

Effects of Leading Edge Enlargement on the Primary Vortex of Blunt-Edged Delta Wing VFE-2 Profile

Mazuriah Said, Shabudin Mat, Shuhaimi Mansor and Airi Ali

Department of Aeronautics, Automotive and Ocean Engineering, School of Mechanical Engineering, Faculty of Engineering, Universiti Teknologi Malaysia, 81310 Skudai, Johor, Malaysia

mazuriah2@graduate.utm.my, shabudin@fkm.utm.my,
shuhaimi@mail.fkm.utm.my, airi.ali@utm.my

Abstract: This paper presents the results obtained from wind tunnel experiments on VFE-2 wing profile model which are differentiated by their leading edge profiles; medium- and large-edged. VFE-2 was established to investigate the effects of Reynolds number, angle of attack, Mach number and leading edge bluntness on vortex properties above-blunt-edged delta wing. The original VFE-2 wing has 4 sets of interchangeable leading edge profile namely sharp, small, medium and large-edge ratio. There were lot of experiments and simulations data in VFE that compares sharp-edged with the medium-edged wings within the VFE campaign. This paper presents the current data on a blunter wing or large-edged wing. These experiments were conducted at UTM - Low Speed Wind Tunnel, Aerolab. The experiments were carried out at speed of 18, 36 and 54 m/s representing Reynolds numbers of 1×10^6 , 2×10^6 & 3×10^6 . Two measurement techniques were employed on the wing, i.e. steady balance and surface pressure measurements. The results obtained from the large-edged wing were compared with the results from medium-edge wing. The results showed that the primary vortex depends on the leading edge bluntness, angle of attack and Reynolds number. The results obtained from steady balance data showed that lift coefficient is sensitive to leading edge bluntness at higher Reynolds numbers. Several important observations were noted on the large-edged wing; i.e. the development of primary vortex has been delayed and the vortex breakdown occurred further aft of the wing. The data obtained provide a better insight into the leading edge effect on the delta-shaped wing and also for the development of Unmanned Combat Air Vehicle (UCAV) which most of them are integrated with delta wing technology.

Keywords: Delta wing; VFE-2 profile; Vortex; blunt leading edge; Reynolds number

1 Blunt-edged Flow Topology

The exploitation of vortex lift on delta wing existed since 1940's [1]. Since then, there are researches that investigate the vortex flows above sharp-edge delta wing [2-5]. On sharp-edged delta wing; primary separation takes place when a stable shear layer is formed from a series of small vortices that shed in the leading edge of the wing. These shear layers form curling up over the wing upper surface into concentrated vortices in a spiral fashion [2-4]. The *primary vortex* is generated and initiated from the wing apex and it grows in strength and size extended towards the wing trailing edge. Underneath this primary vortex, the adverse pressure gradient increases in the region and another spinning vortex is developed in the leading edge. This vortex is called the *secondary vortex*, which rotates in the opposite direction of the primary vortex. [2].

The flow on blunt-edged delta is different from the flow formed on the sharp-edged delta wing. Firstly, the flow separation does not happen in the apex region. The flow is attached to the surface of the wing in a certain chord-wise position. The primary vortex is then developed further aft of the wing that is based on a Reynolds number, angle of attack and leading-edge bluntness [6-8]. This shows that the onset of leading-edge separation was a function of flow conditions such as Reynolds number, Mach number, blockage factors and wing geometry [9]. Another important flow phenomenon is that the primary separation line no longer occurs at the leading edge but somewhere in vicinity of it [6]. This causes the flow on the blunt-edged wing to be complicated and unpredictable.

Therefore, a research group has been established across Europe and USA to further investigate flow phenomena on a blunt-edged delta wing. This group is called as VFE-2 or International Vortex Flow Experiment 2 under AVT-113. The group has the objective to compare the results obtained from numerical calculations with wind tunnel experiments [9]. This group has used the original Chu and Luckring [10] model tested in NASA NTF shown in Figure 1(a) as a generic profile. The NASA original model has a flat plate in the middle with 4 sets of interchangeable leading edge profiles namely as the sharp-edged, small-edged, medium-edged and large-edged. These leading edge profiles were differentiated by its leading edge radii to the wing chord ratio; i.e. 0 for sharp, 0.05 for small-edged, 0.15 for medium-edged and 0.30 for large-edged wing as shown in Figure 1(b).

During the VFE-2 campaign, only the second wing or medium-edged wing was selected for further experiments in several wind tunnels such as Glasgow University [11-13], Tubitak Sage [14], Munich Technical University [15], ONERA [16] and several other wind tunnel facilities. The main objective of the campaign at that time was to study how either the numerical analysis or CFD can well predict the flow on the blunt-edged delta wing. The results obtained from the wind tunnel experiments were compared with the results obtained from Numerical analysis.

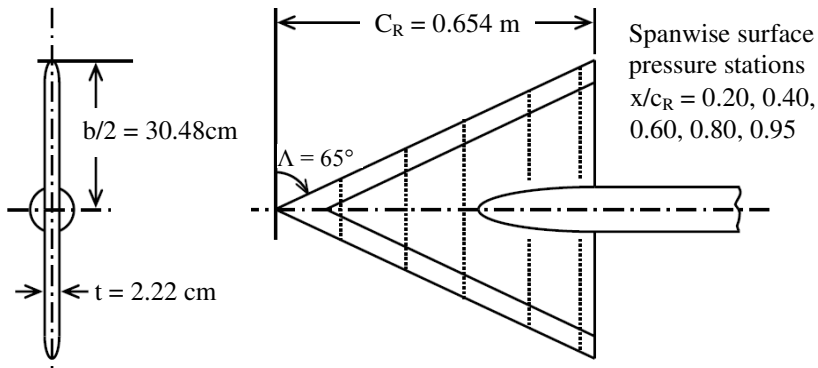


Figure 1(a)

Original NASA and VFE-2 configuration [9]

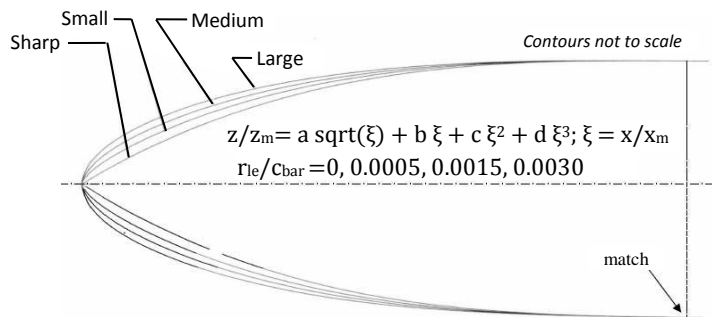


Figure 1(b)

Streamwise leading-edge contours of original NASA and VFE-2 configuration [8]

The sample results taken from VFE-2 [17] campaign of the medium-edged wing are shown in Figure 2 below. The flow on the medium-edge wing is covered by a non-separated flow on the entire wing at a low angle of attack. However, it is unclear whether the vortex is developed further aft of the wing in the trailing edge. No data is available to date [13]. When the attack angle is increased, the primary vortex is formed at a certain chordwise position from the apex as shown in Figure 2 below. From the figure, it can be observed the wing has been covered by two main sections, i.e. the attached flow and the primary vortex. The primary vortex moved forward or backward depending on the angle of attack, the Reynolds number and also the leading edge bluntness. Increasing in angle of attack has caused the primary vortex to move forward; there is no data that can indicate the primary vortex is formed in the apex region if the angle of attack continues to be increased to more than $\alpha = 25^\circ$ to date.

Reduction the Reynolds number has caused the primary vortex to move forward as shown in Figure 2. The comparison here was made between the results at R_{mac}

of 3×10^6 and R_{mac} of 2×10^6 at constant α of 13° and Mach number of 0.4. It cannot be confirm, also, whether the primary vortex will develop in the apex region if the Reynolds number is further decreased. Another factor that influenced the flow is the leading edge bluntness itself; an increase in leading edge bluntness has caused the primary vortex to be delayed. However the data on the blunter wing of large-edged is still limited [8-19]. Another important observation that has been found in the VFE-2 group was that they found another vortex formed inboard of the wing. This vortex is named *inner vortex*. The formation of this vortex also depends on leading edge bluntness, Reynolds number and also the angle of attack. More experiments are necessary to study this vortex at a higher leading-edge radius [13, 15, 17-18, 20].

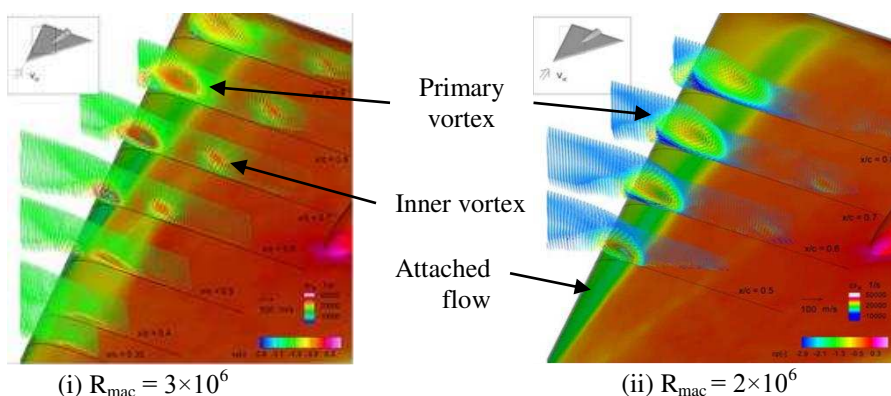


Figure 2

Pressure distributions on VFE-2 configuration at $\alpha = 13^\circ$, $M = 0.4$ on medium-edged wing at (i) R_{mac} of 3×10^6 and (ii) at R_{mac} of 2×10^6 [17]

A research group has been established in UTM to further investigate the influences of leading edge bluntness and Reynolds number on the VFE-2 model. Since the VFE-2 research group has focused on the Medium-edged wing, the team in UTM has decided to focus more on a blunter wing with a large-edged profile. The main purpose of conducting the experiment on the large-edged wing was to further investigate the characteristics of the primary vortex and vortex breakdown at higher leading edge bluntness. In addition, the surface pressure data on this wing is very limited as the VFE-2 group has only focused on the medium-edged wing. In this current paper, the experiments performed at Reynolds number varies from 1×10^6 to 3×10^6 where the flow is strongly influenced by laminar, transition and turbulence. Current data such as drag and detailed surface pressure measurement obtained from the large-edged wing were compared with the medium-edge wing. Therefore, this paper presents the flow characteristics of VFE-2 profile when the leading edge bluntness is increased. Some interesting data will be discussed in the next sub section.

2 Experimental Tests Set-Up

2.1 UTM Aerolab

The experiments were conducted in a closed-circuit UTM-LST wind tunnel facility in Aerolab (refer Figure 5). The dimension of the test section is 2.0 m (W) \times 1.5 m (H) \times 5.8 m (L) with maximum speed of 80 m/s . The average turbulence intensity at the centre of the test section is 0.06% measured at 40 m/s . The boundary layer thickness is about 40 mm at a speed of 40 m/s . The facility was equipped with 3–strut–support system located underneath the test section.

2.2 UTM VFE-2 Model

The original 65° swept angle NASA delta wing model tested in NASA [9] or called as VFE-2 configuration in AVT-113 campaign has been replicated and machined again in UTM under the Malaysian Ministry of Education Research Grant for further experiments at lower Reynolds number. The original NASA model has 4 sets of interchangeable leading edges namely as sharp, small, medium and large radius wing that corresponding to the ratio of leading-edge radii to mean aerodynamics chords r_{LE} of 0, 0.05, 0.15 and 0.3 respectively. In UTM, only two blunter wings, namely medium and large radius wings were built for further experiments. This model is named as UTM VFE-2. The model has a root chord length of, $C_R = 1.311\text{ m}$. The size of the UTM VFE-2 model is 2 times bigger than the original NASA model. This is done in order to get a high Reynolds number ($R_{mac} = 3 \times 10^6$) in a subsonic wind tunnel. The original NASA model was tested in the transonic wind tunnel. The final dimensions of UTM-VFE-2 model and the contours of both leading edges are shown in Figure 3. The UTM VFE-2 model has been machined from three main components. The first component is called a flat-plate delta shaped with fix sharp trailing edge portion. The second components are the leading edges itself, both leading edges will be attached to the flat plate during the experiments. The final component is called as lower surface flat cover. All parts were made from aluminium as shown in Figure 4.

2.3 Measurement Techniques

The experiments were carried out at the Reynolds number of 1×10^6 , 2×10^6 & 3×10^6 corresponding to the speeds of 18 , 36 & 54 m/s base on the mean aerodynamic chord of 0.87 m . The angles of attacks were varied from $0^\circ \leq \alpha \leq 25^\circ$ with 3° increment. The models were attached to six-component external balances located underneath turntable. The models installation is shown in Figure 5. From the figure, the model angle of attack can be created by adjusting the aft support vertically. Two measurements techniques were employed on the model. The first experiment was the steady balance data to obtain the forces and moments in x , y and z . The steady balance data are measured using a heavy capacity external balance located underneath the wind tunnel. This load cell can measure the forces

and moments in 6 axes. For this project, the lift and drag are measured by forces in $-x$ and $-y$ axes while the pitching moment is measured by the moment in the $-z$ axis. The sampling rate for each channels were captured at 100 Hz.

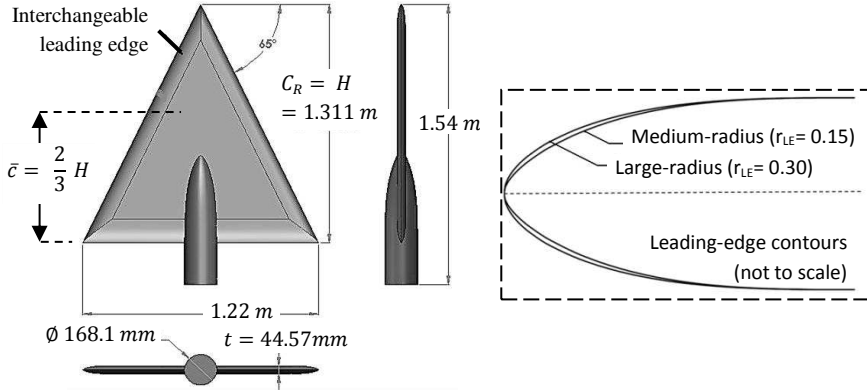


Figure 3
Dimensions of UTM VFE-2 model

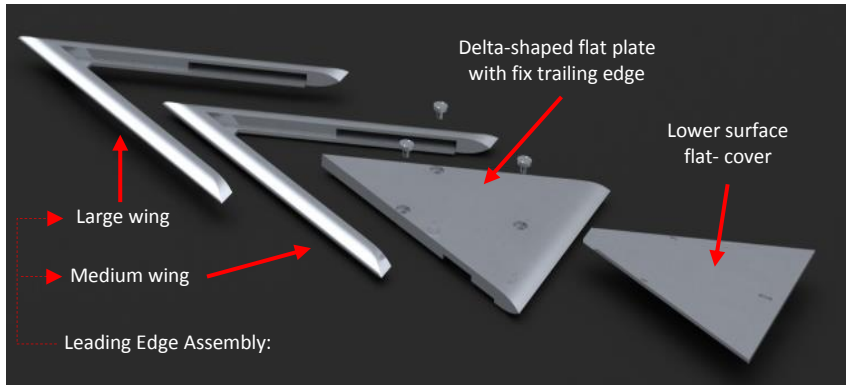


Figure 4
Delta wing model assembly

The final experiment was the surface pressure measurements that were captured on the upper surface of the wing. There were 86 pressure taps located on starboard side of the wing. The diameter of the orifice was $d_{outer} = 1 \text{ mm}$ which located normal to the wing surface. The pressure taps were arranged in 10 different chord-wise stations started from 10% to 97% from the wing apex, i.e. in Y/C_R of 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8 and 0.97 of the wing. In spanwise positions, more pressure taps were placed in the leading edge region. This is done in order to measure the primary vortex that developed in the leading edge. This is shown in Figure 6. During the experiment, the data were captured at 1000 samples in 10 second or the sampling rate was 100 Hz. The installation of UTM-VFE-2 model in the test section of UTM-LST is shown in Figure 5 below.

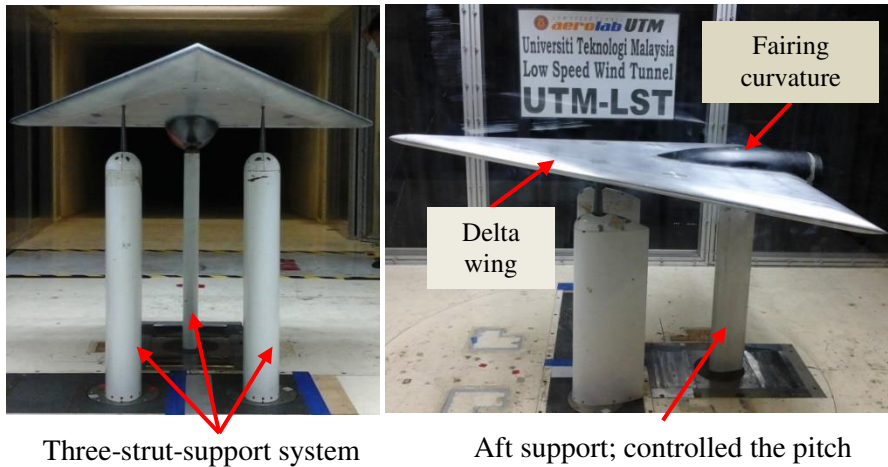
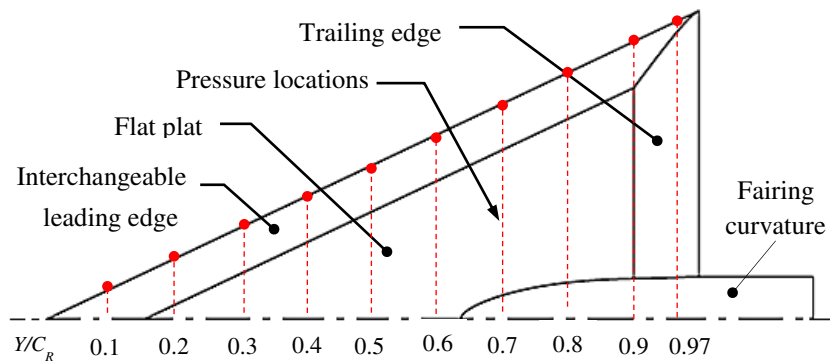


Figure 5
The model installations at UTM-LST



Chord-wise surface pressure locations:
 $Y/C_r = 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.97$

Figure 6
Leading-edge contours (not to scale)

3 Results and Discussion

In order to observe the impact of leading edge bluntness on the primary vortex, the data obtained from the large-edged wing were compared to the medium-edged wing.

3.1 Medium and Large-edged Flow Characteristics

3.1.1 Aerodynamic Coefficients

The data obtained from the external balance data has been analysed and presented in Figure 7. In order to investigate the influence of leading edge bluntness on C_L , C_D and C_M the data obtained from the large-edged wing has been compared with the results obtained from the medium-edged wing. The sample data at R_{mac} of 3×10^6 is shown in Figure 7 below. The results show that both C_L and C_D are reduced if leading edge bluntness is increased. This is consistent with Mat [12] who experimentally showed that the magnitude of lift and drag forces decrease when the leading edge bluntness is increased. This situation occurs because the strength of the primary vortex is decreasing when leading edge bluntness is increased. The reduction in C_D is also caused by the increase in leading edge suction force acting in the leading edge of the wing.

A clear observation in Figure 7(a) shows that the C_L for large-edged wing reduces compared to medium-edged wing. This situation happened starting from $\alpha = 6^\circ$ onwards. This shows that C_L decreases when the leading edge bluntness is increased. This phenomenon is linked with the strength of the primary vortex. The increases in the leading-edge radius have weakened the primary vortex. The primary vortex is weakened because the primary separation has been delayed by the leading edge profile. Another factor that causes C_L to decrease is the attached flow. The large portion of attached flow covered in the apex region of the large-edged wing has reduces the C_L as shown in Figure 7(a). The large fraction of the leading edge suction force act on the large-edged has contributed to this behaviour and thus increased the C_L/C_D ratio shown in Figure 7(c). The results here consistent with Ronoei [21] who experimentally measured the C_L/C_D ratio on a generic span delta wing.

The pitching moment coefficient (C_M) is plotted in Figure 7(d). In general, both wings show to have a nose-down pitching moment. The medium-edged wing has experienced a higher nose-down C_M compared to the large-edged until $\alpha \approx 13^\circ$. This may link to the greater strength of primary vortex that formed earlier on the medium-edged wing compared to the large-edged wing. At angle of attack $\alpha \geq 13^\circ$, the pitching moment has becomes more negative for both cases. At attack angle higher than $\alpha = 13^\circ$, it is notable that C_M for large-edged wing is higher than medium-edged wing. The reason for this is unknown to date and more experiments are needed to verify this.

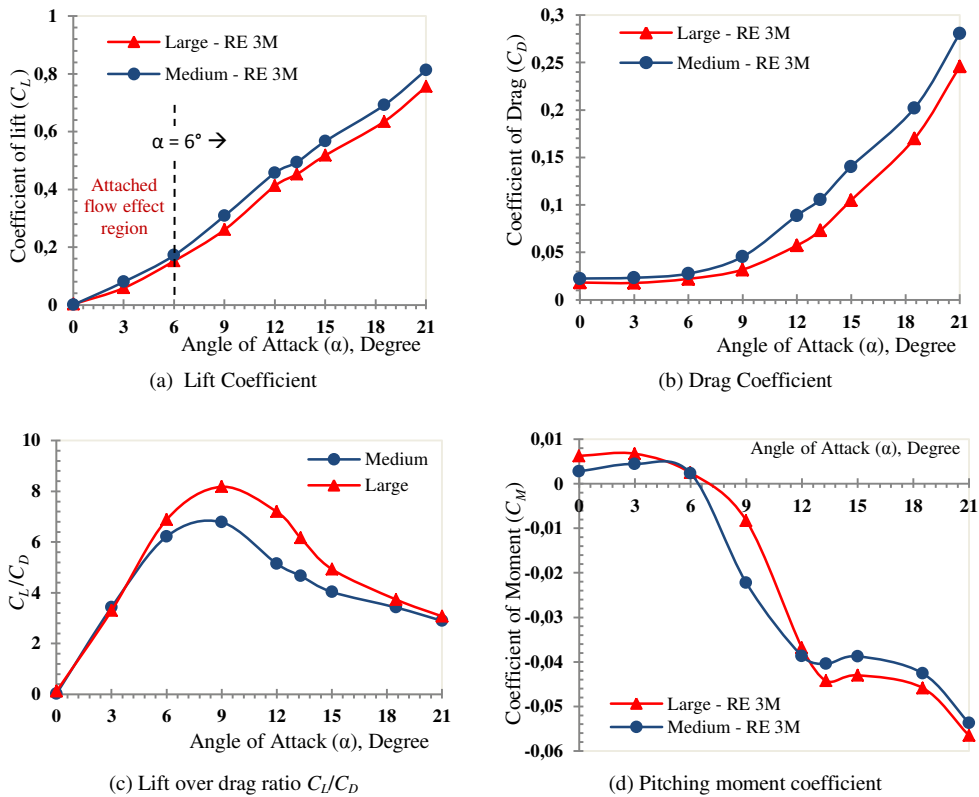


Figure 7
Effects of leading edge bluntness on aerodynamic coefficients at $R_{mac} = 3 \times 10^6$

3.1.2 Surface Pressure Coefficient

This section discusses the results obtained from the surface pressure measurement studies measured on the upper surface of both wings. In order to compare the effects of leading edge bluntness, the surface pressure obtained for large-edged wing has been compared to medium-edged wing. For example, the result at $\alpha = 13^\circ$ and $R_{mac} = 3 \times 10^6$ is compared in Figure 8. The pressure taps were arranged in 10 different chord-wise stations started from 10% to 97% of the wing. For the medium-edged wing, it can be noted that the attached flow is formed from the wing apex to 40% downstream. The primary vortex begins to occur at 50% from the apex. The leading edge effect is obvious here, when the leading edge bluntness is increased, it is notable that the attached flow area on the large-edged has covered about 60% of the wing from the apex. Then, the primary vortex begins at about 70% from the wing apex [19, 22].

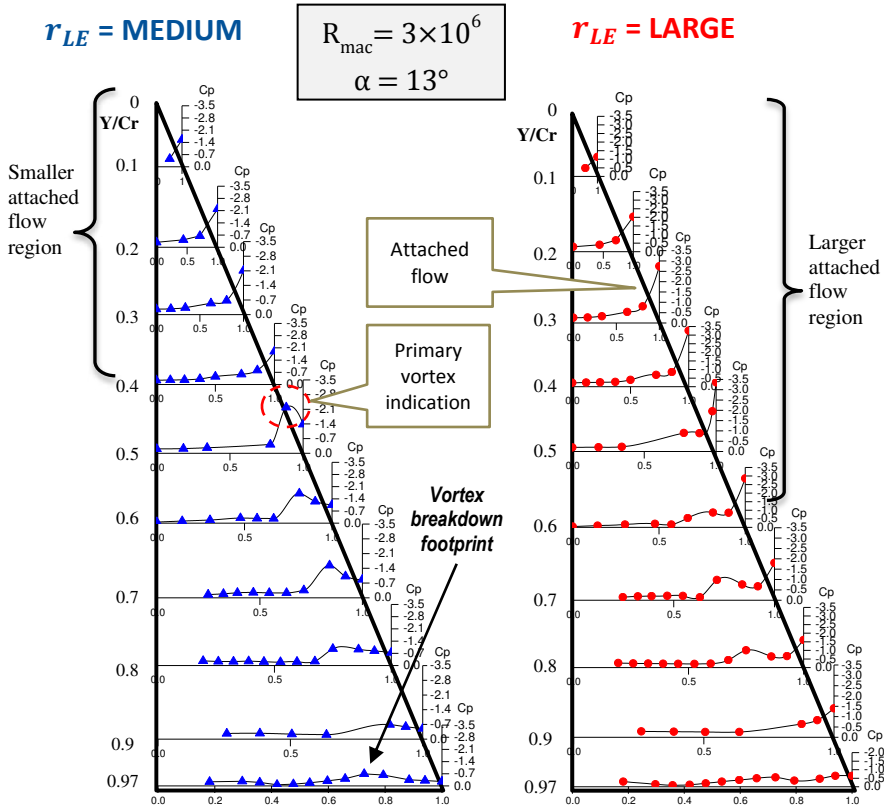


Figure 8

Pressure distribution for Medium- and Large-edged wing at $R_{mac} = 3 \times 10^6$, $\alpha = 13^\circ$

The results at higher angle of attack of $\alpha = 21^\circ$ is shown in Figure 9. The surface pressure for the medium-edged wing has been compared with the data from the large-edged wing in the same figure. From the Figure, it has been noted that the primary vortex has moved forward to 30% from the wing apex on medium-edged wing compared to 40% from the Apex for the large-edged case. This indicates that the upstream progression of the primary vortex has been slowed at a higher angle of attack. In order to observe the strength of the primary vortex on both wings, the pressure coefficient at positions $Y/C_r = 0.3, 0.6 \& 0.7$ were compared in the diagram. From the figure, it can be observed that the peak for the medium-edged wing is relatively higher compared to the large-edged for all positions. This indicates that the strength of the primary vortex decreases if the leading edge bluntness is increased.

The impact of leading edge bluntness on the vortex breakdown is also shown in the figure. A clear observation in the trailing edge area at Y/C_r of 0.8 and below showed that the vortex breakdown is delayed for the large-edged wing compared to the sharper wing of the medium-edged. The stable shear layer on the blunter

wing is suspected to delay the breakdown. By having short run of attached flow in the leading edge region and delay in separation had reduces in the instability of the shear layer on the blunter leading edge [22-23].

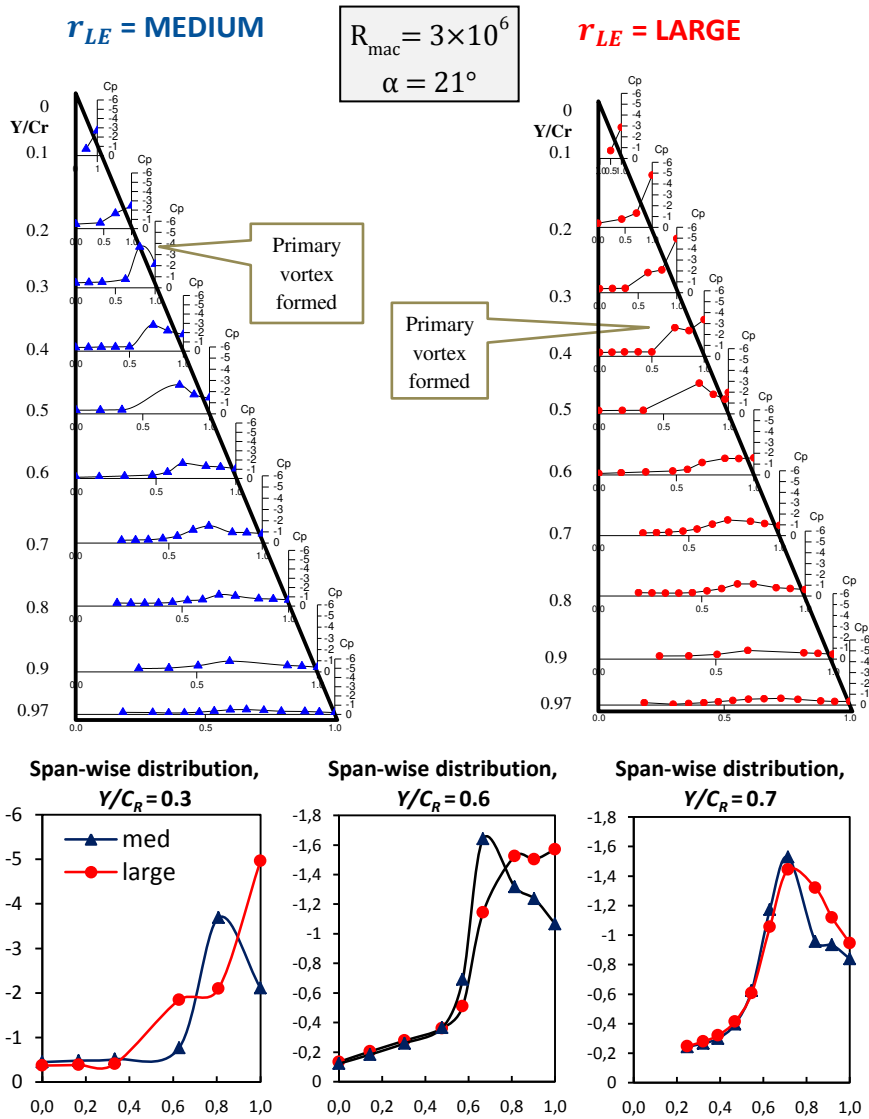


Figure 9
Pressure distribution for Medium- and Large-edged wing at $R_{mac} = 3 \times 10^6$, $\alpha = 21^\circ$

A statistical technique called as Krigging method has been used to obtain the flow topology on the surface of the wing. The sample surface flow topology performed on both wings at $\alpha = 13^\circ$ and R_{mac} of 2×10^6 is shown in Figure 10. The figure shows that the primary vortex is shifted outboard on the blunter wing of the large-edged wing. This again shows that the size and strength of the primary vortex has decreased when the leading edge bluntness is increased.

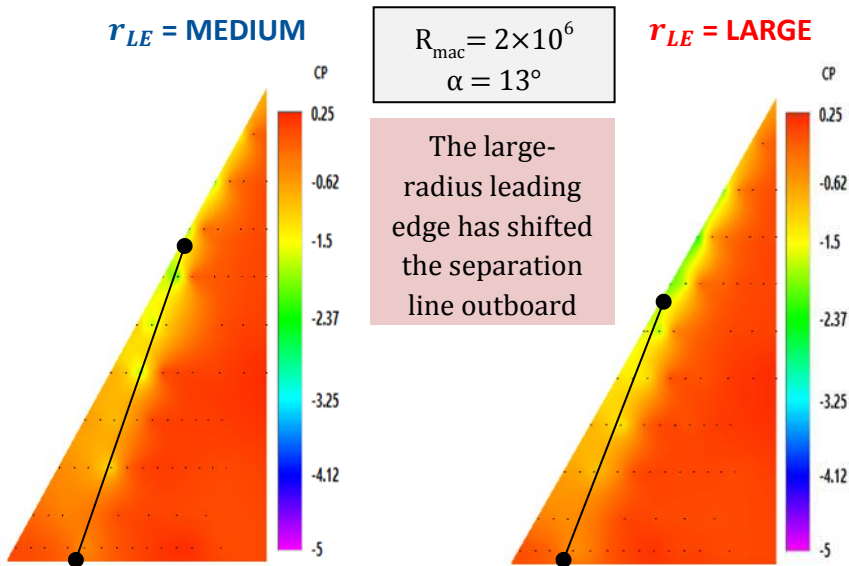


Figure 10

Flow topology comparison at $R_{mac} = 2 \times 10^6$, $\alpha = 13^\circ$

3.1.3 Leading Edge Pressures

The pressure coefficients in the leading edge can be used to predict the onset of the primary vortex on the blunt-edged wing [17-18]. In this case, the leading edge coefficient has been plotted for positions Y/C_r of 0.2, 0.4, 0.5, 0.6 & 0.8 from the apex. This is done in Figure 11 for Reynolds number of 1×10^6 and Figure 12 for Reynolds number of 3×10^6 . At Reynolds number 1×10^6 , the flow remains attached to the surface even at $\alpha = 25^\circ$ for both medium- and large-edged wings at $Y/C_R = 0.2$. The results obtained here contrast with Mat [12-13] who experimentally showed that the primary vortex was developed in the apex area at $R_{mac} = 1 \times 10^6$. An important observation of these current results in the flow is attached to the surface in the apex area as long as the leading edge is blunt [19]. The effects of leading edge bluntness is observed at $Y/C_R = 0.4$ where the flow on the medium-edged wing has separated at $\alpha = 9^\circ$ while it separates at $\alpha = 12^\circ$ for large-edged wing. Similar observation also can be noticed at $Y/C_R = 0.8$ where the separation occurs earlier on the shaper wing. It can be concluded here that the increase in leading bluntness has delayed the formation of vortex above the wing.

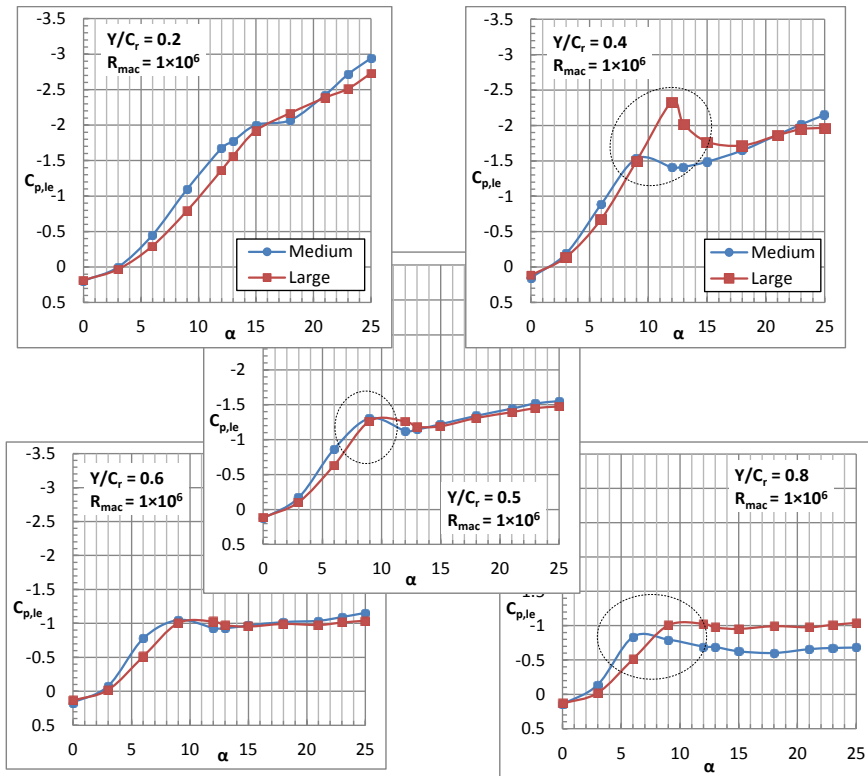


Figure 11

The bluntness effects to leading edge pressures at $R_{mac} = 1 \times 10^6$

The flow characteristics at higher Reynolds of 3×10^6 are shown in Figure 12 below. At this Reynolds number, the effects of leading edge bluntness are more obvious. At position $Y/C_r = 0.2$, the onset of the primary vortex is developed at $\alpha = 12^\circ$ for medium-edged wing, while for the large-edged wing, the primary vortex is still not developed even when the attack angle has been increased to $\alpha = 21^\circ$. At position $Y/C_r = 0.4$, it can be seen that the primary vortex has developed at $\alpha = 9^\circ$ for the medium-edged wing, while it developed at a higher attack angle of $\alpha = 18^\circ$ on the blunter wing. A similar situation happened at Y/C_r of 0.5, 0.6 and 0.8. The current data here showed that the upstream progression of the primary vortex has been delayed if the leading edge bluntness is increased [24].

