# Effects of a Nano-structured Surface Layer on Titanium Implants for Osteoblast Proliferation Activity

**Árpád Joób-Fancsaly[1], Albert Karacs[2], Gábor Pető[2], Kinga Körmöczi[1], Sándor Bogdán[1], Tamás Huszár[1]**

[1]Semmelweis University, Department of Oral and Maxillofacial Surgery, Mária u. 52, H-1085 Budapest, Hungary

[2]Research Institute for Technical Physics and Materials Science, Hungarian Academy of Sciences, P.O. Box 49, H-1525 Budapest, Hungary

E-mail: joob_fancsaly.arpad@dent.semmelweis-univ.hu

*Abstract: The goal of this study is to compare surface morphologies on dental implants created by a range of five different surface-modification technologies and, in addition, cell assays to assess the subsequent cell proliferation on each treated surface. In our research, we surface-treated 5 mm-diameter – 2 mm-thick discs, machined from Grade 1 titanium. We treated the surfaces of the discs with chemical etching, electro-polishing, $Al_2O_3$ sand-blasting, and surface melting with 1 Joule or 3 Joule impulse-energy laser beam. We carried out quantitative as well as qualitative analyses with stereo and scanning electron microscopes (SEM), confocal and atomic force microscopes (AFM) and goniometer. We examined each surface with cell-testing, as a measure of osseointegration. In tests with fibroblasts, the highest cell proliferation occurred on the $Al_2O_3$-roughened surfaces. In the case of osteoblasts, we measured the greatest cell activity on the laser-melted samples with different energy levels.*

*Keywords: surface modification; surface analyses; surface morphology; cell proliferation*

## 1 Introduction

Currently the examination and deliberate structuring of surface morphologies of dental implants, designing new processes for surface manipulation and the assessment of those results, are all at the cutting edge of implantology. Surface treatment of implants influences, to a major extent, the success of integration with bone [1], making these processes crucial to the functional value of implants. To ensure this integration, a suitable morphology must be created on the surface of a given implant [2].

## 1.1    Conditions for Osseointegration

Thirty years ago researchers considered how surface engineering, topography, and the morphology or lack of contamination of implant surface, might affect osseointegration. Since that time, investigators have made substantial progress in making issues of surface shape and contamination central to dental implantology research. In the early 1970s, development has sought to: understand the interaction between the bone and the implant surface, explain the process of osseointegration from the moment of implantation until the phase when secondary stability is attained, and define both the duration of, and the biological mechanisms involved in, each phase [3, 4, 5, 6].

Exact phases of the development of the interface zone are still unknown. The interface layer ensures a connection between the oxide layer on the implant surface and bone proteins. The layer around the osseointegrated implant has a thickness of 2 to 5 nm [7, 8]. By the beginning of the 1970s, Hulbert et al. had already established that porous surfaces enable bone regeneration to occur more quickly and allow osseous tissue to grow [9]. The consensus since the year 2000 has been that rough surfaces are better for osseointegration than smooth or polished surfaces [10-15].

Molecular biology has gradually gained significance in surface treatment. Researchers are developing increasingly sophisticated surface-engineering techniques that take into account biochemical signals and cues involved in how various materials bond with bone tissue. The possibility of applying organic materials (chemical modifications) and several proteins (BMP) on an implants surface, creates new areas of investigation [16]. The interaction of proteins at the nanometric level is emerging as a major decider in the integration of implants.

Up until the end of the 1990s, investigations mostly addressed surface topography and morphology. Since then, attention has increasingly focused on surface-chemistry research. Surfaces are modified with physical-chemical methods to promote faster and more comprehensive osseointegration. High surface energy optimally promotes interaction with the biological environment [17-19]. It is accepted that surface roughness influences alkaline phosphatase (ALP) and osteocalcin (OC) levels, thus indicating the osteoblast activity of the surrounding tissues [20, 21]. On a smooth surface, "osteoblast-like" cells show a sparse and flat distribution, while greater numbers and more dense distributions are found on rough surfaces [22, 23]. Rough surfaces have an impact on cytoskeleton functions as well [24]. Micro-scale surface structures influence cellular functions while nano-scale elements have an effect on sophisticated cascades involved in protein binding (Figure 1) [25-27].
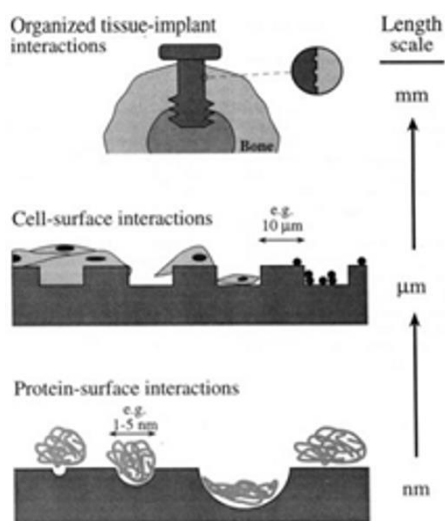
Figure 1

The effect of surfaces with various degrees of structure on the binding of organic elements

According to the current status of the science for osseointegration, it is not a function of quantitative porosity but qualitative porosity [28]. In the case of dental implants the most commonly used and at least accepted parameters are the following: the average height differences of elements raised above the surface, usually represented as "Sa" (in the case of a two-dimensional examination the same value is signified by "Ra"), effectively the average of the wave-lengths (distances), which is named "Scx", and finally the sum of treated and untreated surfaces, which is called Sdr and is described as a "hybrid" value. According to Wennerber and colleagues as well as most of the research community, what can be described as ideal surface parameters are $Sa = 1.4$ µm, $Scx = 11.6$ µm, $Sdr = 1.5$ µm [29-33]. On the basis of the literature a surface can be viewed as rough when the hybrid value is >2 µm. It is medium rough if that given value falls between 1 µm and 2 µm, and smooth, or at least mildly rough if <1 µm [34]. Wennerberg and colleagues showed with animal experiments the 1.5 µm hybrid value is the most favorable and associated with cases of the greatest osseointegration [35]. Lazzara and colleagues showed this and Ivanov was doing human experiments supported with results involving human volunteer subjects [36, 37].

It's possible to characterize the quality of a surface with bone-implant contact (BIC) values, which show what percentage of an implant surface comes into direct contact with bone. After decalcification, with a histology-morphology examination of a tissue slice it is possible to tabulate the percentage ratio of bone-implant connection (BIC) [38]. For sufficient bone integration the lower jaw BIC average is 40.7%, while in the case of the upper jaw this average is only 37.2% [39].

The percentage ratio of bone-implant contact is not easy to calculate. Results recorded as obtained in various journals in the literature vary widely across a very broad scale. They depend on the species of experimental animal, the type of bone the implant was attached to (femur, tibia, jaw), the healing time, and the surface treatment applied to the implant.

Hansson devised a mathematical model describing interdependences between the measure of osseointegration and the surface roughness of an implant. Surfaces of commercially obtainable implants serve as the basis of the model. He established that hollows of depth 1.5 µm and diameter of 3 to 5 µm in the implant surface are advantageous for integration with bone [40].

## 1.2    Surface-Modification Methods

Current machining methods give dental implants unique geometries. Leading manufacturers most often create implants with diameters from 1 to 4.2 mm and lengths from 2 to 14 mm [41]. The geometric structure and surface-preparation of implants plays a central role in primary and secondary stabilization [42]. The most common methods for modifying surface morphology are chemical etching [43-45], sand-blasting [46], these two in combination [47], grit blasting [48], surface melting by laser [49], and anodic oxidation [50, 51].

A freshly machined surface is often used as a reference in experiments for mutual comparison. Literature discussing machined surfaces describes massive bone formation around the implants, which results in a stable connection between implant and bone tissue [52]. Anomalies on the machined surface go to a depth of 5 µm in profile. On the surface at a spacing of 5-8 µm, slightly unevenly arranged but roughly parallel grooves can be observed. The separation and depth of the grooves was equal. With the help of a profileometer the three most characteristic pieces of data were Sa = 0.836 µm, Scx = 8.38 µm and Sdr = 1.3 µm [35]. In the literature there are reports of quite wide intervals (between grooves) on the machined surface with Ra, in the range 0.08 – 4.7 µm [53].

The BIC ranges between 39 and 47% according to Ericsson and colleagues, as opposed to around 50% according to Wennerbeg et al. [54, 55]. In animal-experiment studies the value has been observed at 62% after 3 months [56]. One drawback is that in several places the surface is contaminated by noticeable machining fragments.

On the basis of the profilometer studies by Wennerberg et al (Sa = 1.9 µm; Scx = 12.3 µm; Sdr = 1.42 µm) the parameters of the acid-treated surface closely approach the supposedly ideal values [35]. The results of animal studies support this, showing a BIC value of 88% after three months of healing [56]. According to Baker et al. animal experiments in the case of acid-treated surfaces bone adhesion starts earlier than with other treated surfaces [57]. Human trials carried out by Trisi and colleagues showed that in case of acid-treated surfaces the BIC

percentage value resembles that of freshly-machined surfaces [58]. They regard the payoff from this method to be the chance to load the implant early on. The other big benefit in each acid-treating process is a high level of surface cleanliness, in that the acid removes the outermost layer of the surface.

In the case of sand-blasting the surface can be blasted with particles of different sizes (25, 75, 250 µm). With animal studies Wennerberg and colleagues showed that smaller particles resulted better osseointegration than larger, and suggested using 75 µm sized granules. According to data obtained with the help of an optical profilometer, granules of the size of 25 µm produced a surface characterized by Sa values = 1.13 µm; Scx = 9.78 µm; and Sdr = 1.39 µm, while the surface produced by 75 µm particle blasting has the following values Sa = 1.38 µm; Scx = 11.62 µm; Sdr = 1.47 µm (although they approach the ideal parameters very closely). In the case of the 250 µm particles the values were Sa = 2.15 µm; Scx = 13.54 µm; Sdr = 1.79 µm. The BIC was 62/62% (the mandible/maxilla ratio). During animal-experiment studies this value rose after 3 months to 71% [14, 35, 56, 59]. The disadvantage of the process is that the material used for the sand-blasting can contaminate the surface, and because of this the chemical features of the implant can be unacceptable.

Buser et al examined surface treatments of the titanium used as dental-implant base material. The samples were treated with chemical etching, electro-polishing, or coated with hydroxyapatite. They established that surfaces treated with chemical etching and hydoxyapatite underwent more bone-integration than the samples with electro-polished surfaces [60].

The development of nanotechnology, analyzing the possibilities of implants with nanostructured surfaces, is at the center of numerous current research-and-development projects. Christenson and colleagues looked at organized nanostructures falling into the 1 to 100 nm size range. The "nano" expression can cover crystal structures of material, extracted or built out from the surface layer [61].

Carlos et al completed studies on surface treatments to dental implants made of Grade 4 titanium. They changed the surface morphology with sand-blasting, chemical etching, and anodic oxidation. They prepared SEM images of the surface and then established the level of surface roughness with confocal microscopy. They tested each surface for wettability with contact angle measurements. The surface-modified implants were inserted into live rat tibias and then after 12 weeks removed. They established that chemical etching created a homogeneous surface. Surfaces treated with anodic oxidation measured the lowest while at the same time those implants had the highest screw-out force values [62].

Luiz's research group demonstrated the surface modification and testing possibilities of titanium as a raw material for dental implants. In their work they compared surfaces modified by mechanical polishing, electro-polishing, nano-hydroxyapatite, and sprinkling with titanium dioxide and fluoride granules. With

an atomic force microscope, a SEM and an X-ray photoelectron spectroscope the surface structures of the samples were tested. Experiments carried out on rats (time period: 4 weeks) measured the differently surface-treated implants with pull-out tests. It was observed that implants coated with hydroxyapatite and modified with fluoride had better osseointegration than polished-surface implants. They established that in the several weeks after insertion the nanostructure assisted osseointegration [63].

In the research work of Vinzenz usable surface-treatment techniques for titanium implants are demonstrated. Furthermore he discusses several coating techniques, giving special attention to anodic plasma-chemical treatment of surfaces. Using scanning electron microscopy, X-ray spectroscopy, Raman spectroscopy, and laser-surface-roughness measurements he examined surfaces created by coating. With the results of cell and animal-experimentation examinations he showed that from the point of view of osseointegration the anodic plasma-chemical treated surface performed better than when untreated [64].

Research by Göransson introduces several surface-treatment and examination possibilities for dental implants made of titanium. In the experimental work samples with surfaces modified in 12 different ways were examined. He tested measures of osseointegration using cell experiments, animal experiment models, and human experiments. His result was that the surface structure did not significantly influence osseointegration and further that bioactive coatings also failed to produce significant results aiding improved osseointegration [65].

## 1.3    Methods of Examining the Surface

### 1.3.1    In Vitro Methods

It's usual to categorize methods of examination on the basis of which properties of surface morphology are tested by which method. This creates the primary division into chemical methods of analysis (X-ray Photoelectron Spectroscopy, Auger Electron Spectroscopy, Secondary Ion Mass Spectroscopy) and physical methods of analysis (Atomic Force Microscopy, Scanning Electron Microscopy).

### 1.3.2    In Vivo Methods

An often-used method for quantitative evaluation of the bone-implant contact's load-bearing capacity and the effects of the implant's geometrical characteristics (disregarding some unusual biological and chemical consequences of bone integration) is the pull-out test [66, 67]. After the appropriate healing time has elapsed (in the literature 6 weeks, or else 3 or 6 months) with the help of a special instrument the already bone-embedded implant is removed from the bone and the torque necessary to do this is recorded [68, 69]. Values shown by the instrument can be checked against an already-derived table of standard results to obtain a Ncm value. Removal torque is a very widely used method for animal studies.

In addition to studying implant morphology and with the above progress in mind, we decided to investigate how the nanostructure surface influences the proliferation of osteoblasts. Micro- and nano-morphological analyses revealed some correlations, including the potential biological, biochemical and physiological effects of surface roughness on bone cells adjacent to the implants. In these experiments titanium Grade 1 material was used.

# 2    Materials and Methods

For these experiments titanium Grade 1 (ISO5832 Pt.2 Grade I.) supplied thanks to the PROTETIM® company was made use of, machined discs of diameter 5 mm and thickness 2 mm. The surfaces of the discs were first cleaned with ultrasonic equipment in pure ethanol. The chemical modification of the surfaces of the discs chemical treatment involved passivation in sodium-hydroxide + $H_2O_2$, oxalic acid, nitric acid solution, electropolishing, sandblasting was performed with $Al_2O_3$ particles of 250 μm, along with laser modification (Pulsed Nd/Yag laser (Kvant 1) with a 1 J/pulse as the low-pulse energy laser, and the 3 J/pulse was used for surface treating as high-pulse energy laser). In these experiments, the surface-treated samples were compared with the machined discs as reference standards.

The machined surfaces of the discs were examined and images taken before their various treatments, first with a stereomicroscope (type: Olympus SZX 16), then with a scanning electron microscope (type: Philips XL 30). The discs' surface roughness values were compared quantitatively and qualitatively using a confocal microscope (type: Alicona Infinite Focus). Microscopic measures of the discs' nanostructures were taken (type: Veeco® diInnova™). Comparisons were carried out to measure contact angles and therefore hydrophilic or hydrophobic characteristics of the surfaces of the discs (equipment type: Rame-hart). The measurements were taken at 23°C temperature and 50% humidity. The device released 5-5 droplets (5 µl/droplet). The tests all used distilled water. Following the drop test (by 5 seconds) images were taken and the contact angles were measured using image-analysis software.

## 2.1    Experiments Done with Cells

Before the cell test the discs were placed for 20 minutes into tripsin solution, and for 20 minutes in alcohol, then sterilized in an autoclave. In twenty-four dishes, variously treated discs had NIH3T3 fibroblast and MC3T3 osteoblast cells added to their surfaces. For two days these cells were cultured on the disc surfaces in a cell-culturing medium (DMEM + 10% FBS). At the end of the incubation period discs with identical surface treatments are divided up into separate groups.

## 2.2    Cell-Counting

Cells dissolved from the disc surfaces in tripsin solution were counted using a Bürker camera. After cell-counting, cells were dissolved in a low-temperature 200 μl Triton (30 mM HEPES, 100 mM NaCl, 1 mM EGTA, 20 mM NaF, 1% Triton X-100, 1 mM PMSF, 20 μl/ml protease inhibitory cocktail, 1 mM $Na_3VO_4$) and then with the help of a spectrometer by the Bradford method, the concentration of dissolved protein (Bio-Rad Laboratories) was measured. Cell-counting was completed with a CyQuant DNA proliferation test. After culturing for two weeks CyQuant dye was added to the cells, coloring the DNA of each cell. After 5 minutes incubation the quantity of DNA was read off with a plate reader. Table 1 below lists the surface-treatment types and test methods used.

Table 1
Surface treatments and testing possibilities for Grade 1 titanium

| Type of surface treatment | Tests | Results |
|---|---|---|
| Machined | - stereomicroscope<br>- scanning electron microscope<br>- confocal microscope<br> - contact-angle measurement<br>- cell testing (MC3T3) | - surface morphology<br>- surface roughness<br>- wettability<br>- cell-proliferation |
| Machined + chemical etched | | |
| Machined + electro-polished | | |
| Machined + sandblasted ($Al_2O_3$) | | |
| Machined + 1 Joule impulse-energy laser surface modified | | |
| Machined + 3 Joule impulse-energy laser surface modified | | |

# 3    Results

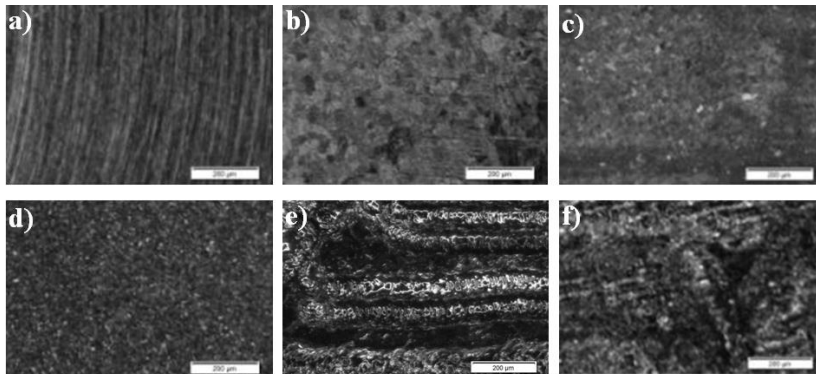In the illustrations below stereomicroscope images of samples are shown (Fig. 2).



Figure 2
Stereo-microscope images of titanium disks machined (a), chemical etched (b), electro-polished (c), $Al_2O_3$ sandblasted (d), 1 J laser surface modified (e), and 3 J laser surface modified

Machined: 3-10 µm width concentric grooves, running parallel to each other, irregular microgrooves. The depth and spacing of the furrows mostly regular. At numerous places on the surface contamination and machine-turning residue observable.

Machined + chemical etched: 1-3 µm width regularly arranged grooves. The primary leaf-form in the grooves run parallel to each other, but with grooves in the adjacent formations nearly perpendicular. This indicates that the groove patterning is not due to the original machining, but is caused by the subsequent surface modification. The leaf-shaped primary structure, about 20 to 25 µm in diameter, is the basic unit of the surface.

Machined + electro-polished: smooth surface, although some traces of grooving from the machined surface are still visible.

Machined + sandblasted ($Al_2O_3$): irregularly-shaped, sharp-contoured raised areas, showing similar geometric depths.

Machined + surface modified with 1 Joule impulse energy laser: 50-70 µm width grooves, regular waves on the surface. The surface is regular and smooth, insofar as microgrooving is smoothed down.

Machined + surface modified with 3 Joule impulse energy laser: 20-50 µm width furrows and 10-20 µm diameter droplets. An orderly spaced, wave like primary structure can be observed with a distance of 30-60 µm from each other. The wave crest projections show 30 to 50 µm. The so-called secondary structure is pear shaped having a size of 10-15 µm.

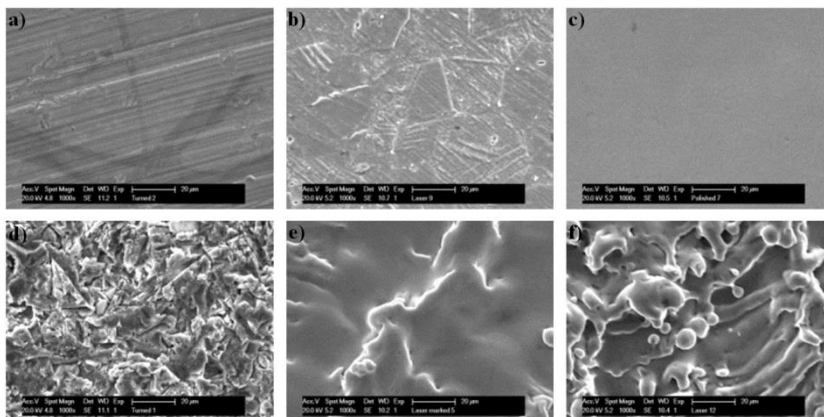On the SEM images 3-10 µm width grooves can be seen on the surface of the machined sample (Figure 3).



Figure 3

Scanning electron microscope images (magnification 1000×) machined (a), chemical etched (b), electropolished (c), $Al_2O_3$ sandblasted (d), 1 J laser surface modified (e), and 3 J laser surface modified

After chemical etching there were ordered structure on the surface of the disk, parallel – though oriented in the original direction of the surface machining – grooves separated by intervals of about 1-3 µm. After electropolishing, a smooth uniform surface – absent the concentric grooving seen on the freshly machined samples – was apparent. The surface bombarded with $Al_2O_3$ sandblasting shows irregular roughly granular zones. The effect of 1 Joule impulse energy laser surface modification is 50-70 µm wide, regular, wave-like grooving. Whereas 3 Joule impulse energy laser surface modification has the effect of creating 20-50 µm width grooves on the surface, in which solidified droplets of diameter 10-20 µm are visible.

With a confocal microscope of 3-3 measurements the average surface roughness (Ra) was established. Measurements were taken at 1500 µm intervals. The different surface morphologies were compared in this case with the machined sample as a reference. The results of these measurements are brought together in Figure 4 below.
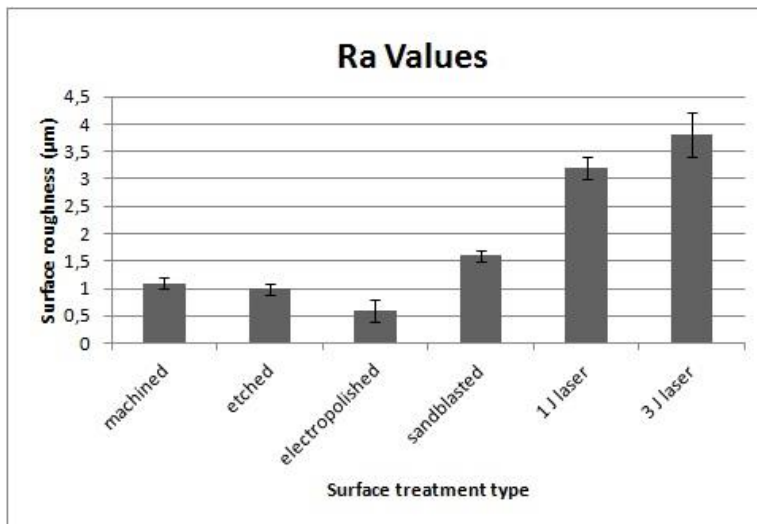


Figure 4

Measured values for the surface roughness (Ra) of the discs

The surface roughness of the sample treated with electropolishing was the lowest (Ra = 0.6 ± 0.2 µm), while the highest value (Ra = 3.8 ± 0.4 µm) was for the discs treated with 3 Joule impulse laser modification. The chemically etched discs' surface roughness (Ra = 1.0 ± 0.1 µm) was slightly less than that of the freshly machined sample (Ra = 1.1 ± 0.1 µm). The titanium discs sandblasted with $Al_2O_3$ particles had a roughness of Ra = 1.6 ± 0.1 µm. The 1 Joule impulse energy laser modified surface had a roughness level (Ra = 3.2 ± 0.2 µm) only slightly less than that created by the 3 Joule laser.

In the case of the measurements by the atomic-force microscope the measured surface was 400 $\mu m^2$. The results of the measurements are collected together in Figure 5 below.
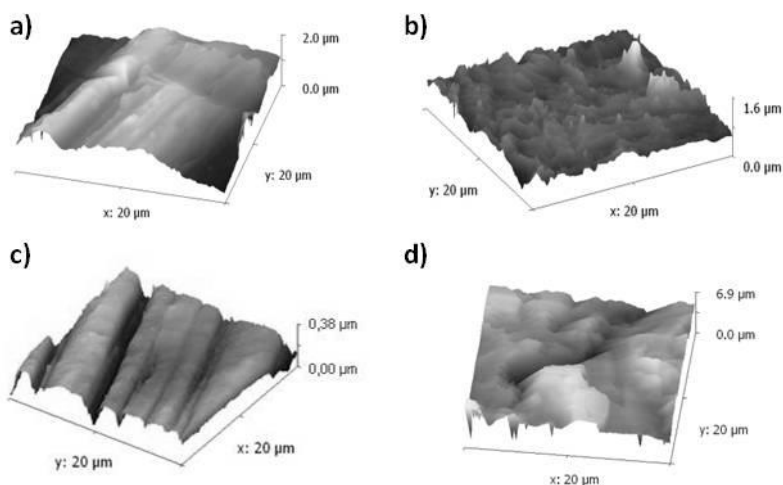


Figure 5

Titanium Grade 1 discs measurement results: machined (a), etched (b), electropolished (c), $Al_2O_3$ sandblasted (d)

The Ra roughness value for various treatments to the surface of the disc was as follows: machined 0.62 µm; chemically etched 0.08µm; electropolished 0.04 µm; $Al_2O_3$ sandblasting 0.27 µm. Samples treated with laser surface-modification did not get subjected to AFM measurements. This is accounted for by physical limitations on the measurement method (extremely large Ra surface roughness, which goes beyond the measurement capacity of the instruments – NB: both the AFM and confocal microscope use methods which establish surface roughness values, but there are variations in the results which deviate from the surfaces examined). On the basis of the AFM measurements it was established that the value of surface roughness Ra was lowest for the electropolished discs and highest for the machined samples.

Contact angles were established, with the $Al_2O_3$ sandblasting $(60 \pm 10°)$ and 1 Joule impulse energy laser surface modification $(51 \pm 3°)$ results close together within margins of error. The electropolished sample surface had the lowest value $(27 \pm 4°)$, while the 3 Joule impulse energy laser treatment gave the highest contact-angle values $(102 \pm 4°)$.

The following quantitative and qualitative characteristics were subjected to cell testing.

## 3.1 Cell-Counting and Protein Concentration Measurement

Results of cell-counting and protein concentration in both the fibroblast and the osteoblast cases showed increased cell proliferation on the roughened surfaces when compared with the machined surface. Tables 2 and 3 present the measured results as averages ± standard deviation. Numbers in the columns are relative values which show what the cell count or protein concentration was on the given surface to the value for the machined surface – used widely in the literature as the reference surface.

Table 2

Changes of cell numbers and protein concentration of NIH3T3 fibroblast cells measured on various modified surfaces compared to the group with a machined surface

| Sample | Cell number | Prot. Conc. |
|---|---|---|
| Machined | 1 | 1 |
| Machined + chemical etched | 1.85 ± 0.15 | 1.34 ± 0.28 |
| Machined + electropolished | 3.24 ± 0.12 | 2.53 ± 0.35 |
| Machined + sandblasted ($Al_2O_3$) | 5.19 ± 0.59 | 3.18 ± 0.41 |
| Machined + 1 Joule impulse energy laser surface modified | 3.17 ± 0.33 | 2.66 ± 0.26 |
| Machined + 3 Joule impulse energy laser surface modified | 2.83 ± 0.36 | 2.31 ± 0.15 |

Table 3

Changes of cell numbers and protein concentration of MC3T3 osteoblast cells measured on various modified surfaces compared to the group with a machined surface

| Sample | Cell number | Prot. Conc. |
|---|---|---|
| Machined | 1 | 1 |
| Machined + chemical etching | 1.55 ± 0.16 | 1.40 ± 0.21 |
| Machined + electropolished | 2.15 ± 0.32 | 2.32 ± 0.41 |
| Machined + sandblasted ($Al_2O_3$) | 3.25 ± 0.33 | 3.05 ± 0.45 |
| Machined + 1 Joule impulse energy laser surface modification | 3.25 ± 0.22 | 2.40 ± 0.37 |
| Machined + 3 Joule impulse energy laser surface modification | 3.83 ± 0.29 | 3.41 ± 0.41 |

**Conclusions**

In the experiments 5 mm diameter, 2 mm thickness discs machined from Grade 1 titanium were used. The chemical treatment to the surface of the discs involved passivation in sodium-hydroxide + $H_2O_2$, oxalic acid, nitric acid solution, electropolishing, $Al_2O_3$ sandblasting, as well as laser pulse surface modification (Pulsed Nd/Yag laser (Kvant 1) with a 1 J/pulse is the low pulse energy laser, and

the 3 J/pulse was used for surface treating as high pulse energy laser). Images were taken of the discs using a stereomicroscope and a scanning electron microscope. On the scanning electron microscope images 3-10 µm width grooving can be seen on the freshly-machined sample surfaces.

After chemical etching, structures appear on the surface of the discs parallel to each other – but these follow the orientation of the structures on the machined surface, 1-3 µm grooves.

After electropolishing we obtain a smooth surface, without the concentric grooves visible on the machined samples.

The $Al_2O_3$ sandblasted surfaces showed irregular, roughly granular zones. The effect of the 1 Joule impulse energy laser surface modification was 50-70 µm width, regular, wave-shaped surface grooving. Whereas the effect of 3 Joule impulse energy laser surface modification was to create 20-50 µm width grooves on the surface, in which solidified droplets of diameter 10-20 µm were visible.

During examination with a confocal microscope, it was established that the surface roughness value of the electropolished samples (Ra = $0.6 \pm 0.2$ µm) was the lowest, while the highest roughness values (Ra = $3.8 \pm 0.4$ µm) were produced on 3 Joule impulse energy laser modified discs. The roughness value for those discs treated with chemically etching (Ra = $1.0 \pm 0.1$ µm) was slightly less than that for the freshly machined samples. The surface roughness of the titanium discs treated with $Al_2O_3$ sandblasting was Ra = $1.6 \pm 0.1$ µm. Roughness values for discs surface-modified by the 1 Joule pulsed energy laser (Ra = $3.2 \pm 0.2$ µm) deviated only slightly from the 3 Joule laser value.

AFM measurements of the samples were carried out. The machined discs had surface roughness values of 0.62 µm, the chemically etched 0.08 µm, the electropolished 0.04 µm, the $Al_2O_3$ sandblasted samples 0.27 µm. It was not possible to take accurate AFM measurements of the laser surface modified surfaces.

Contact angles were established for wetting of the $Al_2O_3$ sandblasted surfaces ($60 \pm 10°$), the 1 Joule pulse energy laser modified surfaces ($51 \pm 3°$), values that were very close together. On the electropolished surfaces these were measured with the lowest values ($27 \pm 4°$), while the highest contact-angle values ($102 \pm 4°$) appeared on the 3 Joule impulse-energy laser modified surfaces. Quantitative and qualitative cell-test characterizations of the surfaces were made. In the case of tests with fibroblasts the greatest degree of cell proliferation was the surfaces sandblasted with $Al_2O_3$ particles. With the osteoblast tests the highest levels of cell activity were for the two types of laser-modified surface.

## References

[1]     Zhe, Q., Xiaohui, R. F., Marco, W., Michael, M., Andreas, S.: The Initial Attachment and Subsequent Behavior Regulation of Osteoblasts by Dental

Implant Surface Modification. Journal of Biomedical Materials Research Part A, 82 (2007:3) 658-668

[2]    Albrektsson, T., Brånemark, P. I., Hansson, H. A., Lindström, J.: Osseointegrated Titanium Implants. Acta Orthop Scand, 52 (1981) 155-170

[3]    Kieswetter, K., Schwartz, Z., Dean, D. D., Boyan, B. D.: The Role of Implant Surface Characteristics in the Healing of Bone. Crit Rev Oral Biol Med, 7 (1996:4) 329-345

[4]    Stadlinger, B., Korn, P., Tödtmann, N., Eckelt, U., Range, U., Bürki, A., Ferguson, S. J., Kramer, I., Kautz, A., Schnabelrauch, M., Kneissel, M., Schlottig, F.: Osseointegration of Biochemically Modified Implants in an Osteoporosis Rodent Model. Eur Cell Mater. 25 (2013:8) 326-40

[5]    Cochran, D. L., Buser, D., ten Bruggenkate, C. M., Weingart, D., Taylor, T. M., Bernard, J. P., Peters, F., Simpson, J. P.: The Use of Reduced Healing Times on ITI Implants with a Sandblasted and Acid-etched (SLA) Surface: Early Results from Clinical Trialson ITI SLA Implants. Clinical Oral Implants Research, 13 (2002) 144-153

[6]    Ramazanoglu, M., Oshida, Y.: Osseointegration and Bioscience of Implant Surfaces - Current Concepts at Bone-Implant Interface, Implant Dentistry - A Rapidly Evolving Practice, Prof. Ilser Turkyilmaz (Ed.), Open Access Publisher: InTech, 2011

[7]    Guéhennec, L. L., Soueidan, A., Layrolle, P., Amouriq, Y.: Surface Treatments of Titanium Dental Implants for Rapid Osseointegration. Dental Materials, 23 (2007) 844-854

[8]    Puleo, D. A., Nanci, A.: Understanding and Controlling the Bone-Implant Interface. Biomaterials, 20 (1999) 2311-2321

[9]    Hulbert, S. F., Morisson, S. F., Klawitter, J. J.: Tissue Reaction to Three Ceramics of Porous and Non-Porous Structures. J Biomed Mater Res 6 (1972) 347-374

[10]   Carlsson, L., Rostlund, T., Albrektsson, T.: Removal Torque for Polished and Rough Titanium Implants. Int J Oral Maxillofacial Implants, 3 (1988) 21-22

[11]   Cochran, D. L.: A Comparison of Endosseus Dental Implant Surface. J Periodontal, 70 (1999) 1523- 1539

[12]   De Assis, A. F., Beloti, M. M., Crippa, E. G., Oliveira, P. T., Morra, M., Roza, A. L.: Development of the Osteoblastic Phenothype in Human Alveolar Bone-derived Cells Grown on a Collagen Type I-coated Titanium Surface. Clin Oral Impl Res, 20 (2009) 240-246

[13]   Wennerberg, A., Albrektsson, T., Andersson, B.: Design and Surface Characteristics of 13 Commercially Available Oral Implant Systems. Int J Oral Maxillofacial Implants, 8 (1993) 622-633

[14]    Wennerberg, A., Ektassobi, A., Albrektsson, T., Johansson, L., Andersson, B.: A 1-Year Follow up of Implants of Differing Surface Roughness Placed in Rabbit Bone. Int J Oral Maxillofac Implants, 12 (1997) 486-490

[15]    Suba, Cs., Lakatos-Varsányi, M., Miko, A., Kovács, L., Velich, N., Kádár, B., Szabó, Gy.: Study of the Electrochemical Behaviour of Ti Osteosynthesis Plates Used in Maxillofacial Surgery. Mater Sci Eng, 447 (2007) 347-354

[16]    Junker, R., Dimakis, A., Thoneick, M., Jansen, J. A.: Effect of Implant Surface Coatings and Composition on Bone Integration: a Systematic Review. Clin Oral Impl Res, 20 (2009) 185-206

[17]    Kasemo, B., Gold, J.: Implant Surface and Interface Processes. Adv Dent Res, 13 (1999) 8-20

[18]    Schwarz, F., Herten, M., Sager, M., Wieland, M., Dard, M., Becker, J.: Bone Regeneration in Dehiscence-Type Defects at Chemically Modified (SLActive) and Conventional SLA Titanium Implants: a Pilot Study in Dogs. J Clin Periodontol, 34 (2007) 78-86

[19]    Velich, N., Kádár, B., Kiss, G., Kovács, K., Réti, F., Szigeti, K., Garagiola, U., Szabó, Gy.: Effect of Human Organism on the Oxide Layer Formed on Titanium Osteosynthesis Plates: A Surface Analytical Study. The Journal of Craniofacial Surgery, 6 (2006) 13-17

[20]    Buser, D., Broggini, N., Wieland, M., Schenk, R. K., Denzer, A. J., Cochran, D. L., Hoffmann, B., Lussi, A., Steinemann, S. G.: Enhanced Bone Apposition to a Chemically Modified SLA Titanium Surface. J Dent Res, 7 (2004) 529-533

[21]    Joób, F. Á., Huszár, T., Divinyi, T., Rosivall, L., Szabó, Gy.: The Effect of the Surface Mikromorphology of Titanium Implants on the Fibro- and Osteoblast Prolifeartion Activity. Fogorvosi Szemle, 97 (2004) 251-255

[22]    Shibli, J. A., Grassi, S., de Figueiredo, L. C., Feres, M., Marcantonio, E. Jr.. Iezzi, G., Piattelli, A.: Influence of Implant Surface Topography on Early Osseointegration: a Histological Study in Human Jaws. Journal of biomedical materials research. Part B, Applied biomaterials, 80 (2007) 377-385

[23]    Soskolne, W. A., Cohen, S., Sennerby, L., Wennerberg, A, Shapira, L.: The Effect Oftitanium Surface Roughness on the Adhesion of Monocytes and their Secretion ofTNF-Alpha and PGE2. Clinical Oral Implants Research, 13 (2002) 86-93

[24]    Klokkevold, P. R., Nishimura, R. D.: Osseointegration Enhanced by Chemical Etching of Titanium Surface. Clin Oral Impl Res, 8 (1997) 442-447

[25] Schwarz, F., Herten, M., Sager, M., Wieland, M., Dard, M., Becker, J.: Bone Regeneration in Dehiscence-Type Defects at Chemically Modified (SLActive) and Conventional SLA Titanium Implants: a Pilot Study in Dogs. J Clin Periodontol, 34 (2007) 78-86

[26] Orsini, G., Assenza, B., Scarano, A., Piatteli, M., Piatteli, A.: Surface Analysis of Machined versus Sandblasted and Acid-etched Titanium Implants. Int J Oral Maxillofac Implants, 15 (2000) 779-784

[27] Joób, F. A., Huszár, T., Divinyi, T., Rosivall, L.: The Effect of The Surface Micromorphology of Titanium Dental Implants on Proliferation Activity of Fibro-Osteoblasts. Paper presented at the 9. Kongress Der Österreichischen Gesellschaft Für Mund-, Kiefer-Und Gesichtschirurgie, Bad Hofgastein, Austria, 2005

[28] Yoshiki O.: Bioscience and Bioengineering of Titanium Materials. Second edition. Oxford: Elsevier, 2013

[29] Gaggl, A., Schultes, G., Müller, W. D., Karchen, M.: Scanning Electron Microscopical Analysis of Laser-treated Titanium Implant Surfaces – a comparative study. Biomaterials, 21 (2000) 1067-1073

[30] Nentwig, G. H., Reichel, M.: Vergleichende Untersuchung zur Mikromorphologie und Gesamtoberfläche enossaler Implantate Z. Zahnärztl Implantologie, 10 (1994) 150-154

[31] Wennerberg, A., Albrektsson, T., Ulrich, H., Krol, J. J.: An Optical Three-Dimensional Technique for Topographical Descriptions of Surgical Implants. Journal of Biomedical Engineering, 14 (1992) 412-418

[32] Wennerberg, A., Albrektsson, T., Andersson, B.: Design and Surface Characteristics of 13 Commercially Available Oral Implant Systems. Int J Oral Maxillofacial Implants, 8 (1993) 622-633

[33] Wennerberg, A., Albrektsson, T.: Bone Tissue Response to Commercially Pure Titanium Implants Blasted with Fine and Coarse Particles of Aluminium Oxide. Int J Oral Maxillofacial Implants, 11 (1996) 38-45

[34] Albrektsson, T., Wennerberg, A.: Die klinische Bedeutung verschidener Oberflacheneigenschafften enossaler. Titanimplantate Implantologie, 3 (1999) 235-246

[35] Wennerberg, A., Albrektsson, T.: Suggested Guidlines for the Topographic Evaluation of Implant Surfaces. Int J Oral Maxillofacial Implants, 15 (2000) 331-334

[36] Lazzara, R. J., Testori, T., Trisi, P., Porter, S. S., Weinstein, R. L.: A Human Histologic Analysis of Osseotite and Machined Surface using Implants with 2 Opposing Surfaces. Int J of Periodontics Restorative Dent, 19 (1999) 117-129

[37]   Ivanov, C. J., Hallgren, C., Widmark, G., Sennerby, L., Wennerberg, A.: Histologic Evaluation of the Bone Integration of TiO2 Blasted and Turned Titanium Microimplants in Humans. Clin Oral Impl Res, 12 (2001) 44-50

[38]   Buser, D.: Titanimplantate mit angerauhter Oberfläche. Implantologie, 3 (1999) 249-268

[39]   Piatteli, A., Corigliano, M., Scarano, A., Quaranta, M.: Bone Reactions to Early Occlusal Loading of Two Stage Titanium Plasma-sprayed Implants: a Pilot Study in Monkeys. Int J Periodontics Restorative Dent, 17 (1997) 162-169

[40]   Hansson, S., Norton, M.: The Relation between Surface Roughness and Interfacial Shear Strength for Bone-anchored Implants. A Mathematical Model. J Biomech, 32 (1999) 829-36

[41]   Miyamoto, I., Tsuboi, Y., Wada, E., Suwa, H., Iizuka, T.: Influence of Cortical Bone Thickness and Implant Length on Implant Stability at the Time of Surgery–Clinical. Prospective, Biomechanical, and Imaging Study, 37 (2005) 776-780

[42]   Xianshuai, C., Longhan, X., Jianyu, C. R., Feilong, D.: Design and Fabrication of Custom-made Dental Implants. Journal of Mechanical Science and Technology, 26 (2012:7) 1993-1998

[43]   Juodzbalys, G., Sapragoniene, M., Wennerberg, A., Baltrukonis, T.: Titanium Dental Implant Surface Micromorphology Optimization. Journal of Oral Implantology, 34 (2007:4) 177-85

[44]   Park, J. Y., Davies, J. E.: Red Blood Cell and Platelet Interactions with Titanium Implant Surfaces. Clinical Oral Implants Research, 12 (2000) 530-539

[45]   Pammer, D., Schindler, Á., Bognár, E.: Chemical Etching of Dental Implants. Gépészeti szám, 60 (2013) 29-32

[46]   Rosa, A. L., Beloti, M. M.: Rat Bone Marrow Cell Response to Titanium and Titanium Alloy with Different Surface Roughness. Clinical Oral Implants Research, 14 (2003:1) 43-48

[47]   Galli, C., Guizzardi, S., Passeri, G., Martini, D., Tinti, A., Mauro, G., Macaluso, G. M.: Comparison of Human Mandibular Osteoblasts Grown on Two Commercially Available Titanium Implant Surfaces. Journal of Periodontology, 76 (2005) 364-372

[48]   Engquist, B., Astrand, P., Dahlgren, S., Engquist, E., Feldmann, H., Grondahl, K.: Marginal Bone Reaction to Oral Implants: a Prospective Comparative Study of AstraTech and Branemark System Implants. Clinical Oral Implants Research, 13 (2002) 30-37

[49]  Gaggl, A., Schultes, G., Muller, W. D., Karcher, H.: Scanning Electron Microscopic Analysis of Laser-treated Titanium Implant Surfaces -- a Comparative Study. Biomaterials, 21 (2000) 1067-1073

[50]  Yamagami, A., Yoshihara, Y., Suwa, F.: Mechanical and Histologic Examination of Titanium Alloy Material Treated by Sandblasting and Anodic Oxidization. The International Journal of Oral Maxillofacial Implants, 20 (2005:1) 48-53

[51]  Anil, S., Anand, P. S., Alghamdi, H., Jansen, J. A.: Dental Implant Surface Enhancement and Osseointegration. Implant Dentistry - A Rapidly Evolving Practice, Prof. Ilser Turkyilmaz (Ed.), Open Access Publisher: InTech, 2011

[52]  Lill, W., Velikogne, W., Danhel-Mayrhauser, M., Haider, R., Plenk, H., Watzek, G.: Histomorphometrische untersuchung der Knochenreaktion um extraoralen Brånemark® - und TPS® Titanschrauben beim Schaf Zeitschrift für Zahnärzliche. Implantologie, 8 (1992) 103-112

[53]  Uitto, V. J., Larjava, H., Peltonen, J., Brunette, D. M.: Expressions of Fibronectin and Integrins in Cultured Periodontal Ligament Epithelial Cells. J. Dent. Res., 71 (1997) 1203-1211

[54]  Ericsson, I., Johansson, C. B., Bystedt, H., Norton, M. R.: A Histomorphometric Evaluation of Bone-to-Implant Contact on Machine-prepared and roughened Titanium Dental Implants. Clin Oral Impl Res, 5 (1994) 202-206

[55]  Wennerberg, A., Ektassobi, A., Albrektsson, T., Johansson, L., Andersson, B.: A 1- Year Follow up of Implants of Differing Surface Roughness Placed in Rabbit Bone. Int J Oral Maxillofac. Implants, 12 (1997) 486-490

[56]  Cordioli, G., Piatteli, A.: Removal Torque and Histomorphometric Investigation of 4 Different Titanium Surface: an Experimental Study in the Rabbit Tibia. Int J Oral and Maxillofacial Implants, 15 (2000) 668-674

[57]  Baker, D. A., London, R. M., Oneal, R. B.: Integration Strength and Speed of Dual-etched Titanium Implants: A Comparative Study in Rabbits. 13[th] Annual Meeting, Academy of Osseointegratio, Atlanta, Georgia, 1998

[58]  Trisi, P., Lazzara, R., Rao, W., Reboudi, A.: Bone-Implant Contact and Bone Quality: of Expected and Actual Bone Contact on Machined and Osseotite Implant Surface. Int J Periodontics Rest Dent, 22 (2002) 535-545

[59]  Wennerberg, A., Hallgren, C., Johansson, C., Danelli, S.: A Histomophometric Evaluation of Screw-shaped Implants Each Prepared with Two Surface Rougness. Clin Oral Impl Res, 9 (1998) 11-19

[60]  Buser, D., Schenk, R. K., Steinemann, S., Fiorellini, J. P., Fox, C. H., Stich, H.: Influence of Surface Characteristics on Bone Integration of Titanium

Implants. A Histomorphometric Study in Miniature Pigs. J Biomed Mater Res, 25 (1991) 889-902

[61]    Christenson, E. M., Anseth, K. S., van den Beucken, J. J., Chan, C. K., Ercan, B., Jansen, J. A., Laurencin, C. T., Li, W. J., Murugan, R., Nair, L. S., Ramakrishna, S., Tuan, R. S., Webster, T. J., Mikos, A. G.: Nanobiomaterial Applications in Orthopedics. J Orthop Res, 25 (2007) 11-22

[62]    Carlos, N. E., Yoshiki, O., José, H. C. L., Carlos, A. M.: Relationship between Surface Properties (Roughness, Wettability and Morphology) of Titanium and Dental Implant Removal Torque. Journal of the Mechanical Behavior of Biomedical Materials, 1 (2008:3) 234-242

[63]    Luiz, M.: On Nano Size Structures for Enhanced Early Bone Formation. PhD-dissertation, Inst of Odontology. Dept of Prosthetic Dentistry, Dental Material Science, Göteborg University, Göteborg, Sweden, 2007

[64]    Frauchiger, V. M.: Anodic Plasma-Chemical Treatment of Titanium Implant Surfaces. PhD-dissertation, Technische Wissenschaften ETH Zürich, Zürich, Switzerland, 2002

[65]    Göransson, A.: On Possibly Bioactive CP Titanium Implant Surfaces. PhD-dissertation, Department of Prosthetic Dentistry, University of Gothenburg, Göteborg, Sweden, 2009

[66]    Ogiso, M., Yamamura, M., Kuo, P. T., Borgese, D., Matsumoto, T.: Comperative Push-Out Test of Dense HA Implants and HA-coated Implants: Findings in a Canine Study. J. Biomed. Mater. Res, 39 (1998) 364-372

[67]    Svehla, M., Morberg, P., Zicat, B., Bruce, W., Sonnabend, D., Walsh, W. R.: Morphometric and Mechanical Evaluation of Titanium Implant Integration: Comparison of Five Surface Structures. J. Biomed. Mater. Res., 51 (2000) 15-22

[68]    Sennersby, L., Thomsen, P., Ericson, L. E.: A Morphometric and Biomechanic Comparison of Titanium Implants Inserted in Rabbit Cortical and Cancellous Bone. Int J Oral Maxillofac Implants, 7 (1992) 62-71

[69]    Cordioli, G., Piatteli, A.: Removal Torque and Histomorphometric Investigation of 4 Different Titanium Surface: an Experimental Study in the Rabbit Tibia. Int J Oral and Maxillofacial Implants, 15 (2000) 668-674

# Performance Modeling of Web-based Software Systems with Subspace Identification

**Ágnes Bogárdi-Mészöly, András Rövid, Shohei Yokoyama**

Department of Automation and Applied Informatics, Budapest University of Technology and Economics, Magyar tudósok krt. 2, 1117 Budapest, Hungary
Institute of Applied Mathematics, Óbuda University, Bécsi út 96/B, 1034 Budapest, Hungary
Department of Computer Science, Shizuoka University, Hamamatsu Campus, Shizuoka ken, Hamamatsu shi, Naka ku, Jouhoku 3-5-1, 432 8011, Japan
agi@aut.bme.hu, rovid.andras@nik.uni-obuda.hu, yokoyama@inf.shizuoka.ac.jp

*Abstract: Performance modeling and prediction of web-based software systems are important and complicated considerations. The goal of our paper is to establish proper mathematical models in the form of difference equations, by subspace identification, in order to model and predict the performance of web-based software systems. First, simulation models have been provided to simulate the behavior of thread pool and queued requests. Second, analytical models have been proposed in form of state space models using subspace identification. In addition, it has been demonstrated that the proposed models can be applied to performance prediction of web-based software systems. The proposed models have been validated and verified. Furthermore, performance factor identification and performance prediction techniques have been proposed based on subspace identification.*

*Keywords: web-based software system; subspace identification; performance modeling; performance factor identification; performance prediction*

# 1   Introduction

Web-based software systems are client-server software applications, in which, clients run in web browsers. Because they serve a large number of users, their performance and efficiency have become of key importance. Their architecture and runtime environment mainly diverge from the previous concepts. Their performance modeling and prediction are active research fields.

Performance evaluation is significant at every stage of the development process. There are three main techniques for performance evaluation: analytical modeling, simulation, and measurement [1].

Web-based software systems access some resources while executing the requests of the clients. Typically, several requests arrive at the same time, thus, a competitive situation is established for the resources. For modeling such situations, queueing model-based approaches are widely used [1] [2] [3].

In our previous work [4] [5] [6], statistical methods have been provided in order to identify and investigate novel performance factors as well as queueing network models and evaluation algorithms have been proposed to model the identified dominant thread pool and queue limit performance factors.

In this work [7] [8], a different approach is investigated for modeling the behavior of thread pool and queued requests, state space models are provided using subspace identification, in order to model and predict the performance.

The paper is organized as follows: Section 2 covers the background and related work. Section 3 presents the contribution, namely, simulation models, state space models, performance prediction, error analysis, and propositions. Finally, the last section reports the conclusions.

## 2   Background and Related Work

This section is devoted to review the background and research efforts related to this work, namely, the concept of thread pool and queued requests, the used SimEvents Toolbox of MATLAB Simulink, and the applied subspace identification method.

### 2.1   Concept of Thread Pool and Queued Requests

In the case of using a thread pool depicted in Fig. 1, when a request arrives, the application adds it to an incoming queue [9]. A group of threads retrieves requests from this queue and processes them. As each thread is freed, another request is executed from the queue.

