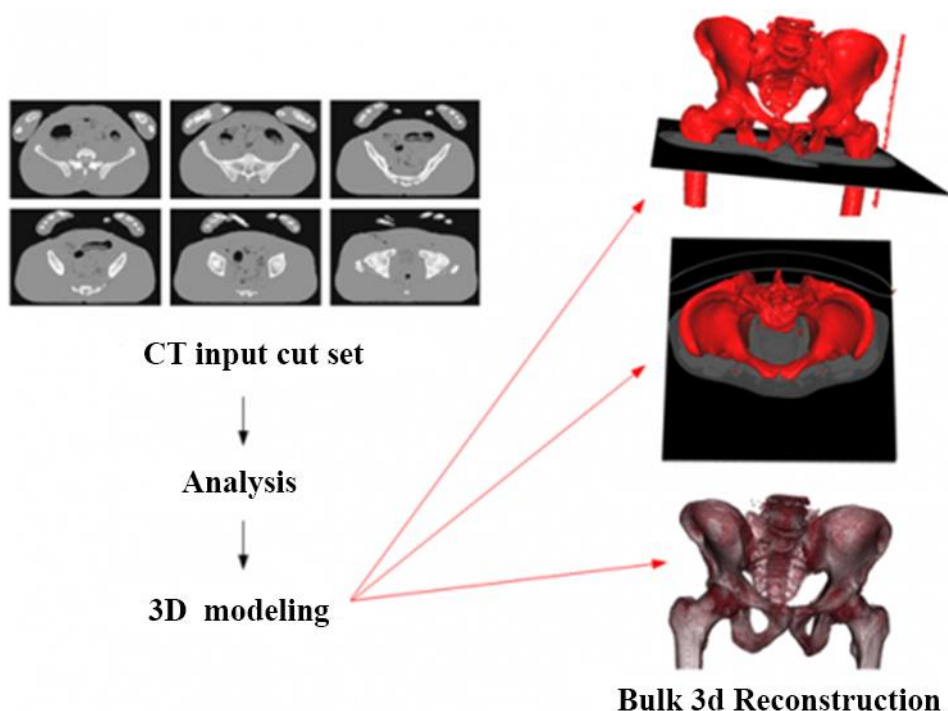


Fig. 1 The procedure for 3D data reconstruction [1]

### III. DATA RETRIEVAL APPARATUS

The first step in creating a 3D model of a human body from medical image data is to obtain it. Various medical devices exist to create these data. In this section we will describe some of the basic devices for these purposes.

Nowadays, the processing of medical data has more and more possibilities due to rapid technology development as well as the advancement of digitalization of information in the world.

In medicine, the worldwide DICOM standard is used for data storage [2]. There are currently several devices that can display image information from the human body.

The most common are X-ray machines, computer tomography, ultrasound devices, as well as magnetic resonance imaging. X-ray machines produce images that are flat and cannot be seen on either the front, the middle, or the back. This problem of spatial imaging has been improved by computer tomography. The human body will be taken from different distances in a row, and thus we can determine the arrangement of individual organs in space. As a result, the human body can be seen from the side, from the front and from the top downwards [3].

### A. Techniques for obtaining flat images

Among the best-known X-ray techniques that enable a flat image to be obtained are Skiaskopia and Skiagrafia. Skiascopy is a radiological examination method that uses X-rays to record a three-dimensional object in a two-dimensional image. Created images are projected in real time, which means that it works as a camera that uses X-rays to display it. The main drawback of this imaging method is a higher radiation load. It is mainly used in the investigation of the gastrointestinal tract, but has significantly decreased its importance by the onset of computed tomography and magnetic resonance.

Skiagraphy or even imaging is a basic radiological imaging method that uses X-rays to photograph parts of the human body. Compared to fluoroscopy, it has a more detailed and detailed picture and produces a lower radiation load on the human body. This method is most often encountered in the human skeleton scan (Fig. 2), but also in the detection of teeth (teeth).



Fig. 2 X-ray images of the left foot

### B. Computer Tomography (CT)

Computed tomography is a radiological examination method that displays the inside of the human body using X-rays. The device performing this examination is a computer tomograph (Fig. 3).

Fig. 3 CT snapshots

The mathematical procedure called tomographic reconstruction is used for recording. It uses inverse radon transformation.

It is a transformation that assigns the original function to Radon's image. The radon image is obtained by means of an integral transformation, which assigns to the real function f defined on the n-dimensional real space another function bearing information about the integrals of the function f through all the affine superficial spaces on which it is defined [4].