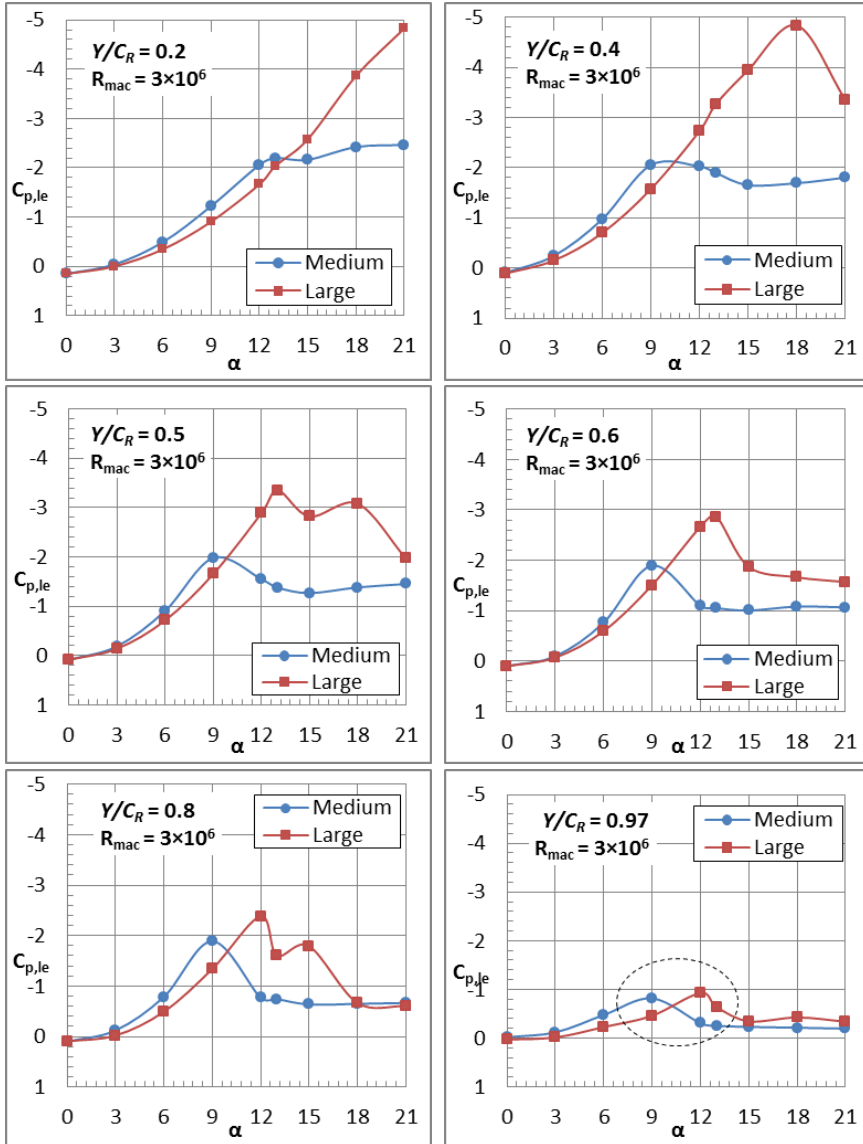


Figure 12
The bluntness effects to leading edge pressures at $R_{mac} = 3 \times 10^6$

3.2 Reynolds Number Effects on Large-edged Wing

Since the data on large-edged wing is limited, this section will further discuss the effects of Reynolds number, angle of attack on this wing. The effects of Reynolds

number at constant angle of attack of $\alpha = 13^\circ$ is presented in Figure 13 below. The data compared the surface pressure on the upper surface at three different Reynolds numbers of 1×10^6 , 2×10^6 and 3×10^6 . The primary vortex developed at about a 20% chordwise distance from the apex ($Y/C_r = 0.2$) at a Reynolds number of 1×10^6 . When the Reynolds number is increased to 2×10^6 and 3×10^6 , the primary vortex shifted further aft of the wing to at about $Y/C_r = 0.4$ and $Y/C_r = 0.6$ of the apex. The results showed that the increase in Reynolds number has slowed down the primary vortex further aft of the wing. The results here were consistent with [13].

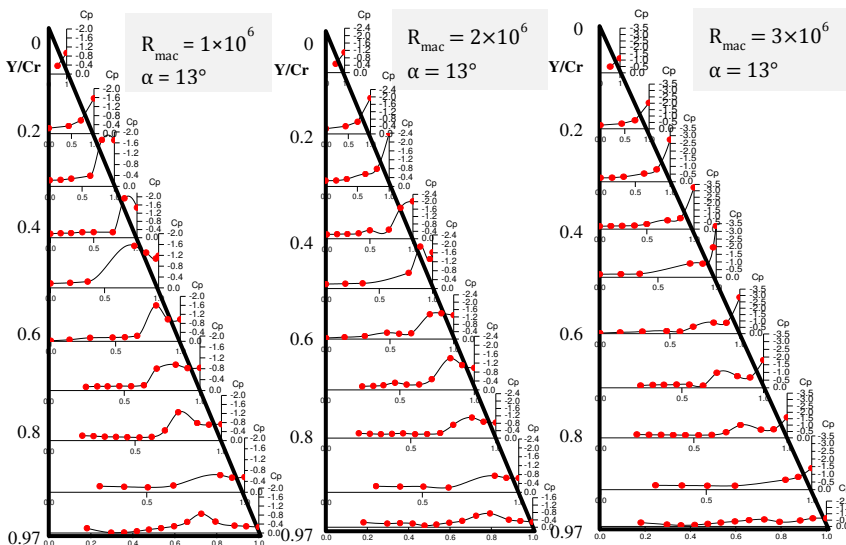


Figure 13

Reynolds number effects on large-edged wing at $\alpha = 13^\circ$

The flow characteristics when the angle of attack is increased to $\alpha = 18^\circ$ is shown in Figure 14. Similar flow physics is observed here where the Reynolds number has delaying the upstream progression of the primary vortex further aft of the wing. The surface flow topology in the second figure showed the primary vortex has been shifted more outboard when the Reynolds number is increased. In addition, the magnitude of pressure topology formed in the leading increases when the Reynolds number is increased. This shows that the primary vortex is stronger when the Reynolds number is increased. The plot of surface pressure in the third figure at Y/C_r at 0.7 also showed that the primary vortex is shifted outboard with the Reynolds number. The characteristics of the flow either being laminar or turbulence, the main factor that leads to these results [6]. At a low Reynolds number where the flow is dominated by laminar flow, the onset of the primary vortex develops earlier. The stronger ability of the turbulent boundary layer at higher Reynolds number has endured the adverse pressure gradient and thus delaying the development of the primary vortex [6, 15].

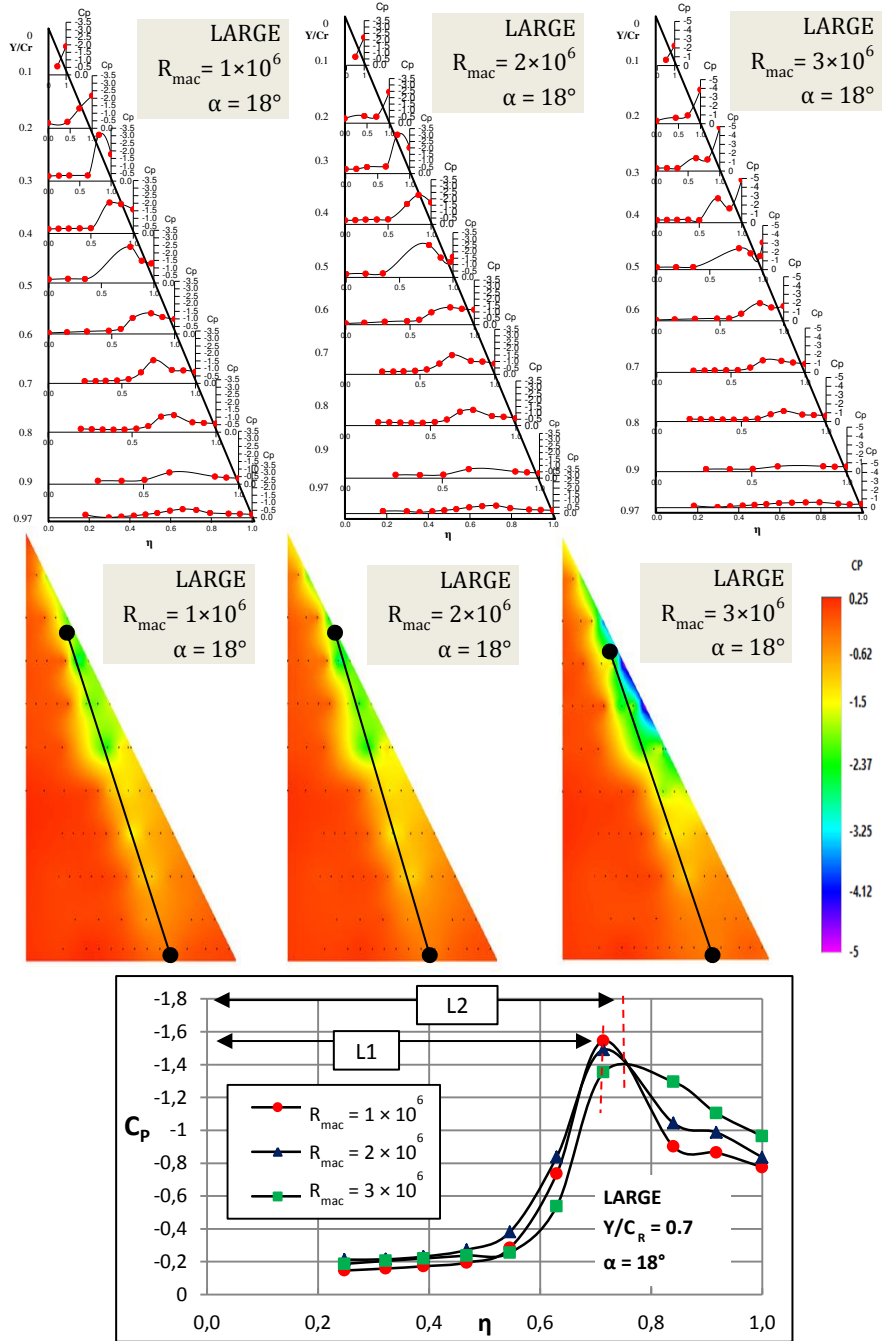


Figure 14

Reynolds number effects on large-edged wing at $\alpha = 18^\circ$

One of the problems that cannot be solved with the VFE-2 for medium-edged wings was to assess the laminar/turbulent status. At lower Reynolds number the flow is laminar and it is turbulent at high Reynolds number, S. Mat [13] in his experiment has shown that the flow at Reynolds number 1×10^6 is dominated by laminar flow. The Figure 15 shows the distribution of pressure coefficients at Reynolds number of 1×10^6 for the large-edged wing. From the figure, it can be observed that the flow is attached to the surface at relatively low attack angles. In addition, it can be noted also that the attached flow still existed even if the attack angle of attack has been increased to $\alpha = 23^\circ$. The boundary layer status is still unverified from this experiment. More experiments are needed to verify this.

4 Further Experiments on Blunt-edged Delta Wing

Delta wing is the best platform for the development of the Unmanned Air Combat Air Vehicle (UACV) aircraft. For most UCAV aircrafts, the wing has been designed with blunt leading edge. The data obtained from this experiment provides a useful knowledge for future UCAV development. In a continuation of the VFE-2 project, another model of delta derivative wing called diamond wing was proposed and is currently fabricated. The interests in this project were an extended research project that initiated from AVT-183 task group, a collaborative task group with AVT-113 under NATO. This research project will focus on understanding the detail interactions between the inboard inner vortex and the primary vortex of blunt-leading-edged vortex separation. The diamond wing was configured with blunt leading-edge of constant airfoil, moderate leading-edge sweep of 53° categorized as non-slender wing, and swept trailing edge as shown in Figure 16.

Besides the ability to induce vortex potential lift, diamond wing configurations with blunt and reduced sweep angle were more relevant to application because it also can enhance aircraft longitudinal static stability [25]. However, diamond wings exhibit more complex vortical flows as compared to slender, sharp-edged wings. The vortices formed on diamond wings are more unsteady and breakdown occurs at a much lower angle of attack than on highly-swept slender wings.

This current research was focused to investigate the inboard inner vortex effects to the onset and progression of leading-edge vortex separation. The leading-edge vortices investigation will have both moderate-sweep effects and blunt-leading-edge effects, coupled together [26]. Several measurement techniques will be employed on the wing suitable to measure the unsteady vortices on the diamond wing configuration as shown in Figure 17.

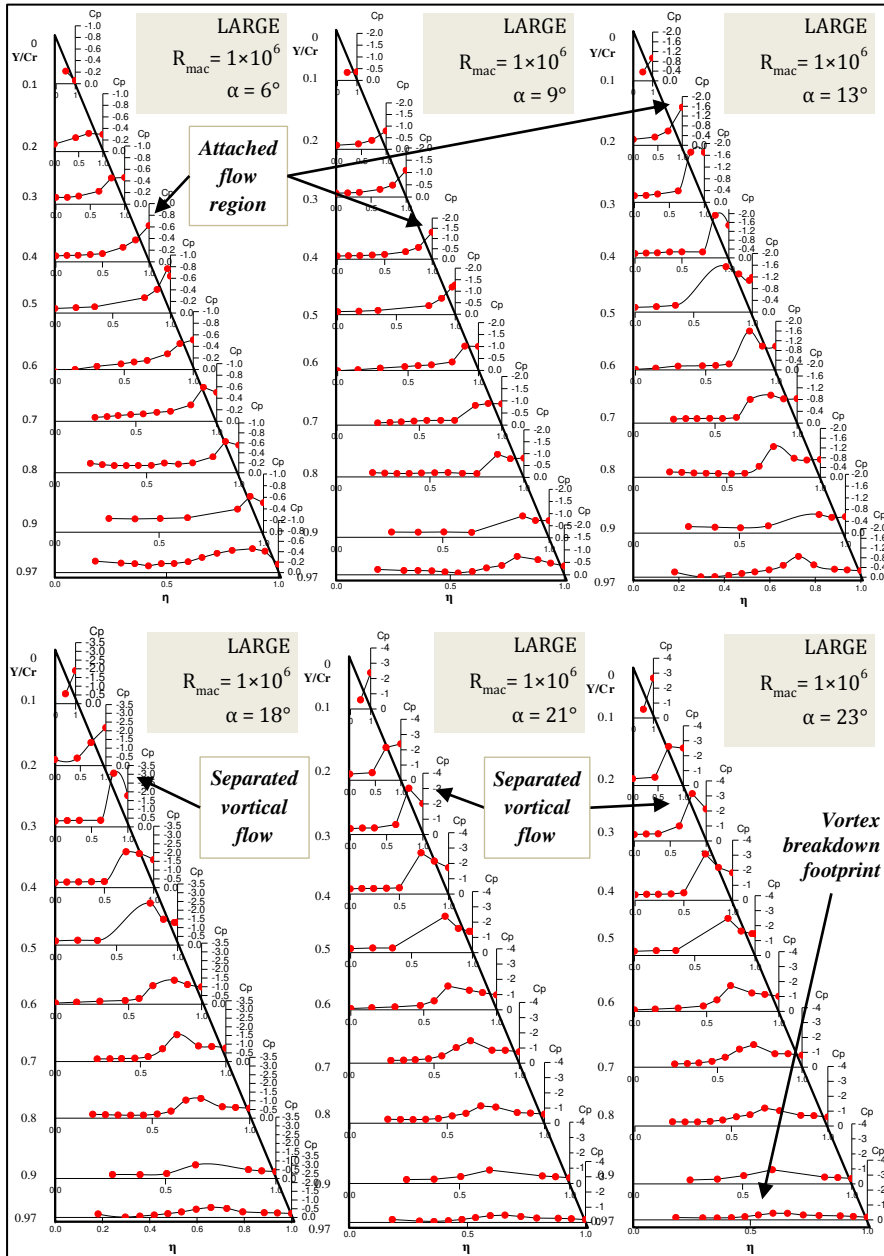


Figure 15

Pressure distribution on large-edge wing at $R_{mac} = 1 \times 10^6$ at various angles of attack

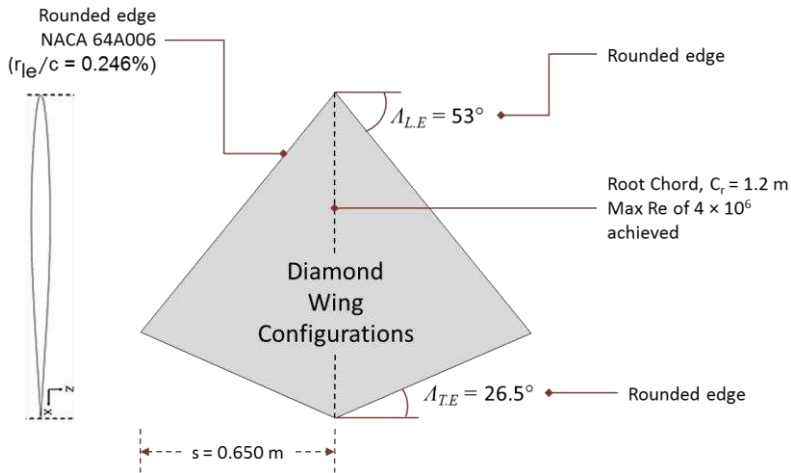


Figure 16
Diamond wing configurations

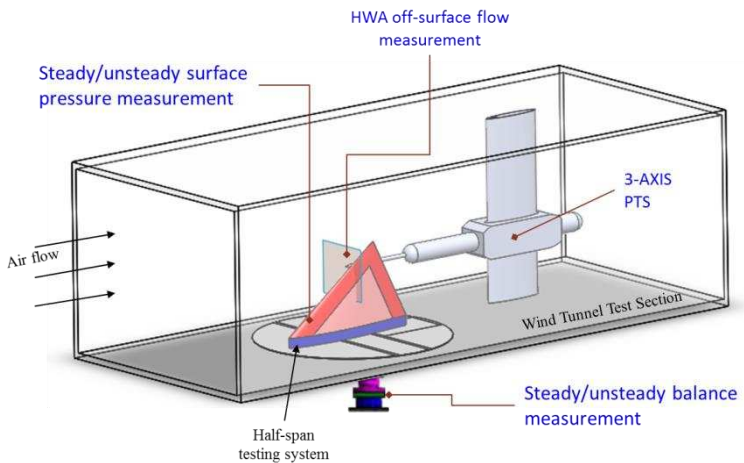


Figure 17
UTM Diamond wing experimental setup and measurements

Conclusions

This paper discusses further the effects of leading edge bluntness on the vortex properties above blunt-edged delta wing. In the VFE-2 campaign, concentrations have been given to a medium-edged wing. A series of experiments were conducted to study the performance of vortex on the blunter wing of large-edged wing. The results obtained from the large-edged wing have been compared with those from the medium-edge wing. The current results here showed that the primary vortex above large-edged delta wing is also dependent on Reynolds

number, leading edge bluntness and angle of attack. The results from steady balance data has showed that the lift/drag ratio is increased if the leading edge bluntness increases. Among the important observations from this study was the area covered by attached flow is enlarged. That the primary vortex also developed further aft of the wing has been shifted outboard to the leading edge area with the bluntness effects. The advantage of the blunter wing also that the formation of the vortex breakdown and its upstream progression has been delayed. Since most of the UCAV aircrafts are in the delta-shaped planform, this paper has highlighted some of the most important considerations in the design stage such as the progression of the primary vortex and vortex breakdown behaviors.

Acknowledgement

This research was funded by a grant from Ministry of Higher Education Malaysia (UTM Research University Grant No. 4F172, 4F718 and 12H06). The data presented, the statement made, and views expressed are solely the responsibility of the author.

References

- [1] Kulfan, R. M. (1979) Wing Geometry Effects on Leading Edge Vortices. Aircraft Systems and Technology Meeting. 20-22 August. New York, 79 – 1872
- [2] Earnshaw, P. B. (1962) An Experimental Investigation of the Structure of a Leading- Edge Vortex. Aeronautical Research Council Reports and Memoranda, No 3281
- [3] Pershing, B. (1964) Separated Flow Past Slender Delta Wings With Secondary Vortex Simulation. El Segundo Technical Operations. TDR–269(4560 – 10) – 4
- [4] Gad-El-Hak, M. and Blackwelder, R. F. (1985) The discrete Vortices from a Delta Wing. Technical Report 1985, Vol. 23, 961-962
- [5] Hummel, D. (1979) On the vortex formation over a slender wing at large incidence. AGARD-CP-247, Paper No. 15
- [6] Hummel, D. (2004) Effects of Boundary layer Formation on the vortical Flow above Slender Delta Wings. RTO specialist Meeting on Enhancement of NATO military Flight Vehicle Performance by Management of Interacting Boundary Layer transition and Separation. Meeting Proceedings RTO-MPAVT- 111. Page 30-1 to 30-2
- [7] Luckring, J. M. (2004a) Compressibility and Leading-Edge Bluntness Effects for a 65° Delta Wing. 42th AIAA Aerospace Science Meeting and Exhibit. 5-8 January, Reno, Nevada. AIAA-2004-0765
- [8] Luckring, J. M. (2004b) Reynolds Number, Compressibility, and Leading Edge Bluntness Effects on Delta Wing Aerodynamics. 24th International

- Congress of the Aeronautical Sciences. 29 – 3 September. Yokohama, Japan
- [9] Luckring, J. M. (2013) Initial Experiments and Analysis of Blunt-edge Vortex Flows for VFE-2 configurations at NASA Langley, USA. *Aerospace Science and Technology*. Vol. 24, Issue 1, pp. 10-21
- [10] Chu, J. and Luckring, J. M. (1996) Experimental Surface Pressure Data Obtained on 650 Delta Wing across Reynolds Number and Mach number Ranges. NASA Technical Memorandum 4645
- [11] Coton, F., Mat, S. B., and Galbraith, R. (2008) Chapter 22 – Experimental Investigations on the VFE-2 Configuration at Glasgow University, United Kingdom. RTO-TR-AVT-113, Pages 22-1 – 22-18
- [12] Mat, S. B., (2011) The analysis of flow on round-edged delta wings. Doctor Philosophy. University Of Glasgow, United Kingdom
- [13] Mat, S. B., Green, R., Galbraith, R., and Coton, F. (2015) The Effect of Edge Profile on Delta Wing Flow. *Proceedings of the Institution of Mechanical Engineers, Part G: Journal of Aerospace Engineering*
- [14] Kurun, S. (2008) Chapter 23 – Experimental Investigations on the VFE-2 Configuration at TUBITAK-SAGE, TURKEY. RTO-TR-AVT-113. Pages 23-1 – 23-18
- [15] Furman, A. and Breitsamter, C. (2013) Turbulent and Unsteady Flow Characteristics of Delta Wing Vortex Systems. *Aerospace Science and Technology*. Vol. 24, Issue 1, pp. 32-44
- [16] Rodriguez, O. (2008) Chapter 20 – Experimental Investigation on the VFE-2 Configuration at ONERA, France. RTO-TR-AVT-113
- [17] Luckring, J. M. and Hummel, D. (2008) Chapter 24 – What Was Learned From The New VFE-2 Experiments. RTO-TR-AVT-113
- [18] Luckring, J. M. and Hummel, D. (2013) What Was Learned From The New VFE-2 Experiments. *Aerospace Science and Technology*. Vol. 24, Issue 1, pp. 77-88
- [19] Tajuddin, N., Mat, S., Said, M. and Mansor, S. (2017) Flow Characteristic of Blunt-edged Delta Wing at High Angle of Attack. *Journal of Advanced Research in Fluid Mechanics and Thermal Sciences*. Vol. 39. Issue 1, pp. 17-25
- [20] Konrath, R., Klein, C., and Schröder, A. (2013) PSP and PIV Investigation on the VFE-2 Configuration in Sub- and Transonic Flow. *Aerospace Science and Technology*. Vol. 24, Issue 1, pp. 22-31
- [21] Ronoie, K. (1996) Low Speed Aerodynamics Characteristics of 60° Rounded Leading-Edge Delta Wing with Vortex Flaps: Part 1: 457 mm Span Delta Wing. Cranfield University. COA Report No. 9611

- [22] Said, M. B., (2016) Effects of Leading Edge Radius, Reynolds number and Angle of Attack on The Vortex Formation above Large-Edged Delta Wing. Master Degree. Universiti Teknologi Malaysia, Malaysia
- [23] Renac, F., Barberis, D. and Molton, P. (2005) Control of Vortical Flow over a Rounded Leading- Edge Delta Wing. *AIAA Journal*. Vol. 43, No. 7, pp. 1409-1417
- [24] Said, M. and Mat, S., (2016) Effects of Reynolds Number on The Onset of Leading Edge Vortex Separation Above Blunt-edge Delta Wing VFE-2 Configurations. 30th Congress of the International Council of the Aeronautical Sciences, 25-30 September, Daejeon, South Korea, 2016_0608
- [25] Hitzel, S. M. (2013) Perform and Survive – Evolution of Some U(M)CAV Platform Requirements. STO AVT Workshop on Innovative Control Effectors for Military Vehicles. Stockholm, Sweden, 20-22 May, No. 1, STO-MP-AVT-215
- [26] Luckring, J. M, Boelens, O. J., Breitsamter, C., Hövelmann, A., Knoth, F., Malloy, D. J., Decke, S. (2016) Objectives, approach, and scope for the AVT-183 diamond-wing investigations. *Aerospace Science and Technology*. Vol. 57, pp. 2-17

Unsupervised Clustering for Deep Learning: A tutorial survey

Artúr István Károly^{1,2}, Róbert Fullér^{3,4}, Péter Galambos¹

¹ Antal Bejczy Center for Intelligent Robotics

Óbuda University, Bécsi út 96/B. H-1034 Budapest, Hungary
info@irob.uni-obuda.hu

² Doctoral School of Applied Informatics and Applied Mathematics

Óbuda University, Bécsi út 96/B. H-1034 Budapest, Hungary

³ Institute for Applied Mathematics

Óbuda University, Bécsi út 96/B. H-1034 Budapest, Hungary

⁴ Department of Informatics,

Széchenyi István University, Egyetem tér 1, H-9026, Győr, Hungary
e-mail: rfuller@sze.hu

Abstract: Unsupervised learning methods play an essential role in many deep learning approaches because the training of complex models with several parameters is an extremely data-hungry process. The execution of such a training process in a fully supervised manner requires numerous labeled examples. Since the labeling of the training samples is very time-consuming, learning approaches that require less or no labeled examples are sought. Unsupervised learning can be used to extract meaningful information on the structure and hierarchies in the data, relying only on the data samples without any ground truth provided. The extracted knowledge representation can be used as a basis for a deep model that requires less labeled examples, as it already has a good understanding of the hidden nature of the data and should be only fine-tuned for the specific task. The trend for deep learning applications most likely leads to substituting as much portion of supervised learning methods with unsupervised learning as possible. Regarding this consideration, our survey aims to give a brief description of the unsupervised clustering methods that can be leveraged in case of deep learning applications.

Keywords: Unsupervised learning; Clustering; Deep learning

1 Introduction

The three primary methods for learning-based systems are supervised, unsupervised and reinforcement learning. Reinforcement learning is applied in fields when an agent takes actions in an environment, and a suitable policy for acting has to be learned [1]. The other two learning methods are used when the output of the system does not influence the inputs in any way. In case of supervised learning the training samples are provided with correct labels, so a ground truth is available. Meanwhile, in unsupervised learning, no such a-priory knowledge of the data is required. Models that are trained in an unsupervised manner only require the collected data for training.

Deep learning is a widely researched topic currently. The majority of deep learning approaches utilize supervised learning [2]. However, given the vast number of trainable parameters of such models, the training process requires numerous labeled examples in order to achieve good generalization [2]. The labeling of the samples is a very resource-intensive and time-consuming process, and usually, the labeling can only be done manually [3]. So naturally arises a need for methodologies that enable the training of such models with less or no labeled examples. This is usually done by applying unsupervised learning first, and then fine-tuning the model with the help of labeled samples and supervised learning [2, 4]. Among others, the future of deep learning is expected to be driven by the development of more sophisticated and accurate unsupervised learning methods [2].

Supervised learning methods always have a well-defined objective, like classifying the inputs into one of the formerly known classes, or the regression of a function between a set of inputs and inspected outputs. In this case, the output features are formerly known, just like the class labels in classification. In unsupervised learning, however, the aim is to discover unknown structures, similarities, and grouping in the data [5].

Clustering is the process when the objective is to create groups of data samples (clusters) based on some kind of similarity measure between the data samples [5]. The difference between classification and clustering is that clustering is carried out in an unsupervised manner, so no class labels are provided, and sometimes even the number of clusters is not known a-priory.

In this survey, we provide a brief introduction of the most significant unsupervised clustering methods and their applicability in the field of deep learning. We aim to give a summary of such clustering techniques that can also be leveraged in deep learning applications, in the aspect of the expected trends of the future development in this field of study [2]. Previous approaches focused on either clustering methods or unsupervised deep learning. A detailed and general description of unsupervised clustering methods can be found in the work of Xu and Wunch, who provide an in-depth survey on clustering algorithms in [5]. They introduce the general description of a clustering procedure, provide multiple measures for similarity that are used

for clustering, and give a detailed explanation of several clustering algorithms and their applicability. Other works that provide a detailed introduction of clustering algorithms are [6] and [7]. Bengio et al. give an exhaustive survey on unsupervised learning in deep learning in [4].

2 Unsupervised clustering methods

In case of unsupervised learning, the task is to uncover hidden relationships, structures, associations or hierarchies based on the data samples provided to the system [5]. This kind of information gives us a better understanding of the data and the underlying processes generating it. The clusters can be constructed based on similarity or distance measures, so similar samples belong to the same cluster. The similarity measure can also be described as a distance measure because the similarity of two samples can be interpreted as the distance between the two samples in the feature space [5].

One approach for unsupervised clustering is to use these similarity measures and construct the regions of the feature space corresponding to the different clusters based on the computed distances between the training samples [5].

Another approach is to extract features with high discrimination power or to find principal directions of the feature space along which the data can be better separated. These features can be any subset of the input data or they can also be constructed by a network architecture [4].

Often, it is more convenient to use the computation mechanism of a supervised learning algorithm, but work it around, so no true labeling is needed. An example of this is anomaly detection, which can be handled as a clustering problem. In case of anomaly detection, one can expect two clusters, one for the normal samples and one for the anomalies. In case of an unsupervised training process for anomaly detection, one can extract the training data from only one cluster, the cluster of normal samples. As the training data, in this case, is expected to come from only one cluster, it can be automatically labeled as belonging to the same class, and the calculation for a supervised method can be used to train a system to classify similar samples as ones that belong to this cluster [8]. However, when an anomaly appears, that is not similar to the examples seen during the training process, then it is classified as a sample not belonging to the cluster of normal samples; thus it is classified as an anomaly.

The automatic generation of multiple labels is also beneficial. However, this method is no longer referred to as unsupervised learning. In the case of simulator-based labeling, for example, there is no need for the time consuming manual labeling of the data [9]. This enables large and very diverse datasets to be annotated. However, these samples are still provided with ground truth, even though they are not manually labeled. In case of the anomaly detection, the ground truth is the same for all of the

samples, so it is more of an assumption on the nature of the data rather than real labeling. That is why we call the use of only one label unsupervised learning, and not the automatic generation of multiple labels.

In case of unsupervised learning for deep learning models, there are two major approaches. Both rely on formulating an output for the neural network that can be used to train it in the usual way with the help of gradient descent [4].

The first one is to try to reconstruct the input of the network on its output. The loss function is computed based on the reconstruction error between the input and the output of the network [10, 4]. This method is expected to extract a meaningful compressed representation of the input from which it can be reconstructed with minimal error. This requires the compressed features to represent features of high discrimination power among the presented training samples. This way, the unsupervised training can be carried out on the whole network or layer-by-layer. After such training, the network is usually trained further with labeled examples, but it requires much less of them because thanks to the unsupervised pre-training a good representation of the input data is already available [2].

The second one is to use two networks in parallel. One of these networks that are called generator is used to generate data that is as similar to the input data of the other network (the discriminator) as possible [11]. The discriminator's objective is to discriminate the generated samples from the real ones. Both networks are trained in parallel. The generator is trained to produce data that can fool the discriminator even better and the discriminator is trained to be able to differentiate between synthetic and real data more accurately. The model itself appends the label for the inputs of the discriminator as a synthetic sample or real sample because it knows which samples come from the generator, while the update of the generator is based on the output response of the discriminator. So the system does not require a ground truth annotation. During the process of training, the discriminator has to develop an understanding of the features of the training dataset and later it can be used for classification as well (in a similar way like the anomaly detectors) [11].

3 Clustering algorithms

In this section we provide a brief description of the clustering algorithms which are especially suitable for deep learning applications. In table 1, we draw a straightforward categorization of the mentioned unsupervised clustering methods 1.

First, we discuss the approaches that are based on a distance measure. These methods define the similarity of data samples by the distance of the samples in the feature space. Because of this property, the simpler variants of these methods fail to cluster data that is not linearly separable. However, with the creative formulation of the similarity measure, or with proper pre-processing of the data, even these techniques can be applied for nonlinear clustering tasks [12].

Class	Methods	Section
Distance measure based	K-means clustering	3.1
	Hierarchical clustering	3.2
	Fuzzy clustering	3.3
	Support vector machines	3.4
	Spectral clustering	3.5
	Decision trees	3.6
Statistical	Expectation maximization algorithm	3.7
Neural networks	Self organizing maps (Kohonen networks)	3.8
	Adaptive resonance theory	3.9
	Autoencoders	3.10
	Co-localization	3.11
	Generative models	3.12

Table 1
Classification of unsupervised clustering methods

3.1 K-means clustering

The k-means algorithm is an iterative (learning) method to discover k number of clusters in the input space. The number k is defined a-priory [13, 14, 15]. Each cluster is represented with a cluster centroid in the feature space. The similarity measure of samples is a simple distance measure. The distance between the samples and the cluster centroids can be computed, and a sample is associated with the cluster with the closest centroid. All of the samples are associated with one and only one cluster, so the clusters cover the whole feature space, and they are disjoint. The iteration process consists of two stages: In the first stage, all the training samples are associated with one of the clusters. Then in the second stage, the location of the cluster centroids is updated. These two steps are repeated until a stopping condition is met. The stopping criteria can be that no further change in the classification of training samples happened after the last update of the cluster centroids, or the distance between the centroids before and after the update is smaller than a specified value (ϵ) [13, 14, 15].