The architecture of ASP.NET environment [5] [10] can be seen in Fig. 2. If a client is requesting a service from the server, the request goes through several subsystems before it is served. From the Internet Information Services (IIS), the requests are placed into the named pipe, which is a global queue between IIS and ASP.NET, its limit is set by the *requestQueueLimit* property. From the named pipe, the requests are placed into an application queue, which is used to maintain the availability of worker and I/O threads, its limit is configured by the *appRequestQueueLimit* property. When the limit is exceeded, the requests are rejected.
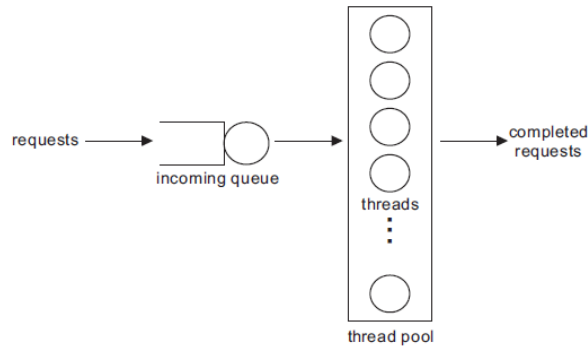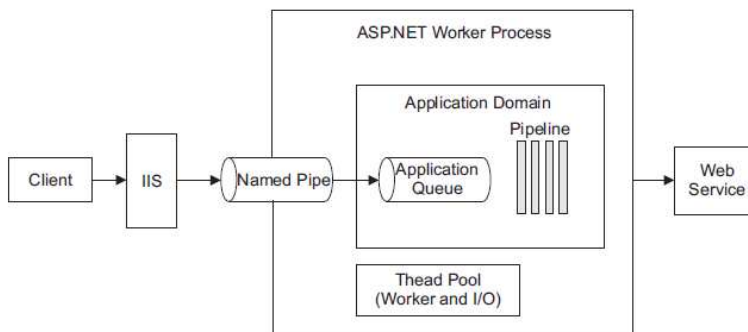
Figure 1
Thread pool and queued requests



Figure 2
Architecture of ASP.NET environment

## 2.2    MATLAB Simulink SimEvents

MATLAB Simulink is used for modeling, simulation, and analysis of dynamic systems. Simulink is designed for time-based simulation, while SimEvents [11] is aimed for event-based simulation. SimEvents extends Simulink behavior with a discrete-event simulation model of computation.

In SimEvents models, during a simulation, entities can pass through a network of queues, servers, switches, etc. Entities can carry data as attributes. Events can change state variables, outputs, occurrences of other events, for example, an entity advances from one block to another, and the service of an entity is completed in a server. Entities can wait in a queue; its capacity limits the simultaneous number of entities in the queue. Servers can serve entities in service time of given distribution; simultaneous number of entities can be finite or infinite. Entities can be routed from/to the selected input/output port using input/output switches. Multiple paths can be merged into a single path using path combiners.

## 2.3   Subspace Identification

*Definition 1.* The state space representation of a deterministic, discrete time, linear, time invariant system is defined by the following difference equations:

$$\mathbf{x}_{k+1}=\mathbf{A}\mathbf{x}_k+\mathbf{B}\mathbf{u}_k, \tag{1}$$

$$\mathbf{y}_k=\mathbf{C}\mathbf{x}_k+\mathbf{D}\mathbf{u}_k, \tag{2}$$

where $\mathbf{x}_k$ represents the state vector, $\mathbf{u}_k$ is the input vector, and $\mathbf{y}_k$ is the output vector at time $k$ as well as $\mathbf{A}$ is the state matrix, $\mathbf{B}$ is the input matrix, $\mathbf{C}$ is the output matrix, and $\mathbf{D}$ is the feedthrough matrix.

The aim is to determine system matrices $\mathbf{A}$, $\mathbf{B}$, $\mathbf{C}$, $\mathbf{D}$ from input-output data by subspace identification. The main thoughts of subspace identification algorithm are demonstrated as follows [12]. Input and output block Hankel matrices ($\mathbf{U}_{1|i}, \mathbf{Y}_{1|i}$) are constructed reflecting the history of input-output data. State sequence ($\mathbf{X}_i$) plays an important role in derivation and interpretation.

$$\mathbf{U}_{1|i}=\begin{bmatrix} \mathbf{u}_1 & \mathbf{u}_2 & \cdots & \mathbf{u}_j \\ \mathbf{u}_2 & \mathbf{u}_3 & \cdots & \mathbf{u}_{j+1} \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{u}_i & \mathbf{u}_{i+1} & \cdots & \mathbf{u}_{i+j-1} \end{bmatrix}, \; \mathbf{Y}_{1|i}=\begin{bmatrix} \mathbf{y}_1 & \mathbf{y}_2 & \cdots & \mathbf{y}_j \\ \mathbf{y}_2 & \mathbf{y}_3 & \cdots & \mathbf{y}_{j+1} \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{y}_i & \mathbf{y}_{i+1} & \cdots & \mathbf{y}_{i+j-1} \end{bmatrix} \tag{3}$$

$$\mathbf{X}_i=\begin{bmatrix} \mathbf{x}_i & \mathbf{x}_{i+1} & \cdots & \mathbf{x}_{i+j-1} \end{bmatrix} \tag{4}$$

State and output equations can be written using extended version of controllability ($\mathbf{\Delta}_i$) and observability ($\mathbf{\Gamma}_i$) matrices as well as lower block triangular Toeplitz matrix ($\mathbf{H}_i$). In geometrical interpretation, output is in the vector space determined by the union of state and input row spaces, state sequence can be estimated by projection of output row space onto orthogonal complement of input row space. Rank can be determined using singular value decomposition.

$$\mathbf{\Gamma}_i=\begin{bmatrix} \mathbf{C} \\ \mathbf{CA} \\ \vdots \\ \mathbf{CA}^{i-1} \end{bmatrix}, \; \mathbf{\Delta}_i=\begin{bmatrix} \mathbf{A}^{i-1}\mathbf{B} & \cdots & \mathbf{AB} & \mathbf{B} \end{bmatrix} \tag{5}$$

$$\mathbf{H}_i=\begin{bmatrix} \mathbf{D} & 0 & 0 & \cdots & 0 \\ \mathbf{CB} & \mathbf{D} & 0 & \cdots & 0 \\ \mathbf{CAB} & \mathbf{CB} & \mathbf{D} & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathbf{CA}^{i-2}\mathbf{B} & \mathbf{CA}^{i-3}\mathbf{B} & \mathbf{CA}^{i-4}\mathbf{B} & \cdots & \mathbf{D} \end{bmatrix} \tag{6}$$

$$\mathbf{X}_{i+1} = \mathbf{A}^i \mathbf{X}_1 + \mathbf{\Delta}_i \mathbf{U}_{1|i}, \ \ \mathbf{Y}_{1|i} = \mathbf{\Gamma}_i \mathbf{X}_1 + \mathbf{H}_i \mathbf{U}_{1|i} \tag{7}$$

System matrices can be estimated in least squares sense from the following equations:

$$\begin{bmatrix} \mathbf{X}_{i+1} \\ \mathbf{Y}_i \end{bmatrix} = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{X}_i \\ \mathbf{U}_i \end{bmatrix} \tag{8}$$

Subspace identification has been successfully applied in various application fields, also together with queueing models [13].

# 3   Performance Modeling

The main contribution of this paper is to establish state space models using subspace identification method in order to model and predict the performance of web-based software systems. Table 1 summarizes the notations used in provided models.

Table 1
Notations of models

| Notation | Input/output | Model | Meaning |
|---|---|---|---|
| $\mathbf{u}_1$ | Input | simulation | number of all users |
| $\mathbf{u}_2$ | Input | simulation | service time |
| $\mathbf{u}_3$ | Input | simulation | number of all dropped requests |
| $\mathbf{u}_4$ | Input | simulation | waiting time in queue |
| $\mathbf{y}$ | Output | simulation | average response time |
| $\mathbf{y}_{measured}$ | Output | simulation | average response time ($\mathbf{y}_{measured} = \mathbf{y}$) |
| $\mathbf{y}_{model}$ | Output | state space | average response time |

## 3.1   Simulation Models

Firstly, for modeling thread pool and queue limit, simulation models (Fig. 3) have been provided using SimEvents of MATLAB Simulink in order to simulate the thread pool and queued requests behavior as well as to obtain input-output data for subspace identification method. Results of other simulation models or measurements can also be used as input-output data.
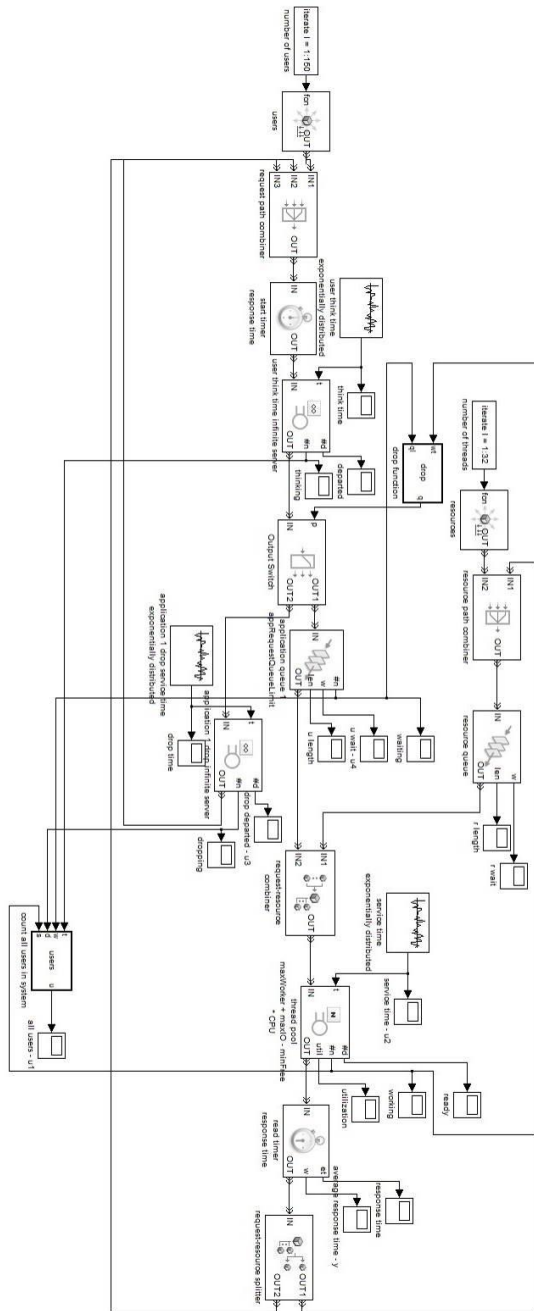
Figure 3
Simulation model

SimEvents models have been designed according to the concept of thread pool and queued requests (see in Section 2.1) as well as the available SimEvents blocks (see in Section 2.2).

Users (requests) and working threads (resources) have been generated by time- and event-based entity generators, combined by entity combiners, split by entity splitters. Feedbacks have been realized by path combiners. Response time has been monitored using start and read timers.

Application and global queues and their limits have been represented by queues and their capacities. User think time and dropping requests have been modeled as infinite server, thread pool as N-server with exponentially distributed think-, drop- and service time. Dropping method has been implemented by output switch and embedded function of its switching criteria according to *appRequestQueueLimit* and *requestQueueLimit* properties.

Simulation models have been provided for a given number of users and increasing number of users. Simulation has been performed over time using lower and higher queue limits. Simulation results are depicted in Fig. 4.

In case of a higher queue limit (5000), all requests can be served; there are no dropped requests (2nd and 4th rows). In case of a lower queue limit (100), some requests are rejected, dropped (1st and 3rd rows).

For a given number of users (150), the number of users in the system is fixed; the average response time tends to a steady value (1st and 2nd rows). For increasing number of users (from 0 to around 1700), the number of users is increasing, in case of a lower queue limit, the same number of users can be served, increasing number of dropped requests causes some slightly increasing overhead in response time (3rd row), however, in case of a higher queue limit, all requests can be served, response time is increasing linearly (4th row).

## 3.2 State Space Models

Secondly, analytical models in form of state space models using subspace identification have been provided examining various input factors and their effects to each other.

For provided simulation models subspace identification method has been applied and implemented in MATLAB. We would like to model and predict the performance, hence a performance metric has been chosen as output like average response time, and various performance factor candidates have been investigated as inputs like number of all users ($\mathbf{u}_1$), service time ($\mathbf{u}_2$), number of all dropped requests ($\mathbf{u}_3$), waiting time in queue ($\mathbf{u}_4$).
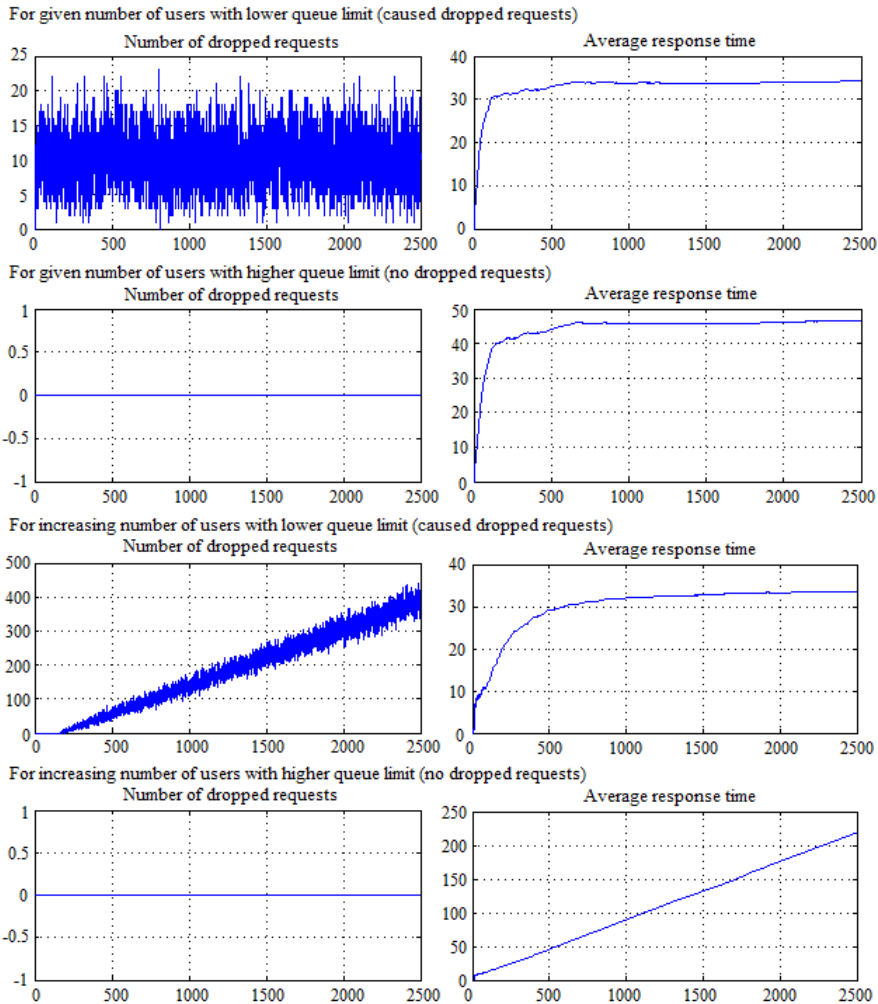
Figure 4
Simulation results (left column: $\mathbf{u}_3$ , right column: $\mathbf{y}$ )

The relationship between the above mentioned metric and the factor candidates, has been modeled by subspace identification (see Fig. 5). In case of input $\mathbf{u}_2$ the relationship is weak, since the service time is exponentially distributed. In case of inputs $\mathbf{u}_1$ and $\mathbf{u}_3$ the relationship is moderately strong (1st row). In case of input $\mathbf{u}_4$ the relationship is extremely strong (bottom right corner). Furthermore, combining inputs $\mathbf{u}_1$ and $\mathbf{u}_3$ – whose individual relationship to output is only moderately strong – the relationship between combined inputs $\mathbf{u}_{1,3}$ and output is stronger (bottom left corner).
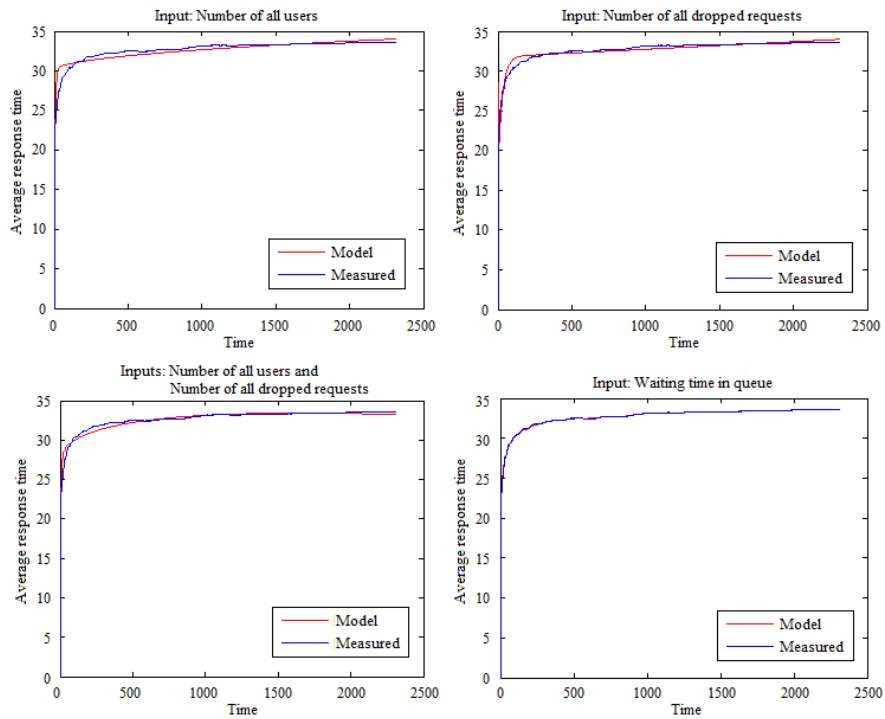
Figure 5

Relationship between input factors and output metric **y** in case of increasing number of users using lower queue limit

The system matrices of the state space model in case of the strongest relationship, namely, in case of input $\mathbf{u}_4$ using lower queue limit for increasing number of users are shown in Fig. 6. Since in this case the system has only one input and one output, **B** is a column vector, **C** is a row vector, **D** is a constant. As the model order 4 was chosen reflecting the best setting.

The detailed results in case of the strongest relationship, namely, in case of input $\mathbf{u}_4$ are depicted in Fig. 7 using lower and higher queue limits for a given and increasing numbers of users. The relationship is extremely strong in all cases.

The proposed method has been validated, and its correctness has been verified by comparing the results of simulation ($\mathbf{y}_{measured}$) and analytical models ($\mathbf{y}_{model}$) depicted in Figs. 5 and 7.

|   |   |   |   |   |
|---|---|---|---|---|
| **A** = | 0.78953 | 0.27977 | 0.47337 | -0.033463 |
|   | -0.80127 | 0.34866 | 0.23873 | 0.13293 |
|   | 0.28777 | -0.78382 | -0.18142 | 0.089506 |
|   | -0.39018 | 0.74561 | -1.0258 | 0.086449 |
| **B** = | -0.54887 |   |   |   |
|   | 0.84371 |   |   |   |
|   | 1.4622 |   |   |   |
|   | -2.5302 |   |   |   |
| **C** = | 4.8106 | -0.28536 | 2.2575 | -0.25992 |
| **D** = | 0 |   |   |   |

Figure 6

System matrices in case of input factor $\mathbf{u}_4$ and output metric $\mathbf{y}$ for increasing number of users using
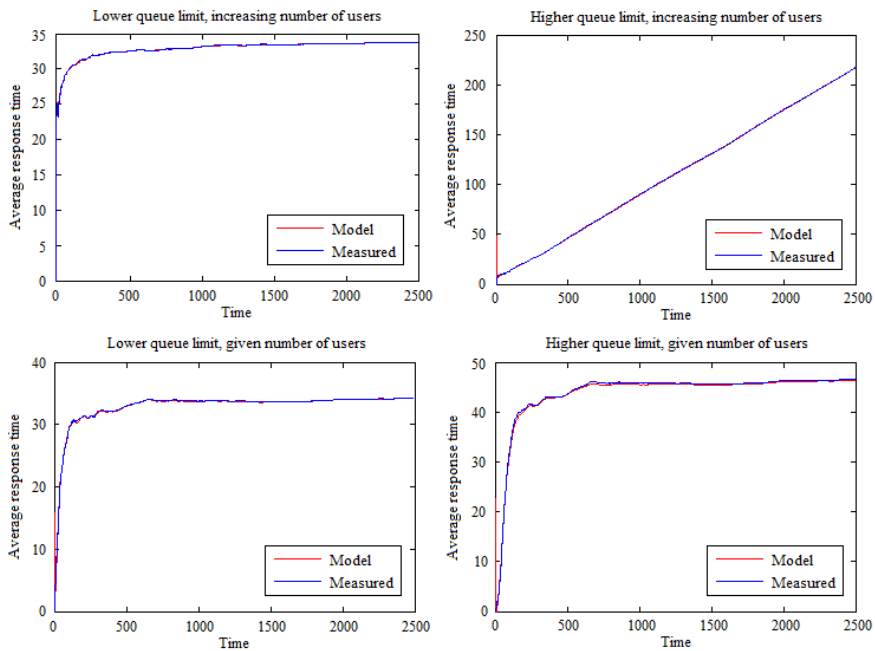lower queue limit



Figure 7

Relationship between input factor $\mathbf{u}_4$ and output metric $\mathbf{y}$

## 3.3    Performance Prediction

In addition, it has been demonstrated that the proposed method can be applied to performance prediction of web-based software systems. The state space model has been determined based on only the first half of the simulated data. For the second half it has been predicted by the proposed state space model.

The results in case of input $\mathbf{u}_4$ are depicted in Fig. 8. The relationship is extremely strong in all cases. The performance prediction facility has been validated and verified, by comparing the results of simulation ($\mathbf{y}_{measured}$) and analytical models ($\mathbf{y}_{model}$). The relationship without and with prediction has around the same strength comparing the results shown in Figs. 7 and 8.
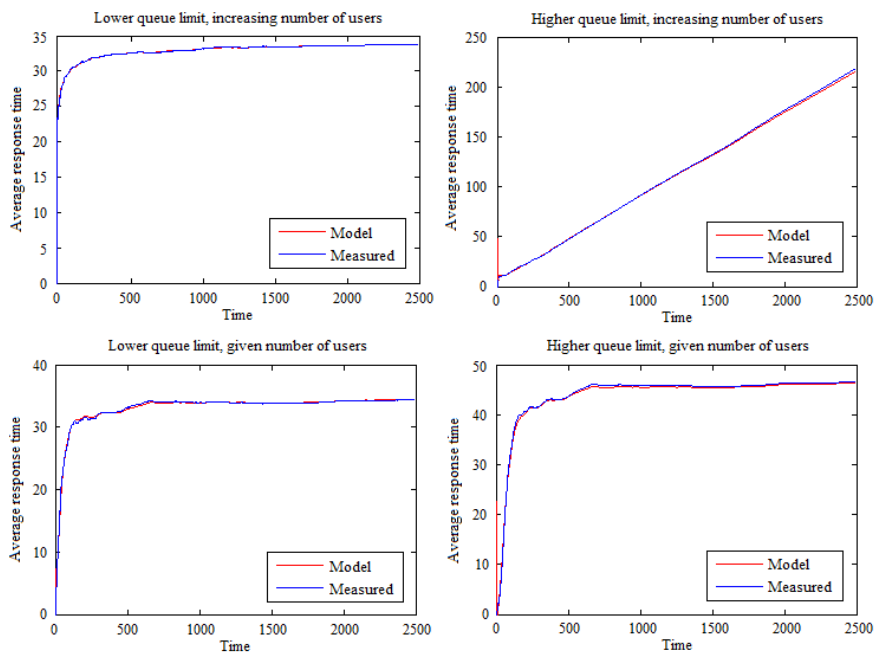


Figure 8

Relationship between input factor $\mathbf{u}_4$ and output metric $\mathbf{y}$ with prediction

## 3.4    Error Analysis

Error analysis has been performed by examining goodness of fit, firstly using lower and higher queue limits for a given and increasing number of users (without prediction), secondly comparing the results without and with prediction.

The model and measured outputs have been compared by the NRMSE (Normalized Root Mean Square Error) fitness value, namely, the goodness of fit.

***Definition 2.*** The goodness of fit is defined as follows where *norm* stands for L2 norm, and the associated notations are in Table 1.

$$\left( 1 - \frac{norm(\mathbf{y}_{measured} - \mathbf{y}_{model})}{norm(\mathbf{y}_{measured} - mean(\mathbf{y}_{measured}))} \right) \times 100 \tag{9}$$

The results (without prediction) are presented in Table 2. Recall that for input $\mathbf{u}_2$ the relationship is weak. In case of a higher queue limit, when there are no dropped requests, for input $\mathbf{u}_3$ the relationship is weak. In case of a lower queue limit, when some requests are dropped, for input $\mathbf{u}_3$ the relationship is moderately strong. For input $\mathbf{u}_1$ the relationship is (moderately) strong. For input $\mathbf{u}_4$ the relationship is extremely strong.

Furthermore, combining inputs, whose individual relationship to output is only moderately strong, the relationship between combined inputs $\mathbf{u}_{1,3}$ and output is stronger.

Column "Increasing number of users" "Queue limit" "100" corresponds to Fig. 5, row "$\mathbf{u}_4$" to Fig. 7.

Table 2
Goodness of fit in percentages for different inputs (without prediction)

| Input | Increasing number of users | | Given number of users | |
|---|---|---|---|---|
| | Queue limit | | | |
| | 100 | 5000 | 100 | 5000 |
| $\mathbf{u}_1$ | 67.39 | 95.74 | 81 | 80.41 |
| $\mathbf{u}_2$ | - | - | - | - |
| $\mathbf{u}_3$ | 69.87 | - | 55.79 | - |
| $\mathbf{u}_4$ | 96.71 | 98.16 | 89.21 | 92.76 |
| $\mathbf{u}_{1,3}$ | 85.35 | 95.01 | 86.4 | 83.32 |

The relationship without and with prediction is extremely strong in all cases, and it has around the same strength comparing the results of goodness of fit shown in Table 3. Column "Without prediction" corresponds to Fig. 7, while column "With prediction" to Fig. 8.

Table 3
Goodness of fit in percentages for input $\mathbf{u}_4$ without and with prediction

| | | | Without prediction | With prediction |
|---|---|---|---|---|
| **Increasing number of users** | **Queue limit** | **100** | 96.71 | 96.73 |
| | | **5000** | 98.16 | 97.54 |
| **Given number of users** | | **100** | 89.21 | 88.35 |
| | | **5000** | 92.76 | 92.72 |

## 3.5   Propositions

To summarize, we can say that the subspace identification method can be applied to performance factor identification and performance prediction.

*Proposition 1.* The subspace identification method can be applied to performance factor identification.

It is shown the way in which the subspace identification method can be applied to identifying performance factors. The input is one (or more) performance factor candidate(s). The output is one (or more) performance metric(s).

The relationship between the factor candidate and the metric is modeled by subspace identification. If the relationship cannot be calculated, in other words, weak, then it is not a performance factor. If the relationship can be calculated, the goodness of fit in percentages is high, the performance factor candidate influences the performance, namely, a novel performance factor is identified.

The related experimental results are demonstrated in Sections 3.2 and 3.4.

*Proposition 2.* The subspace identification method can be used for performance prediction.

It is shown the way in which the subspace identification method can be applied to predicting performance. The input is one (or more) identified performance factor(s) by Proposition 1. The output is one (or more) performance metric(s).

The relationship between the factor candidate and the metric is modeled by subspace identification. The state space model is determined based on only the first part of the measured or simulated data. The other part is predicted by the provided state space model.

The related experimental results are shown in Sections 3.3 and 3.4.

**Conclusions**

Recently, web-based software systems have proliferated. Their performance modeling and prediction are relevant issues. In this work, proper mathematical models have been established in form of difference equations by subspace identification in order to model and predict their performance.

First, simulation models have been provided using SimEvents of MATLAB Simulink to simulate the thread pool and queued requests behavior of web-based software systems as well as to obtain input-output data for subspace identification method. Second, analytical models in form of state space models using subspace identification have been provided investigating various input factors and their effects to each other. In addition, it has been shown that the proposed method can be applied to performance prediction of web-based software systems. The proposed models have been validated and their correctness has been verified by comparing the results of simulation and analytical models, as well as, by error analysis. Above all, performance factor identification and performance prediction techniques have been proposed based on subspace identification.

**Acknowledgement**

**References**

[1]     R. Jain: The Art of Computer Systems Performance Analysis, John Wiley and Sons, 1991

[2]     T.G. Robertazzi: Computer Networks and Systems: Queueing Theory and Performance Evaluation, Springer, Cambridge, 2000

[3]     B. Urgaonkar, G. Pacifici, P. Shenoy, M. Spreitzer, A. Tantawi: An Analytical Model for Multi-tier Internet Services and its Applications, ACM SIGMETRICS Performance Evaluation Review, 2005, Vol. 33, No. 1, pp. 291-302

[4]     Á. Bogárdi-Mészöly, Z. Szitás, T. Levendovszky, H. Charaf: Investigating Factors Influencing the Response Time in ASP.NET Web Applications, Lecture Notes in Computer Science, Springer, 2005, Vol. 3746, pp. 223-233

[5]     Á. Bogárdi-Mészöly: Improved Performance Models for Web-based Software Systems, Modeling Thread Pool and Queue Limit Performance Factors, Lambert Academic Publishing, Saarbrücken, 2010, 132 p.

[6]     Á. Bogárdi-Mészöly, T. Levendovszky: A Novel Algorithm for Performance Prediction of Web-based Software Systems, Performance Evaluation, Elsevier, 2011, Vol. 68, No. 1, pp. 45-57

[7]     Á. Bogárdi-Mészöly, A. Rövid, S. Yokoyama: Subspace Identification for Web-based Software Systems, International Conference on Engineering and Applied Science, Japan, Sapporo, July 20-22, 2015, pp. 119-127

[8]     Á. Bogárdi-Mészöly, A. Rövid, S. Yokoyama: Performance Prediction of Web-based Software Systems with Subspace Identification, International Scientific Conference on Engineering and Applied Sciences, Japan, Naha, July 29-31, 2015, pp. 89-98

[9]     D. Carmona: Programming the Thread Pool in the .NET Framework, in .NET Development (General) Technical Articles, Microsoft Spain, 2002, http://msdn.microsoft.com/en-us/library/ms973903.aspx

[10]   T. Marquardt: ASP.NET Performance Monitoring, and When to Alert Administrators, in ASP.NET Technical Articles, 2003, http://msdn.microsoft.com/en-us/library/ms972959.aspx

[11]   SimEvents, MATLAB, MathWorks, http://www.mathworks.com/help/simevents/

[12]   P. Overschee, B. Moor, Subspace Identification for Linear Systems: Theory – Implementation – Applications, Kluwer Academic Publishers, 2011

[13]   P. Várlaki, T. Vadvári, Queueing Models and Subspace Identification in Logistics, Acta Technica Jaurinensis, 2014, Vol. 8, No. 1, pp. 63-76

# Theory of Acoustic Metamaterials and Metamaterial Beams: An Overview

## Livia Cveticanin, Gyula Mester

Óbuda University, Doctoral School of Safety and Security Sciences
Népszínház utca 8, H-1081 Budapest, Hungary
cpinter.livia@bgk.uni-obuda.hu, gmester@inf.u-szeged.hu

*Abstract: This paper presents a theoretical examination of acoustic metamaterials and their application in vibration absorption. Acoustic metamaterials are concerned as analogy to electromagnetic metamaterials which are suitable for refraction and decline of electromagnetic waves at certain frequencies. Due to the analogy with these materials, the acoustic metamaterials is required to have negative effective (dynamic) mass to enable vibration elimination at the certain frequency. The concept of negative effective mass is explained based on the motion of an externally excited mass-in-mass system where the vibration elimination at the certain frequency is due to the mass-spring unit. Using these vibration absorber units, the acoustic metamaterial beams are made. Depending on the way how the units are attached to the beam, the structure may absorb waves in one-direction (for example, longitudinal waves) or waves in two directions (such as, for instance, transversal and longitudinal waves). Moreover, according to the frequency properties of the absorber units the acoustic metamaterial beams may give one, two or multi-frequency gaps. This work provides an overview of the mathematical models of acoustic metamaterial beams and also contains some suggestions for future work.*

*Keywords: elastic metamaterial; acoustic metamaterial beam; spring-mass systems; negative effective mass*

# 1 Introduction

One of the requirements of environmental and occupant health protection is noise reduction and elimination of sources of sound pollution. Various methods are developed to damp and reduce the acoustic influence on the health of a population. One of the most often applied methods is based on the use of materials for acoustic isolation which absorb the acoustic energy. It is known that acoustic waves which come to a surface generate waves in the surface itself and these waves transmit sound power through the surface to the other side [1]. The fraction of the sound power that is transmitted to the other side is known as the transmission coefficient $T_c$ and the isolation effectiveness is expressed through the transmission loss *TL* which is defined as:

$$TL = 10log_{10}\left(\frac{1}{T_c}\right) \tag{1}$$

where TL is in dB. The frequencies which are of interest to be absorbed are in the low range of human hearing, approximately 100 Hz to 1 kHz, as for these values the irritation to health occurs. Usually, the acoustic wave transmission loss is controlled by the so-called 'mass law': the isolation mechanism is the inertia provided by the mass of the isolating partition. The transmission loss for normally incident wave propagation through a homogenous solid is estimated as:

$$TL = 10log_{10}\left[1 + \left(\frac{\pi\rho h f}{c\rho_a}\right)^2\right] \tag{2}$$

where $\rho$ is the density and $h$ is the thickness of the solid, $f$ is the frequency, $\rho_a$ is the density of air and $c$ is the speed of sound in the air. In [1] it is stated that at a given frequency the level of sound transmitted through a partition will be reduced by 5-6 dB for every doubling of the mass of the partition.

To improve the reduction of the acoustic effect in buildings, machinery, ships and other applications instead of acoustic isolators the novel acoustic absorbers are suggested. Their working mechanism is based on the concept of conventional vibration absorbers. As it is well known, the conventional vibration absorber consists of a lumped mass $m_2$ attached with a linear spring $k_2$ to the mechanical system with mass $m_1$ (see Fig. 1). If the excitation force acts on the mass $m_1$ differential equations of motion are:

$$m_1\ddot{u}_1 + k_2(u_1 - u_2) = F_0\exp(i\omega t) \tag{3}$$

$$m_2\ddot{u}_2 + k_2(u_2 - u_1) = 0 \tag{4}$$

Solutions of (3) and (4) have the form:

$$u_1 = a_1\exp(i\omega t), \quad u_2 = a_2\exp(i\omega t) \tag{5}$$

where:

$$a_1 = \frac{F_0(k_2 - m_2\omega^2)}{(k_2 - m_1\omega^2)(k_2 - m_2\omega^2) - k_2^2}$$

$$a_2 = \frac{F_0 k_2}{(k_2 - m_1\omega^2)(k_2 - m_2\omega^2) - k_2^2} \tag{6}$$

and $i = \sqrt{-1}$ is the imaginary unit. For this model only one local resonance frequency exists. The vibration absorber uses the 1:1 external resonance between the forcing frequency on the main system $\omega$ and the local resonance frequency of the absorber $\omega_2 = \sqrt{k_2/m_2}$ to transform the vibration energy to the absorber and stop the main system motion ($u_1$=0). Based on this conception an idea of a new material, which would be the acoustic absorber is developed. Moreover, motivated by the mathematical analogy between acoustic and electromagnetic waves the investigation were directed toward so-called metamaterials which exhibit exceptional properties not observed in nature or in the constituent materials [2].
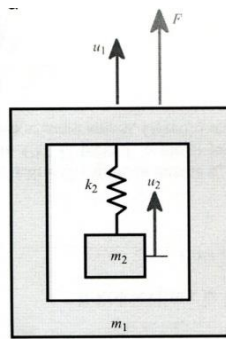
Figure 1
Mass-in-mass model [3]

Acoustic (elastic) metamaterials have to be the counterpart to electromagnetic metamaterials. The main property of electromagnetic metamaterials is that they have negative permittivity and magnetic permeability which result in a negative refractive index [4-7]. According to analogy it has to be asked the acoustic metamaterial to be with negative mass density and negative modulus. The negative effective mass/mass density is not the physical property of the material but is the result of inaccurate modeling of acoustic metamaterials.

## 2    Effective Mass Density for Mass-in-Mass System

Let us consider the mass-in-mass model (Fig.1) as a subunit of a metamaterial which is suggested to be identified with a single mass with effective mass $m_{eff}$ whose motion is the same as that of $m_1$. The effective mass is defined by treating this two-degree-of-freedom system as a one-degree-of-freedom one by assuming the internal absorber being unknown to the observer. In other words, the identity of the internal mass $m_2$ would be ignored and its effect would be absorbed by the introduction of an effective mass $m_{eff}$ as shown in Fig. 2.
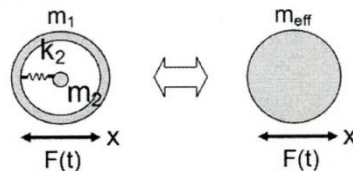


Figure 2
Identity of the mass-in-mass model and of the model with effective mass [8]

If the motion of the mass $m_1$ is $u_1$, the effective mass has also the motion $u_1$. Linear momentums for the both models given in Fig. 2 have to be equal, i.e.:

$$m_{eff}\frac{du_1}{dt} = m_1\frac{du_1}{dt} + m_2\frac{du_2}{dt} \tag{7}$$

Substituting the assumed solution (5) into (7) it is:

$$m_{eff}a_1 = m_1a_1 + m_2a_2 \tag{8}$$

Motion of the mass $m_2$ is given with the equation (4). Substituting the assumed solutions (5) we have

$$-m_2\omega^2a_2 + k_2(a_2 - a_1) = 0 \tag{9}$$

After some modification equations (8) and (9) yield the effective mass [8-10]:

$$m_{eff} = m_1 + m_2\frac{k_2}{k_2-m_2\omega^2} \tag{10}$$

which is for: $\omega_2 = \sqrt{k_2/m_2}$

$$m_{eff} = m_1 + m_2\frac{\omega_2^2}{\omega_2^2-\omega^2} \tag{11}$$

Analyzing the relation (11) it is obvious that the effective mass depends on the ration between the excitation frequency $\omega$ and natural frequency of the system $\omega_2$:

$$m_{eff} = m_1 + m_2\frac{1}{1-\frac{\omega^2}{\omega_2^2}} \tag{12}$$

For the system three modes of motion are evident: 1) acoustic mode when $\omega<\omega_2$, 2) resonant mode when $\omega=\omega_2$ and 3) optic mode when $\omega>\omega_2$. For the acoustic mode the effective mass is positive. In the resonant mode the effective mass is theoretically infinite. In the optical mode the effective mass is negative for:

$$m_1 + m_2\frac{1}{1-\frac{\omega^2}{\omega_2^2}} < 0. \tag{13}$$

Otherwise, it is positive. In Fig. 3, according to (12) the effective mass - frequency diagram is plotted

Differentiating the relation (7) we have:

$$(m_{eff} - m_1)\ddot{u}_1 = m_2\ddot{u}_2 \tag{14}$$

Substituting (14) into (4) we obtain:

$$k_2(u_2 - u_1) = -(m_{eff} - m_1)\ddot{u}_1 \tag{15}$$

Equation (3) and (15) give:

$$-m_{eff}\ddot{u}_1 = F_0\exp(i\omega t) \tag{16}$$

The effective mass is the ratio between the excitation force and acceleration of the mass $m_1$:

$$m_{eff} = \frac{F}{\ddot{u}_1} = \frac{F_0}{-\omega^2a_1} = -\frac{F}{\omega^2u_1} \tag{17}$$
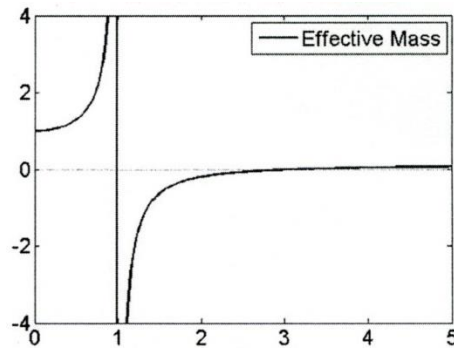
Figure 3
Dimensionless effective mass $m_{eff}/m_1$ as a function of $\omega/\omega_2$ [10]

According to Fig. 3 it is obvious that for $\omega_2=\omega$ the effective mass tends to infinity. For that value the motion of mass $m_1$ is zero and the inertial force of the mass $m_2$ is equal to the excitation force: $F(t) = m_2\ddot{u}_2$. So, the external force is eliminated with the inertia force $-m_2\ddot{u}_2$ through the spring $k_2$. This is the concept of vibration absorbers.

Finally, the following is concluded:

In the acoustic mode, when $\omega<\omega_2$, the effective mass $m_{eff}$ is positive and the motions $u_1$ and $u_2$ are in phase. For the optical mode, when $\omega>\omega_2$, the effective mass $m_{eff}$ may be positive or negative, while the displacements $u_1$ and $u_2$ are $180^0$ out of phase. Then, the absorber works efficiently in the optical mode against the external acting on the mass $m_1$. The excitation is absorbed with the inertial force.

# 3 Mechanical Structure with Negative Effective Mass

Mechanical structures are designed by incorporating of the previously mentioned mechanical subunits in a natural material with the aim to resonate during mechanical wave propagation in it. For the local mechanical resonance the designed structures have negative effective masses. The negative mass behavior is generated by oscillating of resonant structures within the material with $180^0$ out of phase to the acoustic waves which apply to surface. It causes existence of frequency bands where wave propagation is theoretically impossible. These bands are usually called band gaps. The aim of the designed structure is to overcome the limitations of the mass law for solids by creating engineered materials with useful acoustic band gaps, and the key to the generation of these band gaps is an inhomogeneous structure. Huang et al. [10] composed a one-dimensional lattice which contains mass-in-mass lattices (Fig. 4). The model is based on those with negative mass as explained in [7] and [8]. Equations of motion for the unit cell are:

$$m_1^{(j)} \ddot{u}_1^{(j)} + k_2\left(u_1^{(j)} - u_2^{(j)}\right) + k_1\left(2u_1^{(j)} - u_1^{(j-1)} - u_1^{(j+1)}\right) = 0$$

$$m_2^{(j)} \ddot{u}_2^{(j)} + k_2\left(u_2^{(j)} - u_1^{(j)}\right) = 0 \tag{18}$$

where $k_1$ is the rigidity of connection.
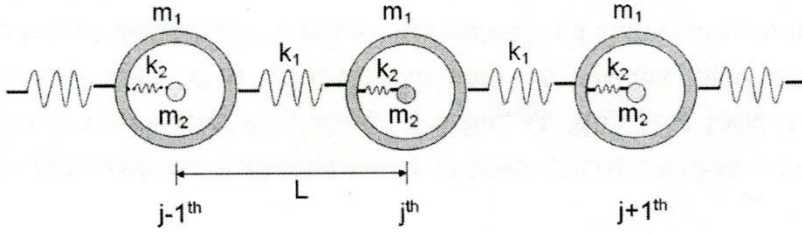


Figure 4

Model of subunits connected in lattice [11]

The harmonic wave solution for (18) is assumed in the form:

$$u_1^{(j)} = U_1 \exp(i\beta x - i\omega t)$$

$$u_2^{(j)} = U_2 \exp(i\beta x - i\omega t)$$

$$u_1^{(j+1)} = U_1 \exp(i\beta x - i\omega t)\exp(i\beta L)$$

$$u_1^{(j-1)} = U_1 \exp(i\beta x - i\omega t)\exp(-i\beta L) \tag{19}$$

Substituting (19) into (18) two homogenous equations for $U_1$ and $U_2$ follow which give the dispersion equation:

$$m_1 m_2 \omega^4 - [(m_1 + m_2)k_2 + 2m_2 k_1(1 - \cos(\beta L))]\omega^2 + 2k_1 k_2(1 - \cos(\beta L)) = 0 \tag{20}$$

In Fig. 5, both branches of the band structure (20) which correspond to the optical mode, when $\omega > \omega_0$, and to the acoustic mode, when $\omega < \omega_0$, are plotted. The parameter values are $m_2/m_1=9$, $k_2/k_1=0.1$ and $\omega_0=\sqrt{k_2/m_2}=149.07 \text{s}^{-1}$. The frequency distribution is given as a function of the wave number $\beta L$.

The displacements in (19) are functions of $\exp(i\beta x)$. For the case when the wave number has a complex form:

$$\beta L = \gamma + i\alpha \tag{21}$$

it is

$$u \propto \exp\left(\frac{i\gamma x}{L}\right)\exp\left(-\frac{\alpha x}{L}\right). \tag{22}$$
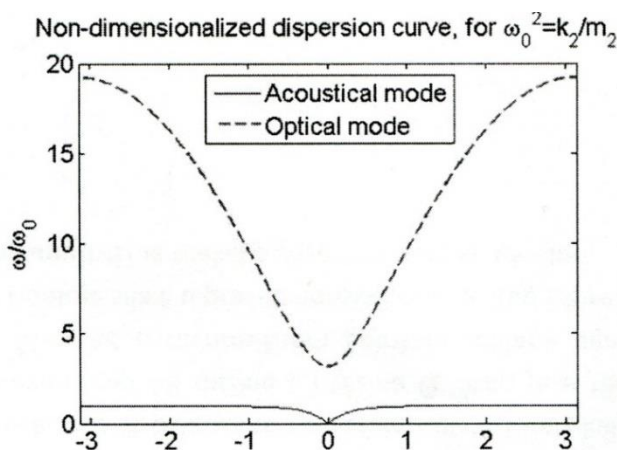
Figure 5
Nondimensional dispersion curve for the mass-in-mass lattice model [10]

The amplitude of the displacement decays as the exponential function $\exp(-\frac{\alpha x}{L})$ if the attenuation factor $\gamma$ is positive. It is of interest to analyze the case when wave frequency approaches the local resonance frequency $\omega_0$ and the attenuation factor $\gamma$ theoretically becomes unbounded.

If the lattice system is reduced to a homogenous mono-atom lattice system in which only effective masses $m_{eff}$ are connected by springs with rigidity $k_1$, the homogenous lattice system has the dispersion equation:

$$\omega^2 = \frac{2k_1}{m_{eff}}(1 - \cos(\beta L)) \tag{23}$$

The suggested model is equivalent to the original mass-in-mass system if their dispersion equations (20) and (23) are identical. Substituting (23) into (20) the effective mass is obtained:

$$m_{eff} = m_{st} + \frac{m_2(\frac{\omega}{\omega_0})^2}{1-(\frac{\omega}{\omega_0})^2} \tag{24}$$

where $m_{st} = m_1 + m_2$. Comparing (24) with the result (12) it is seen that they are identical and that the diagram shown in Fig. 3, is also valid here. Let us extend the discussion of the mentioned diagram. Especially the vibrations out of resonant region are considered. Thus, for frequencies far below the local resonance, the internal masses oscillate in phase with the solid in which they are embedded. But as the excitation frequencies passes through the resonance and the effective mass becomes negative near the local resonance frequency $\omega_0$, the phase angle of the response changes by close to $180^0$. That means that near this frequency the acceleration of the resonant structures within the metacomposite will have a component whose direction is opposite to that of the force of pressure applied to the surface of the metacomposite.

However, the most important rule is the resonant one. Analyzing relations (23) and (24) it is obvious that the effective mass is negative only when $(1-\cos(\beta L))$ is negative and when $k_1$ and $\omega^2$ are positive. It means that the dimensionless wavenumber $\beta L$ has to be complex. Thus, frequencies corresponding to the negative mass are in the stopping band. In other words, a negative effective mass in the equivalent mass-spring lattice system yields spatial attenuation in wave amplitude. For a material with a negative mass density the speed of sound is imaginary, and therefore, only evanescent waves, which decay exponentially from the surface, can exist.
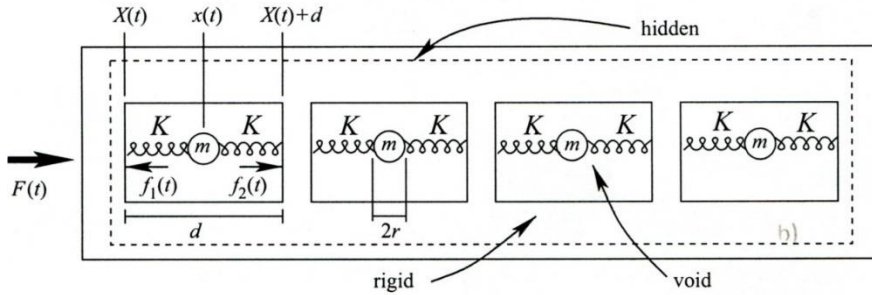


Figure 6
A one-dimensional material where the mass depends on the frequency $\omega$ and can be negative [8].

Based on this theoretical consideration, Milton and Willis [8], developed a bar-like acoustic metamaterial with heterogeneous material properties which are valid for waves of wavelengths much longer than the sizes of subunits. Milton and Willis [8] concluded that solids containing many identical small resonators [8] exhibit band gap behavior near their resonance frequency, even though the size and spacing of the resonators were over a 100 times smaller than the wavelength at the band gap frequency. Furthermore, lack of order of periodicity in the arrangement of the resonators did not materially affect the results, and the band gap frequencies could be tuned through changes to the resonators natural frequencies. These locally resonant sonic materials (LRSM), shown in Fig. 4, are considered to be a type of acoustic absorbers [11-13]. The model of the structure is one-dimensional [8]. From a bar, made of rigid material, cylindrical cavities of length $d$ have been carved out. Each cavity is modeled as a sphere of mass $m$ and radius $r$. The sphere is attached to the ends of the cavity with two possibly viscoelastic springs each having the same complex spring constant $K$. It is assumed that everything is varied harmonically with time with frequency $\omega$. The model is plotted in Fig. 6.

In Fig. 7 the fabricated lattice beam is shown. The model is suitable to describe the propagation of the dispersive wave in the lattice.
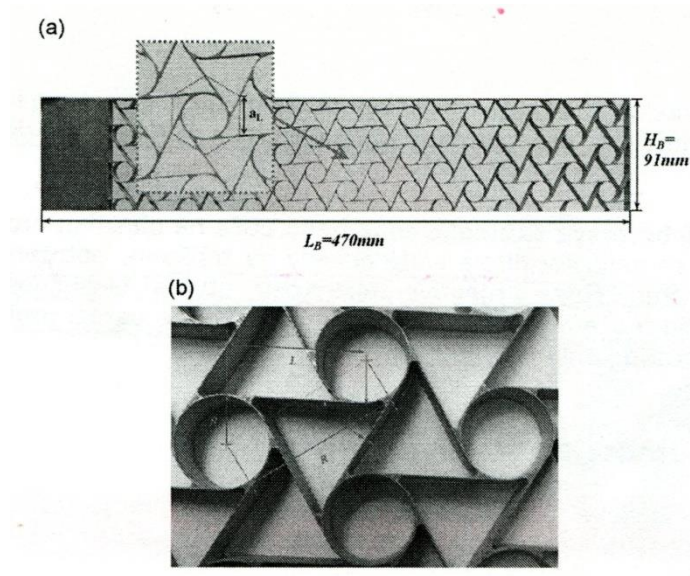
Figure 7
a) The fabricated chiral lattice beam and its zoomed unit cell and b) the topology of the hexagonal chiral lattice [14]

This structure is non-homogenous. If it is required that the metamaterial behave like a homogeneous material described by its averaged material properties, its subunits must be much smaller than the shortest wavelength of waves propagating in it. The averaging may result in the existence of a useful but mysterious phononic stop-band that allows no waves within that frequencies range to propagate forward, and most current designs of acoustic metamaterials are based on the stop-band effect [10]. For manufacturing such metamaterials with tiny subunits in order to have stop-bands, expensive manufacturing techniques are required including micro and nano-manufacturing technologies.

# 4    Design of Metamaterial Beams for Broadband Absorption and Isolation

To eliminate the lack of the previously mentioned beam-absorbers made of LRSM and of the chiral lattice beam, acoustic absorber based on a metamaterial beam is developed. The metamterial beam consists of a uniform isotropic beam with many small spring-mass-damper subsystems integrated at separated locations along the beam to act as vibration absorbers (Fig. 8). The spring-mass-damper subsystems create a stop-band in which no elastic waves in this frequency range can propagate forward. The concepts of negative effective mass and stiffness is applied for metamaterial design.
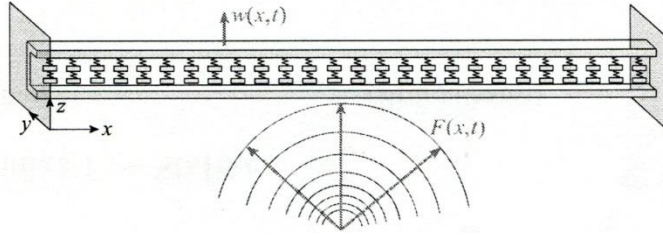
Figure 8
Model of a metamaterial beam for vibration absorption [2].