*C. Magnetic Resonance Imaging (MRI)*

Magnetic resonance imaging is a radiological imaging method that uses a strong and homogeneous magnetic field. Unlike CT, MRI does not use X-rays but ionizing radiation. Despite good X-ray control (the risk of exposure), MRI scanning is now seen as a better choice than CT.
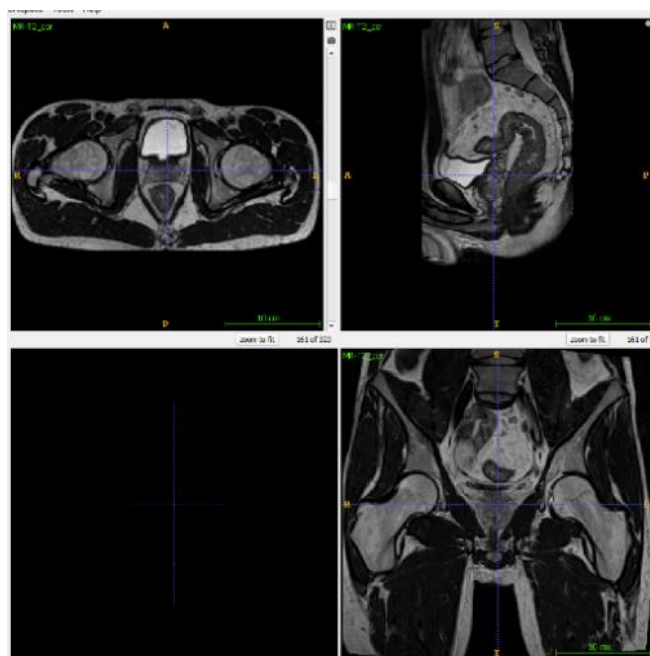


Fig. 4 MRI images

Unlike CT scanners, MR is noisier and can cause discomfort in patients because sensing is more time-consuming, is louder, and is usually inserted into a narrow tube. It is also not suitable for people who have a medical implant or other non-removable metal parts inside the body. MRI testing may not be safe for them if a high magnetic field device is used. At that time, the metal parts can be attracted by the magnet and can cause burns. This risk also arises when a

patient has a tattoo or permanent makeup that contains metal particles. Magnetic resonance imaging is suitable for recording changes and abnormalities in soft tissue, but it can also show the bone tissue of a patient well, see. Fig. 4 As with CT instruments, it has a high resolution and has no problem capturing small details [6].

### D.  *Ultrasonography (USG)*

Ultrasonography is a diagnostic imaging technique that uses ultrasound to visualize the human body. The ultrasonography sends high frequency ultrasound waves to the patient. Tissues in the body have different acoustic impedances, so the ultrasonic wave propagation rate is not the same in all tissues, and it is therefore possible to capture reflected waves. The ultrasound is transmitted in millisecond pulses and the instrument registers the intensity of the reflected signals and the time it takes to return the pulse to the sensor. Ultrasonic wave properties also include the fact that its intensity decreases exponentially, therefore the detected signal is further adjusted for good visualization. The adjustment consists in amplifying the signal in proportion to the time elapsed since the signal was transmitted. Ultrasonography can be performed differently. We know for example A mode, B mode, M mode, but also 2D and 3D mode. A mode is a one-dimensional view in which the amplitudes of the reflected signals are displayed on the screen and seen as a curve at the output. In this mode, distances are accurately measured. In medicine, a B mode is often used, which is a one-dimensional image in which the reflected signal pulses are displayed as different shades of gray. Then we can observe a line segment consisting of pixels with different degrees of brightness. This mode is the basis for other display modes. M mode represents the display of data from B mode in succession over time. It allows moving structures to be displayed and is most often used for cardiac examination. 2D mode is basic. One-dimensional images in mode A, B and M can also be obtained from this image. This image is obtained as a row next to one-dimensional lines of one-dimensional representation in B mode. It is used for the examination of human internal organs. The 3D mode is the latest kind of display that is created as a reconstruction of a series of two-dimensional cuts in a row. For this reconstruction information on the location of the cuts is needed. An example of images displayed in 2D and 3D mode can be seen in Fig. 5
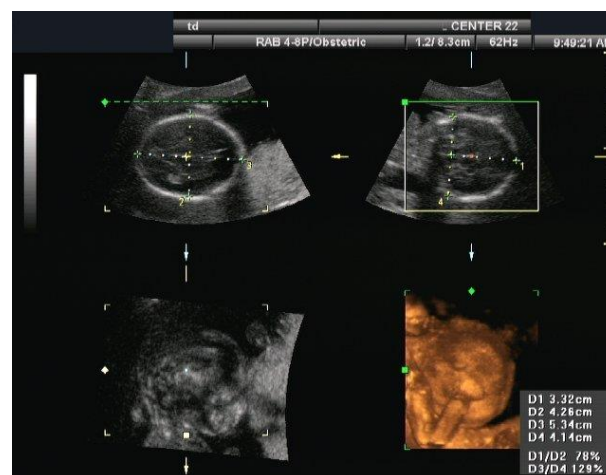
Fig. 5 Ultrasonograph images

### E.  *Positive emission tomography*

Positron-emission tomography (PET) belongs to imaging methods of nuclear medicine. It is based on the principle of annihilation, radiation detection and subsequent assembly of the cut in a given plane.