Before the iteration process, there are some preliminary considerations to make that affect the result of the clustering. These considerations are the number of clusters k , the starting positions for the centroids of the clusters and the stopping condition.

A simple measure of the distance between the samples and the centroids is the Euclidean distance. The Euclidean distance of the two points is the Euclidean norm

of the vector pointing from one of the points to the other. The computation of the Euclidean norm of a vector $\mathbf{x} \in \mathbb{R}^n$ can be seen in equation 1.

$$\|\mathbf{x}\| = \sqrt{x_1^2 + \dots + x_n^2} \quad (1)$$

The Euclidean distance of a sample $\mathbf{x} \in \mathbb{R}^n$ and a centroid $\mathbf{c} \in \mathbb{R}^n$ is $\|\mathbf{x} - \mathbf{c}\|$. Apart from the Euclidean distance, other metrics can be used like the Manhattan or the Mahalanobis distance [5].

The update of the centroids of clusters is carried out after each element of the training set has been associated with one of the clusters. The new cluster centroids can be computed as like a center of mass for all the samples associated with that given cluster.

The starting location of the cluster centroids can be randomly placed in the feature space, or random samples can be selected as initial cluster centroids to ensure that they are located in the feature space from where the data is drawn.

An algorithm for the k-means clustering method is presented in algorithm 1. In the algorithm, the squared Euclidean distance is used for distance measurement, in line 13.

Algorithm 1 : K-means clustering

- 1: **Definitions:**
 - 2: Let k be the number of clusters,
 - 3: \mathbf{X} , the set of training samples,
 - 4: $\mathbf{x}_i \in \mathbf{X}$, the i^{th} sample, $i \in \{0, 1, \dots, n\}$,
 - 5: $\mathbf{x}_i \Rightarrow j^{th}$ cluster, means that the i^{th} sample belongs to the j^{th} cluster,
 - 6: \mathbf{c}_j , the j^{th} cluster centroid, $j = \{1, 2 \dots k\}$,
 - 7: s_j , the number of samples associated to the j^{th} cluster
 - 8: **Initialization:**
 - 9: Let \mathbf{c}_j^+ be a random sample drawn from \mathbf{X} , $\forall j$
 - 10: **Iteration:**
 - 11: **repeat**
 - 12: $\mathbf{c}_j = \mathbf{c}_j^+$, $\forall j$
 - 13: $\forall i$, **find** j for which $\|\mathbf{x}_i - \mathbf{c}_j\|^2$ is minimal $\forall j$ and set $\mathbf{x}_i \Rightarrow j^{th}$ cluster
 - 14: $\forall j$, $\mathbf{c}_j^+ = \frac{1}{s_j} \sum_{k=1}^{s_j} \mathbf{x}_k$, where $\mathbf{x}_k \Rightarrow j^{th}$ cluster
 - 15: **until** $\mathbf{c}_j^+ - \mathbf{c}_j < \epsilon$, or no samples are re-associated to an other cluster
-

It can be seen that the k-means algorithm is heavily constrained. Its performance is profoundly affected by the proper selection of preliminary parameters, like the number of clusters or the initial location for the centroids. There are methods for determining a good set of parameters for a given training set [16]; however these methods usually require the construction of several clustering systems with different

parameters and evaluating their results. It is computationally intensive to calculate the distance of each sample and each centroid, so for a large number of samples, the basic algorithm has to be altered [17]. Also, if the number of clusters has to be changed, the whole iteration process has to be done from scratch, so methods for enabling on-line k-mean clustering were also developed [18, 19]. The basic method also fails to reveal appropriate clusters for non-linearly separable data. This can be worked through by formulation a similarity measure that operates in a transformed space (which is usually higher dimensional) where the samples are linearly separable [12].

3.2 Hierarchical clustering

Unlike the k-means algorithm, hierarchical clustering does not propose disjoint clusters. It builds a hierarchical structure of clusters, that can be represented as a dendrogram [5]. The leaves of the dendrogram structure are the samples themselves (each belonging to its own class), and the root of the structure is the cluster that includes all of the samples. Thus cutting the dendrogram at different levels of hierarchy results in a different number of clusters. Unlike the k-means clustering algorithm, hierarchical clustering does not require the a-priory declaration of the number of clusters, however doing so can serve as a stopping condition, resulting in faster computation for the proposed clusters. The dendrogram can be built from the leaves to the root (agglomerative method) or from the root toward the leaves (divisive method) [5].

Both agglomerative and divisive methods are based on distance measures like the Euclidean distance to compute the similarity of clusters. This measure in the context of hierarchical clustering is called dissimilarity measure, which is the basis of the four major clustering strategies [5, 20].

Single linkage clustering defines the similarity of two clusters with the help of the minimum of all pairwise dissimilarities between the elements of the two clusters [20]. The complete linkage clustering strategy defines the similarity of two clusters as the maximum of the pairwise dissimilarities between the elements of the two clusters [20]. If the similarity between the clusters is defined by the average of the pairwise dissimilarities of the samples in the two clusters, then it is group-average clustering [20]. Finally, the clusters can also be given centroids (computed from the samples belonging to the clusters), just like the clusters in the k-means algorithm. The centroid clustering strategy defines cluster similarity with the dissimilarity measure between the centroids of the clusters [20].

Agglomerative methods start with assigning a cluster for all samples. Then a cluster for the two most similar clusters is created. This process is repeated until all samples belong to a single cluster (root) or until a certain amount of clusters (k) is discovered [5, 20].

Divisive methods start from one cluster (the root of the dendrogram) that holds

all the samples. The cluster is divided into two sub-clusters, but due to the large number of possible splits, the evaluation of all of the possible splits would be too computationally expensive. Usually, a good split is done by finding the two elements of the cluster with the highest dissimilarity and grouping the other samples to the element of the selected two, that is more similar to the given sample [20]. The created clusters can also be split into two, while all the resulted clusters contain only one sample, or a formerly given number of clusters is discovered.

The divisive method is harder to implement, but it can extract more meaningful clusters than the agglomerative approach because the latter tends to construct clusters based on local similarities without the knowledge of the global distribution of the data, while divisive methods have global information from the beginning [20].

3.3 Fuzzy clustering

The previously introduced methods of k-means and hierarchical clustering consider clusters with hard margins, meaning that a sample either belongs to a given cluster or not. In case of fuzzy clustering methods, a sample has a degree of membership in a given cluster, which is a continuous value rather than a binary. As clustering has a close relation to set theory, most of the clustering algorithms have a fuzzy implementation. In this section, we introduce a fuzzy equivalent of the k-means algorithm, the fuzzy c-means algorithm [21, 22, 23].

The fuzzy c-means algorithm is very similar to the k-means algorithm. The applied objective function to be minimized can be seen in equation 2. Where N is the number of samples, C is the number of clusters, \mathbf{x}_i is the i^{th} sample, $i \in \{1, 2 \dots N\}$, \mathbf{c}_j is the j^{th} cluster centroid, $j \in \{1, 2 \dots C\}$, μ_{ij} is the degree of membership of \mathbf{x}_i in cluster j , $\|\cdot\|$ is any norm for measuring distance (like the Euclidean norm) and m is a coefficient to control fuzziness $1 \leq m \leq \infty$ [21, 22].

$$J_m = \sum_{i=1}^N \sum_{j=1}^C \mu_{ij}^m \|\mathbf{x}_i - \mathbf{c}_j\|^2 \quad (2)$$

The function J_m is minimized with an iteration process, during which the degrees of membership for each sample and each clusters are updated. After the update, the new centroids for the clusters are computed. The degrees of membership can be determined according to the equation 3 [21, 22].

$$\mu_{ij} = \frac{1}{\sum_{k=1}^C \left(\frac{\|\mathbf{x}_i - \mathbf{c}_j\|}{\|\mathbf{x}_i - \mathbf{c}_k\|} \right)^{\frac{2}{m-1}}} \quad (3)$$

$$i \in \{1, 2 \dots N\}, j \in \{1, 2 \dots C\}$$

The centroids of the clusters are calculated like in equation 4 [21, 22].

$$\mathbf{c}_j = \frac{\sum_{i=1}^N \mu_{ij}^m \cdot \mathbf{x}_i}{\sum_{i=1}^N \mu_{ij}^m} \quad (4)$$

$$j \in \{1, 2, \dots, C\}$$

The preliminary steps before the iterative algorithm are to define the number of clusters (C), the coefficient for fuzziness (which is usually set to 2) and to assign an initial degree of membership for all training sample for each cluster. This is usually done by filling a matrix U of size $N \times C$ with random values for μ_{ij} . The stopping condition can be formulated exactly like in the k-mean algorithm [21, 22].

After these preliminary steps, the cluster centroids are computed with the help of equation 4 based on the training samples and the given U matrix containing the degree of membership of each training sample in each cluster. Then the elements of the matrix U are modified according to equation 3. These two steps are repeated until the stopping condition is met.

It can be seen from equation 3 that in the marginal case, if m is set to be $m = 1$, the degrees of membership converge to either zero or one, making it a crisp clustering method, like k-means.

According to this approach, most of the clustering methods can be fuzzified by assigning a degree of membership to the samples.

3.4 Support vector machines

The SVM-based clustering is usually referred to as support vector clustering described in detail in [24]. The idea behind the support vector clustering method is based on the work of Schölkopf et al. [25] and Tax and Dunin [26], who introduced methods to carry out the support vector description [27] of data structures in a high dimensional space with the help of kernel functions [12].

Rather than separating the data samples from each other directly in feature space, the kernel function enables to formulate a distance measure of a higher dimensional space called feature space and use this kernel function to design the separation of the data samples [12]. Such a separation can lead to highly nonlinear, complex decision boundaries in the data space.

The complete mathematical description of the support vector clustering (SVC) method can be found in [24]. In this survey, we only explain the core idea behind this technique.

In case of SVC, the training samples are mapped to a high dimensional feature space utilizing a Gaussian kernel function. The data in feature space is enclosed in a hypersphere of center \mathbf{a} and a radius of R . A penalty parameter is added to control the allowed number of outliers. An outlier \mathbf{x} is a sample in data space for which $\|\Phi(\mathbf{x}) - \mathbf{a}\|_2^2 > R^2 + \xi$. Where $\Phi(\cdot)$ is the kernel function that maps the sample \mathbf{x} from the data space to the feature space, and ξ is the slack variable to enable soft margin [24].

The contour of the hypersphere forms boundaries in the data space that separates points of the data space that are inside and those that are outside of the hypersphere when mapped to the feature space, with the given kernel function. These boundaries can be non-convex and can even form disjoint sets of points in the data space. The shape of the decision boundary depends on the parameters of the kernel function and the penalty coefficient for outliers [24]. The proper tuning of these parameters depends on the noise and overlap of structures in the provided data, and it is detailed in [24]. If the parameters are all set to suitable values, then smooth disjoint boundaries should form in the data space.

The clusters are marked by the disjoint sets in the data space [24]. So two samples in the data space \mathbf{x}_1 and \mathbf{x}_2 are said to belong to different clusters if any path that connects these two points in the data space exits the hypersphere in the feature space. In [24] this criterion is checked numerically for twenty points of a connecting line between \mathbf{x}_1 and \mathbf{x}_2 .

A more straightforward approach for unsupervised clustering with support vector machines is to use the one-class support vector machine (OCSVM) [25, 26]. The OCSVM method operates as the basis of the SVC algorithm. If only two clusters are expected, like in anomaly detection [8], there is no need for the cluster assignment method proposed in [24] so the system can be simplified.

3.5 Spectral clustering

Spectral clustering is used for graph partitioning [28] by analyzing graphs with methods of linear algebra. The spectral clustering algorithm is also based on a similarity measure. The training data can be represented as a similarity graph, which is an undirected graph, with the training samples as the vertexes and the edges associated with a weight of the similarity between the two vertexes they connect. From the similarity graph, the graph Laplacian is computed. The different kinds of similarity graphs and graph Laplacians can be found in [28].

The graph Laplacian matrix is used to split the data into clusters. Given a required number of clusters noted by k , the first k eigenvectors with the largest corresponding eigenvalues of the graph Laplacian are selected [28]. These eigenvectors are used as centroids for the clusters. The data samples are then associated with one of the clusters with the help of the k-means method.

The implementation and interpretation of the spectral clustering method are de-

scribed from several aspects in [28].

3.6 Decision trees

A decision tree is a tree-like graph structure. In order to assign a sample to a cluster, the data is fed to the root node of the structure. At the nodes of the structure, the data is inspected according to one of its given features and decided in which branch the given sample should be propagated. So a splitting node has branches for all possible values of the given feature towards the leaves. The leaves of the structure correspond to different clusters, meaning that data with similar features fall to the same cluster [29].

The structure is constructed based on a training set. During the training, the objective is to find the best way to split the data, so to select an appropriate inspection feature for all of the nodes. This method is called a greedy algorithm to find the splitting feature [29]. The number of clusters can be controlled by the branching factor of the structure [29].

The appropriate splitting feature is selected based on a measure of their discrimination power. A split based on a feature that can discriminate the data will result in sets that are more homogeneous than that before the split. So often, the homogeneity measure is used that can also be computed based on the distance metric between the samples [30].

Basak and Krishnapuram proposed a method for unsupervised clustering with decision trees [30]. They introduced two homogeneity measures in their paper that are based on the distance metric between the samples. During the construction of the structure, at each node, a clustering of the training set is carried out based on a single or a group of features [30].

The appropriate feature to split by in a given node can be selected by removing features from the feature space and computing its effect on the homogeneity of the data, or by simply observing the homogeneity of the data along each feature separately [30].

Let the number of data points be N . If the similarity of two data points \mathbf{x}_i and \mathbf{x}_j is computed like μ_{ij} $i, j \in \{1, 2 \dots N\}$ in equation 5, where d_{ij} is the distance between \mathbf{x}_i and \mathbf{x}_j by a distance measure that is not necessarily the Euclidean distance and d_{max} is the maximum distance of all pairwise distances between the data points [30].

$$\mu_{ij} = g \left(1 - \frac{d_{ij}}{d_{max}} \right) \quad (5)$$

$$i, j \in \{1, 2 \dots N\}$$

The function g is computed according to equation 6 [30].

$$g(x) = \begin{cases} x & \text{if } 0 \leq x \leq 1 \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

The discrimination power of a feature can be calculated according to equation 7, if the discrimination power is measured by removing the feature from the feature set, and according to equation 8, if the discrimination power is measured exclusively for a given feature [30].

$$H_f = - \left(\sum_{i,j} \mu_{ij} (1 - \mu_{ij}^f) + \mu_{ij}^f (1 - \mu_{ij}) \right) \quad (7)$$

$$\hat{H}_f = - \sum_{i,j} \hat{\mu}_{ij}^f (1 - \hat{\mu}_{ij}^f) \quad (8)$$

$$i, j \in \{1, 2 \dots N\}$$

The expression μ_{ij}^f is the similarity of \mathbf{x}_i and \mathbf{x}_j with the feature f removed from the feature set and $\hat{\mu}_{ij}^f$ is the similarity of \mathbf{x}_i and \mathbf{x}_j only along the feature f . In both cases the feature with the highest discrimination power is selected for splitting feature at a given node [30].

Basak and Krishnapuram found that with this method, a well interpretable hierarchical clustering of the data can be made, which can also be translated into clustering rules [30]. They also found that it is better to select a single feature as a splitting criterion than a set of features.

3.7 Expectation maximization algorithm

The expectation maximization algorithm is an iterative process to find the parameters of statistical models that describe the probability distribution of the data [31]. In case of clustering, we suppose that the data can be distributed into k different clusters and data in each cluster have its own probability distribution with its given parameters. The type of the distribution of the clusters is based on assumption [31].

If it is known which data points belong to which clusters, the estimation for the parameters of their probability distribution can be computed. If the parameters of the probability distributions are given, the probability for each data sample coming from the given distribution, belonging to a specific cluster can be computed. This boils down to a k-means-like iteration process, where the initial parameters for the probability distribution of the clusters are initialized, the probability of samples belonging to one cluster thus can be computed, and the parameters of the probability

distributions can be refined based on these calculated probabilities [31]. Then this process is repeated until a specified stopping condition is met.

So the expectation maximization algorithm is very much like the fuzzy c-means clustering, but with a stochastic aspect [31]. Instead of the degree of membership of the samples in each cluster, the probability of the samples of belonging to the clusters is used and the parameters to be updated are the parameters of the assumed statistical model [31].

3.8 Self organizing maps (Kohonen network)

The Kohonen network is a fully connected artificial neural network with only an input layer and an output layer [32, 33]. The input vectors presented to the network are associated with one of the k different clusters. In the output layer of the network, there is an output neuron for every cluster. So the output layer has k number of neurons. The neuron with the highest activation decides the cluster a sample belongs to [32, 33].

The network can be trained with the winner-take-all method [32, 33]. Let the weight vectors of the output neurons be \mathbf{w}_i where $i \in \{1, 2, \dots, k\}$ and the input vectors to be \mathbf{x} . The j^{th} output neuron is selected as the winner neuron, and the sample is associated with its cluster, if

$$\|\mathbf{x} - \mathbf{w}_j\| = \min_{i=1\dots k} \|\mathbf{x} - \mathbf{w}_i\|$$

The weight vector with the minimum distance from the sample can also be found by finding the maximum of the scalar products of the input vector and the weight vectors if the weight vectors are all normalized [32, 33]. This can be seen in equation 9.

$$\langle \mathbf{w}_j, \mathbf{x} \rangle = \max_{i=1\dots k} \langle \mathbf{w}_i, \mathbf{x} \rangle \quad (9)$$

Equation 9 only holds if the weight vectors of the network are normalized, so $\|\mathbf{w}_i\| = 1 \forall i$. The scalar product of the input vector and the weight vectors is the activation of the output neurons. This is why the winning neuron is selected as the one with the highest activation.

The interpretation of the scalar product is the projection of x in the direction of \mathbf{w}_i . So if the projection is greater in a given direction, that means the input vector is more similar to the normal weight vector pointing to that direction.

In case of the winner-take-all method, only the weight vector of the winner neuron is modified during the training process [32, 33]. The objective is to minimize the squared distance between the input vector and the winning weight vector. This can be done by computing the gradient of the objective function and use gradient descent

to modify the weights of the winning neuron [32, 33].

The gradient of the objective function with respect to the weights can be seen in equation 10.

$$\frac{d\|\mathbf{x} - \mathbf{w}_j\|^2}{d\mathbf{w}_j} = -2(\mathbf{x} - \mathbf{w}_j) \quad (10)$$

According to equation 10 the weight vector \mathbf{w}_j should be modified in the direction of $\mathbf{x} - \mathbf{w}_j$. The modified weights can be calculated according to equation 11.

$$\mathbf{w}_j := \mathbf{w}_j + \eta(\mathbf{x} - \mathbf{w}_j) \quad (11)$$

After the weight update, the weight vector \mathbf{w}_j has to be normalized again.

If the training data is linearly separable, the weight vectors of the Kohonen network will converge to point to the center of mass of the clusters [32, 33]. The number of output neurons must be larger than the number of clusters even if the number of clusters is not known a-priory. The appropriate neurons to cluster by then can be selected, by inspecting the direction of their weight vectors during and after training and the ones that are not necessary can be omitted [32, 33].

3.9 Adaptive resonance theory

The adaptive resonance theory (ART) [34] model for neural networks is very similar to the Kohonen network, but it includes some other functionality [35]. The structure of the ART model is the same as the Kohonen network with the exception that the output neurons also implement lateral inhibition. So the activation of an output neuron decreases the activation of the others. The objective is the same, to find a weight vector that is similar to the input vector. After the classification of the input vector, the output activations are compared to a vigilance parameter [35]. If the activation of the winning neuron is higher than the vigilance parameter, the training continues like it was a Kohonen network. However, if the vigilance parameter is larger than the winning neuron's activation, it means that the presented input vector is out of an expected range around the weight vector. In this case, the winning neuron is turned off, and the prediction is made again. This is done until one of the weight vectors overcome the vigilance parameter. If it is not overcome in any trials, then such a neuron is selected that does not represent a cluster yet, and its weights are modified towards the input vector [35].

With the tuning of the vigilance parameter, ART models can control the smoothness of classification [35]. High values of the vigilance parameter result in fine clusters and a lower value of the vigilance parameter results more general, smooth clusters [35].

3.10 Autoencoders

Autoencoders are artificial neural network structures that try to reconstruct their input on their output [36]. The loss function is computed from the reconstruction error of the network. As only the inputs and the computed outputs are used for the loss function, there is no need for labeling, and thus the network can be trained in an unsupervised manner [36].

The training of the autoencoder structure is carried out with the error backpropagation algorithm. A simple autoencoder structure can be seen in Figure 1.

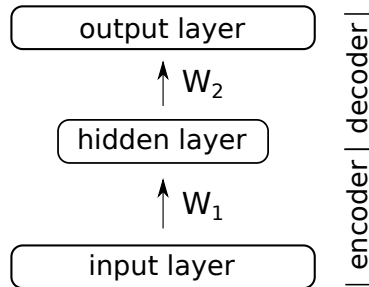


Figure 1
Simple autoencoder structure

In the structure in Figure 1, the network has fully connected layers. Let the input vector be $\mathbf{x} \in \mathbb{R}^n$, so the number of neurons in the input layer is n . The objective is to reconstruct the input in the output layer, so the output layer also consists of n number of neurons. The number of neurons in the hidden layer is chosen to be m where $m < n$. The matrix W_1 is the weight matrix between the input and the hidden layers. Each row of the W_1 matrix is a weight vector of a neuron in the hidden layer. So W_1 is a matrix of $\mathbb{R}^{m \times n}$. The rows of W_2 are the weight vectors of the output neurons, so $W_2 \in \mathbb{R}^{n \times m}$.

The hidden layer activations are computed according to equation 12, where $f(\cdot)$ is the activation function of the hidden layer neurons, and it is applied element-wise to the result of the matrix-vector multiplication. It is important to note that the bias weights are not treated separately in this equation.

$$\mathbf{h} = f(W_1 \mathbf{x}) \quad (12)$$

The output of the network can be computed similarly to the hidden layer activations.

$$\mathbf{y} = f(W_2 \mathbf{h})$$

As the hidden layer has fewer neurons than the input layer and the objective of the

network is to reconstruct the input on its output, the encoder part is responsible for compression of the data [36]. The compressed representation of the data must be informative about the input so it can be decoded with high accuracy. So the encoder part of the network tries to extract features from the input that describe the data well. These features can be used for clustering because they represent directions in the feature space in which the input data can be well-separated [36].

The training of the structure can be accomplished by forming a loss function from the difference of the input and the output, such as $\|\mathbf{x} - \mathbf{y}\|^2$ and minimizing this function with respect of the weights of the network [36]. The gradient of the loss function with respect of the weights can be computed, so the error-back-propagation algorithm can be used to train the network without the need for labeled examples [36].

3.11 Co-localization

Image co-localization is an unsupervised method for object discovery in images. A set of images that all contain a common object are provided as training samples [37]. The task of image co-localization is to find a bounding box for the common object across all images [37].

Wei et al. proposed an image co-localization method, DDT (deep descriptor transformation) that leverages the features extracted by pre-trained convolutional neural networks [37]. A deep descriptor is a component vector of the output volume of a convolutional layer. If the output volume has the dimensions of w for width, h for height and d for depth, then a deep descriptor at the index i, j is the vector $\mathbf{x}_{i,j} \in \mathbb{R}^d$, $i \in \{1, \dots, h\}$ and $j \in \{1, \dots, w\}$.

The pre-trained model is presented with N number of images that has a common object on them to be localized. The deep descriptors for each image are collected $X^n = \{\mathbf{x}_{i,j}^n \in \mathbb{R}^d\}$, where $n \in \{1, \dots, N\}$ and the mean vector for all of the deep descriptors is calculated as in equation 13 [37].

For equations 13,14 and 15, $i \in \{1, \dots, h\}$, $j \in \{1, \dots, w\}$ and $K = w \cdot h \cdot N$.

$$\bar{\mathbf{x}} = \frac{1}{K} \sum_n \sum_{i,j} \mathbf{x}_{i,j}^n \quad (13)$$

After calculating the mean vector, Wei et al. compute the covariance matrix according to equation 14 [37].

$$\text{Cov}(\mathbf{x}) = \frac{1}{K} \sum_n \sum_{i,j} (\mathbf{x}_{i,j}^n - \bar{\mathbf{x}})(\mathbf{x}_{i,j}^n - \bar{\mathbf{x}})^T \quad (14)$$

The eigenvector of $\text{Cov}(\mathbf{x})$ that corresponds to the largest eigenvalue is noted as ξ_1

and the first principal component of a deep descriptor at an index of i, j for a given image is described as in equation 15 [37].

$$p_{i,j}^1 = \xi_1^T (\mathbf{x}_{i,j} - \bar{\mathbf{x}}) \quad (15)$$

The first principal component can be calculated for all values of i and j , and the result can be organized into a matrix $P^1 \in \mathbb{R}^{h \times w}$. The elements of P^1 for a given image with positive value represent positive correlation across all the N number of images for that descriptor, so it is likely to belong to the common object [37]. Thus the matrix P^1 is thresholded at zero for all of the images and the location of the largest connected positive region is sought. As the dimension of the P^1 matrix is the same as the feature map of the convolutional layer ($w \times h$), the location in the P^1 matrix can be reflected the image. So a region on the image can be found that correlates across all of the images [37]. A minimal enclosing bounding box can be formed for the proposed regions, thus solving the task of image co-localization. Also, if the P^1 matrix does not contain any element with a positive value, it means that the image does not contain the common object [37].

3.12 Generative models

In generative adversarial networks (GANs), there are two networks trained at the same time [11], a generator network and a discriminator network. The generator network generates data from random vectors, and the discriminator network tries to tell apart the generated and synthetic data samples. The objective of the generator network is to fool the discriminator. So it develops its internal parameters in a way that it can generate data seemingly coming from the same domain as the real training samples. The discriminator has to develop an understanding of the essential features of the data in order to be able to discriminate between the synthetic ones [11].

After both the networks are trained, the transfer of the output of the generator network can be examined by interpolating in the input space [11]. The results show that the transfer between two input vectors with input space interpolation is smooth. This implies that the discriminator network also possesses a smooth continuous representation of the feature space, which means that such a model can be used for extracting robust general low-level features even in case of training sets that are discontinuous in the feature space [11].

Mathematical operations with the input vectors also show that the generator deals with the features in the sense of similarity [38]. In case of image generation for human faces, for example, let an input vector \mathbf{x}_{sf} yield an output of a smiling female face, \mathbf{x}_{nf} a neutral female face and \mathbf{x}_{nm} result in a neutral male face. The input vector $\mathbf{x}_{sm} = \mathbf{x}_{sf} - \mathbf{x}_{nf} + \mathbf{x}_{nm}$ will result in a smiling male face. This also implies, that the discriminator also has a sense of similarity, like smiling faces are similar, female faces are similar, male faces are similar etc. moreover, this knowledge can be utilized for clustering as well [38].

4 Applications

In this section, we introduce some examples of how the different unsupervised clustering techniques can be leveraged in deep learning applications.

The k-means clustering can be used to learn low-level filters for convolutional neural networks [39, 40, 41]. Socher et al. introduced a convolutional-recursive deep learning structure for 3D object recognition from RGB-D data [39]. A single convolutional layer first processed both the RGB and the depth modalities. They proposed an unsupervised method to build the filters based on the k-means clustering algorithm. They compared their proposed method to other models introduced in [42, 43, 44, 45]. Their experiments show, that their model, with an accuracy of 86.8 ± 3.3 , was able to outperform all other methods except for the one introduced in [45], which had a 0.7% higher accuracy, but required five times more memory.

Coates and Ng introduced an unsupervised method for learning deep hierarchical features with the help of k-means clustering [46]. In their paper they describe the main considerations and limitations to perform multiple layers hierarchical feature representation with k-means clustering. They also show that this method, with an accuracy of 82%, can achieve the performance of the state-of-the-art unsupervised deep learning methods such as vector quantization (81.5% accuracy) and convolutional DBN (78.9% accuracy) on the full CIFCAR-10 dataset, but with easier implementation (only one hyperparameter 'k') and better scalability. A similar approach can be seen in [47].

Reducing the dimensionality of the data is an essential task for both visualizing high dimensional data for better understanding and for clustering, as the distance measures become simpler in reduced dimensional spaces. Yang et al. proposed a method to optimize the dimensionality reduction and the clustering method together, to construct a meaningful representation of the data in a reduced dimensional space [48]. They used a deep neural network as the dimensionality reduction system and trained it in respect to the clustering method (k-means clustering). In order to avoid trivial solutions where the network maps any input to such latent space that it can be trivially separated, a loss for the reconstruction of the input was also introduced, like in the autoencoder structure. This way the network was able to create a latent representation of the input with well separable clusters that are evenly scattered around the cluster centroids.

For graph partitioning Tian et al. showed that the reconstruction of the similarity matrix with autoencoders is a suitable alternative for the traditional matrix calculations in case of spectral clustering, in large-scale clustering problems, where the input space is very high dimensional [49]. The hidden layer activations can be used for the k-means clustering directly, instead of calculating the eigenvectors of the graph Laplacian to place cluster centroids. Based on this result Vilcek proposed an autoencoder structure for unsupervised detection of communities in social networks [50].

The connections between two neural network layers decide which features of the first layers affect which features of the second layer. In case of fully connected neural networks, all the neurons in the first layer can affect the activation of each neuron in the second layer. It is decided during training, which connections are neglected and which are of greater importance, by tuning the weights associated with the connections. Connections get neglected because not all first layer features are necessary for the computation of a given second layer feature [51, 41]. For unsupervised learning, the tuning of the weights this way is not always possible, and it is also computationally ineffective. Unsupervised clustering can be used to design the connections between neural network layers [51, 41].

Bruna et al. proposed a generalization of deep convolutional architectures (locally connected networks) based on analogies with graph theory and introduced a hierarchical and a spectral construction method for convolutional structures [52]. Experiments were carried out on a downsampled MNIST dataset, where the proposed method was able to achieve equal or lower classification error than a fully connected network that had more than twice the amount of parameters. In [53] another study on the topic of spectral methods for deep learning is presented.

Fuzzy rules can be extracted from the collected data with the help of deep learning [54]. In [54] a method is proposed for extracting fuzzy rules from the data by feeding it to a restricted Boltzmann machine and applying a probability based clustering method (similar to the expectation maximization algorithm) to form the fuzzy rules.

DFuzzy is a deep learning based fuzzy clustering approach for graph partitioning [55]. DFuzzy enables vertexes to belong to multiple clusters with different degree. An autoencoder structure is used to create graph partitions that can be mapped to vertexes by the decoder. An initial clustering of the graph is performed with the PageRank algorithm [56].

The NDT (neural decision tree) is a hybrid architecture of a decision tree and multilayer neural networks [57]. At each node in the decision tree, the splitting is implemented by a neural network. Describing the structure as a whole and assuming shared weights, enable the optimization of the whole architecture globally. The authors compared the test set accuracy of the NDT, a decision tree and a neural network on 14 different datasets, and found that none of these methods have a significant advantage over the other, but the NDT model accuracy is in the top two on 13 of the 14 datasets.

Patel et al. proposed a probabilistic framework for deep learning [58]. Their proposed model, the Deep Rendering Mixture Model (DRMM) can be optimized with the expectation maximization algorithm. The method was introduced as an alternative for deep convolutional neural networks and their optimization with backpropagation. The model can be trained in an unsupervised and semi-supervised manner. The experiments show that the best performing DRMM architectures were able to achieve a test error rate of 0.57%, 1.56%, 1.67% and 0.91% in a semi-supervised scenario for the MNIST dataset with 100, 600, 1000 and 3000 labeled examples respectively.

Most of the convolutional and generative models for comparison had error rates that are nearly two times greater than these.

In [59] a spatial mixture model was proposed for the unsupervised identification of entities in the input. In the mixture model, each entity is described by a neural network with a given set of parameters. Based on the expectation maximization algorithm, an unsupervised clustering method is introduced, that enables optimization by differentiation. The basic idea behind this approach is similar to the neural decision tree [57], but instead of a decision tree, a mixture model is created.

A detailed description of the autoencoder structure, its role in unsupervised learning and place in deep learning, along with different types of autoencoders can be found in [36].

Radford et al. introduced the Deep Convolutional Generative Adversarial Network structure (DCGAN) [60]. In their work, they showed that deep convolutional models, as generative networks, can extract useful features from the presented images in an unsupervised manner. Their result shows that both the generator and the discriminator network can be trained to extract general features and thus they can be used for other purposes as well, for example as feature extractors.

Summary

The current results show that deep learning can benefit a lot from unsupervised clustering methods. The applications that utilize unsupervised learning in the process of deep learning, perform well in many cases [39, 41, 46, 48, 52, 57, 58] and can have other advantages, like fast training and inference, smaller memory needs and easy implementation due to the lack of labeling. This paper promotes the use of unsupervised techniques in the field of deep learning and argues that a significant aspect of deep learning research should be to find ways to exploit the information provided by the sheer data better, rather than acquiring more and more data in order to build even more complex models to enhance performance.

Acknowledgement

Róbert Fullér has been partially supported by FIEK program (Center for Cooperation between Higher Education and the Industries at the Széchenyi István University, GINOP-2.3.4-15-2016-00003) and by the EFOP-3.6.1-16-2016-00010 project.

References

- [1] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge, 1998.
- [2] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436, 2015.
- [3] Allan Hanbury. A survey of methods for image annotation. *Journal of Visual Languages & Computing*, 19(5):617–627, 2008.

- [4] Yoshua Bengio, Aaron C Courville, and Pascal Vincent. Unsupervised feature learning and deep learning: A review and new perspectives. *CoRR*, *abs/1206.5538*, 1:2012, 2012.
- [5] Rui Xu and D. Wunsch. Survey of clustering algorithms. *IEEE Transactions on Neural Networks*, 16(3):645–678, May 2005.
- [6] P. Berkhin. A Survey of Clustering Data Mining Techniques. In Jacob Kogan, Charles Nicholas, and Marc Teboulle, editors, *Grouping Multidimensional Data: Recent Advances in Clustering*, pages 25–71. Springer Berlin Heidelberg, Berlin, Heidelberg, 2006.
- [7] Anil K Jain, M Narasimha Murty, and Patrick J Flynn. Data clustering: a review. *ACM computing surveys (CSUR)*, 31(3):264–323, 1999.
- [8] A. I. Károly, J. Kuti, and P. Galambos. Unsupervised real-time classification of cycle stages in collaborative robot applications. In *2018 IEEE 16th World Symposium on Applied Machine Intelligence and Informatics (SAMII)*, pages 97–102, Feb 2018.
- [9] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In *Intelligent Robots and Systems (IROS), 2017 IEEE/RSJ International Conference on*, pages 23–30. IEEE, 2017.
- [10] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [11] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- [12] Bernhard Schölkopf. The kernel trick for distances. In *Advances in neural information processing systems*, pages 301–307, 2001.
- [13] James MacQueen et al. Some methods for classification and analysis of multivariate observations. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, volume 1, pages 281–297. Oakland, CA, USA, 1967.
- [14] Xindong Wu, Vipin Kumar, J Ross Quinlan, Joydeep Ghosh, Qiang Yang, Hiroshi Motoda, Geoffrey J McLachlan, Angus Ng, Bing Liu, S Yu Philip, et al. Top 10 algorithms in data mining. *Knowledge and information systems*, 14(1):1–37, 2008.
- [15] Tapas Kanungo, David M Mount, Nathan S Netanyahu, Christine D Piatko, Ruth Silverman, and Angela Y Wu. An efficient k-means clustering algorithm:

- Analysis and implementation. *IEEE transactions on pattern analysis and machine intelligence*, 24(7):881–892, 2002.
- [16] Catherine A Sugar and Gareth M James. Finding the number of clusters in a dataset: An information-theoretic approach. *Journal of the American Statistical Association*, 98(463):750–763, 2003.
- [17] Aristidis Likas, Nikos Vlassis, and Jakob J Verbeek. The global k-means clustering algorithm. *Pattern recognition*, 36(2):451–461, 2003.
- [18] Edo Liberty, Ram Sriharsha, and Maxim Sviridenko. An algorithm for online k-means clustering. *CoRR*, abs/1412.5721, 2014.
- [19] Isis Bonet, Adriana Escobar, Andrea Mesa-Múnera, and Juan Fernando Alzate. Clustering of Metagenomic Data by Combining Different Distance Functions. *Acta Polytechnica Hungarica*, 14(3), October 2017.
- [20] Hinrich Schütze, Christopher D Manning, and Prabhakar Raghavan. *Introduction to information retrieval*, volume 39. Cambridge University Press, 2008.
- [21] J. C. Dunn. A fuzzy relative of the isodata process and its use in detecting compact well-separated clusters. *Journal of Cybernetics*, 3(3):32–57, 1973.
- [22] Wang Peizhuang. Pattern recognition with fuzzy objective function algorithms (james c. bezdek). *SIAM Review*, 25(3):442, 1983.
- [23] Ernesto Moya-Albor, Hiram Ponce, and Jorge Brieva. An Edge Detection Method using a Fuzzy Ensemble Approach. *Acta Polytechnica Hungarica*, 14(3):20, 2017.
- [24] Asa Ben-Hur, David Horn, Hava T Siegelmann, and Vladimir Vapnik. Support vector clustering. *Journal of machine learning research*, 2(Dec):125–137, 2001.
- [25] Bernhard Schölkopf, Robert C Williamson, Alex J Smola, John Shawe-Taylor, and John C Platt. Support vector method for novelty detection. In *Advances in neural information processing systems*, pages 582–588, 2000.
- [26] David MJ Tax and Robert PW Duin. Support vector domain description. *Pattern recognition letters*, 20(11-13):1191–1199, 1999.
- [27] Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine learning*, 20(3):273–297, 1995.
- [28] Ulrike von Luxburg. A tutorial on spectral clustering. *CoRR*, abs/0711.0189, 2007.
- [29] S Rasoul Safavian and David Landgrebe. A survey of decision tree classifier methodology. *IEEE transactions on systems, man, and cybernetics*, 21(3):660–674, 1991.