The proposed metamaterial beam is based on the concept of conventional mechanical vibration absorbers. It uses the incoming elastic wave in the beam to resonate the integrated spring-mass-damper absorbers to vibrate in their optical mode at frequencies close to but above their local resonance frequencies to create shear forces and bending moments to straighten the beam and stop the wave propagation. Metamaterials usually designed as metacomposites have unusual response to elastic waves.

## 4.1   Metamaterial Beams for Broadband Absorption of Longitudinal Elastic Waves

In Fig. 9, a metamaterial beam for acoustic absorption is presented. In the longitudinal beam spring-mass subsystems are integrated. The beam has to absorb low and high frequency elastic waves [3]. The absorption of the longitudinal elastic waves is based on the negative mass effect. The concept of the negative effective mass and vibration absorbers with two-degrees-of-freedom, mass-in-mass system, forced with harmonic excitation is demonstrated in Section 2. It is shown that the negative effective mass density gives the stop band when the incoming wave frequency is slightly higher than the local frequency.



Figure 9
Wave propagation in a metamaterial beam [12].

Using the concept of vibration absorbers, the unit cell model is shown in Fig. 10. The equations for a unit cell of an infinite metamaterial beam can be derived using the extended Hamiltonian principle.

$$\int_0^L (\delta T - \delta U + \delta W_{nc}) = 0 \qquad (25)$$

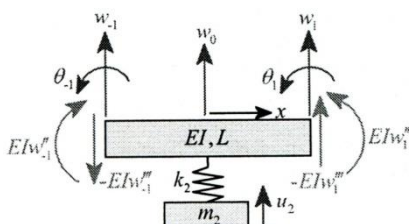where $T$ is kinetic energy, $U$ is the elastic energy and $W_{nc}$ is the work of the nonconservative forces.

Figure 10
The unit cell model

For the unit cell the kinetic energy is:

$$\delta T = -\int_{-\frac{L}{2}}^{\frac{L}{2}} \rho A dx \frac{\partial^2 w}{\partial t^2} \delta w - m_2 \frac{\partial^2 u_2}{\partial t^2} \delta u_2 \tag{26}$$

the elastic energy is:

$$\delta U = k_2(u_2 - w_0)\delta(u_2 - w_0) + \int_{-\frac{L}{2}}^{\frac{L}{2}} EI \frac{\partial^4 w}{\partial x^4} dx \delta w$$

$$+EI\left(\frac{\partial^2 w_1}{\partial x^2}\right)\delta\left(\frac{\partial w_1}{\partial x}\right) - EI\left(\frac{\partial^2 w_{-1}}{\partial x^2}\right)\delta\left(\frac{\partial w_{-1}}{\partial x}\right)$$

$$- EI\left(\frac{\partial^3 w_1}{\partial x^3}\right)\delta w_1 + EI\left(\frac{\partial^3 w_{-1}}{\partial x^3}\right)\delta w_{-1} + EI\left(\frac{\partial^3 w_{0+}}{\partial x^3}\right)\delta w_0 + EI\left(\frac{\partial^3 w_{0-}}{\partial x^3}\right)\delta w_0 \tag{27}$$

and the work of the nonconservative forces is:

$$\delta W_{nc} = EI\left(\frac{\partial^2 w_1}{\partial x^2}\right)\delta\left(\frac{\partial w_1}{\partial x}\right) - EI\left(\frac{\partial^2 w_{-1}}{\partial x^2}\right)\delta\left(\frac{\partial w_{-1}}{\partial x}\right)$$

$$- EI\left(\frac{\partial^3 w_1}{\partial x^3}\right)\delta w_1 + EI\left(\frac{\partial^3 w_{-1}}{\partial x^3}\right)\delta w_{-1} \tag{28}$$

where $\left(\frac{\partial w}{\partial x}\right) = \theta$. It is $EI\left(\frac{\partial^3 w_{0+}}{\partial x^3}\right) \neq EI\left(\frac{\partial^3 w_{0-}}{\partial x^3}\right)$ because the vibration absorber creates a concentrated shear force at $x=0$. Substituting $(26) - (28)$ into $(25)$ and separating the terms with $\delta w$ and $\delta u_2$ the following equations are obtained:

$$0 = -\rho A \frac{\partial^2 w}{\partial t^2} - EI \frac{\partial^4 w}{\partial x^4} + \left[k_2(u_2 - w_0) + EI\left(\frac{\partial^3 w_{0+}}{\partial x^3}\right) + EI\left(\frac{\partial^3 w_{0-}}{\partial x^3}\right)\right]\delta(x)$$

$$m_2 \frac{\partial^2 u_2}{\partial t^2} + k_2(u_2 - w_0) = 0 \tag{29}$$

where $\delta(x)$ is the Dirac function. Due to periodicity along $x$ direction, flexural wave propagation through the infinite periodic beam can be expressed in a harmonic form:

$$w(x,t) = Wexp(i\beta x - i\omega t), \qquad w_0(x,t) = W_0 exp(-i\omega t)$$

$$u_2(x,t) = U_2 exp(-i\omega t) \tag{30}$$

where $i=\sqrt{-1}$ is the imaginary unit, $\beta$ is the wave number and $\omega$ is the vibration frequency.

If the metamaterial is treated as a homogenized uniform beam the integration over the whole length gives:

$$0 = k_2(u_2 - w_0) - \int_{-\frac{L}{2}}^{\frac{L}{2}} \rho A \frac{\partial^2 w}{\partial t^2} dx - EI\left(\frac{\partial^3 w_1}{\partial x^3}\right) + EI\left(\frac{\partial^3 w_{-1}}{\partial x^3}\right) \tag{31}$$

Using the harmonic wave solution and the assumption that the beam segment is a lumped mass $\tilde{m}$ and a spring $\tilde{k}$ system, we have:

$$\tilde{m}\frac{\partial^2 w_0}{\partial t^2} + \tilde{k}w_0 + k_2(u_2 - w_0) = 0 \tag{32}$$

where:

$$\tilde{m} = -\frac{2\rho A sin\left(\frac{\beta L}{2}\right)}{\beta}, \qquad \tilde{k} = -2EI\beta^3 sin\left(\frac{\beta L}{2}\right) \tag{33}$$

The values $\tilde{m}$ and $\tilde{k}$ depend on the wavenumber $\beta$. Substituting (30) into (29) and (32), and equating the determinant of corresponding parameters with zero, we have:

$$\begin{bmatrix} -\tilde{m}\omega^2 + \tilde{k} + k_2 & -k_2 \\ -k_2 & -m_2\omega^2 + k_2 \end{bmatrix} = 0 \tag{34}$$

The dispersion equation follows as:

$$(-\tilde{m}\omega^2 + \tilde{k} + k_2)(-m_2\omega^2 + k_2) - k_2^2 = 0 \tag{35}$$

To each specific value of $\omega$ the positive value of $\beta$ has to be determined. If no positive values for $\beta$ exist, the value of $\omega$ is the stop band. Namely, if $\beta=i\alpha$ and $\alpha>0$, and:

$$w(x,t) = W exp(-\alpha x)exp(-i\omega t) \tag{36}$$

there is an evanescent non-propagating wave and the stop band exists.

## 4.2    Metamaterial Beams with Two Spring-Mass Systems for Broadband Absorption of Longitudinal Elastic Waves

The metamaterial previously considered has single-mass absorbers and produces a stop band at the high-frequency side. Pai et al., [12] presented a new metamaterial beam based on multi-frequency vibration absorbers for broadband vibration absorption.

The metamaterial beam consists of a uniform isotropic beam and small two-mass spring-mass-damper subsystems at many locations along the beam to act as multi-frequency vibration absorber. This type of absorber produces two stop-bands.

The existence of stop-bands is caused by spring-mass-damper subsystems and their existence is explained according to negative effective mass and effective stiffness.

For an incoming wave with a frequency in one of the two stop-bands, the absorbers are excited to vibrate in their optical modes to create shear forces to straighten the beam and stop the wave propagation. For an incoming wave with a frequency outside of but between the two stop-bands, it can be damped by the damper with the second mass of each absorber. So, the stop-bands are connected into a wide stop-band.

To show the use of multi-frequency vibration absorbers to design broadband metamaterials Pai et al. [2] consider the three-degree-of-freedom oscillating system shown in Fig. 11. On the mass $m_1$, connected with the spring $k_1$ and damping $c_1$, a harmonic excitation force acts. The vibration absorber uses two lumped masses $m_2$ and $m_3$ connected with springs $k_2$ and $k_3$ and damping coefficients $c_2$ and $c_3$.



Figure 11
Three-degree-of-freedom oscillatory system [2]

We describe the motion of the system with a system of coupled differential equations

$$m_1\ddot{u}_1 + (c_1 + c_2)\dot{u}_1 - c_2\dot{u}_2 + (k_1 + k_2)u_1 - k_2u_2 = F_0\exp(i\omega t)$$

$$m_2\ddot{u}_2 + (c_2 + c_3)\dot{u}_2 - c_2\dot{u}_1 - c_3\dot{u}_3 + (k_2 + k_3)u_2 - k_2u_1 - k_3u_3 = 0$$

$$m_3\ddot{u}_3 + c_3(\dot{u}_3 - \dot{u}_2) + k_3(u_3 - u_2) = 0 \qquad (37)$$

The aim of the absorber is to make $u_1=0$ by adjusting one of two local natural frequencies equal to the excitation frequency $\omega$. For the undamped system for which $u_1=0$ equations (37) simplify to:

$$m_2\ddot{u}_2 + (k_2 + k_3)u_2 - k_3u_3 = 0$$

$$m_3\ddot{u}_3 + k_3(u_3 - u_2) = 0 \qquad (38)$$
where the natural frequencies are:

$$\omega_{1n}, \omega_{2n} = \omega_2 \sqrt{\frac{\bar{m} + \bar{k} + \bar{m}\bar{k} \pm \sqrt{(\bar{m} + \bar{k} + \bar{m}\bar{k})^2 - 4\bar{m}\bar{k}}}{2\bar{m}}} \tag{39}$$

where $\bar{m} = \frac{m_3}{m_2}, \bar{k} = \frac{k_3}{k_2}, \omega_2 = \sqrt{\frac{k_2}{m_2}}$. Around these two frequencies two stop-bands exist.

The goal is to design metamaterial beam based on the vibration absorber shown in Fig. 12 that can stop propagation of elastic waves. The unit cell model is plotted in Fig. 12a. Using the relation (25), equations for a unit cell model are derived.



Figure 12
Model of an infinite metamaterial beam: a) a unit cell, b) an infinite beam [2]

The kinetic energy, elastic energy and the work of the non-conservative forces are:

$$\delta T = -\int_{-\frac{L}{2}}^{\frac{L}{2}} \rho A dx \frac{\partial^2 w}{\partial t^2} \delta w - m_2 \frac{\partial^2 u_2}{\partial t^2} \delta u_2 - m_3 \frac{\partial^2 u_3}{\partial t^2} \delta u_3 \tag{40}$$

$$\delta U = k_2(u_2 - w_0)\delta(u_2 - w_0) + \int_{-\frac{L}{2}}^{\frac{L}{2}} EI \frac{\partial^4 w}{\partial x^4} dx \delta w$$

$$+ EI\left(\frac{\partial^2 w_1}{\partial x^2}\right)\delta\left(\frac{\partial w_1}{\partial x}\right) - EI\left(\frac{\partial^2 w_{-1}}{\partial x^2}\right)\delta\left(\frac{\partial w_{-1}}{\partial x}\right) + k_3(u_3 - u_2)\delta(u_3 - u_2)$$

$$- EI\left(\frac{\partial^3 w_1}{\partial x^3}\right)\delta w_1 + EI\left(\frac{\partial^3 w_{-1}}{\partial x^3}\right)\delta w_{-1} + EI\left(\frac{\partial^3 w_{0+}}{\partial x^3}\right)\delta w_0 + EI\left(\frac{\partial^3 w_{0-}}{\partial x^3}\right)\delta w_0 \tag{41}$$

$$\delta W_{nc} = EI \left( \frac{\partial^2 w_1}{\partial x^2} \right) \delta \left( \frac{\partial w_1}{\partial x} \right) - EI \left( \frac{\partial^2 w_{-1}}{\partial x^2} \right) \delta \left( \frac{\partial w_{-1}}{\partial x} \right)$$

$$- EI \left( \frac{\partial^3 w_1}{\partial x^3} \right) \delta w_1 + EI \left( \frac{\partial^3 w_{-1}}{\partial x^3} \right) \delta w_{-1} \tag{42}$$

where $\left( \frac{\partial w}{\partial x} \right) = \theta$. It is $EI \left( \frac{\partial^3 w_{0+}}{\partial x^3} \right) \neq EI \left( \frac{\partial^3 w_{0-}}{\partial x^3} \right)$ because the vibration absorber creates a concentrated shear force at $x=0$. Substituting (40) – (42) into (25) and separating the terms with $\delta w$, $\delta u_2$ and $\delta u_3$ the following three equations are obtained:

$$0 = -\rho A \frac{\partial^2 w}{\partial t^2} - EI \frac{\partial^4 w}{\partial x^4} + \left[ k_2 (u_2 - w_0) + EI \left( \frac{\partial^3 w_{0+}}{\partial x^3} \right) + EI \left( \frac{\partial^3 w_{0-}}{\partial x^3} \right) \right] \delta(x) \tag{43}$$

$$m_2 \frac{\partial^2 u_2}{\partial t^2} + k_2 (u_2 - w_0) + k_3 (u_2 - u_3) = 0 \tag{44}$$

$$m_3 \frac{\partial^2 u_3}{\partial t^2} + k_3 (u_3 - u_2) = 0 \tag{45}$$

where $\delta(x)$ is the Dirac function. Due to periodicity along $x$ direction, wave propagation through the infinite periodic beam can be expressed in a harmonic form:

$$w(x,t) = W exp(i\beta x - i\omega t), \qquad w_0(x,t) = W_0 exp(-i\omega t),$$

$$u_2(x,t) = U_2 exp(-i\omega t), \qquad u_3(x,t) = U_3 exp(-i\omega t) \tag{46}$$

where $i=\sqrt{-1}$ is the imaginary unit, $\beta$ is the wave number and $\omega$ is the vibration frequency. If the metamaterial is treated as a homogenized uniform beam the integration over the whole length (see Fig. 12b) gives:

$$0 = k_2 (u_2 - w_0) - \int_{-\frac{L}{2}}^{\frac{L}{2}} \rho A \frac{\partial^2 w}{\partial t^2} dx - EI \left( \frac{\partial^3 w_1}{\partial x^3} \right) + EI \left( \frac{\partial^3 w_{-1}}{\partial x^3} \right) \tag{47}$$

Using the harmonic wave solution and the assumption that the beam segment is a lumped with mass $\tilde{m}$ and spring rigidity $\tilde{k}$, we obtain the equation (32) with explanation (33). The values $\tilde{m}$ and a spring $\tilde{k}$ depend on the wavenumber $\beta$. Substituting (46) into (32), (44) and (45), and separating the terms with $W_0$, $U_2$ and $U_3$ the determinant of the system is formed. The solution of the system is non-trivial if the determinant is zero, i.e.,

$$\begin{bmatrix} -\tilde{m}\omega^2 + \tilde{k} + k_2 & -k_2 & 0 \\ -k_2 & k_2 + k_3 - \omega^2 m_2 & -k_3 \\ 0 & -k_3 & k_3 - \omega^2 m_3 \end{bmatrix} = 0 \tag{48}$$

Developing the determinant (48) the dispersion equation is obtained:

$$\left( -\tilde{m}\omega^2 + \tilde{k} + k_2 \right) \left[ \left( -\tilde{m}\omega^2 + \tilde{k} + k_2 \right) (k_3 - \omega^2 m_3) - k_3^2 \right] - k_2^2 (k_3 - \omega^2 m_3) = 0 \tag{49}$$

For the solution of the dispersion equation in which the wave number $\beta=i\alpha$ and $\alpha>0$, the function $w$ has the form (36). It is obvious that the amplitude decreases and tends to zero, i.e., the stop-band exists.

## 4.3 Multi-Flexural Band Gaps in an Euler-Bernoulli Beam with Lateral Local Resonators Transformation of the Flexural into Longitudinal Vibrations

In the previous text, single and double band gap metamaterials are shown, which are not suitable for the multi-frequency wave suppression. Further, in the configurations the local resonators are attached to continuum beams to generate band gaps for stopping the propagation of waves in one (longitudinal or lateral or torsional) direction. Wang et al. [15] suggested an acoustic metamaterial beam which is based on the multi-frequency vibration absorption and transformation of the flexural waves into longitudinal, i.e., the flexural vibration is attenuated into another direction in a beam. For theoretical consideration of the dynamic characteristics of the flexural wave propagation in an Euler-Bernoulli beam with lateral local resonators (LLR) the mathematical model is formed. Namely, it is of interest to suppress the flexural vibration in such a beam. The LLR structures substructures partially transform the flexural waves into longitudinal waves and block the wave propagation in another direction. In Fig. 13 a simple model of an Euler-Bernoulli beam with periodical LLR substructures in $x$ direction is plotted. One LLR consists of two lateral resonators with spring and mass constant of $k_2$ and $m_2$ and a vertical resonator with spring and mass constant $k_1$ and $m_1$ and a four link mechanism with rigid and massless trusses. The beam and the vertical resonator vibrate in $z$ direction and the lateral resonators vibrate in $x$ direction with displacement $w$, $u_1$ and $u_2$, respectively.
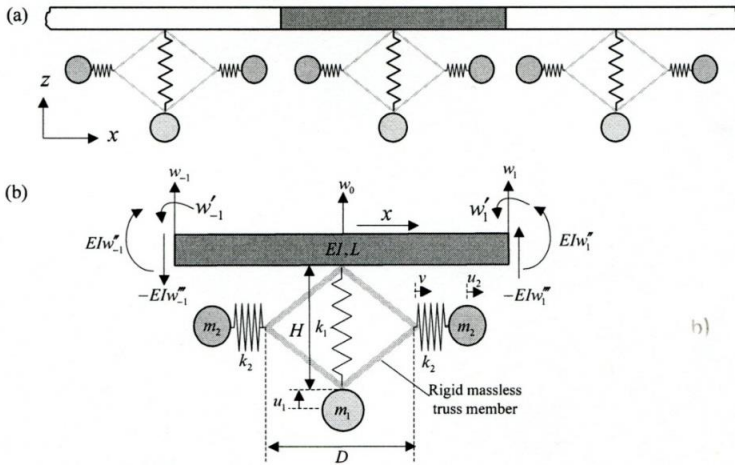


Figure 13
Construction of metamaterial beam: a) an infinite beam, b) a typical unit cell [15]

The unit cell of an infinite periodic metamaterial beam is shown in Fig. 13b. Inertial force in $z$ direction of the elementary beam is:

$$\rho A dx \frac{\partial^2 w}{\partial t^2}$$

Inertial forces of masses are: $m_1 \frac{\partial^2 u_1}{\partial t^2}$ and $m_2 \frac{\partial^2 u_2}{\partial t^2}$. Then, the elementary kinetic energy is:

$$\delta T = -m_1 \frac{\partial^2 u_1}{\partial t^2}\delta u_1 - 2m_2 \frac{\partial^2 u_2}{\partial t^2}\delta u_2 - \int_{-L/2}^{L/2}\rho A dx \frac{\partial^2 w}{\partial t^2}\delta w \tag{50}$$

The elastic force of the unit beam is $EI\frac{\partial^4 w}{\partial x^4}dx$ and the elastic forces in springs are $k_1(u_1 - w_0)$ and $k_2(u_2 - v)$ where $w_0$ is the flexural displacement of the center of beam and $v$ is the displacement of the truss end connected the lateral resonators. Based on the assumption of small displacements, we have:

$$v = -\frac{H}{2D}(w_0 - u_1) \tag{51}$$

Introducing the boundary conditions for the elementary unit, the elementary elastic energy of the system is:

$$\delta U = k_1(u_1 - w_0)\delta(u_1 - w_0) + \int_{-\frac{L}{2}}^{\frac{L}{2}}EI\frac{\partial^4 w}{\partial x^4}dx\delta w + 2k_2(u_2 + \frac{H}{2D}(w_0 -$$

$$u_1))\delta(u_2 + \frac{H}{2D}(w_0 - u_1)) + EI\left(\frac{\partial^2 w_1}{\partial x^2}\right)\delta\left(\frac{\partial w_1}{\partial x}\right) - EI\left(\frac{\partial^2 w_{-1}}{\partial x^2}\right)\delta\left(\frac{\partial w_{-1}}{\partial x}\right)$$

$$- EI\left(\frac{\partial^3 w_1}{\partial x^3}\right)\delta w_1 + EI\left(\frac{\partial^3 w_{-1}}{\partial x^3}\right)\delta w_{-1} + EI\left(\frac{\partial^3 w_{0+}}{\partial x^3}\right)\delta w_0 + EI\left(\frac{\partial^3 w_{0-}}{\partial x^3}\right)\delta w_0 \tag{52}$$

The elementary work of the non-conservative forces is:

$$\delta W_{nc} = EI\left(\frac{\partial^2 w_1}{\partial x^2}\right)\delta\left(\frac{\partial w_1}{\partial x}\right) - EI\left(\frac{\partial^2 w_{-1}}{\partial x^2}\right)\delta\left(\frac{\partial w_{-1}}{\partial x}\right)$$

$$- EI\left(\frac{\partial^3 w_1}{\partial x^3}\right)\delta w_1 + EI\left(\frac{\partial^3 w_{-1}}{\partial x^3}\right)\delta w_{-1} \tag{53}$$

Because of a concentrated shear force created by LLR structure at $x=0$, it is $EI\left(\frac{\partial^3 w_{0+}}{\partial x^3}\right) \neq EI\left(\frac{\partial^3 w_{0-}}{\partial x^3}\right)$. Using the extended Hamilton's principle (25) with (51)-(53) and separating the terms with $\delta w$, $\delta u_1$ and $\delta u_2$ the following equations are obtained:

$$0 = -\rho A \frac{\partial^2 w}{\partial t^2} - EI\frac{\partial^4 w}{\partial x^4} + \left[k_1(u_1 - w_0) - \frac{H}{D}k_2\left(u_2 + \frac{H}{2D}(w_0 - u_1)\right) +\right.$$

$$\left. EI\left(\frac{\partial^3 w_{0+}}{\partial x^3}\right)\delta w_0 + EI\left(\frac{\partial^3 w_{0-}}{\partial x^3}\right)\delta w_0\right]\delta(x) \tag{54}$$

$$-m_1 \frac{\partial^2 u_1}{\partial t^2} - k_1(u_1 - w_0) + \frac{H}{D}k_2\left(u_2 + \frac{H}{2D}(w_0 - u_1)\right) = 0 \tag{55}$$

$$m_2 \frac{\partial^2 u_2}{\partial t^2} + k_2\left(u_2 + \frac{H}{2D}(w_0 - u_1)\right) = 0 \tag{56}$$

Using the wave propagation function (46) and treating the system as a homogenized uniform metamaterial beam, the integration over the whole length gives:

$$0 = -\int_{-\frac{L}{2}}^{\frac{L}{2}} \rho A \frac{\partial^2 w}{\partial t^2} dx - EI\left(\frac{\partial^3 w_1}{\partial x^3}\right) + EI\left(\frac{\partial^3 w_{-1}}{\partial x^3}\right) + k_1(u_1 - w_0) - \frac{H}{D} k_2 \left(u_2 + \right.$$

$$\left. \frac{H}{2D}(w_0 - u_1)\right) \qquad (57)$$

$$\tilde{m} \frac{\partial^2 w_0}{\partial t^2} + \tilde{k} w_0 + k_1(u_1 - w_0) - \frac{H}{D} k_2 \left(u_2 + \frac{H}{2D}(w_0 - u_1)\right) = 0 \qquad (58)$$

where:

$$\tilde{m} = -\frac{2\rho A sin\left(\frac{\beta L}{2}\right)}{\beta}, \qquad \tilde{k} = -2EI\beta^3 sin\left(\frac{\beta L}{2}\right)$$

Combining (58) with (55) and (56) we obtain:

$$\begin{bmatrix} -\tilde{m}\omega^2 + \tilde{k} - k_1 - \frac{1}{2}\left(\frac{H}{D}\right)^2 k_2 & k_1 + \frac{1}{2}\left(\frac{H}{D}\right)^2 k_2 & -\frac{H}{D} k_2 \\ k_1 + \frac{1}{2}\left(\frac{H}{D}\right)^2 k_2 & m_1\omega^2 - k_1 - \frac{1}{2}\left(\frac{H}{D}\right)^2 k_2 & \frac{H}{D} k_2 \\ -\frac{H}{2D} k_2 & \frac{H}{2D} k_2 & m_2\omega^2 - k_2 \end{bmatrix} \times \begin{Bmatrix} W_0 \\ U_1 \\ U_2 \end{Bmatrix} = 0 \quad (59)$$

To obtain the non-trivial solution, the determinant of the coefficient matrix should be set to 0. If the solution of the determinant is a function of the wavenumber which is non-positive or even imaginary the wave propagation is stopped. So, the flexural waves are partially transformed into longitudinal waves and then totally blocked. It stimulates the lateral resonance to create inertial forces to counterbalance the shear forces resulting in wave suppression in the other directions. This type of beams is promising to be applied in the flexural absorber and isolator for vibration and noise control.

**Conclusions**

It can now be concluded:

1)  The theory of a conventional vibration absorber (spring - mass system) is suitable for explanation of the basic concept of acoustic metamaterial with negative effective (dynamic) mass.

2)  Acoustic metamaterial beams contain an isotropic beam with built-in spring-mass vibration absorbers. Depending on the way how the units are attached to the beam, the structure may absorb waves in one-direction (such as, for example, longitudinal waves) or waves in two directions (such as, for instance, transversal and longitudinal waves).

3)  Usually, acoustic metamaterial beams may produce one or two frequency gaps. At the moment, it is of interest to produce beams for multi-frequency gaps. The stop-bands are rather small. They need to be enlarged.

4)  We suggest introducing the construction of the acoustic metamaterial beams units which would have springs with nonlinear elastic properties instead of linear ones. Some investigation are already done (see [16-20]). The distinctive feature of this non-linear absorber would be that the vibration energy

transferred to the absorber, would be either localized or dissipated internally and would not re-enter the main system, even if the excitation of the main system is halted.

**Acknowledgement**

**References**

[1]     E. P. Calius, X. Bremaud, B. Smith, A. Hall, Negative mass sound shielding structures: Early results, Physica Status Solidi B 246 (9), 2009. pp. 2089-2097

[2]     H. Sun, X. Du, P. F. Pai, Theory of metamaterial beams for broadband vibration absorption, Journal of Intelligent Material Systems and Structures, 21, July, 2010, pp. 1085-1101

[3]     P. F. Pai, H.Peng, S. Jiang, Acoustic metamaterial beams based on multi-frequency vibration absorbers, International Journal of Mechanical Sciences 79, 2014, pp. 195-205

[4]     V. G. Veselago, The electrodynamics of sustances with simultaneously negative values of ε and μ, Sov. Phys. Usp 10, 1968, pp. 509-514

[5]     J. B. Pendry, Negative refraction makes a perfect lens, Phys. Rev. Lett. 85, 2000, pp. 3966-3969

[6]     J. Valentine, S. Zhang, T. Zentgraf, E. Ulin-Avila, D. A. Genov, G. Bartal, X. Zhang, Three-dimensional optical metamaterial with a negative refractive index, Nature 455, 2008, pp. 376-379

[7]     P. Sheng, X. X. Zhang, Z. Liu, C.T. Chan, Locally resonant sonic materials, Physica B 338, 2003, pp. 201-205

[8]     G. W. Milton, J. R.Willis, On modifications of Newton's second law and linear continuum elastodynamics, Proceedings of the Royal Society A 463, 2007, pp. 855-880

[9]     G. W. Milton, New metamaterials with macroscopic behavior outside that of continuum elastodynamics, New Journal of Physics 9, 359, 2007, pp. 1-13

[10]    H. H. Huang, C. T. Sun, G. I Huang, On the negative effective mass density in acoustic metamaterials, International Journal of Engineering Science 47, 2009, pp. 610-617

[11]    P. F. Pai, Metamaterial-based broadband elastic wave absorber, Journal of Intelligent Material Systems and Structures 21(5), 2010, pp. 517-528

[12]    P. Sheng, X. X. Zhang, Z. Liu, C. T. Chan, Locally resonant sonic materials, Physica B 338, 2003, pp. 201-205

[13]   P. Sheng, J. Mei, Z. Liu, W. Wen, Dynamics mass density and acoustic metamaterials, Physica B 394, 2007, pp. 256-261

[14]   R. Zhu, X. N. Liu, G. K. Hu, C. T. Sun, G. L. Huang, A chiral elastic metamaterial beam for broadband vibration suppression, Journal of Sound and Vibration 333, 2014, pp. 2759-2773

[15]   T. Wang, M. P. Sheng, Q. H. Qin, Multi-flexural band gaps in an Euler-Bernoulli beam with lateral local resonators, Physics Letters A, 380, 2016, pp. 525-529

[16]   P. F. Pai, B. Wen, A. S. Naser, M. J. Schulz, Structural vibration control using PZT patches and nonlinear phenomena, Journal of Sound and Vibration 215(2), 1998, pp. 273-296

[17]   P. F. Pai, M. J. Schulz, A refined nonlinear vibration absorber, International Journal of Mechanical Sciences 42, 2000, pp. 537-560

[18]   P. F. Pai, B. Rommel, M. J. Schulz, Nonlinear vibration absorbers using higher-order internal resonances, Journal of Sound and Vibration 234(5), 2000, pp. 799-817

[19]   L. Cveticanin, M. Kalami-Yazdi, H. Askari, Analytical approximations to the solutions for a generalized oscillator with strong nonlinear terms, Journal of Engineering Mathematics 77(1), 2012, pp. 211-223

[20]   M. M. Maaza, A. Arab, M. Belkhatir, S. Hammoudi, M. P. Luong, A. Benaissa, Non-linear behavior of sands under longitudinal resonance testing, Acta Polytechnica Hungarica 9(2), 2012, pp. 209-220

[21]   L. Cveticanin, Gy. Mester, I. Biro, Parameter influence on the harmonically excited Duffing oscillator, Acta Polytechnica Hungarica 11(5), 2014, pp. 145-160

# A New Class of Cascade Orthogonal Filters based on a Special Inner Product with Application in Modeling of Dynamical Systems

## Nikola B. Danković, Dragan S. Antić, Saša S. Nikolić, Staniša Lj. Perić, Marko T. Milojković

University of Niš, Faculty of Electronic Engineering, Department of Control Systems, Aleksandra Medvedeva 14, 18000 Niš, Serbia
nikola.dankovic@elfak.ni.ac.rs, dragan.antic@elfak.ni.ac.rs,
sasa.s.nikolic@elfak.ni.ac.rs, stanisa.peric@elfak.ni.ac.rs,
marko.milojkovic@elfak.ni.ac.rs

*Abstract: A class of cascade filters, orthogonal with respect to a new inner product, is presented in this paper. A sequence of generalized Malmquist orthogonal rational functions is used for design of these filters. In addition, by using these functions Müntz polynomials which are orthogonal in respect to a special inner product were derived. Obtained Müntz polynomials are applied in determination of outputs of the proposed filters. Depending on whether the design of the filters is performed in the s-domain or complex z-domain, we can derive a class of analogue or digital filters, respectively. Outputs from these filters are orthogonal with respect to the two different inner products. Both classes of filters are practically realized and their application in modeling of continuous-time and discrete-time dynamical systems is given. Obtained results show that there are great agreements between the outputs of models and real dynamical systems.*

*Keywords: orthogonal filters; inner product; dynamical systems modeling; Malmquist functions; Müntz polynomials*

## 1  Introduction

Around the middle of the last century, a new class of orthogonal rational functions were developed [1–3], and later they were used for synthesis of orthogonal filters. Also, classes of generalized Legendre polynomials and generalized orthogonal polynomials of Szegö were derived [4, 5]. Further generalization of classical orthogonal polynomials and orthogonal rational functions is given in [6, 7]. Derivation of Müntz orthogonal polynomials from these functions is described, as well in [8]. An overview of the theory of such orthogonal systems and some problems in applications of orthogonal polynomials are given in [9]. A few years later, the paper [9] was translated to English [10].

A method for obtaining a sequence of orthogonal rational functions is given in [11–14]. The method is based on the pole–zero and zero–pole mapping by using symmetric transformation. In this way obtained rational functions were used to design of new classes of orthogonal filters, quazi-orthogonal filters [15], almost orthogonal filters [16-21] and finally, generalized orthogonal filters with complex poles [22]. Using the same method, new classes of Malmquist orthogonal functions can also be obtained. These classes involve already known classes of Malmquist rational functions [23]. The method for obtaining Müntz orthogonal polynomials from sequence of orthogonal Malmquist functions is presented in [13, 14].

In this paper, generalized Malmquist functions are used for derivation of Müntz orthogonal polynomials (associated Müntz polynomials). By using these polynomials, a new class of filters, orthogonal with respect to a special inner product, is designed. Orthogonal rational function was derived by using the special symmetric transformation for pole–zero mapping and vice versa. Rational functions are used for structure design of filters, and appropriate Müntz polynomials are used to determine the outputs of these filters. Two subclasses of these filters are designed: subclass of analogue filters and subclass of digital filters. Practical realization of these filters was performed and their applications in modeling continual and discrete systems are given.

## 2    The Sequence of Orthogonal Rational Functions and Associated Müntz Orthogonal Polynomials

Let us introduce the sequence of rational functions:

$$W_n\left(s\right) = \frac{\prod\limits_{k=0}^{n-1}\left(s - \alpha_k^*\right)}{\prod\limits_{k=0}^{n}\left(s - \alpha_k\right)}, \quad n = 1, 2, 3, \dots \tag{1}$$

where the zeroes $\alpha_k^*$ are obtained by mapping the poles $\alpha_k$, and the poles $\alpha_k$ are obtained by mapping the zeroes $\alpha_k^*$, using the following symmetric transformation:

$$\alpha_k^* = f\left(\alpha_k\right), \ \alpha_k = f\left(\alpha_k^*\right). \tag{2}$$

Substituting (2) in (1), transformed sequence is obtained:

$$\overline{W_n(s)} = \frac{\prod\limits_{k=0}^{n-1}(s-\alpha_k)}{\prod\limits_{k=0}^{n}(s-\alpha_k^*)} \,. \tag{3}$$

As we can see, the zeroes of the transformed sequence are equal to the poles of $W_n(s)$ and vice versa.

Now, let us consider the inner product:

$$J_{n,m} = \left(W_n(s), \overline{W_m(s)}\right) = \frac{1}{2\pi i} \oint\limits_{C_P} W_n(s)\overline{W_m(s)}ds \,, \tag{4}$$

where $C_p$ involves all the poles of $W_n(s)$. In the cases of: $m > n$ and $m < n$, poles in the integrand (4) inside the contour $C_p$ are annulled with appropriate zeroes. Applying Cauchy theorem, we have $J_{n,m} = 0$. In the case of $m = n$, there exists one first-order pole inside the contour $C_p$, so $J_{n,n} = \oint\limits_{C_P} W_n(s)W_n(s)ds = N_n^2$. Hence, applying symmetric transformation of poles to zeroes of $W_n(s)$, we obtain the orthogonal sequence of the rational functions.

Let us apply symmetric transformation which maps the zeroes to poles and the poles to zeroes of $W_n(s)$ in the following way: $\alpha_k^* = \dfrac{b}{\overline{\alpha_k}}$, i.e., $\alpha_k = \dfrac{b}{\overline{\alpha_k^*}}$, to the sequence (1). In this way, we obtain two sequences orthogonal with respect to inner product (4):

$$W_n(s) = \frac{\prod\limits_{k=0}^{n-1}\left(s - \dfrac{b}{\overline{\alpha_k}}\right)}{\prod\limits_{k=0}^{n}(s-\alpha_k)}, \quad \overline{W_m(s)} = \frac{\prod\limits_{k=0}^{n-1}(s-\alpha_k)}{\prod\limits_{k=0}^{n}\left(s - \dfrac{b}{\overline{\alpha_k}}\right)}, \tag{5}$$

where $W_n(s)$ and $\overline{W_n(s)}$ are generalized Malmquist functions.

Using (4) and (5), we obtain:

$$\left(W_n(s)\overline{W_m(s)}\right) = \frac{b^n}{\left(|\alpha_n|^2 - b\right)\prod\limits_{k=1}^{n-1}|\alpha_k|^2}\delta_{n,m}\,. \tag{6}$$

One class of Müntz polynomials, orthogonal on $(0, 1)$, derived from the sequence of orthogonal rational function is given in [8, 11]. In the similar way, Müntz polynomials [9] can be obtained from the orthogonal sequence (5) as follows:

$$P_n(x) = \frac{1}{2\pi i} \oint_{C_P} W_n(s) x^s ds .$$
(7)

Using (5) and (7), we obtain:

$$P_n(x) = \sum_{k=0}^{n} A_{n,k} x^{\alpha_k} ,$$
(8)

where: $A_{n,k} = \dfrac{\prod_{j=0}^{n-1}\left(\alpha_k - \dfrac{b}{\alpha_j}\right)}{\prod_{j=0, j\neq k}^{n}\left(\alpha_k - \alpha_j\right)}, \ k = 0,1,2,...,n .$

A new operator for product over monoms $x^\alpha$ and $x^\beta$ is defined in the following manner [9]:

$$x^\alpha \ \square \ x^\beta = x^{\alpha\beta} .$$
(9)

Then, we define a new product of two Müntz polynomials $P_n(x) = \sum_{k=0}^{n} p_k x^{\alpha_k}$ and

$P_m(x) = \sum_{j=0}^{m} q_j x^{\alpha_j}$ in the following way:

$$\left(P_n \ \square \ P_m\right)(x) = \sum_{k=0}^{n}\sum_{j=0}^{m} p_k q_j x^{\alpha_k \alpha_j} .$$
(10)

On the basis of the product of polynomials defined above, the new inner product can be defined as:

$$\left(P_n, P_m\right)_\square = \int_0^1 \left(P_n \ \square \ \bar{P}_m\right)(x) \frac{dx}{x^2} .$$
(11)

Applying this inner product to Müntz polynomials (8), we obtain:

$$\left(P_n(x), P_m(x)\right)_\square = \frac{b^n}{\left(|\alpha_n|^2 - b\right)\prod_{k=0}^{n-1}|\alpha_k|^2} \delta_{n,m} .$$
(12)

In this way, we showed that Müntz polynomials obtained by using a class of generalized Malmquist functions are orthogonal with respect to the new inner product (12). A connection between orthogonal sequences of rational functions

(generalized Malmquist functions) and Müntz orthogonal polynomials is established (see Eq. 6).

# 3    Design of a New Class of Orthogonal Filters

When $W_n(s)$ have real poles, then associated Müntz polynomials have real exponents. In this case, substituting $x = e^t$ into $P_n(x)$, we obtain the sequence of exponential functions:

$$\varphi_n(t) = P_n(e^t) = \sum_{k=0}^{n} A_{n,k} e^{\alpha_k t} . \tag{13}$$

These functions are orthogonal with respect to the new inner product:

$$\left( \varphi_n(t), \varphi_m(t) \right)_\square = \int_0^\infty \varphi_n(t) \square \, \varphi_m(t) e^{-t} dt , \tag{14}$$

where: $\varphi_n(t) \square \; \varphi_m(t) = \sum_{k=0}^{n} \sum_{j=0}^{m} A_{n,k} A_{m,j} e^{\alpha_k \alpha_j t}$ .

The sequence of orthogonal rational functions (5) provides for the design of a new class of cascade orthogonal filters. Let us suppose that these functions have real poles. Now, we can design a cascade filter with real poles shown in Figure 1.



Figure 1

Block diagram of proposed orthogonal filter

Outputs from the filter in the time domain are:

$$\varphi_l(t) = L^{-1}\{W_l(s)\} = \frac{1}{2\pi i} \oint_{C_p} W_l(s) e^{st} ds , \tag{15}$$

where: $W_l(s) = \prod_{k=0}^{l} \dfrac{s - \dfrac{b}{\alpha_k}}{s - \alpha_k}$ , $l = 0,1,2,...,n$ .

Therefore, comparing (7), (8), and (15), we can notice that filter outputs are determined when we introduce substitution $x = e^t$ in Müntz polynomials (8),

obtained using the sequence of rational functions $W_l(s)$. Thereby, the contour $C_p$ involves all poles of $W_l(s)$, and all zeroes of $W_l(s)$ lie outside this contour.

Outputs $\varphi_l(t)$ are orthogonal in the time domain on the interval $(0, \infty)$ with respect to the new inner product (14). The diagram in Figure 1 is used for design of two new classes of orthogonal filters: analogue and digital version. The class of analogue orthogonal filters is obtained when we introduce *s* as an operator for differentiating. These filters are orthogonal in the complex domain on the contour which involves all poles of the filter. In the time domain, the filter outputs $\varphi_l(t)$ are orthogonal with respect to the inner product (14).

Digital orthogonal filters are obtained after moving in *z*-domain using the transformation $s = \dfrac{1}{T}\ln z$, where *z* presents the operator of prediction for one sample period in the time domain. These filters are also orthogonal in the complex domain, while outputs in the time domain are orthogonal in the classical sense.

# 4    Practical Realization of the New Orthogonal Filters

The obtained filter block scheme given in Figure 1 is suitable for practical implementation. For the purpose of analogue filter realization, we will write the transfer function in the following form:

$$W_n(s) = \frac{\displaystyle\prod_{k=0}^{n-1}\left(s + \alpha_k^*\right)}{\displaystyle\prod_{k=0}^{n}\left(s + \alpha_k\right)} = \frac{1}{s + \alpha_0}\prod_{k=1}^{n}\frac{s + \alpha_{k-1}^*}{s + \alpha_k}, \quad \alpha_k^* = \frac{b}{\alpha_k}, \quad \alpha_k \in R, \quad \alpha_k \geq 0. \tag{16}$$

In such a way the modified filter is orthogonal both in the complex and time domain. We have practically realized proposed filter in our Laboratory for modeling, simulation, and control systems and it is shown in Figure 2. The setup includes a printed circuit board with the realized new type of orthogonal filter, a microprocessor unit, an acquisition unit, and power supply. The realized filter has real and adjustable poles.
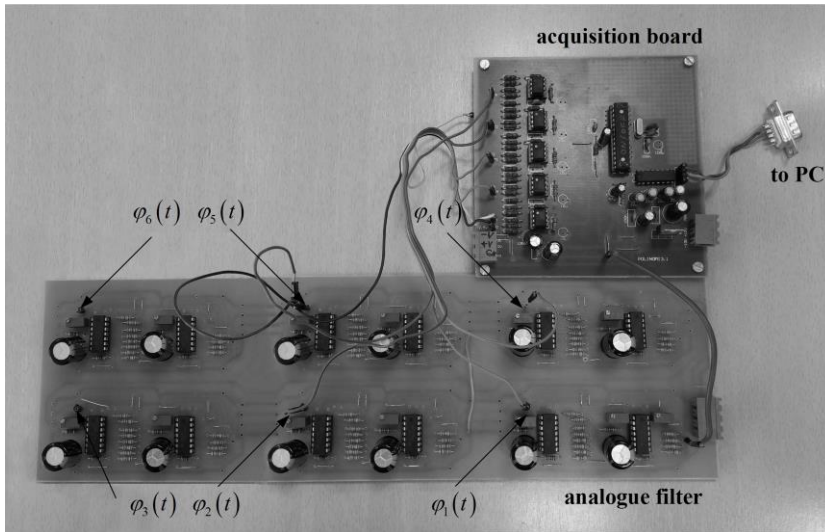
Figure 2
The analogue orthogonal filter, a printed circuit board

For the illustrative purpose, a sequence of functions on the outputs of the first five cascades of the proposed analogue filter (16) for the following values $b = -1$, $\alpha_0 = 1$, $\alpha_1 = 2$, $\alpha_2 = 3$, $\alpha_3 = 4$, $\alpha_4 = 5$ is obtained mathematically and given by:

$$\varphi_0(t) = e^{-t}, \ \varphi_1(t) = 3e^{-2t} - 2e^{-t},$$

$$\varphi_2(t) = 7e^{-3t} - \frac{15}{2}e^{-2t} + \frac{3}{2}e^{-t}, \ \varphi_3(t) = \frac{65}{4}e^{-4t} - \frac{70}{3}e^{-3t} + \frac{35}{4}e^{-2t} - \frac{2}{3}e^{-t}, \quad (17)$$

$$\varphi_4(t) = \frac{77}{2}e^{-5t} - \frac{1105}{16}e^{-4t} + \frac{455}{12}e^{-3t} - \frac{105}{16}e^{-2t} + \frac{5}{24}e^{-t}.$$

The outputs from the realized filter are shown in Figure 3. Orthogonality with respect to the new inner product (14) of these outputs can be easily verified.

As we have already said, when we get into $z$-domain (i.e., the operator $s$ is substituted with $z$), a digital filter can be realized on the basis of the modified scheme by using the following transfer function:

$$W_n(z) = \frac{z}{z - \alpha_0} \prod_{k=1}^{n} \frac{z - \alpha_{k-1}^*}{z - \alpha_k}, \ \alpha_k^* = \frac{b}{\alpha_k}, \ \alpha_k \in R \ . \quad (18)$$

Practical realization of the proposed digital filter, which is also performed in our laboratory, is shown in Figure 4. It too, has real and adjustable poles.

Figure 3
Outputs $\varphi_i(t)$ from the analogue filter, signals sensed on a printed circuit board



Figure 4
The digital orthogonal filter, a printed circuit board

A sequence of functions on the outputs of cascades of the proposed digital filter (18) is obtained mathematically for the following values $b = 1$, $\alpha_0 = \dfrac{1}{2}$, $\alpha_1 = \dfrac{1}{3}$, $\alpha_2 = \dfrac{1}{4}$, $\alpha_3 = \dfrac{1}{5}$, and the first few outputs are given by:

$$\varphi_0(K) = \left(\frac{1}{2}\right)^K, \quad \varphi_1(K) = 10\left(\frac{1}{3}\right)^K - 9\left(\frac{1}{2}\right)^K,$$

$$\varphi_2(K) = 231\left(\frac{1}{4}\right)^K - 320\left(\frac{1}{3}\right)^K + 90\left(\frac{1}{2}\right)^K, \tag{19}$$

$$\varphi_3(K) = 9576\left(\frac{1}{5}\right)^K - 17325\left(\frac{1}{4}\right)^K + 8800\left(\frac{1}{3}\right)^K - 1050\left(\frac{1}{2}\right)^K.$$

The same outputs are obtained by simulation in the Matlab/Simulink software package (Figure 5) and from the practically realized digital filter (Figure 6).

Figure 5

Outputs $\varphi_l(K)$, $l = 0, 1, 2, 3$ from the digital filter, signals sensed in Matlab

Figure 6

Outputs $\varphi_l(K)$, $l = 0, 1, 2, 3$ from the digital filter, signals sensed on a printed circuit board

Channel 1 in Figure 6 represents a step input signal, while channel 2 represents an appropriate signal from the filter section. Signals were recorded by using GW INSTEC digital storage oscilloscope (series GDS-3254).

The outputs from realized filter are orthogonal with respect to the inner product (14):

$$J_{n,m} = \left( \varphi_m(K), \varphi_n(K) \right) = \sum_{k=1}^{\infty} \varphi_m(K)\varphi_n(K) = N_n^2 \delta_{n,m} . \tag{20}$$

# 5 Case Study—Application in Modeling Continuous and Discrete Systems

Analogue and digital versions of the proposed orthogonal filters will be applied by modeling one continuous and one discrete system, both well known and commonly used in practice.

## 5.1 Modeling of a Protector Cooling System

The analogue version of the newly designed cascade orthogonal filter has been applied in modeling of one technological process in the tyre industry. It is a process of tyre strip cooling, more precisely protector (outer part of a tyre) cooling [15, 22, 24]. This is a complex electromechanical and thermodynamic system

which usually consists of 5–16 cascades about 15$m$ long. There are about several hundred systems like this in the world. The model is obtained by successively modeling cascades, starting with the first one. Modeling is performed applying a genetic algorithm. We used the adjustable model shown in Figure 7. Thereby, the output from $i$-th cascade in the model is described by:

$$y_{M,i} = \sum_{k=0}^{3} b_k \varphi_k(t), \quad i = 1, 2, ..., N,$$  (21)

where $N$ is an order number of the system cascade, and $b_j$ are summation coefficients.



Figure 7
Block diagram of an adjustable model with the proposed orthogonal analogue filter

For determining model of the first cascade, step response is used, while for the models of other cascades, previous ones are used (for modeling $i$-th cascade, ($i$-1)-th cascade is used).

The outputs from the first cascade of the process, $y_S(t)$, and the model, $y_M(t)$, are shown in Figure 8.

Using a genetic algorithm with minimization of the mean squared error $J = \dfrac{1}{N} \displaystyle\int_0^N \left( y_S(t) - y_M(t) \right)^2 dt$, we obtained poles of the process $\alpha_i$, summation coefficients $b_i$, and mapping parameter $b$. In our previously study, we have already used the genetic algorithm for the adjustment of parameters and minimization of criteria function $J$ [14, 17, 22].

The specific genetic algorithm used in experiments has the following parameters: initial population of 1000 individuals, a number of generations of 600, a stochastic uniform selection, a reproduction with ten elite individuals, and Gaussian mutation with shrinking. The used structure of chromosome was with 8 parameters coded by real numbers: $\alpha_0$, $\alpha_1$, $\alpha_2$, $\alpha_3$, $b_0$, $b_1$, $b_2$ and $b_3$. The main goal of the experiments was to obtain the best model of the unknown system in regard to the criteria function, i.e. mean squared error.

Figure 8

Outputs from the first cascade of a protector cooling system and the adjustable model

The model of the first cascade has the following form:

$$W(z) = \frac{b_0}{s + \alpha_0} + \sum_{k=1}^{3} b_k \frac{s + \dfrac{b}{\alpha_k}}{s + \alpha_k} \,, \tag{22}$$

where $b = 0.62$ and obtained values for poles and summation coefficients are given in Table 1. The obtained value for the mean squared error is: $J_{\min} = 0.77954 \cdot 10^{-3}$.

Table 1

Obtained parameters of the orthogonal analogue model

| k | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| $\alpha_k$ | 1.68 | 1.21 | 4.53 | 6.53 |
| $b_k$ | 2.08 | -0.22 | 3.35 | 3.35 |

The transfer function of the model (22) after substituting the values from Table 1 can be rewritten in the following way:

$$W(s) = \frac{8.57s^3 + 51.31s^2 + 103.12s + 76.38}{s^4 + 13.96s^3 + 63.65s^2 + 108.17s + 60.25} \,. \tag{23}$$

In order to verify quality of the model based on the proposed orthogonal filter, comparison with the model based on orthogonal Legendre's filter was performed (zeroes of the filter, $\alpha_k^*$ are shifted for $l$ related to poles $\alpha_k$ [21, 22]). Filter shown in Figure 7 is now modified according to the mapping function $\alpha_k^* = \alpha_k + l$, and

outputs from both models are given in Figure 8. The following results were obtained: $l$=0.7, $\alpha_0$= 5.82, $\alpha_1$=4.20, $\alpha_2$=0.54, $\alpha_3$=1.31, $b_0$=6.48, $b_1$=-0.38, $b_2$=0.18, $b_3$=-0.08, and $J_{min} = 5.3457 \cdot 10^{-3}$. Hence, in this case, mean squared error is much bigger.

## 5.2    Modeling of a Linear Part of DPCM System

On the other hand, a cascade orthogonal digital filter has been applied in modeling of the linear part of differential pulse code modulation transmission system [25]. Differential Pulse Code Modulation—DPCM is a well known and commonly used technique for signal transmission in telecommunications. An estimation, i.e. a prediction of the present value of the input signal is based on knowledge of its earlier values [26, 27]. This is why one of the most important parts of every DPCM and ADPCM (Adaptive Differential Pulse Code Modulation) is a predictor (a linear part of the system). The linear part in DPCM encoder will be modeled. In the encoder, the predictor is situated in the direct branch of a positive feedback loop as opposed to the decoder [27].