Annihilation is a process in elementary particle physics. This process can occur when a particle collides with its antiparticle, whereby the original particle disappears and its mass is transformed into some form of energy. In PET, this process is used to fuse the positron with the electron [7]. This device records inflammatory tissues very well and localizes tumor deposits. The most common use of this device is in collaboration with a CT device, when the light generated is precisely identified in the human body. On Fig. 6 we can see the difference between CT and PET images as well as the resulting image after joining these data [8].
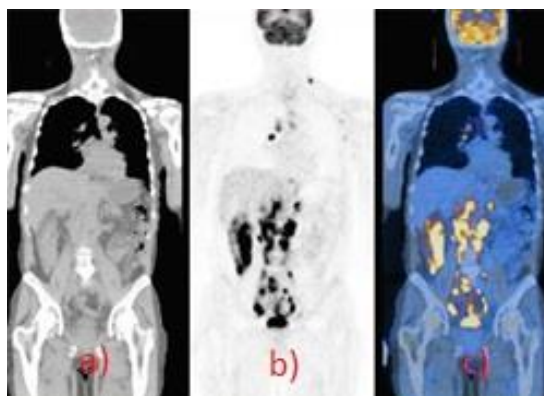


Fig. 6 a) CT scan b) PET scan c) CT + PET scan

## IV. DATA RETENTION FORM

We need to keep the acquired data in some data form. The data we receive from medical devices is not effectively stored in basic formats (JPEG, PNG) because they contain additional information that cannot be written to them. Data standards and advanced file formats are used.

Data standards are important tools for maintaining information systems. This enables information and communication technologies to be used effectively across different health areas.

The standard is essentially a set of rules linked to the creation, development and use of information systems, which includes characteristics, methods, procedures and conditions, especially in terms of security and integrability with other systems. The issue of standards and standards for health care in Slovakia falls within the competence of the Ministry of Health of the Slovak Republic.

### A. DICOM

DICOM [9] is an abbreviation of Digital Imaging and Communications in Medicine and is the communications standard for displaying, storing, distributing, and printing medical data worldwide. DICOM is used by all healthcare facilities, but each in a different range.

In particular, the standard was created to unify the transmission of image and data information between workplaces dealing with different diagnostic methods to display X-ray, MR, CT and PET. The need to standardize the creation and display of medical data has increased mainly due to the transfer of information and images when using devices from different manufacturers. The standard is very extensive.

In addition to data, the DICOM file format also includes a header with additional information describing not only the image itself, but also the patient information, type of examination, time, date of origin, and so on. Thus, the header ensures that the data cannot be exchanged and is uniquely identified. On a similar principle, JPEG also works, which can also store additional information. In each file, however, there is only one data object containing the pixels themselves.

DICOM has expanded over the years to include several guidelines and standards that describe not only image data, but also the transmission, printing and reporting of these data, as unambiguous information objects, among which it also defines interconnections. Supports network operations using TCP / IP protocol. Works with offline media that use industry standards such as CD-R or MOD. It also determines the match level, t. j. declares to the standard implementer the exact structure of the compliance document where it describes each applied function, as well as the properties of the objects. DICOM also specifies the semantics of commands and associated data. It also uses logical file systems FAT 16, ISO 9660, etc. [9].

DICOM consists of services that allow data transmission over network channels. Basic services include: Store - Send object to PACS, Storage Commitment - Retention certificate, Query / Retrieve - Search, Modality Worklist - List of scheduled tasks (examination) with object, Print - Print X-ray image, Off-line media how to store on removable media and more [2].

*B. RAW*

Raw Image File RAW is a "raw" image file from the camera, with only minimal data from the image sensor. These data are not yet ready for printing or editing with the bitmap graphics editor. Normally, the image is processed by a converter in a wide-range internal color space, where precise adjustments can be made before converting to a "positive" format such as JPEG, which is suitable for storage, printing or other manipulation. This color space is often dependent on the device being used.

They are often attributed to the term "digital negatives", although they are not negative images, but perform the same role as negatives from a film camera. The essence is that the negative is not directly usable as an image, but has all the information needed to create the image.

Their task is to capture the best possible radiometric characteristics of the scene (light intensity information and scene color) in the best possible sensor performance. Most often, this information is stored in pixels.

Raw format has several advantages compared to compressed format. For intensity, they have 12 or 14 channels, giving 4096 to 16384 shades compared to 8 bits that use compressed formats. It is because of this feature that they are also very often used to store medical data, as imaging devices provide shades of gray, so multiple objects in the image can be distinguished. This format also provides data transformation, e.g. increase exposure. An important feature of the RAW format is that the changes that are made to the file are non-destructive, that is, only the metadata that controls the rendering changes to different output versions, leaving the original unchanged [10].

## V. Segmentation

When data is processed in certain formats, segmentation algorithms are run on them. These will create 3D images from 2D images so that the Z-axis is also appended to the points on the X-axis and Y-axis, so that the individual frames are stacked on each other to obtain a spatial view. The snapshot order defines the process by which the data was collected. They are stored according to the scan direction of the device. Here's some basic approaches to image data segmentation.

According to [11], the main task of segmentation is to divide the image into parts or categories. We know complete and partial segmentation. When fully segmented, the image is divided into disjoint parts that have a high correlation with real world objects or regions

displayed in the image. They are only partially responsible for partial segmentation. In this kind of segmentation, we divide the image into homogeneous parts in terms of the property we have chosen, e.g. brightness, reflectance, texture, and so on. Most often, we will see segmentation in solving the problem of background image separation from the foreground. If the object being scanned is not on the entire frame, we say that we will focus segmentation on the area of interest. In doing so, the target objects may not always be clear, and the individual types of segmentation must also take into account undesirable effects such as, for example, bugs in the form of noise that we were unable to remove when preparing images before segmentation or parts (e.g., other objects) that are not in the interest of our investigation.