- [30] Jayanta Basak and Raghu Krishnapuram. Interpretable hierarchical clustering by constructing an unsupervised decision tree. *IEEE transactions on knowledge and data engineering*, 17(1):121–132, 2005.
- [31] Todd K Moon. The expectation-maximization algorithm. *IEEE Signal processing magazine*, 13(6):47–60, 1996.
- [32] Teuvo Kohonen and Timo Honkela. Kohonen network. *Scholarpedia*, 2(1):1568, 2007.
- [33] Teuvo Kohonen. *Self-organization and associative memory*, volume 8. Springer Science & Business Media, 2012.
- [34] Stephen Grossberg. Adaptive resonance theory. Technical report, Boston University Center for Adaptive Systems and Department of Cognitive and Neural Systems, 2000.
- [35] Gail A Carpenter, Stephen Grossberg, and David B Rosen. Art 2-a: An adaptive resonance algorithm for rapid category learning and recognition. *Neural networks*, 4(4):493–504, 1991.
- [36] Pierre Baldi. Autoencoders, unsupervised learning, and deep architectures. In *Proceedings of ICML Workshop on Unsupervised and Transfer Learning*, pages 37–49, 2012.
- [37] Xiu-Shen Wei, Chen-Lin Zhang, Jianxin Wu, Chunhua Shen, and Zhi-Hua Zhou. Unsupervised object discovery and co-localization by deep descriptor transforming. *CoRR*, abs/1707.06397, 2017.
- [38] Tom White. Sampling generative networks: Notes on a few effective techniques. *CoRR*, abs/1609.04468, 2016.
- [39] Richard Socher, Brody Huval, Bharath Bath, Christopher D Manning, and Andrew Y Ng. Convolutional-recursive deep learning for 3d object classification. In *Advances in Neural Information Processing Systems*, pages 656–664, 2012.
- [40] Eugenio Culurciello, Jordan Bates, Aysegul Dundar, José Antonio Pérez-Carrasco, and Clément Farabet. Clustering learning for robotic vision. *CoRR*, abs/1301.2820, 2013.
- [41] Aysegul Dundar, Jonghoon Jin, and Eugenio Culurciello. Convolutional clustering for unsupervised learning. *CoRR*, abs/1511.06241, 2015.
- [42] K. Lai, L. Bo, X. Ren, and D. Fox. A large-scale hierarchical multi-view rgb-d object dataset. In *2011 IEEE International Conference on Robotics and Automation*, pages 1817–1824, May 2011.
- [43] L. Bo, X. Ren, and D. Fox. Depth kernel descriptors for object recognition. In *2011 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 821–826, Sept 2011.

- [44] M. Blum, Jost Tobias Springenberg, J. Wülfing, and M. Riedmiller. A learned feature descriptor for object recognition in rgb-d data. In *2012 IEEE International Conference on Robotics and Automation*, pages 1298–1303, May 2012.
- [45] Liefeng Bo, Xiaofeng Ren, and Dieter Fox. Unsupervised feature learning for rgb-d based object recognition. In *Experimental Robotics*, pages 387–402. Springer, 2013.
- [46] Adam Coates and Andrew Y Ng. Learning feature representations with k-means. In *Neural networks: Tricks of the trade*, pages 561–580. Springer, 2012.
- [47] M. Dundar, Q. Kou, B. Zhang, Y. He, and B. Rajwa. Simplicity of kmeans versus deepness of deep learning: A case of unsupervised feature learning with limited data. In *2015 IEEE 14th International Conference on Machine Learning and Applications (ICMLA)*, pages 883–888, Dec 2015.
- [48] Bo Yang, Xiao Fu, Nicholas D. Sidiropoulos, and Mingyi Hong. Towards k-means-friendly spaces: Simultaneous deep learning and clustering. *CoRR*, abs/1610.04794, 2016.
- [49] Fei Tian, Bin Gao, Qing Cui, Enhong Chen, and Tie-Yan Liu. Learning deep representations for graph clustering. In *AAAI*, pages 1293–1299, 2014.
- [50] Alexandre Vilcek. Deep learning with k-means applied to community detection in network. Project Report CS224W-31, Stanford University Center for Professional Development, 2014.
- [51] Eugenio Culurciello, Jonghoon Jin, Aysegul Dundar, and Jordan Bates. An analysis of the connections between layers of deep neural networks. *CoRR*, abs/1306.0152, 2013.
- [52] Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann LeCun. Spectral networks and locally connected networks on graphs. *CoRR*, abs/1312.6203, 2013.
- [53] L. Shao, D. Wu, and X. Li. Learning deep and wide: A spectral method for learning deep networks. *IEEE Transactions on Neural Networks and Learning Systems*, 25(12):2303–2308, Dec 2014.
- [54] Erick De la Rosa and Wen Yu. Data-driven fuzzy modeling using deep learning. *CoRR*, abs/1702.07076, 2017.
- [55] Vandana Bhatia and Rinkle Rani. Dfuzzy: a deep learning-based fuzzy clustering model for large graphs. *Knowledge and Information Systems*, pages 1–23, 2018.
- [56] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The pagerank citation ranking: Bringing order to the web. Technical Report 1999-66, Stanford InfoLab, November 1999. Previous number = SIDL-WP-1999-0120.

-
- [57] Randall Balestriero. Neural decision trees. *CoRR*, abs/1702.07360, 2017.
- [58] Ankit B Patel, Minh Tan Nguyen, and Richard Baraniuk. A probabilistic framework for deep learning. In *Advances in neural information processing systems*, pages 2558–2566, 2016.
- [59] Klaus Greff, Sjoerd van Steenkiste, and Jürgen Schmidhuber. Neural expectation maximization. *CoRR*, abs/1708.03498, 2017.
- [60] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *CoRR*, abs/1511.06434, 2015.

Design of a Single-Master/Multi-Slave Nonlinear Teleoperation System through State Convergence with Time Varying Delays

Umar Farooq^{1,3}, Muhammad Usman Asad¹, Jason Gu¹, Ghulam Abbas⁴, Valentina E. Balas², Marius M. Balas²

¹Department of Electrical and Computer Engineering Dalhousie University, Halifax, N.S. B3H 4R2, Canada

²Department of Automatics and Applied Software, “Aurel Vlaicu” University of Arad, Romania

³Department of Electrical Engineering, University of The Punjab, Quaid-e-Azam Campus, Lahore, 54590 Pakistan

⁴Department of Electrical Engineering, The University of Lahore, Pakistan

umar.farooq@dal.ca, usman.asad@dal.ca, jason.gu@dal.ca,
ghulam.abbas@ee.uol.edu.pk, valentina.balas@uav.ro, marius.balas@uav.ro

Abstract: This paper presents the design of a nonlinear teleoperation system which is comprised of a single master and multiple slave (SM/MS) units. The interaction between these units follows the extended state convergence architecture which allows multiple linear master units to influence multiple linear slave units. However, in this study, the nonlinear dynamics of the master and slave units is considered and the resulting nonlinear teleoperation system is analyzed in the presence of time delays. To be specific, the following objectives are defined: (i) the nonlinear teleoperation remains stable in the presence of time varying delays, (ii) the slave units follow the position commands of the master unit and (iii) the operator receives a force feedback proportional to the interaction forces of the slaves with their environments. Towards this end, Lyapunov-Krasovskii theory is utilized which provides guidelines to select the control gains of the extended state convergence architecture such that the aforementioned objectives are achieved. The efficacy of the proposed scheme is finally verified through simulations in MATLAB/Simulink environment by considering a two degrees-of-freedom (DoF) single-master/tri-slave nonlinear teleoperation system.

Keywords: Teleoperation; nonlinear dynamics; state convergence; MATLAB

1 Introduction

Teleoperation refers to the control of a distant process and has found diverse applications ranging from miniaturized medical procedures to large-scale industrial processes. It is usually accomplished through the use of master and slave robotic devices which are connected through a communication channel. Based on the number of these robotic devices, teleoperation systems can be classified as either bilateral or multilateral systems. In a typical bilateral teleoperation system, human operator drives the master manipulator and the resulting motion commands are transmitted across the communication channel towards the slave manipulator which performs the desired task at the remote site. A force feedback is also provided by the slave manipulator to improve human's perception of the remote environment. By deploying more than one slave manipulator, the task can be carried out more efficiently. The teleoperation system in such a setting is known as single-master/multi-slave system and is one of the topologies in a broader class of multilateral systems. Other arrangements in this category include dual user systems for training tasks, and multi-master/single-slave and multi-master/multi-slave systems for collaborative missions [1]-[3].

All these forms of teleoperation need an effective control system to achieve the required task. An ideal control algorithm should be able to ensure that the teleoperation system remains stable against the time delays of the communication channel while providing a superior position and force tracking performance under systems' uncertainties. This is a challenging task since stability and transparency (the position and force tracking requirement is collectively referred as transparency) are two conflicting objectives and the presence of uncertainties complicates the problem further. Many research efforts have been directed to address these performance issues in teleoperation systems. Passivity schemes are popular in research community as they transform the delay-vulnerable system into a delay-resilient one [3]-[11]. However, transparency of the teleoperation system is sacrificed during this transformation process especially when large time delays exist, for which some modifications to the standard passivity algorithms have also been proposed [12]. To ensure that the teleoperation system performs well under uncertainties, non-passive algorithms based on H_∞ [13], [14], sliding mode [15]-[18] and adaptive control [19]-[21] theories are also proposed. However, time delay appears to be a limiting factor in the complete success of these algorithms. The use of intelligent control techniques such as fuzzy logic [22]-[26] and neural networks [27], [28] has also been investigated. Encouraging results are reported based on the combination of neural networks and passivity algorithms [29], [30].

State convergence is another novel scheme which provides an elegant design procedure to determine control gains for bilateral teleoperation systems modeled on state space [31]. It was originally proposed for linear systems with small time delay in the communication channel. Later, the applicability of the scheme to nonlinear teleoperation systems suffering from time-varying delays was

demonstrated through the use of Lyapunov-Krasovskii functional [32]. In our earlier work, we have employed the state convergence scheme to control a nonlinear teleoperation system which can be approximated by a class of Takagi-Sugeno fuzzy models. We have also extended this scheme to design controllers for multiple linear one DoF teleoperation systems [33].

This paper builds on our earlier framework of [33] and discusses the design of a multi-DoF SM/MS nonlinear teleoperation system in the presence of time varying delays. To the best of authors' knowledge, state convergence based design of SM/MS nonlinear teleoperation system has not been discussed in the literature. Further, the earlier methodology on the control of nonlinear bilateral teleoperation system through state convergence [32] has become a special case of the proposed multilateral controller. To proceed, we first define the control objectives to be the synchronization of master and slave position signals along with the mixed force reflection to the operator from the slave environments. Then, to achieve these objectives, Lyapunov-Krasovskii theory is utilized to design the control gains of the extended state convergence architecture following the lines of [32]. The proposed methodology is finally verified through MATLAB simulations on a 2-DoF single-master/tri-slave nonlinear teleoperation system in the presence of time delays.

This paper is structured as follows: We start by presenting the modeling of SM/MS teleoperation system in Section 2. Preliminaries are included in Section 3 while control objectives are listed in Section 4. Stability analysis and control design is discussed in Section 5. Simulation results are presented in Section 6. Finally, conclusions are drawn in Section 7.

2 Modeling of the SM/MS Teleoperation System

We consider a nonlinear teleoperation system which is comprised of n -DoF single master and l -slave manipulators/units with the following dynamics:

$$M_m(q_m)\ddot{q}_m + C_m(q_m, \dot{q}_m)\dot{q}_m + g_m(q_m) = \tau_m + F_h \quad (1)$$

$$M_s^i(q_s^i)\ddot{q}_s^i + C_s^i(q_s^i, \dot{q}_s^i)\dot{q}_s^i + g_s^i(q_s^i) = \tau_s^i - F_e^i, \forall i = 1, 2, \dots, l \quad (2)$$

Where $(M_m, M_s^i) \in \mathbb{R}^{n \times n}$, $(C_m, C_s^i) \in \mathbb{R}^{n \times n}$, $(g_m, g_s^i) \in \mathbb{R}^{n \times 1}$, $(q_m, q_s^i) \in \mathbb{R}^{n \times 1}$, $(\dot{q}_m, \dot{q}_s^i) \in \mathbb{R}^{n \times 1}$, $(\ddot{q}_m, \ddot{q}_s^i) \in \mathbb{R}^{n \times 1}$ denote inertia matrices, coriolis/centrifugal matrices, gravity vectors, joint positions, joint velocities, joint

accelerations, and input torques of master and slave units respectively. Operator's forces are assumed to be constant while environments are assumed to be passive in this study. It is also assumed that the environments can be modeled as spring-damper systems, i.e. $F_e^i = K_e^i q_s^i + B_e^i \dot{q}_s^i$ where $K_e^i, B_e^i \in \mathbb{R}^{n \times n}$ are positive definite diagonal matrices.

Now, the communication between the master and slave units is established through the use of extended state convergence architecture. This communication framework is shown in Figure 1 and is comprised of the following parameters:

$K_m = \begin{bmatrix} K_{m1} & K_{m2} \end{bmatrix} \in \mathbb{R}^{n \times 2n}$: This parameter is the stabilizing feedback gain matrix for the master unit. Each of its constituent members ($K_{m1}, K_{m2} \in \mathbb{R}^{n \times n}$) is an unknown but negative definite diagonal matrix and will be found through Lyapunov-Krasovskii based design procedure.

$K_s^i = \begin{bmatrix} K_{s1}^i & K_{s2}^i \end{bmatrix} \in \mathbb{R}^{n \times 2n}$: This parameter is the stabilizing feedback gain matrix for the i^{th} slave unit. Each of its constituent members ($K_{s1}^i, K_{s2}^i \in \mathbb{R}^{n \times n}$) is an unknown but negative definite diagonal matrix and will be found through Lyapunov-Krasovskii based design procedure.

$R_s^i = \begin{bmatrix} R_{s1}^i & R_{s2}^i \end{bmatrix} \in \mathbb{R}^{n \times 2n}$: This parameter models the influence of the master unit's motion onto the i^{th} slave unit's motion. Each of its constituent members ($R_{s1}^i, R_{s2}^i \in \mathbb{R}^{n \times n}$) is an unknown but positive definite diagonal matrix and will be found through Lyapunov-Krasovskii based design procedure.

$R_m^i = \begin{bmatrix} R_{m1}^i & R_{m2}^i \end{bmatrix} \in \mathbb{R}^{n \times 2n}$: This parameter models the influence of the i^{th} slave unit's motion onto the master unit's motion. Each of its constituent members ($R_{m1}^i, R_{m2}^i \in \mathbb{R}^{n \times n}$) is an unknown but positive diagonal matrix and will be found through Lyapunov-Krasovskii based design procedure.

$G_2^i \in \mathbb{R}^+$: This parameter models the influence of the operator's applied force onto the i^{th} slave unit.

$T_{mi}(t) \in \mathbb{R}^+$: This represents the time delay on the link which connects the i^{th} slave unit to the master unit. In this study, only the upper bounds on these time delays (T_{mi}^+) are known.

$T_{si}(t) \in \mathbb{R}^+$: This represents the time delay on the link which connects the master unit to the i^{th} slave unit. In this study, only the upper bounds on these time delays (T_{si}^+) are known.

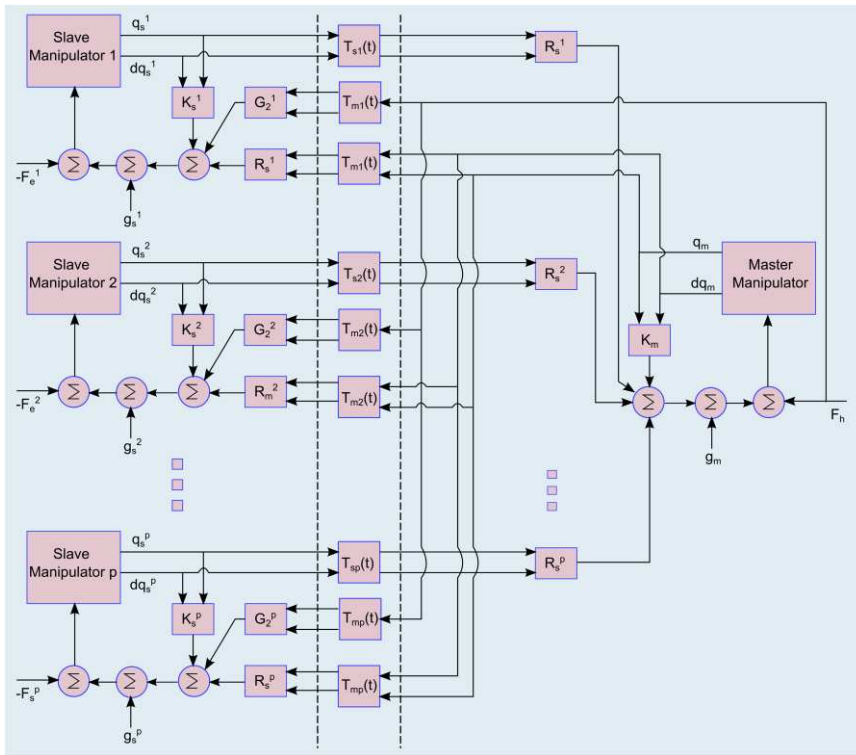


Figure 1

Single-master/multi-slave teleoperation system through state convergence

3 Preliminaries

3.1 Properties of Master/Slave Units

The master and slave units as modeled by (1),(2) possess the following properties:

(P1) The inertia matrices are symmetric, positive definite and bounded, i.e. there exist positive constants β_l and β_u such that $0 < \beta_l I < M(q) < \beta_u I < \infty$.

(P2) A skew-symmetric relation exists between the inertia and coriolis/centrifugal matrices such that $x^T \left(\dot{M}(q) - 2C(q, \dot{q}) \right) x = 0, \forall x \in \mathbb{R}^n$.

(P3) The coriolis/centrifugal force vectors are bounded i.e., there exists positive constant β_f such that $\left\| C(q, \dot{q})\dot{q} \right\| \leq \beta_f \left\| \dot{q} \right\|$.

(P4) If the joint variables q and \dot{q} are bounded, then the time derivative of coriolis/centrifugal matrices is also bounded.

3.2 Lemmas

For any vectors $x, y \in \mathbb{R}^n$, positive definite diagonal matrix $F \in \mathbb{R}^{n \times n}$, scalar $\gamma > 0$ and variable time delay $T_i(t)$ having upper bound T_i^+ , the following inequalities hold:

$$(L1) \quad -2 \int_0^{t_f} x^T F \int_0^{T_i(t)} y(t-\sigma) d\sigma dt \leq \gamma \int_0^{t_f} x^T F x dt + \frac{T_i^{+2}}{\gamma} \int_0^{t_f} y^T F y dt$$

$$(L2) \quad x(t - T_i(t)) - x(t) = \int_0^{T_i(t)} \dot{x}(t - \sigma) d\sigma \leq T_i^{+ \frac{1}{2}} \left\| \dot{x} \right\|_2$$

4 Control Objectives

Besides establishing the stability, we intend to achieve the following objectives in SM/MS nonlinear teleoperation system:

Control Objective # 1: During the free motion, the joint positions of all the slave units should converge to the corresponding joint positions of the master unit in steady state i.e. $\lim_{t \rightarrow \infty} \|q_s^i(t) - q_m(t)\| = 0, \forall i = 1, 2, \dots, l$

Control Objective # 2: During the contact motion, operator should feel a force proportional to the aggregated environmental forces, i.e. $F_h \propto \sum_{i=1}^l F_e^i$

5 Stability Analysis and Control Design

Consider the SM/MS teleoperation system of Fig. 1 with time varying delays in the communication channel. The control inputs for the master and slave units in this time-delayed teleoperation system are:

$$\tau_m = g_m(q_m) + K_{m1}q_m + K_{m2}q_m + \sum_{i=1}^l R_{m1}^i q_s^i(t - T_{si}(t)) + \sum_{i=1}^l R_{m2}^i q_s^i(t - T_{si}(t)) \quad (3)$$

$$\begin{aligned} \tau_s^i = & g_s^i(q_s^i) + K_{s1}^i q_s^i + K_{s2}^i q_s^i + R_{s1}^i q_m(t - T_{mi}(t)) + R_{s2}^i q_m(t - T_{mi}(t)) + \\ & G_2^i F_h(t - T_{mi}(t)), \forall i = 1, 2, \dots, l \end{aligned} \quad (4)$$

By plugging (3) in (1) and (4) in (2), we obtain the closed loop dynamics of the master and slave units as:

$$M_m \ddot{q}_m + C_m \dot{q}_m = K_{m1}q_m + K_{m2}q_m + \sum_{i=1}^l R_{m1}^i q_s^i(t - T_{si}(t)) + \sum_{i=1}^l R_{m2}^i q_s^i(t - T_{si}(t)) + F_h \quad (5)$$

$$\begin{aligned} M_s^i \ddot{q}_s^i + C_s^i \dot{q}_s^i = & K_{s1}^i q_s^i + K_{s2}^i q_s^i + R_{s1}^i q_m(t - T_{mi}(t)) + R_{s2}^i q_m(t - T_{mi}(t)) + \\ & G_2^i F_h(t - T_{mi}(t)) - F_e^i, \forall i = 1, 2, \dots, l \end{aligned} \quad (6)$$

In equilibrium points for master and slave units, we have:

$$\begin{aligned} q_m = q_m(t - T_{mi}(t)) = \bar{q}_m, q_m = q_m = 0 \\ q_s^i = q_s^i(t - T_{si}(t)) = \bar{q}_s^i, q_s^i = q_s^i = 0 \end{aligned} \quad (7)$$

Considering the environmental models and evaluating (6), (7) at equilibrium, we have:

$$\begin{aligned} 0 = & K_{m1} \bar{q}_m + \sum_{i=1}^l R_{m1}^i \bar{q}_s^i + F_h \\ 0 = & K_{s1}^i \bar{q}_s^i + R_{s1}^i \bar{q}_m + G_2^i F_h(t - T_{mi}(t)) - K_e^i \bar{q}_s^i, \forall i = 1, 2, \dots, l \end{aligned} \quad (8)$$

Let us now analyze the closed loop teleoperation system of (5), (6) in a new coordinate system formed by the variables q_m, q_s^i and their time delayed versions $q_m(t - T_{mi}(t)), q_s^i(t - T_{si}(t))$ as defined below:

$$q_m = q_m - \bar{q}_m \quad (9)$$

$$q_m(t - T_{mi}(t)) = q_m(t - T_{mi}(t)) - \bar{q}_m \quad (10)$$

$$q_s^i = q_s^i - \bar{q}_s^i \quad (11)$$

$$q_s^i(t - T_{si}(t)) = q_s^i(t - T_{si}(t)) - \bar{q}_s^i \quad (12)$$

By using (9)-(12) with (5)-(8), we obtain the transformed teleoperation system as:

$$M_m \ddot{q}_m + C_m \dot{q}_m = K_{m1} q_m + K_{m2} \dot{q}_m + \sum_{i=1}^l R_{m1}^i q_s^i(t - T_{s1}^i(t)) + \sum_{i=1}^l R_{m2}^i \dot{q}_s^i(t - T_{s1}^i(t)) \quad (13)$$

$$M_s^i \ddot{q}_s^i + C_s^i \dot{q}_s^i = K_{s1}^i q_m^i + K_{s2}^i \dot{q}_m^i + R_{s1}^i q_m^i(t - T_{m1}^i(t)) + R_{s2}^i \dot{q}_m^i(t - T_{m1}^i(t)) - K_e^i q_s^i - B_e^i \dot{q}_s^i, \forall i = 1, 2, \dots, l \quad (14)$$

Now we study the asymptotic stability and position coordination behavior of the time-delayed teleoperation system in Theorem 1 while the force reflection behavior is discussed in Theorem 2.

Theorem 1: The origin of the transformed time-delayed teleoperation system (13), (14) is asymptotically stable and the position coordination between the master and slave units is achieved in free motion when the control gains of (15), (16) are used and $l+1$ inequalities in (17), (18) are also satisfied.

$$K_{m1} = -lK, K_{m2} = -(l+1)K_1 - \sum_{i=1}^l K_{md}^i$$

$$K_{s1}^i = -K, K_{s2}^i = -2K_1 - K_{sd}^i, \forall i = 1, 2, \dots, l \quad (15)$$

$$R_{m1}^i = R_{s1}^i = K, R_{m2}^i = 2K_{md}^i, R_{s2}^i = 2K_{sd}^i, \forall i = 1, 2, \dots, l$$

$$K_{md}^i = \left(1 - T_{sj}^i(t)\right) K_1, K_{sd}^i = \left(1 - T_{mj}^i(t)\right) K_1 \quad (16)$$

$$K_1 - \frac{\gamma_{mj}}{2} K - \frac{T_{sj}^{i+2}}{2\gamma_{sj}} K > 0, \forall i = 1, 2, \dots, l \quad (17)$$

$$K_1 - \sum_{i=1}^l \frac{\gamma_{sj}}{2} K - \sum_{i=1}^l \frac{T_{mj}^{i+2}}{2\gamma_{mj}} K > 0 \quad (18)$$

Where, γ_{sj}, γ_{mj} are positive constants, $K, K_1 \in \mathbb{R}^{n \times n}$ are positive definite diagonal matrices, T_{sj}^i, T_{mj}^i are the time derivatives of communication delays which are assumed to be less than unity. Therefore, $K_{sd}^i, K_{md}^i \in \mathbb{R}^{n \times n}$ are also positive definite diagonal matrices.

Proof: Consider the following Lyapunov-Krasovskii functional candidate:

$$\begin{aligned}
V \left(q_m, q_s^i, q_s^i - q_m, q_s^i \right) &= \frac{1}{2} q_m^T M_m q_m + \frac{1}{2} \sum_{i=1}^l q_s^{iT} M_s^i q_s^i + \\
\frac{1}{2} \sum_{i=1}^l \left(q_s^i - q_m \right)^T K \left(q_s^i - q_m \right) &+ \sum_{i=1}^l \int_{t-T_{mj}(t)}^t q_m^T(\eta) K_1 q_m(\eta) d\eta + \\
\sum_{i=1}^l \int_{t-T_{sj}(t)}^t q_s^{iT}(\eta) K_1 q_s^i(\eta) d\eta &+ \frac{1}{2} \sum_{i=1}^l q_s^{iT} K_e^i q_s^i
\end{aligned} \tag{19}$$

By taking the time derivative of (19) along the trajectories of the teleoperation system (13), (14) and using the property P2 of the master and slave units, we obtain:

$$\begin{aligned}
\dot{V} &= q_m^T \left(K_{m1} q_m + \sum_{i=1}^l R_{m1}^i q_s^i(t-T_{si}(t)) + K_{m2} q_m + \sum_{i=1}^l R_{m2}^i q_s^i(t-T_{si}(t)) \right) + \\
\sum_{i=1}^l q_s^{iT} \left(K_{s1}^i q_s^i + K_{s2}^i q_s^i + R_{s1}^i q_m(t-T_{mi}(t)) + R_{s2}^i q_m(t-T_{mi}(t)) - K_e^i q_s^i - B_e^i q_s^i \right) &+ \\
\sum_{i=1}^l \left(q_s^{iT} K q_s^i - q_s^{iT} K q_m - q_m^T K q_s^i + q_m^T K q_m + q_s^{iT} K_e^i q_s^i \right) &+ \\
\sum_{i=1}^l \left(\frac{d}{dt} q_m^T K_1 q_m(t-T_{mi}(t)) \left(1 - T_{mi}(t) \right) K_1 q_m(t-T_{mi}(t)) \right) &+ \\
\sum_{i=1}^l \left(\frac{d}{dt} q_s^{iT} K_1 q_s^i(t-T_{si}(t)) \left(1 - T_{si}(t) \right) K_1 q_s^i(t-T_{si}(t)) \right)
\end{aligned} \tag{20}$$

By grouping the terms in (20) and using the definition of the time varying matrices (16), we have:

$$\begin{aligned}
\dot{V} &= q_m^T \left(K_{m1} + lK \right) q_m + \sum_{i=1}^l q_m^T \left(R_{m1}^i q_s^i(t-T_{si}(t)) - K q_s^i \right) + \\
\sum_{i=1}^l q_s^{iT} \left(K_{s1}^i + K \right) q_s^i + \sum_{i=1}^l q_s^{iT} \left(R_{s1}^i q_m(t-T_{mi}(t)) - K q_m \right) &+ \\
\sum_{i=1}^l \left(q_m^T \left(\frac{K_{m2}}{l} + K_1 \right) q_m + q_m^T R_{m2}^i q_s^i(t-T_{si}(t)) - q_s^{iT} \left(t-T_{si}(t) \right) K_{md}^i q_s^i(t-T_{si}(t)) \right) &+ \\
\sum_{i=1}^l \left(q_s^{iT} \left(K_{s2}^i + K_1 - B_e^i \right) q_s^i + q_s^{iT} R_{s2}^i q_m(t-T_{mi}(t)) - \frac{d}{dt} q_m^T(t-T_{mi}(t)) K_{sd}^i q_m(t-T_{mi}(t)) \right)
\end{aligned} \tag{21}$$

Let us now define the following position error signals:

$$\begin{aligned} e_{q_m}^i &= q_m - q_s^i(t - T_{si}(t)) \\ e_{q_s}^i &= q_s^i - q_m(t - T_{mi}(t)) \end{aligned} \quad (22)$$

By substituting the control gains of (15) in (21) and using the time derivative of (22) in the resulting expression, we obtain:

$$\begin{aligned} \dot{V} &= \sum_{i=1}^l q_m^{\square T} K \left(q_s^i(t - T_{si}(t)) - q_s^i \right) - \sum_{i=1}^l e_{q_m}^{\square T} K_{md}^i e_{q_m}^i + \sum_{i=1}^l q_s^{\square T} K \left(q_m(t - T_{mi}(t)) - q_m \right) - \\ &\quad \sum_{i=1}^l e_{q_s}^{\square T} K_{sd}^i e_{q_s}^i - q_m^{\square T} K_1 q_m - \sum_{i=1}^l q_s^{\square T} K_1 q_s^i - \sum_{i=1}^l q_s^{\square T} B_e^i q_s^i \end{aligned} \quad (23)$$

By integrating (23) over the time interval $[0, t_f]$, rewriting first and third terms in integral form and finally using lemma L1, we have:

$$\begin{aligned} \int_0^{t_f} \dot{V} d\eta &\leq \sum_{i=1}^l \left(\frac{\gamma_{si}}{2} \int_0^{t_f} q_m^{\square T} K q_m d\eta + \frac{T_{si}^{+2}}{2\gamma_{si}} \int_0^{t_f} q_s^{\square T} K q_s^i d\eta \right) + \\ &\quad \sum_{i=1}^l \left(\frac{\gamma_{mi}}{2} \int_0^{t_f} q_s^{\square T} K q_s^i d\eta + \frac{T_{mi}^{+2}}{2\gamma_{mi}} \int_0^{t_f} q_m^{\square T} K q_m d\eta \right) - \\ &\quad \sum_{i=1}^l \int_0^{t_f} e_{q_m}^{\square T} K_{md}^i e_{q_m}^i d\eta - \sum_{i=1}^l \int_0^{t_f} e_{q_s}^{\square T} K_{sd}^i e_{q_s}^i d\eta - \\ &\quad \int_0^{t_f} q_m^{\square T} K_1 q_m d\eta - \sum_{i=1}^l \int_0^{t_f} q_s^{\square T} K_1 q_s^i d\eta - \sum_{i=1}^l \int_0^{t_f} q_s^{\square T} B_e^i q_s^i d\eta \end{aligned} \quad (24)$$

The simplification of (24) leads to:

$$\begin{aligned} \int_0^{t_f} \dot{V} d\eta &\leq - \int_0^{t_f} q_m^{\square T} \left(K_1 - \sum_{i=1}^l \frac{\gamma_{si}}{2} K - \sum_{i=1}^l \frac{T_{mi}^{+2}}{2\gamma_{mi}} K \right) q_m d\eta - \\ &\quad \sum_{i=1}^l \int_0^{t_f} q_s^{\square T} \left(K_1 - \frac{\gamma_{mi}}{2} K - \frac{T_{si}^{+2}}{2\gamma_{si}} K \right) q_s^i d\eta - \\ &\quad \sum_{i=1}^l \int_0^{t_f} e_{q_m}^{\square T} K_{md}^i e_{q_m}^i d\eta - \sum_{i=1}^l \int_0^{t_f} e_{q_s}^{\square T} K_{sd}^i e_{q_s}^i d\eta - \sum_{i=1}^l \int_0^{t_f} q_s^{\square T} B_e^i q_s^i d\eta \end{aligned} \quad (25)$$

$$\begin{aligned}
 V(t_f) - V(0) \leq & -\mu \left(K_1 - \sum_{i=1}^l \frac{\gamma_{si}}{2} K - \sum_{i=1}^l \frac{T_{mi}^{+2}}{2\gamma_{mi}} K \right) \left\| q_m \right\|_2^2 - \\
 & \sum_{i=1}^l \mu \left(K_1 - \frac{\gamma_{mi}}{2} K - \frac{T_{si}^{+2}}{2\gamma_{si}} K \right) \left\| q_s^i \right\|_2^2 - \\
 & \sum_{i=1}^l \mu(K_{md}^i) \left\| e_{q_m^i} \right\|_2^2 - \sum_{i=1}^l \mu(K_{sd}^i) \left\| e_{q_s^i} \right\|_2^2 - \sum_{i=1}^l \mu(B_e^i) \left\| q_s^i \right\|_2^2
 \end{aligned} \tag{26}$$

Where $\mu(X)$ denotes the minimal Eigen value of X and the notation $\|x(t)\|_2$ represents the L_2 norm of the signal $x(t)$ in the time interval $[0, t_f]$. Now, if the inequalities in (17), (18) are satisfied and the time derivative of the communication delays remains less than unity, then the right hand side of (26) remains negative. Taking the limit as $t_f \rightarrow \infty$, it can be concluded that the signals