For the purpose of modeling a linear part of the encoder in DPCM system (in further text DPCM linear part), we use an adjustable model shown in Figure 9, which is based on the new orthogonal digital filter. In this particular case, we use filter with six sections and it has real and adjustable poles.

The output from the orthogonal model can be calculated as:

$$y_M(K) = \sum_{k=0}^{n} b_k \varphi_k(K), \tag{24}$$

where $K$ represents the number of samples.



Figure 9

Block diagram of an adjustable model with the proposed orthogonal digital filter

The desired model of DPCM linear part is obtained by adjusting the following parameters: orthogonal filter poles $\alpha_k$ ($k$=0,1,…,5), coefficients $b_k$ ($k$=0,1,…,5), and the mapping parameter $b$. In the case of modeling a particular unknown system, parameters of the model should be adjusted in such a way that the model

in Figure 9 corresponds as closely as possible to the unknown system. The process of modeling is performed in the well-known manner by introducing the same input to the system itself, as well as to its adjustable model based on the new cascade orthogonal digital filter (Figure 10) [15, 17, 18, 20].



Figure 10
The input of DPCM linear part and the adjustable model

The next step is measuring the outputs from the system $y_S(t)$ and the filter $y_M(t)$, and calculating the mean squared error (criteria function) as in the previous experiment: $J = \dfrac{1}{N}\sum_{K=0}^{N}\left(y_S(K) - y_M(K)\right)^2$. Unknown parameters are obtained by minimization of $J$.

The genetic algorithm used in the experiment has same values for initial population, the number of generations like in previous one, and reproduction with six elite individuals. Also, we used the stochastic uniform selection and Gaussian mutation with shrinking. The used structure of chromosome was with 12 parameters coded by real numbers: $\alpha_0, \alpha_1, \alpha_2, \alpha_3, \alpha_4, \alpha_5, b_0, b_1, b_2, b_3, b_4$ and $b_5$.

The original signal (output from DPCM linear part) and the signal obtained using the adjustable model based on the orthogonal digital filter are given in Figure 11.

Figure 11

Outputs from DPCM linear part and the adjustable model

The results obtained for the optimal values of the parameters of the adjustable model are presented in Table 2. The obtained value for the mean squared error is: $J_{min} = 13.182 \cdot 10^{-3}$.

Table 2

Obtained parameters of the orthogonal digital model

| k | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| $\alpha_k$ | 0.89880 | 0.80415 | 0.95560 | 0.71041 | -0.16135 | 0.86723 |
| $b_k$ | 0.80535 | 0.55137 | 0.30523 | 0.05137 | -0.10907 | 0.02849 |

We can notice a high level of matching between signals from DPCM linear part and the proposed orthogonal digital filter from the Figure 11.

Now, we have the model of DPCM linear part in the following form:

$$W_M(z) = \sum_{k=0}^{5} b_k \prod_{i=0}^{k} \frac{z - \alpha_{i-1}^{*}}{z - \alpha_i}, \ \alpha_{-1}^{*} = 0 \ , \tag{25}$$

where $\alpha_k^{*} = 0.82/\alpha_k$ and appropriate values of parameters are given in Table 2.

**Conclusion**

This paper presents a new class of cascade orthogonal filters based on the special inner product. A new method is applied in obtaining orthogonal Müntz polynomial from Malmquist rational functions and generalised Malmquist functions. Müntz polynomials obtained in this way are orthogonal with respect to the new inner product. Generalised Malmquist orthogonal functions were used for design of two new classes of orthogonal filters, for continuous (analogue) and discrete systems (digital filters). Müntz polynomials are used to prove orthogonality in the time domain, as well as to determine the outputs of these filters.

New filters are in the complex domain orthogonal on the contour which surrounds all the poles of the filters, while all zeroes lie outside this contour. Outputs of the filter are orthogonal with respect to the new inner product. Both analogue and digital filter were realized in our laboratory. Great matching between outputs from these filters and outputs obtained mathematically, by using Müntz polynomial, is shown.

The effectiveness of new classes of cascade orthogonal filters, analogue and digital, is demonstrated in the cases of determining a model of complex technological process in the tyre industry and for modeling the linear part of DPCM system, respectively. These filters can be applied in case of modeling dynamical systems when we adopt different mean squared errors (criteria functions) between the output of the process being modeled and the output from the adjustable filter (e.g. when the mean squared error is given in dB).

The class of these filters with complex poles (both analogue and digital) can be also a subject for consideration in some future works.

**Acknowledgement**

**References**

[1]    N. Akhiezer, *Theory of Approximation*, New York, Dover Publications, 1956

[2]    M. M. Djrbashian: Orthogonal Systems of Rational Functions on the Circle with Given Set of Poles, Dokl. Akad. Nauk SSSR, Vol. 147, pp. 1278-1281, 1962

[3]    M. M. Djrbashian: Orthogonal Systems of Rational Functions on the Circle, Akad. Nauk Armyan. SSR, Vol. 1, pp. 3-24, 106-125, 1996

[4]    G. Szegö, *Orthogonal Polynomials,* American Mathematical Society, Colloquium Publications, 23, Providence, 1975

[5]     G. V. Badalyan: Generalisation of Legendre Polynomials and Some of Their Applications, Akad. Nauk. Armyan. SSR Izv. Ser. Fiz.-Mat. Estest. Tekhn. Nauk, Vol. 8, No. 5, pp. 1-28, 1955

[6]     A. K. Taslakyan: Some Properties of Legendre Quasi-polynomials with Respect to a Müntz System, Mathematics, Erevan University, Erevan, Vol. 2, pp. 179-189, 1984

[7]     G. Mastroianni and G. Milovanović, *Interpolation Processes – Basic Theory and Applications*, Springer-Verlag, Berlin, Heidelberg, 2008

[8]     P. B. Borwein, T. Erdelyi, and J. Zhang: Müntz Systems and Orthogonal Müntz-Legendre Polynomials, Trans. Amer. Math. Soc. Vol. 342, No. 2, pp. 523-542, 1994

[9]     M. M. Djrbashian: A Survey on the Theory of Orthogonal Systems and Some Open Problems, In P. Nevai (Ed.), *Orthogonal Polynomials*, pp. 135-146, Springer, Netherlands, 1990

[10]    K. Müller and A. Bultheel: Translation of the Russian paper ˝Orthogonal Systems of Rational Functions on the Circle˝, Report TW253, Katholieke Univ. Leuven, Leuven, 1997

[11]    P. C. McCarthy, J. E. Sayre, and B. L. R. Shawyer: Generalized Legendre Polynomials, Journal of Mathematical Analysis and Applications, Vol. 177, No. 2, pp. 530-537, 1993

[12]    B. Danković, G. V. Milovanović, and S. Rančić: Malmquist and Müntz Orthogonal Systems and Applications, in *Inner Product Spaces and Applications*, T. M. Rassias, eds., Addison-Wesley Longman, Harlow, pp. 22-41, 1997

[13]    G. V. Milovanović, B. Danković, and S. Rančić: Some Müntz Orthogonal Systems, Journal of Computational and Applied Mathematics, Vol. 99, No. 1-2, pp. 299-310, 1998

[14]    S. B. Marinković, B. Danković, M. S. Stanković, and P. M. Rajković: Orthogonality of Some Sequences of the Rational Functions and the Müntz Polynomials, Journal of Computational and Applied Mathematics, Vol. 163, No. 2, pp. 419-427, 2004

[15]    D. Antić, Z. Jovanović, V. Nikolić, M. Milojković, S. Nikolić, and N. Danković: Modeling of Cascade-connected Systems using Quasi-orthogonal Functions, Electronics and Electrical Engineering, Vol. 18, No. 10, pp. 3-8, 2012

[16]    B. Danković, S. Nikolić, M. Milojković, and Z. Jovanović: A Class of Almost Orthogonal Filters, Journal of Circuits, Systems, and Computer, Vol. 18, No. 5, pp. 923-931, 2009

[17]    M. Milojković, D. Antić, M. Milovanović, S. S. Nikolić, S. Perić, and M. Almawlawe: Modeling of Dynamic Systems using Orthogonal Endocrine

Adaptive Neuro-fuzzy Inference Systems, Journal of Dynamic Systems, Measurement, and Control, Vol. 137, No. 9, pp. DS-15-1098, 2015

[18]  D. Antić, B. Danković, S. Nikolić, M. Milojković, and Z. Jovanović: Approximation Based on Orthogonal and Almost Orthogonal Functions, Journal of the Franklin Institute, Vol. 349, No. 1, pp. 323-336, 2012

[19]  M. Milojković, S. Nikolić, B. Danković, D. Antić, and Z. Jovanović: Modelling of Dynamical Systems based on Almost Orthogonal Polynomials, Mathematical and Computer Modelling of Dynamical Systems, Vol. 16, No. 2, pp. 133-144, 2010

[20]  D. Antić, S. Nikolić, M. Milojković, N. Danković, Z. Jovanović, and S. Perić: Sensitivity Analysis of Imperfect Systems using Almost Orthogonal Filters, Acta Polytechnica Hungarica, Vol. 8, No. 6, pp. 79-94, 2011

[21]  S. Nikolić, D. Antić, B. Danković, M. Milojković, Z. Jovanović, and S. Perić: Orthogonal Functions Applied in Antenna Positioning, Advances in Electrical and Computer Engineering, Vol. 10, No. 4, pp. 35-42, 2010

[22]  S. S. Nikolić, D. S. Antić, S. L. Perić, N. B. Danković, and M. T. Milojković: Design of Generalised Orthogonal Filters: Application to the Modelling of Dynamical Systems, International Journal of Electronics, Vol. 103, No. 2, pp. 269-280, 2016

[23]  P. Heuberger, P. Van den Hof, and B. Wahlberg, *Modelling and Identification with Rational Orthogonal Basis Functions*, Springer-Verlag, London, 2005

[24]  D. Trajković, V. Nikolić, D. Antić, and B. Danković: Analyzing, Modeling and Simulation of the Cascade Connected Transporters in Tire Industry using Signal and Bond Graphs, Machine Dynamics Problems, Vol. 29, No. 3, pp. 91-106, 2005

[25]  N. S. Jayant and P. Noll, *Digital Coding of Waveforms, Principles and Applications to Speech and Video, Chapter 6*, Prentice-Hall, Englewood Cliffs NJ, USA, 1984

[26]  N. B. Danković and Z. H. Perić: A Probability of Stability Estimation of DPCM System with the First Order Predictor, Facta Universitatis, Series: Automatic Control and Robotics, Vol. 12, No. 2, pp. 131-138, 2013

[27]  N. Danković, Z. Perić, D. Antić, D. Mitić, and M. Spasić: On the Sensitivity of the Recursive Filter with Arbitrary Order Predictor in DPCM System, Serbian Journal of Electrical Engineering, Vol. 11, No. 4, pp. 609-616, 2014

# The Special Characteristics of Stepping Motor Drives and a New Type of Classification

**László Számel, Tibor Vajsz**

Budapest University of Technology and Economics, Department of Electric
Power Engineering, Egry József utca 18, H-1111 Budapest, Hungary
E-mail: szamel.laszlo@vet.bme.hu, vajsz.tibor@vet.bme.hu

*Abstract: Stepping motor drives are widely used for positioning applications due to their easy controllability and straightforward connectivity to digital electronics. One of their greatest advantages is the possibility to perform positioning without requiring a closed-loop position control system. Stepping motor drives have special characteristics and therefore are considered as special types of electric drives. In this paper, these special characteristics are presented, along with a new type of classification. This classification is based on both the presented special characteristics (construction, etc.) and a newly derived equation called the fundamental equation of stepping motor drives. It is also shown in this paper that based on the new classification stepping motor drives can be divided into two subcategories: synchronous-type and asynchronous-type steppers. These subcategories are based on their similarities to traditional synchronous- and asynchronous (induction) motor drives. Also, the basic equations of the presented stepping motor drives will be derived from the fundamental equation of stepping motor drives. Therefore, this new type of classification is well applicable for both scientific analysis and educational purposes.*

*Keywords: stepping motor; stepper motor; step motor; special characteristic; classification; the fundamental equation of stepping motor drives; variable reluctance stepping motor; permanent magnet stepping motor; hybrid stepping motor*

# 1 Introduction

## 1.1 The Basic Characteristics of Stepping Motor Drives

Stepping motor drives have been a subject of investigation for a long time. This is due to the fact that stepping motor drives have a very wide variety of construction. They are frequently used in low-power applications. The number of phases is usually two, three, four, or five. The arc of a full step is between 0,72° and 15° in most cases. The number of steps per revolutions for a 0,72° full-step motor in half-stepping mode is 1000, which is in the lower range of the resolution of the digitalized-output incremental encoders. There are three types of stepping motors: variable reluctance, permanent magnet, and hybrid [1].

The magnetic field established in the stator can be unidirectional or bidirectional, depending on the type of the stepping motor used. In the case of variable reluctance stepping motors the application of unidirectional magnetic field is sufficient, while in the case of permanent magnet stepping motor drives the bidirectional magnetic field can be well-used for increasing the electromagnetic torque of the motor.

All of the stepping motors are used in an open-loop fashion, which means that current-control is independent of the position of the rotor and therefore the electromagnetic torque of the machine is not flat. This results in a reduced loadability. This is similar to the V/F controlled permanent magnet synchronous motor drive, where a load-torque step causes heavy speed-oscillations and therefore the loadability of the drive must be reduced in order to prevent falling out of synchronism.

Nowadays, other solutions are becoming more and more popular, where current-control of a stepping motor depends on the actual position of the rotor. Therefore, the loadability of the motor is improved [2], [3].

## 1.2 The Fundamental Equation of Stepping Motor Drives

An electric machine produces electromagnetic torque of non-zero mean-value if the stator- and the rotor magnetic fields are at standstill relative to each other [4]. This requirement is often called the **speed-requirement** or the **frequency-requirement** and is widely used in Hungary for the classification of conventional electric machines [5]. The speed-requirement can be expressed as follows:

$$\omega_{stator-field,stator} = \omega_{rotor-field,rotor} + \omega_{rotor} \tag{1}$$

Where:

$\omega_{stator-field,stator}$: the angular speed of the stator-field relative to the stator

$\omega_{rotor-field,rotor}$: the angular speed of the rotor-field relative to the rotor

$\omega_{rotor}$: the angular speed of the rotor

In the case of synchronous machines the rotor-field is fixed to the rotor and therefore:

$$\omega_{rotor-field,rotor} = 0$$

$$\omega_{stator-field,stator} = \omega_{rotor} \tag{2}$$

In the case of DC-machines the stator-field is fixed to the stator and therefore:

$$\omega_{stator-field,stator} = 0$$

$$\omega_{rotor-field,rotor} = -\omega_{rotor}. \tag{3}$$

In the case of asynchronous (induction) machines, neither is the stator-field fixed to the stator, nor is the rotor-field fixed to the rotor. Therefore, in (1) nothing is necessarily zero, which means that the case of asynchronous machines is the most general among the three basic machine types.

The classification of stepping motor drives is not obvious. This is because in the case of stepping motors each phase conducts occasionally and therefore a movement of the stator- and the rotor magnetic fields occur only when the conducting stator phase is changing. It must be mentioned that stepping motor drives are basically considered synchronous motor drives, because in steady-state the speed of the motor can be controlled by the switching frequency. However, some of the stepping motor drives are supplied by direct currents and this makes their classification difficult. Therefore, stepping motor drives are rather classified as special electric motor drives.

The following equation can be derived from (1) for stepping motor drives:

$$\Delta\alpha_{stator-field,stator} = \Delta\alpha_{rotor-field,rotor} + \Delta\alpha_{rotor} \tag{4}$$

Where:

$\Delta\alpha_{stator-field,stator}$: the change in the position of the stator-field relative to the stator for one step-period

$\Delta\alpha_{rotor-field,rotor}$: the change in the position of the rotor-field relative to the rotor for one step-period

$\Delta\alpha_{rotor}$: the change in the position of the rotor for one step-period

This is the **fundamental equation of stepping motor drives**, which will be used for analyzing the properties of all stepping motor drives and will be the basis for their new classification.

Similarly in the case of conventional machine drives, different types of stepping motor drives can be defined based on (4). For synchronous-type stepping motor drives the following equations hold:

$$\Delta\alpha_{rotor-field,rotor} = 0,$$
$$\Delta\alpha_{stator-field,stator} = \Delta\alpha_{rotor} \tag{5}$$

Based on (3) a DC-motor type stepping motor would mean that:

$$\Delta\alpha_{stator-field,stator} = 0$$
$$\Delta\alpha_{rotor-field,rotor} = -\Delta\alpha_{rotor} \tag{6}$$

However, there are no stepping motor drives that would satisfy this equation because in all of the stepping motor drives the stator-field is moving relative to the stator due to the switching between the stator-phases. Therefore, DC-motor type steppers do not exist. Although some of the stepping motor drives are supplied by DC-currents, but based on (6) they cannot be called DC-motor type steppers.

In the case of asynchronous-type stepping motor drives – similarly to the case of conventional asynchronous (induction) motor drives – in (4) nothing is zero, which means that the case of asynchronous-type steppers is the most general among stepping motor drives.

# 2 Variable Reluctance Stepping Motor Drives

## 2.1 Basic Types

Variable reluctance (VR) stepping motors have teeth both on the stator and the rotor [1]. They do not contain any permanent magnet. Only the stator has winding. The rotor is made of soft iron material. Figure 1 shows a variable reluctance stepping motor.



Figure 1
A variable reluctance stepping motor [6]

The principle of the operation is very simple. The excited stator-phase magnetizes the stator-teeth belonging to it. The magnetized stator-teeth attract the rotor-teeth closest to themselves. Because both north- and south-poles attract the iron, the direction of the stator magnetic field is irrelevant. Therefore, the application of a unidirectional magnetic field is sufficient.

Although stepping motors are basically considered synchronous-type motors, asynchronous-type and synchronous-type variable reluctance stepping motors can be distinguished. **Unlike in the literature and for the sake of simplicity, the stator- and the rotor-teeth will be represented by lines.**

Figure 2 shows a synchronous-type VR stepping motor. The motor has three phases, six stator-teeth and two rotor-teeth. In the case of VR-steppers (or more generally, in the case stepping motors utilizing unidirectional stator magnetic fields) the minimal number of phases required to construct a symmetrical machine is three.



Figure 2
A synchronous-type variable reluctance stepping motor [7]

There are an even number of teeth belonging to each phase (on Figure 2 two teeth belong to each phase) and thus, the resultant of the forces effecting the shaft will be zero. Also, this arrangement provides a magnetically symmetrical construction. Thus, the number of stator-teeth ($Z_1$) is as follows:

$$Z_1 = 2pm^*$$ (7)

Where:

$p$: the number of pole-pairs (per phases)

$m^*$: the number of phases


## 2.2 Synchronous-Type Variable Reluctance Stepping Motors

In the case of synchronous-type variable reluctance stepping motors the magnetic field of the rotor is fixed to the rotor. Therefore, for the synchronous-type VR-stepper on Figure 2, from equation (5):

$$\Delta\alpha_{stator-field,stator} = 60°, \Delta\alpha_{rotor-field,rotor} = 0°, \Delta\alpha_{rotor} = 60°$$

The number of steps per revolutions ($S$) is:

$$S = m^*Z_2$$ (8)

Where:

$Z_2$: the number of rotor-teeth

A general relationship can be derived from equation (4) if $\Delta\alpha_{rotor}$ is measured in revolutions:

$$\frac{1}{Z_1} = \frac{1}{S} \tag{9}$$

This means that:

$$Z_1 = S \tag{10}$$

Substituting equations (7) and (8) into equation (10) gives:

$$Z_2 = 2p \tag{11}$$

This means that the number of rotor-teeth is equal with the number of poles. It must be noted that synchronous-type VR-motors are generally not used as stepping motors. Instead, they are mainly applied as high-speed switched reluctance motors [7], [8], [9], [10].

## 2.3    Asynchronous-Type Variable Reluctance Stepping Motors

Figure 3 and Figure 4 show two asynchronous-type VR-steppers. The motor on Figure 3 has 6 teeth on the stator and 4 teeth on the rotor, therefore it can be called a 6/4 asynchronous-type VR-stepper. For similar reasons, the motor on Figure 4 can be called a 6/8 asynchronous-type VR-stepper.



Figure 3

A 6/4 asynchronous-type VR-stepper [7]

Figure 4
A 6/8 asynchronous-type VR-stepper [7]

For the 6/4 stepper, the components of equation (4) are:

$\Delta\alpha_{stator-field,stator} = 60°, \Delta\alpha_{rotor-field,rotor} = 90°, \Delta\alpha_{rotor} = -30°$

Similarly, for the 6/8 stepper:

$\Delta\alpha_{stator-field,stator} = 60°, \Delta\alpha_{rotor-field,rotor} = 45°, \Delta\alpha_{rotor} = 15°$

It can be concluded that the magnetic field of the rotor is moving relative to the rotor and therefore, these steppers are similar to the asynchronous (AC induction) motors.

Similarly to the case of synchronous-type VR-steppers, general relationships can be derived from equation (4):

$$\frac{1}{Z_1} = \frac{1}{Z_2} \pm \frac{1}{S} \tag{12}$$

After a few equivalent transformations we get:

$$S = \frac{Z_1 Z_2}{|Z_2 - Z_1|} \tag{13}$$

Substituting equations (7) and (8) into equation (12) we get:

$$Z_2 = Z_1 \pm 2p \tag{14}$$

It can be concluded from (14) that there are two rotor-configurations belonging to a given stator-configuration. In general, the rotor-configuration with more teeth is chosen in order to increase the number of steps per revolutions.

Like the asynchronous (AC induction) motors, the asynchronous-type VR-steppers have "slip", too. The "slip"-definition of an asynchronous-type VR-stepper can be derived from the slip-definition of an asynchronous (AC induction) motor.

Therefore:

$$"slip" = \frac{\omega_{stator-field,stator} - \omega_{rotor}}{\omega_{stator-field,stator}} = \frac{\omega_{rotor-field,rotor}}{\omega_{stator-field,stator}} = \frac{\frac{1}{Z_2}}{\frac{1}{Z_1}} = \frac{Z_1}{Z_2} \qquad (15)$$

The "slip" of the 6/4 stepper on Figure 3 is 1.5. This means that the rotor is rotating in the opposite direction of the stator-field and with a speed, which is half of that of the stator-field. Similarly, for the 6/8 stepper on Figure 4 the "slip" is 0.75. This means that the rotor is rotating in the same direction as the stator-field and with a speed, which is 25% of that of the stator-field.

Although the motors described above are asynchronous-type VR-steppers, these motors also have the basic characteristic of conventional synchronous-type motors, which means that the speed of these motors is independent of the load-torque and an increase in the load-torque will cause an increase in the so-called load angle. Also, because of the unidirectional stator magnetic field, there is no need for changing the direction of the current-flow and therefore, the supply is DC, like in the case of DC-motors. Thus, it can be concluded that asynchronous-type VR-steppers have the characteristics of all basic machine types: synchronous, asynchronous (AC induction), and DC.

## 2.4    Increasing the Number of Steps per Revolutions

There are several solutions for increasing the number of steps per revolutions ($S$). One of the simplest methods is multiplying the number of teeth. This means that the magnetic field of the rotor travels a distance of several rotor-teeth instead of a distance of only one rotor-tooth, if a switching from one stator-phase to another one takes place. Let us mark the rotor-teeth multiplication coefficient with ($k + l$). This coefficient is made of two components because the stator-teeth multiplication coefficient ($k$) is not necessarily the same as the rotor-teeth multiplication coefficient. This issue can be explained by a technological reason: an increase in the number of stator-teeth will cause a decrease in the size of the stator-slots, which makes the winding-process more difficult.

Similarly to the previous cases, general relationships can be derived from equation (4) ($Z_1$ marks the number of wound stator-teeth without stator-teeth multiplication):

$$\frac{1}{Z_1} = \frac{k+l}{Z_2} \pm \frac{1}{S} \qquad (16)$$

After a few equivalent transformations we get:

$$S = \frac{Z_1 Z_2}{|Z_2 - (k+l)Z_1|} \qquad (17)$$

Substituting equations (7) and (8) into equation (16) we get:

$$Z_2 = (k + l)Z_1 \pm 2p \qquad (18)$$

This means that there are two rotor-configurations belonging to a given stator-configuration. Figure 5 shows an example for teeth-multiplication. In practice, the teeth-multiplication coefficients can be higher.



Figure 5
Teeth-multiplication [7]

In this example $m^* = 3$, $p = 1$, $k = 2$, $l = 0$, $Z_1 = 6$, $Z_2 = 12 \pm 2$. On Figure 5 there are 10 rotor-teeth and according to (15) the "slip" is:

$$"slip" = \frac{(k+l)Z_1}{Z_2} = \frac{2*6}{10} = 1.2 \tag{19}$$

## 2.5    The Speed of Variable Reluctance Stepping Motors

In steady-state synchronous motors rotate on the same speed as the rotating magnetic field of stator. This speed is called the synchronous speed. It can be expressed as follows.

$$n_1 = \frac{60f_1}{p} \tag{20}$$

Where:

$n_1$: the synchronous speed

$f_1$: the frequency of the fundamental component of the stator voltages

$p$: the number of pole-pairs

In the case of variable reluctance stepping motors a similar expression applies:

$$n_1 = \frac{60f_0}{S} = \frac{60f_1}{Z_2} \tag{21}$$

Where:

$n_1$: the speed of the motor in steady-state

$f_0$: the stepping frequency

$f_1$: the switching frequency for each of the stator phases, $f_0 = m^* f_1$

$Z_2$: the number of rotor-teeth

It can be concluded from these equations that variable reluctance stepping motors are basically considered synchronous-type motors. The correspondent quantity of the number of pole-pairs is the number of rotor-teeth.

# 3    Permanent Magnet Stepping Motor Drives

## 3.1    Basic Construction

Figure 6 shows a permanent magnet (PM) stepping motor. This motor has teeth with windings on the stator and permanent magnet on the rotor.



Figure 6
A permanent magnet stepping motor [6]

Contrary to the case of variable reluctance stepping motors, the direction of the stator magnetic field is relevant, because the stator phases see constantly alternating magnetic poles before themselves. Therefore, bidirectional magnetic field is required in order to operate the motor.

It should be noted that VR-steppers are driven by only the so-called reluctance torque because the rotor is not excited and the motor is magnetically asymmetrical. The reluctance torque can be expressed as follows [8], [9], [10]:

$$M = \frac{1}{2} i^2 \frac{\partial L}{\partial \propto} \tag{22}$$

Where:

$i$: the stator phase-current

$L$: the self-inductance of the stator

$\propto$: the mechanical angle of the rotor

Equation (22) assumes that only one phase is energized and the motor is operated in the linear region (the iron cores are not saturated).

In the case of PM-steppers the motor is driven by both the excitation torque and the reluctance torque. The excitation torque is (if only one phase is excited) [11]:

$$M = c\psi_r i \sin \alpha_e \tag{23}$$

Where:

$c$: a constant

$\psi_r$: the amplitude of the rotor fluxvector

$\alpha_e$: the electrical angle between the stator fluxvector and the rotor fluxvector

Equation (23) suggests that there is a linear relationship between the stator current of one phase and the electromagnetic torque. If only one stator phase is excited with nominal current, then for a PM-stepper the excitation torque is significantly greater than the reluctance torque. This means that $M \sim i$ and therefore the application of microstepping is possible, whereas in the case of VR-steppers.

Permanent magnet stepping motors – like variable reluctance stepping motors – have two basic types: synchronous-type and asynchronous-type permanent magnet stepping motors.

## 3.2    Synchronous-Type Permanent Magnet Stepping Motors

Figure 6 shows a synchronous-type permanent magnet stepping motor. In order to find out the basic equations of synchronous-type permanent magnet stepping motors, the fundamental equation of stepping motor drives (equation (4)) must be utilized. For a synchronous-type stepping motor $\Delta\alpha_{rotor-field,rotor} = 0$ and $\Delta\alpha_{stator-field,stator} = \Delta\alpha_{rotor}$. This means that for the PM-stepper on Figure 6 $\Delta\alpha_{stator-field,stator} = \Delta\alpha_{rotor} = 90°$.

Like for synchronous-type VR-steppers, equation (9) can be derived from equation (4). For a synchronous-type permanent magnet stepping motor the following equations hold:

$$Z_1 = 2p_1 m^* \tag{24}$$

$$S = 2p_2 m^* \tag{25}$$

Where $p_1$ and $p_2$ are the number of stator- and rotor pole-pairs respectively. If we substitute equations (24) and (25) into equation (9), after equivalent transformations we get:

$$p_1 = p_2 \tag{26}$$

This means that for synchronous-type PM-steppers, the number of stator pole-pairs and rotor pole-pairs is equivalent, like in the case of traditional motor drives (AC induction motors, permanent magnet synchronous motors, etc.). However, this is not true for all of the permanent magnet stepping motor drives, as we shall see for asynchronous-type PM-steppers.

### 3.3    Asynchronous-Type Permanent Magnet Stepping Motors



Figure 7
An asynchronous-type PM-stepper [6]

Figure 7 shows an asynchronous-type permanent magnet stepping motor. Like in the case of synchronous-type PM-steppers the fundamental equation of stepping motor drives will be utilized (equation (4)). For the PM-stepper on Figure 7 the components of equation (4) are as follows: $\Delta\alpha_{stator-field,stator} = 90°$, $\Delta\alpha_{rotor-field,rotor} = 120°$, $\Delta\alpha_{rotor} = -30°$.

Like in the case of asynchronous-type VR-steppers, $\Delta\alpha_{rotor-field,rotor} \neq 0$. However, in the case of asynchronous-type VR-steppers this is true because the magnetic field of the rotor is moving relative to the rotor. The same is impossible in the case of asynchronous-type PM-steppers because the rotor-field is established by the permanent magnets fixed to the rotor. The question is then how it is possible that $\Delta\alpha_{rotor-field,rotor} \neq 0$. The answer is that during a switching between stator phases the attracted pole-pairs on the rotor are changing. This means that $\Delta\alpha_{rotor-field,rotor} \neq 0$ is true because the torque-forming pole-pairs on the rotor are changing.

A similar relationship to equation (12) can be derived from the fundamental equation of stepping motor drives. This is:

$$\frac{1}{Z_1} = \frac{1}{p_2} \pm \frac{1}{S} \tag{27}$$

Where $p_2$ is the number of rotor pole-pairs. After a few equivalent transformations we get:

$$S = \frac{Z_1 p_2}{|p_2 - Z_1|} \tag{28}$$

This equation is similar to equation (13). It can be concluded that the number of rotor-teeth in the case of VR-steppers has the same role as the number of rotor pole-pairs in the case of PM-steppers. If we take it into consideration that equation (25) applies for asynchronous-type PM-steppers as well, then by substituting equation (25) into equation (27), the following expression can be derived:

$$p_2 = Z_1 \pm p_1 \tag{29}$$

According to equation (29) the number of rotor pole-pairs differs from the number of stator pole-pairs in the case of asynchronous-type PM-steppers. This is a significant difference from conventional electric motor drives like AC induction motor drives, permanent magnet synchronous motor drives, etc.

The "slip"-definition of asynchronous-type PM-steppers is similar to that of asynchronous-type VR-steppers. Based on equation (15), the "slip" of an asynchronous-type PM-stepper is:

$$\text{"slip"} = \frac{Z_1}{p_2} \tag{30}$$

## 3.4   Comparison of VR-Stepper and PM-Stepper Drives

Variable reluctance stepping motors are cheaper than permanent magnet stepping motors. This is because they have a much simpler construction and they do not contain any permanent magnet. Also, VR-steppers are mechanically more robust than PM-steppers. While PM-steppers are exposed to the danger of demagnetization due to excessive currents, vibration, etc., VR-steppers do not suffer of this problem, because they do not contain any permanent magnet.

However, PM-steppers have a much higher power-density. A consequence of this is that for the same nominal power and for the same stator current PM-steppers produce more torque than VR-steppers.

In the case of VR-steppers the number of steps per revolutions ($S$) that can be maximally achieved is much higher than in the case of PM-steppers. In fact, the main problem with PM-steppers is that they have a relatively low $S$. This problem can be reduced by three different ways.

The first solution is the application of microstepping. This is applicable because $M \sim i$ in the case of PM-steppers. However, microstepping complicates the control of the motor [12]. The other problem with microstepping is that it does not necessarily increase the accuracy of positioning because the error in the load-angle caused by the load-torque can be significantly higher than the resolution of microstepping.

The second solution is the increasing of both the number of stator phases and the number of rotor pole-pairs. The former is very expensive, the latter is very limited. The third solution is the application of a hybrid stepping motor.

# 4   Hybrid Stepping Motor Drives

## 4.1   Basic Construction

Hybrid stepping motors combine the advantages of VR- and PM-stepper motor drives while eliminating their problems [1]. Figure 8 shows a hybrid stepping motor.

The stator is the same as in the case of VR-steppers, including teeth-multiplication. The rotor has a special construction, however. There are two cups which are made of soft-iron material mounted on the permanent magnet rotor. The permanent magnet has a magnetic field of axial-direction and therefore the two cups have opposite magnetic polarity. Also, the two cups have teeth.

According to Figure 9 there is a tooth offset between the two cups, which means a 180° offset in electrical degrees. This is necessary because the magnetic field of the stator has the same direction at both of the cups and without the 180° offset in electrical degrees (which is in fact a one-pole offset) the resultant torque would be zero.

It is an important characteristic of the motor that – in contrast to the case of permanent magnet synchronous motors and PM-steppers – the permanent magnet on the rotor is not exposed to the danger of demagnetization due to excessive currents. This is because the stator generates a magnetic field of radial-direction while the permanent magnet on the rotor has a magnetic field of axial-direction.

Figure 8
A hybrid stepping motor [6]



Figure 9
The offset between the two cups [6]

Another important characteristic is that the number of steps per revolutions ($S$) can be as high as in the case of VR-steppers because the achievable maximum of the number of rotor-teeth can be relatively high.

However, the torque-forming of a hybrid stepper is more similar to that of a PM-stepper: the excitation torque is dominant while the effect of the reluctance torque is less significant [13], [14]. Therefore, $M \sim i$ is true for hybrid steppers, too.

There are two types of hybrid stepping motors: synchronous-type and asynchronous-type. Synchronous-type hybrid stepping motors are not used in practice because the achievable maximum of the number of steps per revolutions is lower than in the case of asynchronous-type hybrid steppers. Therefore, only asynchronous-type hybrid steppers will be discussed.

## 4.2   Asynchronous-Type Hybrid Stepping Motors

There are two possibilities for the excitation of the stator windings of an asynchronous-type hybrid stepper: bipolar and unipolar. Bipolar excitation is used when establishing of a bidirectional magnetic field is necessary, otherwise unipolar excitation (which establishes a unidirectional magnetic field) is sufficient.

The relationships for the number of steps per revolutions ($S$) in the case of unipolar excitation are the same as in the case of VR-steppers including teeth-multiplication, notably equations (8), (17) and (18) are still valid. Figure 10 shows a simple hybrid stepper that is applicable for unipolar excitation. The stator winding requires only unidirectional magnetic field because the rotor-teeth always get aligned with the currently excited stator-teeth. However, in the case of hybrid steppers that are made for bipolar excitation, this is not true.

Figure 11 shows a bipolar hybrid stepper in two different positions. In the first position the required polarity of the lower stator-tooth is north because otherwise stability-problems would occur in standstill. In the second position, when the lower stator-tooth is aligned with a rotor-tooth, the required polarity of the lower stator-tooth is south, because in an aligned position the attraction between the opposite poles is required in order to hold the motor in a firm position.



Figure 10

A hybrid stepper for unipolar excitation

Figure 11
A hybrid stepper for bipolar excitation, a) first position, b) second position

This means that in the case of bipolar hybrid steppers both the attraction between the opposite magnetic poles and the repulsion between the identical magnetic poles are utilized for torque-production. This is a significant difference from conventional electric motor drives, where the attraction of the opposite magnetic poles is utilized only.

Let us derive the basic equations of bipolar hybrid steppers. The number of steps per revolutions is:

$$S = 2m^* Z_2 \tag{31}$$

If we compare this equation with equation (8) then we can conclude that $S$ has been doubled. This is because a bidirectional magnetic field can be established instead of a unidirectional one. Substituting equation (31) into the fundamental equation of stepping motor drives (equation (4)), we get:

$$\frac{1}{Z_1} = \frac{(k+l)}{Z_2} \pm \frac{1}{2m^* Z_2} \tag{32}$$

Where $(k + l)$ is the rotor-teeth multiplication coefficient. After a few equivalent transformations the following equation can be derived:

$$Z_2 = (k + l)Z_1 \pm p_1 \tag{33}$$

Where $p_1$ is the number of stator pole-pairs. Substituting this back to equation (32) and after a few equivalent transformations we get:

$$S = \frac{Z_1 Z_2}{|Z_2 - (k+l)Z_1|} = \frac{Z_1 Z_2}{p_1} \tag{34}$$

The "slip"-definition of an asynchronous-type hybrid stepping motor is the same as in the case of a VR-stepper including teeth-multiplication, notably equation (19) is true for asynchronous-type hybrid steppers as well.

It must be mentioned that there are hybrid stepping motors that are suitable for both unipolar and bipolar control. In this case the application of bipolar excitation is more advantageous because the achievable maximum torque can be doubled. Figure 12 shows a hybrid stepper that is suitable for both unipolar and bipolar excitation.



Figure 12

A hybrid stepper for both unipolar and bipolar excitation [15]

If this motor is considered a four-phase (unipolar) stepper (phases "a", "b", "c", "d" are separate phases) then $m^* = 4$, $p = 1$, $k = 2$, $l = 0$, $Z_1 = 8$ and by substituting into equation (18) we get $Z_2 = 16 \pm 2$. If the motor is consider a two-phase (bipolar) stepper (phases "a" and "c" together form one phase and the same applies for phases "b" and "d") then $m^* = 2$, $p = 2$, $k = 2$, $l = 0$, $Z_1 = 8$ and by substituting into equation (33) we get $Z_2 = 16 \pm 2$. This means that this motor is suitable for both unipolar and bipolar excitation.

## 5    The New Classification of Stepping Motor Drives

In the previous chapters, a new type of classification of variable reluctance, permanent magnet and hybrid stepping motor drives has been made. This new type of classification is based on the fundamental equation of stepping motor drives, the construction and the mode of excitation. Figure 13 summarizes this new type of classification.

As it can be seen on Figure 13, synchronous-type VR-steppers are used as switched reluctance motor (SRM) drives. In the case of switched reluctance motor drives teeth-multiplication is not applied, because it would make it more difficult to synchronize the stator currents to the rotor position. However, asynchronous-type VR-steppers are widely used as stepping motor drives. Teeth-multiplication can be applied in order to increase the number of steps per revolutions. In the case

of all VR-steppers, the excitation is unipolar, because only a unidirectional magnetic field is required in order to operate the motor.

Stepping motor drives

Variable reluctance (excitation: unipolar)　　Permanent magnet (excitation: bipolar)　　Hybrid

Synchronous type (used as SRM, no teeth-multiplication)　　Asynchronous type　　Synchronous type　　Asynchronous type　　Synchronous type (not used)　　Asynchronous type

With teeth-multiplication　　Without teeth-multiplication

Unipolar　　Bipolar　　Bipolar & Unipolar

Figure 13
The new classification of stepping motor drives

In the case of permanent magnet stepping motor drives, both synchronous-type and asynchronous-type steppers are used. The excitation is always bipolar, because a bidirectional magnetic field is required in order to operate the motor.

In the case of hybrid stepping motor drives only asynchronous-type steppers are used. The excitation can be unipolar or bipolar. There are certain constructions that are suitable for both unipolar and bipolar control.

**Conclusions**

This paper has presented the special characteristics of stepping motor drives, along with a new type of classification. In order to make this new type of classification the fundamental equation of stepping motor drives has been derived. Based on this equation synchronous-type and asynchronous-type stepping motor drives have been distinguished. The basic equations of both types have been derived in all practical cases for all basic construction types (variable reluctance, permanent magnet and hybrid). It is an interesting conclusion of this paper that asynchronous-type steppers are far more frequently used than synchronous-type steppers.

**References**

[1]     T. A. Khan, T. A. Taj, I. Ijaz: Hybrid Stepper Motor and its Controlling Techniques a Survey, Proceedings of the 2014 IEEE NW Russia Young Researchers in Electrical and Electronic Engineering Conference (ElConRusNW), St. Petersburg, 2014, IEEE, pp. 79-83

[2]     N. Dahm, M. Huebner, J. Becker: FPGA System-on-Chip Solution for a Field Oriented Hybrid Stepper Motor Control, 9[th] International Multi-Conference on Systems, Signals and Devices (SSD), Chemnitz, 2012, IEEE, pp. 1-6

[3] W. Kim, C. Yang, C. C. Chung: Design and Implementation of Simple Field-Oriented Control for Permanent Magnet Stepper Motors Without DQ Transformation, IEEE Transactions on Magnetics, Vol. 47, No. 10, October 2011, pp. 4231-4234

[4] G. Müller: Theorie elektrischer Maschinen, VCH Verlagsgesellschaft mbH, D-69451 Weinheim (Bundesrepublik Deutschland), 1995, pp. 30-33

[5] Gy. Retter: Villamosenergia-átalakítók I. (In English: Electric Energy Converters I.), Műszaki Könyvkiadó, Budapest, 1986, pp. 120-126.

[6] Microchip Technology Inc.: Microchip Webseminars: Introduction to Stepper Motors, Part 1: Types of Stepper Motors, 2007, pp. 12-25

[7] L. Szamel: The special characteristics of stepping motor drives, ENELKO, Alba Iulia, Romania, 2012, pp. 150-154

[8] J. Borka, K. Lupan, L. Szamel: Control Aspects of Switched Reluctance Motor Drives, Proceedings of the IEEE Inernational, Symposium on Industrial Electronics (ISIE'93), Budapest, 1993, pp. 1-3

[9] L. Szamel: Optimal Control of Transistor SRM Converters with Reduced Number of Switching Element, 12[th] International Power Electronics and Motion Control Conference (EPE-PEMC 2006), Portoroz, Slovenia, 2006, IEEE, pp. 1-4

[10] L. Szamel: Model Reference Adaptive Control of Ripple Reduced SRM Drives, Periodica Polytechnica: Electrical Engineering Vol. 46, No. 3-4, 2002, pp. 163-174

[11] S. Derammelaere, B. Vervisch, F. Verbelen, K. Stockman: Torque ripples in stepping motor driven systems, 17[th] European Conference on Power Electronics and Applications (EPE'15 ECCE-Europe), Geneva, 2015, IEEE, pp. 1-6

[12] B. S. Somesh, A. Mukherjee, S. Sen, P. Karmakar: Constant Current Control of Stepper Motor in Microstepping Mode using PIC16F877A, 2[nd] International Conference on Devices, Circuits and Systems (ICDCS), Combiatore, 2014, IEEE, pp. 1-4

[13] S. Derammelaere, C. Debruyne, F. De Belie, K. Stockman, L. Vandevelde: Load angle estimation for two-phase hybrid stepping motors, IET Electric Power Applications, 2014, pp. 257-266

[14] S. Derammelaere, B. Vervisch, J. Cottyn, B. Vanwalleghem, P. Cox, F. De Belie, K. Stockman, L. Vandevelde, G. Van Den Abeele: The Efficiency of Hybrid Stepping Motors: Analyzing the Impact of Control Algorithms, IEEE Industry Applications Magazine, 2014, pp. 50-60

[15] I. Schmidt, Gyné Vincze, K. Veszprémi: Villamos szervo- és robothajtások (In English: Electric servo- and robot drives), Műegyetemi Kiadó, Budapest, 2000, pp. 210-212

# A Type-2 Fuzzy-based Approach to the Minnesota Code

**Norbert Sram, Márta Takács**

Óbuda University
Bécsi út 96/b, H-1034 Budapest, Hungary
E-mail: sramm.norbert@phd.uni-obuda.hu, takacs.marta@nik.uni-obuda.hu

*Abstract: Cardiovascular diseases are still among the most common causes of death. Online and automated diagnostic solutions are considered as a possible remedy. In this paper, the authors present one such solution, which is based on the Minnesota Code and type-2 fuzzy for identifying cardiovascular diseases. The presented diagnostic system is case studied on various ECG data sets and compared to diagnostic results provided by physicians.*

*Keywords: fuzzy; type-2 fuzzy; interval fuzzy; Minnesota Code*

## 1 Introduction

Cardiovascular diseases are some of the most common causes of death. Based on the World Health Organization reports [1], about 30% of deaths are caused by either Ischemic Heart Disease or stroke. The prediction of sudden cardiac deaths is still a major concern and remains mostly unsolved [2, 3]. It is now well-established that classifications based on clinical circumstances can be misleading and often impossible because 40% of sudden deaths may be unwitnessed [4]. There are approaches that try to improve this situation through online automated and expert system based setups [20, 21].

### 1.1 Minnesota Code

The Minnesota Code [5] is used as a starting point for providing automated monitoring and diagnostic systems for predicting cardiovascular diseases. The Minnesota Code [5] is a classification system for the electrocardiogram that utilizes a defined set of measurement rules to assign specific numerical codes according to the severity of the ECG (Electrocardiography) findings. As for the definition, the Minnesota Code is a structured list of rules that examines certain characteristics of ECG waveforms. The input for the Minnesota Code diagnostic

system is a heartbeat cardiac cycle (ECG signal shown on Figure 1) and the related waveform annotations (P, Q, R, S, T) for all the 12 ECG leads (I, II, III, V1, V2, V3, V4, V5, V6, aVL, aVF, aVR). The diagnostic system studies the various aspects (e.g. length, amplitude) of the waveforms to produce a result. The analysis of the inputs is done with diagnostic rules. One such rule is shown on Figure 2, where the two aspects (Q/R amplitude, Q duration) of the ECG signal are compared to the optimal values (1/3, 0.03) defined by the Minnesota Code.



Figure 1
ECG waveform with the corresponding annotations

## Anterolateral site (leads I, aVL, V$_6$)

1-1-1  Q/R amplitude ratio $\geq 1/3$, plus Q duration $\geq 0.03$ sec in lead I or V$_6$.

Figure 2
Minnesota Code diagnostic rule definition

The outputs of the Minnesota diagnostic system are "true"/"false" statements for various diagnostic rules. The boolean values of various diagnostic rules are used to produce the diagnosis (a cardiovascular disease).

The Minnesota Code combines three major elements: a set of measurement rules, a classification system for reporting ECG findings and a set of exclusion rules. The relationship between the three major sets is vaguely defined. Although the Minnesota Code has some known weaknesses, the authors have been investigating various approaches in the past to improve those, such as incorporating fuzzy logic [5] and representing the decision rule set as Ontology [6]. In this paper, the authors apply an interval type-2 fuzzy logic based approach for handling uncertainties in the diagnostic system since it has been shown to be a solution in

other domains [16, 17, 19]. A case study is also performed to show that it is a viable approach in terms of automated diagnostic solutions.

# 1 Interval-based Fuzzy Sets

Compared to regular (type-1) fuzzy sets, interval type-2 fuzzy set can have intervals as a result of their membership value, which defines the constraints of the membership:

$$A: X \rightarrow E([0,1])$$

Membership functions defined in this way are called interval based fuzzy sets, which can be represented using two functions, an upper and a lower curve as shown in Figure 3.

Figure 3
An interval based fuzzy set

Using interval based fuzzy sets, an uncertainty factor can be paired to the members of the set [8]. The reliability of a system increases with the usage of interval based sets, while the precisions decreases. The usage of interval sets also has a significant impact on the computational complexity [9].

Interval based fuzzy sets can be further generalized if fuzzy sets are used instead as intervals. Each interval itself can be a fuzzy set and it can assign another fuzzy set to each member of the fuzzy set as the value of membership. Fuzzy sets defined in this way are 2nd order fuzzy sets, also called type-2 sets (see Figure 4). These 2nd order fuzzy sets can be further generalized to 3rd and higher order fuzzy sets. In the case of 3rd order fuzzy sets, the membership values are 2nd order fuzzy sets.

Figure 4

A general type-2 fuzzy set, where *FOU* stands for "field of uncertainty" [10]

As already stated, the interval type-2 fuzzy is a specialized case of the general type-2 fuzzy, where the 3rd dimension does not contain fuzzy based values. The interval can be described with an upper and lower bound as shown in Figure 5.



Figure 5

Interval Type-2 fuzzy function

The Type-2 fuzzy set seen in Figure 2 can be represented as:

$$f(x) = \begin{cases} f_{lm}(x) = trapmf(x, \qquad b, c, d, e) \\ f_{um}(x) = trapmf(x, a, c, d, f) \end{cases}$$

Where *f(x)* is the value of uncertainty in the interval of $[f_{lm}(x), f_{um}(x)]$ for the variable *x*. The representation of uncertainty in the following way means that the length of the interval is the indicator of uncertainty. The longer the interval, the more uncertain the value is.

## 2    Interval Type-2 Fuzzy Definitions of the Minnesota Code

The *if-then* rules are present in the case of type-2 fuzzy as well, however the antecedent and consequent sets are type-2 sets. Type-2 fuzzy logic can be applied when there is an uncertainty factor present. With this approach, crisp values that constitute the diagnostic logic are represented as type-2 fuzzy sets. Similar to the type-1 fuzzy logic, the type-2 fuzzy also has a fuzzification step, rules, an inference step and an output processor. The output processor contains both the type-reducer and a defuzzification step. Based on the type-2 fuzzy logic, the crisp values defined by the Minnesota diagnostic code are represented with an upper and lower membership function. A unified approach is applied for the fuzzification of all crisp values. In each case, a 10% tolerance zone is applied to the original value [6].          The tolerance zone is defined with two membership functions based on the Minnesota code crisp values. For modeling the Minnesota Code diagnostic rules with type-2 fuzzy, the authors used interval type-2 fuzzy sets, where, for a specific waveform input, the output is an interval defined by the lower and upper membership functions. The lower membership function of a type-2 definition represents the value that is 5% below the optimal value, while the upper membership function is above the optimal value by 5%. Thus, the 2 membership functions form a range, where the optimal value means no risk. As an input diverges from the optimal value in the defined tolerance zone, the length of the interval representing the uncertainty factor also increases.



Figure 6

The fuzzy type-2 interval definition of one possible Q waveform state

The interval type-2 fuzzy function shown in Figure 6 has the following definition:

$$f(x) = \begin{cases} f_{lm}, & smf(x, 0.035, 0.05) \\ f_{um}, & smf(x, 0.04, 0.05) \end{cases}$$

where the *smf* function is the type-1 S-shaped fuzzy function.

The interval type-2 fuzzy function shown in Figure 7 can be defined as:

$$f(x) = \begin{cases} f_{lm}, & trapmf(x, 0.03, 0.04, 0.05, 0.06) \\ f_{um}, & trapmf(x, 0.035, 0.04, 0.05, 0.55) \end{cases}$$

The used *trapmf* function corresponds to the *trapezoid* type-1 fuzzy membership function. The above definitions show that the type-1 fuzzy membership functions can be used to define the interval type-2 fuzzy functions, where $f_{lm}$ represents the lower membership function, while $f_{um}$ the upper membership function.



Figure 7
The Fuzzy Type-2 interval definition of one possible Q waveform state

## 2.1 Defuzzification

The defuzzification of a type-2 membership function does not produce a crisp value that can be used by the Minnesota code. In the case of an interval type-2 function, the result of the function is an interval, as show in Figure 8.

Figure 8
The result of an interval type-2 fuzzy function

In order to be able to execute the diagnostic steps outlined by the Minnesota Code, the intervals need to be reduced to a single truth value. This can be done by applying a type-reducer to the given interval. There are various approaches and type-reducers that one can choose from [18]. The chosen type-reduction method affects the diagnostic output. Another possibility is to introduce a new type-reduction algorithm, which is defined for the specific use.