In segmentation, the individual exploration areas can be represented by a closed boundary, and each boundary that is enclosed describes an area. Segmentation methods can be categorized by segmentation method:

### A. Thresholding Segmentation

Tresholding is the basis for a number of segmentation methods, mainly for its simple implementation and computational simplicity. The basic step in thresholding is to choose the right threshold. Later, the image is divided into two groups. The first group contains pixels whose degree of brightness is greater than the threshold value and the second one whose degree is less than or equal to the threshold value. Subsequently, the image is converted into binary form, i.e. black and white image.

According to the number of thresholds per image, we divide the threshold into global and local. Global thresholding is the easiest and most convenient to use when the background image brightness and foreground object intensity is vastly different. At that time, a single threshold is sufficient to divide all the points in the image. The threshold value is obtained from a histogram in which the brightness of each pixel is recorded. This may be automatic or manual. If the contrast between background and objects is low, local thresholding is used. In doing so, the image is divided into smaller parts based on the width of the surroundings we want to work with and the threshold is determined separately for each area depending on the surroundings. Local thresholding is mainly used for images that contain noise. Noise can be created especially when the object is not evenly distributed

### B. Area-based segmentation

Area-based segmentation method must ensure that the criterion for full segmentation is met:

$$R = \bigcup_{i=1}^{S} R_i \quad R_i \bigcap R_j = 0 \quad i \neq j \tag{2.1}$$

and conditions of maximum homogeneity:

$$H(R_i) = \text{TRUE, pre } i = 1, 2, \ldots, S \tag{2.2}$$

$$H(R_i \cap R_j) = \text{FALSE pre } i \neq j \text{ a } R_i \text{ is adjacent to } R_j \tag{2.3}$$

where R is the region designation, H is the homogeneity mark and S is the number of all regions. This method uses 3 basic approaches:
1. Connecting areas
2. Partitioning
3. Partitioning and joining areas

*C. Edge-based segmentation*

The basic principle of segmentation based on edge detection is the discontinuity of image data. Subsequently, borders are detected between regions. In this method, it is important to distinguish two basic terms, which are edge and boundary. An edge is only a property of a pixel that represents its degree (brightness, gradient). Unlike the boundary that is already the result of segmentation, it shows the contour of the segmented area. The boundaries are detected by gradient operators (1st and 2nd derivatives). In places where the level of brightness, color or texture is changed, a local edge is created. The drawback of this method is that it cannot work with noise. This causes detection of local edges even where there is no global boundary and vice versa, sometimes they are missing where the global boundary exists. An important factor is therefore the choice of the threshold height for the local edges, provided that the thicker edges are part of the boundary. In particular, the method achieves valid results when the contrast between the subject and the background is large enough. But this may not apply to local edges, which gives us incomplete boundaries [11].

*D. Clustering Segmentation*

The clustering method is characterized in that the objects are placed in clusters based on the similarity of data properties. The basic clustering algorithms include K-means and K-medoids.

K-means is a two-step method. In the first step, the centroids are created. These are assigned at random positions. The number of centroids created is the same as the resulting number of clusters we require. This point can be moved, unlike a data point that is static. The data point is assigned to the cluster, based on the distance from the centroid. When the assignment of data points to clusters is complete, the second point continues to calculate the centroid distances and they move to the center of the cluster to which they are assigned. This process of assigning data points to clusters and recalculating centroid coordinates is repeated until their position is stable (Fig. 7). The centroids are then removed and only clusters with data points whose properties are similar remain [12].

K-medoids is a method that is an enhancement of K-means. Compared to the previous one, it has less sensitivity to disturbing and inappropriate data that a dataset may contain. Also, the distance is not calculated with respect to the average value, but to the median (nearest value to average value). In the first step, like K-means, we determine the random objects, that is, the number of resulting clusters. We assign data points to clusters based on the smallest distance. In the second step we then select a random median and try to replace it with the current one. The method is similar to the first method, until the median position is stabilized [12].
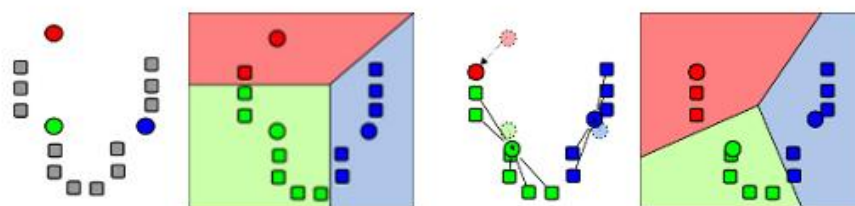