$\left\{ q_m, q_s^i, q_s^i - q_m, q_s^i \right\} \in L_\infty$ and $\left\{ q_m, q_s^i, e_{q_m^i}, e_{q_s^i} \right\} \in L_2$. The boundedness of the

signals $\left\{ q_s^i - q_m, q_s^i \right\}$ implies that q_m is also bounded and therefore $q_m \in L_\infty$. Now,

we study the boundedness of the signals $\left\{ q_m, q_s^i \right\}$. Towards this end, we rewrite

(13) and (14) as:

$$\dot{q}_m = -(M_m)^{-1} \left[C_m q_m - K_{m1} q_m - \sum_{i=1}^l R_{m1}^i q_s^i(t - T_{si}(t)) - K_{m2} q_m - \sum_{i=1}^l R_{m2}^i q_s^i(t - T_{si}(t)) \right] \tag{27}$$

$$\dot{q}_s^i = -(M_s^i)^{-1} \left[C_s^i q_s^i - K_{s1}^i q_s^i - R_{s1}^i q_m(t - T_{mi}(t)) + K_e^i q_s^i - K_{s2}^i q_s^i - R_{s2}^i q_m(t - T_{mi}(t)) + B_e^i q_s^i \right] \tag{28}$$

In (27) and (28), boundedness of the signals, $q_m - q_s^i(t - T_{si}(t)), q_s^i - q_m(t - T_{mi}(t))$ needs to be established in order to draw conclusions on the boundedness of the perturbed acceleration signals. These position error signals can be written as:

$$q_m - q_s^i(t - T_{si}(t)) = \overbrace{q_m - q_s^i}^1 + \overbrace{q_s^i - q_s^i(t - T_{si}(t))}^2 \tag{29}$$

$$q_s^i - q_m^i(t - T_{mi}(t)) = \overbrace{q_s^i - q_m^i}^1 + \overbrace{q_m - q_m(t - T_{mi}(t))}^2 \quad (30)$$

The first terms in (29), (30) have already been shown to be bounded. The second terms in these relations can be re-written using lemma L2 as:

$$\begin{aligned} q_s^i - q_s^i(t - T_{si}(t)) &= \int_0^{T_{si}(t)} q_s^i(t - \sigma) d\sigma \leq T_{si}^{\frac{1}{2}} \left\| \ddot{q}_s^i \right\|_2 \\ q_m - q_m(t - T_{mi}(t)) &= \int_0^{T_{mi}(t)} q_m(t - \sigma) d\sigma \leq T_{mi}^{\frac{1}{2}} \left\| \ddot{q}_m \right\|_2 \end{aligned} \quad (31)$$

It can now be concluded from (31) that the signals $\{q_m - q_s^i(t - T_{si}(t)), q_s^i - q_m^i(t - T_{mi}(t))\} \in L_\infty$. Using the properties P1, P3 of the manipulators and the combined result, $\left\{ \ddot{q}_m, \ddot{q}_s^i, \ddot{q}_s^i - \ddot{q}_m^i, \ddot{q}_m^i - \ddot{q}_s^i(t - T_{si}(t)), \ddot{q}_s^i - \ddot{q}_m^i(t - T_{mi}(t)) \right\} \in L_\infty$, it is established that the perturbed acceleration signals of master and slave units are bounded, i.e. $\left\{ \ddot{q}_m, \ddot{q}_s^i \right\} \in L_\infty$. By Barbalat's lemma, this boundedness of the transformed

acceleration signals in conjunction with the result $\left\{ \ddot{q}_m, \ddot{q}_s^i \right\} \in L_2$ leads to the zero convergence of the perturbed velocity signals, i.e. $\lim_{t \rightarrow \infty} \dot{q}_m = \lim_{t \rightarrow \infty} \dot{q}_s^i = \lim_{t \rightarrow \infty} e_{q_m^i} = \lim_{t \rightarrow \infty} e_{q_s^i} = 0$. Next, we analyze the time derivative of (27) and (28):

$$\begin{aligned} \dot{q}_m &= -\frac{d}{dt} (M_m)^{-1} \begin{bmatrix} C_m \dot{q}_m - K_{m1} q_m - \sum_{i=1}^l R_{m1}^i q_s^i(t - T_{si}(t)) \\ -K_{m2} q_m - \sum_{i=1}^l R_{m2}^i q_s^i(t - T_{si}(t)) \end{bmatrix} \\ &\quad - (M_m)^{-1} \frac{d}{dt} \begin{bmatrix} C_m \dot{q}_m - K_{m1} q_m - \sum_{i=1}^l R_{m1}^i q_s^i(t - T_{si}(t)) \\ -K_{m2} q_m - \sum_{i=1}^l R_{m2}^i q_s^i(t - T_{si}(t)) \end{bmatrix} \end{aligned} \quad (32)$$

$$\begin{aligned} \ddot{q}_s^i = & -\frac{d}{dt} \left(M_s^i \right)^{-1} \begin{bmatrix} C_s^i \dot{q}_s^i - K_{s1}^i q_s^i - R_{s1}^i q_m \left(t - T_{mi} \left(t \right) \right) + K_e^i \dot{q}_s^i \\ -K_{s2}^i \dot{q}_s^i - R_{s2}^i q_m \left(t - T_{mi} \left(t \right) \right) + B_e^i \dot{q}_s^i \end{bmatrix} \\ & - \left(M_s^i \right)^{-1} \frac{d}{dt} \begin{bmatrix} C_s^i \dot{q}_s^i - K_{s1}^i q_s^i - R_{s1}^i q_m \left(t - T_{mi} \left(t \right) \right) + K_e^i \dot{q}_s^i \\ -K_{s2}^i \dot{q}_s^i - R_{s2}^i q_m \left(t - T_{mi} \left(t \right) \right) + B_e^i \dot{q}_s^i \end{bmatrix} \end{aligned} \quad (33)$$

The derivative terms involving inertia matrices in (32), (33) are computed as:

$$\begin{aligned} \frac{d}{dt} \left(M_m \right)^{-1} &= - \left(M_m \right)^{-1} \left(C_m + C_m^T \right) \left(M_m \right) \\ \frac{d}{dt} \left(M_s^i \right)^{-1} &= - \left(M_s^i \right)^{-1} \left(C_s^i + C_s^{iT} \right) \left(M_s^i \right) \end{aligned} \quad (34)$$

The properties P1 and P3 of the master and slave units along with the earlier result

$\left\{ q_m, \dot{q}_s^i, \ddot{q}_m, \dot{q}_s^i \right\} \in L_\infty$ dictate the boundedness of the derivative terms in (34). The

remaining derivative terms in (32), (33) also turn out to be bounded following the application of properties P1, P3, P4 and the earlier result:

$$\left\{ \begin{array}{l} q_m, \dot{q}_s^i, \ddot{q}_s^i - q_m, \dot{q}_s^i, q_m - \dot{q}_s^i \left(t - T_{si} \left(t \right) \right), \\ q_s^i - q_m \left(t - T_{mi} \left(t \right) \right), \ddot{q}_m, \dot{q}_s^i \end{array} \right\} \in L_\infty. \text{ Since all the terms on right hand}$$

sides of (32), (33) are bounded, we have $\left\{ \ddot{q}_m, \ddot{q}_s^i \right\} \in L_\infty$. By using the results,

$$\lim_{t \rightarrow \infty} \ddot{q}_m = \lim_{t \rightarrow \infty} \ddot{q}_s^i = 0 \text{ and } \left\{ \ddot{q}_m, \ddot{q}_s^i \right\} \in L_\infty, \text{ it can be concluded that } \lim_{t \rightarrow \infty} \dot{q}_m = \lim_{t \rightarrow \infty} \dot{q}_s^i = 0.$$

With the zero convergence of perturbed velocity and acceleration signals, the closed loop teleoperation system of (13), (14) in combination with (15) becomes:

$$\lim_{t \rightarrow \infty} \sum_{i=1}^l \left\| q_m - \dot{q}_s^i \left(t - T_{si} \left(t \right) \right) \right\| = 0 \quad (35)$$

$$\lim_{t \rightarrow \infty} \left\| \dot{q}_s^i - q_m \left(t - T_{mi} \left(t \right) \right) \right\| = -K^{-1} K_e^i \dot{q}_s^i \quad (36)$$

The time delay terms in (35) and (36) can be written as:

$$\begin{aligned}
q_s^i(t - T_{si}(t)) &= q_s^i - \int_{t - T_{si}(t)}^t q_s^i d\eta \\
q_m^i(t - T_{mi}(t)) &= q_m^i - \int_{t - T_{mi}(t)}^t q_m^i d\eta
\end{aligned} \tag{37}$$

Since $\lim_{t \rightarrow \infty} q_m^i = \lim_{t \rightarrow \infty} q_s^i = 0$, then the integral terms in (37) disappear. By using this result in (35), (36) and considering the free motion behavior of the teleoperation system, it can be concluded that the perturbations in joint position errors converge to zero, i.e. $\lim_{t \rightarrow \infty} q_m^i = \lim_{t \rightarrow \infty} q_s^i = 0$. Thus the origin of the transformed teleoperation

system $\left\{ q_m^i, q_s^i, q_m^i, q_s^i \right\}$ is asymptotically stable. This further implies that

$\lim_{t \rightarrow \infty} q_m^i = \overline{q_m^i}, \lim_{t \rightarrow \infty} q_s^i = \overline{q_s^i}$. By using these results in the original time-delayed teleoperation system (5), (6), it is found that the position error between the master and slave units is vanished in the absence of operator and environmental forces and the control objective #1 is achieved \square

Remark 1: In case of SM/MS teleoperation system with time varying delays in the communication channel, the control gains for the joint velocities of master and slave units depend on the derivative of time delays as can be seen from (16). These gains are unrealizable since no information about the trajectories of time delays is available except for their upper bounds. In order to overcome this limitation, we transmit extra ramp signals across the communication channel and their time derivatives are computed to realize the velocity control gains as:

$$\begin{aligned}
K_{md}^i &= r(t - T_{sj}(t)) K_1, \forall i = 1, 2, \dots, l \\
K_{sd}^i &= r(t - T_{mj}(t)) K_1, \forall i = 1, 2, \dots, l
\end{aligned} \tag{38}$$

Theorem 2: During the contact motion of the teleoperation system under the control gains of (15), static force is reflected to the operator which is proportional to the aggregated environmental force.

Proof: Consider the steady state behavior of the teleoperation system (1), (2) in the presence of operator and environmental forces. By plugging the control gains (15) in (5), (6), we have:

$$F_h = K \sum_{i=1}^l (\overline{q_m^i} - \overline{q_s^i}) \tag{39}$$

$$F_e^i = K (\overline{q_m^i} - \overline{q_s^i}) + G_2^i F_h \tag{40}$$

By taking the summation ($i = 1$ to $i = l$) on both sides of (40) and using (39), we obtain:

$$F_h = \frac{\sum_{i=1}^l F_e^i}{1 + \sum_{i=1}^l G_2^i} \quad (41)$$

It can be seen from (41) that the operator experiences a static force which is a scaled version of the aggregated environmental force. If all the slave units experience the same force while working in their environments and the force transmission coefficients (G_2^i) are all unity, then $F_h = \frac{l}{l+1} F_e \approx F_e$ for large l .

This completes the proof.

Remark 2: By setting l as unity in (41), the earlier state convergence based nonlinear bilateral controller [32] becomes a special case of the proposed multilateral controller.

6 Simulation Results

In order to validate the proposed scheme, a single-master/tri-slave nonlinear teleoperation system is setup in MATLAB/Simulink environment. The master and slave units forming this teleoperation system are all two link manipulators with two degrees-of-freedom motion. Their dynamical system representation is given by (1), (2) with the following matrix/vector entries:

$$M(q) = \begin{bmatrix} 3ml^2 + 2ml^2 \cos(q_2) & ml^2 + ml^2 \cos(q_2) \\ ml^2 + ml^2 \cos(q_2) & ml^2 \end{bmatrix} \quad (42)$$

$$C(q, \dot{q}) = \begin{bmatrix} -\dot{q}_2 ml^2 \sin(q_2) & -(\dot{q}_1 + \dot{q}_2) ml^2 \sin(q_2) \\ \dot{q}_1 ml^2 \sin(q_2) & 0 \end{bmatrix} \quad (43)$$

$$g(q) = \begin{bmatrix} a_g ml \sin(q_1 + q_2) + 2a_g ml \sin(q_1) \\ a_g ml \sin(q_1 + q_2) \end{bmatrix} \quad (44)$$

Where, $m_1 = m_2 = m$ denotes the mass of the links, $l_1 = l_2 = l$ denotes the link lengths and $a_g = 9.8ms^{-2}$ is the acceleration due to gravity. The numerical values of these parameters are chosen as $\{m_m = 2.0, l_m = 1.0\}$, $\{m_{s1} = 10.0, l_{s1} = 1.5\}$,

$\{m_{s2} = 5.0, l_{s2} = 2.0\}$ and $\{m_{s3} = 8.0, l_{s3} = 2.5\}$. The other parameters of the teleoperation system are the stiffness and damping of the remote environments which are selected as $\{K_e^1 = \text{diag}(100,100), B_e^1 = \text{diag}(20,20)\}$, $\{K_e^2 = \text{diag}(50,50), B_e^2 = \text{diag}(10,10)\}$ and $\{K_e^3 = \text{diag}(20,20), B_e^3 = \text{diag}(5,5)\}$.

Let us now consider that the time varying delays exist in the communication channel between the master and slave units as shown in Figure 2. We assume that the upper bound on these time delays is 0.8. We further assume that all the gamma constants are unity. The inequalities in (17), (18) are then solved which leads to the selection of decisive matrices as $K = \text{diag}(20,10)$ and $K_1 = \text{diag}(60,30)$. With the knowledge of these matrices, control gains of the teleoperation system are found to be:

$$\begin{aligned}
 K_{m1} &= \begin{bmatrix} -60 & 0 \\ 0 & -30 \end{bmatrix}, K_{m2} = \begin{bmatrix} -80 & 0 \\ 0 & -40 \end{bmatrix} - \sum_{i=1}^3 r(t-T_{si}(t)) \times \begin{bmatrix} 20 & 0 \\ 0 & 10 \end{bmatrix} \\
 K_{s1}^i &= \begin{bmatrix} -20 & 0 \\ 0 & -10 \end{bmatrix}, K_{s2}^i = \begin{bmatrix} -40 & 0 \\ 0 & -20 \end{bmatrix} - r(t-T_{mi}(t)) \times \begin{bmatrix} 20 & 0 \\ 0 & 10 \end{bmatrix}, \forall i = 1, 2, 3 \\
 R_{m1}^i &= \begin{bmatrix} 20 & 0 \\ 0 & 10 \end{bmatrix}, R_{m2}^i = 2r(t-T_{si}(t)) \times \begin{bmatrix} 20 & 0 \\ 0 & 10 \end{bmatrix}, \forall i = 1, 2, 3 \\
 R_{s1}^i &= \begin{bmatrix} 20 & 0 \\ 0 & 10 \end{bmatrix}, R_{s2}^i = 2r(t-T_{mi}(t)) \times \begin{bmatrix} 20 & 0 \\ 0 & 10 \end{bmatrix}, \forall i = 1, 2, 3
 \end{aligned} \tag{45}$$

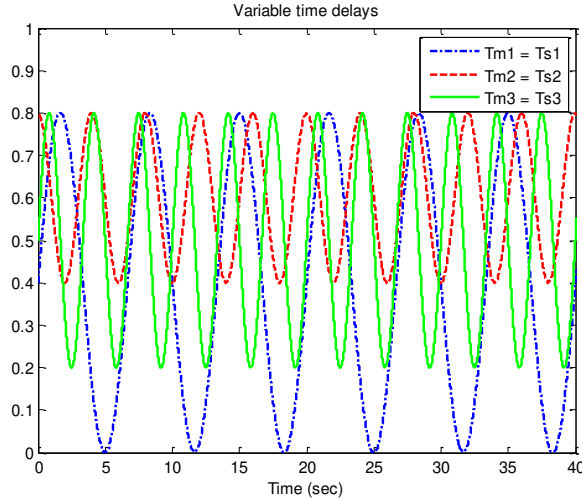


Figure 2
Time delays between master and slave units

We first analyze the time delayed teleoperation system in free motion when the operator applies a constant force as shown in Figure 3. The resultant joint positions of the master and slave units are shown in Figures 4-6. It can be seen that the slave units are following the master unit in the presence of time varying delays and the joint positions converge when the applied force becomes zero. This shows that the free motion behavior of the teleoperation system is stable.

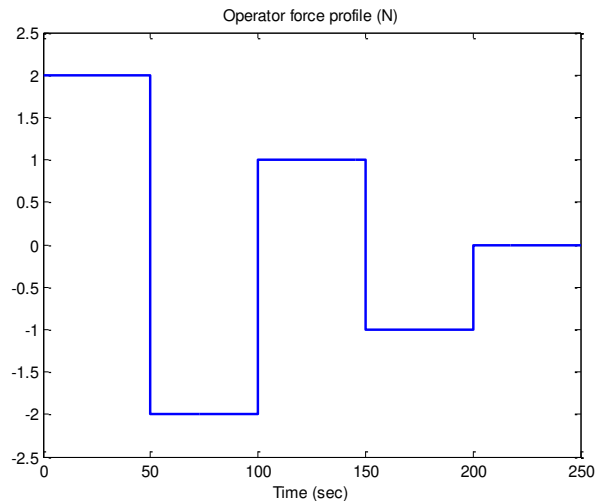
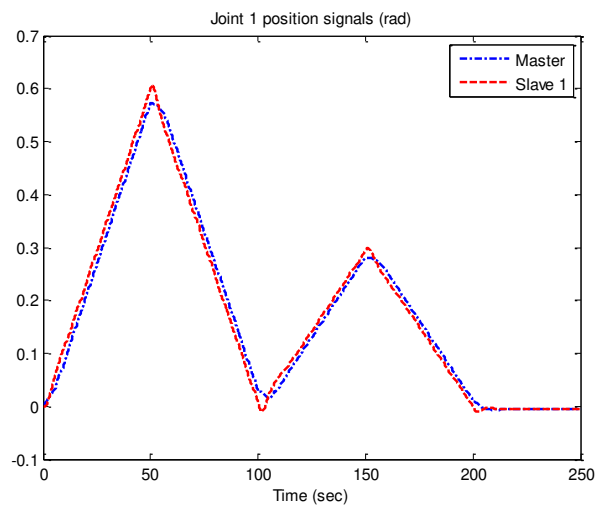
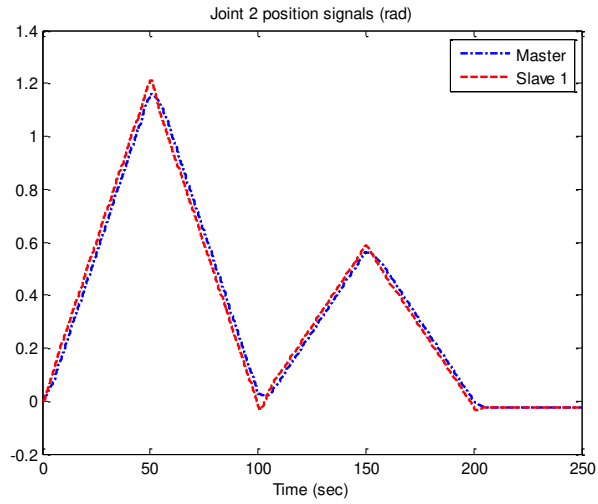


Figure 3
Operator's force profile



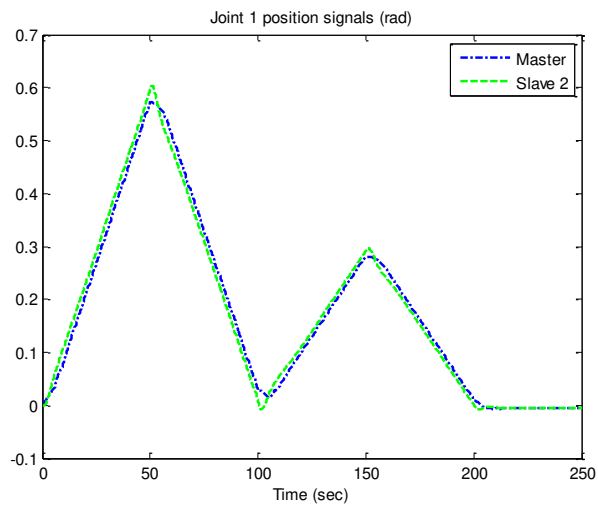
(a)



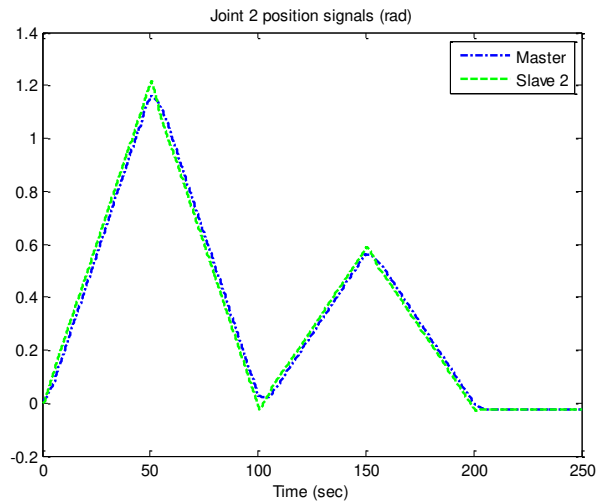
(b)

Figure 4

Free motion of Slave # 1 (a) Joint 1 position (b) Joint 2 position



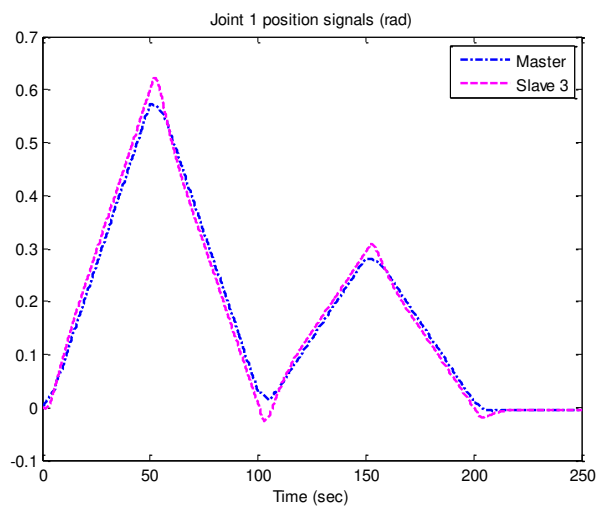
(a)



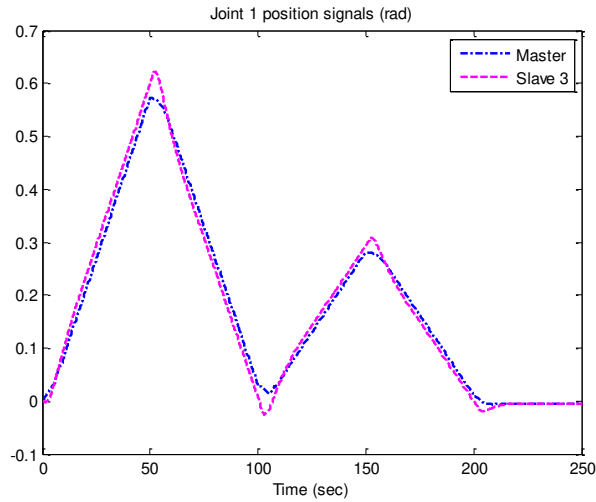
(b)

Figure 5

Free motion of Slave # 2 (a) Joint 1 position (b) Joint 2 position



(a)

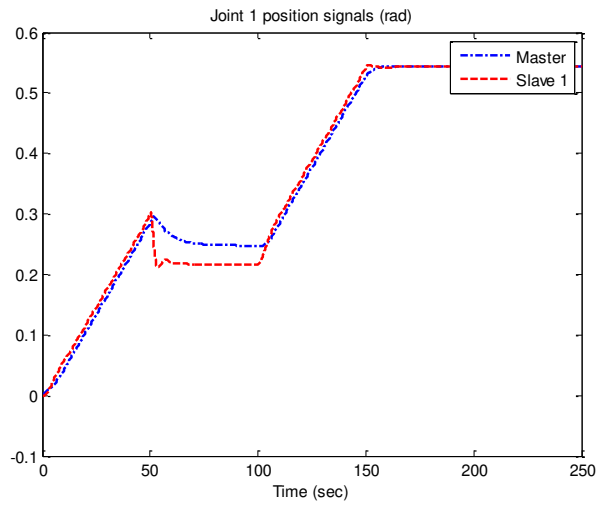


(b)

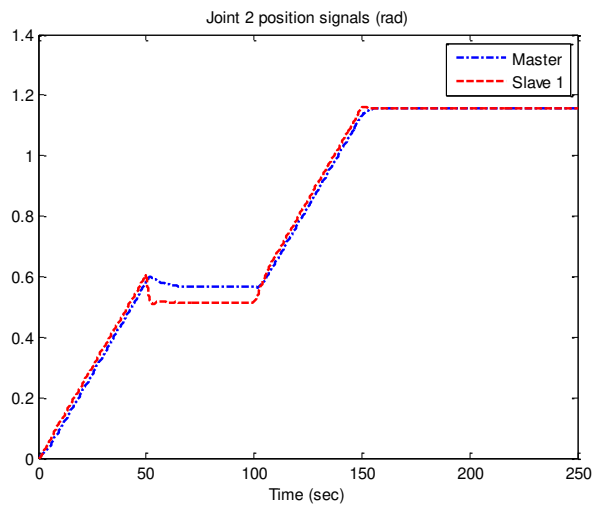
Figure 6

Free motion of Slave # 3 (a) Joint 1 position (b) Joint 2 position

We now consider the contact motion of the slaves when the operator exerts a constant force of 1 N that lasts for 150 s. The contact motion of all the slave units starts at $t = 50s$ and ends at $t = 100s$. The resultant position trajectories of the master and slave units are depicted in Figures 7-9. It is evident that the joint positions of the slave units are tracking the corresponding joint positions of the master unit and all the position signals remain bounded which implies that the teleoperation system is stable during both the free and contact motion cases. The force reflection ability of the time delayed teleoperation system is also analyzed. Theoretical result of (41) indicates that the static force reflection should be $0.25 \times (F_e^1 + F_e^2 + F_e^3)$, which is confirmed through simulation results on force reflection as shown in Figure 10.



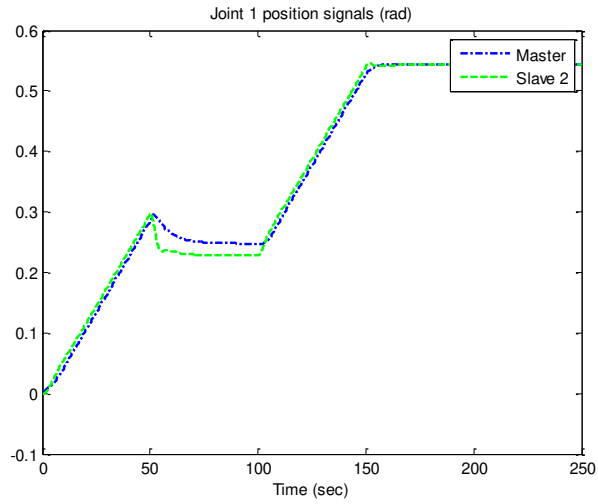
(a)



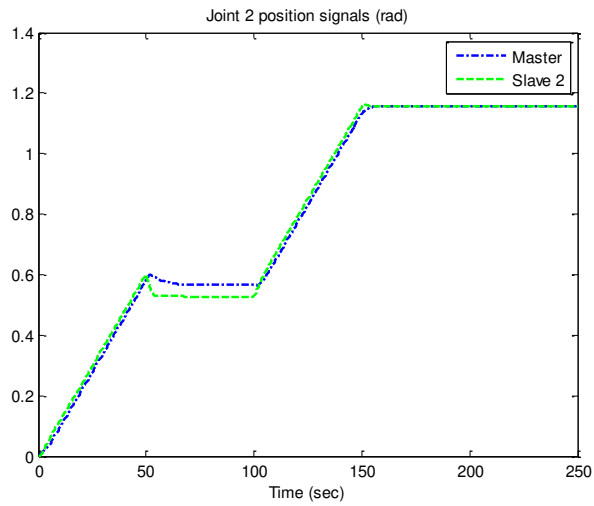
(b)

Figure 7

Free plus contact motion of Slave # 1 (a) Joint 1 position (b) Joint 2 position



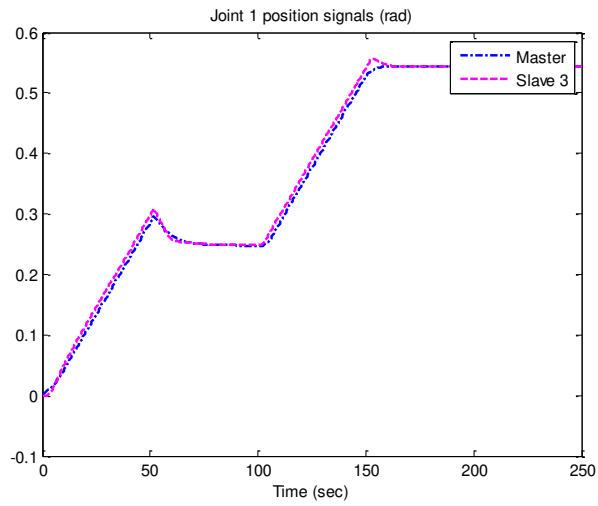
(a)



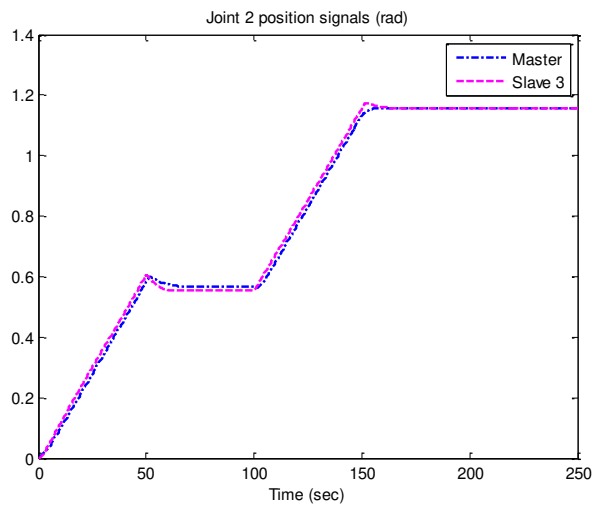
(b)

Figure 8

Free plus contact motion of Slave # 2 (a) Joint 1 position (b) Joint 2 position



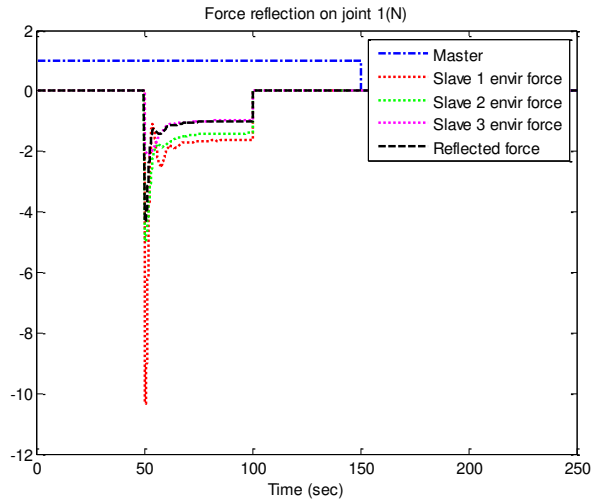
(a)



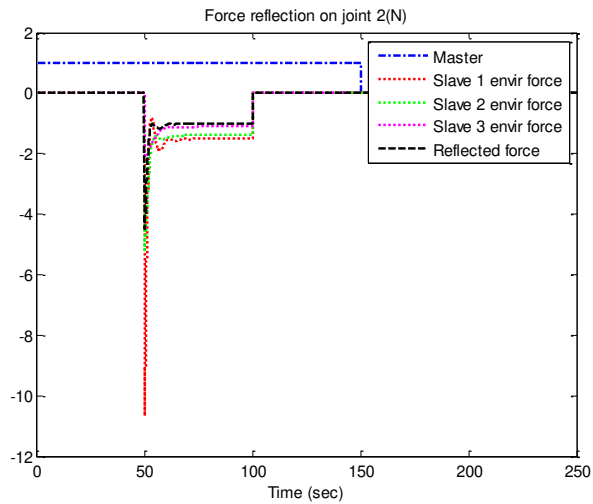
(b)

Figure 9

Free plus contact motion of Slave # 3 (a) Joint 1 position (b) Joint 2 position



(a)



(b)

Figure 10
Free plus contact motion of Slave # 3 (a) Joint 1 position (b) Joint 2 position

Conclusions

The design of a SM/MS nonlinear teleoperation system is presented in this paper. The proposed design builds upon our earlier work on the multilateral linear teleoperation systems, but considers the nonlinear dynamics of the master and slave units as well as the asymmetric time delays of the communication channel. With the help of Lyapunov-Krasovskii control theory, it is shown that the origin of the teleoperation system is asymptotically stable and the slave units track the position commands of the master unit. It is also shown that the proposed teleoperation system possess force reflection ability. To validate the theoretical findings, MATLAB simulations are performed on a single-master/tri-slave nonlinear teleoperation system where each master/slave unit has two DoF motion. It is found that all the control objectives including stability, position synchronization and force reflection are achieved. Future work involves the real time implementation of the proposed scheme.

Acknowledgement

The work of Dr. Umar Farooq was supported by Nova Scotia Graduate Scholar Program.