In the case of a *MISO (Multi input single output)* inference system, the rules are of the type $if\ x_1\ is\ A_1\ and\ x_2\ is\ A_2, \dots x_n\ is\ A_n\ then\ y\ is\ B$ where the input is $x_1^*, \dots x_n^*$. The firing rate for an interval is $\left[f_{lm}^k(x_k^*), f_{um}^k(x_k^*)\right] (k = 1, n)$. As stated earlier, the computational complexity with the general type-2 reduction approaches is high. Because of this, the authors introduced the following 2 methods for the specific use case of the Minnesota Code. The following two approaches are used for aggregating and reducing the multi parameter diagnostic rules:

- **1st method:** In the first case, the type-reduction step is applied to produce a crisp value then produced crisp values aggregated.
- **2nd method:** The second method aggregates the result intervals of the Type-2 membership functions, after this step, a type-reduction is applied to produce a crisp value.

## 2.2    Reduced Value Aggregation (1st method)

The method of reduced value aggregation provides a value that can be used by the Minnesota code by applying type-reduction to the intervals and then aggregating the results of the type-reductions.

In the process of evaluating a waveform parameter, the degree of truth will be the average of the interval formed by the upper and lower membership functions. If the diagnostic rule requires the processing of multiple parameters (for example,

the Q waveform length and R amplitude), the diagnostic rules degree of truth is determined by processing each input separately and aggregating the produced averages. In this case, the type-reduction is performed for each parameter and then the aggregation applied to the partial results formed by the type-reduction. One of the characteristics of this method (the 1[st] method) is that a single parameter can have a significant impact on the evaluation of the diagnostic rule. This is especially the case, when the minimum operator is used for aggregation.

$$Dof = T(\frac{u_1 + l_1}{2}; \frac{u_2 + l_2}{2})$$

$$l_1 \qquad u_1 \qquad l_2 \qquad u_2$$

## 2.3    Type-Reduction of Aggregated Intervals (2[nd] method)

The given $\left[f_{lm}^k(x_k^*), f_{um}^k(x_k^*)\right] \in R$ intervals are aggregated to produce an uncertainty interval. The type-reduction of the produced uncertainty interval provides the firing rate.

For the given $\left[f_{lm}^k(x_k^*), f_{um}^k(x_k^*)\right] \in R$ intervals, the following relationships are possible:

$$l_1 < u_1 < l_2 < u_2 \qquad Dof = \frac{u_1 + l_2}{2}$$

$l_1 \qquad u_1 \qquad l_2 \qquad u_2$

$$l_1 < l_2 < u_1 < u_2 \qquad Dof = \frac{l_2 + u_1}{2}$$

$l_1 \qquad l_2 \quad u_1 \ u_2$

$$l_1 = l_2 < u_1 = u_2 \qquad Dof = \frac{u_1 + l_1}{2} = \frac{u_2 + l_2}{2}$$

$l_1 = l_2 \qquad u_1 = u_2$

$$l_1 < l_1 < u_2 < u_1 \qquad Dof = \frac{l_2 + u_2}{2}$$

$l_1 \ l_2 \qquad u_2 \ u_1$

The diagnostic results can be produced without the application of a step-by-step type-reduction. In the case of complex diagnostic rules, type-2 fuzzy sets are used. In this case the given intervals are aggregated until a single interval is produced. The aggregation of intervals happens based on the relationship between the two aggregated intervals. In case there is an overlap between the intervals representing

the diagnostic parameter, the overlapping interval will be the result of the aggregation. If there is no overlap, the interval formed by the distance between the two diagnostic parameter intervals will be the result of the aggregation. The distance is formed by the smaller intervals upper bound and the larger intervals lower bound. The type-reduction of the aggregation forms the result of the diagnostic rule.

## 2.4   Consistency Level

The outputs of the Minnesota diagnostic system are *"true"/"false"* statements for various diagnostic rules. The *boolean* values of various diagnostic rules are used to produce the diagnosis (a cardiovascular disease). Because of this, the degree of truth values produced from the type-2 fuzzy sets need to be converted to *boolean* values. The authors introduced a consistency level for performing the conversion. For a given membership function $A$, with a consistency level of $\alpha$, the output set $D$ consists of all samples from $A$, for which the predicate $A(x)$ is greater than $\alpha$ holds true. The elements of set $D$ are considered to have a *true* value when mapped to a *boolean* domain. The consistency level value $\alpha$ is an input property of the fuzzy based diagnostic system. By raising the value of $\alpha$, the strictness of the fuzzy based diagnostic system increases, by decreasing, it allows more uncertainty.

### 2.4.1   General Definition of the Consistency Level

Let $I1 \in [0,1]$ be the output of the reference decision system and $I2 \in [0,1]$ the output of the experimental decision system. The consistency level is the number $a \in [0,1]$ that is used in the following way to classify the results:

- If $a < |I2 - I1|$ than the results given by the two decision systems are inconsistent, meaning they differ in the outcome. In such a case a new decision method or decision system parameter tuning is advised to have an unambiguous outcome.
- If $a > |I2 - I1|$ that the two results are consistent and the outcome of the reference decision system is accepted.

The generalized consistency domains can be defined by a $[0,1] \times [0,1]$ domain in the following way (Figure 9):

Figure 9
Generalized consistency domains

## 2.5   Extension of the Minnesota Code ontology with Type-2 Definitions

Type-2 fuzzy definitions of various waveform parameters can be stored as ontology annotations coupled to the ontology concepts representing various waveform states. Possible solutions for adding fuzzy values to ontologies are outlined in the papers [11, 12] and previous work of the authors. The ontology annotations can contain both the type-1 and type-2 fuzzy value definitions using an XML structure that differentiates between the two. Figure 10 shows the relevant type-2 part of a definition, where type-1 membership function definitions are reused to define the interval type-2 set.

Figure 10

Interval type-2 definition of the Q waveform state

# 3 Case Study

For evaluating the effectiveness of various methodologies, the authors used three freely available datasets from PhysioNet [13]. These are the T-Wave Alternans Challenge Database (*TWA database*) [14], the PTB Diagnostic ECG Database (*PTB database*) [15] and St.-Petersburg Institute of Cardiological Technics 12-lead Arrhythmia Database (*Incart database*). These datasets contain ECG signal recordings for various anonymized patients along with the diagnosis produced by their physician. The datasets can serve as an input for evaluating the effectiveness of the proposed type-2 based diagnostic system. The evaluation of the samples found in these ECG databases has been conducted with the classical Minnesota Code expert system and the two type-2 based fuzzy solutions. The results of the classical expert system based approach are used as a reference. For all ECG recordings, the Q and QS rule groups of the Minnesota Code have been executed. During the evaluation, the results of a single sample are compared to the results produced by the mentioned methodologies. The authors evaluated the impact of the cutoff point on the diagnostic results of the fuzzy based solutions in the case of each sample. The values of the cutoff point were chosen from a predefined set of values in order to have defined points of comparison.

## 3.1 Results of the Incart Database

The evaluation of the Incart database produced 5265 entries for the Q and QS diagnostic rule groups. Of these entries, 28% back are "*true*" statements according to the classical approach. This means that in 28% of the cases the diagnostic rules fired since the defined criteria has been met. These 1484 (28%) entries are a sign of irregularities in the patient, a symptom of a cardiovascular disease.



Figure 11

Comparison of the "true" results for the Incart database

Figure 11 displays the number of diagnostic rules that fired ("*true*") for the fuzzy methodologies compared to the classic approach (baseline). The horizontal axis is the value of the consistency level. The vertical axis shows the number of diagnostic rules that fired for a specific consistency level. The figure shows the effect that the value of the consistency level has on the diagnostic output of the fuzzy based methodologies. For the evaluated fuzzy methodologies, the difference is negligible above the consistency level of 0.9. The comparison of *"false"* diagnostic rules is presented in Figure 12. As the images show, the difference between the 2 fuzzy based diagnostic methods shows significant decreases if the value of the consistency level is above 0.5. By further increasing the value of the consistency level, the difference continues to decrease, as well.

Figure 12
Comparison of the "false" results for the Incart database

## 3.2  Results of the TWA Database

The processing of the TWA database results in 4980 diagnostic rule entries. Of all entries, 30% are *"true"*, meaning a third of the diagnostic rules fired. This is the reference value produced by the classical, expert system based approach.

The behavior of the diagnostic methodologies for the TWA database is very similar to the one exhibited in the case of the Incart database. The difference between the results of the fuzzy methodologies decreases significantly at the cutoff value of 0.5. At higher cutoff values this difference becomes negligible as seen in Figure 13 and Figure 14.

Figure 13
Comparison of the "true" results for the TWA database



Figure 14
Comparison of the "false" results for the TWA database

## 3.3    Results of the PTB Database

The largest of the three used databases, the evaluation of the PTB database produces 37959 diagnostic rule entities. According to the classical system, 11507 (~30%) of these diagnostic rule entities are "true", hence, patients are showing the symptoms of some form of hearth diseases.

Figure 15
Comparison of the "true" results for the PTB database



Figure 16
Comparison of the "false" results for the PTB database

The behavior of the fuzzy diagnostic methodologies corresponds to the one observed in the case of the Incart and TWA databases. Figure 15 and Figure 16 underline this.

It can be deduced that the $2^{nd}$ method (type-reduction of aggregated intervals) provides a more graceful handling for scenarios where one of the inputs of the diagnostic rules heavily influences the output towards the *"false"* spectrum. In these cases, the usage of interval distances acts as a type of counteract against bias where 1 input would determine the output. The $2^{nd}$ method is more tolerant since the processing is done on the distance formed by the 2 intervals. This results in

coupling higher truth factors to the diagnostic rule outputs. With this approach, small deviations are less likely to influence the diagnostic output. Of course, the value of the consistency level still plays a significant role in the tolerance of the system. The 1$^{st}$ method (aggregated, type-reduced intervals) is closer to the original approach of the Minnesota Code, which evaluates the inputs, without those biasing one another in any way.

## 3.4    Analysis of Diagnostic Outputs

Apart from the ECG data, the diagnostic conclusions of the physicians are also available for the PTB database. This information can be used to compare the results of various diagnostic methodologies for each patient to the one provided by their physician. Using the Minnesota Code diagnostic rule results, a diagnostic outcome can be provided based on a predefined mapping table. A part of such a mapping table can be seen in Table 1. With the diagnostic rules that fired, one can match a corresponding hearth disease. For example, if the results of a patient state that the diagnostic rule *1-1-1* is true, that given patient has a case of "*Q wave Myocardial Infraction*". Using this, one can check the effectiveness of the present diagnostic system by comparing the diagnostic results provided by physicians.

Table 1
Prediction of hearth diseases based on the Minnesota Code rules

| ECG Categories Associated With Myocardial Infarction / Ischemia | | |
|---|---|---|
| *Definition and Description* | | *Minnesota Code* |
| *Q wave MI* | *Q wave MI; major Q waves with or without ST-T abnormalities* | *1-1-x* |
| | *Q wave MI; moderate Q waves with ST-T abnormalities* | *1-1-1 plus 4.1, 4.2* |
| *Isolated minor Q and ST-T abnormalities* | *Minor Q waves without ST-T abnormalities* | *1-3-x* |
| | *Minor ST-T abnormalities* | *4-3, 4-4, 5-3, 5-4* |

In this section, the authors introduce some of the more insightful diagnostic cases of the PTB database that were identified by using the algorithms described in [11].

The first case to analyze is *Patient 01*. According to the diagnostic description, the patient had *Myocardial infarction*. The Minnesota Code states that *Myocardial infarction* or one of its specific cases can be diagnosed if the diagnostic rules *1.1.x, 1.2.x or 1.3.x* are true (Table 1). The result of the *1.1.x* rules are the same for classical and fuzzy based methodologies, however the *1.2.x* and *1.3.x* rule group results differ.

Table 2

Diagnostic differences in the case of Patient 01

| Rule | Classic approach | Aggregated Type-Reduced Intervals | Aggregated intervals |
|------|------------------|-----------------------------------|----------------------|
| 1.2.1 | False | 0.7 | 0.8954 |
| 1.3.1 | False | 0.7 | 0.9 |

Table 2 shows the type-2 fuzzy based results for *Patient 01*. This provides additional information related to the diagnostic outcome. Besides identifying that the patient has a case of *Myocardial infarction,* it also provides a hint at a more specific case. Using the result of the fuzzy diagnostic solution, it can be inferred that *Patient 01* might have a case of *"Q wave Myocardial infarction"* or *"Minor Q waves abnormalities"*.

In the case of *Patient 03,* there are contradictions. The documented diagnostic outcome for the patient is *"Myocardial infarction"*. With the Minnesota Code, this diagnostic outcome holds, if one of the *1.3.x* rules fire. In the case of the rule *1.3.1,* there are differences between the classical approach and the fuzzy based one. As can be seen in Table 3, the classical approach does not produce the same diagnostic output as the physicians'. On the other hand, the fuzzy diagnostic results are in line with the diagnosis reached by the physicians with an additional associated risk factor.

Table 3

Diagnostic differences in the case of Patient 03

| Rule | Classic approach | Aggregated Type-Reduced Intervals | Aggregated intervals |
|------|------------------|-----------------------------------|----------------------|
| 1.3.1 | False | 0.65 | 0.95 |
| 1.3.1 | False | 0.7 | 0.9 |

In the case of *Patient 18,* there are numerous contradictions between the type-2 and the classical approach. From the available 227 processed samples, there is a difference in 39 cases. That is 17% difference for the processing of a single diagnostic rule group. Table 4 shows that most of the type-2 fuzzy diagnostic results differ from the classic one and can be considered borderline cases [11]. The last 2 highlighted rows of Table 4 are an example for the weakness of the classical methodology.

Table 4

Highlighted diagnostic rule results of Patient 18

| Rule | Classic approach | Aggregated Type-Reduced Intervals | Aggregated intervals |
|------|------------------|-----------------------------------|----------------------|
| 1.2.2 | False | 0.7 | 0.7 |
| 1.1.1 | False | 0.9 | 0.975 |
| 1.1.1 | False | 0.78 | 0.94 |

| 1.2.2 | False | 0.85 | 0.85 |
|-------|-------|------|------|
| 1.1.1 | False | 0.975 | 0.99 |
| 1.3.3 | False | 0.55 | 0.85 |
| 1.2.1 | False | 0.97 | 0.99 |
| 1.1.1 | False | 0.975 | 0.97 |
| 1.3.3 | False | 0.55 | 0.85 |
| 1.2.1 | False | 0.97 | 0.99 |
| 1.1.1 | False | 0.975 | 0.97 |
| 1.1.2 | False | 0.985 | 0.985 |
| 1.3.1 | False | 0.85 | 0.949 |
| **1.2.1** | **False** | **1.0** | **1.0** |
| **1.2.2** | **False** | **1.0** | **1.0** |

In order to examine the cause of the difference in the last 2 cases, one needs to inspect the measured values (inputs) that caused this behavior. In the case of the highlighted *rule 1.2.1,* the measured value for the Q and R amplitude division is 0.34, which meets the expectation of the diagnostic rule, which states that *Q/R amplitude ratio must be greater than or equal to 1/3*. The other parameter used by *1.2.1* is the Q waveform length. According to the rule definition, *the Q waveform length needs to be greater than or equal to 0.02s and less than 0.03s*. The measured value of the waveform length is 0.03s, which does not meet the requirements of the rule definition. The advantages of the fuzzy methodologies are shown in these kinds of borderline cases.

Similar to the other patients, *Patient 18* was also diagnosed with "*Myocardial infarction*". As stated earlier, in order to have the diagnostic conclusion of "*Myocardial infarction"* in the Minnesota Code, rule *1.3.x* must meet the requirements. However, in the case of the classical approach this is not met. With the fuzzy based approach and a reasonable consistency level (0.8-0.9), the same results are achieved as those provided by the physicians.

**Conclusions**

The fuzzy type-2 approach to the Minnesota Code algorithm provides the needed flexibility for the system and can act as guidance for the physicians in providing the diagnostic outcome for patients. Because of the hierarchical setup of the Minnesota code, various diagnostic decision trees might be affected by value of the consistency level. One of the advantages of the presented method is that with the modification of the consistency level, the potential diagnostic outcomes can be investigated, as well. The risk associated with a diagnostic result is inversely proportionate to the value of the consistency level.

The presented type-2 fuzzy based diagnostic method is an approach for improving the weakness of the original Minnesota code. Further improvements include the incorporation of decision trees based on the results of the various algorithms.

Predicting and using the right methodology for a given patient and diagnostic can be taken a step forward by including statistical steps.

## References

[1]     World Health Organization, "Annex Table 2: Deaths by Cause, Sex and Mortality Stratum in WHO Regions, Estimates for 2002", The world health report, 2004

[2]     Engelstein ED, Zipes DP., "Sudden Cardiac Death", The Heart, Arteries and Veins. New York, NY: McGraw-Hill; 1998:1081-1112.4

[3]     Myerburg RJ, Castellanos A., "Cardiac Arrest and Sudden Death.", Heart Disease: A Textbook of Cardiovascular Medicine. Philadelphia, Pa: WB Saunders; 1997:742-779

[4]     de Vreede Swagemakers JJM, Gorgels APM, Dubois-Arbouw WI, van Ree JW, Daemen MJAP, Houben LGE, Wellens HJJ. "Out-of-Hospital Cardiac Arrest in the 1990's: a Population-based Study in the Maastricht Area on Incidence, Characteristics and Survival", J Am Coll Cardiol. 1997; 30:1500-1505

[5]     Prineas, Ronald J., Crow, Richard S., Zhang, Zhu-ming, "The Minnesota Code Manual of Electrocardiographic Findings", ISBN: 978-1-84882-777-6

[6]     Sram, N., Takacs, M., "Minnesota Code: A Fuzzy Logic-based Approach", Proc. of the 11$^{th}$ International Symposium on Computational Intelligence and Informatics (CINTI), pp. 233-236, 2010, Budapest, Hungary, 2010

[7]     Sram, N., Takacs, M., "An Ontology Model-based Minnesota Code", Acta Polytechnica Hungarica Vol. 12, No. 4, 2015

[8]     Dongrui Wu, Jerry M. Mendel, „Uncertainty Measures for Interval Type-2 Fuzzy Sets", Elsevier Information Sciences, Volume 177, Issue 23, December 2007

[9]     Kóczy T. László, Tikk Domonkos, "Fuzzy rendszerek", Typotex, ISBN 963 9132 55 1, 2001

[10]    https://en.wikipedia.org/wiki/Type-2_fuzzy_sets_and_systems

[11]    Sram, N., Takacs, M.,"Analysis of Fuzzy Logic Assisted Evaluation of the Minnesota Code", Computational Cybernetics (ICCC), 2013 IEEE 9$^{th}$ International Conference, pp. 121-124

[12]    Fernando Bobillo, Umberto Straccia, "Fuzzy Ontology Representation using OWL 2", International Journal of Approximate Reasoning, 2011, pp. 1073-1094

[13]    Goldberger AL, Amaral LAN, Glass L, Hausdorff JM, Ivanov PCh, Mark RG, Mietus JE, Moody GB, Peng C-K, Stanley HE. PhysioBank, PhysioToolkit, and PhysioNet: Components of a New Research Resource

for Complex Physiologic Signals. Circulation 101(23):e215-e220 [Circulation Electronic Pages; http://circ.ahajournals.org/cgi/content/full/101/23/e215]; 2000 (June 13)

[14] Moody GB, „The PhysioNet / Computers in Cardiology Challenge 2008: T-Wave Alternans", Computers in Cardiology 35:505-508; 2008

[15] Bousseljot R, Kreiseler D, Schnabel, A. Nutzungm "EKG-Signaldatenbank CARDIODAT der PTB über das Internet", Biomedizinische Technik, Band 40, Ergänzungsband 1 (1995) S 317

[16] Nagy K., Takacs M.,"Type-2 Fuzzy Sets and SSAD as a Possible Application", Acta Polytechnica Hungarica Vol. 5, No. 1, 2008

[17] Nagy, K., Divéki, S., Odry, P., Sokola, M., Vujičić, V.,"A Stochastic Approach to Fuzzy Control", Acta Polytechnica Hungarica Vol. 9, No. 6, 2012

[18] J. M. Mendel, R. I. John, and F. Liu.,"Interval Type-2 Fuzzy Logic Systems Made Simple", Trans. Fuz Sys. 14, 6 (December 2006), 808-821

[19] Hao Ying, "General Interval Type-2 Mamdani Fuzzy Systems are Universal Approximators", Fuzzy Information Processing Society, 2008. NAFIPS 2008

[20] E. Tóth-Laufer, "Soft Computing-based Techniques in Real-Time Health Monitoring Systems", in Proc. of the International Engineering Symposium at Bánki, Efficiency, Safety and Security (IESB 2013), Budapest, Hungary, November 19, 2013

[21] J. Min Kang, T. Yoo, H. Chan Kim, "A Wrist-Worn Integrated Health Monitoring Instrument with Tele-Reporting Device for Telemedicine and Telecare", IEEE Transactions on Instrumentation and Measurement, Vol. 55, No. 5, October 2006, pp. 1655-1661, DOI: 10.1109/TIM.2006.881035

# A Method for Quantitative Comparison of 2D Skeletons

## Gábor Németh[1], György Kovács[2], Attila Fazekas[3] and Kálmán Palágyi[1]

[1]Department of Image Processing and Computer Graphics, University of Szeged
Árpád tér 2, 6720 Szeged, Hungary
{gnemeth, palagyi}@inf.u-szeged.hu

[2]Analytical Minds Ltd.
Árpád út 5, 4933 Beregsurány, Hungary
gykovacs@analyticalminds.hu

[3]Department of Computer Graphics and Image Processing,
University of Debrecen
Kassai u. 26, 4028 Debrecen, Hungary
attila.fazekas@inf.unideb.hu

*Abstract: Skeletons are widely used shape descriptors which summarize the general form of binary objects. There exist numerous skeletonization techniques that produce various skeleton-like features for the same object. Despite of the fact, that some researchers have made efforts to compare skeletons and evaluate skeletonization algorithms, we propose a new similarity measure that is based on the concept of normalized distance maps. In addition, a novel method for the quantitative comparison of skeletons is also presented. The reported method uses a high resolution dataset containing pairs of elongated objects and their expected skeletons. Our method is validated with the help of generalized morphological skeletons driven by neighborhood sequences. Based on the proposed method, we compared and ranked nineteen existing 2D thinning algorithms.*

*Keywords: skeleton; comparison of skeletons; generalized morphological skeleton; neighborhood sequences*

## 1 Introduction

Skeleton is a region-based shape descriptor which represents the general form of objects. It plays important role in various applications in image processing and pattern recognition. The skeleton of a 2D continuous object can be defined as the set of the centers of all maximal inscribed (open) disks [1]. A disk is maximal inscribed if it is included in the considered object, but it is not covered by any other inscribed disk.

Skeletonization means a process for producing an approximation to the skeleton of a discrete/digital object. There exist various skeletonization techniques that produce different skeleton-like features for the same object [2]. For example Németh and Palágyi presented 21 new algorithms in a single paper [3].

Some researchers have made efforts to compare skeletons and evaluate 2D skeletonization algorithms [4] [5] [6] [7]. They proposed some similarity measures between two skeletal sets that do not take the original elongated objects into account. The only exception is the measure of reconstructibility [5], but it may view numerous sets of points as "best" skeletons of an object. This is why we propose some new types of similarity measures that are based on normalized distance maps.

In this paper we propose a novel method, for the quantitative comparison of skeletons. The two key components of our method are a specific similarity measure for skeletons and the created gold standard image database containing 55 pairs of reference 2D images and their expected skeletons. The proposed method is validated with the help of generalized morphological skeletons driven by comparable neighborhood sequences. According to our experiments, the reported method can be used for evaluating arbitrary skeletonization algorithms.

Note that, our first attempt at this was published in a conference paper [8]. In that work the generalized morphological skeletons driven by neighborhood sequences were compared by using a small test database (containing just ten pairs of images) and we applied a similarity measure that ignore the original images.

The rest of the paper is organized as follows. Section 2 provides a method for creating a gold standard image database for comparison of skeletons. In Section 3, we propose some new similarity measures to give to the distance between two kinds of skeletons extracted from the same object. Section 4 reports the generalized morphological skeletons are combined with neighborhood sequences, furthermore we validate the proposed method with the help of generalized morphological skeletons driven by comparable neighborhood sequences. Section 5 compares 19 existing 2D thinning algorithms. Finally, we round off the paper with some concluding remarks.

## 2   Creation of Gold Standard Images

In this section a technique is reported for creation of gold standard images that are suitable for quantitative comparison of skeletons. It involves the following steps:

1) The selection of base images

2) The creation of reference skeletons

3) The generation of reference images

These steps will now be described in more detail.

## 2.1    Selection of Base Images

We collected 55 high resolution binary images of different shapes. Here the selected images are called base images. Note that there are some collections of binary images (say the Kimia 216 dataset), but they contains rather small silhouettes with several thin parts. Hence all skeletonization algorithms are obliged to produce similar results for those images.

## 2.2    Creation of Reference Skeletons

Skeletonization algorithms need to take the following requirements into account:

- Force the "skeleton" to retain the topology of the original image (i.e., skeletonization must be a topology-preserving reduction [9])
- Force the "skeleton" to be in its geometrically correct position (i.e., in the "center" of the object)
- Produce a minimal structure (i.e., the desired "width" of the "skeleton" is one point)

Our aim was to create the kind of reference skeletons from the base images that would meet these three conditions. In order to fulfill the first requirement, we extracted a topologically correct raw skeleton from each base image using the topology-preserving thinning algorithm $AK^2$ [10]. These raw skeletons may include some unwanted side branches. So as to satisfy the other two conditions, raw skeletons were corrected. This pruning process could be performed automatically [11], but we edited the 55 raw skeletons manually. As a result, our reference skeletons satisfy all of the three conditions listed above.

Note that method for generating reference skeletons from the base images is not significant, since any topologically correct skeletonization algorithms would do.

## 2.3    Generation of Reference Images

It must not be assumed that a reference skeleton is the expected skeleton of the corresponding base image. Hence we constructed reference images to replace base images.

The first step is to calculate (error free) Euclidean distance maps from the white (i.e., non-object) points of base images, where each element $p$ has a value that gives the Euclidean distance to the nearest white point [12] [13]. The Euclidean distance map is defined as follows:

$$DM_{wBI}(p) = \min_{q \in wBI} d_E(p, q), \tag{1}$$

where $wBI$ and $d_E(p, q)$ denote the set of white points in base image $BI$ and the Euclidean distance between points $p$ and $q$, respectively. Note that $DM_{wBI}$ is stored in an array of floating point numbers.

Figure 1
Creating a pair of reference skeleton and reference image. A 115×90 base image of an elephant (a); its
raw skeleton (b); Euclidean distance map calculated from the white points of the base image (c);
reference skeleton (d); the reference image (here we used the raw skeleton) (e); the difference image
(i.e., "base image" XOR "reference image") (f).

The set of object points *RI* in the reference image is generated as follows:

$$RI = \bigcup_{p \in RS} \Delta_E\big(p, DM_{wBI}(p)\big) \tag{2}$$

where *RS* is the set of skeletal points in the corresponding reference skeleton, $DM_{wBI}$ is the Euclidean distance map calculated from the white points in the base image *BI*, and $\Delta_E(p,r)$ denotes the "best" discrete approximation to the Euclidean disk of radius $r \in \mathbb{R}$ centred at point $p \in \mathbb{Z}^2$, that is,

$$\Delta_E(p,r) = \{q \mid q \in \mathbb{Z}^2, d_E(p,q) \leq r\} \tag{3}$$

In other words, the generated reference image *RI* is the union of disks that are centred at skeletal points in *RS* and the radii of these disks are determined by using the Euclidean distance map $DM_{wBI}$.

Figure 1 provides an illustrative example of creating a pair of reference skeleton and reference image. One may say that the procedure of reference skeleton construction introduces a strong bias. These reference skeletons are subjective indeed. That is why we do not consider reference skeletons as expected ones of the base images. Reference images paired with reference skeletons differ from the corresponding base images (see Figure 1f). We assumed that the reference skeleton *RS* is the expected discrete skeleton of the reference image *RI*. Note that it is not guaranteed, that for each $p \in RS$, disk $\Delta_E(p, DM_{wBI}(p))$, is a maximal inscribed one in *RI*, but we insist that reference skeletons satisfy all the three conditions for skeletonization methods.

All of the 55 pairs of reference images and reference skeletons, are available at
`https://www.inf.u-szeged.hu/~gnemeth/compskel/`

# 3 Similarity Measures for Skeletons

## 3.1 Existing Similarity Measures

If we have a gold standard (i.e., reference skeleton images associated with reference images of elongated objects), then measuring the goodness of skeletons produced by an algorithm seems to be a fairly simple task. Numerous measures have been proposed to define the similarity/distance between two sets of points [4] [5] [6].

Let us consider the frequently applied Hausdorff distance between two (arbitrary) sets of points $P$ and $Q$, which may be defined as follows [14]:

$$H(P,Q) = \max\{\ \max_{p \in P} \min_{q \in Q} d_E(p,q),\ \max_{q \in Q} \min_{p \in P} d_E(p,q)\ \},\ =$$
$$\max\{\ \max_{p \in P} DM_Q(p),\ \max_{q \in Q} DM_P(q)\ \} \tag{4}$$

where $DM_P(q)$ denotes the value of the Euclidean distance map calculated from the set of points $P$ at position $q$.

In our first attempt, we sought to make a comparison of the skeleton $S$ (extracted from a reference image) with the corresponding reference skeleton $RS$ using the similarity measure $H(S,RS)$, but it did not work. Just one salient point (like an endpoint of an unwanted line segment) in $S$ may determine $H(S,RS)$, hence it is not a fair assessment of a method.

Lee, Lam, and Suen [5] proposed a sophisticated similarity measure between two skeletons $P$ and $Q$ which is defined by the following Formula (5):

$$C(P,Q) = \left( \frac{1}{\#(P)} \sum_{p \in P} \frac{1}{DM_Q(p)^2 + 1} + \frac{1}{\#(Q)} \sum_{q \in Q} \frac{1}{DM_P(q)^2 + 1} \right) / 2 \tag{5}$$

where $\#(P)$ denotes the number of points in set $P$.

Similarly to the Hausdorff distance and others proposed in some studies [4] [5] [6], measure $C$ does not take into account the original (elongated) object. Hence we do not regard these similarity measures as acceptable for evaluating skeletons. Skeletons should be treated as special kinds of sets of points.

Lee, Lam, and Suen [5] proposed an additional measure that takes the original object into account. This measure of reconstructibility is defined by the formula

$$\alpha(S,I) = \frac{\#\left( \bigcup_{p \in S} \Delta_E(p, DM_{wI}(p)) \right)}{\#(I)} \tag{6}$$

where $S$ is a "skeleton" of object $I$. The measure takes values from the interval $[0,1]$, since $\bigcup_{p \in S} \Delta_E(p, DM_{wI}(p)) \subseteq I$. They say that: $\alpha(S,I)=1$ means that $S$ is identical to the "best" skeleton of image $I$. Unfortunately, this is not always the case. There is no guarantee that an Euclidean disk included in $I$ and centred at $p \in S$ with radius $\Delta_E(p, DM_{wI}(p))$ will be a maximal inscribed one. One can construct various sets of point $S \subseteq I$ such that $\bigcup_{p \in S} \Delta_E(p, DM_{wI}(p)) = I$. In

addition, it is not hard to see that $\alpha(I,I)=1$ (for any object *I*), but an elongated object may not be treated as the "best" skeleton of itself.

Note that Couprie and Bertrand [15] also proposed some measures (i.e., spuriousness factor, reconstruction error, thickness factor) between 3D curve-skeletons, however these measures can be calculated in a complex way furthermore do not consider the thickness of different parts of objects. In addition, their approach does not yield a fully automated method.

Sobieczki et. al. [16] [17] also investigated some similarity measures to compare 3D mesh-contraction-based curve-skeletonization algorithms. Unfortunately, they assumed mesh representation, hence their method cannot be applied for pixel-based images.

Some shape matching algorithms are based on skeletal graphs. Skeletal graphs are derived from 3D curve-skeletons or 2D centerlines in which endpoints and junction points represents the set of nodes/vertices, and there is an edge between two nodes if the corresponding pixels/voxels are connected by a skeletal path. These methods consider pruned skeletons (i.e., some unwanted branches are removed [11]), and they are based on some time consuming graph matching methods [18] [19] [20] [21] [22]. Unfortunately, the similarity measures that are used in graph matching methods are not skeleton-specific ones, they assume general sets of points, and do not take the original object into consideration.

Note that chamfer matching [23] [24] could also yield a similarity measure, where the query and the target contours are the two skeletons to be compared. Unfortunately, the original object would be also ignored, and chamfer distances are to approximate the Euclidean metric with integers (or rational numbers).

Despite the wealth of previously proposed similarity measures, we looked for new ones. It should be mentioned that in our previous paper [8], we applied five kinds of similarity measures which also ignored the original objects.

## 3.2   Distance Map Normalization

From the results of our experiments, we came to realize that not every skeletal point is equally important (i.e., positioning error of a certain size in a "thin" part is much more serious than the same error in a "thick" segment). This is why we propose a normalized distance map which is defined as:

$$\overline{DM}_{S,wI} = DM_S/(DM_S + DM_{wI}) \qquad\qquad (7)$$

where *S* is a set of skeletal points that is extracted from the image *I* by a skeletonization algorithm. (Note that "/" and "+" merely denote the point-by-point division and addition of two arrays of floating point numbers which have the same size, respectively.) Figure (2) shows normalized distance maps for five kinds of skeletons (see Section 4).

<div align="center">

$RI$          $\overline{DM}_{RS,wRI}$          $\overline{DM}_{SS(RI,<1>),wRI}$

$\overline{DM}_{SS(RI,<2>),wRI}$      $\overline{DM}_{SS(RI,<1,2>),wRI}$      $\overline{DM}_{SS(RI,\mathcal{A}_{opt}),wRI}$

</div>

Figure 2
A reference image RI (see Fig. 1e); the corresponding normalized distance maps of its reference skeleton RS (see Fig. 1d) and the four kinds of skeletons shown in Fig. 5

It can readily be seen that the following three properties hold:

- $0 \leq \overline{DM}_{S,wI}(p) \leq 1$ for each point $p \in wI$
- $\overline{DM}_{S,wI}(p) = 0$ if and only if $p \in S$
- $\overline{DM}_{S,wI}(p) = 1$ if and only if $p \in wI$

## 3.3 A New Similarity Measure Based on Normalized Distance Maps

Let us consider the following measure between a skeleton $S$ of image $I$ and a normalized distance map $\overline{DM}_{S,wI}$:

$$D_{avg}(S, \overline{DM}_{S,wI}) = \frac{1}{\#(S)} \sum_{p \in S} \overline{DM}_{S,wI}(p) \qquad (8)$$

We are now ready to introduce a new similarity measure that is recommended for comparing two skeletons:

$$AA_I(S_1, S_2) = \left( D_{avg}\left(S_1, \overline{DM}_{S_2,wI}\right) + D_{avg}\left(S_2, \overline{DM}_{S_1,wI}\right) \right)/2 \qquad (9)$$

where $S_1$ and $S_2$ are two skeletal sets of points that are extracted from the same image $I$ ($S_1, S_2 \subseteq I$). In addition the following three properties hold for the similarity measure $AA$ (for any $S1$, $S2$, and $I$):

- $0 \leq AA_I(S_1, S_2) \leq 1$
- $AA_I(S_1, S_2) = AA_I(S_2, S_1)$
- $AA_I(S_1, S_1) = 0$

We should stress here, that the smaller value means a better similarity between the two skeletons in question.

## 3.4 Goodness of Similarity Measures

Consider two skeletonization techniques $T_1$ and $T_2$ that produce skeletal sets of points $T_1(I)$ and $T_2(I)$ for image $I$. Suppose that it is known that $T_1$ is better than $T_2$ (i.e., $T_1$ can produce more reliable skeletons than $T_2$). Let $(RI,RS)$ be a pair of reference image and its reference skeleton.

We say that the similarity measure *SM* is *reasonable* for $(RI,RS)$ if

$$SM_{RI}(T_1(RI), RS) \leq SM_{RI}(T_2(RI), RS) \tag{10}$$

The purpose of our experiments was to show that the proposed similarity measures is reasonable. The question is: How to find such comparable skeletonization techniques $T_1$ and $T_2$?

# 4 Validation

## 4.1 Comparable Skeletons

In this section a new family of skeletons called *sequence skeletons* are introduced. These skeletons are not competitive with others produced by some existing skeletonization algorithms but we can validate our comparison method with the help of them.

Mathematical morphology, developed by Matheron and Serra [25], is a powerful tool for image processing and image analysis. Its operators can extract relevant topological and geometrical information from images by using structuring elements (i.e., geometric templates to probe some properties of interest) of various shapes and sizes. We use the fundamental concepts and notions of mathematical morphology as reviewed by Gonzalez and Woods [26].

### 4.1.1 Neighborhood Sequences and their Disks

The aim of this subsection is twofold. First, notions and results related to neighborhood sequences and the derived discrete distances will be reviewed in brief. Second, the disks corresponding to the neighborhood sequences will be formally expressed in terms of dilations (i.e., fundamental morphological operations) [26].

We will now present some basic notions and results concerning neighborhood sequences.

Let $n,m \in \mathbb{N}$ with $m \leq n$. Two points $p=(p_1, ..., p_n)$ and $q=(q_1, ..., q_n)$ in $\mathbb{Z}^n$ are said to be m-adjacent if both of the following conditions are satisfied:

- $|p_i - q_i| \leq 1$ $(i \in \{1,2,...,n\})$
- $\sum_{i=1}^{n} |p_i - q_i| \leq m$

Note that these relations are reflexive and symmetric. In the case of $n=2$ (i.e., the 2-dimensional orthogonal grid), 1– and 2–adjacency relations are often referred to as 4– and 8–adjacencies, respectively [9].

The sequence $\mathcal{A}=\langle A(1),A(2),...\rangle$ is called an $nD$–*neighborhood sequence* if $A(i)\in\{1, 2, ..., n\}$ for all $i\in\mathbb{N}$. If for some $t\in\mathbb{N}$, we have $A(i+t)=A(i)$ for all $i\in\mathbb{N}$, then the neighborhood sequence $\mathcal{A}$ is called periodic with period $t$. For simplicity, let $\mathcal{A}=\langle A(1),...,A(t)\rangle$ stand for a periodic neighborhood sequence having a period $t$.

Let $\mathcal{A}=\langle A(1),A(2), ...\rangle$ be an $nD$–neighborhood sequence. The sequence of points $\langle r_0,..., r_l\rangle$ $(r_j\in\mathbb{Z}^n, j\in\{0,..., l\})$ is an $\mathcal{A}$-*path* of length $l$ $(l \geq 0)$ from point $p$ to point $q$ if $p=r_0$, $q=r_l$, and $r_{j-1}$ and $r_j$ are $A(j)$-adjacent for all $j$ $(j\in\{1, ..., l\})$.

Let $d_\mathcal{A}(p, q)$ stand for the $\mathcal{A}$-*distance* between two points $p$ and $q$. It is defined as the length of the shortest $\mathcal{A}$-path(s) between $p$ and $q$.

As we are considering 2D binary images, we shall now examine 2D neighborhood sequences. According to the definitions above, 2D–neighborhood sequences may contain two kinds of elements, namely "1" and "2". Notice that distances $d_{\langle 1\rangle}$, $d_{\langle 2\rangle}$, and $d_{\langle 1,2\rangle}$ correspond to cityblock, chessboard, and octagonal distances, respectively [27]. It can readily be seen that there exist an infinite number of possible neighborhood sequences. The trick is to choose the neighborhood sequence which gives the best approximation to the Euclidean distance. The existence of the best approximating neighborhood sequence was proved in [28]. This non-periodic sequence is:

$$\mathcal{A}_{opt} = \langle 2,1,1,1,2,1,2,1,1,2,1,1,...\rangle \tag{11}$$

Let $d_E(p, q)$ represent the Euclidean distance between the two points $p$ and $q$ in $\mathbb{Z}^2$. A natural partial ordering relation "$\preccurlyeq$" can be defined for 2D–neighborhood sequences $\mathcal{A}_1$ and $\mathcal{A}_2$.

If $|d_{\mathcal{A}1}(p, q) - d_E(p, q)| \leq |d_{\mathcal{A}2}(p, q) - d_E(p, q)|$ (for any two points $p$ and $q$), then $\mathcal{A}_1 \preccurlyeq \mathcal{A}_2$ (i.e., $\mathcal{A}_1$ is better than $\mathcal{A}_2$).

The following conditions hold for the four neighborhood sequences under comparison [28]:

$$\mathcal{A}_{opt} \preccurlyeq \langle 1,2\rangle \preccurlyeq \langle 1\rangle$$
$$\mathcal{A}_{opt} \preccurlyeq \langle 1,2\rangle \preccurlyeq \langle 2\rangle \tag{12}$$

$$\Delta_{<1>} \qquad \Delta_{<2>} \qquad \Delta_{<1,2>}$$

Figure 3
Sample disks corresponding to the cityblock, chessboard, and octagonal distances of radii up to 4. The points denoted by a "0" are the origin and each point denoted by $r \leq k$ belongs to a disk of radius $k$.



$$<1> \qquad <2> \qquad <1,2> \qquad \mathcal{A}_{opt}$$

Figure 4
Approximations of the Euclidean disk of radius 96 (represented as a black circles) considering four neighborhood sequences <1>, <2>, <1,2>, and $\mathcal{A}_{opt}$.

$$\mathcal{A}_{opt} \preccurlyeq\; < 1,2 > \;\preccurlyeq\; < 1 >$$

$\mathcal{A}_{opt} \preccurlyeq\; < 1,2 > \;\preccurlyeq\; < 2 >$ Let us consider the discrete distance $d_{\mathcal{A}}$ based on the neighborhood sequence $\mathcal{A}$. The corresponding discrete disk of radius $k$ ($k=0,1,\dots$) centred at the origin $\mathcal{O}$ is defined by

$$\Delta_{\mathcal{A}}(p) = \{\, p \mid d_{\mathcal{A}}(\mathcal{O},p) \leq k \,\} \tag{13}$$

Figure 3 shows some discrete disks derived from the three periodic discrete distances $d_{<1>}$, $d_{<2>}$, and $d_{<1,2>}$.

It is well known that the neighborhood sequences <1> and <2> (which are composed of only one kind of adjacency relation) are diamond–shaped and square–shaped, respectively, and we can get various octagon–shaped discrete disks if both relations are combined.

In order to get discrete disks based on neighborhood sequences in terms of dilations, we will assign structuring elements to adjacency relations.

Let us consider the $m$-adjacency in $\mathbb{Z}^2$ ($m=1,2$). The structuring element $Y(m)$ for the $m$-adjacency is defined by:

$$Y(p) = \{\, p \mid p \in \mathbb{Z}^2 \text{ such that } p \text{ is } m\text{-adjacent to } \mathcal{O} \,\} \tag{14}$$

Since $m$-adjacency is a reflexive and symmetric relation over $\mathbb{Z}^2$, the structuring element $Y(m)$ contains the origin and it is symmetric, i.e., if $p=(p1,p2) \in Y(m)$, then $-p=(-p1,-p2) \in Y(m)$.

It can readily be seen that the discrete disk $\Delta_{\mathcal{A}}(k)$ can be expressed in terms of dilations (denoted by "$\oplus$" [26]) as follows

$$\Delta_{\mathcal{A}}(k) = \begin{cases} \{\mathcal{O}\} & \text{if } k = 0 \\ \Delta_E(k-1) \oplus Y\big(A(k)\big) & \text{otherwise} \end{cases}$$

$$= \Big(...\big(\{\mathcal{O}\} \oplus Y\big(A(1)\big)\big) \oplus ...\Big) \oplus Y\big(A(k)\big) \tag{15}$$

Approximations of Euclidean disks with the structuring elements (or discrete disks) derived from four kinds of neighborhood sequences are illustrated in Figure 4.

## 4.1.2 Generalized Morphological Skeletons Driven by Neighborhood Sequences

The skeleton of a discrete binary image can be characterized via morphological operations. In this section first the conventional morphological skeleton that just uses one structuring element will be reviewed. Then we will focus on the generalized morphological skeletons that are driven by neighborhood sequences.

The 2D morphological skeleton $S(X,Y)$ of a discrete set of points $X \subset \mathbb{Z}^2$ (i.e., object points in a 2D binary image) determined by a structuring element $Y$ consists of the centers of all maximal inscribed discrete disks of radius $k$ ($k=0,1,...$) [26]. With this approach, the structuring element $Y$ is assumed to be the unit disk (i.e., a disk of radius 1) and the discrete disk $Y^k$ of radius $k$ is derived from $Y$ by successive dilations:

$$Y^k = \begin{cases} \{\mathcal{O}\} & \text{if } k = 0, \\ Y^{k-1} \oplus Y & \text{otherwise} \end{cases}$$

$$= \underbrace{\Big(...\big(\big(\{\mathcal{O}\} \oplus Y\big) \oplus Y\big) \oplus ...\big) \oplus Y}_{k-\text{times}} \tag{16}$$

A point $p \in X$ is the center of a maximal inscribed discrete disk of radius $k$ ($k = 0,1,...$) if $p \in X \ominus Y^k$ and $p \notin (X \ominus Y^{k+1}) \oplus Y$, where "$\ominus$" denotes the erosion (i.e., a fundamental morphological operation that is dual to dilation) [26].

For this reason, the morphological skeleton $MS$ of a set $X$ determined by a structuring element $Y$ is defined by:

$$MS(X,Y) = \bigcup_{k=0}^{K} MS_k(X,Y)$$

$$= \bigcup_{k=0}^{K} (X \ominus Y^k) - [(X \ominus Y^{k+1}) \oplus Y] \tag{17}$$

where $K$ is the radius of the largest inscribed disk. In other words,

$$K = \max \ \{ \, k \mid X \ominus Y^k \neq \emptyset \, \} \tag{18}$$

According to the formulation defined by (17) and (18), the morphological skeleton is the union of the disjoint skeletal subsets, where $MS_k(X, Y)$ contains the centers of all maximal inscribed disks of radius $k$ ($k=0,1, ..., K$). An interesting property of

the morphological skeleton is that, a set *X*, can be exactly reconstructed from the *K+1* skeletal subsets:

$$X = \cup_{k=0}^{K} MS_k(X,Y) \oplus Y^k \tag{19}$$

The main limitation of a morphological skeleton is that its construction is based on "disks" of the form $Y^k$. If the chosen structuring element $Y=Y(m)$ (*m*=1,2) (see (14)), then the discrete disk $Y^k=\Delta_{<m>}(k)$ (see (15)). Discrete disks $Y^k=\Delta_{<1>}(k)$ and $Y^k=\Delta_{<2>}(k)$ do not give "good" approximations to the Euclidean disks (see Figure 4), hence we suspect that morphological skeletons are probably not "close" to the expected skeleton.

In order to reduce the shortcomings of the conventional morphological skeleton, Maragos proposed generalized morphological skeleton transforms that allows us to use varying structuring elements in different steps [29]. In his approach, the structuring element $Y^k$ (i.e., a discrete disk of radius *k*, see (16)) can be replaced by $\{O\} \oplus Y_1 \oplus . . . \oplus Y_k$, where $<Y_1, ..., Y_k>$ is the prefix of length *k* of an arbitrary sequence structuring elements ($k = 0, 1, …$).

These generalized morphological skeletons can be combined with neighborhood sequences by using the sequence of structuring elements $< Y(A(1)), Y(A(2)), ... >$ which are related to the neighborhood sequence $\mathcal{A} =< A(1), A(2), ... >$.

This *sequence skeleton* makes use of discrete disks $\Delta_{\mathcal{A}}(k)$ ($k = 0, 1, …$) (see (15)).

The sequence skeleton *SS* of a $X \subseteq \mathbb{Z}^2$ driven by a neighborhood sequence $\mathcal{A}$ is defined by

$$SS(X,\mathcal{A}) = \cup_{k=0}^{K} SS_k(X,\mathcal{A}) \tag{20}$$

where

$$SS_k(X,\mathcal{A}) = \left( X \ominus \Delta_{\mathcal{A}}(k) \right) - \left[ \left( X \ominus \Delta_{\mathcal{A}}(k+1) \right) \oplus Y\left( A(k+1) \right) \right] \tag{21}$$

and *K* is the radius of the largest inscribed disk; that is

$$K = \max \ \{ k \mid X \ominus \Delta_{\mathcal{A}}(k) \neq \emptyset \} \tag{22}$$

With this formulation defined by (20) and (22), the sequence skeleton is the union of disjoint skeletal subsets, where $SS_k(X,\mathcal{A})$ contains the centers of all maximal inscribed disks $\Delta_{\mathcal{A}}(k)$ ($k$=0, 1, . . . ,$K$).

It is easy to see that in the $\mathcal{A} = <m>$ case (*m*=1,2)

$$SS(X,\mathcal{A}) = MS\left(X, Y(m)\right) \tag{23}$$

thus the conventional morphological skeleton is a special case of sequence skeletons. Some illustrative examples of sequence skeletons are given in Fig. 5.

$$X \qquad\qquad SS(X,<1>) = MS\big(X,Y(1)\big) \qquad SS(X,<2>) = MS\big(X,Y(2)\big)$$

$$SS(X,<1,2>) \qquad\qquad SS\big(X,\mathcal{A}_{\text{opt}}\big)$$

Figure 5
A $115 \times 90$ image of an elephant and its sequence skeletons driven by four kinds of neighborhood sequences. Notice that the first two sequence skeletons are also conventional morphological skeletons.

The reconstruction formula from the sequence skeleton $SS(X,\mathcal{A})$ is analogous to (19), hence:

$$X = \bigcup_{k=0}^{K} SS_k(X,\mathcal{A}) \oplus \Delta_{\mathcal{A}}(k) \,. \tag{24}$$

This means that like the conventional morphological skeletal subsets, the subsets of the sequence skeleton also fully represent the original set of points [29] [30].

Next, note that the connectivity of the conventional morphological skeletons and sequence skeleton is not guaranteed (i.e., these skeletons are not connected and topologically correct for numerous connected objects).

It is known that the non–periodic neighborhood sequence $\mathcal{A}_{\text{opt}}$ (see (11)) provides the best approximation to the Euclidean distance and the Euclidean disk [28]. Hence we can assume that $SS(X, \mathcal{A}_{\text{opt}})$ is the best sequence skeleton for any $X$ (i.e., it gives the best approximation to the expected skeleton).

## 4.2   Validation with Sequence Skeletons

In this section we validate the proposed method for the quantitative comparison of skeletons with the help of neighborhood sequences.

We examined our gold standard image database, containing 55 pairs of reference images and reference skeletons, the four suggested similarity measure $AA$, and the four metrical neighborhood sequences <1>, <2>, <1, 2>, and $\mathcal{A}_{\text{opt}}$ (see (11)).

For each pair of ($RI$, $RS$) we calculated the followings:

- The four sequence skeletons driven by the four neighborhood sequences in question:
  $S1 = SS(RI, <1>)$

$S2 = SS(RI, <2>)$

$S3 = SS(RI, <1, 2>)$

$S4 = SS(RI, \mathcal{A}_{opt})$ (see (20))

- The five normalized distance maps corresponding to the reference skeleton and the four sequence skeletons:

$\overline{DM}_{RS,wRI}$

$\overline{DM}_{S_1,wRI}$

$\overline{DM}_{S_2,wRI}$

$\overline{DM}_{S_3,wRI}$

$\overline{DM}_{S_4,wRI}$ (see (7) and Fig. 2)

- The four values of similarity measures:

$AA_{RI}(S_1,RS)$

$AA_{RI}(S_2,RS)$

$AA_{RI}(S_3,RS)$

$A_{RI}(S_4,RS)$ (see (9))

All the measures for the 55 pairs of reference images and reference skeletons, are presented in the following website:

`https://www.inf.u-szeged.hu/~gnemeth/compskel/`

Observe that a smaller value in a row, means a better similarity of the sequence skeleton and the reference skeleton.

We know that $\mathcal{A}_{opt} \leqslant <1, 2> \leqslant <1>,<2>$ (see (12)), hence the following inequalities:

$$SM_{RI}\big(SS(RI, \mathcal{A}_{opt}), RS\big) \leq$$
$$SM_{RI}(SS(RI, <1,2>), RS) \leq$$
$$SM_{RI}(SS(RI, <1>), RS)$$

should hold for a reasonable similarity measure *SM*, for each pair of reference image and reference skeleton (*RI,RS*) (see (10)). We should add that the similarity measure *AA* satisfies it for high resolution images (in the case of Kimia dataset the image resolution is too low to measure big differences), hence, it is judged a reasonable similarity measure. Note, as well, that none of the five types of similarity measures applied in our previous paper [8] are reasonable.