Fig. 7 Method of K-means Method

*E. Livewire Segmentation*

Livewire [13] (magnetic lasso) is a kind of segmentation that anchors segmentation points to edges in the input image. The algorithm works with volume data in sections. The segmentation area of interest is bounded by a boundary. In determining a boundary, a price function is used that describes the two adjacent pixels the price of the edge between them. Segmentation begins by setting the starting point so that the user clicks on the image pixel. We consider this pixel an anchor. Subsequent movement of the mouse creates the smallest path to the pixel where the

current mouse is located. If a user wants to select the path that this movement has to show, the user will click the image again. For a more accurate object contour, it is possible to add another control vertex (another anchor). After joining all the anchors with the edges, an image contour is created that resembles a magnetic lasso (Fig. 8).
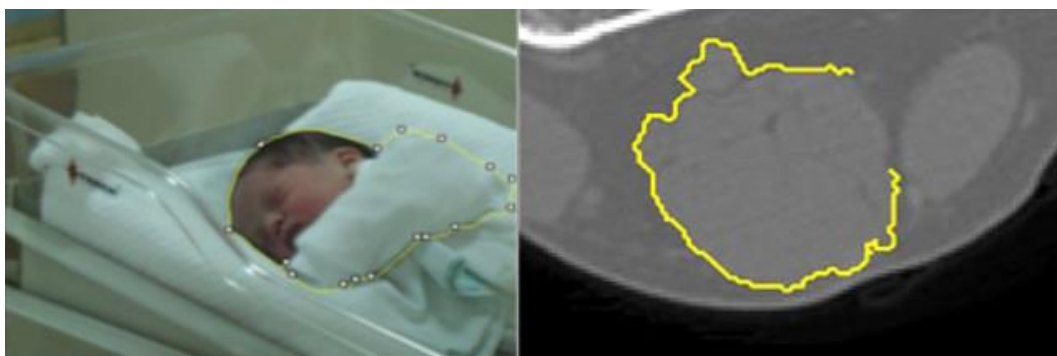


Fig. 8 Livewire selection

## VI. METHODOLOGY OF SEGMENTATION ALGORITHMS EVALUATION

For the segmentation of image data, not only medical ones, there are several basic approaches that we have described in the previous chapter. By implementing these approaches, it is possible to obtain specific algorithms that exhibit different properties when segmented. In this work we focus on evaluation of properties of these algorithms. To this end, it is necessary to design a methodology that will allow us to do so.

Medical image data segmentation algorithms are implemented in multiple successful software tools, and therefore their implementation would be rather unnecessary. We prefer to focus on existing software and explore the algorithms implemented in it. We will evaluate the work with existing software as well as the resulting segmentation images of individual algorithms on different data.

For tools, we will first monitor whether the tool is free or is paid for. Another aspect will be open-source code, number of implemented segmentation algorithms, support for input / output formats, availability and dossier development, and finally the complexity of tooling. In segmentation algorithms, we will evaluate the rate of segmentation and the resulting 3D model, using the algorithm. Our goal will be to segment the support system. We will observe whether the resulting segmentation contains other human tissues, but also whether the segmentation has captured the entire area of the support system shown in the images (empty bones). Based on this, we will evaluate the individual algorithms according to Tab. I.

TABLE I
METHODOLOGY OF SEGMENTATION ALGORITHMS EVALUATION

| Good | The 3D model does not contain any (or very few) tissues other than the tissues of the support system. Object contours are clear and objects in the image are easily recognizable. The algorithm speed is appropriate to the results. |
|---|---|
| Average | The 3D model contains tissues other than the support system tissue, but the contours of the objects are not quite clear and the objects in the image are more difficult to recognize. Segmentation required more time with respect to the result. |
| Poor | The 3D model contains other tissues, such as tissue of the support system. Object contours are not clear and objects in the image are difficult to recognize. Segmentation required too much time due to the result. |

## VII. Software Tools for Medical Data Segmentation

The implementation of segmentation algorithms is dealt with by several companies or institutions. It has been found that effective segmentation of image data can greatly aid in the treatment of the patient, therefore the emphasis is on the most confidential display of reality. This largely affects the correct selection of the algorithm. In larger projects, several methods are used simultaneously for segmentation to achieve the best result. The 3D models thus obtained are often more sophisticated than models that were created using only one segmentation method. Many tools that implement visualization and segmentation algorithms are freely accessible. The best known are Medviso [14], 3D Slicer [15], ITK-SNAP [16], MeVisLab [17], TurtleSeg [18], ImageJ [19]. However, there are many tools using different combinations of algorithms. In this work, we focused on 2 tools, namely ITK-SNAP and MeVisLab, where ITK-SNAP is a tool that is also suitable for a user who does not have great knowledge of advanced data preparation and subsequent segmentation. MeVisLab is a representative of an advanced tool and some knowledge is needed to work with it.

### A. ITK-SNAP

ITK-SNAP [16] is a software application that is used to segment structures in 3D medical images. It was created in collaboration with the University of Pennsylvania, Utah University and Paul Yushkevich. This tool is multi-platform and free. ITK-SNAP provides semi-automatic segmentation using active contour methods as well as manual delimitation and image navigation. In addition to basic features, it also offers other support tools.