References

- [1] G. Niemeyer, C. Preusche, G. Hirzinger: Springer Handbook of Robotics, Springer-Verlag, New York, 2008
- [2] M. Ferre, M. Buss, R. Aracil, C. Melchiorri, C. Balaguer: Advances in Telerobotics, Springer, 2007
- [3] M. Shahbazi, S. Atashzar, R. Patel: A Systematic Review of Multilateral Teleoperation Systems, IEEE Transactions on Haptics, Early Access Article, 2018
- [4] R. Anderson, M. W. Spong: Bilateral Control of Teleoperators with Time Delay, IEEE Transactions on Automatic Control, Vol. 34, No. 5, 1989, pp. 494-501
- [5] G. Niemeyer, J. J. E. Slotine: Stable adaptive teleoperation, IEEE Journal of Oceanic Engineering, Vol. 16, No. 1, 1991, pp. 152-162
- [6] J. Ryu, D. Kwon, B. Hannaford: Stable Teleoperation with Time-Domain Passivity Control, IEEE Transactions on Robotics and Automation, Vol. 20, No. 2, 2004, pp. 365-373
- [7] Y. Ye, Y. J. Pan, Y. Gupta, J. Ware: A Power-Based Time Domain Passivity Control for Haptic Interfaces, IEEE Transactions on Control Systems Technology, Vol. 19, No. 14, 2011, 874-883
- [8] T. Kanno, Y. Yokokohji: Multilateral Teleoperation Control over Time-Delayed Computer Networks using Wave Variables, Proc. IEEE Haptics Symposium, 2012, pp. 125-131

-
- [9] H. Van Quang, J. H. Ryu: Stable Multilateral Teleoperation with Time Domain Passivity Approach, Proc. IEEE/RSJ International Conference on Intelligent Robots and Systems, 2013, pp. 5890-5895
- [10] Z. Chen, Y. J. Pan, J. Gu, S. Forbrigger: A Novel Multilateral Teleoperation Scheme with Power Based Time Domain Passivity Control, Transactions of the Institute of Measurement and Control, 2016
- [11] Z. Cheng, F. Huang, W. Song, S. Zhu: A Novel Wave Variable Based Time-Delay Compensated Four Channel Control Design for Multilateral Teleoperation Systems, IEEE Access, Vol. 6, 2018, pp. 25506-25516
- [12] D. Sun, F. Naghdy, H. Du: Application Of Wave-Variable Control to Bilateral Teleoperation Systems: A Survey, Annual Reviews in Control, Vol. 38, No. 1, 2014, pp. 12-31
- [13] J. Yan, S. E. Salcudean: Teleoperation Controller Design using H_{∞} Optimization with Application to Motion Scaling, IEEE Transactions on Control Systems Technology, Vol. 4, No. 3, 1996, pp. 244-258
- [14] Y. Kawai, T. Miyoshi, M. Fujita: Written communication system based on multilateral teleoperation using robust control, Proc. IEEE International Conference on Advanced Intelligent Mechatronics, 2017, pp. 657-662
- [15] W. Shen, J. Gu, Z. Feng: A Stable Tele-Robotic Neurosurgical System Based on SMC, Proc. IEEE International Conference on Robotics and Biomimetics, 2007, pp. 150-155
- [16] R. Kikuuwe, K. Kanaoka, M. Yamamoto: Phase-Lead Stabilization of Force-Projecting Master-Slave Systems With a New Sliding Mode Filter, IEEE Transactions on Control Systems Technology, Vol. 23, No. 6, 2015, pp. 2182-2194
- [17] M. Shahbazi, S. F. Atashzar, H. A. Talebi, F. Towhidkhan, M. J. Yazdanpanah: A Sliding Mode Controller for Dual User Teleoperation with Unknown Constant Time Delays, Robotica, Vol. 31, 2013, pp. 589-598
- [18] H. S-Reyes, L. G. G-Valdovinos, H. J-Hernandez, T. S-Jimenez, L. A. G-Zarco: Higher Order Sliding Mode Based Impedance Control for Dual-User Bilateral Teleoperation Under Unknown Constant Time Delay, Proc. IEEE/RSJ International Conference on Intelligent Robots and Systems, 2015, pp. 5209-5215
- [19] N. Chopra, M. W. Spong, R. Lozano: Adaptive Coordination Control of Bilateral Teleoperators with Time Delay, Proc. IEEE Conference On Decision and Control, 2004, pp. 4540-4547
- [20] Z. Chen, Y. J. Pan, J. Gu: A Novel Adaptive Robust Control Architecture for Bilateral Teleoperation Systems Under Time Varying Delays, International Journal of Robust and Nonlinear Control, Vol. 25, No. 17, 2015, pp. 3349-3366
-

-
- [21] Z. Chen, Y. J. Pan, J. Gu: Integrated Adaptive Robust Control for Multilateral Teleoperation Systems Under Arbitrary Time Delays, *International Journal of Robust and Nonlinear Control*, Vol. 26, No. 12, 2016, pp. 2708-2728
- [22] Y. Yang, C. Hua, X. Guan: Adaptive Fuzzy Finite-Time Coordination Control for Networked Nonlinear Bilateral Teleoperation System, *IEEE Transactions on Fuzzy Systems*, Vol. 22, No. 3, 2014, pp. 631-641
- [23] Z. Li, Y. Xia and F. Sun: Adaptive Fuzzy Control for Multilateral Cooperative Teleoperation of Multiple Robotic Manipulators Under Random Network-Induced Delays, *IEEE Transactions on Fuzzy Systems*, Vol. 22, No. 2, 2014, pp. 437-450
- [24] X. Yang, C. C. Hua, J. Yan, X. P. Guan: A New Master-Slave Torque Design for Teleoperation System by T-S Fuzzy Approach, *IEEE Transactions on Control Systems Technology*, Vol. 23, No. 4, 2015, pp. 1611-1619
- [25] U. Farooq, J. Gu, M. El-Hawary, M. U. Asad, G. Abbas: Fuzzy Model Based Bilateral Control Design of Nonlinear Tele-Operation System Using Method Of State Convergence, *IEEE Access*, Vol. 4, 2016, pp. 4119-4135
- [26] U. Farooq, J. Gu, M. El-Hawary, V. E. Balas, M. U. Asad, G. Abbas: Fuzzy Model Based Design of a Transparent Controller For A Time Delayed Bilateral Teleoperation System Through State Convergence, *Acta Polytechnica Hungarica*, Vol. 14, No. 8, 2017, pp. 7-26
- [27] C-C. Hua, Y. Yang, X. Guan: Neural Network Based Adaptive Position Tracking Control for Bilateral Teleoperation Under Constant Time Delay, *Neurocomputing*, Vol. 113, No. 3, 2013, pp. 204-212
- [28] Z. Li, Y. Xia, D. Wang, D-H. Zhai, C-Y. Su, X. Zhao: Neural Network Based Control of Networked Trilateral Teleoperation with Geometrically Unknown Constraints, *IEEE Transactions on Cybernetics*, Vol. 46, No. 5, 2016, pp. 1051-1064
- [29] C. Yang, X. Wang, Z. Li, Y. Li, C-Yi Su: Teleoperation Control Based on Combination of Wave Variable and Neural Networks, *IEEE Transactions on Systems, Man and Cybernetics*, Vol. 47, No. 8, 2017, pp. 2125-2136
- [30] D. Sun, F. Naghdy, H. Du: Neural Network Based Passivity Control of Teleoperation System Under Time Varying Delays, *IEEE Transactions on Cybernetics*, Vol. 47, No. 7, 2017, pp. 1666-1680
- [31] J. M. Azorin, O. Reinoso, R. Aracil, M. Ferre: Generalized Control Method by State Convergence of Teleoperation Systems with Time Delay, *Automatica*, Vol. 40, No. 9, 2004, pp. 1575-1582

- [32] J. C. Tafur, C. Garcia, R. Aracil, R. Saltaren: Stability Analysis of Teleoperation System by State Convergence with Variable Time Delay, Proc. American Control Conference, 2013, pp. 5696-5701
- [33] U. Farooq, J. Gu, M. E. El-Hawary, M. U. Asad, J. Luo: An Extended State Convergence Architecture for Multilateral Teleoperation Systems, IEEE Access, Vol. 5, 2017, pp. 2063-2079

ANN-based Classification of Urban Road Environments from Traffic Sign and Crossroad Data

Zoltán Fazekas¹, Gábor Balázs^{1,2}, Péter Gáspár¹

¹Institute for Computer Science and Control (MTA SZTAKI), Kende u. 13-17, H-1111 Budapest, Hungary
e-mails: zoltan.fazekas@sztaki.mta.hu, peter.gaspar@sztaki.mta.hu

²Zukunft Mobility GmbH, Ruppertswies 14, D-85092 Kösching, Germany
e-mail: gabor.balazs@zf.com

Abstract: A method that distinguishes between urban road environment types, based on traffic sign (TS) and crossroad (CR) data is presented in this paper. The types and the along-the-route locations of the TSs and the CRs – encountered during car trips – are recorded either by a human data entry assistant, or by an advanced driver assistance subsystem that has been enhanced for the purpose. A feed-forward artificial neural network (ANN) – trained in a supervised manner – carries out the classification tasks. ANNs with different topologies and training regimes are considered and tested for the purpose. These ANNs are characterized by different degrees of modularity ranging from fully modular to non-modular networks. The fully modular ANN consists of three functional modules. Two of these three were trained initially as standalone ANNs, to infer the actual road environment type solely from the TS and the CR data, respectively. The outputs of these two modules are combined via the third module. Further synapses supplement the module-level connections in the less modular ANNs. During the training of the full ANN, the TS and the CR modules are kept relatively intact, while the weights and the biases within the merger module can evolve. Test results for the considered ANNs are provided and compared.

Keywords: detection of the road environment; artificial neural networks; traffic sign recognition systems

1 Introduction

The amount of real-time data measured, gathered and processed on-board high-end road vehicles, increases continually from year to year. Such data gathering and data processing are carried out, for instance, by the real-time traffic sign recognition system (TSR) presented in [1]. The system described therein relies on an RGB image sequence, and a depth image sequence and GPS location, odometer

and map data. The depth image sequence is used for the selection of the regions-of-interest, while template matching is applied to the color-segmented RGB image sequence. The latter locates the traffic signs within the image sequence. The system repeatedly performs a particle filter based localization of the ego-vehicle based on these data. The real-time implementation is achieved via the use of multicore processors and a number of graphics processing units. As a more recent example of on-board real-time data gathering and high-volume data processing, the object classification system proposed in [2] could be mentioned. The system fuses image data and the LIDAR point cloud data. Furthermore, it uses a deep convolutional neural network that is fed with pixel-level depth data – obtained via point cloud up-sampling – and RGB color image data. In our view, the increasing trend of the real-time data gathered and processed on-board is partly explainable with – but also permits and calls for – a richer set of on-board environment perception capabilities, much richer than it was customary, say, a decade ago.

The need for a more detailed perception of the environment arises in a wide range of road environments and traffic situations. Notably, it arises in urban environments. For instance, the narrow and possibly blocked streets – typical in historical town-centers – require due attention, clear perception and fast decision making [3]. Furthermore, a profound understanding/model of the road environment is indispensable in the ever-growing and increasingly hectic road traffic at major urban junctions for human drivers and for smart/semi-autonomous/autonomous road vehicles alike. In case of smart road vehicles, the information concerning the key components of the road and vehicle environment is made available to the drivers in the form of advanced driver assistance functions; while in case of semi-autonomous/autonomous road vehicles, it is utilized by the vehicles' own control system.

The data describing the surroundings of high-end vehicles – data typically collected and processed on-board – includes the positions, dimensions and velocities of the pedestrians and the vehicles using the road, as well as the positions, shapes and dimensions of the markings, traffic signs, traffic lights and other objects on, along and in the vicinity of the roads. For instance, a computer vision solution that classifies the road environment into urban, rural, or off-road categories based on color and texture features was proposed in [4]. These features were derived from the color and texture distributions extracted from various regions of interest. The features were then combined using a trained classifier approach to resolve two road-type classification problems. The first was the determination of the off-road/on-road situations. The second was the multiclass road environment problem of determining the actual road type, namely off-road, urban, major/trunk road and multilane motorway/carriageway.

From various data mentioned above, computer vision and artificial intelligence units on-board 'guess' (calculate, estimate, determine) current traffic conditions, the actual road and lane geometry, as well as, the traffic control status (e.g., actual speed limit, green light).

Herein, the urban road environment surrounding an ego-car is classified into one of the three predefined road environment categories based on traffic sign (TS) and crossroad (CR) data. These three road environment categories are as follows: downtown (Dt), industrial/commercial (Ind), and residential (Res) areas. The classification is carried out by an artificial neural network (ANN). Though, the collected data is processed in a post-collection manner, the trained ANNs could well be installed on-board of smart cars and operated in a real-time manner within the advanced driver assistance systems (ADAS).

ANNs of different degrees of modularity are proposed and tested for the purpose of urban road environment classification in the present paper. The recognition performances of these ANNs are evaluated and compared.

Modularity is a desired characteristic of systems of any practical purpose. Such systems include computerized systems and computing/ processing machinery, as well. ANNs are no exceptions to this rule. Modularity serves many of the well-founded engineering demands during the life-cycle of systems (i.e., during system design, maintenance, validation and system renewal/update). It also promotes the traceability of the systems, and makes easier for the system developers and users to understand, follow and verify the computations carried out within.

Modularity also promotes the reusability of programs, applications, methods, modules, subsystems, and even – as it is the case here – the tested and properly functioning weights and biases associated with interneuron synapses within the ANNs and the neurons themselves. These are to be used within modules, or subsystems in a more complex computing network. But modularity comes at a price, e.g., the processing could require several stages/layers and the precision could be slightly impacted. Herein, we look at ANNs exhibiting different levels of modularity, evaluate and compare their performance when these are employed in the given application context.

To cope with the different levels of modularity, furthermore to maintain consistency among the various ANNs used in our experiments, several training regimes – using the analogy of the vibrations of excited nodes within a coupled mechanical network – were devised. These training regimes when applied to an ANN, retain – as much as possible – the weights and biases in certain well-defined parts (e.g., within a module, or within a subsystem) of the network. It is expected that the reasonably good starting values for these weights and biases – if retained or modified with care – shorten the necessary training effort in respect of the network.

The aforementioned three road environment types are rather different from a traffic safety point of view, as it was pointed out by the authors of [5]. Because of these differences, the human drivers, the semi- and the fully- autonomous road vehicles need to look out for very different hazardous traffic situations within these environments. Some important accident and crash data, for different urban road environments in the City of San Antonio, Texas, USA, are given in the cited

paper. The data presented in the article, as well as, similar accident and crash data from other cities, e.g., accident and crash data from Xi'an (China) that was presented and analyzed in [6], underline the need for the ADAS function proposed herein.

In a smart car driven by a human driver, the output of the road environment detection ADAS subsystem could be simply displayed to the driver as a short message (e.g., “You are now probably driving in a downtown area.”), so that the driver can adjust to expected traffic conditions and any foreseeable safety hazards. In case of semi-autonomous/autonomous vehicles, for instance, different voluntary speed-limits could be set for different road environments and these limits could be deferred to by the vehicular control system.

The identified road environment type, on the other hand, could be utilized within the mentioned cooperative system architecture. For instance, the customary/standard size of the TSs may vary in different road environments (e.g., nowadays some very small TSs are also deployed in downtown of Budapest), and the TS size for the given environment category, could be beneficial for cross-checking the detected TS candidates. As another example for supporting, with actual road environment information, the processing carried out within a TSR ADAS subsystem, one could mention the different occurrence probabilities of the various TSs. If several candidate TS types are identified for a particular TS encountered along the route, then for choosing the most likely type, these probabilities could be taken into consideration.

In an advantageous implementation, the TS and CR data is gathered, processed, combined and used by ADAS subsystems (e.g., by a camera-based TSR system, or by a lane keeping assistant). These subsystems may rely on data provided by on-board measurement devices (e.g., LiDARs).

The rest of the paper is structured as follows. The first subject in Section 2 is the need for a more robust cooperation amongst the ADAS subsystems in the context of road environment detection. Then achievements in the field of road environment perception, modeling, interpretation and representation are touched upon and some precursors to the proposed ADAS subsystem are mentioned therein. Finally, still in Section 2, the socio-economic relevance of the road and traffic related data is pointed out with reference to an interesting large-scale application of such data. In Section 3, a summary of the car-based TS and CR data collection work carried out in the present research is presented. In Section 4, the ANN-based urban road environment classifiers – exhibiting different degrees of modularity, used in our investigations, are described and their constituent modules/subnetworks, as well as, the two-phase training regime used in conjunction with the classifiers are also described therein. In Section 5, the test results are presented, graphically and also in tabular form, for a particular test route. The results are also discussed therein. Conclusions are drawn and the future work is outlined at the end of this paper.

2 Related Work

2.1 Making Good Use of the Subsystem-Level Cooperation within Advanced Driver Assistance Systems

A biologically inspired system architecture that supports environmental perception capabilities and makes extensive use of subsystem-level cooperation is presented in [7]. The authors of the paper argue that within the majority of the ADAS – at least then, i.e., in 2011 – the ADAS subsystems have clearly defined functionalities, and work fairly independently to complete their specific tasks. In other words, the subsystems do not make extensive use of the results produced by other ADAS subsystems. Although, the loosely integrated ADAS – built from ‘individualistic’ subsystems – show good performance at the implemented functions, such systems are stuck at a low level of abstraction and are unable to handle complex scenes and traffic situations in a generic way. Furthermore, the extension and the modification of the implemented ADAS capability is far from being straightforward. The biologically inspired system architecture proposed in their paper – which incorporates a module responsible for static domain-specific tasks, pathways for object recognition and location/distance computations, as well as a module for environmental interaction – warrants a higher abstraction level and eliminates the mentioned drawbacks, they claim.

2.2 Road Environment Detection

Road environment detection, perception, modeling, interpretation and representation have been targets of research for some time [7] [8]. The former paper has been discussed above. The authors of the latter paper list a number of current reliable working ADAS subsystems in production cars. In their view, the research activities worldwide have just started to address them, and are definitely addressing now, the driver assistance related environment perception problems for inner-city scenarios.

The aim of these activities is to provide substantial and reliable information about the actual driving situation even in such complex spatial environments. As the authors point out in their paper, this task necessitates the application of multi-sensor systems and advanced sensor fusion technologies. The most frequently used data fusion methods are either object-based, or occupancy grid-based. The authors present an occupancy grid-based fusion concept that is optimized for the environment perception task facing and posed by road vehicles. Their system initially separates the static and the dynamic information coming from the on-board environment sensors (i.e., stereo cameras, radars).

An interesting view, with considerable insight in the field of data fusion, is presented in [9]. The authors of the paper forecast that future environment perception systems, including those on-board road vehicles, will rely on model-free grid-based representations and, at the same time, on model-based object tracking solutions. This is because only a combination of both will meet the requirements of the complex ADAS.

In regards of our present topic, the multi-layer representation of stationary inner-city intersections presented in [3], is of great interest. The authors of the cited paper, use the term ‘multi-layer’ both in a geometrical sense (to distinguish between ground and raised features) and in an algorithmic/computer architectural sense (to distinguish between the sensor- and data-specific layer, the tristate abstraction layer and the fusion layer). They differentiate between and/or indicate explicit and implicit free spaces, the ground texture, the curbs, the elevated occupancies, the texture+elevated occupancies, and multiple occupancies within the intersection.

The parametric free space (PFS) map – a novel generic 2D environment representation suitable for automotive applications – was introduced in [10]. The representation proposed there is more compact than those based on common grid-based models, and therefore, it is could be used for the purpose of automotive CAN transmission. The PFS map maintains explicit information about relevant free spaces, while setting aside data describing irrelevant free spaces. Using the PFS map, an arbitrarily fine spatial evaluation can be carried out in a sensor-principle independent and real-time manner. The authors consider radar and stereo camera data streams in their experiments, but claim that additional sensors could be included in the map generation. The generation process, however, is computationally more demanding than pure grid mapping.

In the papers discussed above, the term ‘road environment’ is used in a fairly narrow, more or less, geometrical sense. The term refers to the immediate/close surroundings of the ego-car, e.g., to the physical extent of an inner-city intersection/roundabout, of a multi-lane road segment, or of a road segment by a construction site. The authors of [11], on the other hand, use a slightly different term, namely ‘driving environment’, and they study the ‘critical changes’. The examples given there are related to the sudden changes of illumination (e.g., tunnel entry, tunnel exit, shadow of an overpass), but still refers to the immediate/close surroundings of the ego-car.

In our view, the above papers tackle the highly relevant and practical spatial perception, modeling, interpretation and representation problems that arise in urban spaces, but miss out on dealing with the urban environment – around the ego-car – at a somewhat larger scale. Herein, the term ‘road environment’ is used in a more socio-economic sense and refers to larger urban spaces; this interpretation is used in [12]. The authors of the cited paper also define and use a number of socio-economic measures characterizing the urban form and everyday life. These include the residential density, employment density, land use diversity,

intersection density, size of working-age population within a (short) driving distance, number of jobs accessible by car, size of working-age population within commute distance and the number of jobs accessible within a short transit commute. Looking at maps of a given city/town with the above measures noted, one gets a fairly good understanding of how everyday urban activities are conducted in that settlement. One could easily define the road environment categories Dt and Res in a quantitative manner based on these and similar socio-economic measures, and could even combine these measures to come up with practical definitions. Herein, however, a fairly simple and a somewhat subjective categorization of these urban environments are used, which considers observable traits related to the residential and intersection densities.

In respect of the third urban environment type (i.e., Ind), we refer to [13]. The author of this doctoral dissertation examines the spatial distribution of commercial activities in Montréal, Canada. In particular, the impact of spatial determinants pertaining to street permeability and street centrality on the character and spatial distribution of retail activities is investigated. Among a wide range of analyses carried out in respect of the 'urban tissue', the author investigates the spatial characteristics of commercial streets using a morphological approach. To account for the different types of intersections on a commercial street, mainly T-types and +-types, the distances between these intersections were measured for both sides of the street, and the averages are produced. Such an approach could be utilized in precisely defining the Ind road environment category, and also the road environment data collection could be improved in this manner, as presently we do not distinguish between the left-hand-side and right-hand-side environments of the road.

2.3 Statistical Inference and ANN-based Methods for Road Environment Detection

Road environment detection (RoED) – in urban areas and in the above detailed socio-economic sense – was tackled in [14], [15] and [16]. The methods presented in these papers rely on TS and/or CR input data. Two different approaches of RoED, namely statistical inference and ANN-based classification, were presented in these publications.

It should be emphasized that this algorithmic dichotomy is not unique to the problem at hand. The statistical and the ANN-based methods commonly applied in transportation research are surveyed in [17]. The authors look at the differences between and the similarities of these methods and provide insights on how to choose one from the available algorithmic palette.

A shallow ANN was utilized to identify the actual urban road environment based on TS data in [16]. The urban RoED method proposed herein builds upon the results presented there and to a lesser extent, upon those published in [14] and

[15]. In fact, the ANN – described in [16] – that was trained and used for processing TS data has been re-applied herein as a functional module/subnetwork. In the following, this functional unit will be referred to as the TS processing module/ subnetwork. It is augmented with a similar module/subnetwork that inputs CR data and with another module that merges the outputs of the former two.

Multiple functionally independent subnetworks (modules), as a part of the whole ANN, were used and experimented with, in [18], the paper also presented various training techniques for such ANNs. Some of the training techniques described were tested in regard to the RoED system proposed herein.

Modularity within ANNs serves a variety of design objectives, see [19] for a good overview; the main motivations for turning to a modular design in this case were the availability of a trained ANN (i.e., reuse of a readily available software component), the expected reduction of the required training effort, and the increased understandability/traceability of the data processing.

2.4 Road Environments, Traffic Patterns, Composition of the Traffic and their Socio-Economic Relevance

Clearly, urban RoED, even in the socio-economic sense of the term as used here, can be carried out in a number of different ways and from different kinds of input data. The most obvious option is to use a navigational device and stored map data, such as provided by the OpenStreetMap [20], for the purpose and to rely on the urban area categorization given there. Nonetheless, a similar argument could be brought forward in conjunction with TSR, and still we see that there are many camera-based TSR systems on the market, and these ADAS subsystems either use, or do not, the ego-car’s geographical position and stored map data for data fusion and data corroboration purposes.

We opted for using the TS and CR data as inputs to the RoED ADAS function as the TS data is readily available on-board of high-end road vehicles through the TSR function, while the CR data can be generated from LIDAR point clouds. We note here that the LIDAR-based ADAS solutions have gained popularity in the recent years [21] [22], and we expect to see CR detection ADAS functions in smart cars in the near future. Having said that, it is true that instead of looking at TSs and CRs, or in deed the built environment in general, RoED could be accomplished via observing traffic intensities, traffic patterns and the composition of the traffic (e.g., with respect to road vehicle types, brands, and models).

In this context, the work presented in [23] could be mentioned here. The authors of the cited paper turn to the deep learning methodology for estimating the socio-economic characteristics of approximately two-hundred US regions. Their aim, however, is not related to driver assistance. They extract, using deep learning-

based computer vision techniques, the motor vehicles encountered by the Google cars taking Street View images. In total, several millions of Street View images were processed for the purpose of this major endeavor. From the extracted image segments/blobs, the makes, the models, and the production years of these vehicles were sought and identified. This data is then used for estimating various socio-economic characteristics of the administrative regions of the studied country.

3 Collection of Traffic Sign, Crossroad and Urban Road Environment Data

In this section, a brief description of the data collection work in respect of TSs, CRs and urban road environments is given. More details, including route maps, as well as photographs of typical scenes can be found in [14] [15] [16]. In respect of the TS data, several car-based data collection trips were made to three urban settlements within Hungary, namely to Csepel¹, Százhalombatta² and Vác³. These locations were chosen to cover urban settlement population sizes characteristic for the country. It was our aim to include industrial, cultural, commercial centers for data collection, while trying to include both historic and modern settlements, as well as settlements with garden suburbs and green residential areas. For convenience reasons, we chose destinations that are not too far from Budapest, where our research institute is located.

An Android application was used for recording the relevant TS data. The application automatically records the car-trajectory, and provides means to log the TSs. The data logged for each TS includes the TS type – the types used by the proposed RoED are shown in Figure 1 – and the TS location along the route covered, and the actual road environment type. As it was mentioned in the Introduction, three predefined road environment types were considered, namely

¹ Csepel is the 21st district of Budapest. Some decades ago, it was a working-class borough with many factories. Today, Csepel contains housing estates, as well as middle-class garden suburbs. It has approximately 85,000 inhabitants. (Excerpts from <https://en.wikipedia.org/wiki/Csepel>)

² Százhalombatta is situated about 30 km from Budapest, along national motorway No. 6. It has about 18,000 inhabitants. Looking at the modern industrial town today, it is hard to believe that the settlement has about 4,000-year-old history. (Excerpts from <http://www.1hungary.com/info/szazhalombatta/>)

³ Vác is a town in Pest county with about 35,000 inhabitants. The town is located 35 kilometres north of Budapest on the eastern bank of the Danube river. Its history dates back to the days of the Roman Empire. Later, in the middle ages, the town became a Roman Catholic bishopric. Nowadays, Vác is an educational, cultural, commercial, industrial and religious center of Pest county, as well as a popular summer resort. (Excerpts from <https://en.wikipedia.org/wiki/Vác>)

the downtown (Dt), industrial/commercial (Ind), and residential (Res) areas. A data entry assistant entered the data describing the TSs seen along the route, as well as categorized and recorded the actual road environment according to their best judgment.



Figure 1

The TS types used by the proposed RoED system; the different shades of the background signify the Dt, Ind and Res road environments, respectively, in which the given signs prevail

Eight TS types were identified in [14] as occurring frequently in urban areas and prevailing in one of the three urban environment types mentioned above. In Fig. 1, these TS types are shown in groups corresponding, from left to right, to Dt (dark grey background), Ind and Res road environment types (black and light grey backgrounds, respectively). Frequent occurrences of the two Dt TSs suggest that one is driving in a Dt area, those of the four Ind TSs indicate with some probability an Ind area, while the two Res TSs are indicative, again with some probability, of a Res area.

For the purpose of our study, the CR data was added in a post-trip manner to the trajectory data. It was done by collating the recorded trajectory with the intersection data extracted from the road layer of a public geographical information system provided by [20].

Five CR categories were considered for the purpose of RoED, namely the T-shaped CRs, the +-shaped CRs, the complex CRs, the roundabouts (all these without traffic lights), and any CRs controlled by traffic lights. Examples of these CRs are shown, from left to right, in Figure 2. These CR categories had been used in [15] for detecting change, with properly tuned Page-Hinkley change detectors, in the character of the urban road environment sweeping past the ego-car collecting data.



Figure 2

Instances of CR categories used by the proposed RoED system, namely instances of T-shaped, +-shaped, complex CRs, roundabouts – all these without traffic lights – and CRs controlled by traffic lights

It should be noted that in an advantageous implementation of the proposed RoED system, the TS and CR data should be gathered, combined and processed by ADAS subsystems. The processing could easily be carried out in real-time, both for the change detection approach and for ANN-based RoED.

4 Urban Road Environment Classifiers Making Use of ANNs Exhibiting Different Degrees of Modularity

4.1 TS-Processing Module/Subnetwork

An ANN with one hidden layer was chosen for detecting the type of the actual urban environment solely from TS data in [14]. The ANN presented therein, as well as, its variants described and used herein, were implemented, trained and tested in the simulation environment Simbrain. The features and capabilities of this simulation environment are presented in [24]. Simbrain has a number of pre-defined ANN-models, including the backpropagation model used in the present work, and instances of these models can be easily created, trained and tested.

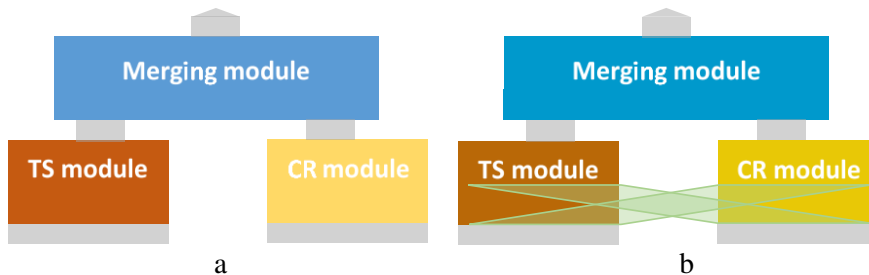


Figure 3

The modular RoED system with separate TS and CR modules (a) and the non-modular RoED system with TS and CR subnetworks connected via a merging module and via additional synapses (b)

The trained ANN is reused herein, as a module and a subnetwork, see Figures 3a and 3b, respectively, for the purpose of TS-processing within the full ANNs that consider both TS and CR data.

The input features of the TS-processing ANN and also of the TS-processing module/subnetwork, within the full ANNs, are the average distances between consecutive relevant TSs (of any sort) over the last 250, 500, 1000 and 2000 meters of the trajectory, and the number of occurrences of the typical TSs pertaining to each of the considered three road environments, again over the mentioned path-lengths. That is, there are in total 16 neurons in the input layer of the TS-processing module as can be verified in Figure 4. (The same number of neurons are used in the non-modular ANN by the TS-processing subnetwork.) Figure 4 shows the inner structure of the (trained) modular ANN for RoED and the three modules shown in Figure 3a, can be identified therein.

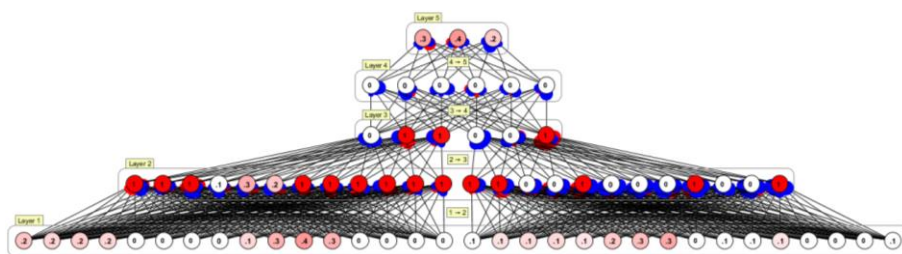


Figure 4

The inner structure of the (trained) modular ANN for RoED

The input features are calculated for consecutive route-segments of 50 meters of the car trajectory. These features are used both in the individual training and testing of the TS-processing module/subnetwork, and in training and testing of the full ANN. The network weights and biases computed in the standalone phase are retained for the full ANN and only a limited and controlled modification is allowed during the final training. After appropriate training, the TS-processing ANN – used herein as module/subnetwork, could achieve a 67.3% agreement with the ground truth data, for a particular route in Csepel.

4.2 CR-Processing Module/Subnetwork

An ANN with the same topology as the TS-processing ANN described above (cf. the left and right modules formed by the lower layers of the full ANN shown in Figure 4) that also relies on the same multiscale approach embodied in the input features was set up for detecting the type of the actual urban environment solely from CR data.

The CR data consists of the type and the location of each intersection along the route. The input features to the CR-processing ANN are the average distances between consecutive CRs (of any sort) over the mentioned path-lengths, and the number of occurrences of the typical intersection types for each of the road environment types again over the above path-lengths.

The T-shaped crossings are ‘slightly typical’ to Res areas, the +-shaped and the *-shaped (i.e., more than four-legged) crossings are ‘slightly typical’ to Dt areas, while the roundabouts and the traffic light-controlled crossings are ‘slightly typical’ to Ind areas. The above described CR-processing ANN is used herein as a module and a subnetwork – see Figures 3a and b, respectively – for the purpose of CR-processing within the full ANNs.

Similarly to the training of the TS-processing module/subnetwork, the CR module/subnetwork was trained separately, using a supervised learning approach, via backpropagation. The resultant weights and biases within the module/subnetwork were retained for the full ANN and only a limited and controlled

modification was allowed during its final training. After appropriate training, the CR-processing ANN, used here as module/subnetwork, could achieve a 59.7% agreement with the ground truth data, for a test route in Csepel.

4.3 ANNs for Processing TS and CR Data Jointly

In Figures 3a and b, the inner structures of the two main types of the full ANNs that were used in our experiments are presented schematically. A supervised learning approach was used, as both input data (i.e., TS and CR data) and the desired output (i.e. the actual urban road environment type) were known, and the desired output was available as reference data. In Figure 4, the trained version of the full modular ANN, i.e. the one sketched in Figure 3a, is shown using the Simbrain simulation environment.

Each of the two, full ANNs, comprises three functional modules/subnetworks: two of these feed into the third one. The two feeding modules/subnetworks process, exclusively/primarily, the TS and the CR data, respectively. The input signals, marked collectively with grey rectangles in Figures 3a and b, are fed into the ANNs at the bottom, while the classification results, marked with grey rectangles and arrows, concerning the current urban environment types, appear at the top of the modules/subnetworks.

The TS and CR data logs provide two separate ‘views’ on the actual urban road environment, while the third module combines the outputs of the other two modules/subnetworks, and produces a final road environment guess. In Figures 3a and 4, the TS and CR modules interact only through the merger module, while in Figure 3b, apart from the interaction via the merger module, there are also synapses between the neurons of the TS and CR processing subnetworks. These synapses are collectively marked by the two green parallelograms in the figure. The detailed graphical representation of the non-modular ANN is omitted from this paper.

Initially, each of the full ANNs in Figure 3a and b comprise only an individually trained TS-processing module and an individually trained CR-processing module. In each of the two networks, the aforementioned modules are augmented with a merger module, which is then trained through backpropagation, but without modifying the TS and CR modules. This training of the merger module is referred to as initial training of the full ANN. After this initial training, the full ANN achieved a 63.7% agreement with the ground truth road environment type data on a particular test route in Csepel. This agreement value falls between the agreement values⁴ computed for the individual modules; i.e., it falls between 67.3% for the TS-processing module/ANN and 59.7% for the CR-processing module/ANN.

⁴ These agreement values were given in Subsections 4.1 and 4.2, respectively.

5 Urban Road Environment Detection Results

Following the initial training, a two-phase training regime was carried out. The training parameters, the training errors and the test agreements for the modular full ANN, shown in Figures 3a and 4, trained according to this regime are given in Table 1.

Table 1
First-stiff-then-loose training regime used for the modular full ANN

Parameter setting	Learning rate	Momentum	Training error	Agreement	Test stripe shown below
A1	0.0250	0.0900	11.7%	56.1%	
A2	0.2500	0.9000	12.1%	56.1%	
B1	0.0050	0.0180	15.5%	49.6%	
B2	0.0500	0.1800	14.2 %	57.9 %	
C1	0.0010	0.0036	13.9%	59.3%	
C2	0.0050	0.0180	14.8%	55.8%	
D1	0.0002	0.0072	3.5%	71.9%	✓
D2	0.0010	0.0036	3.1%	68.7%	

According to this regime, a stiff training phase is followed by a loose training phase. The stiff training phase modifies all three modules (i.e., the TS-processing, the CR-processing and the merger modules) in a controlled manner, which is then followed by a loose training phase that modifies only the weights and biases within the merger.

The stiff training phases with different parameter settings are referred to as A1, B1, C1 and D1. See Table 1 for the concrete parameters. Each of these concrete training phases were then followed by loose training phases, namely A2, B2, C2 and D2, respectively. The parameter settings for these concrete training phases are also given in Table 1.

Similar data for the full ANN shown in Figure 3b trained according to the mentioned two-phase training regime are presented in Table 2. The rows of the table correspond to full non-modular ANNs with increasing percentages (i.e., 20%, 40%, 60% and 80%) of synapses, with the new synapses randomly added to the existing ones, between the neurons of the TS-, and CR-processing subnetworks.

Table 2
First-stiff-then-loose training of the ANN with additional synapses

TS-CR synapses	Stage	Learning rate	Momentum	Training error	Agreement	Test stripe shown below
20%	stiff	0.0025	0.009	07.9%	69.4%	✓
20%	loose	0.0250	0.090	06.3%	68.3%	
40%	stiff	0.0025	0.009	03.1%	73.4%	
40%	loose	0.0250	0.090	06.5%	73.0%	✓

60%	stiff	0.0025	0.009	05.6%	74.1%	
60%	loose	0.0250	0.090	14.0%	43.2%	
80%	stiff	0.0025	0.009	3.0%	76.3%	✓
80%	loose	0.0250	0.090	19.1%	63.7%	

In Figure 5, the test results for the full ANNs, corresponding to the parameter settings and training stages tagged in Tables 1 and 2, are compared to the ground truth data in respect of a particular test route in Csepel. (Note that all the lower test stripes within the pairs shown in the figure are the same, i.e., the ground truth road environment categorization of the test route. It is just repeated so that the test results are easier compared with the ground truth.)