# 5 Results

In this section we compare and rank nineteen thinning algorithms. The similarity measure *AA* has been computed for the $4\times55=220$ morphological skeletons driven by the four neighborhood sequences $<1>$, $<2>$, $<1,2>$ and $\mathcal{A}_{opt}$. In each case, $\mathcal{A}_{opt}$ provided the best result.

Evaluation is based on three different ranking methods:

1) *Sum of ranks*: For each test image, the similarity measure *AA* has been calculated and the algorithms have been sorted according the *AA* value. The scores have been summarized for each algorithm. The winner has the lower score value. Table 1 summarizes the result according to this ranking method.

2) *Sum of AA values*: The score of an algorithm has been computed as the sum of *AA* values for each test image. The winner algorithm has the lowest score. Table 2 shows the result according to this ranking.

3) *Tournament*: In a competition two algorithms play matches for each test image. If an algorithm is better (i.e., has a lower *AA* value) for more images than the other one in a competition, then it wins 1 point. In the tournament, the algorithms plays competitions pairwise. The best algorithm wins the most competitions. Table 3 presents the result of this ranking.

According to the our quantitative and fully automated comparison with the three types of rankings, we can state, that the 2*2-subiteration parallel thinning algorithm SI-<NE,SW,NW,SE>-E3, proposed by Németh Kardos and Palágyi [31], is the best choice, among the nineteen thinning algorithms compared.

**Conclusions**

A novel method for quantitative comparison of skeletons was presented herein. The proposed method is based on a new similarity measure and a gold standard 2D image database, containing pairs of reference images with elongated objects and their expected skeletons. Our method is validated using generalized morphological skeletons, driven by neighborhood sequences. According to the experiments, the proposed method can be used for evaluating arbitrary 2D skeletonization algorithms. Based on our method, the quantitative comparison of nineteen 2D thinning algorithms were presented, as well. In future work, we plan to extend our method to evaluate 3D skeletonization techniques.

Table 1
"Sum of ranks" method

| Rank | Algorithm | Ref. | Type | Sum of ranks for each image |
|---|---|---|---|---|
| 1 | SI-<NE,SW,NW,SE>-E3 | [31] | sub iteration-based | 535 |
| 2 | SI-<NE,SW,NW,SE>-E2 | [31] | sub iteration-based | 547 |
| 3 | BM99 | [32] | fully parallel | 916 |
| 4 | H89 | [33] | fully parallel | 1235 |
| 5 | SI-<NE,SW,NW,SE>-E1 | [31] | sub iteration-based | 1268 |
| 6 | GH92C | [34] | fully parallel | 1344 |
| 7 | GH89A1 | [35] | sub iteration-based | 1656 |
| 8 | FP-E3 | [31] | fully parallel | 1685 |
| 9 | FP-E2 | [31] | fully parallel | 1701 |
| 10 | PAV81 | [36] [37] | fully parallel | 1862 |
| 11 | EM93 | [38] | fully parallel | 1965 |
| 12 | FP-E1 | [31] | fully parallel | 2071 |

| 13 | AK$^2$ | [10] | fully parallel | 2131 |
|---|---|---|---|---|
| 14 | SI-<N,E,S,W>-E2 | [31] | sub iteration-based | 2965 |
| 15 | SI-<N,E,S,W>-E3 | [31] | sub iteration-based | 2965 |
| 16 | ZSLW | [39] | sub iteration-based | 3146 |
| 17 | SI-<N,E,S,W>-E1 | [31] | sub iteration-based | 3570 |
| 18 | RUT66 | [40] | fully parallel | 4267 |
| 19 | CWSI87 | [41] | fully parallel | 4759 |

Table 2
Sum of AA values" methods

| Rank | Algorithm | Ref. | Type | Sum of *AA* values |
|---|---|---|---|---|
| 1 | SI-<NE,SW,NW,SE>-E3 | [31] | sub iteration-based | 3.3108 |
| 2 | SI-<NE,SW,NW,SE>-E2 | [31] | sub iteration-based | 3.3122 |
| 3 | SI-<NE,SW,NW,SE>-E1 | [31] | sub iteration-based | 3.7920 |
| 4 | BM99 | [32] | fully parallel | 3.8911 |
| 5 | H89 | [33] | fully parallel | 4.0083 |
| 6 | PAV81 | [36] [37] | fully parallel | 4.0200 |
| 7 | GH92C | [34] | fully parallel | 4.0210 |
| 8 | FP-E3 | [31] | fully parallel | 4.0587 |
| 9 | FP-E2 | [31] | fully parallel | 4.0602 |
| 10 | EM93 | [38] | fully parallel | 4.0604 |
| 11 | AK$^2$ | [10] | fully parallel | 4.1016 |
| 12 | GH89A1 | [35] | sub iteration-based | 4.1181 |
| 13 | FP-E1 | [31] | fully parallel | 4.1665 |
| 14 | SI-<N,E,S,W>-E3 | [31] | sub iteration-based | 5.1786 |
| | SI-<N,E,S,W>-E2 | [31] | sub iteration-based | 5.1786 |
| 16 | ZSLW | [39] | sub iteration-based | 5.3166 |
| 17 | SI-<N,E,S,W>-E1 | [31] | sub iteration-based | 5.7857 |
| 18 | RUT66 | [40] | fully parallel | 6.9958 |
| 19 | CWSI87 | [41] | fully parallel | 9.0067 |

Table 3
"Tournament" methods

| Rank | Algorithm | Ref. | Type | Number of winner single combats |
|---|---|---|---|---|
| 1 | SI-<NE,SW,NW,SE>-E3 | [31] | sub iteration-based | 18 |
| 2 | SI-<NE,SW,NW,SE>-E2 | [31] | sub iteration-based | 17 |
| 3 | SI-<NE,SW,NW,SE>-E1 | [31] | sub iteration-based | 16 |
| 4 | BM99 | [32] | fully parallel | 15 |
| 5 | H89 | [33] | fully parallel | 14 |
| 6 | GH92C | [34] | fully parallel | 13 |
| 7 | GH89A1 | [35] | sub iteration-based | 12 |
| 8 | FP-E3 | [31] | fully parallel | 11 |
| 8 | FP-E2 | [31] | fully parallel | 10 |
| 10 | FP-E1 | [31] | fully parallel | 9 |
| 11 | PAV81 | [36] [37] | fully parallel | 8 |
| 12 | EM93 | [38] | fully parallel | 7 |

| 13 | AK$^2$ | [10] | fully parallel | 6 |
|----|--------|------|----------------|---|
| 14 | SI4-<N,E,S,W>-E3 | [31] | sub iteration-based | 4 |
|    | SI4-<N,E,S,W>-E2 | [31] | sub iteration-based | 4 |
| 16 | ZSLW | [39] | sub iteration-based | 3 |
| 17 | SI-<N,E,S,W>-E1 | [31] | sub iteration-based | 2 |
| 18 | RUT66 | [40] | fully parallel | 1 |
| 19 | CWSI87 | [41] | fully parallel | 0 |

## References

[1]  P. Giblin and B. B. Kimia: A formal classification of 3d medial axis points and their local geometry, *IEEE Transactions on Pattern Analysis and Machine Intelligence,* vol. 26, 2004, pp. 238-251.

[2]  K. Siddiqi and S. Pizer, Medial representations − Mathematics, algorithms and applications, vol. 37, Springer, 2008.

[3]  G. Németh and K. Palágyi: Topology preserving parallel thinning algorithms, *International Journal of Imaging Systems and Technology,* vol. 21, 2011, pp. 37-44.

[4]  L. Lam and C. Y. Suen: Automatic comparison of skeletons by shape matching methods, *Int. Journal on Pattern Recognition and Artificial Intelligence,* vol. 7, 1993, pp. 1271-1286.

[5]  W. Lee, L. Lam and C. Y. Suen: A systematic evaluation of skeletonization algorithms, *Int. Journal on Pattern Recognition and Artificial Intelligence,* vol. 7, 1993, pp. 1203-1225.

[6]  R. Plamondon, C. Y. Suen, M. Bourdeau and C. Barriére: Methodologies for evaluating thinning algorithms for character recognition, *Int. Journal on Pattern Recognition and Artificial Intellelligence,* vol. 7, 1993, pp. 1247-1270.

[7]  W.-P. Choi, K.-M. Lam and W.-C. Siu: Extraction of the Euclidean skeleton based on a connectivity, *Pattern Recognition,* vol. 36, 2003, pp. 721-729.

[8]  A. Fazekas, K. Palágyi, G. Kovács and G. Németh: Skeletonization Based on Metrical Neighborhood Sequences, *Proceedings of the 6th Int. Conf. Computer Vision Systems), LNCS*, vol. 5008, A. Gasteratos, M. Vincze and J. K. Tsotsos, (eds.), Springer, 2008, pp. 333-342.

[9]   T. Y. Kong and A. Rosenfeld: Digital topology: Introduction and survey," *Computer Vision, Graphics, and Image Processing,* vol. 48, 1989, pp. 357-393.

[10]  G. Bertrand and M. Couprie: New 2D parallel thinning algorithms based on critical kernels, *Proceedings of the 11th Int. Workshop Combinatorial Image Analysis*, *LNCS,* vol. 4040, R. Reulke, U. Eckhardt, B. Flach, U. Knauer and K. Polthier (eds.), Springer, 2006, pp. 45-59.

[11]  D. Shaken and A. Bruckstein: Pruning Medial Axes, *Computer Vision and Image Understanding,* vol. 69, no. 2, 1998, pp. 156-169.

[12]  G. Borgefors: Distance transformations in arbitrary dimensions, *Computer Vision, Graphics, and Image Processing,* vol. 27, 1984, pp. 321-345.

[13]  R. Fabbri, L. F. Costa, J. C. Torelli and O. M. Bruno: 2D Euclidean distance transform algorithms: A comparative survey, *ACM Computing Surveys,* vol. 40, no. 2, 2008, pp. 1-44.

[14]  R. Klette és A. Rosenfeld: Digital Geometry – Geometric Methods for Digital Picture Analysis, Morgan Kaufmann Publisher, 2004.

[15]  M. Couprie and G. Bertrand: Asymmetric parallel 3D thinning scheme and algorithms based on isthmuses, *Pattern Recognition Letters,* vol. 76, 2016, pp. 21-31.

[16]  A. Sobieczki, H. Yassan, A. Jalba and A. Telea: Qualitative comparison of contraction-based curve skeletonization methods, *Proceedings of 11th International Symposium on Mathematical Morphology*, Springer-Verlag, 2013, pp. 425-439.

[17]  A. Sobieczki, A. Jalba and A. Telea: Comparison of curve and surface skeletonization methods for voxel shapes, *Pattern Recognition Letters,* vol. 47, 2014, pp. 147-156.

[18]  C. Aslan, A. Erdem, E. Erdem and S. Tari: Disconnected skeleton: Shape at its absolute scale, *IEEE Trans. Pattern Analysis and Machine Intelligence,* vol. 30, no. 12, 2008, pp. 2188-2203.

[19]  X. Bai and L. J. Latecki: Path similarity skeleton graph matching, *IEEE Trans. Pattern Analysis and Machine Intelligence,* vol. 30, 2008, pp. 1282-1292.

[20]  J. Tschirren, G. McLennan, K. Palágyi, E. A. Hoffman and M. Sonka: Matching and anatomical labeling of human airway tree, *IEEE Trans. Medical Imaging,* vol. 24, 2005, pp. 1540-1547.

[21]  A. Brenneke és T. Isenberg: 3D shape matching using skeleton graphs,

*Proceedings of Simulation and Visualisation*, 2004, pp. 299-310.

[22] H. Sundar, D. Silver, N. Gagvani and S. Dickinson: Skeleton based shape matching and retrieval, *IEEE*, Washington, DC, USA, 2003.

[23] H. Barrow, J. Tenenbaum, R. Bolles and H. Wolf: Parametric condespondence and chamfer matching: Two new techniques for image matching, *Proceedings of the 5th international joint conference on Artificial intelligence*, Morgan Kaufmann Publishers Inc., 1977, pp. 659-663.

[24] A. Thayanantan, B. Stenger, P. S. Torr and R. Chipolla: Shape context and chamfer matching in cluttered scenes, *Proceedings of Computer Vision and Pattern Recognition*, 2003, pp. 127-133.

[25] J. Serra: Image Analysis and Mathematical Morphology, Academic Press, 1982.

[26] R. C. Gonzalez and R. E. Woods: Digital Image Processing (3rd Edition), Prentice Hall, 2008.

[27] A. Rosenfeld and J. L. Pfaltz: Distance functions on digital pictures," *Pattern Recognition,* vol. 1, 1968, pp. 33-61.

[28] A. Hajdu and L. Hajdu: Approximating the euclidean distance by digital metrics, *Discrete Mathematics,,* vol. 238, 2004, pp. 101-111.

[29] P. Maragos: Unified Theory of Translation-Invariant Systems with Applications to Morphological Analysis, and Coding of Images. PhD Thesis, Atlanta, GA: School of Elect. Engineering, Georgia Inst. of Technology, 1985.

[30] R. Kresch and D. Malah: Skeleton-based morphological coding of binary images, *IEEE Transactions on Image Processing,* vol. 7, 1998, pp. 1387-1399.

[31] G. Németh, P. Kardos and K. Palágyi: 2D parallel thinning and shrinking based on sufficient conditions for topology preservation, *Acta Cybernetica,* vol. 20, 2011, pp. 125-144.

[32] T. Bernard and A. Manzanera: Improved low complexity fully parallel thinning algorithm., *Proceedings of the 10th Int. Conf. on Image Analysis and Processing*, Venice, IEEE, 1999, pp. 215-220.

[33] R. W. Hall: Fast parallel thinning algorithms: parallel speed and connectivity preservation, *Communications of the ACM,* vol. 32, no. 1, 1989, pp. 124-131.

[34] Z. Guo and R. W. Hall: Fast fully parallel thinning algorithms., *Computer Vision, Graphics, and Image Processing,* vol. 55, no. 3, 1992, pp. 317-328.

[35] R. W. Hall: Parallel connectivity-preserving thinning algorithms, *Topological Algorithms for Digital Image Processing* , Elsevier, 1996, pp. 145-179.

[36] T. Pavlidis: A flexible parallel thinning algorithm., *Proceedings of IEEE Comp. Soc. Conf. Pattern Recognition, Image Processing*, 1981, pp. 162-167.

[37] T. Pavlidis: An asynchronous thinning algorithm., *Computer Graphics and Image Processing,* vol. 20, no. 2, 1982, pp. 133-157.

[38] U. Eckhardt and G. Maderlechner: Invariant thinning., *Pattern Recognition and Artificial Intelligence,* vol. 7, 1993, pp. 1115-1144.

[39] H. Lü and P. Wang: A comment on "A fast parallel algorithm for thinning digital patterns", *Communications of the ACM,* vol. 29, 1986, pp. 239-242.

[40] D. Rutovitz: Pattern recognition, *Journal of the Royal Statistical Society,* vol. 129, 1966, pp. 504-530.

[41] R. Chin, H. Wan, D. Stover and R. Iverson: A one-pass thinning algorithm and its parallel implementation., *Computer Vision, Graphics, and Image Processing,* vol. 40, no. 1, 1987, pp. 30-40.

[42] M. Couprie: Note on fifteen 2D parallel thinning algorithms. Internal Report., Université de Marne-la-Vallée, IGM2006-01, France, 2006.

[43] T. Sebastian, P. Klein and B. Kimia: Recognition of shapes by editing their shock graphs, *IEEE Trans. Pattern Analysis and Machine Intelligence,* vol. 20, no. 5, 2004, pp. 550–571.

[44] N. Cornea, D. Silver and P. Min: Curve-skeleton properties, applications and algorithms, *IEEE Transactions on Visualization and Computer Graphics,* vol. 13, no. 3, 2007, pp. 530-548.

# Evaluation and Choosing of Recycling Technologies by Using FAHP

## Pavlović Aleksandar[1*], Tadić Danijela[2], Arsovski Slavko[2], Jevtić Dragan[3], Pavlović Milan[4]

[1] Faculty of Economics and Engineering Management in Novi Sad
University Business Academy Novi Sad
Cvećarska 2, 21000 Novi Sad, Serbia, a.pavlovic@pfb-design.rs

[2] Technical Faculty "Mihajlo Pupin" in Zrenjanin, University of Novi Sad,
Djure Djakovica bb 23000 Zrenjanin, Serbia, pmilan@sbb.rs

[3] Faculty of Engineering, University of Kragujevac
6 Sestre Janjic Str. 34000 Kragujevac, Serbia, galovic@kg.ac.rs, cqm@kg.ac.rs

[4] Ministry of Spatial Planning, Civil Engineering and Ecology
Government of Republic of Srpska
Trg Republike Srpske 1, 78000 Banja Luka, Bosnia and Herzegovina,
d.jevtic@mgr.vladars.net

*Abstract: The evaluation and ranking of recycling technologies for each treated waste with respects to many different criteria has important results for the management team of any recycling center. Improvement of business strategy is based on the obtained rank of recycling technologies. It represents a key success factor for a recycling center in dealing with crisis. Uncertainties in: relative importance of evaluation criteria and priority of recycling technologies under each criterion are described by triangular fuzzy numbers. Relative importance of evaluation criteria is stated by fuzzy pair-wise comparison matrices. Determining of elements values of these matrices can be considered as a fuzzy group decision making problem. Aggregation of individual opinions into group consensus is performed by using fuzzy averaging method and Fuzzy Ordered Weighted Aggregation (FOWA,) Operator Fuzzy Analytic Hierarchy Process (FAHP) is used for determination of rank of recycling technologies with respects to evaluation criteria and its weights. Proposed model is tested by example with real life data.*

*Keywords: waste; recycling technologies; fuzzy set; FAHP*

# 1   Introduction

Theory and management of good practice has shown a strong interest in the domain of waste management, regulated by laws and standard ISO 14000 in the large number of developed and developing countries. During the last few decades, it has been frequently employed in order to establish mechanisms for environmental protection, reducing usage of natural resources, profit increase and better competitive positioning of any enterprise. Consequently, it is possible to realize economic sustainability of every country, meeting ecology standards, by using recyclable materials and applying these different recycling technologies.

Many and varied types of uncertainty exist in a treated problem. The term uncertainty implies that in a certain situation, a person does not have a tendency, indication or they lack ability to analyze information which quantitatively and qualitatively is appropriate to describe, prescribe or predict deterministically and numerically a system, its behavior or other characteristic [23]. It is assumed that these uncertainties are far better judged by using linguistic expressions than by representing them in terms of precise numbers. It is very useful in situations, which are too complex or not well defined to be reasonably described in conventional quantitative expressions [23].

The main contribution of this paper is the introduction of a fuzzy group decision making approach and fuzzy AHP for ranking of recycling technologies for each selected type of waist. In the literature, there are many developed approaches for handling FAHP. Chang introduced a new approach with use of triangular fuzzy numbers for pair-wise comparison scale of FAHP, and use of the extent analysis method for synthetic extent value of the pair-wise comparison [5]. Use of a developed approach does not involve cumbersome mathematical operations, and it has the ability to capture the vagueness of human thinking style. With respects to opinion, the authors of this paper suggested that fuzzy AHP is appropriate for evaluation and selection of recycling technologies with respect to numerous evaluation criteria and its relative importance in an uncertain environment.

The paper is organized in the following way. The literature review is presented in Section 2. The model framework, modelling of uncertainties and base of fuzzy Analytic Hierarchical Process which is introduced in [5] is given in Section 3. In Section 4, a proposed model is illustrated by an example with real-life data. Conclusions are presented in Section 5.

# 2   Literature Review

In the literature, many papers used to describe recycling technologies of different waste types can be found in [10]; [20]; [7]; [3]. Evaluation criteria of recycling technologies of each waste type are determined by experts and stakeholders of

each RC. It is assumed that choice of recycling technology for each type of waste is based on knowledge, experience of waste management and stakeholders of RCs.

Evaluation and selection of recycling technologies presents one of the most important management tasks of a recycling center (RCMT).   By respecting selected recycling technology, necessary finances for supply of recycling equipment, could be determined by review of personal capacity, quantity of recycling material, etc. Strategy for the increase of quality of management of RCs (one of the requirements of ISO 9000:2008), and strategy of sustainable regional development may be based on the obtained results. This approach is used for explanation importance of the problem.

Solutions of the decision making problems which belong to different research areas are obtained by using proposed FAHP [5]. For instance, Kahraman et al. considered selecting of the location facility [12]. Erensal et al. have determined key capabilities in technology management [9]. Environmental risk management as part of risk analysis is treated in [21]. Chan and Kumar used FAHP in selection and ranking the best global supplier for a manufacturing firm to supply one of its most critical parts used in assembling process [6]. The priorities of organizational capital measurement indicators are determined in [4]. Assessment of water management plans in the one region of Brazil is considered in [18].

By comparing papers that propose a modified FAHP, certain differences can be noted, which are further described. This analysis, at the same time, shows advantages of the proposed model.

Selection evaluation criteria can be given according to literature data (by analogy Kahraman et al., [12]; Erensal et al. [9]; Seçme et al. [17]; Srđević and Medeiros, [18]) or results of good practice (by Tesfamariam and Sadiq []21; Chan and Kumar [6]; Bozbura and Beskese [4]). In this paper, the assumption is introduced that evaluation criteria may be selected according to assessment of stakeholders.

Many authors suggest that rating of the relative importance of criteria and priorities of attribute should be stated as fuzzy group decision making problem. Aggregation of individual opinions into a group consensus can be performed by applying different methods, for instance by using the method of fuzzy average value (Kaya and Kahraman, [13]; Tadić et al. [19]) and the fuzzy ordered weight method (FOWA) (Merigó and Gil-Lafuente [16]) in [1]. In this paper, aggregation of opinions of RCMTs is performed by fuzzy averaging method. The aggregated assessment of stakeholder is given by FOWA. With respect to nature of considered problem, these operators are used in this paper.

The modelling of the relative importance of criteria and priority of alternative is based on the fuzzy set theory. These decision variables are presented by triangular fuzzy numbers (TFNs) in all analyzed papers. Respecting this fact, in this paper is the introduced assumption that uncertainty should be modelled by TFNs. In the considered papers, criteria weights are determined by using FAHP which is

proposed in [5]. In this paper, FAHP (by analogy [5]) is applied for determining criteria weights.

The main difference and significant contribution of the proposed FAHP that has been analyzed in this section, is calculation of values of elements of fuzzy pair-wise comparison matrices. In the authors' opinion, introduced modifications of FAHP provide significantly more reliable information to decision makers than the proposed FAHP which can be seen in the literature. Therefore, RCMTs may determine the best recycling technology for each waste type.


# 3    Methodology

The evaluation and choosing of recycling technologies for each identified waste type is based as multi-criteria optimization task under uncertainties. The assessments of decision makers are described by linguistic expressions. It is assumed that it is a closer to human way of thinking that assessments are made by decision makers are represented by using linguistic expressions than precise numbers. In this paper, existing uncertainties are modelled by using fuzzy set theory [15]; [23]. In the literature, the TFNs are widely used for modelling different uncertainties. TFNs offer a good compromise between descriptive power and computational simplicity. In this paper, with respect to the type and size of the considered problem and results of investigators, five linguistic expressions are used, at the most, assigned to the existing linguistic variables. Ranking of recycling technologies for each identified waste type with respect to pre-defined evaluation criteria and their weights is given by applied FAHP, which is proposed in [5]. In literature this method is used in many papers where different problems of fuzzy multi-tasking decision making are described [17]; [22]; [13]; [19].


## 3.1    The Mathematical Formulation of Treated Problem

In this paper, reverse logistic chain (RLC) consists of a few recycling centers which can be presented by set indices $\Omega=\{1...j...J\}$. The total number of recycling centers (RCs) of considered RLC is denoted as J, and j, j=1,…,J is index of recycling center. In general, waste is presented by set of indices $I=\{1...i...I\}$, where the index of waste type is denoted as i and I is the total number of identified waste types. Each waste type i, i=1,..,I at the each RC j, j=1,..,J can be recycled by applying different recycling technologies. Recycling technology is determined for each type of waste, separately. Data base of recycling theory and their characteristics (capacity, economy characteristics, ecology characteristics, (Data bases are made by domestic and international waste management associations) employment, etc.), exists. Recycling technology of each waste type are defined by RCMT. Assessment of RCMT is based on experience, knowledge, and the initial prediction of RCMT. The recycling technologies which can be used for recycling

of waste type i, i=1,..,I are presented by set of indices $\Theta_i = \{1,\ldots,t,\ldots T_i\}$. The index of recycling technology of each waste type i, i=1,..,I is denoted as $t, t = 1,..,T_i$ and $T_i, i = 1,\ldots,I$ is the total number of defined recycling technologies for waste type i, i=1,..,I. Evaluation of recycling technologies for any identified waste types is performed with respects to many criteria. The kind and number of evaluation criteria are defined by literature sources. Formally, the evaluation criteria are presented by set indices $\chi = \{1,\ldots,k,\ldots K\}$, where k is index of evaluation criteria and K is the total number of evaluation criteria. It can be assumed that all the criteria for evaluating selected recycling technologies of each waste type are usually not of the same relative importance, and do not depend on the recycling technologies.

The fuzzy rating of the relative importance of each pair of evaluation criteria is performed by stakeholders at the recycling center level (general managers of recycling centers, managers of local administrations, and the main technology staff of recycling centers). It is supposed that stakeholders at the recycling center level make decisions by consensus. It is assumed that stakeholders of different RCs are not equal, as far as, relative importance. This assumption is introduced because of economy development, and quantity of waste and its morphology generated in different regions are not the same. The aggregation of judgements of stakeholders of considered RCs is performed by Fuzzy Ordered Aggregated Operation (FOWA) [16]. The fuzzy rating of each pair of recycling technologies under each pre-defined criterion is performed by RCMT (main manager, technology staff, manager of ecology, and financial manager). The aggregation of individual opinions of RCMTs into group consensus can be given by fuzzy averaging method.

## 3.2   Modelling of Uncertainties

Fuzzy pair-wise comparison matrix of the relative importance of evaluation criteria and treated recycling technological are stated by analogy [14]; [1].

In this paper, the fuzzy rating of stakeholders of each RC is described by linguistic expressions which can be represented as TFN $\tilde{W}_{kk'}^{j} = \left(x; l_{kk'}^{j}, m_{kk'}^{j}, u_{kk'}^{j}\right)$. Value 1 marks the evaluation criterion k over evaluation criterion k´, k, k´=1, … , K´; k ≠ k´, has lower importance. On the other hand, the value 9 denotes that evaluation criterion k over evaluation criterion k´, k, k´=1, … , K´; k ≠ k´ has the most importance. If strong relative importance of criterion k´ over criterion k holds, then pair-wise comparison scale can be represented by the TFB

$$\underset{kk'}{\overset{j}{\sim}}{w} = \left( \underset{k'k}{\overset{j}{\sim}}{w} \right)^{-1} = \left( \frac{1}{\underset{k'k}{\overset{j}{u}}}, \frac{1}{\underset{k'k}{\overset{j}{m}}}, \frac{1}{\underset{k'k}{\overset{j}{l}}} \right).$$ If k = k′ (k, k′ = 1, … , K) then relative

importance criterion k over criterion $k^{'}$ is represented by single point 1 which is a triangular fuzzy number (1,1,1). Similarly,, the preference of recycling technology t over recycling technology t′, t, t′ = 1, … , $T_i$; t ≠ t′ are assessed by RCMT e, e=1,..,E by using pre-defined linguistic expressions. These linguistic expressions are modeled by triangular fuzzy numbers which are given in the following way:

*Very low level importance / preferred (VL)- (x; 1. 1. 5.5)*

*Low level importance / preferred (L)- (x; 1, 3, 9)*

*Middle level importance / preferred (M)- (x;1, 5, 9)*

*High level importance / preferred (HV)- (x; 1, 7, 9)*

*Very high level importance / preferred (VH) – (x; 4.5, 9, 9)*

## 3.3 The Proposed Algorithm

The proposed procedure can be realized through steps.

*Step 1*. The fuzzy pair-wise comparison matrix of the relative importance of evaluation criteria for each RC are stated:

$$\left[ \underset{kk'}{\overset{j}{\tilde{W}}} \right]_{K \times K} \quad k, k^{'} = 1,...,K; k \neq k^{'}; \ j = 1,..,J$$

*Step 2*. Aggregation of individual judgements of stakeholders of RCs is obtained by applying FOWA:

$$\left[ \tilde{W}_{kk'} \right]_{K \times K} \quad k, k^{'} = 1,...,K; k \neq k^{'}; \ j = 1,..,J$$

where: $\tilde{W}_{kk'} = \sum_{j=1}^{J} \omega_j \overset{j}{\tilde{W}}_{kk'}$ and $\omega_j$ is the relative importance of the stakeholder of the recycling center j, j=1,..,J

*Step 3*. The aggregated fuzzy rating of the priority of each pair of recycling technologies is given by using the fuzzy averaging method:

$$\tilde{W}_{t,t'} = \frac{1}{E} \overset{e}{\tilde{W}}_{t,t'}$$

*Step 4*. Calculate the criteria priority weights for each RC j, j=1,…,J and recycling technologies priority preferences [5]. The weights vector represented is obtained by applying the method for fuzzy numbers comparison [2]; [8].

In a similar way, the normalized priorities vector of recycling technologies for waste type i, i=1,..,I under each identified evaluation criteria and for each waste type is:

$$V_{t\,k}^{i} = \left( v_{1\,i}^{i},..,v_{t\,k}^{i},...,v_{T_i\,k}^{i} \right) t = 1,...,T_i; i = 1,..,I; k = 1,..,K$$

*Step 5*. The fuzzy composite priorities of the recycling technology for waste type i, i=1,..,I is given:

$$z_{t,i} = \sum_{k=1}^{K} w_k v_t^i k$$

*Step 6*. Organize all $z_{t,i}$ in descending sequence. Recycling technology first in sequence could be considered the best for that type of waste o i, i=1,.,I.

# 4   Case Study

Developed model is tested on real life data which are obtained from the region of West Balkan. There are six recycling centers (see Fig. 1). Recycling centers j=1 and j=2 are in region of Zvornik, in the region Brcko is the recycling center denoted as j=3, recycle centers j=4 and j=5 are in Tuzla and recycling center j=6 is in Trebinje.

The proposed model is tested by data from the social industry. Participation of building industry in social product of the considered region of West Balcan is about 15%-20%. It could be said that building industries are taking a greater part in the definition of the regional development strategy. The waste is classified according to Schedule of categories of waste (Officia lGazzete Bosnia and Hercegovina No. 9/05). Choosing waste types treated in the paper is based on stakeholder's assessment and forecasted waste quantities. The six waste types: concrete (I=1), brick (I=2), tile (I=3), rubber (I=4), plastics (I=5), and ash (I=6). Choosing of these waste types is performed with respect to two criteria: (1) estimated waste qualities which are stored and recycling in the treated RCs, and (2) demand for recycles which are obtained by applying recycling technologies. Recycling technologies for each specific type of waste are defined according to the existing data base of recycling technologies and presented in Table 1.

Table 1

Possible recycling technologies for every type of the waste

| Type of waste | Possible recycling technologies |
|---------------|--------------------------------|
| i=1 | mobile recycling technologies; mobile technologies for recycling and separation; mobile recycling technologies, selection and separation |
| i=2 | cleaning for reusing; cleaning for reusing and selection; processing bricks for reuse |
| I=3 | cleaning for reusing; cleaning for reusing and selectin; processing bricks for reuse |
| i=4 | grinding process; grinding process and separation; grinding and pyrolysis |
| i=5 | selection; grinding process, process of pressing; |
| i=6 | ash separation; recycling of flying ash; slug recycling; recycling and pressing |

The number and kind of criteria used to determine mark of recycling technologies for each waste type are defined by administrations of Bosnia and Hercegovina, provinces, city administrations, owner of each considered recycling center, recycling equipment manufactories, etc. The evaluation criteria are: employment level (k=1), quantities of waste (k=2), environment impact (k=3), sustainable development of city and province (k=4), social cohesiveness of province (k=5), enrolment of people from different social categories (k=6), level of usage of waste material (k=7), compliance with European Union Standards considering waste management (k=8), support of domestic industry (k=9), level of soil usage (k=10), innovation ability of local providers (k=11), level of dependence based on imported material (k=12), and relation of price between recycled resource and price of resource on market (k=13).

With respect to number of populations in a considered region, kinds of industrial enterprises which exists in provinces, it is assumed that the relative importance of evaluation criteria may be defined for a group of RCs. The treated RCs are divided into three groups: the first group consists of three RCs (j=1; j=2; j=6). The recycling center (j=3) presents the second group of recycling centers. The third group of recycling centers consists of two recycling centers (j=4; j=5).

The normalized weights vector of evaluation criteria is:

W = ( 0.103 , 0.085, 0.085, 0.107, 0.066, 0.077, 0.106, 0.095, 0.075, 0.035, 0.033, 0.034, 0.100)

The priorities of recycling technologies of each considered waste type under each evaluation criterion is calculated and presented in the following tables.

Table 2

The priorities of the recycling technologies under each evaluation criterion

| | Vector priorities of recycling technologies for waste type i=1 | Vector priorities of recycling technologies for waste type i=2 and i=3 | Vector priorities of recycling technologies for waste type i=4 |
|---|---|---|---|
| k=1 | $[0.191,0.348,0.461]^{T}$ | $[0.333,0.333,0.333]^{T}$ | $[0.207,0.387,0.406]^{T}$ |
| k=2 | $[0.179,0.247,0.574]^{T}$ | $[0.202,0.388,0.410]^{T}$ | $[0.175,0.396,0.429]^{T}$ |
| k=3 | $[0.258,0.334,0.408]^{T}$ | $[0.333,0.333,0.333]^{T}$ | $[0.447,0.297,0.256]^{T}$ |
| k=4 | $[0.158,0.392,0.479]^{T}$ | $[0.307,0.336,0.362]^{T}$ | $[0.296,0.307,0.398]^{T}$ |
| k=5 | $[0.161,0.375,0.463]^{T}$ | $[0.291,0.338,0.371]^{T}$ | $[0.180,0.371,0.449]^{T}$ |
| k=6 | $[0.151,0.378,0.471]^{T}$ | $[0.204,0.320,0.476]^{T}$ | $[0.189,0.373,0.439]^{T}$ |
| k=7 | $[0.122,0.352,0.526]^{T}$ | $[0.333,0.333,0.333]^{T}$ | $[0.301,0.329,0.370]^{T}$ |
| k=8 | $[0.233,0.323,0.444]^{T}$ | $[0.282,0.332,0.385]^{T}$ | $[0.279,0.350,0.371]^{T}$ |
| k=9 | $[0.215,0.278,0.507]^{T}$ | $[0.215,0.271,0.513]^{T}$ | $[0.275,0.354,0.371]^{T}$ |
| k=10 | $[0.210,0.280,0,510]^{T}$ | $[0.333,0.333,0.333]^{T}$ | $[0.419,0.373,0.208]^{T}$ |
| k=11 | $[0.291,0.332,0.377]^{T}$ | $[0.241,0.305,0.454]^{T}$ | $[0.333,0.333,0.333]^{T}$ |
| k=12 | $[0.333,0.333,0.333]^{T}$ | $[0.333,0.333,0.333]^{T}$ | $[0.333,0.333,0.333]^{T}$ |
| k=13 | $[0.258,0.326,0.416]^{T}$ | $[0.557,0.370,0.073]^{T}$ | $[0.333,0.333,0.333]^{T}$ |

By using the proposed algorithm (Step 5 and Step 6), the fuzzy composite priorities of the recycling technology and their rank is presented.

Applying $(t_{3\,1})$ we get have a concrete of a different granulation which has a higher market value than a recycled one, obtained by using an additional two other methods. Placed first, in the ranking is recycling technology $(t_{3\,1})$. Choosing recycling technology corresponds certain rank. The main constrain of applying this recycling technology regards to RC financial resources.

Table 3

The priorities of the recycling technologies under each evaluation criterion (continue)

| | Vector priorities of recycling technologies for waste type i=5 | Vector priorities of recycling technologies for waste type i=6 |
|---|---|---|
| k=1 | $[0.247,0.345,0.408]^{T}$ | $[0.170,0.200,0.319,0.311]^{T}$ |

| | | |
|---|---|---|
| k=2 | $[0.187, 0.370, 0.443]^T$ | $[0.250, 0.250, 0.250, 0.250]^T$ |
| k=3 | $[0.378, 0.348, 0.274]^T$ | $[0.118, 0.221, 0.323, 0.318]^T$ |
| k=4 | $[0.247, 0.338, 0.415]^T$ | $[0.179, 0.372, 0.224, 0.224]^T$ |
| k=5 | $[0.215, 0.363, 0.421]^T$ | $[0.161, 0.386, 0.227, 0.227]^T$ |
| k=6 | $[0.150, 0.346, 0.501]^T$ | $[0.351, 0.147, 0.207, 0.237]^T$ |
| k=7 | $[0.206, 0.343, 0.451]^T$ | $[0.228, 0.193, 0.282, 0.297]^T$ |
| k=8 | $[0.282, 0.316, 0.402]^T$ | $[0.205, 0.192, 0.189, 0.415]^T$ |
| k=9 | $[0.205, 0.348, 0.447]^T$ | $[0.164, 0.184, 0.314, 0.338]^T$ |
| k=10 | $[0.242, 0.260, 0.5123]^T$ | $[0.084, 0.192, 0.317, 0.340]^T$ |
| k=11 | $[0.177, 0.357, 0.465]^T$ | $[0.161, 0.081, 0.321, 0.336]^T$ |
| k=12 | $[0.306, 0.331, 0.363]^T$ | $[0.186, 0.196, 0.229, 0.389]^T$ |
| k=13 | $[0.225, 0.311, 0.464]^T$ | $[0.219, 0.202, 0.287, 0.929]^T$ |

Table 4

Rank of technology used for concrete recycling (i=1):

| Type of technology | Total priority coefficient | Rank |
|---|---|---|
| Mobile recycling technology ($t_{1\,1}$) | 0.2004 | 3 |
| Mobile recycling technology and selection ($t_{2\,1}$) | 0.3345 | 2 |
| Mobile recycling technology, selection and separation ($t_{3\,1}$) | 0.4685 | 1 |

Table 5

Rank of technologies used for brick recycling (i=2) and tile recycling (i=3)

| Type of technology | Total priority coefficient | Rank |
|---|---|---|
| Cleaning for brick/tile reuse ($t_{1\,2}$)/ ($t_{1\,3}$) | 0.3124 | 3 |
| Cleaning for reuse and brick/tile ($t_{2\,2}$)/ ($t_{2\,3}$) | 0.3357 | 2 |
| Brick/tile processing for reuse($t_{3\,2}$)/ ($t_{3\,3}$) | 0.3532 | 1 |

With respect to calculated values of priority coefficient it can be concluded that each of considered recycling technologies for brick as well for tile could be in first place. It is expected because the advantage of each technology is overruled by its disadvantage. For example, applying technology $(t_{1\,2})$/ $(t_{1\,3})$ technological level of process is lower (therefore, using this type of technology is much cheaper) but amount of labor is higher (usually unskilled). Choice of recycling technology for these two types of waste should be based on results of realized cost-benefit analysis. If wasted brick/tile is broken, then $(t_{1\,2})$/ $(t_{1\,3})$, or $(t_{2\,2})$/ $(t_{2\,3})$ will overpower $(t_{3\,2})$/ $(t_{3\,3})$. Otherwise, the best recycling technology is $(t_{3\,2})$/ $(t_{3\,3})$ which is placed first in the ranking.

Table 6
Rank of technologies used for recycling of rubber (i=4):

| Type of technology | Total priority coefficient | Rank |
|---|---|---|
| Grinding process $(t_{1\,4})$ | 0.2815 | 3 |
| Grinding processes and separation $(t_{2\,4})$ | 0.3183 | 2 |
| Grinding process, separation and pyrolysis $(t_{3\,4})$ | 0.3712 | 1 |

The technology $(t_{34})$ is placed first in the ranking. With respect to the given results it can be concluded that this recycling technology has the most priority compared to the other two technologies. Basic hydrocarbons obtained by using the process of pyrolysis, represents raw materials found in production of rubber and plastic products. Recycled material which we obtain by using $(t_{34})$ have a higher value in the market of reused raw materials than ones obtained by already analyzed technologies. It can be concluded that business efficiency RC could be significantly increased by applying $(t_{34})$. The main task of RCMT can be defined as taking over management initiatives (e.g., a continuous supply to RCs of sufficient quantities of rubber which leads to application of $(t_{34})$ will be economically justified. The main disadvantage of the technology $(t_{34})$ compared to $(t_{14})$ and $(t_{24})$ are a higher cost of recycling. The cost of pyrolysis devices represents the majority of recycling costs by applying the technology $(t_{34})$.

Table 7
Rank of technologies used for recycling of plastic (i=5):

| Type of technology | Total priority coefficient | Rank |
|---|---|---|
| selection $(t_{1\,5})$ | 0.2169 | 3 |
| grinding process $(t_{2\,5})$ | 0.3182 | 2 |
| pressing processes $(t_{3\,5})$ | 0.4253 | 1 |

With respect to the rank of recycling technologies for (i=5), it can be concluded that the dominant technology is $(t_{35})$. Applying $(t_{35})$ we get recycle material, which represents the final product which can be used in different types of industry.

Market cost of recyclate acquired by using technology ($t_{35}$) is much higher than recycle materials obtained by applying technologies ($t_{15}$) and ($t_{25}$). Recycle material could be differently granulated which impacts its market value. Recycle material obtained by using technology ($t_{25}$) cannot be used like a final product; it is used like a raw material in production of different products in building and other branches of industry.

Table 8
Rank of technologies used in recycling of ash (i=6):

| Type of technology | Total priority coefficient | Rank |
|---|---|---|
| ash separation ($t_{1\,6}$) | 0.2015 | 4 |
| recycling of flying ash ($t_{2\,6}$) | 0.2293 | 3 |
| recycling of bottom ash ($t_{3\,6}$) | 0.2660 | 2 |
| recycling and pressing ($t_{4\,6}$) | 0.2981 | 1 |

Importance of recycling this type of waste can be illustrated by forecasting quantities (about 360 000 t per year) and numerous and varied sources for this type of waste. Fist in the ranking is technology ($t_{46}$). This technology should not have been realized if another technology is not considered. Technology ($t_{46}$) should be used where the final product could be transported to factories where ash is used like a production raw material. Usually, recycled ashes are used in production of cement because ash is one of the basic raw materials found in cement production.

**Conclusion**

Based on the results of good practice in developed countries, it is known that well organized, existing RCs have a high influence on realization of state development strategy. One of the management problems of RCMT is selection technologies which can be applied at each RC. The criteria toevaluate recycling technologies are defined by the stakeholders of RCs. It is assumed that stakeholders of treated RCs have different relative importance. The assessment and selection of recycling technologies may be introduced through identification of waste type. Assessment of defined recycling technologies priorities is performed by RCMT at the RC level. Solution of the considered problem is obtained in an exact way because the solution is less burdened by the subjective judgments of the decision makers.

The priority of selected recycling technologies with respect to all evaluation criteria and their weights is obtained by using fuzzy AHP. The elements of fuzzy pair-wise comparison matrix of the relative importance of evaluation criteria are calculated by using FOWA operator. The element of fuzzy pair-wise comparison matrix of the priority of the selected recycling technologies under each type of waste is given by using fuzzy averaging method.

With respects to economic aspect, the worst technology is the one who gives the recycle materials numerous restrictions, insufficient funds, unskilled staff, unsuitable capacity of RC etc. Respecting constrains, choosing recycling technologies is very important for RCMTs because of that existing strategy and development of RC is based on obtained results. On other hand, improvement strategy (requirements of ISO 9001:2008, and ISO 144000) of the recycling processes should be based on obtained results.

The proposed procedure is illustrated by real-life data from RCs in the West Balcan region. Some of the possible strategies that may be employed for improving values of named RFs are: creation of partner relationships with current suppliers, increasing warehouse potential related to resources if market conditions are impacted by an unstable politic and economic situation, and development of safety–critical systems.

Besides the advantages, the proposed model has certain constraints, which are: the number of type of waste, available capacity of RC, change of number of recycling technologies for one or all considered types of waste, change of political and economic environment, etc. For a set period of time (in this case period of one year is realistic) it could be considered that selected technology has a higher priority for treated RCs.

The focus of future research should be set on a case study with a large sample of type of wastes in each RC. All of these modifications can be easily and quickly incorporated into the proposed model and do not increase the complexity of the mathematical computation. In addition, the software solution could be expanded with additional functionalities for better management of RCs.

### Acknowledgement

### References

[1]     Aleksić, A., Stefanović, M., S. Arsovski, and D. Tadić (2013) An Assessment of Organizational Resilience Potential in SMEs of the Process Industry, a Fuzzy Approach, Journal of Loss Prevention in the Process Industries, 26, 1238-1245

[2]     Bass, S.M., Kwakernak, H. (1979) Rating and Ranking of Multiple-Aspect Alternatives Using Fuzzy Sets, Automatica, Vol. 3, 47-58

[3]     Bolden, J., Abu-Lebdeh, T., Fini, E. (2013). Utilization of recycled and waste materials in various construction applications, American Journal of Environmental Science, Vol. 9, No. 1, 14-24.

[4]     Bozbura, S. B., Beskese, A. (2007) Prioritization of Organizational Capital Measurement Indicators Using Fuzzy AHP, International Journal of Approximate Reasoning, 44, 124-147

[5]     Chang, D., Y., (1996) Applications of the Extent Analysis Method on Fuzzy AHP. European J. of Operational Research, 95 649-655

[6]     Chan, S. T. F., Kumar, N. (2007) Global Supplier Development Considering Risk Factors Using Fuzzy Extended AHP-based Approach, Int. J. of Production Research, 46, 417-431

[7]     de Brito, J., Robles, R. (2010) Recycled Aggregate Concrete (RAC) Methodology for Estimating its Long-Term Properties, Indian Journal of Engineering & Materials Science, 17, 449-492

[8]     Dubois, D., Prade, H. (1979) Decision-Making under Fuzziness. In Advances in Fuzzy Set Theory and Applications (ed. R. R. Yager), Ed.-North-Holland, 279-302

[9]     Erensal, Y. C., Oncant, T., Dermican, M. L. (2006) Determining Key Capabilities in Technology Management Using FAHP: A Case Study of Turkey. Information Science, 176, 2755-2770

[10]    Hanneqaurt, C.: Good Practice Guide on Waste Plastics Recycling (2004) A Guide by and for Local and Regional Authorities. Association of Cities Regions for Recycling (ACRR), Belgium

[11]    Hwang, C. and Masud, A. (1979) Multi Objective Decision Making-Methods and Applications: a State-of-the Art Survey. Springer, Berlin

[12]    Kahraman C (2008) Introduction: Fuzzy Theory and Technology. Multiple Valued Logic and Soft Computing, 15(2-3) 103-115

[13]    Kaya, T., Kahraman, C. (2011) Multicriteria Decision Making in Energy Planning Using a Modified Fuzzy TOPSIS methodology, Expert Systems with Applications, 38, 6577-6585

[14]    Kelemenis, A, Askounis, D., (2010) A New TOPSIS-based Multi-criteria Approach to Personal Selection. Expert System with Applications, 7 (37), 4999-5008

[15]    Klir, G. J., Folger, T., Fuzzy Sets, Uncertainty, and Information. Prentice Hall, Upper Saddle River, NJ., USA, 1988

[16]    Merigó, J. M., Gil-Lafuente, M. A., (2011) Fuzzy Induced Generalized Aggregation Operators and its Application in Multi-Person Decision Making. Expert System with Applications, 38 (8), 9761-9772

[17]    Seçme, Y. N., Bayrakdaroğu, Kahraman, C. (2009) Fuzzy Performance Evaluation in Turkish Banking Sector using Analytic Hierarchy 11709. Process and TOPSIS, Expert Systems with Applications, 36, 11699-11709

[18]    Srđević, B., Medeiros, Y. D. P. (2008) Fuzzy AHP Assessment of Wather Management Plants, Wather Resources Management, Vol. 22, 877-894

[19]　Tadić, D., Gumus, T. A., Arsovski, S., Aleksić, A., Stefanović, M. (2013) An Evaluation of Quality Goals by Using Fuzzy AHP and Fuzzy TOPSIS Methodology. Journal of Intelligent & Fuzzy Systems, 25, 547-556

[20]　Tam, V. W. Y., Tam, T. C. M. (2006) A Review on the Viable Technology

for Construction Waste Recycling, Resources. Conservation and Recycling, 47 209-221

[21]　Tesfamariam, T. S., Sadiq, S. R. (2006) Risk-based Environmental Decision-Making Using Fuzzy Analytic Hierarchy Process (FAHP) Stoc. Environmetal Research Risk Assessment, 21, 35-50

[22]　Torfi, F., Farahani, Z., R., Rezapour, S. (2010) Fuzzy AHP to Determine the Relative Weights of Evaluation Criteria and Fuzzy TOPSIS to Rank the Alternatives, Applied Soft Computing, 10, 520-528

[23]　Zimmermann, H. J. (2001) Fuzzy set Theory and its Applications. Kluwer Nijhoff Publising: Boston

# The Impact of Creative Construction Tasks on Visuospatial Information Processing and Problem Solving

## Bernadett Babály[1,2], Andrea Kárpáti[3]

[1] Szent István University, Ybl Miklós Faculty of Architecture and Civil Engineering, Thököly út 74, H-1146 Budapest, Hungary
[2] Eötvös Loránd University, Doctoral School of Education, Kazinczy u. 23-27, H-1075 Budapest, Hungary
[3] Eötvös Loránd University, Faculty of Science, Pázmány Péter sétány 1/A, 7/7 25, H-1117 Budapest, Hungary
E-mail: babaly.bernadett@ybl.szie.hu; andrea.karpati@ttk.elte.hu

*Abstract: This study explores the potential of using creative 3D modeling for the development of spatial abilities. We investigate the efficiency of spatial training programs with a focus on differences in spatial information processing in real and virtual environments. Participants were architecture and civil engineering students in the first and second study year. The standardized Spatial Ability Test by Séra, Kárpáti and Gulyás (2002) was used for the assessment of relevant skill components: spatial perception, visualization and mental manipulation. In order to analyze visuospatial information processing and problem solving, we documented the phases of planning and modeling and revealed problems and motivating factors of the design process through student surveys. We discuss factors influencing the perception and interpretation of space and showed strategies of engineering students in solving spatial problems. The effectiveness of the program was unrelated to gender, specialization, secondary level studies and learning environments (real and virtual spaces). Post-test results of the experimental groups were significantly higher (t[226]=-4.70, p<0.001) and the effect size of the developmental program was d=1.07. Research has proven that an appropriately constructed set of creative problem solving tasks in modeling and construction, results in significant development of spatial skills and are as effective as traditional drawing tasks. Creative modeling is an activity with high motivation value and can be utilized to develop spatial abilities that are basic for the professional development of engineers and architects.*

*Keywords: spatial abilities; creative problem solving; 3D modeling; project pedagogy; virtual learning environment*

# 1    Introduction

Creative problem solving is highly appreciated by educators, but rarely introduced in school curricula. In a teacher opinion survey by the European Commission on creativity, over 95% of respondents agreed that creativity is a fundamental competence that can be applied to every domain of knowledge and to every school subject, and therefore developed by every discipline at school [8]. The *Partnership for 21ˢᵗ Century Skills*, an American organization of worldwide recognition that advocates the fusion of the three R's and four C's (critical thinking and problem solving, communication, collaboration and creativity and innovation), also emphasizes the importance of research-based examples to realize these objectives [36].

Visuospatial information processing is present in a wide range of everyday activities, from gardening or sewing to building self-representations on a social website. Construction challenges for knowledge building, develops problem solving skills and prepares for flexible retrieval and utilization of information in the world of work. Tens of thousands of people around the world organize peer-learning communities to acquire its technology and aesthetics as it is rarely taught in schools [12]. Image production (both in two- or three- dimensional formats) fosters the creation of accurate mental representations [27] [31]. Design and construction is associated with genetic forms of knowledge building: exploration, trial and play. The design process involves instinctive, spontaneous phases that may lead to inspiring detours and incongruences that result in the discovery of new solutions [1]. In the project reported here, we want to prove that open-ended, creative design and construction tasks may develop spatial skills with at least the same efficacy as traditional methods based on drill-like exercises in representational conventions. We also show that creative tasks inspire students at a university of technology, to produce high quality work that illustrates the development of their creativity, as well as, their spatial skills.