Manual segmentation provides three orthogonal planes at a time. It also supports various input formats (DICOM - image series, DICOM - single image, GE, GIPL, MetaImage, NIfTI, Raw, Siemens Vision, VoxBo CUB, VTK Image), but also outputs (JPEG, PNG, TIFF, VTK Image etc) . The app contains a linked cursor that makes 3D navigation easier. 3D images can be color coded and explored using a function that can cut a plane. The graphical interface preview of this tool is shown in Fig. 9.
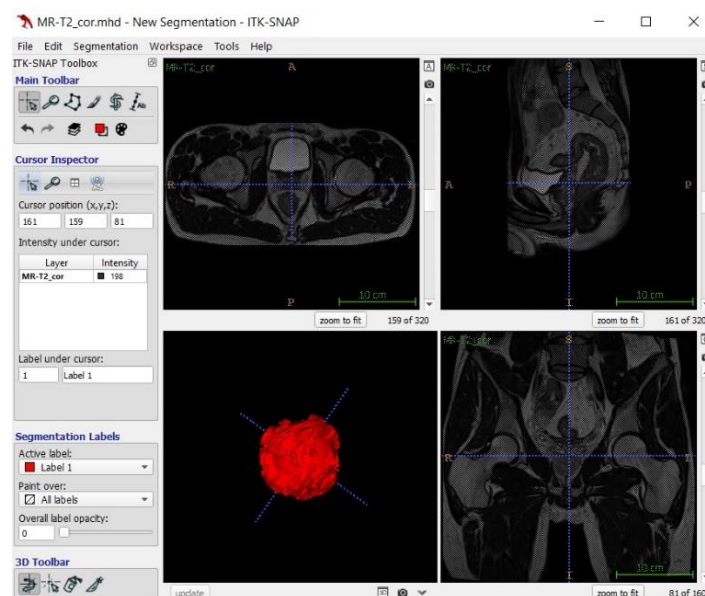


Fig. 9 ITK-SNAP tool

ITK-SNAP has implemented 4 kinds of segmentation algorithms that can be used for semi-automatic as well as manual segmentation.

The segmentation methods they use are:
1. Thresholding - Provides two-page thresholding (Segmentation Threshold is closed from top and bottom) or only from one side (Top / Bottom Threshold setting only)
2. Clustering - ability to determine the number of resulting clusters (parameter k)
3. Edge Based Segmentation - Edge smoothing parameter
4. Classification

ITK-SNAP developers focused primarily on image segmentation. The other unrelated functions are developed only minimally. Application orientation is intuitive, also provides training videos for easier work and user experience. It also offers a variety of color options to display the segmentation area.

### B.  MeVisLab

MeVisLab [17] is developed by MeVis Medical Solutions AG in collaboration with Frauhofer MEVIS. The application is a powerful research and development tool for image processing with a special focus on medical imaging. Enables rapid integration and testing of new algorithms. It is also used to develop clinical applications prototypes.
The application includes advanced modules for segmentation, registration, volumetry, as well as quantitative morphological and functional analysis. It is often used for neuro-imaging, dynamic image analysis, operation planning as well as cardiovascular analysis.
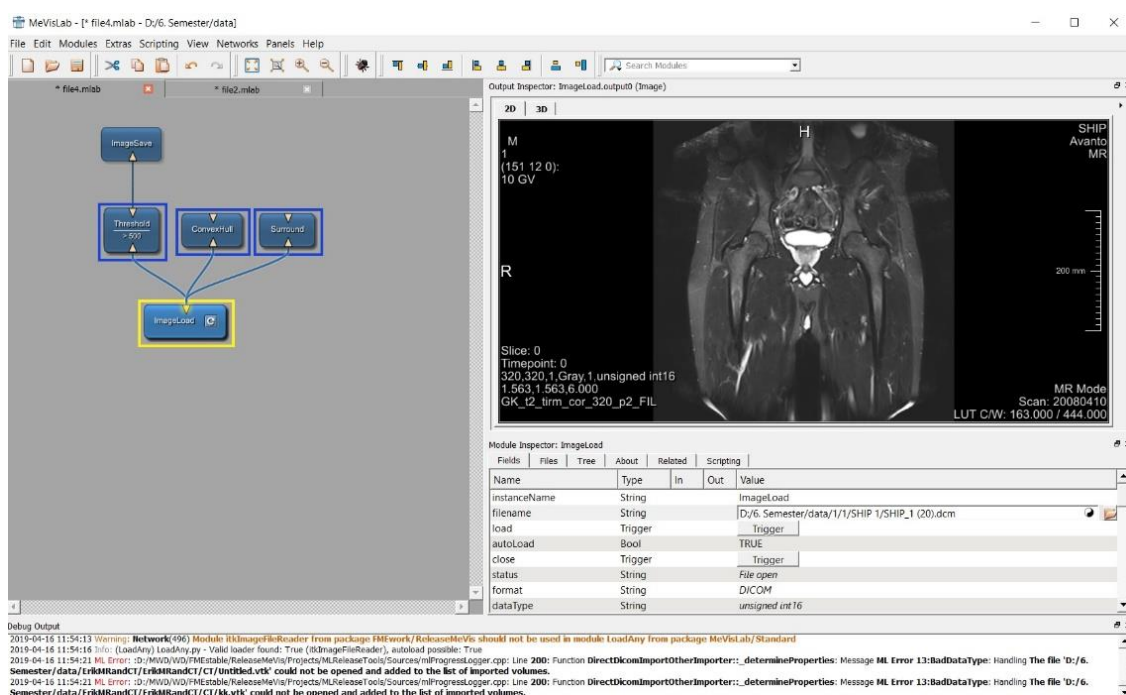


Fig. 10 MeVisLab tool

MeVisLab (Fig. 10) uses known libraries as well as third party technologies. It uses the QT library, OpenGL standard, for the graphical interface and data visualization.