From the top to the bottom, the test stripes for the following parameter settings and training stages appear in the figure: parameter setting D of the modular ANN after the stiff training phase, the non-modular ANN with 20% of the possible TS-CR synapses after the stiff training phase, the non-modular ANN with 40% of the possible TS-CR synapses after stiff training followed by loose training phase, the non-modular ANN with 80% of the possible TS-CR synapses after stiff training phase.

The results shown in Tables 1 and 2 indicate that the agreements with the ground truth road environment types may considerably improve, via the application of the proposed two-phase training regime, compared to the agreements achieved by the full ANNs after just the initial training. Still, one finds that even the best agreement values are around 75%, i.e., they are not that high. Clearly, better results can be gained via using GPS and reliable map data used together. This evident RoED approach was mentioned, together with some other RoED alternatives, in Subsection 2.4.

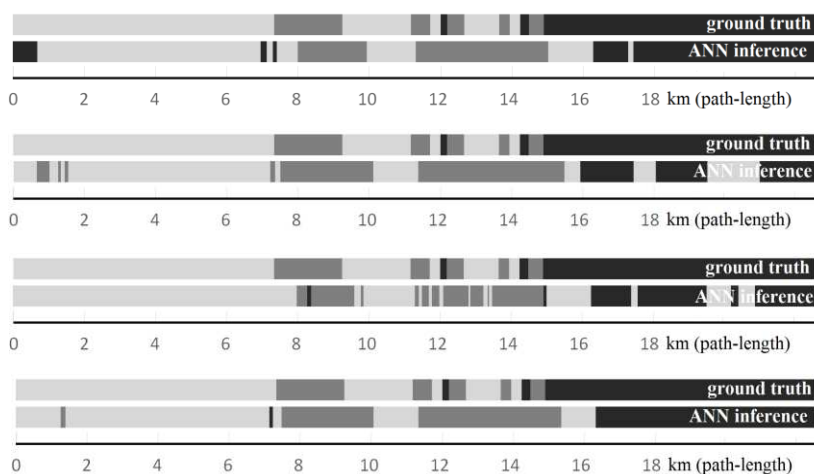


Figure 5

The road environment types manually recorded along a test route in Csepel, and those inferred by the full ANNs in different training phases (see the text for details)

In our view, when one evaluates the agreement values shown in Tables 1 and 2, one should bear in mind what sort of raw input data is available, in the presented application scenario, for the road environment classification. The raw input data series, generated by the car-based TS and CR data collection in respect of the road infrastructure, can be thought of as a realization of a marked random point process with TS and CR types appearing as marks for the along-the-route locations. The locations are characterized only by the path-length covered by the ego-car. Based on the raw input data series, the average distances between TSs and CRs, as well as the number of TS and CR occurrences over certain path-lengths, see Subsections 4.1 and 4.2 for details, are calculated by some circuitry (not shown in the figures). This derived/aggregated data are used then as input data by the TS-, and CR-processing modules/subnetworks. The marked random point process generating the raw input data is modeled with a marked Poisson point process in [15].

In our view, the TS and CR data described above is fairly basic, it lacks important road details (e.g., number of lanes and road widths are not known) that could be beneficial in determining the actual road environment type, but which are much more time-, and resource-consuming to collect, and are impractical to use on-board, at the present time. Considering the simplicity of the raw data series, the resulting agreement values seem reasonable, perhaps even surprisingly high.

Conclusions and further work

The problem of identifying the actual urban road environment type, sweeping past an ego-car, based on TS and CR data was tackled in this paper. ANNs exhibiting different degrees of modularity and having been trained according to a two-phase regime were considered and tested for the purpose. In the modular case, shown in Fig. 3a, the TS and CR modules were first trained individually and then were put together with the help of a merger module. The aim was to complement the capabilities of the processing modules and bring about improved classification results for the full network.

The considered non-modular ANN layout is shown in Fig. 3b. In this case, some additional synapses between the TS-processing and the CR-processing subnetworks are activated in a random manner and used in the processing. In both cases, improvements of the initial classification performance were achieved for certain parameter settings.

The input data used, was gathered in a small-scale data collection exercise. It is certain that a larger collection of TS and CR data from diverse regions and countries would be essential for real automotive application of the RoED system. Defining and using more urban road environment types could also extend the usability of the approach. Also, adding more CR types and corresponding input features could also be valuable for future research.

Acknowledgement

The work presented herein was supported by the Higher Education Excellence Program of the Ministry of Human Capacities in the frame of Artificial Intelligence Research Area of Budapest University of Technology and Economics (BME FIKPMI/FM).

References

- [1] K. Par, O. Tosun: Real-time Traffic Sign Recognition with Map Fusion on Multicore/Many-core Architectures. *Acta Polytechnica Hungarica*, 9, 231-250, 2012
- [2] H. Gao, B. Cheng, J. Wang, K. Li, J. Zhao, D. Li: Object Classification using CNN-Based Fusion of Vision and LIDAR in Autonomous Vehicle Environment. *IEEE Transactions on Industrial Informatics*, 14, 4224-4231, 2018
- [3] J. Rieken, R. Matthaei, M. Maurer: Toward Perception-Driven Urban Environment Modeling for Automated Road Vehicles, In: *IEEE Int. Conference on Intelligent Transportation Systems*, 731-738, 2015
- [4] I. Tang, T. P. Breckon: Automatic Road Environment Classification, *IEEE Trans. on Intelligent Transportation Systems*, 12, 476-484, 2011
- [5] E. Dumbaugh, R. Rae: Safe Urban Form: Revisiting the Relationship between Community Design and Traffic Safety, *Journal of the American Planning Association*, 75, 309-329, 2009
- [6] Y. G. Wang, S. S. Huang, W. S. Xiang, Y. L. Pei: Multipattern Road Traffic Crashes and Injuries: a Case Study of Xi'an City. *Acta Polytechnica Hungarica*, 8, pp. 171-181, 2011
- [7] R. Kastner, T. Michalke, J. Adamy, J. Fritsch, C. Goerick: Task-Based Environment Interpretation and System Architecture for Next Generation ADAS, *IEEE Intelligent Transportation System Magazine*, 3, 20-33, 2011
- [8] T. N. Nguyen, M. M. Meinecke, M. Tornow, B. Michaelis: Optimized Grid-Based Environment Perception in Advanced Driver Assistance Systems, In: *IEEE Intelligent Vehicles Symposium*, 425-430, 2009
- [9] C. Glaser, T. P. Michalke, L. Burkle, F. Niewels: Environment Perception for Inner-City Driver Assistance and Highly-Automated Driving. In *IEEE Intelligent Vehicles Symposium*, 1270-1275, 2014
- [10] M. Schreier, V. Willert, J. Adamy: From Grid Maps to Parametric Free Space Maps – a Highly Compact, Generic Environment Representation for ADAS, In: *IEEE Intelligent Vehicles Symposium*, 938-944, 2013
- [11] C. Y. Fang, S. W. Chen, C. S. Fuh: Automatic Change Detection of Driving Environments in a Vision-Based Driver Assistance System, *IEEE Trans. on Neural Networks*, 14, 646-657, 2003

- [12] K. Ramsey, J. Thomas: EPA's Smart Location Database: A National Dataset for Characterizing Location Sustainability and Urban Form, Draft Technical Report, Washington, DC, USA, Environmental Protection Agency, Office of Sustainable Communities, 1-27, 2012
- [13] J. Villain: The Impact of the Urban Form on the Spatial Distribution of Commercial Activities in Montréal, doctoral dissertation, Concordia University, Montreal, Quebec, Canada, 1-203, 2011
- [14] Z. Fazekas, G. Balázs, L. Gerencsér, P. Gáspár: Inferring the Actual Urban Road Environment from Traffic Sign Data Using a Minimum Description Length Approach, *Transportation Research Procedia*, 27, 516-523, 2017
- [15] Z. Fazekas, G. Balázs, L. Gerencsér, P. Gáspár: Detecting Change in the Urban Road Environment Along a Route Based on Traffic Sign and Crossroad Data. In: *Intelligent Transport Systems - From Research and Development to the Market Uptake*, Lecture Notes of the Institute for Computer Sciences, Social Informatics and Telecommunications Engineering (222) Springer, Cham, Switzerland, 252-262, 2018
- [16] Z. Fazekas, G. Balázs, P. Gáspár: Identifying the Urban Road Environment Type from Traffic Sign Data Using an Artificial Neural Network. In: *the Proceedings of the International Scientific Conference on Modern Safety Technologies in Transportation*, Herlány, Slovakia, 42-49, 2017
- [17] M. G. Karlaftis, E. I. Vlahogianni: Statistical Methods versus Neural Networks in Transportation Research: Differences, Similarities and Some Insights, *Transportation Research Part C: Emerging Technologies*, 19, 387-399, 2011
- [18] T. Caelli, L. Guan, W. Wen: Modularity in Neural Computing, *Proceedings of the IEEE*, 87(9), 1497-1518, 1999
- [19] R. Rojas: *Neural Networks: A Systematic Introduction*, Springer Science & Business Media, Berlin, Germany, p. 502, 2013
- [20] OpenStreetMap contributors: Road Network in Hungary, URL: <http://planet.openstreetmap.org>, 2015, last accessed: 23 Feb, 2018
- [21] A. Asvadi, C. Premebida, P. Peixoto, U. Nunes: 3D Lidar-based Static and Moving Obstacle Detection in Driving Environments: An Approach Based on Voxels and Multi-region Ground Planes. *Robotics and Autonomous Systems*, 83, 299-311, 2016
- [22] F. Jiménez, J. E. Naranjo, J. J. Anaya, F. García, A. Ponz, J. M. Armingol: Advanced Driver Assistance System for Road Environments to Improve Safety and Efficiency. *Transportation Research Procedia*, 14, 2245-2254, 2016
- [23] T. Gebru, J. Krause, Y. Wang, D. Chen, J. Deng, E. Lieberman Aiden, L. Fei-Fei: Using Deep Learning and Google Street View to Estimate the Demographic Makeup of Neighborhoods across the United States. *Proceedings of the National Academy of Sciences*, 114, 13108-13113, 2017
- [24] Z. Tosi, J. Yoshimi: Simbrain 3.0. *Neural Networks*, 83, 1-10, 2016

Assessment of the Investment in Real Estate through Innovative Funding Mechanisms

Rita Remeikiene¹, Ligita Gaspareniene², Romualdas Ginevicius³

^{1,2} Lithuanian Institute of Agrarian Economics, V. Kudirkos str. 18-2, 03105 Vilnius, Lithuania; e-mail: rita.remeikiene@laei.lt; ligita.gaspareniene@laei.lt

³ Vilnius Gediminas Technical University, Sauletekio av. 11, 10223 Vilnius, Lithuania; e-mail: romualdas.ginevicius@vgtu.lt

Abstract: The main purpose herein, is to assess the popularity of the innovative funding mechanisms when acquiring land at the national level. The results of the research suggest that the innovative funding mechanisms are neither popular nor available when making the investment in land in Lithuania. The unavailability of the innovative real estate funding mechanisms was determined by the lessons of the financial crisis of 2008. Now land plots are mainly acquired with personal funds as loans issued by credit unions. The future prospects are not very favorable, since a very small part of the transactions in the land market in Lithuania are funded with bank loans due to the Scandinavian banking policy which is currently responding to the crisis in the Scandinavian asset market. Issuance of corporate bonds serves as another source of funding. Although, theoretically, the mechanism of “crowd funding” could also be employed, the cases of its employment, thus far, have not been registered at the national level. The novelty of this article lies in the provision of a comprehensive approach to the innovative land funding mechanisms since scientific literature, thus far, has lacked the research on innovative real estate funding mechanisms.

Keywords: real estate; funding mechanisms; investment; land; innovative funding

1 Introduction

As real estate plays one of the major roles in modern economics, the mechanisms of its funding are under scrutiny at both European and global scales. The last global financial crisis has forced to look for the ways to reduce welfare expenditures. National governments have started reducing their investments and in many cases, completely stopped funding the improvement of public infrastructures (road maintenance, renovation of buildings, etc.).

Real estate markets are increasingly being treated as beneficial and able to flexibly respond to consumer needs, and so promote the revival of general economics.

Since 2008, advanced European countries have focused on the innovation of real estate funding mechanisms: they are looking for the ways to optimize the role of lending institutions, balance the debt-to-asset ratio in real estate development projects, rationally assess the risks of these projects, establish reasonable debt repayment requirements and promote public-private partnership. Currently, Europe possesses a wide variety of the innovative real estate funding mechanisms. The development of the global financial system conditions the establishment of the institutions searching for any asset accumulation opportunities. The visions of these institutions are linked to the investment in land through the innovative funding mechanisms that serve as an excellent basis for the new interest in land. As it was noted by Knuth (2015), the innovative real estate funding systems open the way for global, cross-regional and local investment in land. While analyzing the current policies of large-scale land transactions (co-called land grabbing), it seems reasonable to research land investment expediency, tendencies and funding opportunities.

Although, the current scientific literature has been rich in the studies focused on real estate funding issues (real estate funding forms were analyzed by Haffner and Boelhouwer (2006), Kemp (2007), Griggs and Kemp (2012), Squires et al. (2016) and others), the opportunities to invest in land through the innovative funding mechanisms have hardly been covered, in particular, at the national level.

The novelty of this article. Scientific literature lacks the studies on the innovative real estate value determinants and funding mechanisms. Hence, we find it purposeful to introduce a comprehensive approach to the innovative land funding mechanisms.

The main purpose of this article is to assess the popularity of the innovative funding mechanisms when acquiring land at the national level. For implementation of the defined purpose, the following **objectives** were developed:

- 1) To research the theoretical aspects of the innovative real estate funding mechanisms,
- 2) To select and introduce the methodology of the research,
- 3) To assess the popularity of the innovative funding mechanisms when making the investment in land as an investment object at the national level.

The methods of the research include systematic and comparative literature analysis and expert evaluation.

2 Bank- and Market-based Financial Systems. The Innovative Real Estate Funding Mechanisms

As large-scaled real estate projects in Europe are funded by employing the innovative funding mechanisms, it is extremely important to understand the character and features of these mechanisms so as to assess their role in real estate development.

The character of the innovative real estate funding mechanisms in different countries to a large scale depends on the national institutional environment and regulation of the financial system. Tiwari & White (2014) distinguish between bank-based financial systems and market-based financial systems. The differences between these two types of financial systems are observed since each of them differently accumulates savings from households, businesses and governments, differently selects and monitors the investment, and differently manages the risk. The role of any financial system, in this context, is to expand the mechanisms that would allow to effective funding of investments in real estate, with consideration of the class and characteristics.

So, which funding mechanisms can be provided by bank-based and market-based financial systems to redirect the accumulated savings to the development of real estate and property investment? Although traditionally, real estate is funded by employing bank lending mechanisms, over the last two decades more innovative funding mechanisms which enhance the role of market-based financial systems have also been developed.

Under the conditions of the modern economy, the investment in real estate is commonly funded through loans and subsidies (Bilal & Kratke, 2013). Marseguerra and Cortelezzi (2009), who researched the effects of debt-financing on real estate investment decisions, found that debt financing induces agents to invest earlier than in the case of pure equity financing. Nevertheless, the former type of funding is relatively inflexible: issuance of a bank loan is a long process, especially in terms of preparation of a project, documentation, submission of guarantees and deposits, the period of consideration, etc. In addition, if a state follows the policy of loan issuance limitation (such policies were established in weaker economies, including Lithuania, after the global financial crisis of 2007-2008 before which unreasonable availability of loans (often even without verification of a debtor's solvency) had caused painful problems of insolvency and numerous cases of foreclosure)), there are no guarantees that a loan will be issued at all. The process of subsidisation is even longer since a person (natural or juridical) who applies for subsidies has to prepare a project and pass the procedures of tendering (subsidisation in a state is commonly limited and granted only to most promising and relevant projects). Considering the facts explicated above, it can be stated that unlike the traditional real estate funding mechanisms, the innovative mechanisms allow to accumulate the funds, share

the investment risks among the participants of a project (an owner, an investor, a state, financial institutions, international organisations, etc.) and benefit from greater flexibility of funding (Carter, 2006).

Risk is one of the major components that must be considered when developing the innovative funding mechanisms (Bartke, 2013). Apart from the opportunities provided to large-scale projects, the innovative funding mechanisms incorporate many factors of risk that must be assessed along with the potential financial returns. As it was noted by Squires et al. (2016), the majority of the modern real estate funding mechanisms, such as value capture bond mechanisms, cover the innovation risks which are not high, but call for reasonable consideration. Although in general sense the development of real estate can be risky, the risks can be mitigated by increasing the degree of diversification of financial investments (Havard, 2013). As it was noted by Beracha et al. (2017), short-fall risk in a decumulation portfolio decreases with substantial allocations to real estate either public or private. What is more, an investor's strive to earn profits acts as an incentive for further development whatever methods of risk management are employed (Weber, 2002). The rational measures of risk management (i.e., consideration of the variety of risk and profit factors) ensure the efficient and stable long-term financial development of long-scale real estate projects.

The investment in real estate can also be funded through solidarity, public-private partnership as well as loan and bond innovation mechanisms (Grishankar, 2009) (see Fig. 1).

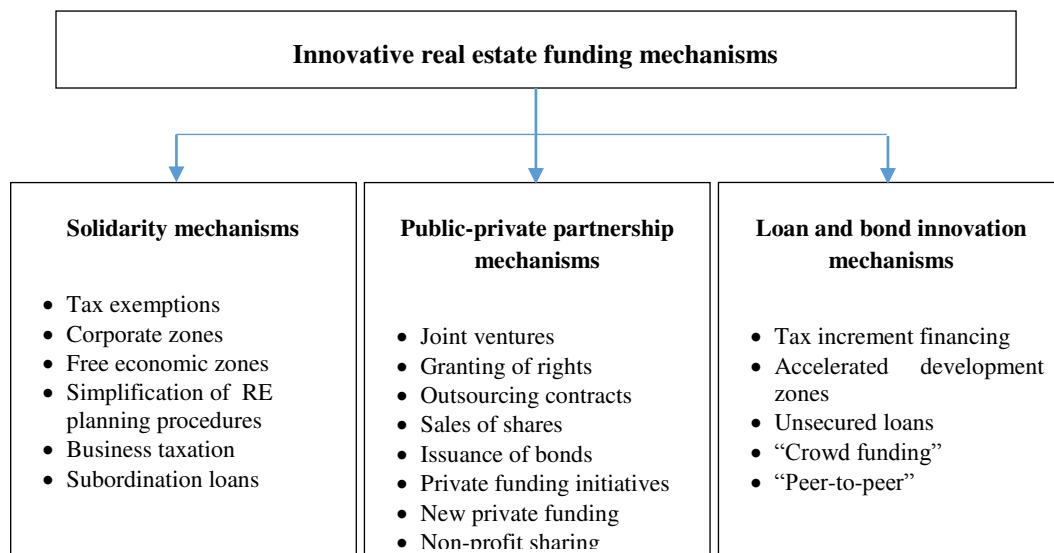


Figure 1

Innovative real estate funding mechanisms (source: compiled by the authors)

Public-private partnership is one of the most modern real estate funding mechanisms. The partnership of this kind facilitates budget constraints, contributes to the improvement of public service quality, promotes innovations and optimizes risk sharing (Liu & Wilkinson, 2014). Public-private partnership covers a wide variety of agreements starting with private financial initiatives, joint ventures and granting of rights, and ending with outsourcing contracts and sales of shares (McQuaid & Scherrer, 2008). Public-private partnerships are commonly employed in the countries where the schemes of real estate development, in the private sector, are based on long-term commitments.

Iblher & Lucius's (2003) study revealed that the demand for joint ventures in Germany is increasing, and this increase is linked to the growing demand for real estate in the largest German cities. The mechanism of joint ventures is often employed when co-operating institutions build close long-term relationship. The main advantage of this mechanism is sharing of responsibilities.

The "Lammenschans" real estate development project in Leiden (the town with the population of 120000 people) in the Western Netherlands can serve as an example of the employment of the innovative real estate funding mechanisms. The real estate was located in the southern part of Leiden, near the railway station. The town's municipality had developed the strategy following which the area of the project was divided into some complexes for different constructions, and the use of land in these complexes was restricted. The territory had to be rearranged into mixed-purpose land plots for a school, residential housing (dormitories, apartments), retail centers, parks and squares, industrial buildings and service centers. The owners of the land as well as the developers of the project took the initiative to protect their finances and agree on land ownership restrictions. The innovative approach in this project was funding of the real estate in the land plots after rearrangement (van der Krabben & Needham, 2008).

The above-described mechanism is the instrument of public planning employed for reduction of a real estate project risk and generation of the potential returns to the owners and project developers (van der Krabben & Heurkens, 2014). By employing this mixed real estate funding mechanism, land owners reformed their ownership rights, i.e. they acquired the right to change the purpose of the land (to use the land for construction in accordance with the terms of the project). Later, the ownership rights were transferred to the municipality. This financial innovation allowed the owners to participate in the publically-controlled project of regional development. After the "Lammenschans" real estate development project, this innovation has become widely-used in the projects of regional development as a measure that allows land owners and real estate developers to revive stagnant areas.

Public-private partnership initiatives are sometimes funded through issuance of bonds which are commonly linked to particular indices (the bonds of this type are called index-linked debts) and put bond holders at a potentially high inflation risk

(The European Public-Private Partnership Expertise Centre, 2010). The practice of land plot re-parceling, in some cases employed by municipalities to revive particular plots of land, can serve for public-private partnership initiatives (van der Krabben & Heurkens, 2014).

When analyzing the mechanisms of public-private partnership, the partnerships between foreign institutional asset funds and private development loan providers should also be mentioned. According to Squires *et al.* (2016), the latter mechanism is suitable for the projects of the residential property for rent, large-scaled real estate development projects and infrastructural projects, i.e. a substantial number of institutions could be attracted to invest in a particular asset class on condition they were offered a stable long-term low-risk return non-correlated to the return on the investment in other assets. Regardless of whether the real estate market is going up or down, large-scale projects possess longer stages of implementation, so it may take time to eliminate possible market inefficiencies even when the stimuli of the investment have been recognized. Financial innovations may help to “lock” the value of assets in any stage of a project (e.g., when acquiring land, starting-up or finishing construction works, etc.). From this point of view, the above-described partnerships are less dependent on the changes in real GDP, interest rates, inflation components, money supply and stock market returns – the factors that are recognized to significantly affect real estate fund returns (Delfim & Hoesli, 2016). It can be argued that distribution of funds for different stages of a project is a “long deal”, and long-term investments are more favorable for implementation of large-scale projects funded by institutional investors. Nevertheless, institutional funding is commonly employed when the relationships with banks become complicated and long-term lending restrictions are unreasonably strict.

Private funding initiatives, such as employment of the innovative funding mechanisms, are quite controversial as private funding initiatives are linked to high construction risks which are transferred to the private sector just following the argument that the private sector is capable of managing this type of risk (Adair *et al.*, 2011). Services are provided on the basis of a contract between a private consortium and an authorized public institution. As in the case of public-private partnership, private funding initiatives can offer the value for money and increase the overall efficiency of the private sector (Wall & Connolly, 2009). Private funding initiatives are arranged to ensure the full coverage of the consortium costs and generate extra returns on the borrowed capital (i.e., the returns on investment) (Greenhalgh & Squires, 2011). Despite the initial difficulties, such as insufficient operational flexibility, governments sometimes introduce the modifications of private funding initiatives so as to offset plausible shortcomings of this type of funding (HM Treasury, 2012). One of the most innovative forms of private funding initiatives is Private Funding 2 (PF2). Non-profit sharing in Scotland has already pushed out private funding initiatives. Non-profit sharing is reported to bring extra benefits: it generates limited returns so that the return surplus could be

reinvested in the public sector and thus would allow meeting the public interest (Scottish Futures Trust, 2013).

The most recent loan and bond innovations include tax increment financing (TIF) and accelerated development zones. The latter have become extremely important for the projects of infrastructure (The British Property Federation, 2008; Webber, 2010). For instance, the “BatterSea” power station project in London was initiated following the mixed scheme of the residential and commercial real estate development. The project was funded by employing the traditional debt-asset principle, i.e. the funds were obtained from foreign investors, pension funds and international banks. Although the development of the transport infrastructure was funded together with the construction of buildings, the former was based on the mechanism of public-private partnership which at that time seemed to be quite innovative. The case of the “BatterSea” showed that the project was not successful: it was stopped in 1983. Later on, the power station saw the changes in its owners and bankruptcy administrators till in 2006 the area of 750000 square meters with 600 residential and commercial premises was sold to the “Treasury Holdings”, the Irish real estate development company, for 400 million pounds. The project was terminated in 2011 when the “Treasury Holdings” took over its administration and appointed the National Asset Management Agency (NAMA) a sole proprietor of this real estate. In September 2012, the property was transferred to the “SP Setia”, the Malaysian consortium, which initiated the employment of the innovative real estate funding mechanisms. Having an unconditional right of ownership, the “SP Setia” initiated the campaign of the “BatterSea” power station development and started-up the works of the area reviving (the works were completed in stages and lasted for nearly 12 years). The value of the real estate reviving works amounted to approximately 8 billion pounds (Squires et al., 2016). The innovative funding mechanisms included tax increment financing as namely this mechanism allowed to balance the costs and returns by increasing the total value of the property and fixing the value of any improvements. This model of funding was largely directed towards the residential buildings with the aim to sell these buildings and thus raise the funds for the financing of the other stages of the project (currently the revenues from the sales and exploitation of residential buildings can be earned due to the vitality of the housing market in London). For funding of the transport infrastructure, the Public Sector Loans Board granted a loan of one billion pounds. It is expected that this expenditure will be covered after fixing of influence taxes – this way, the Northern underground line will be extended (Squires et al., 2016).

Coleman and Grimes (2009) and Medda et al. (2012) focused on a betterment tax and an accessibility increment contribution. Increment is an increase in the value of any property determined by public decisions and interventions. For instance, increment can be determined by abolition of land-use restrictions, changes in the purpose of a land plot after modification of the general land-use plan, improvement of the transportation or utility infrastructure, etc. To prevent an owner from receiving unearned increment and compensate a state or a community

a part of the spending on the infrastructure improvement, a betterment tax is imposed. Coleman and Grimes (2009) discuss two scenarios of increment: improvement of the infrastructure and changes in the purpose of a land plot (e.g., arable land is transformed into residential land). According to Coleman and Grimes (2009), both a regular land tax and an increment contribution should be levied against such increase in the value of property.

Loans can be issued by employing crowd funding systems. Crowd funding is a method to fund a project or an activity when funds are collected from a large number of people. Funds can be collected by mail, during various events, through online intermediaries, etc. The crowd funding model is based on participation of three main agents: an initiator of a project, supporters and intermediaries (intermediary platforms which link project initiators and supporters). The statistics show that in 2015 over 34 billion dollars worldwide were raised for funding of different projects by employment the method of crowd funding (Barnett, 2015). The model peer-to-peer is a method of debt financing which allows individuals to lend and borrow money without applying to financial institutions as intermediaries (Steinisch, 2012). Money is lent or borrowed online by combining the interests of lenders and debtors. The advantage of this method to lenders is that the peer-to-peer offers higher interest than the traditional lending methods (for instance, keeping money in bank accounts). For debtors, it is an opportunity to raise funds for different projects or activities that could hardly be funded through other channels. The main disadvantage of the peer-to-peer fund raising method is that a lender has practically no guarantees of a debtor's credibility. For this reason, lenders in some cases may demand higher interest for higher risks (Kennard, Bond, 2011).

Summarizing, it can be stated that the importance of public-private partnerships between foreign institutional asset funds and development loan providers is rising, in particular, when implementing real estate projects that require wide infrastructures. Employment of different methods of funding in particular stages of the real estate cycle can be considered as an extremely flexible form of real estate financing. Fixing of value as well as concentration on an increase in the total value of property are efficient supplements of the innovative real estate funding systems. Although the financing based on the expectations of the value increase in the future is comparatively risky, the appropriate measures of risk management (e.g., fixing of value in different stages of project implementation) can mitigate this risk. The variety of the real estate funding mechanisms is an important determinant of the real estate market development since it allows a reduction in the costs and risk of the investment in Real Estate, diminishes the compulsory volumes of investment per person and shorten the time of investment. In other words, the innovative real estate funding mechanisms make preconditions for effective and well-structured real estate transactions. Nevertheless, the innovative real estate funding mechanisms not always can be considered as an equivalent alternative to the traditional forms of funding. In some cases, the innovative

mechanisms serve as extra measures promoting the development of the real estate market.

In Lithuania, real estate is mainly funded by issuing bank loans or bonds. The other real estate funding methods are not popular due to the imperfections of the country's legal framework, poor subsidization of housing acquisition and a lack of VAT exemptions for real estate buyers.

3 Research Methodology

Literature research revealed that within the area of real estate financing, the methods of descriptive statistics that quantitatively describe and/or summarize the features of a collection of information are prevalent (see Table 1).

Table 1
Review of previous research methodologies applied in the area of real estate financing (source: compiled by the authors)

Research methods	Author(s), year
Literature review	Kane, 2001; Breuer, Kreuz, 2011; Vicent, 2015; Olsson, 2015
Citation analysis	Breuer, Kreuz, 2011
Statistical data review	“Knight Frank”, 2017; “JLL”, 2017
Case analysis	Squires, 2015; “JLL”, 2017
Desk research	Squires, 2015
Secondary data analysis	Ezimuo et al., 2014; Olsson, 2015
Sample surveys	Ogedengbe, Adesopo, 2003; Nkyi, 2012; Mwathi, 2013; Ezimuo et al., 2014
Interviews	Ogedengbe, Adesopo, 2003; Iblher, Lucius, 2003; Nkyi, 2012; Ezimuo et al., 2014; Squires, 2015
Trend analysis	“Knight Frank”, 2017
Univariate regression	Lasfer, 2007
Multiple regression	Gonenc, 2005; Lasfer, 2007; Abor, 2007; Nkyi, 2012
T-test	Brown et al., 1996; Redman et al., 2002; Ali et al., 2006; Nkyi, 2012
Correlation analysis	Nkyi, 2012
Chi-square	Brown et al., 1996; Nkyi, 2012
Factor analysis	Nkyi, 2012
SWOT analysis	Acquah, 2011; Nkyi, 2012
Chow test	Gonenc, 2005

As it can be seen from the review in Table 1, the most substantial part of previous studies adopt sample surveys (survey questionnaires) and interviews, whereas some of the studies rely on the sources of secondary data (financial statement,

articles from companies, press releases). The use of statistical tools demonstrates a considerable drift from basic (percentages, ratios) to more complicated tools (univariate and multiple regressions, T-test, correlation, Chi-square, factor analysis, Chow test). SWOT analysis is employed to research the impact of real estate financing related internal (strengths and weaknesses) and external (opportunities and threats) factors.

Expert evaluation is one of the most popular insight methods applied in different areas of research (Baležentis & Žalimaitė, 2011). According to Rudzkiene & Burinskiene (2007), expert evaluation can be treated as a generalized opinion of a group of experts. It is a procedure that allows combining different opinions and having an insight in the general understanding. Expert evaluation is commonly employed for the research of a particular problem, process or phenomenon that requires specific knowledge and abilities. The results of this research are submitted as reasoned conclusions and recommendations (Rudzkiene & Burinskiene, 2007). According to Makridakis et al. (1998), expert evaluation should involve 10 – 100 experts depending on the purpose of the research and expert competence in the area under consideration. Other scientists submit slightly different recommendations. For instance, Augustinaitis et al. (2009) recommends inclusion of at *least 5 experts* to ensure the accuracy and reliability of the research results. While conducting the empirical research, the authors of this work, followed the latter methodological recommendation in order to keep the focus on the expert competence, their specific knowledge of the real estate market and understanding of the conditions and problems of business environment rather than the scale of the questionnaire survey. First, 8 experts were included in the study, however, due to the split of opinions by filling out the questionnaire and improving the meaning of Cronbach's alpha and Kendall Concordance, three of them were removed from the expert evaluation.

Following the above-described recommendations, the group of the experts included 5 people:

- *Marius Dubnikovas*, who is currently in charge of Business Development Manager position at “Compensa Life Vienna Insurance Group SE”, with more than 15 years of professional and practical experience in the areas of real estate valuation and finance. He started his career as the President of Lithuanian Financial Brokers Association, and subsequently followed the position of Client Investment Manager at “Finasta Ltd.”. The expert is also the Chairman of the Tax Committee, Lithuanian Business Confederation. The financial analyst is particularly active with his speeches and insights into the trends of the real estate market in media;
- *Saulius Vagonis*, who is the Head of Valuation and Analysis Department in “OBER-HAUS Real Estate Ltd.”. He has acquired his experience in working with real estate for over 20 years. During the expert's career, more than 3000 asset evaluations and about 100 outsource market studies and analyses have been conducted. Saulius Vagonis is a board member of

Lithuanian Association of Property Valuers and Lithuanian Association of Property and Business Valuation Enterprises, and the Chairman of the Commission on Science and Education. He actively participates in real estate conferences (e.g. Real Estate Conference 2016 and 2017, organized by the Bank of Lithuania);

- *Dr. Vytautas Azbainis*, who has gained his experience in drawing up real estate investment projects and land plot detailed plans during 13 years of professional career. Since 2005 he has held the position of the director of “Vilnius Namas Ltd.”. In 2014, he defended the dissertation on the topic “Real Estate Market Cycle Management and Modeling”;
- *Romualdas Paulauskas*, who has accumulated more than 15 years of experience in the real estate sector. Currently, he is the Head of “OBERHAUS Real Estate Ltd.”, Panevėžys Department. His professional insights are published in popular Lithuanian newspapers “Verslo žinios”, “Lietuvos rytas”, “Vakarų ekspresas”, etc.;
- *Emilijus Gedvilas*, who is a broker at “Akorus Real Estate”. The expert has been purposefully working with land investment, purchase and sales of real estate, and the development of real estate objects for about 4 years.

During the empirical research, the experts were asked to assess the land investment funding mechanisms in Lithuania. With reference to the results of scientific literature analysis, the globally-practiced real estate funding mechanisms were classified into three main categories: the traditional mechanisms of private investment, the traditional mechanisms of public investment and the innovative real estate funding mechanisms. The experts were asked to indicate which mechanisms are most available when funding the investment in land in Lithuania. The main purpose of the questions was to identify the most available real estate funding mechanisms in the country and promote the need to develop the network of more innovative sources of real estate funding.

The experts were asked to evaluate particular mechanisms and statements on a scale from 1 to 5 (i.e. from 1 – “I completely disagree / It is completely irrelevant” to 5 – “I completely agree / It is completely relevant”). In accordance with the strength of their agreement / disagreement, the experts could select the intermediate numerical values 2, 3 or 4.

The data was processed with SPSS (Statistical Package for Social Sciences) and “Microsoft Excel” software.

In general, reliability of expert evaluations depends on the number of experts and the level of their knowledge. Presuming that experts are accurate assessors, it can be stated that an increase in the number of experts contributes to higher reliability of an expertise. The degree of an expert’s competence (i.e. the degree of an expert’s qualification in the area under consideration) is quantitatively measured by employing the coefficient of competence. However, it was not employed for this research.

The special attention should be drawn to possible interpretations of *Cronbach's alpha* coefficient when developing the conclusions of the expert evaluation. *Cronbach's alpha* indicates whether a questionnaire reflects an object under consideration with appropriate accuracy. Some scientists, for instance, Nunnally & Bernstein (1994), state that *Cronbach's alpha* has to be higher than 0.7, while others, for instance, Malhotra & Birks (2003), propose that the lowest marginal value of a questionnaire's reliability is 0.6. Hence, the selection of the lowest marginal value of a questionnaire's reliability is a subjective matter that may depend on the nature and qualitative aspects of a particular study. For this empirical research, the authors of this article selected 0.7 as the lowest marginal value of *Cronbach's alpha* coefficient.