# 2    The Role of Creative Construction Tasks in the Development and Assessment of Visuospatial Information Processing and Spatial Representation

Visuospatial information processing and spatial representation (referred to as spatial skills in the subsequent part of this paper) are key areas of development in many educational disciplines. STEM (Science, Technology, Engineering and Mathematics) educators claim that spatial ability influences knowledge acquisition in their disciplines [7] [35] [45]. If one spatial skill component is developed, it may affect the level of others [39]. Different types of learning environments have special effects on the development of spatial skills [40] [42] [43].

Several research findings suggest that construction activities positively affect spatial skills as well as learning achievement in mathematics and science. Moreover, learning deficits may be revealed and individually treated in early childhood through construction game based tools [38]. In a project by *Verdine*, *Golinkoff*, *Hirsh-Pasek*, *Newcombe*, *Filipowicz* and *Chang* [47] studied the emergence of construction skills in relation to verbal literacy, numeracy, gender and social background. Both the development and the assessment were performed through building tasks with LEGO bricks: children had to build a construction out of 2-4 bricks according to a drawing. *McKnight* and *Mulligan* [29] saw the role of construction toys in activating intuitive, informal modes of knowledge acquisition and thus reveal the level of differentiated thinking. They used open-ended construction tasks for studying spatial problem solving skills. A similar research approach was employed by *Ferrara*, *Hirsh-Pasek*, *Newcombe*, *Golinkoff* and *Lam* [11] who developed spatial skills through three methods: (1) free play with building blocks; (2) guided play; (3) play with prefabricated constructions. No tasks were prescribed, but the researchers analyzed how children interacted with the educational environment through assembling and disassembling, constructing and reconstructing the blocks.

The major influential factor in the enhancement of spatial skills is teaching methodology. The *developmental effects of construction,* was first utilized systematically in education by *Friedrich Froebel*. Developed in the early 1800s for his Kindergarten, the Froebel Gifts appear deceptively simple but elicit a wide range of sophisticated operations described in detailed manuals. With a wide range of construction activities, he clearly aimed at the development of visuospatial information processing and through this, cognitive development: "From objects to pictures, from pictures to symbols, from symbols to ideas, leads the ladder of knowledge." (*Froebel* [13] quoted by *Marenholtz-Buelow* [30] p. 36). Most of the Froebel Gifts were monochrome: a clear indication of the emphasis on structural relationships among geometric elements, and not the imitation of real-life objects. *Maria Montessori* employed construction as a free experimentation method through which sensory organs may be developed. In the fifties, when the effects of construction on cognitive development became widely studied, her "Child Size Tools for Small Hands" series and her creative tasks based on experimentation inspired educational toys worldwide. The US Government encouraged the inclusion of building and construction in several curricular areas to promote learning in mathematics and science [16].

The meta-analysis by *Nath* and *Szűcs* [34] of studies that indicate the *effects of building and construction activities on performance in geometry and general spatial ability tests*, claim that most studies lack data about background variables that might have influenced (and explained) results. In studies with sufficient background information, the effects of gender, social environment and verbal expression were revealed. Boys perform better in spatial tasks already at age 4, largely due to social conditioning by offering them building blocks and

construction kits [25]. *Richardson*, *Jones*, *Croker* and *Brown* [38] are among the first to focus on educational implications and provide a system of tasks based on their difficulty and complexity. This model was used to construct a diagnostic test of cognitive development, based on LEGO toys. The test, validated for children aged 7-11 and adults, requires the building of a structure out of 4-11 LEGO bricks, based on isometric charts (similar to those used by IKEA for the assembly of furniture), within pre-set time limits.

There is no significant difference between the objectives and results of *spatial development projects performed in real or virtual spaces*. Both approaches are beneficial for skills enhancement and the reduction or even elimination of gender differences. Effect size is influenced by the frequency and quality of developmental programs. As the importance of digital technologies in education increases, the role of digital tools in the development of spatial skills is likely to increase as well. More authentic and research based educational methods will be employed to increase student motivation, provide individualized learning and also contribute to the development of 21st Century skills [28].

Developmental programs, however, often neglect *creative problem solving*. In contemporary art education, examinations are required to include such tasks as they are necessary for reliable and authentic assessment of competences related to creation and perception of the arts. A modernized version of the portfolio of the 19th Century art academies, the *process folio* that includes a logbook with documentation of research and variations on the theme leading to the final solution is a sensitive and flexible tool for the assessment of spatial intelligence [14]. Carefully chosen creative tasks have high intercultural validity as proven by the International Baccalaureate for the Arts [4]. A large scale Dutch-Hungarian study showed that *project based portfolio assessment* of creative work by trained jurors through detailed, illustrated evaluation criteria diminishing juror bias is a valid and reliable tool in education [18]. The reliability of assessment of visual skills may be increased through the introduction of creative design tasks and standardized competence tests at the same time [21].

Creative tasks as competence evaluation tools have a long tradition in higher education in the arts. Architectural or technical plans and artworks are used for final, summative assessment of professional skills at art and design academies and universities of technology as well. However, using open-ended, creative tasks for formative assessment regularly is far from being standard practice in engineering education as they are not considered objective measures of competence. In this paper, we show criteria for creative tasks to be used for the assessment of spatial creation and perception and prove that such tasks are not only motivating, but also valid measures of this basic for engineering work competence.

High student dropout numbers in technology education necessitated competence assessment of first and second year students at the Faculty of Architectural Studies at Szent István University, Budapest. The Mental Cutting Test was used to detect

deficiencies in this area [5] [6] [22]. Similar studies preceded remedial skills development projects at other universities where deficient spatial skills resulted in learning problems [3] [33]. All three large-scale studies revealed difficulties in the development of spatial skills of 18-23-year-olds. University students could not completely overcome learning handicaps resulting from low levels of spatial perception and analysis. Many of them had to quit studies because of their inability to imagine or create three-dimensional objects represented in two-dimensional plans and projections and vice versa. The results of several Hungarian studies correlate, with similar findings, showing the high impact of spatial skills deficiencies, in the failure in engineering studies [15] [24] [31] [42].

# 3   Objectives and Methods of the Research Project

We investigated factors influencing the development of spatial abilities with a focus on differences in spatial information processing in real and virtual spaces. In this paper we discuss the effects of methods based on the transfer effects of creative tasks in space on formal spatial representations.

Research questions:

Q1 Do creative problem solving tasks in modeling and construction result in significant development of spatial skills and are they as effective as traditional two-dimensional drawing tasks?

Q2 Is there significant difference in effect sizes between the two learning environments: modeling in real space and constructing models virtually?

Q3 Which factors influence the natural growth and educational enhancement of spatial abilities?

Q4 Which pedagogical methods and strategies are most suitable for the development of spatial abilities in the age range 18-22 years?

Related hypotheses:

H1 Gender influences spatial performance. Male students will score higher on the pre- and post-tests than female students.

H2 University level training type influences spatial performance. Civil engineering students would score higher on the pre- and post-tests than architects.

H3 Secondary school level training type influences spatial performance. Vocational secondary school graduates will outperform secondary grammar school graduates in the pre- and post-tests.

## 3.1    Sample

The experiment was performed at the Faculty of Architectural Studies at Szent István University, Budapest, with students in the first and second study year. We targeted freshmen of architecture (n=198) and civil engineering (n=74) as both professions require, besides spatial perception (observation and interpretation of spatial relations, mental processing of visualizations, etc.) also creative skills (representation of shapes, spatial organization, modeling and construction based on an analysis of interrelationships of material, structure and form, etc.). The developmental program was introduced as an elective seminar. Students of the same study years who did not attend the spatial development program constituted the control groups. Pre-tests of spatial skills showed no significant difference between the experimental and control groups.

## 3.2    Assessment Tools

Spatial development experiments usually employ formal, standardized tests only for pre- and post- testing [46]. As our program focused on the effects of open, creative tasks, we also used process folios and a jurying method described in the papers surveyed in the introductory section for this paper. In order to document phases of visuospatial information processing, we documented phases of the creative process through films and photo sequences and on-site observation by external experts. A student questionnaire yielded information on skills and activities with potential influence on the development of spatial skills like engagement in visual arts, design or sports. Members of the experimental groups filled out satisfaction surveys as well. They evaluated the tutor's and their own performance and motivating factors as well as problems with working on an open ended art project (an unusual component in engineering studies).

The learning process was supported by an e-learning environment (Moodle) where study materials were shared and discussions about tasks documented. Problems with task solution were also managed through individual mentoring in private e-mails. The analysis of questionnaires, on-site observations and documentation provided by the virtual learning focused on a wide variety of aspects of the learning process: difficulties with understanding the tasks, quality of solution versions, correction and mentoring work and student motivation.

The formal assessment of development of spatial skills of 19-21 year old students was conducted by the Hungarian Spatial Ability Test by *Séra*, *Kárpáti* and *Gulyás* [41]. This test was successfully employed in large scale assessment projects both in secondary [19] and in higher education [10]. The two identical versions of the tool can be used for pre- and post-testing purposes (Version A, with 56 items and Version B with 47 items have a Cronbach's coefficient 0.81 and 0.93, respectively). Besides task to be solved through selecting from alternative

solutions through mental operations, the test also contains drawing tasks and can be used for the assessment of "thinking by drawing".

The test measures two large skills clusters, as described by Tóth (2014): A) Basic mental operations: mental analysis (observation of hidden spatial structures); mental synthesis (compositions); and B) Complex mental operations: mental rotation and transformation and construction of mental spatial images.

The test contains both *object processing* and *spatial processing* items – components that have been identified by the Mental Imagery and Human-Computer Interaction Laboratory at Harvard University as representing the two dominant cerebral processing mechanisms of space [2] [23]. *Those who use spatial processing,* a dorsal visual pathway function, rely on spatial relations for orientation, while those who use object processing, a ventral visual pathway function, rely more on the characteristics of the objects in space (shape, color, texture and size). Tests including tasks for both orientation types are more reliable assessment instruments.

In this study, we describe the developmental program, show its effects on students of architecture and civil engineering and discuss interrelations among performance and background variables.

## 3.3 The Experimental Program

The experimental course was carried out during the Fall and Spring semesters of the academic year 2014/15. Students and tutors met thirteen times during a semester for 120-minute sessions. In the Fall semester, we formed three groups: a *control group* with students performing the usual spatial tasks of the course program only. (This program includes two-dimensional drawing tasks: perspective representation of arrangements of geometric shapes, furniture and interiors; reconstruction based on the Monge system; and structural drawing of objects.) (Figure 1)
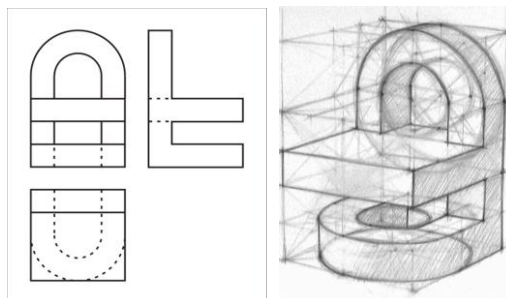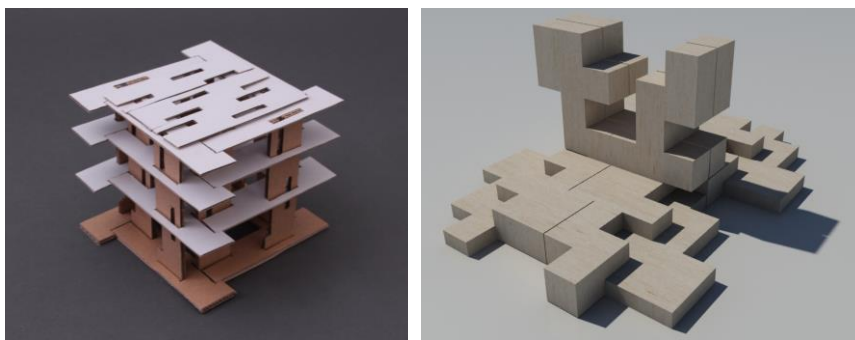


Figure 1
A two-dimensional reconstruction task for the control group studying in the traditional training program

Our *two experimental groups* worked on similar three-dimensional modeling tasks but in different design environments: in real space, using different materials (cardboard, wood, textiles, glass) for building models and in virtual space, using computer-aided design – CAD – software for modeling. Constructing tasks range from simple to complex arrangements and focus on the relations of elements, the spatial structure and transformations like flexion, truncation, and proportion changes. The investigation of visual effects like plane and spatial form characters, rhythms, transparency, physical properties of materials, light and shadow effects are the most important components. The developmental program is intended to inspire analytic thinking and representation using spatial experiences and creative imagination that also reflects the designer's identity and personality. (Figs. 2, 3)



Figures 2 and 3

Spatial modeling developed according to tutor-defined criteria by Balázs Fehér, architect student, 2015 (in real space) and by Zsolt Vittay, architect student, 2015 (in virtual space)

Table 1

Representation methods, task types, experiment periods and sample sizes

| Groups | Two- or three-dimensional type and techniques of spatial representation | | Tutor-defined / open task | Semester in the study year 2014/2015 | Sample size |
|---|---|---|---|---|---|
| Control group | 2D | drawing techniques | tutor-defined | Fall | 178 |
| 1. Experimental group | 3D | cardboard models in real space | tutor-defined | Fall and Spring | 39 |
| 2. Experimental group | 3D | computer-assisted design in virtual space | tutor-defined | Fall and Spring | 35 |
| 3. Experimental group | 3D | models in real and virtual space with self-selected materials and techniques | open tasks | Spring | 20 |

In the Spring semester, we repeated the course and also organized a *third experimental group* that worked on self-selected spatial problems through building three-dimensional models in real and virtual space. Members in this group, were also allowed to select materials and techniques for their models. (Figure 4)



Figure 4
An open, creative task of spatial modeling by Balázs Veres, architect student, 2015

We used a *process-oriented project method* in all three experimental groups, organizing work around a design problem. This approach revealed problem solving strategies that usually remain hidden when task-oriented methods are used. In Groups 1 and 2, the tutor defined spatial problems to be solved through model building, and explained the phases of modeling and the technology to be used. Spatial problems for tasks were suggested by design tutors and represented real-life issues in architecture and civil engineering. Table 1 summarizes the educational methods and other data of the experimental and control groups.

# 4   Results

Results of the pre- and post-tests show normal distribution, although some of the task groups of the Spatial Ability Test were easy for the architecture and Civil Engineering students. Pre-test average is 63.04% (skewness: -0.44, kurtosis: 0.12), post-test average is 72.01% (skewness: -0.99, kurtosis: 1.17), cf. Figures 5 and 6.



Figures 5 and 6

Distribution of results in the pre- and post-tests of the experimental and control groups (N=272)

In the Fall semester, there were no significant differences between pre-test result of the students in the experimental and control groups. ($M_{cont}$=62.16%, $M_{exp}$=64.90% t[226]=-1.28, p<0.205). During the semester, the spatial skills of all groups developed significantly, (t[178]=-7.01, p<0.001; t[48]=-9.56, p<0.001), but performance growth in the control group was only 6.48%, while that of the experimental groups was 15.44%. Post-test results of the experimental groups were significantly better ($M_{cont}$=68.63%, $M_{exp}$=80.33%; t[226]=-4.70, p<0.001). Test results are summarized on Table 2.

Table 2

Student performance on the Spatial Ability Test in the first experimental phase: a comparison of results of the experimental and control groups

| Groups | Pre-test (%) | | Post-test (%) | | Pre- and post-test (%) |
|---|---|---|---|---|---|
| | Mean | Standard deviation | Mean | Standard deviation | Paired-samples t-test |
| Control (N=178) | 62.15 | 14. 92 | 68.63 | 17.09 | t=-7.01, p<0.001 |
| Experimental (N=48) | 64.90 | 12.73 | 80.33 | 14.79 | t=-9.56, p<0.001 |
| Independent-samples  t-test | t[226]=-1.28, p<0.205 n.s.* | | t[226]=-4.70, p<0.001 | | - |

*not significant

The effect size of the developmental program was d=1.07, is more than 1 standard deviation. When we correct this value with the effect size of the development of the control group (d=0.44), the resulting effect size is still d=0.63, a high value that shows the efficacy of the developmental programs based on creative three-dimensional tasks.

On Figures 7 and 8, individual results of students are shown. Pre- and post-test scores are projected on top of each-other. Student whose performance was the same in both tests are shown on the diagonal. Above the diagonal, we can see students who performed better on the post test and under it, those who underperformed. On the left, results of members of the control group are almost evenly distributed above and below the diagonal. On the right, results of members of the experimental groups show two peculiar features. First, only those performed worse on the post-test, whose pre- and post-test results were very high. Second, there are many students with outstanding performance growth. Irrespective of higher or lower spatial skills levels at start, the experimental project resulted in substantial development.



Figures 7 and 8

Results of the experimental and control groups in the pre- and post-tests (N=272)

The most important difference among the control and the experimental groups was the method of development through two- and three-dimensional tasks respectively. In the second phase of the experiment, we wanted to find out if it was the *spatial representation modality* (dimension) of the tasks or the *spatial problems* involved (the content of the developmental programs) that resulted in significant student development. Therefore, we organized an experimental group with a new type of training program, where students were allowed to select the spatial problems they wanted to solve and also the materials and three-dimensional techniques employed. The performance of the two different types of experimental groups, are summarized on Table 3.

Table 3
Results of the tutor-defined and the open task groups in the pre- and post-tests

| Type of experimental treatment | Pre-test (%) | | Post-test (%) | | Pre- and post-test (%) |
|---|---|---|---|---|---|
| | Mean | Standard deviation | Mean | Standard deviation | Paired-samples t-test |
| Tutor-defined modeling (N=74) | 64.30 | 12.66 | 80.01 | 13.00 | t=-11.60, p<0.001 |
| Open modeling (N=20) | 66.25 | 12.70 | 72.45 | 13.63 | t=-1.82, p<0.085 n.s.* |
| Independent-samples t-test | t[94]=-0.61, p<0.542 n.s.* | | t[94]=2.22, p<0.034 | | - |

*not significant

Results of the two experimental groups showed no significant difference in the pre-test (t[94]=-0.61, p<0.542), and the distribution of results also had a similar pattern. (Levene-test=0.17, p=0.68). The average performance growth of the open modeling task group was 6.20%, but the development was not significant. The performance of the tutor-defined task group was significantly higher on the post-test, then on the pre-test (M=15.72%; t=-11.60, p<0.001). The methodology of the developmental program (task types, structure of the teaching units etc.) contributed to student growth considerably (t[94]=2.22, p<0.034). However, as the sample size for the open modeling group was relatively low (n=20), we cannot decide if student development was the result of the methodology of the training program or just the better representational functions of the 3D environment. Further research is needed to establish the significance of the learning content and teaching methods in the development of spatial abilities through creative tasks.

In the Spring semester, we repeated the tutor-defined modeling program too. Results were similar to those in the Fall semester: experimental groups scored significantly higher than control groups, and the effect size of the tutor-defined modeling program was the same in both phases.

Table 4
Results of groups working in real and virtual space in the pre- and post-tests

| Type of modeling environment | Pre-test (%) | | Post-test (%) | | Pre-and post-test (%) |
|---|---|---|---|---|---|
| | Mean | Standard deviation | Mean | Standard deviation | Paired-samples t-test |
| Material model (in real space, N=39) | 65.18 | 12.81 | 78.62 | 15.21 | t=-7.14, p<0.001 |
| Digital model (in virtual space, N=35) | 63.31 | 12.61 | 81.57 | 10.00 | t=-9.67, p<0.001 |
| Independent-samples t-test | t[74]=-0.63, p<0.531 n.s.* | | t[74]=-1.00, p<0.322 n.s.* | | - |

*not significant

In both semesters, one of the experimental groups worked in virtual space, using CAD software. On Table 4 we compare the effects of spatial skills development in real versus in virtual space. There is no significant difference in effect sizes between the two learning environments: modeling in real space with tangible materials equally enhances spatial skills as constructing models virtually.

# 5    Background Variables Affecting Student Results

In this part of the paper, the effects of three background variables included in the student questionnaire will be discussed: gender, specialization (studies for an architecture or civil engineering degree), and secondary level education in a technical vocational school or secondary grammar school.

Table 5 shows test results by gender. Research on sex differences, in the level of spatial ability; almost always indicate the supremacy of boys and men [26]. We also hypothesized that both pre- and post-test results of male students would be higher. This assumption was justified in the present study as well. Another result, however, that indicates similar degree of development for men and women shows that spatial abilities of female engineering students with much lower pre-test scores than males may equally be enhanced through creative, three-dimensional tasks.

Table 5

Results of male and female students in the pre- and post-tests of the experimental and control groups

| Gender | Pre-test (%) | | Post-test (%) | | Development Pre-and post-test (%) | |
|---|---|---|---|---|---|---|
| | Mean | Standard deviation | Mean | Standard deviation | Mean | Standard deviation |
| Male (N=161) | 65.43 | 14.71 | 73.61 | 16.54 | 8.18 | 12.74 |
| Female (N=111) | 59.57 | 12.71 | 69.69 | 16.34 | 10.13 | 13.39 |
| Independent samples t-test | t[272]=3.50, p<0.001 | | t[272]=1.93, p<0.055 n.s.* | | t[272]=-1.21, p<0.226 n.s.* | |

*not significant

Male students scored significantly higher in the pre-test, (t[272]=3.50, p<0.001), but in the post-test, their better performance is not so significant (t[272]=1.93, p<0.055). In further studies, we will compare the level of development in relation to treatment types of the spatial skills development programs for female students to increase their performance in skill components that they seem to have a genetic or cultural handicap in. Female students showed lower performance in tasks that require mental manipulations like mental rotation and transformation. They achieved similar results in the recognition and interpretation of two- and three-dimensional shapes. We could not detect gender-related differences in tasks that appear regularly in the university training programs (e.g. reconstruction).

We also investigated the effects of specialization (studies for an architecture or civil engineering degree) on the results of spatial abilities development during the experimental courses. In the B.Sc. program for Civil Engineering, Mathematics and Technology is more pronounced than in the Bachelor in Architecture program. Therefore, we presumed that results of civil engineering students would be higher on both the pre- and post-tests. Contrary to our expectations, future architects performed significantly better than future civil engineers on both tests. (Cf. Table 6 for a comparison of results). There is no significant difference in their pace of spatial skills development either. The participants of our sample were all in their first or second year of study in a four-year degree course. It would be worthwhile to test the spatial skills of graduating students as well to verify this result and redesign courses if needed.

Table 6

Results of students of Architecture and Civil Engineering in the pre- and post-tests of the experimental and control groups

| Specialization | Pre-test (%) | | Post-test (%) | | Development Pre-and post-test (%) | |
|---|---|---|---|---|---|---|
| | Mean | Standard deviation | Mean | Standard deviation | Mean | Standard deviation |
| Architect (N=198) | 66.34 | 12.76 | 76.06 | 13.03 | 9.72 | 12.46 |
| Civil engineer (N=74) | 54.20 | 14.16 | 61.19 | 19.94 | 6.99 | 14.31 |
| Independent-samples t-test | t[272]=6.46, p<0.001 | | t[272]=5.96, p<0.001 | | t[272]=-1.45, p<0.150 n.s.* | |

*not significant

University staff often report better performance of those students who are graduates of vocational secondary schools for whom it is easier to understand and solve tasks at the beginning of their training, than for students coming from secondary grammar schools. Based on this experience, we expected vocational secondary school graduates to outperform secondary grammar school graduates in the pre- and post-tests. Again, the hypothesis based on previous experiences turned out to be wrong: the type of secondary education had no effect on test results (cf. Table 7 for details). Tóth [44] compared the spatial ability of 14-18 year old students, of secondary grammar and vocational schools. Two important conclusions of the study: (1) some courses of the vocational training schools (like technical drawing) effectively develop spatial skills, (2) but the general cognitive abilities of secondary grammar school students are more advanced than those of vocational school students, and these positively affect their spatial performance. Consequently, spatial abilities of students from the two types of secondary education institutions are similar in the higher grades. It would be important to study the role of cognitive skills in future studies conducted with engineering students as well.

Table 7

Results of graduates from vocational secondary schools and secondary grammar schools in the pre-
and post-tests of the experimental and control groups

| Type of secondary school | Pre-test (%) | | Post-test (%) | | Development Pre- and post-test (%) | |
|---|---|---|---|---|---|---|
| | Mean | Standard deviation | Mean | Standard deviation | Mean | Standard deviation |
| Secondary grammar school | 63.04 | 14.35 | 71.38 | 16.06 | 8.34 | 12.50 |
| Vocational secondary school | 63.03 | 14.03 | 73.02 | 17.37 | 9.99 | 13.80 |
| Independent-samples t-test | t[272]=0.01, p<0.994 n.s.* | | t[272]=-0.78, p<0.436 n.s.* | | t[272]=-1.00, p<0.319 n.s.* | |

*not significant

Creative and constructive task types have substantially developed the visual skill components studied. These tasks were especially useful for enhancing the ability to solve mental transformations and spatial operations requiring the integration of different viewpoints (construction of mental spatial images). They were also beneficial for the development of the other spatial skill components where this treatment also resulted in significant performance increase (1 t=-2,928, p<0,005; 2 t=-5,134, p<0,001; 3 t=-5,462, p<0,001; 4 t=-7,204, p<0,001). (Figure 9)



Figure 9

Results by spatial skill components in the pre- and post-tests (N=74, 1st and 2nd experimental groups)

Female students performed lower in tasks requiring mental rotation and transformation, but in the integration of viewpoints and in the interpretation of shapes their performance was similar to those of male students. No gender-related differences were found in those tasks that are continuously present in the university training program (for example, reconstruction), - a result that indicates the possibility to reduce and eventually eliminate factors that result in weaker

spatial performance of women. Results were not influenced by response type (choosing from alternatives or drawing the correct solution). Task difficulty was inherent in mental visualisation (creation of mental images). Tasks without response options represented as images resulted in more false answers.

Table 8 shows results for two issues raised in the survey: ranking problems encountered during the creative process and in the use of visual language by difficulty level. According to the process-oriented project method, tasks were organized in a logical sequence. This arrangement supported students in the identification of design problems and the elaboration of appropriate solutions. During the design experiment we could observe how students become more and more conscientious of spatial problems, how they learn to analyze them correctly and reflect with growing self-assurance on solutions of their own as well as those of peers.

Table 8

Frequency of responses on questions about the experimental program for the development of spatial skills (experimental groups)

| Question | Items ranked | Frequency of highest ranking of this item (%) |
|---|---|---|
| What was the most difficult aspect of solving the creative task? Rank the answer choices in order of preference, 5 being the highest and 1 the lowest! | Understanding the task | 5 |
| | Generating creative, original design ideas | 26 |
| | Finding the most appropriate technology (in a 3D modeling task, selection of the most appropriate material and tool, in computer aided design tasks, the best software function, etc.) | 26 |
| | Preparation of the model (production of an appropriate number of alternatives of good quality) | 35 |
| | Corrections performed according to suggestions by the tutor | 8 |
| What were the major problems in using visual language? Rank the answer choices in order of preference, 6 being the highest and 1 the lowest! | Mental imaging and analysis of spatial arrangements and their representation | 17 |
| | Choice of color | 12 |
| | Design and interpretation of objects in space and their connections | 10 |
| | Texture design | 7 |
| | Visual effects design | 17 |
| | Composition | 37 |

When ranking tasks by difficulty, students ranked the creation of model variations in adequate number and quality first or second. Tutors reported that students were often unable to plan the stages of task completion and considered their first design

idea an accomplished task. They also found it difficult to come up with "creative and original" ideas. (In the second phase of the experiment, we divided the tasks in smaller subtasks and thus defined the process of activities in order to develop design skills more effectively).

First year students (freshmen) found the selection of appropriate techniques and processes for modeling difficult. In the use of visual language (e. g. creating different visual effects), students found tasks that required an exact interior image of the spatial situation to represent most difficult. Interpretation of spatial relations and design of visual effects were found difficult for the same reason: they both require mental modeling and complex mental manipulations (interpretation of spatial positions and directions, mental rotation and mirroring) before the completion of the task in real or virtual space. Even in the groups studied, first and second year students whose training focuses on the development high level spatial skills, it is difficult to harmonize different viewing angles. This deficit made students rank "composition" high on the list of difficult tasks.

The attitude survey of students participating in our experiment showed high level motivation for attendance in the creative modeling workshop. Students considered the enhancement of their general design competence, highly important for their future profession (and not the development of spatial skills) as their major benefit from the course. (Preference score given on a five-level Likert scale: 4,23).

**Conclusions**

Spatial ability is a basic set of skills for engineering students. It is traditionally developed through two-dimensional tasks in a geometric representational system, copying or drawing arrangements of geometric or organic objects. These activities are moderately motivating, even for art students. Therefore, our results for the successful enhancement of spatial skills, through creative modeling and construction tasks has a methodological significance. We have proven that an appropriately constructed set of tasks, supported by face-to-face and online mentoring, results in a significant developmental improvement.

Improved performance in spatial tests of the experimental groups was unrelated to gender, specialization and secondary level studies. Male, as well as, female students and future architects, as well as, engineers benefitted from the program, irrespective of being graduates of secondary vocational or grammar schools. The effect size of the development was close to the effect sizes reported in similar studies of near transfer between very similar contexts [9] [32] [37].

We found no difference between the test results of students working in real and virtual environments. Significant differences were only identified between those working with two-dimensional representations only (the control group) and the experimental groups working in real or virtual three-dimensional environments. Another influential factor was the developmental program: its quality impacted the level of the development of spatial ability. In a previous study, we tested five

art education methodologies for students aged 10-14 [17]. Spatial ability was evaluated through standardized tests and revealed the primacy of the group that engaged in building and construction of objects in real space. Those groups, whose developmental program included only two-dimensional drawing tasks or perceptual tasks (analysis of geometric and artistic representations of space), developed slower. Similar findings were reported by János Katona [20] who found modeling tasks in a virtual, three-dimensional environment more effective for developing spatial skills. These findings and our research results reported here need further confirmation in different sociocultural and educational settings. It is probable that three-dimensional tasks like modeling and construction enhance the perception, processing and interpretation of spatial relations more effectively than traditional methods of two-dimensional representation.

Satisfaction surveys conducted after the two iterations of our experimental courses revealed a substantial "motivating effect" of the creative design and construction tasks. More frequent inclusion of such tasks in engineering programs is encouraged by the research. Free experimentation with design solutions and realization of plans in three-dimensional models are skills required of Architects and Civil Engineers.

Our future research will be focused on developmental differences of spatial ability of students of architecture and civil engineering that is present at the beginning of their training. Is it still there at the end of their training, and if so, what can we do to furnish future Civil Engineers with high level spatial skills? Another research problem is manifest in the lower developmental potential of freely selected modeling tasks. Work produced in this experimental group is of high quality, still this method does not seem to result in a similar level of skills enhancement as tutor-defined tasks. A more detailed analysis of test results (for example, the analysis of the strength of item to item correlations) may reveal which skill components are developed by the open, creative tasks develop, and which may be better enhanced through tutor-defined tasks with some creative options. We intend to introduce new testing measures that may assess a wider set of spatial skill components, to see which element of spatial ability the open, creative tasks develop that standardized spatial tests eventually neglect.

Apart from tests, we also evaluated the performance of students through portfolio assessment of their collections of design alternatives and finished products. Our future work will compare spatial ability development of students in the different experimental and control groups, based on a comparison of test results and portfolio assessment results. The current study has proven the effects of near transfer of a creative program. We intend to further improve this motivating and effective methodology and make it a research-grounded option for Universities of Technology.

## References

[1]     Aden, M.: *Risiken und Nebenwirkungen einer kompetenzorientierten Kunstpädagogik. Ein kritischer Forschungsbericht*. Universität Bremen, Bremen. http://nbn-resolving.de/urn:nbn:de:gbv:46-00102369-13 Last accessed: 25 January 2016 (2011)

[2]     Blazhenkova, O. and Kozhevnikov, M.: The New Object-Spatial-Verbal Cognitive Style Model: Theory and Measurement. *Applied Cognitive Psychology*, 23(5), pp. 638-663 (2009)

[3]     Bárdné Feind, T.: *Építészhallgatók térszemléletének fejlődése és fejlesztése az Ábrázoló geometria tantárgy keretében*. (Development and fostering of spatial abilities of students of architecture in a Descriptive Geometry seminar), Doctoral dissertation. Budapest University of Technology and Economics, Faculty of Architecture (2001)

[4]     Boughton, D: *Assessment of Performance in the Visual Arts: What, How, and Why*. In: Karpati, A. and Gaul, E. (Eds.): From Child Art to Visual Culture of Youth - New Models and Tools for Assessment of Learning and Creation in Art Education. Intellect Press, Bristol (2013) pp. 119-142

[5]     Bölcskei, A., Gál-Kállay, Sz., Kovács, A. Zs. and Sörös, Cs.: Development of Spatial Abilities of Architectural and Civil Engineering Students in the Light of the Mental Cutting Test. *Journal for Geometry and Graphics*, 16(1) pp. 103-115 (2012)

[6]     Bölcskei, A., Kovács, A. Zs. and Kusar, D.: New Ideas in Scoring the Mental Rotation Test. *Ybl Journal of Built Environment*, 1(1) pp. 59-69 (2013)

[7]     Cheng, Y. L. and Mix, K. S.: Spatial Training Improves Children's Mathematics Ability. *Journal of Cognition and Development*, 15(1) pp. 2-11 (2014)

[8]     *Creativity of Schools in Europe – a Survey of Teachers*. European Commission, Paris (2009)

[9]     Csíkos, Cs., Szitányi, J. and Kelemen, R.: The Effects of Using Drawings in Developing Young Children's Mathematical Word Problem Solving: A Design Experiment with Third-Grade Hungarian Students. *Educational Studies in Mathematics*, 81(1) pp. 47-65 (2012)

[10]    Csíkos, Cs. and Kárpáti, A.: Connections between Spatial Ability and Visual Imagery Preferences. *Journal of Engineering* (Submitted)

[11]    Ferrara, K., Hirsh-Pasek, K., Newcombe, N. S., Golinkoff, R. M. and Lam, W. S.: Block Talk: Spatial Language during Block Play. *Mind, Brain, and Education*, 5(3) pp. 143-151 (2011)

[12]    Freedman, K., Hejnen, E., Kallio-Tavin, M., Kárpáti, A. and Papp, L.: Visual Culture Networks for Learning: What and How Students Learn in

Informal Visual Culture Groups. *Studies in Art Education*, 54(2) pp. 103-115 (2013)

[13]   Froebel, F.: *Die Grundsätze, der Zweck und das innere Leben der allgemeinen deutschen Erziehungsanstalt in Keilhau bei Rudolstadt*. Müller Verlag, Erfurt. (1821)

[14]   Gardner, H.: *Art Education and Human Development*. The Getty Center for Education in the Arts, Los Angeles (1990)

[15]   Güven, B. and Kosa, T.: The Effect of Dynamic Geometry Software on Student Mathematics Teachers' Spatial Visualization Skills. *Development*, 2 (3D) (2008)

[16]   Hewitt, K.: Blocks as a Tool for Learning: A Historical and Contemporary Perspective. *Young Children*, 56(1) pp. 6-14 (2001)

[17]   Kárpáti, A.: The Leonardo Program. In: Kauppinen, H. and Dicket, M. (Eds.): *International Trends in Art Education.* NAEA, Washington. pp. 82-96 (1995)

[18]   Kárpáti, A., Zempléni, A., Verhelst, N. V., Velduijzen, N. H. and Schönau, D. W.: Expert Agreement in Judging Art Projects – a Myth or Reality? *Studies in Educational Evaluation*, 24(4) pp. 385-404 (1998)

[19]   Kárpáti, A. and Gaul, E.: From Child Art to Visual Language of Youth: The Hungarian Visual Skills Assessment Study. *International Journal of Art Education*, 9(2) pp. 108-132 (2011)

[20]   Katona, J.: *A geometriai térszemlélet számítógéppel támogatott fejlesztése a műszaki felsőoktatásban* (Computer-supported Development of Geometric Perception of Space), Doctoral dissertation. Debrecen University (2012)

[21]   Koch, D. S.: *The Effects of Solid Modeling and Visualization on Technical Problem Solving*. Doctoral dissertation, Virginia Polytechnic Institute and State University (2006)

[22]   Kovács, A. Zs. and Németh, L.: *Development of Spatial Ability according to Mental Rotation Test at SKF and YBL. Ybl journal of built environment*, 2(1), pp. 18-29 (2014)

[23]   Kozhevnikov, M. and Garcia, A.: *Visual-Spatial Learning and Training in Collaborative Design in Virtual Environments.* In Wang, X. and Jen-Hung J. (Eds.): Collaborative design in virtual environments - Intelligent Systems, Control and Automation: Science and Engineering, Springer Netherlands, 48, pp. 16-27 (2011)

[24]   Leopold, C.: Geometry Education for Developing Spatial Visualisation Abilities of Engineering Students. *Journal of Polish Society for Geometry and Engineering Graphics,* 15, pp. 39-45 (2005)

[25]   Levine, S. C., Ratliff, K. R., Huttenlocher, J. and Cannon, J.: Early Puzzle Play: a Predictor of Pre-Schoolers' Spatial Transformation Skill. *Developmental psychology*, 48(2) pp. 530-542 (2012)

[26]   Linn, M. C., and Petersen, A. C.: Emergence and Characterization of Sex Differences in Spatial Ability: A Meta-Analysis. *Child development*, 56(6) pp. 1479-1498 (1985)

[27]   Mc Kim, R. H.: *Experiences in Visual Thinking*. MA: PWS Publishers, Boston (1980)

[28]   McClarty, K. L., Orr, A., Frey, P. M., Dolan, R. P., Vassileva, V. and McVay, A.: A Literature Review of Gaming in Education. *Gaming in Education.* http://www.pearsonassessments.com/research. Last accessed: 25 January 2016 (2012)

[29]   McKnight, A. and Mulligan, J.: Teaching Early Mathematics "Smarter not Harder": Using open-ended tasks to build models and construct patterns. *Australian primary Mathematics Classroom*, 15(3), pp. 4-9 (2012)

[30]   Marenholtz-Buelow, B. Von: *The Child and Child Nature*. Forgotten Books, London (2013)

[31]   Mohler, J. L., and Miller, C. L.: Improving Spatial Ability with Mentored Sketching. *Engineering Design Graphics Journal*, 72(1), pp. 19-27 (2009)

[32]   Molnár, Gy.: Playful Fostering of 6- to 8-year-old Students' Inductive Reasoning. *Thinking Skills and Creativity,* 6(2), pp. 91-99 (2011)

[33]   Nagyné Kondor, R.: *Dinamikus geometriai rendszer bevezetése a gépészmérnök hallgatók műszaki ábrázolás oktatásába* (Introduction of dynamic geometric systems in technological drawing studies of mechanical engineering students), Doctoral dissertation. Debrecen University (2007)

[34]   Nath, S. and Szűcs, D.: Construction Play and Cognitive Skills Associated with the Development of Mathematical Abilities in 7-year-old Children. *Learning and Instruction*, 32, pp. 73-80 (2014)

[35]   Newcombe, N. S.: Seeing relationships: Using Spatial Thinking to Teach Science, Mathematics, and Social Studies. *American Educator*, 37(1), pp. 26-31 (2013)

[36]   Partnership for 21. Century Learning: *What We Know About CREATIVITY*. Part of the 4Cs Research Series. http://www.p21.org/storage/documents/docs/Research/P21_4Cs_Research_ Brief_Series_-_Creativity.pdf Last accessed: 25 January 2016 (2015)

[37]   Pásztor, A.: Lehetőségek és kihívások a digitális játék alapú tanulásban: egy induktív gondolkodást fejlesztő program hatásvizsgálata (Potentials and challenges of digital game based learning: usability study of a software for the development of inductive thinking) *Magyar Pedagógia*, 114(4), pp. 281-302 (2014)

[38] Richardson, M., Jones, G., Croker, S. and Brown, S. L.: Identifying the Task Characteristics that Predict Children's Construction Task Performance. *Applied Cognitive Psychology*, 25(3), pp. 377-385 (2011)

[39] Salat, A. E. and Séra, L: A téri vizualizáció fejlesztése transzformációs geometriai feladatokkal (The development of spatial visualisation skills through transformative geometry tasks). *Magyar Pedagógia*, 102(4), pp. 459-473 (2002)

[40] Sandstrom, N. J., Kaufman, J. and Huettel, S. A.: Males and Females use Different Distal Cues in a Virtual Environment Navigation Task. *Cognitive Brain Research*, 6(4), pp. 351-360 (1998)

[41] Séra, L., Kárpáti, A. and Gulyás, J.: *A térszemlélet. A vizuális-téri képességek pszichológiája, fejlesztése és mérése* (Spatial ability. Psychology, development and assessment of visuospatial skills) Comenius Kiadó, Pécs (2002)

[42] Sorby, S. A.: Developing 3-D Spatial Visualization Skills. *Engineering Design Graphics Journal*, 63(2) pp. 21-32 (2009)

[43] Sutton, K., Heathcote, A. and Bore, M.: Measuring 3-D Understanding on the Web and in the Laboratory. *Behavior Research Methods*, 39(4) pp. 926-939 (2007)

[44] Tóth, P.: *Téri képességek, mentális műveletek* (Spatial Abilities, Mental Operations) In: Buda, A. and Kiss, E. (*Eds.*): Interdiszciplináris pedagógia és a fenntartható fejlődés: A VIII. Kiss Árpád Emlékkonferencia előadásainak szerkesztett változata (Interdisciplinary pedagogy and sustainable development) Kiss Árpád Archívum Könyvtára - DE Neveléstudományok Intézete, Debrecen (Library of the Árpád Kiss Archive – Debrecen University, Institute of Educational Research) (2014) pp. 76-91

[45] Uttal, D. H. and Cohen, C. A.: Spatial Thinking and STEM Education: When, why and how. *Psychology of learning and motivation*, 57, pp. 147-181 (2012)

[46] Uttal, D. H., Meadow, N. G., Tipton, E., Hand, L. L., Alden, A. R., Warren, C. and Newcombe, N. S.: The Malleability of Spatial Skills: a Meta-Analysis of Training Studies. *Psychological bulletin*, 139(2) pp. 352-402 (2013)

[47] Verdine, B. N., Golinkoff, R. M., Hirsh-Pasek, K., Newcombe, N. S., Filipowicz, A. T. and Chang, A.: Deconstructing Building Blocks: Pre-Schoolers' Spatial Assembly Performance Relates to Early Mathematical Skills. *Child development*, 85(3) pp. 1062-1076 (2013)

# Multidisciplinary Optimization of Journal Bearings, using a RVA Evolutionary Type Optimization Algorithm

## Ferenc János Szabó

Institute of Machine- and Product Design, University of Miskolc, Faculty of Mechanical Engineering and Informatics, Miskolc-Egyetemváros, H-3515 Miskolc, Hungary, E-mail: machszf@uni-miskolc.hu

*Abstract: In this paper, the optimum geometry of a journal bearing is calculated for minimum friction coefficient and for maximum load carrying capacity. The optimized versions can be compared, which makes it possible to draw important conclusions concerning the necessary constructional changes in journal bearings if we want to increase the load carrying capacity or to decrease the energy loss due to friction. It is also interesting to see the differences in the load carrying capacity when the friction coefficient is minimal or in the friction coefficient when the load carrying capacity is maximal. During the investigations the basic equation of the THD (Thermo-Hydrodynamic) state of hydrodynamic journal bearings is solved by using the finite difference technique, while for the optimization the RVA (Random Virus Algorithm) is used. As the result of the optimization process, the load carrying capacity can be increased by more than 28% or the friction coefficient in the oil film can be decreased by 29% compared to the starting design.*

*Keywords: friction coefficient; journal bearing; load bearing capacity; optimization; RVA*

## 1    Introduction

Hydrodynamic sliding and journal bearings are commonly used in many fields of mechanical and energy engineering [1]. The efficiency and performance of such bearings are determined by their load carrying capacity and frictional coefficient, or friction force. Decreasing the friction force in the bearings makes it easier to maintain the motion, which will decrease the energy (fuel) consumption, resulting in the possibility of significant cuts in operational costs and environmental pollution.

Finding the maximum load carrying capacity or the minimum frictional coefficient needs optimization techniques, while the presence of the lubricant (most often oil) and the effects of the temperature will enlarge the analysis process

into a multi-physics or multi-disciplinary analysis process. Therefore the whole optimization process will be an example of Multi-physics Optimization or Multidisciplinary Optimization (MDO) [2]. The disciplines involved in this complex process are: fluid flow, heat transfer, solid mechanics, elasticity and tribology. The complexity of these analysis processes makes it necessary to use several numerical methods (finite difference, finite element), which can sometimes be time consuming and takes a large amount of computing capacity. Therefore very efficient and quick optimization algorithms are needed for the Multidisciplinary Optimization of hydrodynamic bearings, in order to avoid overwhelming calculations and excessively long calculation times.

Over the last 2-3 decades, evolutionary type optimization algorithms have provided the best ways to solve MDO problems, because of their efficiency, robustness and quick convergence. The basic idea of these algorithms came from the study of the behavior and reproduction of several natural systems [3] such us genetic engineering (Genetic Algorithm GA [4]), evolution of biological populations (Evolutionary Programming EP [5], or Evolutionary Strategies ES [6]), Reproduction of Bacteria (Bacterial Foraging Algorithm, BFA [7]), behavior of natural swarms (Particle Swarm Optimization, PSO [8], or Virus-Evolutionary Particle Swarm Optimization, VEPSO [9]), behavior of animal colonies (Ant Colony Algorithm, ACA [10]), or behavior and reproduction of viruses (Random Virus Algorithm, RVA [11]).

In this paper, the Random Virus Algorithm (RVA) is used for the optimization of hydrodynamic journal bearings. For the numerical analysis of the hydrodynamic bearings in each step the finite difference technique is used. The temperature dependence of the lubricant characteristics (density, viscosity) is taken into consideration by iterative steps during the numerical solution of the governing partial differential equation. Two optimized geometries are compared: in first case, the geometry of the bearing is optimized for maximum load bearing capacity. In the second case, the bearing is optimized for minimum frictional coefficient in the lubricant film. Both optimization processes start from the same starting design and are compared each to another and to the original design. On the basis of the comparisons interesting conclusions can be drawn concerning necessary constructional and geometrical changes in order to increase the load carrying capacity of the bearing or to decrease the friction coefficient.

This paper is organized as follows: Section 2 describes the finite difference calculation used for determining the pressure distribution in the lubricant film. Section 3 shows the details of the optimization problems and RVA optimization algorithm. Section 4 gives the results of the optimization procedures and Section 5 contains conclusions.

## 2    Pressure Distribution in the Lubricant Film

The applied numerical method is applicable to any problem that can be described by linear partial differential equations [12]; in this work it is used for solving the pressure distribution $p(x,z)$ in the fluid film of hydrodynamic journal bearings, for a given gap shape function $h(x,z)$. The governing equation of this problem is the Reynolds equation:

$$\frac{\partial}{\partial x}\left(h^3 \frac{\partial p}{\partial x}\right) + \frac{\partial}{\partial z}\left(h^3 \frac{\partial p}{\partial z}\right) - 6\eta U \frac{\partial h}{\partial x} - 12\eta \frac{\partial h}{\partial t} = 0$$

(1)

In Equation (1) the relative velocity of the sliding surfaces is denoted by $U$, and $\eta$ means the absolute viscosity of the lubricant. Geometry of the bearing is shown in Fig. 1.
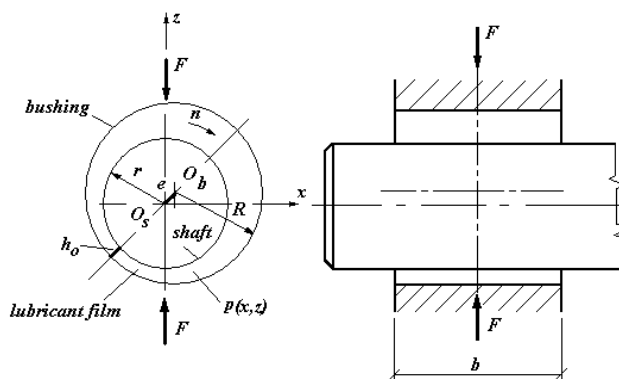


Figure 1

The geometry and dimensions of a hydrodynamic journal bearing

Equation (1) can be written into matrix form, after the discretization of the fluid film domain between the sliding surfaces, by using finite differences, as shown in equation (2). The vector **p** collects the nodal values of the pressure function, and elements of matrix **K** depend on the nodal values of the gap shape function: During the finite difference solution of the equation (1), $h_{i,j}$ represent the nodal values of the gap function and $p_{i,j}$ are the nodal values of the pressure function in the fluid film. By using this notation, the Reynolds- equation can be written in nodal points marked by i,j  [15].

In case of a finite difference mesh having  $u \: X \: v$ nodes, the matrix **K** will have a bandwidth of $2v - 3$, after the applications of the boundary conditions.

**Kp + g = 0**                                                                                                (2)

The nodal form of the Reynolds equation:

$$A_{i,j}p_{i+1,j} + B_{i,j}p_{i-1,j} + C_{i,j}p_{i,j+1} + D_{i,j}p_{i,j-1} + E_{i,j}p_{i,j} + G_{i,j} = 0$$

where

$$A_{i,j} = \frac{1}{4dx^2}\left(h_{i+1,j} - h_{i-1,j}\right) + \frac{h_{i,j}}{3dx^2} = -B_{i,j}$$

$$C_{i,j} = \frac{1}{4dz^2}\left(h_{i,j+1} - h_{i,j-1}\right) + \frac{h_{i,j}}{3dz^2} = -D_{i,j}$$

$$E_{i,j} = -\frac{2h_{i,j}}{3}\left(\frac{1}{dx^2} + \frac{1}{dz^2}\right) \quad ; \quad G_{i,j} = \frac{\eta U}{dx}\left(\frac{1}{h_{i+1,j}} + \frac{1}{h_{i-1,j}}\right)$$

The density and the viscosity of the lubricant is the function of the operating temperature of the bearing. This is taken into account by an iteration during this numerical solution. At the beginning, an approximate temperature is supposed and the equation is solved with characteristics calculated for this temperature. On the basis of the results, new and more accurate temperature value can be determined. The whole calculation will be repeated with lubricant characteristics calculated with this new temperature value. Several trial calculations and experiences show that after three or four iteration cycles, the difference between the temperature values before and after a calculation step will be smaller than 1°C, which is enough accurate for the further calculations. The elastic deformation of the shaft and housing could be checked by finite element model after the solution (quasi-TEHD state), this could be effective if these deformations are small comparing to the gap size (for example in case of steel shaft and steel bushing).

Once we have the solution of this process for the nodal values of the pressure function, the load carrying capacity of the surface pairs $F_n$ can be calculated by numerical integration, using the characteristic sizes ($r$, $R$, $b$, $h_o$, $e$) of the bearing, according to Fig. 1.

$$F_1 = \int_{\varphi=-(\beta-\varphi_1)}^{\varphi_1}\int_{z=-\frac{b}{2}}^{b/2} p\, r\, d\varphi\, dz \cos\varphi \qquad F_2 = \int_{\varphi=-(\beta-\varphi_1)}^{\varphi_1}\int_{z=-\frac{b}{2}}^{b/2} p\, r\, d\varphi\, dz \sin\varphi$$

$$ \; \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad ; $$

$$F_n = \sqrt{F_1^2 + F_2^2} \qquad\qquad\qquad\qquad\qquad\qquad\qquad . \qquad (3)$$

The friction force, which is the force needed for the relative motion between the shaft and bushing can be determined as follows:

$$F_f = \int_{x=-r(\beta-\varphi_1)}^{r\varphi_1}\int_{z=-\frac{b}{2}}^{b/2}\left(\frac{1}{2}\frac{\partial p}{\partial x}h - \eta\frac{r\omega}{h}\right)dx\, dz$$

$$ (4) $$

In equation (4) the angular velocity $\omega = 2\,\Pi\,n$, if the unit of the angular velocity is radians per seconds and the $n$ rotation speed is in rotations per seconds.