The core component of MeVisLab is the object-oriented MeVis Image Processing Library (ML), which provides general image processing support. A distinctive feature of ML is demand-based image processing control, which is based on multiple multithreaded pages. It uses priority site management to avoid unnecessary data recalculation.

It currently has more than 5000 image processing modules, including image filtering, segmentation, and statistical analysis. They can be divided into 5 main groups: Filters, Segmentation, Transformation, Statistics, and more. Each group contains several hundred

modules. It supports various input data formats (DICOM, TIFF, RAW, JPG, BMP).

Segmentation offers the following implemented methods:

1. Thresholding
2. Segmentation based on areas
3. Livewire segmentation
4. Manual Contours

In this work we used modules mainly from segmentation. Tab. II shows an overview of implemented segmentation algorithms in individual tools.

TABLE II
OVERVIEW OF SEGMENTATION ALGORITHM USAGE

| | MeVisLab | ITK-SNAP |
|---|:---:|:---:|
| **Thresholding** | ✓ | ✓ |
| **Area-based segmentation** | ✓ | ✗ |
| **Edge-based segmentation** | ✗ | ✓ |
| **Clustering** | ✗ | ✓ |
| **LiveWire** | ✓ | ✗ |
| **Classification** | ✗ | ✓ |

## VIII. PRACTICAL TESTING OF SEGMENTATION ALGORITHMS

We used CT and MRI data to test algorithms and compare their properties. The CT of the device shows images of the left knee and pelvic area. The knee images consist of 26 and 28 frames (different viewing angles) where the distance of each frame is 4.4 and 3.9 mm. The pelvic imagery consists of 40 images that are 6mm apart. The frame distance determines the width of the cut. MRI images of the pelvic area. These images are 160 (front-to-back view) and 320 (side-to-side and top-down view) The images we use are DICOM or RAW.

The first tool for segmentation in which we performed semgneting was ITK-SNAP. Navigation in it is very simple and intuitive. At least minimal knowledge of the English language is assumed when using this tool. The calculations themselves were quite fast, but this depends on the segmentation method used and the number of iterations triggered. In this case, the iterations are considered to be a sequence of steps, that is, how many times the given algorithm is run at a certain parameter setting over the data. This means that we can combine segmentation algorithms and run them on the same data several times, while the individual segmentation algorithms can be re-calculated multiple times on a single run. In general, the more iterations we run, the more accurate segmentation is. The resulting evaluation of segmentation algorithms is recorded in Tab. III.

From the table we can see that the tools that were selected to compare segmentation algorithms were representatives of two groups. ITK-SNAP was a tool that is simple for a user who does not have a deeper knowledge of segmentation. Best results were obtained with CT images using a threshold algorithm where the threshold was set manually. Comparable results were also observed when using the clustering and classification algorithm in MRI images.

MeVisLab is an advanced tool for image data segmentation and visualization that is very comprehensive and provides a large number of modules that offer data pre-processing, data processing and data editing. Some knowledge of medical image segmentation is needed to work with it. Just because of this, it is already suitable for an advanced user. In MeVisLab, we investigated the threshold method and the effect of threshold value on each frame type. The 3D model, which was very close to reality, came from images created by a CT machine. The threshold value was set to 450.

The segmentation algorithms used showed different results using different data. The main

factors that influence the resulting 3D model include the quality of data collected from medical imaging devices. With CT knee shots, we can observe poor results even with manual but also automatic threshold values. This is probably due to the low number of frames, because using both the CT threshold method and the MRI in the ITK-SNAP tool, we achieved average and good segmentation model results, with the aforementioned CT and MRI data being known to be of higher quality and greater as knee images.

TABLE III
SEGMENTATION PERFORMED STATS

| Segmentation order | Tool | Segmentation algorithm | Data | Iteration count | Rating |
|---|---|---|---|---|---|
| 1. | ITK-SNAP | Thresholding | MRI | 4x150 | Average |
| 2. | ITK-SNAP | Thresholding | MRI | 4x150 | Poor |
| 3. | ITK-SNAP | Zhlukovanie | MRI | 4x1500 | Good |
| 4. | ITK-SNAP | Edge Segmentation | MRI | 4x150 | Poor |
| 5. | ITK-SNAP | Classification | MRI | 4x3500 | Good |
| 6. | ITK-SNAP | Thresholding | CT | 2x150 | Poor |
| 7. | ITK-SNAP | Thresholding | CT | 5x900 | Good |
| 8. | ITK-SNAP | Thresholding | CT | 4x200 | Poor |
| 9. | MeVisLab | Thresholding | MRI | 1 | Poor |
| 10. | MeVisLab | Thresholding | MRI | 1 | Average |
| 11. | MeVisLab | Thresholding | MRI | 1 | Average |
| 12. | MeVisLab | Thresholding | MRI | 1 | Poor |
| 13. | MeVisLab | Thresholding | CT | 1 | Poor |
| 14. | MeVisLab | Thresholding | CT | 1 | Average |
| 15. | MeVisLab | Thresholding | CT | 1 | Good |
| 16. | MeVisLab | Thresholding | CT | 1 | Good |
| 17. | MeVisLab | Thresholding | CT | 1 | Average |
| 18. | MeVisLab | Thresholding | CT | 1 | Good |
| 19. | MeVisLab | Thresholding | CT | 1 | Good |
| 20. | MeVisLab | Thresholding | CT | 1 | Good |

In MeVisLab, we can also see that thresholding achieves different results with the same data (CT images of knee and CT scan of pelvis and chest). These differences can be attributed to data quality (number of frames, accuracy of the device). After the experiments performed, it cannot be stated that CT scan images are more suited to segmenting 3D models than images from an MRI instrument because we achieved good results for both CT and MRI images.

Segmentation algorithms achieved different results, but none of the algorithms used created a 100% 3D model of the human body support system. This problem could be solved by a combination of individual segmentation algorithms that would be used sequentially for different segmentation regions.

## IX. CONCLUSION

The aim of this work was to investigate and compare selected algorithms for medical image data segmentation. In the work we are familiar with the procedure for creating a 3D model of the human body.

The first step is to obtain data from medical devices. We described the scanning procedure of X-ray apparatus, computed tomography, magnetic resonance, ultrasonograph and positron emission tomography. The data produced by these devices is stored in different formats. In addition, we have described the basic medical data storage formats. Most of this data is in the DICOM standard, which also makes data transfer easier. The acquired data is recorded in 2D images from which a 3D model is created using segmentation algorithms. There are several

basic principles for creating a segmentation algorithm. These include thresholding, edge-based segmentation, area-based segmentation, clustering segmentation, and modern livewire segmentation. The described principles implement various software tools. We chose ITK-SNAP and MeVisLab for our research. ITK-SNAP is a segment-oriented tool that is also suitable for users with minimal knowledge of image data segmentation. MeVisLab is an advanced tool that allows you to work more with your data. I.e. not just segmentation but also visualization, transformation, various statistics and others. Both tools have shown good and bad results. These arose depending on different data and setting of input parameters of individual segmentation algorithms. In ITK-SNAP, the manual threshold thresholding algorithm for CT scans, but also the MRI image clustering algorithm, yielded the best results. In MeVisLab, a similar thresholding algorithm with a manually set threshold using CT scan images also achieved similar results. In general, it cannot be argued that the device from which we obtained the data affects the resulting segmentation because we have achieved good results in both CT and MRI images. The implemented segmentation algorithms achieved good results, but this does not mean that they did not have minor deficiencies. These deficiencies could be eliminated by using multiple segmentation methods on a single 3D model.

## REFERENCES

[1] Kamencay, Hudec, 3d rekonstrukcia medicinskych dat, [Online] http://vedanadosah.cvtisr.sk/3d-rekonstrukcia-medicinskych-dat
[2] Majerník, Švída, Majerníková. Medicínska Informatika.
[3] Ako vyzerá naše telo z vnútra [online] https://zdu.uniza.sk/prednasky/nase-telo-vnutri
[4] Radonova transformácia [Online] http://www.veda.sk/?science=6&pojem=Radonova_transform%C3%A1cia
[5] Stoppler, CT Scan [Online] https://www.medicinenet.com/cat_scan/article.htm#ct_scan_facts
[6] Weis, Bořuta. Úvod do magnetickej rezonancie
[7] Šmída. Nadbytek elektronů v kosmickém záření [Online] http://www-hep2.fzu.cz/~smida/www/smida-elektrony+pozitrony.pdf
[8] PET-CT [Online] https://www.massgeneral.org/imaging/services/procedure.aspx?id=2250
[9] DICOM [Online] https://www.dicomstandard.org/ 13
[10] Raw formát [Online] https://digi-foto.sk/zakladne-pojmy/raw-format/ 14
[11] Šonka, Hlaváč, Boyle. Image processing, Analysis and Machine Vision
[12] Michálek, Zhlukovacie algoritmy [online] http://www2.fiit.stuba.sk/~kapustik/ZS/Clanky0809/michalek/index.html#Segmentacne
[13] Alexandre X Falcão, Jayaram K Udupa, and Flavio Keidi Miyazawa. Anultra-fast user-steered image segmentation paradigm: live wire on thefly.
[14] Medviso [Online] http://medviso.com/
[15] 3D-Slicer [Online ] https://www.slicer.org/
[16] ITK-SNAP [Online] http://www.itksnap.org/pmwiki/pmwiki.php
[17] MeVisLab [Online] https://www.mevislab.de/mevislab/
[18] TurtleSeg [Online] http://www.turtleseg.org/
[19] ImageJ [Online] https://imagej.nih.gov/ij/