The frictional coefficient $\mu$ can be calculated as $\mu = F_f / F_n$. Lubrication angle $\beta$ is shown in Fig. 2 together with the angle $\varphi$ marking a general position of the gap function $h(\varphi)$. In general position the thickness of the lubricant film (gap function) can be calculated as:

$$h(\varphi) = R - r - e\cos\varphi \tag{5}$$

This calculation method has been verified and compared to the analytical solutions for infinite width bearings given by Szota and Döbröczöni [12], optimized for maximum load carrying capacity, and good agreement was found between the theoretical and numerical results [15]. Another verification of the method was in the case of finite sliding bearings [2] where the results of this finite difference based code were compared with those calculated by the ANSYS-FLUENT [13] program system and once again good agreement was detected.
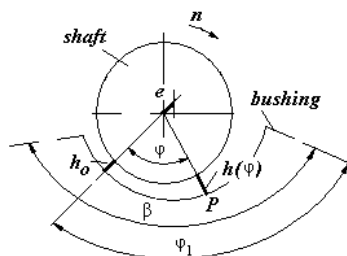


Figure 2
Characteristic angles of the bearing in general case

On the basis of these comparisons it can be concluded that the finite difference based calculation method proposed here can be applied for further investigations of THD state of hydrodynamic sliding surface pairs with finite or infinite width and for hydrodynamic sliding bearings and journal bearings. This calculation method will be integrated in this work with the RVA optimization algorithm for the optimization of a finite width hydrodynamic journal bearing. Two optimization processes will be compared: in the first case the bearing is optimized for maximum load carrying capacity (the objective function is $F_n$), in the second case, the same bearing will be optimized for minimum friction coefficient (the objective function is $\mu$) all the other parameters of the investigations will be the same. Optimal geometrical parameters are compared for this two objective functions in order to draw some useful conclusions for the manufacturers, designers or users of this type of bearings about efficient ways to increase the load carrying capacity or decrease the friction resistance of the bearing by modifying only the geometrical sizes.

# 3    Description of the Optimization Problem and the Random Virus Algorithm

As a starting design of the bearing optimization, a hydrodynamic bearing of an electric generator, is selected. The input data of the bearing are the following: input power: *P = 1300 kW*, the minimum required load bearing capacity: *F = 31400 N*, rotational speed: *n = 1000 rpm*, width to diameter ratio: *b/d = 1.3*, environmental temperature: 20°C, lubrication angle*: β = 180°*, maximum permissible operational temperature: $T_{max} = 80\ ^oC$, material of the shaft: structural steel with yield stress $R_{eH} = 275\ MPa$, material of the bushing: structural steel, with white alloy lining, maximum permissible value of the average pressure in the fluid film: 1 MPa and average surface roughness value: on the shaft - 0.16 μm, on the bushing - 0.32 μm.

The design variables are the nodal coordinates of the finite difference mesh keypoints. For the meshing 40 key nodes are used with variable coordinates (these are the optimization variables) and remaining nodes are placed depending on the keypoints in order to make higher density mesh. For the first optimization problem the objective function is the load carrying capacity of the bearing $F_n$, which is to be maximized. In the second optimization problem the friction factor $\mu$ is minimized. For both optimization problems the generator bearing is used with the given input data as the starting design.

Size constraints:    0 [mm] < $r$ < 500 [mm]   ,    0 [mm] < $R$ < 500 [mm]   ,
0 [mm] < $e$ < 10 [mm]. Implicit constraints: the pressure function should fulfill the Reynolds equation (1) of hydrodynamic surface pairs; the shaft diameter should be higher than the minimum required diameter given in equation (6); the average pressure in the fluid film should be smaller than the maximum permissible average pressure as it is shown in equation (6); and the minimum gap distance $h_o$ should be considerably higher than the sum of the maximum roughness of the surfaces. Maximum permissible operation temperature of the bearing is 80 °C.

$$r \geq 0.5\sqrt{\frac{F}{p_{adm}b/d}}\ , \qquad h_0 \geq 4.5\left(R_{a1} + R_{a2}\right), \qquad \bar{p} \leq \bar{p}_{adm} \tag{6}$$

According to the logic of the RVA optimization algorithm, the first step is to create the first (or starting) population of the possible solutions fulfilling the constraints (Fig. 3). Once the starting population has been generated, each member of the population will reproduce, creating three new members each. This process is stronger than a nuclear explosion, so in the remaining part of the optimization the selection of the best members and elimination of members without good enough objective function values will be very important. At least 60% of the new and of the total members should be eliminated after each population in order to avoid overwhelming calculations. The members that survive this strict selection procedure will give the second population. The programming

of the RVA algorithm is very simple, easy to carry out in any programming language or in macro languages of finite element program systems, if available.
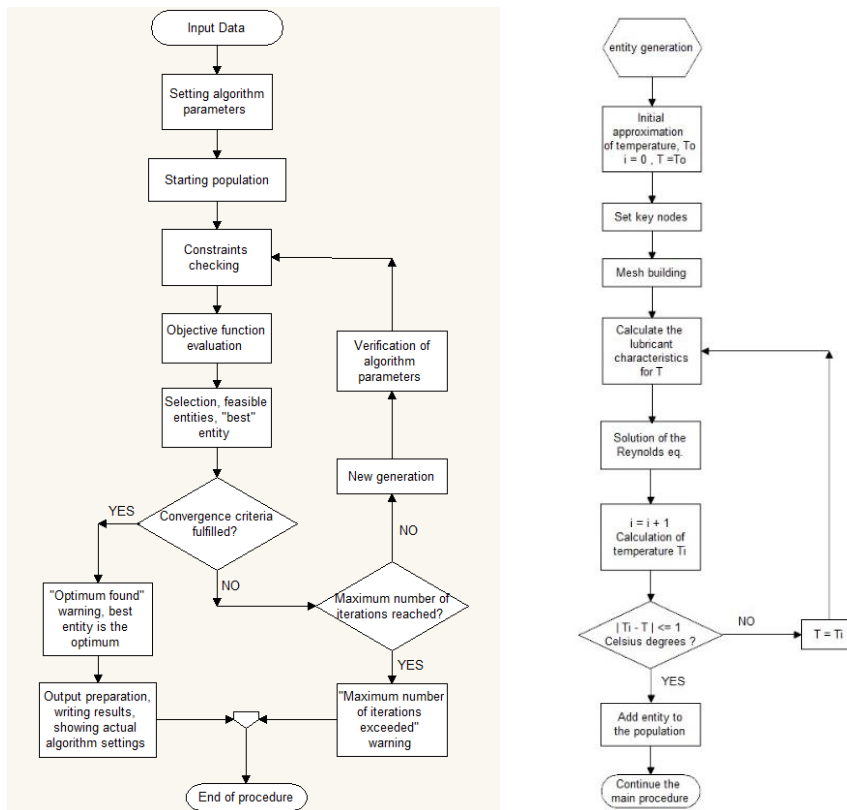


Figure 3
Flow chart of the RVA optimization algorithm and entity generation

This procedure will continue until the pre-defined optimum conditions are fulfilled. Several benchmark problem runs and numerical experiments have shown that the algorithm is very efficient: in the optimization problem investigated in this work 6 populations were enough to find the optimum. The total computation time required for a complete run was 35 minutes on an Intel core i5 desktop computer. The $j^{th}$ population: $P_j = \{x_i\}_j$ ; the reproduction formula:

$$y_k = x_k + R_k q^* \left( up_i - lw_i \right)$$  (7)

Where $y_k$ means the $k^{th}$ variable value of the new member, $q^*$ is the spreading parameter, and $R_k$ is a random number between 0 and 1, simulating the possibility of random mutations. Setting the spreading parameter properly is also very important, because it can have an important effect on the efficiency of the

algorithm. This needs a great deal of experimentation and unique fine-tuning work for each optimization problem. For this optimization process the best value for the spreading parameter was 0.8 in the first three populations and 0.25 afterwards. If the maximum number of iterations is reached without fulfilling the convergence criteria, it means that the search procedure needs more iterations and so the optimization is stopped, but during the results display a warning will say that there is a danger of a local optimum and possibly a new run will be necessary with other parameters or with a higher maximum number of iterations permitted.

# 4    Results of the Optimizations

As numerical example a hydrodynamic journal bearing of an electric generator [14] has been optimized by using the multi-disciplinary optimization (MDO) procedure described in the section 2 and 3. Two calculations have been made: in the first one the bearing is optimized for minimum friction factor in the lubricant film, which gives minimal force necessary to maintain the relative motion (turning) of the shaft. In the second study the same starting design of the bearing was optimized for the maximum load carrying capacity. The two resulting optimized version can be compared in order to draw conclusions for the further design, fabrication and operation of the bearings. Table I shows all the important parameters of the bearing, using the optimum results of the design variables for the calculation of the geometrical dimensions of the bearing. In the table it can be seen that important achievements were made as results of the optimizations: The load carrying capacity of the bearing was increased by more than 28%, while the friction factor was decreased by 29%.

Optimization results show that the increase of the load carrying capacity was realized by changing the shaft radius from 80 mm to 95 mm and changing the bushing radius from 80.13 mm to 95.16 mm. The eccentricity was increased from 79.86 μm to 101.72 μm. As the result of these changes the minimum gap $h_o$ increased from 50.54 μm to 59.77 μm. In the case of maximum load carrying capacity the average value of the pressure in the lubricant was decreased comparing to the starting design from 0.9435 MPa to 0.8630 MPa, but the friction factor remains the same, at 0.003. The temperature of the lubricant $T$ is the active constraint, 79.68$^{o}$C while the permissible temperature is 80$^{o}$C. The joint quality of the bearing remains unchanged.

Regarding the optimum results, for minimum friction factor, it can be seen in the Table I that compared to the starting design the shaft diameter remains the same, the bushing radius decreased from 80.13 mm to 80.104 mm and the eccentricity decreased from 79.86 μm to 68.90 μm. As the result of these changes, the minimum gap decreased from 50.54 μm to 35.10 μm.

Table I

Optimization results for two different objective functions

| Parameters | Starting | Min $\mu$ | Max $F_n$ |
|---|---|---|---|
| r [mm] | 80 | 80 | 95 |
| R [mm] | 80.130 | 80.104 | 95.161 |
| e [mm] | 0.0799 | 0.069 | 0.1017 |
| $\mu$ | 0.00305 | 0.002163 | 0.00305 |
| $F_n$[N] | 31400 | 31400 | 40500 |
| T [$^o$C] | 74.95 | 58.88 | 79.68 |
| $\overline{p}$ [MPa] | 0.9435 | 0.9435 | 0.8630 |
| Decrease in $\mu$ [%] | - | - 29.24 | 0 |
| Increase in $F_n$ [%] | - | 0 | + 28.98 |
| $h_o$ [μm] | 50.54 | 35.10 | 59.77 |
| Joint (ISO) | H7/a9 | H7/b8 | H7/a9 |

In the case of the minimum friction factor the load carrying capacity (31400 N) and average pressure in the lubricant (0.9435 MPa) remain the same. The temperature of the lubricant decreased to 58.88$^o$C from the original 74.95$^o$C. In this case the active constraint is the average pressure. The tolerance of the shaft is stricter (narrower) than for the starting design. The relative position of the shaft diameters and the diameters of the bushing can be compared in Fig. 4. The figure shows that in order to increase the load carrying capacity of the bearing it is necessary to increase the shaft diameter and the bushing diameter comparing to the starting design, and the eccentricity and the minimum gap size should be also increased.
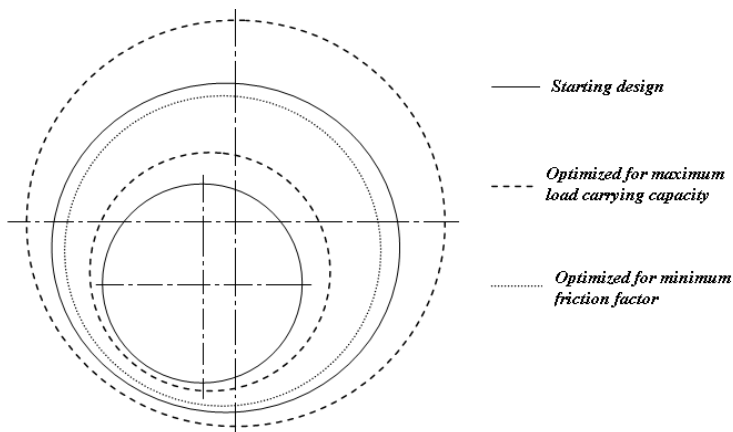


Figure 4

Schematic position of the most important diameters, comparing optimized versions for different objective functions

It can be seen from Fig. 4 that in order to decrease the friction factor, one should decrease the bushing diameter, the shaft diameter should remain unchanged, and the eccentricity and minimum gap size should be decreased compared to the starting design. Fig. 4 shows only the relative position of the diameters (it is possible to see only which one to increase, which one to decrease, and not the real dimensions, because of very small differences).

**Conclusions**

In this paper, a cylindrical hydrodynamic journal bearing (THD state) has been optimized for two different objective functions, all other parameters and constraints are the same. The starting design for the optimizations is the hydrodynamic journal bearing of an electric generator. During the first optimization, the objective function is the load carrying capacity and its maximum is determined. The second optimization study, is the minimization of the friction factor in the bearing.

The pressure distribution is determined by numerical solution of the Reynolds equation, using a finite difference computational code and the algorithm of the optimization is the RVA algorithm. During both of the optimizations the design variables are the nodal coordinates of the keypoints used for the finite difference mesh. The implicit constraints are:

- The pressure field in the lubricant film should fulfill the Reynolds equation

- Shaft diameter should be higher than the minimum necessary shaft diameter

- Maximum admissible value of the average pressure is 1MPa in the lubricant film

- Minimum gap distance between the shaft and the bushing should be higher than the sum of the maximum roughness of the surfaces

- Maximum permissible operation temperature in the lubricant is $80^{\circ}$C

- Temperature dependence of the lubricant characteristics is taken into account by an iterative process during entity generation

Final results of the optimizations show significant achievements: a 29% decrease in the friction coefficient, and a 28% increase in the load carrying capacity. The decrease in the friction coefficient can be very encouraging in terms of operation costs (since a smaller amount of energy is needed for the motion, this allows a large amount of fuel to be saved), and environmental protection (a smaller amount of fuel leads to lower levels of pollution). Higher load carrying capacity can be of interest to designers and/or manufacturers, as this can improve the market position of the factory or decrease the manufacturing costs.

The final optimal results are collected in table (Table I), showing all the important parameters of the starting design, the optimal version for the minimum friction coefficient and the optimal version for maximum load carrying capacity. The schematic position of the most imported sizes (shaft radius, bushing radius, eccentricity) can be seen in one figure (Fig. 4) together, in order to compare more easily the positions and relations of these sizes regarding all three versions (starting design, minimum friction coefficient and maximum load carrying capacity).

Comparison of the numerical results of the optimizations leads to the following conclusions:

- Increasing the load carrying capacity requires increasing all the principal sizes. As a result of these changes, the minimum gap distance will also increase. In this case the active constraint is the temperature of the lubricant, while the friction factor remains the same as it was in case of the starting design. The average value of the pressure in the lubricant decreased by approximately 10%.

- Decreasing the friction factor, requires decreasing the bushing radius and the eccentricity, while the shaft radius remains the same. The minimum gap distance also decreases. In this case the active constraint is the maximum permissible average pressure in the lubricant, and the temperature decreases by approximately 20%. The load carrying capacity remains the same as it was in the starting design.

- The changes arising from optimizations will have an effect the ISO quality of the joint between the shaft and the bushing. In the case of maximum load carrying capacity, the joint can be the same as it was in the starting design (H7/a9), but for the minimum friction coefficient it should meet higher standard: H7/b8, which will need a finer surface for the shaft.

- Regarding the manufacturing costs for the changes needed, the optimal versions, the necessary modifications for the minimum friction coefficient optimization seem to be easier and cheaper, because in this case the shaft diameter does not need to be altered, although finer surface treatment will be necessary, and the bushing diameter can be decreased slightly by methods such as the application of some coatings. The maximum load carrying capacity alternations may be more expensive, because a higher shaft diameter will be necessary (this can be realized by changing the shaft or applying a sleeve on the shaft) and a higher bushing diameter will be necessary, which may require a cutting process and could have further costs.

- It is interesting to see that the two different objective functions need changes, to the starting design, which are totally in contrast: maximum load carrying capacity requires increasing the sizes, while the minimum friction coefficient needs these parameters to be decreased. Therefore, it is advised, to consider carefully, the selection, as the objective function, in a real case.

In further investigations, more parameters could be included as design variables or objective functions (oil viscosity, surface roughness of the shaft and bushing), which will allow the calculation in a more realistic way, the costs of some changes in the design variables.

**Acknowledgements**

**References**

[1]     Kirankumar, B.M., Prajapati, J. M.: A Brief Review on Optimum Design of A Journal Bearings for I.C. Engine. *International Journal of Engineering Research Technology (IJERT),* Vol. 1. No**. 4**. 2012, pp. 1- 9. ISSN: 2278-0181.

[2]     Szabó, F. J.: Edge Shape Optimization of Finite Width Sliding Bearings. *Comput. Sci. Appl.,* Vol. 2., No **1**. 2015. pp.29- 35. USA.

[3]     Abraham, et al.: *Foundations of Computational Intelligence*, Vol. 3. Springer, 2009. 528 p. ISBN 978-3-642-01085-9

[4]     Goldberg, D. E.: *Genetic Algorithms in Search, Optimization and Machine Learning*, Addison- Wesley, Massachusetts, USA, 1989.

[5]     Fogel, L. J.*: Intelligence through Simulated Evolution: Forty Years of Evolutionary Programming*, John Wiley, Chichester, 1999.

[6]     Sheel, A.: *Betrag zur Theorie der Evolutionsstrategie*. Dissertation, TU Berlin, Germany, 1985.
ISSN: 2333-9071.

[7]     Das, S. et al.: On Stability of the Chemotactic Dynamics in Bacterial Foraging Optimization Algorithm. *IEEE Transaction on Systems, Man and Cybernetics, Part A*.: Systems and Humans, Vol. 39. Issue **3**. pp. 670-679. 2009. ISSN: 1083-4427. DOI: 10.1109/TSMCA.2008.2011474.

[8]     Eberhart, R., Kennedy, J.: New Optimizer Using Particle Swarm Theory. In: *Proceedings of VI. International Symposium on Micro Machine Human Science* 1995. pp. 39- 43.

[9]     Gao, F. et al.: Virus- Evolutionary Particle Swarm Optimization Algorithm. In: L. Jiao et al.*: ICNC 2006, Part II., LNCS 4222*, pp. 156- 165, 2006. Springer- Verlag, Berlin- Heidelberg, 2006.

[10]    Martens, D. et al.: Classification with Ant Colony Optimization. *IEEE Transactions on Evolutionary Computation*, Vol. 11. No**. 5**. pp. 651-665, 2007.

[11]    Szabó, F. J.: Multidisciplinary optimization of a structure with temperature dependent material characteristics, subjected to impact loading. *International Review of Mechanical Engineering*, 2 (**3**) pp. 499- 505. (2008).

[12]    Szota, Gy., Döbröczöni, Á.: Influence of the adhesiveness of lubricants on load carrying capacity and coefficient of friction of hydrodynamic sliding surface pairs. In: Stachowiak, G. W.(ed.): *AUSTRIB'94 International Tribology Conference*, Perth, Australia, 5-8 December, 1994. pp. 135-139.

[13]    ANSYS Inc.; SAS IP Inc. (2011): *ANSYS Mechanical APDL Technology Demonstration Guide*, Southpointe, 275 Technology Drive, Canonsburg, PA 15137, USA.

[14]    Szota, Gy.: *Design of Journal Bearings (in Hungarian)*, Műszaki Könyvkiadó, Budapest, Hungary, 1974. p. 258.

[15]    Kovács, B., Szabó, F. J.,Szota, Gy.: A generalized shape optimization procedure for the solution of linear partial differential equations with application to multidisciplinary optimization. *Structural and Multidisciplinary Optimization*, 21, No. **4**, pp.327-331, Springer, 2001.

## List of Symbols

| Name of symbol | Unit | Short description |
|---|---|---|
| $h(x,z)$, $h(\varphi)$ | [mn] | Gap function of the bearing (can be function of coordinates x, z or of angle $\varphi$ . |
| $p(x,z)$, $p(\varphi)$ | [MPa] | Pressure function in the lubricant film. |
| x, y, z | | Axis of the global coordinate system. |
| $\eta$ | [Pas] | Absolute viscosity of the lubricant. |
| U | [mm/s] | Velocity of the relative motion. |
| t | [s] | Time. |
| $F$, $F_n$ | [N] | Load of the bearing, normal load. |
| $h_o$ | [μm] | Minimum gap distance. |
| e | [mm] | Eccentricity between the shaft and the bushing. |

| n | [rpm] | Rotation speed of the shaft. |
|---|---|---|
| r | [mm] | Radius of the shaft. |
| R | [mm] | Radius of the bushing. |
| b | [mm] | Width of the bearing. |
| $O_b$ | | Center point of the bushing. |
| $O_s$ | | Center point of the shaft. |
| $h_{i,j}$ | [mm] | Nodal values of the gap function. |
| $p_{ii}$ | [MPa] | Nodal values of the pressure function. |
| **K** | [1/mm] | Coefficient matrix containing nodal values of gap function. |
| **p** | [MPa] | Vector of nodal pressure values. |
| **g** | [N/mm$^3$] | Vector of constants. |
| $A_{i,j}$ , $B_{i,j}$ , $C_{i,j}$ $D_{i,j}$ , $E_{i,j}$ , $G_{i,i}$ | | Auxiliary parameters. |
| $F_1$ | [N] | Load component in direction $\varphi = 0$ . |
| $F_2$ | [N] | Load component perpendicular to F1. |
| Name of symbol | Unit | Short description |
| | | |
| $F_f$ | [N] | Friction force. |
| $\omega$ | [rad/s] | Angular velocity. |
| $\varphi$ | [$^o$] | Angle describing the position where the gap is measured. |
| $\beta$ | [$^o$] | Lubrication angle. |
| P | [kW] | Input power. |
| $d = 2r$ | [mm] | Shaft diameter. |
| $T_{max}$ | [$^o$C] | Maximum permissible operational temperature. |
| $R_{eH}$ | [MPa] | Yield stress of the shaft material. |
| $\bar{p}_{max}$ | [MPa] | Maximum permissible value of the pressure in the lubricant film. |
| $\bar{p}$ | [MPa] | Average pressure in the lubricant film. |

$R_{a1}$, $R_{a2}$          [µm]          Average surface roughness of the shaft and bushing.

$P_j$                                  The $j^{th}$ population in the RVA algorithm.

$\{x_i\}_j$                            Variables of the $j^{th}$ member of the population

$y_k$                                  $k^{th}$ variable of the "new" member.

$x_k$                                  $k^{th}$ variable of the "old" member.

$R_k$                                  Random number having a value between 0 and 1.

$q*$                                   Spreading parameter

$up_i$                                 Upper limit for the explicit constraint of the $i^{th}$ design variable.

$lw_i$                                 Lower limit of the explicit constraint of the $i^{th}$ design variable.

$\mu$                                  Friction factor

# Prediction of Surface Roughness Parametres by New Experimentally Validated Modelling Algorithm under Abrasive Condition

## István Barányi[1], Róbert Keresztes[2], Zoltán Szakál[2], Gábor Kalácska[2]

[1] Donát Bánki Faculty of Mechanical and Safety Engineering, Óbuda University, Népszínház u. 8, H-1081 Budapest, Hungary,
E-mail: baranyi.istvan@bgk.uni-obuda.hu

[2] Faculty of Mechanical Engineering, Szent István University, Páter Károly u. 1, H-2100 Gödöllő, Hungary,
E-mail: keresztes.robert@gek.szie.hu, szakal.zoltan@gek.szie.hu, kalacska.gabor@gek.szie.hu

*Abstract: In the initial stage of the abrasive wear process, the microtopography of steels drastically changes until it reaches a stable state. This stage can be described by the changes of the 3 dimensional roughness parameters. To predict these parameters, a simulation algorithm has been developed which was validated by wear tests. The applied load and the length of the wear path varied in wide interval during the experiments. The 1.1191 steel specimen was worn using an ISO 6344:1998 abrasive paper with a mesh size of 1200. It was found that both the experimental and the simulation results can be modelled in the same manner, meaning that we use a single function involving different parameters, which depends on the problem to be solved.*

*Keywords: algorithm; surface roughness measurement; wear; abrasion; steel; design of experiment*

## 1    Intorduction

Continuously increasing industrial production demands call forth the tribological design of component surfaces. When providing technical specifications for contacting surfaces, attempts should be made early in the design phase to optimize friction losses of components and to properly define lubrication states during operation. Surface microtopography design – optimal development thereof – is one of the most dynamically developing areas of research from the tribology and manufacturing points of view [1].

Kubiak et al. [2] stated that the roughness parameter Ra of machined surfaces influenced wear intensity and the friction coefficient. In their measurements of tests involving high roughness values, they experienced decreasing friction coefficients and increasing wear, while in case of specimens of lower roughness, they experienced increasing friction coefficients and decreasing wear intensity.

Sedlacek et al. [3] examined the wear process of steel plates in case of various Ra values and manufacturing technology processes. In the course of tests, the relationship between parameters Rku, Rsk, Rpk, and Rvk, the friction coefficient, and wear was defined in case of dry and wet friction. Measurements verified the fact that roughness parameters exert an influence on the friction coefficient.

Sukumaran et al. [4] [5] emphasize the importance of online and offline measurements in their works, and they apply roughness measurements and a high-speed imaging system to describe the wear process.

Czifra et al. [6] and Palásti-Kovács et al. [7] searched for techniques to describe topography in their work in order to describe manufactured and abraded surfaces. In addition to traditionally applied parameters, their work involves evaluations based on fractal, slicing algorithm, peak curvature, and steepness distribution and they examine the applicability of each technique.

Reizer et al. [8] develop a model based on topography measurements to describe wear on honed surfaces, then they certify their simulation results by a model experiment and propose a topography to be realized by the honing technique.

Based on the literature it can be stated that the phase of initial wear, to be characterized by changes in microtopography parameters applied in roughness measurements, is not fully elaborated in the literature. Therefore, a new special truncation algorithm has been developed to model the wear process, validated by measurement results. The model can be used for forecasting changes in roughness parameters Sa, Sq, Ssk, and Sku subject to given abrasion conditions until the total destruction of the microtopography produced.

# 2    Experimental

## 2.1    Truncation Algorithm

In case of microtopographies, the truncation algorithm can be used for examining the progressive destruction of the peak zone by roughness parameters. After reading in the machined surface, the algorithm developed by the first author (Figure 1) specifies the material volume to be removed for the total destruction of the microtopography.
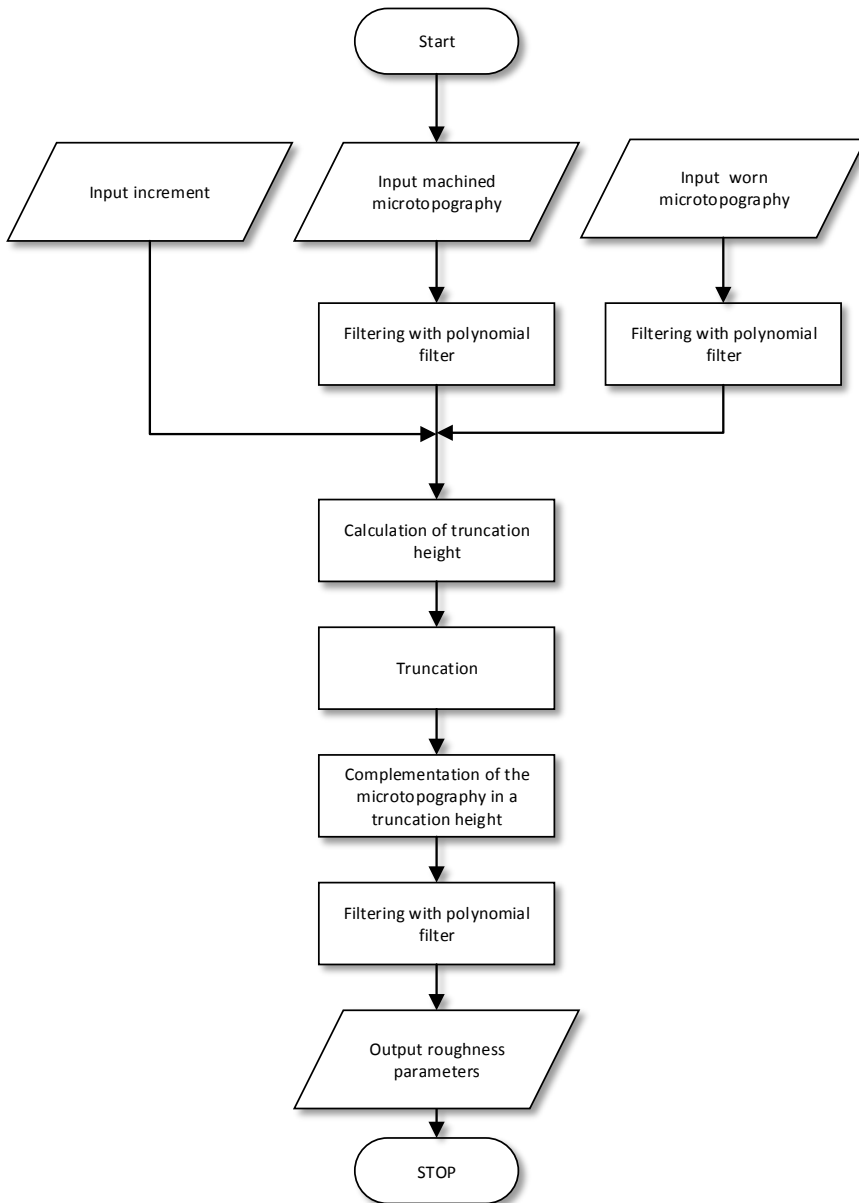
Figure 1
The flowchart of the truncation algorithm

Then the simulation model determines the truncation heights. The truncation heights can be determined by two method (Figure 2):

- according to linear scale,
- according to Abbott – Firestone curve.

Using the second method is a truncation can be implemented where the removed material volume fractions are constant in each step.
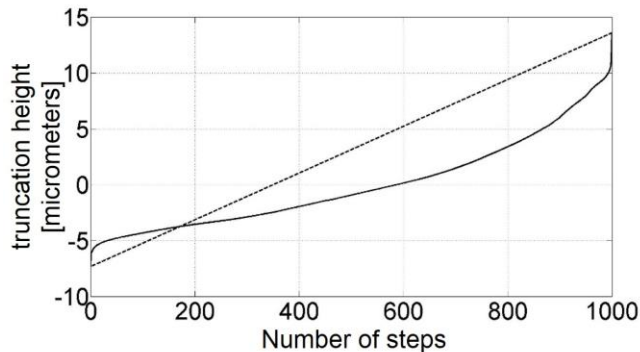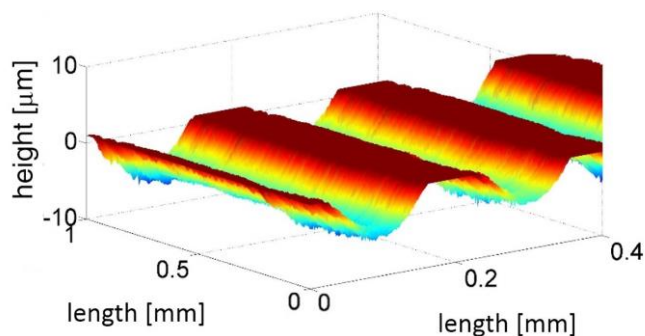


Figure 2

The relationship between the truncation step and the truncation height in the case of linear and constant removed material ratio method

The points removed cause ruptures in the microtopography. These ruptures are substituted by plane or by fully destroyed microtopography points (Figure 3).

In the case of plane patching the results of the algorithm cannot be used perfectly, because the end of the running in stage wear process the Sa and Sq parameters tends to 0 micron, the Ssk parameter trends to minus infinity, and the Sku parameter trends to minus infinite value in the case of oriented microtopography.
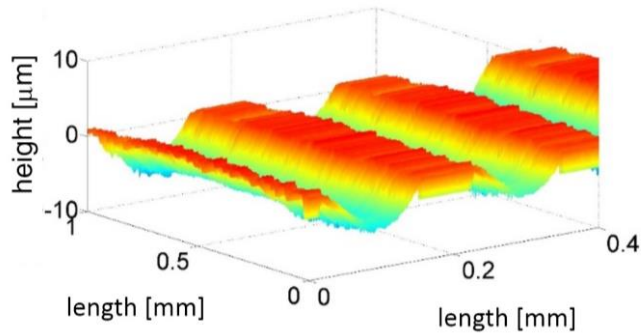
Figure 3

Microtopography patching at a given height by plane and abrasion scratches

The "virtual" microtopography (which generated from measurement results, truncated by Abbott – Firestone curve and patched by abrasion scratches) was filtered by a first-degree polynomial in both directions, retaining both the waviness and roughness characteristics of the measurements.

Afterwards, the results of the algorithm were evaluated using the parameters Sa, Sq, Ssk, and Sku in function of the destroyed volume portion, using equations (1), (2), (3), and (4).

$$Sa = \frac{1}{MN} \sum_{j=1}^{M} \sum_{i=1}^{N} |z(x, y)| \tag{1}$$

$$Sq = \sqrt{\frac{1}{MN} \sum_{j=1}^{M} \sum_{i=1}^{N} (z(x, y))^2} \tag{2}$$

$$Ssk = \frac{1}{MNS_q^3} \sum_{j=1}^{M} \sum_{i=1}^{N} (z(x, y))^3 \tag{3}$$

$$Sku = \frac{1}{MNS_q^4} \sum_{j=1}^{M} \sum_{i=1}^{N} (z(x, y))^4 \tag{4}$$
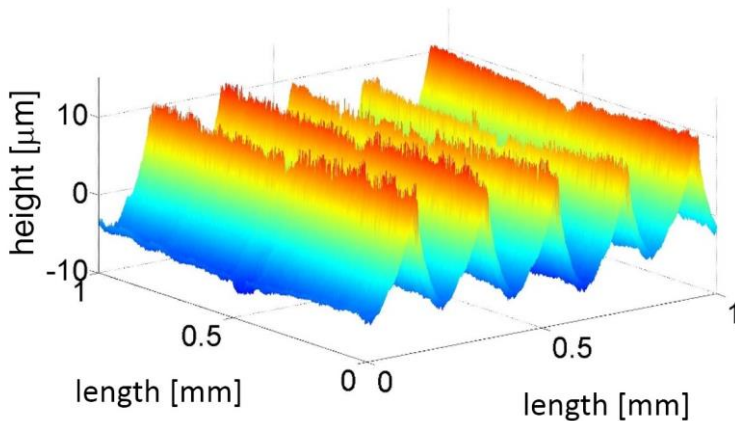
## 2.2   Test Pieces and the Wear Experiment

During the tests – where the initial stage of wear was described from several aspects – 1.0503 steel test pieces (supplier: BÖHLER-UDDEHOLM Hungary Kft.) were used in normalized heat treatment state. A revolving knife technology (rpm: 400, feed 0.2 mm/revolution, flan ID: DCMT 070202-HMP, supplier: Kolroy Inc.) was applied for surface machining. The average general roughness of

the manufactured surface was Sa=3.2 micrometers; the geometrical average of deviations was Sq=4 micrometers. The average values of skewness and kurtosis parameters to characterize the distribution of points in the direction of the vertical axis were Ssk=0.739 and Sku = 2.72, respectively.

Afterwards the tribology test specimens produced were destroyed in a pin-on-plate arrangement in the first step, at 150 mm length and 25 mm/s velocity, with a contact surface of 30 mmx30mm, subject to 600 N normal force, using an abrasion cloth of 1200 fineness (type: CK721X, procurement: Fk-Technika Kft.) perpendicularly to microtopography orientation. The end of the process was defined at 10,800 mm wear route length by the total destruction of roughness valleys located in the waviness valleys produced during manufacturing. In the second half of the experiments, the force compressing the surfaces was changed between 200 N and 600 N by 100N steps; the wear route length was increased by 600 mm steps up to 4800 mm, and then by 1200 mm steps to 10,800 mm, subject to 25 mm/s wear velocity.

The surface microtopography of test specimens was measured after each test using a Mahr stylus instrument for roughness measurement. A FRW750 instrument was applied for measurements, with a nose angle of 90°, and a rounding radius of 5 micrometers. Measurements were performed on a 1mm x 1mm surface, with a 2 micrometer step in each direction. The results yielded were taken into consideration without filtering in the course of further evaluation (Figure 4).
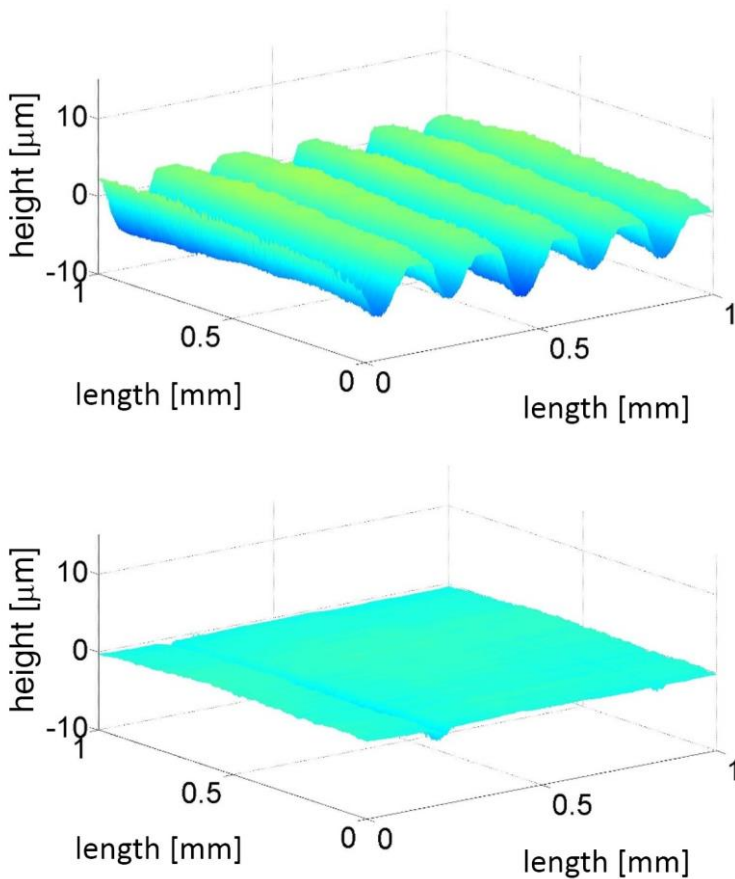
Figure 4
Surface microtopography in three different states: manufactured, abraded with a force of 200 N along 4800 mm, abraded with a force of 600 N along 10,800 mm

# 3    Results

## 3.1    Simulation Model Results

When using the simulation model, the microtopography was divided into 100 parts in the height direction where the volume of the material quantity removed was constant. Figure 5 shows changes in the parameters defined by the algorithm.
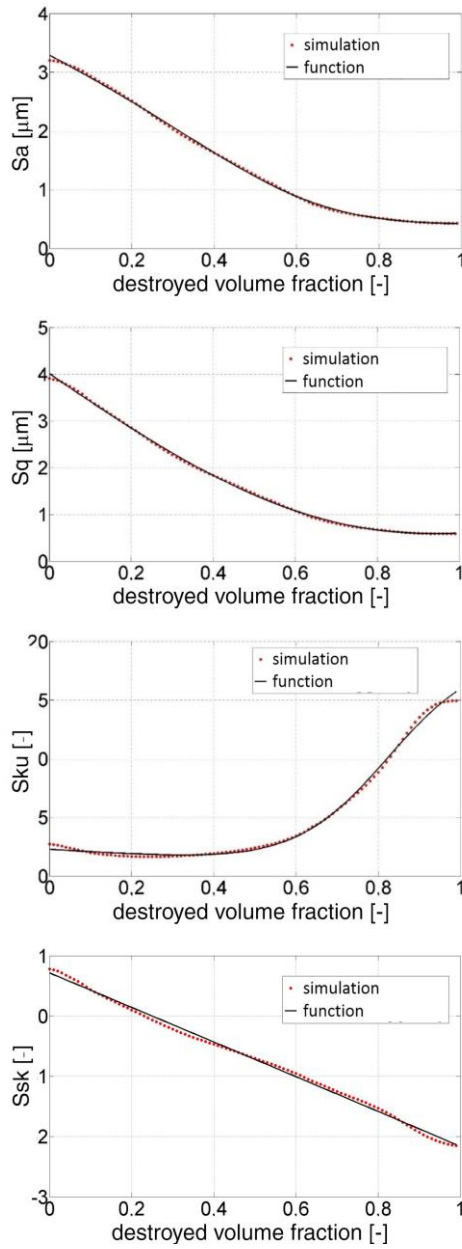
Figure 5
Changes in parameters Sa, Sq, Sku, and Ssk as a function of the volume portion destroyed

For the mathematical description of the results yielded, a modified logistic function (5) was used in the form below by white box method [9]:

$$f = \frac{ax+b}{1+e^{cx+d}} + e \tag{5}$$

Table 1 summarizes the results of functional approximations and the value of variance.

Table 1
The coefficints of the logistic functions

|                  | a      | b      | c      | d      | f     | $R^2$  |
|------------------|--------|--------|--------|--------|-------|--------|
| Sa [mikrometers] | -3.271 | 2.954  | 5.827  | -3.273 | 0.444 | 0.9995 |
| Sq [mikrometers] | -6.086 | 2.176  | 2.875  | -2.321 | 2.031 | 0.9995 |
| Sku [-]          | -2.137 | -16.88 | 8.75   | -7.144 | 19.14 | 0.9973 |
| Ssk [-]          | -6.58  | -1.204 | -0.011 | 0.2644 | 1.24  | 0.9962 |

## 3.2   Test Plan Results

A full factor test series was performed in order to validate results. Experiments were intended to answer the question whether the truncation algorithm represented a suitable approximation to the results yielded by model experiments. In order to compare results, the modified logistic function was extended into three dimensions in the following form:

$$f(x, y) = \frac{(ax+b)}{1+e^{(cy+d)}} + \frac{(fy+g)}{1+e^{(hx+i)}} + j \tag{7}$$

The equation form was used in accordance with equation (8) for function approximation in respect of force, course, and roughness parameters:

$$S = \frac{\left(\dfrac{a \cdot F}{1000}+b\right)}{1+e^{\left(\frac{c \cdot s}{1000}+d\right)}} + \frac{\left(\dfrac{f \cdot s}{1000}+g\right)}{1+e^{\left(\frac{h \cdot F}{1000}+i\right)}} + const \tag{8}$$

If the value of force (F) is substituted into the equation in [N], and course (s) in [mm], then the unit of measurement of roughness parameters depending on amplitude will be micrometers, and statistical parameters will be dimensionless. The constant value of equation (8) defines the displacement in the direction of the vertical axis applied in function approximation, its value equals to the value of totally destroyed microtopography.

The function approximation of measurement results is shown in Figure 6 and the calculated coefficients summarized in Table 2.
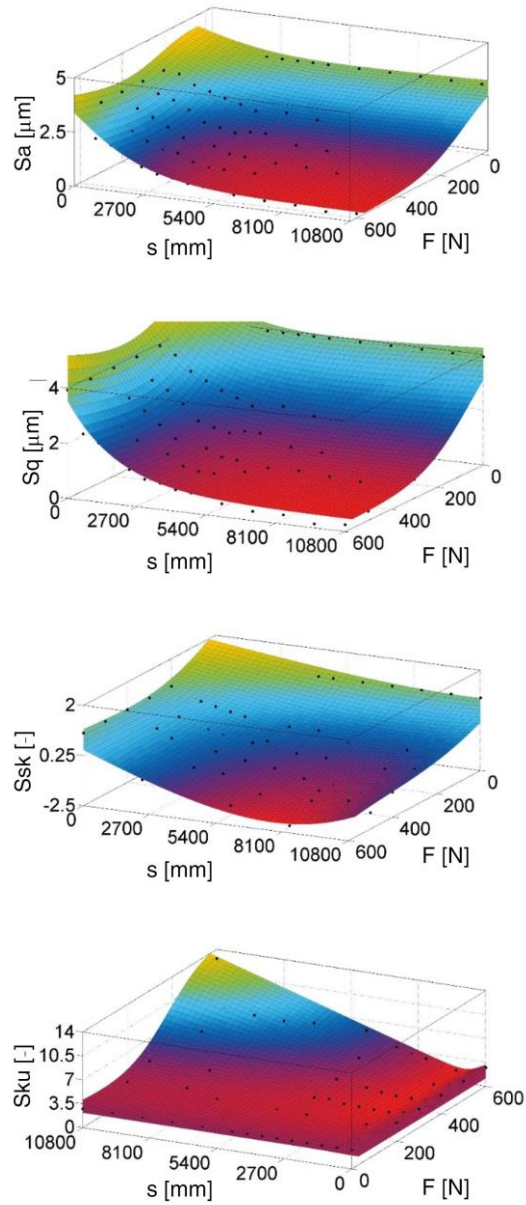
Figure 6
Changes in parameters Sa, Sq, Sku, and Ssk as a function of normal force and wear route length

Table 2

Logistic function coefficients specified from the test plan

|                  | a      | b       | c      | d      | f       |
|------------------|--------|---------|--------|--------|---------|
| Sa [mikrometers] | 7.375  | 2.601   | 0.6106 | 0.3725 | 0.03795 |
| Sq [mikrometers] | 72.5   | 31.81   | 0.5326 | 3.044  | 0.06937 |
| Ssk [-]          | -19.42 | -3.657  | 0.3096 | -4.58  | -0.7837 |
| Sku [-]          | -6.063 | -0.0491 | -3.14  | 1.418  | 1.286   |

|                  | g      | h      | i      | const  | $R^2$  |
|------------------|--------|--------|--------|--------|--------|
| Sa [mikrometers] | 2.875  | 10.02  | -1.516 | 0.4294 | 0.9397 |
| Sq [mikrometers] | 3.524  | 10.01  | -1.224 | 0.5928 | 0.9465 |
| Ssk [-]          | 27.58  | -2.853 | 1.359  | -0.715 | 0.852  |
| Sku [-]          | 0.8471 | -10.8  | 4.791  | 2.7376 | 0.9398 |

**Conclusions**

Using the algorithm developed, a simulation model was developed where – by defining the relationship between slicing height and roughness parameters – the impact of abrasion scratches in the peak zone can also be taken into consideration.

A modified logistic function was used to define the connection of amplitude-dependent roughness parameters with the normal force and the wear route length, under given measurement conditions. Within the abrasion wear model system, the coefficients of the 9-parameter function were determined for roughness index numbers Sa, Sq, Ssk, and Sku in accordance with the full factor test plan.

Based on the results yielded by the algorithm and measurements, respectively, it can be stated that the physical process occurs in accordance with a function of identical shape.

**References**

[1]     Bruzzone, A. A. G., Costa, H. L., Lonardo, P. M. and Lucca, D. A.: Advances in Engineered Surfaces for Functional Performance, CIRP Annals – Manufacturing Technololgy 57, 2008, pp. 750-769

[2]     Kubiak, K. J., Liskiewicz, T. W. and Mathia, T. G.: Surface Morphology in Engineering Applications: Influence of Roughness on Sliding and Wear in Dry Fretting, Tribology International. 44, 2011, pp. 1427-1432

[3]     Sedlaček, M., Podgornik, B. and Vižintin, J.: Influence of Surface Preparation on Roughness Parameters, Friction and Wear, Wear 266, 2009, pp. 482-487

[4]    J. Sukumaran, M. Ando, P. De Baets, V. Rodriguez, L. Szabadi, G. Kalacska, V. Paepegem: Modelling Gear Contact with Twin-Disc Setup, Tribology International. 49, 2012, pp. 1-7

[5]    J. Sukumaran, S. Soleimani, P. De Baets, V. Rodriguez, K. Douterloigne, W. Philips, M. Ando: High-Speed Imaging for Online Micrographs of Polymer Composites in Tribological Investigation, Wear 296, 2012, pp. 702-712

[6]    Czifra, Á., Goda, T. & Garbayo, E.: Surface Characterisation by Parameter-based Technique, Slicing Method and PSD Analysis, Measurement 44, 2011, pp. 906-916

[7]    Palásti-Kovács, B., Néder, Z., Czifra, Á., Váradi, K.: Microtopography Changes in Wear Process, Acta Polytechnica Hungarica 1, 2004, pp. 108-119

[8]    Reizer, R., Pawlus, P., Galda, L., Grabon, W. & Dzierwa, A.: Modeling of Worn Surface Topography Formed in a Low Wear Process, Wear 278-279, 2012, pp. 94-100

[9]    Pokorádi L: Introduction to Mathematichal Diagnostics I.: Theoretical Background, Debreceni Műszaki Közlemények VI: (1), 2007, pp. 65-80

# 2016 Reviewers

Abonyi, János
Ádám, Norbert
Alexik, Mikulas
Ancza, Erzsébet
Ando, Matyas
Anh, Anh Le
Antal, Margit
Babic, Frantisek
Balas, Valentina Emilia
Bartak, Roman
Batyrshin, Ildar
Benczúr, András
Beneda, Károly
Benesova, Wanda
Berenguel, Manuel
Biro, Miklos
Bitó, János
Bobak, Martin
Borsa, Judit
Brown, John N. A.
Bucko, Jozef
Bukovics, István
Bundzel, Marek
Butka, Peter
Búza, Antal
Cádrik, Tomáš
Csabai, István
Csákány, Rita
Csernátonyi, Zoltán
Czifra, Árpád
Do, Tien V.
Dobos, László
Dobránszky, János
Dömötör, Ferenc
Dragan, Antic
Drexler, Dániel András

Drotár, Peter
Duchon, Frantisek
Eigner, György
Ekler, Péter
Elmenreich, Wilfried
Ferenci, Tamás
Firstner, István
Fogarassy-Vathy, Ágnes
Fridli, Sándor
Fullér, Róbert
Gazda, Juraj
Giang, Nguyen
Haraksim, Rudolf
Holik, Ildikó
Horváth, Csaba
Horváth, Richárd
Hreno, Jan
Hudec, Ladislav
Jadlovská, Slávka
Jeszenszky, Peter
Józsa, Lajos
Jung, Jason
Kádár, István
Kajtár, László
Kaptay, George
Karlovitz, Tibor János
Kasanicky, Tomas
Kiss, Rita
Klešč, Marián
Koncsik, Zsuzsanna
Kósi, Krisztián
Kovács, Attila
Kovács, László
Kovács, Levente
Kovács, Szilveszter
Kovács-Coskun, Tünde

Kovarova, Alena
Krajci, Stanislav
Krómer, István
Kucera, Markus
Kuzsella, László
Laki, Sándor
Lazányi, Kornélia
Liščinský, Pavol
Luca, Mihaela
Machado, José
Markopoulos, Angelos
Márkus, Ferenc
Maros-Berkes, Mária
Marosi, Ildikó
Mátray, Péter
Mertinger, Valéria
Michelberger, Pál
Micsik, Andras
Milosavljevic, Cedomir
Molnár, Bálint
Morva, György
Nagy, Dénes Ákos
Nagy, István
Ocelíková, Eva
Palik, Mátyás
Palúch, Stanislav
Paralic, Jan
Piazolo, Marc
Pintér, István
Plavka, Ján
Pleva, Matus
Pokorádi, László
Porkoláb, Zoltán
Precup, Radu-Emil
Puheim, Michal
Rácz, Ervin
Radac, Mircea-Bogdan
Reicher, Regina
Restas, Agoston

Réveszová, Libuša
Sali, Attila
Sallai, Gyula
Sánta, Imre
Sarnovsky, Martin
Semanisin, Gabriel
Sergyán, Szabolcs
Sidló, Csaba István
Simonak, Slavomir
Simsik, Dusan
Škrinárová, Jarmila
Sobota, Branislav
Son, Hua Nam
Szabó, Imre
Szabó, József Zoltán
Szabó, László Zsolt
Szabó, Tamás
Szénási, Sándor
Szlivka, Ferenc
Szűcs, Gábor
Takács-György, Katalin
Tar, József
Tejfel, Máté
Tomoriova, Beata
Tóth, Péter
Tóth-Laufer, Edit
Trinh, Anh Tuan
Trohák, Attila
Uj, József
Vajk, István
Vámossy, Zoltán
Varga, Béla
Varga, Péter János
Varga, Viorica
Vas, László
Vassileva, Bistra
Vilanova, Ramon
Windisch, Gergely
Zdešar, Andrej