4 The Results of the Expert Evaluation

To accomplish the main purpose of the empirical research, the results of the expert evaluation were systematized. A concept, factor or any other aspect under consideration was treated as important if its average rank was equal to or exceeded 3.5. The value of Cronbach's alpha coefficient estimated for the questionnaire was equal to 0.98, which proposed that the questionnaire reflected the dimension under research with appropriate accuracy.

Availability of the funding mechanisms when making the investment in land in Lithuania. The value of Kendall's coefficient of concordance was equal to 0.578 ($p = 0.000$). The experts unanimously indicated that the most available funding mechanisms when making the investment in land in Lithuania are money, loans and mortgages attributable to the category of the traditional mechanisms of private investment. The other funding mechanisms, such as the traditional mechanisms of public investment or the innovative real estate funding mechanisms, are not popular in Lithuania, in particular when it concerns the investment in land. This tendency can be related not only to the imperfections of the legal framework in the country, but also to poor attractiveness of the domestic real estate market in comparison to foreign real estate markets. The results are summarized in Table 1.

Table 1
Unpopular funding mechanisms while making the investment in land in Lithuania

Funding mechanism	Mean	Minimum	Maximum	SD
The traditional mechanisms of private investment				
Use of a part of the share capital	3.40	3	4	0.548
RELPs	2.00	1	3	1.000
CREFs	1.80	1	3	1.095
REITs	2.20	1	3	1.095
Real estate mutual funds	2.20	1	3	1.095

SWFs	2.00	1	3	1.000
The traditional mechanisms of public investment				
Bond issuance	2.60	1	4	1.517
IPOs	2.60	1	4	1.517
Sale and leaseback	3.00	1	5	1.871
ABSs	2.20	1	3	1.095
CMBSs	2.20	1	3	1.095
Real-estate related derivatives (property index certificates, forwards)	2.20	1	3	1.095
Public subsidization of real estate projects	2.60	1	4	1.140
Lower VAT tariffs for housing	2.00	1	3	1.000
The innovative real estate funding mechanisms				
Corporate zones	3.40	2	4	1.140
Subordination loans	3.20	3	4	0.447
Joint ventures	3.40	3	4	0.548
Granting of rights	3.20	2	4	0.837
Outsourcing contracts	2.60	1	4	1.517
Transfers of holdings	3.40	1	5	1.517
Private funding initiatives (PFI or PFI2)	2.80	1	4	1.643
Reinvestment of return surplus	2.80	1	5	1.789
Tax increment financing (TIF)	2.80	1	5	1.789
Accelerated development zones	2.80	1	4	1.304
Unsecured loans	2.40	1	4	1.517
“Crowd funding”	2.80	1	4	1.643
“Peer-to-peer” (P2P) model	2.80	1	4	1.643

Source: compiled by the authors with reference to the results of the expert evaluation

The results of the empirical research have revealed that in the group of the innovative real estate funding mechanisms, *tax increment financing* (with mean rank equal to 3.6), *free economic zones* (with mean rank equal to 3.8), *simplification of real estate planning procedures* (with mean rank equal to 3.6.) and *additional business taxation* (with mean rank equal to (3.6) are considered to be significant, although less available land funding mechanisms. According to the experts, the other mechanisms from the same group are not popular when making the investment in land in Lithuania. *Summarizing, it can be stated that money (cash), loans and mortgages are the most available traditional funding mechanisms commonly employed when making the investment in land in Lithuania.*

When Saulius Vagonis and Marius Dubnikovas were asked the question “What are the sources of land investment funding if land is treated as a kind of real estate?”, the experts clarified that with reference to the data of the Bank of Lithuania, only an insignificant share of transactions in the land market are funded by employing loans. The Scandinavian bank policy was indicated as the main reason for this tendency: as the Scandinavian real estate market is overcoming the period of crisis, the excessive requirements are imposed on land funding. Issuance of corporate bonds serves as another source of land investment funding. Although theoretically the mechanism of “crowd funding” could also be employed, the cases of its employment have not been registered thus far. According to Saulius Vagonis, the funding of land investment in Lithuania is limited. This limitation was partly caused by the financial crisis of 2008, when banks used to lend money to the investors who wanted to acquire land plots driven by unreasonable expectations, but did not have any debt repayment capabilities. In addition, before the beginning of the crisis, many investors had acquired illiquid land plots which later turned out to be out of demand and caused the average price of land plots to decrease by 80 percent. After the burst of the real estate price bubble, banks stopped funding the investment in land plots. Now land plots are mainly acquired with personal funds as well as with loans or mortgages issued by banking institutions (in many cases, by credit unions).

Conclusions

Based on the analysis of scientific literature, it is possible to make the generalizations, that the main characteristics of innovative real estate funding mechanisms is a combination of loans and guarantees, which helps to add value to both land and buildings. The most significant innovations in the area of real estate funding cover solidarity mechanisms, public-private partnerships and the mechanisms of loans and bonds. The funds accumulated in the form of various taxes and fees have become a part of the funding of real estate development, like fixing the land value or tax exemptions which promise investors lower tax tariffs.

Public-private partnerships are invoked when the schemes of the long-term obligations in the private sector need to be matched up to long-term assets. The private funding of real estate (with consideration of both public-private partnership projects and private funding initiatives) allows for the pooling of funds from the public and private sectors, for a mutual purpose and allocate the investment-related risks. Finally, the innovative loan and bond mechanisms help to attract private investments and direct them to well-organized capital markets. The integration between transportation and land value fixing can also be considered as a part of the real estate funding innovations. By following the approach of refunding, the stages of the development of particular assets can be modified, and the development of the assets can be ensured, by promoting the sales and the flows of commercial revenues.

The results of the expert evaluation have revealed that the innovative funding mechanisms are neither popular nor available, when making an investment in land in Lithuania. The unavailability of the innovative real estate funding mechanisms was determined by the lessons of the financial crisis of 2008, before which, banks used to lend the money to the investors, who wanted to acquire land plots driven only by unreasonable expectations, but did not have any debt repayment capabilities. Now land plots are mainly acquired with personal funds, as well as, with loans or mortgages issued by banking institutions (in many cases, credit unions).

References

- [1] Abor, J. (2007) Corporate governance and financing decisions of Ghanaian listed firms. *Corporate Governance: The International Journal of Business in Society*, 7(1), 83-92
- [2] Acquah, P. (2011) Residential development and borrowing in Ghana: a challenge for banks and private estate developers. Retrieved from <http://ir.knust.edu.gh/xmlui/bitstream/handle/123456789/6434/Aquah%20Patrick.pdf?sequence=1>
- [3] Adair, A., Bery, J., Gulati, M., Haran, M., Hutchison, N., Kashyap, A., McCord, M., McGreal, S. Oyedele, J., & Tiwari, P. (2011) *The future of private finance initiative and public private partnership. RICS Research*. London: Royal Institution of Chartered Surveyors
- [4] Ali, Z., McGreal, W. S., Adair, A. S., & Webba, J. R. (2006) Corporate real estate strategy in the UK and Malaysia. *Journal of Corporate Real Estate*, 8(4), 168-17
- [5] Augustinaitis, A., Rudzkienė, V., Petrauskas, R. A., Dagytė, I., Martinaitytė, E., Leichteris, E., Malinauskienė, E., Višnevskā, V., & Žilionienė, I. (2009) *Lietuvos e. valdžios gairės: ateities išvalgų tyrimas*. Vilnius: Mykolas Romeris University
- [6] Baležentis, A., & Žalimaitė, M. (2011) Ekspertinių vertinimų taikymas inovacinių veiksnių plėtros analizėje: Lietuvos inovatyvių įmonių vertinimas [Application of expert evaluations in the analysis of the development of the innovative factors]. *Management theory and studies for rural business and infrastructure development*, 3(27) 23-31
- [7] Barnett, C. (2015) Trends show crowdfunding to surpass VC in 2016. Retrieved from <https://www.forbes.com/forbes/welcome/?toURL=https://www.forbes.com/sites/chancebarnett/2015/06/09/trends-show-crowdfunding-to-surpass-vc-in-2016&refURL=&referrer=>
- [8] Bartke, S. (2013) Science for environment policy – thematic issue: Brownfield regeneration. Retrieved from <http://ec.europa.eu/environment/integration/research/newsalert/pdf/39si.pdf>

-
- [9] Beracha, E., Downs, D. H., & MacKinnon, G. (2017) The 4% rule: does real estate make a difference? *Journal of Property Research*, 34(3), 181-201, <https://doi.org/10.1080/09599916.2017.1293134>
- [10] Bilal, S., & Kratke, F. (2013) Blending loans and grants for development: an effective mix for the EU? Retrieved from <http://ecdpm.org/publications/blending-loans-grants-for-development-effective-mix-eu/>
- [11] Breuer, W., & Kreuz, C. (2011) Real estate and real estate finance as a research field – an international overview. Retrieved from https://papers.ssrn.com/sol3/papers.cfm?abstract_id=1775760
- [12] British Property Federation (2008) Tax increment financing: a new tool for funding regeneration in the UK? Retrieved from <https://www.bpf.org.uk/sites/default/files/resources/Tax-Increment-Finance-tool-for-funding-regeneration.pdf>
- [13] Brown, K. C., Harlow, W. V., & Starks, L. T. (1996) Of tournaments and temptations: an analysis of managerial incentives in mutual fund industry. *The Journal of Finance*, 51(1), 85-110
- [14] Carter, A. (2006) Building an investment market for economic development. *Local Economy*, 21(1), 65-72
- [15] Coleman, A., & Grimes, A. (2009) Fiscal, distributional and efficiency impacts on land and property taxes. Retrieved from <https://treasury.govt.nz/sites/default/files/2017-11/tfr-lpt-1nov09.pdf>
- [16] Delfim, J. C., & Hoesli, M. (2016) Risk factors of European non-listed real estate fund returns. *Journal of Property Research*, 33(3), 190-213. <https://doi.org/10.1080/09599916.2016.1199590>
- [17] Ezimuo, P. N., Onyejiaka, C. J., & Emoh, F. I. (2014) Sources of real estate finance and their impact on property development in Nigeria: a case study of mortgage institutions in Lagos metropolis. *British Journal of Environmental Research*, 2(2), 35-58
- [18] Greenhalgh, B., & Squires, G. (2011) *Introduction to building procurement*. Abingdon: Taylor and Francis
- [19] Gonenc, H. (2005) Comparison of debt financing between international and domestic firms: evidence from Turkey, Germany and UK. *International Journal of Managerial Finance*, 1(1), 49-68
- [20] Griggs, J., & Kemp, P. A. (2012) Housing allowances as income support: comparing European welfare regimes. *International Journal of Housing Policy*, 12(4), 391-412
- [21] Grishankar, N. (2009) Innovating development finance: from financing sources to financial solutions. Retrieved from

- http://siteresources.worldbank.org/CFPEXT/Resources/CFP_Working_Paper_No1.pdf
- [22] Haffner, M. E. A., & Boelhouwer, P. (2006) Housing allowances and economic efficiency. *International Journal of Urban and Regional Research*, 30(4), 944-959
- [23] Havard, T. (2013) *Financial feasibility studies for property development: theory and practice*. London: Routledge
- [24] HM Treasury (2012) *A new approach to public private partnerships*. London: HM Treasury Publishing
- [25] Iblher, F., & Lucius, D. I. (2003) Innovative real estate financing in Germany – a financial desert? *Property Management*, 21(1), 82-96, <https://doi.org/10.1108/02637470310464490>
- [26] JLL (2017) Financing China's real estate: pragmatism and creativity will prevail. Retrieved from <http://www.joneslanglasalle.com.cn/china/en-gb/Research/china-real-estate-financing-report-en.pdf>
- [27] Kane, M. J. (2001) Equity investment in real estate development projects: a negotiating guide for investors and developers. *The Real Estate Finance Journal*, Spring, 5-9
- [28] Kemp, P. (2007) *Housing allowances in comparative perspective*. Bristol: The Policy Press
- [29] Kennard, M., & Bond, S. (2011) Interest soars in US peer-to-peer lending 2016. Retrieved from <https://www.ft.com/content/2345e94a-0bb1-11e1-9a61-00144feabdc0#ixzz21aKdQtwL>
- [30] Knight Frank. (2017) Analysis of institutional funding in real estate. Report 2017. Retrieved from <https://kfcontent.blob.core.windows.net/research/782/documents/en/analysis-of-institutional-funding-in-real-estate-4526.pdf>
- [31] Knuth, S. E. (2015) Global finance and the land grab: mapping twenty-first century strategies. *Canadian Journal of Development Studies*, 36(2), pp. 163-178, <https://doi.org/10.1080/02255189.2015.1046373>
- [32] Lasfer, M. (2007) On the financial drivers and implications of leasing real estate assets: the Donaldsons-Lasfer's Curve. *Journal of Corporate Real Estate*, 9(2), 72-96, <https://doi.org/10.1108/14630010710828090>
- [33] Liu, T., & Wilkinson, S. (2014) Large-scale public venue development and the application of public-private partnerships (PPPs). *International Journal of Project Management*, 32(1), 88-100
- [34] Makridakis, S., Wheelwright, S. C., & Hyndman, R. J. (1998) *Forecasting: methods and applications*. New York: John Wiley & Sons

- [35] Malhotra, N. K., & Birks, D. F. (2003) *Marketing research: an applied approach*. Harlow: Pearson Education
- [36] Marseguerra, G., & Cortelezzi, F. (2009) Debt financing and real estate investment timing decisions. *Journal of Property Research*, 26(3), 193-212, <https://doi.org/10.1080/09599911003669625>
- [37] McQuaid, R. W., & Scherrer, W. (2008) Public private partnerships in the European Union: experiences in the UK, Germany and Austria. *Uprava*, 6(2), 7-34
- [38] Medda, F. R., Caschili, S., & Modelewska, M. (2012) Financial mechanisms for historic city core regeneration and Brownfield redevelopment. In G. Licciardi, & R. Amirtahmasebi (Eds.), *The Economics of Uniqueness: Investing in Historic City Cores and Cultural Heritage Assets for Sustainable Development* (pp. 213-240)
- [39] Mwathi, J. K. (2013) The effect of financing sources on real estate development in Kenya. Retrieved from http://erepository.uonbi.ac.ke/bitstream/handle/11295/58809/Mwathi_The%20effect%20of%20financing%20sources%20on%20real%20estate%20development%20in%20Kenya.pdf?sequence=3
- [40] Nkyi, B. A. (2012) Strategies for financing real estate development in Ghana. Retrieved from <http://ir.knust.edu.gh/bitstream/123456789/4699/1/Benjamin%20Appiagyeyi%20Nkyi.pdf>
- [41] Nunnally, J. C., & Bernstein, I. H. (1994) *Psychometric theory* (3rd ed.) New York: McGraw-Hill
- [42] Ogedengbe, P. S., & Adesopo, A. A. (2003) Problems of financing real estate development in Nigeria. *Journal of Human Ecology*, 14(6), 425-431
- [43] Olsson, F. (2015) Alternative financing options of corporate real estate. Retrieved from <http://www.diva-portal.org/smash/get/diva2:839664/FULLTEXT01.pdf>
- [44] Redman, A. L., Tanner, J. R., & Manakyan, H. (2002) Corporate real estate financing methods: a statistical study of corporations' choices. *Journal of Corporate Real Estate*, 4(2), 169-186
- [45] Rudzkiene, V., & Burinskienė, M. (2007) *Plėtros kryptių vertinimo ir valdymo informaciniai modeliai [The models of the development direction assessment and management]*. Vilnius: Technika
- [46] Scottish Futures Trust. (2013) Revenue funded infrastructure investment in Scotland. Retrieved from https://www.scottishfuturestrust.org.uk/files/publications/Pipeline_of_Revenue_Funded_Projects_NPD_and_hub_April_2013.pdf

- [47] Squires, G. (2015) Innovative financing in real estate development for urban regeneration. Retrieved from https://www.researchgate.net/publication/268742084_Innovative_Financing_in_Real_Estate_Development_for_Urban_Regeneration
- [48] Squires, G., Hutchinson, N., Adair, A., Berry, J., McGreal, S., & Organ, S. (2016) Innovative real estate development finance – evidence from Europe. *Journal of Financial Management of Property and Construction*, 21(1) 54-72, <http://dx.doi.org/10.1108/JFMPC-09-2015-0036>
- [49] Steinisch, M. (2012) Peer-to-peer lending survey. Retrieved from https://www.consumer-action.org/news/articles/2012_p2p_lending_survey/
- [50] The European Public-Private Partnership Expertise Centre. (2010) Capital markets in PPP financing: where we were and where are we going? Retrieved from http://www.eib.org/attachments/epec/epec_capital_markets_en.pdf
- [51] Tiwari, P., & White, M. (2014) *Real estate finance in the new economy*. New Jersey: Wiley-Blackwell
- [52] van der Krabben, E., & Heurkens, E. (2014) Netherlands: a search for alternative public-private development strategies from neighbouring countries. In G. Squires, & E. Heurkens (Eds.), *International Approaches to Real Estate Development* (pp. 66-81)
- [53] van der Krabben, E., & Needham, B. (2008) Land readjustment for value capturing, a new planning tool for urban redevelopment. *Town Planning Review*, 79(6), 651-672
- [54] Vicent, S. (2015) Real estate financing techniques and sources. Retrieved from https://www.researchgate.net/publication/303881343_Real_Estate_Financing_Techniques_and_Sources
- [55] Wall, A., & Connolly, C. (2009) The private finance initiative. *Public Management Review*, 11(5), 707-724
- [56] Webber, C. (2010) *Tax increment financing*. London: Centre for Cities
- [57] Weber, R. (2002) Extracting value from the city: neoliberalism and urban redevelopment. *Antipode*, 34(3), 519-540

Experimental Testson the Plasticity and Deformability Characteristics of Several Stainless Steel Grades used for Hydro–Pneumatic Equipment's Manufacturing

Vasile Alexa, Imre Kiss

University Politehnica Timișoara, Faculty of Engineering Hunedoara
Department of Engineering and Management
5, Revolutiei, 331128 Hunedoara, Romania
e-mails: vasile.alexu@fih.upt.ro,imre.kiss@fih.upt.ro

Abstract: Many of the hydraulic and pneumatic devices are made from high quality stainless steels, through complex and elaborated manufacturing technologies. Thus, not infrequently the semi-finished product used for obtaining the pneumatic and hydraulic equipment is subjected to different kinds of strains in the manufacturing process. Obviously, the units that are currently producing pneumatic and hydraulic equipment should focus on the manufacturing of products mostly requested in a market economy. This requires the modernization of existing production capacities in line with the EU requirements, followed by the update of technologies to the standards applied in the EU economy. The knowledge about the characteristics of deformability has for the technologist, as well as for the designer and researcher, a great practical significance, because they are important elements in establishing a correct technological process. The change of deformation conditions existing in the industrial process, such as the temperature and rate of deformation, are difficult to consider for correcting the deformability determined by testing. In this paper, through "deformability" we understand all the properties characterizing the deformation behavior of the metals and alloys, and the „deformation resistance" of the metals is expressed through the unit strain required to produce a certain degree of plastic deformation, under the conditions of a particular diagram of tensions, deformations and deformation rates, in the absence of external friction forces. This study includes the results of the experimental tests conducted to find the plasticity and deformability characteristics of several stainless steel grades: one martensitic stainless steel (grade X46Cr13), one ferritic stainless steel (grade X6Cr17) and one austenitic stainless steel (grade X5CrNi18–10).

Keywords: plasticity; deformability; stainless steel grades; temperature; heating; tests

1 Introductory Remarks

In many industries and in many types of technical operations, the hydraulic equipment requires steel to withstand high operating temperatures combined with the corrosive action of the environment. These requirements cannot be met without the proper development of the high-alloy and quality steel manufacture, including the thermostable stainless steels. [1-3, 5-6, 16]

Currently, we know various types of stainless steels, which have multiple features and properties, designed to withstand corrosive environments, various working conditions, and weathering, thus providing safety conditions in enterprises, longer life in constructions and hygiene in everyday life. [1-3] The stainless steels are used in all industries today: mechanical engineering, metallurgical fields, medical equipments and instruments, ship-, automotive- and aviation-building, food processing, energy and power engineering, chemical and petrochemical, traffic engineering, construction, etc.

The knowledge about the characteristics of deformability has for the technologist, as well as for the designer and researcher, a great practical significance, because they are important elements in establishing a correct technological process. [2-6, 10-11] The change of deformation conditions existing in the industrial process, such as the temperature and rate of deformation, are difficult to consider for correcting the deformability determined by testing. [2, 5-7, 10,11, 13]

In view of this, the deformability is the ability of a material to be plastically deformed without the occurrence of undesired conditions (cracking or tearing of the material during the plastic deformation, inadequate quality of the surface, wrinkling or curling of the stamped steel sheets, coarse structure, difficulty of material flowing when filling the moulds, or other commercially-imposed conditions). [2, 6, 10,11, 13, 16]

The stainless steels can undergo structural changes under the action of the following technological processes: [5, 10, 11]

- a heat treatment(required by the manufacturing process);
- a cold plastic deformation(austenitic steels);
- annealing,after cold deformation;
- a high temperature thermo-mechanical treatment (e.g. required for hot rolledsteel or subjected to mechanical stress at high temperature).

Regardless of the adopted method for deformability determination, when the technological process are decided, the people involved should bear in mind that, the results have a relative value, i.e. they are significant only in comparison with other steels, whose plastic deformation behavior as deformability indices are already known. [2-6, 10-11]

The processing of metals and alloys via plastic deformation is based on the property of plasticity, which defines their ability to acquire permanent deformations under the action of external forces. [2-6, 8-12, 14] When processing by plastic deformation, the shape modification of a semi-finished product is made by redistributing its elementary volumes under the action of external forces; therefore, unless some unavoidable losses due to equipment imperfection, the processing takes place without any removal of material.

The deformability of metals and alloys characterizes their ability to permanently deform without breaking the internal structural bonds. [2-5] As the deformability of a material is expressed by the degree of deformation to which the first cracks appear, i.e. its tearing resulting from a standard mechanical test or from one specific to the industrial deformation process, it should be pointed out that the breaking process, for all industrial processes of plastic deformation, as well as for the materials plastically deformed in these processes, appears in the form of ductile fracture. [2-6, 10, 11]

The main factors that influence the deformability can be grouped into two categories: [2-6, 10, 11]

- material related factors: composition, structure, purity, metallurgical development, localization of the deformation;
- process related factors: deformation temperature, deformation rate, state of stress and strain, hydrostatic pressure, friction between the tool and workpiece, geometry of the tool and workpiece.

In determining the hot deformability of steels in the laboratory, in general, but especially those stainless, the following conditions in which the plastic deformation takes place under industrial conditions must be taken into account: [5, 6, 10, 11]

- steel heating temperature;
- deformation temperature;
- tensions scheme where the deformation occurs;
- steel-tool contact friction;
- steel structure at the deformation temperature;
- steel deformation rate.

There are several methods for determining the deformability of the steels, such as: [2-6,10]

- compression, rolling and forging (taking account of friction);
- tensile, bending and torsion (without taking account of friction).

The above mentioned methods enable that, besides the determination of deformability characteristics (plasticity and deformation resistance, depending on temperature), to study the influence of the deformation conditions (rate of heating, holding time at heating temperature, friction with the tools, rate of deformation, structural changes in terms of deformation, rate of recrystallization, etc.). [5, 10, 11]

2 Determination of the Stainless Steel Deformability by Torsion

This method is the only one that allows obtaining large deformations along the length of the specimen, so it is mainly used to determine the characteristics of large deformations. [5, 10, 11]

Since the shear strains play an important role in the process of rolling and forging, the deformability caused by torsion reflects quite accurately the steel behavior at hot plastic deformation, and due to the fact that the specimen can be maintained in the oven during deformation, we can ensure the stability of temperature. By this method, the hot deformability of the stainless steel is determined by subjecting to torsion a cylindrical specimen maintained at the deformation temperature in a tubular oven. [5, 10, 11]

The size of the required moment for torsion the specimen expresses the resistance to deformation, and the number of torsions before failure expresses the plasticity limit of that steel. [5, 10, 11]

There are several methods for determining the deformability by hot torsion, such as: [5, 10, 11]

- torsion by maintaining the specimen at constant length;
- torsion by tensioning the specimen;
- torsion with free change of the specimen length.

2.1 Torsion by Maintaining the Specimen at Constant Length

This variant of the method for determining the hot deformability, which does not imply a deformation through pure shearing, is preferable because the rate of deformation can be easily maintained constant throughout the specimen. [5, 10, 11]

The torsional deformation on the specimen surface is:

$$\gamma = r \cdot \frac{\theta}{l} [\text{rad}] \quad (1)$$

in which: r – the specimen radius, [m]; θ – torsion angle,[rad]; l – the specimen length, [m].

Then, the deformation rate on the specimen surface is:

$$v_{\gamma} = \frac{d\gamma}{dt} = \frac{d\left(\frac{r}{l}\theta\right)}{d\theta} \cdot \frac{d\theta}{dt} = \frac{r}{l} \cdot \frac{d\theta}{dt} = \frac{r}{l} \frac{2 \cdot \pi \cdot n}{60} [\text{rad/sec}] \quad (2)$$

in which: $d\theta$ – the specimen's torsion speed, [-]; t – time, [sec], n – number of rotations, [rot/min].

From the value of the torque, we can calculate the shear resistance on the specimen surface:

$$\tau = \frac{(3-k) \cdot M}{2 \cdot \pi \cdot r^3} [\text{N/m}^2] \quad (3)$$

in which: M – the torque moment, [N·m]; r – the specimen radius, [m]; k – sensitivity coefficient expressing resistance to the speed of deformation, for a particular steel, at a certain temperature, [-].

This method, although it does not imply a deformation through pureshearing, it is preferable because the rate of deformation can be easily maintained constant throughout the specimen.

2.2 Torsion by Tensioning the Specimen

The specimen heated in the oven is subjected to a constant tensile strain during the torsion test. Having no calibrated portion, the deformation length depends on the heated length of the test specimen and, therefore, the reproducibility is poor and the deformability characteristics are only informative, being not suitable for scientific processing and interpretation.[5, 10, 11]

Because the specimen is maintained under constant tensile strain, we cannot speak about the balancing of the axial force, which tends to shorten the specimen during cooling and hence its dimensions are changing during the test, resulting the fact that the deformation on the specimen surface does not occur at a constant speed, and due to changes in diameter we cannot calculate the shear resistance on the specimen surface.

2.3 Torsion with Free Change of the Specimen Length

For taking into account the results of the hot torsion test measurements in determining the resistance to deformation, the specimen is necessary to do not miss its straightness, and the deformation to be uniform throughout the specimen; in this case, instead of the axial force, we measure the length of the specimen during its torsion.[5, 10, 11]

Based on the law of specimen volume constancy before and after the deformation, we can write:

$$r_0^2 \cdot l_0 = r^2 \cdot l \quad (4.1)$$

or

$$r = r_0 \sqrt{\frac{l_0}{l}} [\text{m}] \quad (4.2)$$

in which: r – the specimen radius, before and after the deformation, [m]; l – the specimen length before and after the deformation, [m].

As at high temperatures, the deformation resistance is a function of the deformation rate and not of the degree of deformation (due to recrystallization), its value at the specimen surface is:

$$\tau = \frac{1}{2\pi} \left[3 \cdot \frac{M}{r^3} + v \cdot \gamma \cdot \frac{\partial \left(\frac{M}{r^3} \right)}{\partial v \gamma} \right] [\text{N/m}^2] \quad (5)$$

in which: M – the torque moment, at a time, [N·m]; r – the specimen radius, at a time, [m]; v – the deformation speed at the specimen surface.

The radius variation in time and the torsion moment variation with the deformation speed can be obtained only by logarithmisation. Therefore, it is preferred to plot diagrams for representing the torsion moment variation versus the number of torsions and the torsion rate.

3 The Research Methodology

The experimental equipment used to study the stainless steel deformability by hot torsion belongs to the Faculty of Engineering Hunedoara, University Politehnica Timișoara.

The facility is provided with a central shaft on which two side discs and an intermediate disc are mounted in the central area. Spacer bushes have been mounted between the left side disc and the intermediate one, as well as between the intermediate disc and the right side one, capable of keeping the discs at a distance, for fixing the experimental samples, the specimens.

The so-equipped central shaft is connected to an electric motor which provides its rotation along with the specimens. At the top of the facility, above the central area of the shaft (where the experimental sample sare fixed), is placed an electric oven which provides the sample heating in the range 20-1300°C. The temperature is maintained at the desired value by means of a control box, and the speeds can be changed by attaching to the electric motor of astatic frequency converter.

The ensemble of the experimental equipment used to study the stainless steel deformability by hot torsion, with and without the heating oven, is shown in Fig. 1.

The specimens for hot torsions were mechanically taken from $\Phi 20$ mm hot-rolled steel bars, having the form and dimensions presented in Figure 2. The test specimens are typically cylindrical, with a calibrated small-diameter central portion, having the ration $\frac{l}{d} = 5$ in the point of deformation.



a.



b.

Figure 1

The experimental facility for determining the hot deformability of the stainless steels
a. without the heating oven; b. with the heating oven

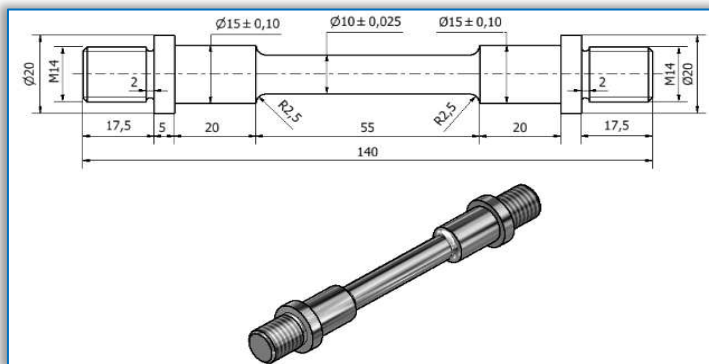


Figure 2

Sample for determining the hot deformability by torsion

The ends are screwed, and the specimen must have a shoulder in the continuation of the thread, to prevent further screwing during the torsion.

The choice of heating regime is currently mostly based on the practical experience of the companies; therefore, the process of establishing the hot processing technology for these steels is primarily related to the definition of heating conditions, according to their technological characteristics.

4 Results and Discussions

For the experimental tests, we used several stainless steel grades. This study includes the results of the tests conducted to find the plasticity and deformability characteristics of the martensitic stainless steel, grade X46Cr13 (Table 1 and Table 2), the ferritic stainless steel, grade X6Cr17 (Table 3 and Table 4) and the austenitic stainless steel, grade X5CrNi18–10 (Table 5 and Table 6).

Table 1

The results of the tests conducted to find the plasticity characteristics of the martensitic stainless steel (hardenable stainless steel, grade X46Cr13)

No.	Deformation temperature [°C]	Maximum torque moment [daN·cm]				
		1	2	3	4	Average
1.	800	274	300	240	250	266
2.	850	242	268	250	280	260
3.	900	266	276	258	269	267.25
4.	950	194	191	188	174	186.75
5.	1000	156	158	134	134	145.50
6.	1050	127	121	119	118	121.25
7.	1100	101	83	112	112	102
8.	1150	90	93	98	88	92.25
9.	1200	69	69	60	40	57
10.	1250	47	48	45	–	46.66

Table 2

The results of the tests conducted to find the deformability characteristics of the martensitic stainless steel (hardenable stainless steel, grade X46Cr13)

No.	Deformation temperature [°C]	The number of torsions before failure [–]				
		1	2	3	4	Average
1.	800	6	5	7	5	5.75
2.	850	5	7	10	10	8
3.	900	9	10	8	9	9
4.	950	10	10	11	13	11
5.	1000	13	11	12	12	12
6.	1050	13	13	12	13	12.75
7.	1100	14	13	14	13	13.75
8.	1150	7	14	14	14	12.25
9.	1200	8	8	8	8	8
10.	1250	8	9	7	–	8

Table 3

The results of the tests conducted to find the plasticity characteristics of the ferritic stainless steel (non-hardenable stainless steel, grade X6Cr17)

No.	Deformation temperature [°C]	Maximum torque moment [daN·cm]				
		1	2	3	4	Average
1.	800	135	135	140	–	136
2.	850	126	126.5	123	–	125.16
3.	900	108	112	111	–	110.33
4.	950	94	73	77	–	81.33
5.	1000	63	57	57	–	59
6.	1050	9	22	38	–	23
7.	1100	83	36	41	–	53.33
8.	1150	29	28	29	–	28.66
9.	1200	21	21	20	–	20.66
10.	1250	18	–	14	–	16

Table 4

The results of the tests conducted to find the deformability characteristics of the ferritic stainless steel (non-hardenable stainless steel, grade X6Cr17)

No.	Deformation temperature [°C]	The number of torsions before failure [–]				
		1	2	3	4	Average
1.	800	31	34	42	–	35.66
2.	850	29	22	26	–	25.66
3.	900	27	17	29	–	24.33
4.	950	34	33	28	–	31.66
5.	1000	35	36	48	–	39.66
6.	1050	62	58	68	–	62.66
7.	1100	15	71	75	–	53.66
8.	1150	105	69	94	–	89.33
9.	1200	43	57	134	–	78
10.	1250	460	–	425	–	443

Table 5

The results of the tests conducted to find the plasticity characteristics of the austenitic stainless steel (non-magnetic stainless steel, grade X5CrNi18–10)

No.	Deformation temperature [°C]	Maximum torque moment [daN·cm]				
		1	2	3	4	Average
1.	800	200	392	340	390	330.5
2.	850	192	362	314	359	306.75
3.	900	276	326	306	316	306
4.	950	194	230	220	227	218.25

5.	1000	170	176	156	194	174
6.	1050	144	130	142	130	136.5
7.	1100	133	123	127	129	128
8.	1150	98	–	100	80	92.66
9.	1200	97	83	84	77	81
10.	1250	58	44	61	–	53.33

Table 6

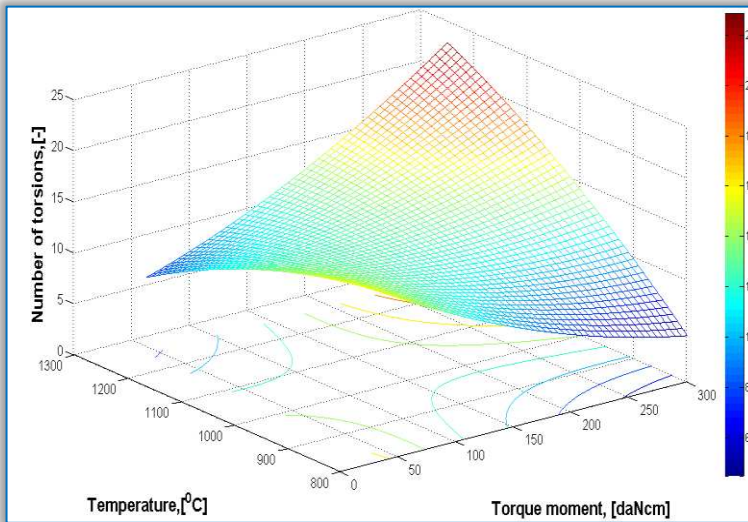
The results of the tests conducted to find the deformability characteristics of the austenitic stainless steel (non-magnetic stainless steel, grade X5CrNi18–10)

No.	Deformation temperature [°C]	The number of torsions before failure [–]				
		1	2	3	4	Average
1.	800	2	2	2	2	2
2.	850	3	3	4	2	3
3.	900	4	2	3	3	3
4.	950	6	5	8	4	5.75
5.	1000	4	6	7	3	5
6.	1050	9	8	8	7	8
7.	1100	10	8	15	12	11.25
8.	1150	9	–	9	9	9
9.	1200	9	12	6	6	9
10.	1250	7	8	6	–	7

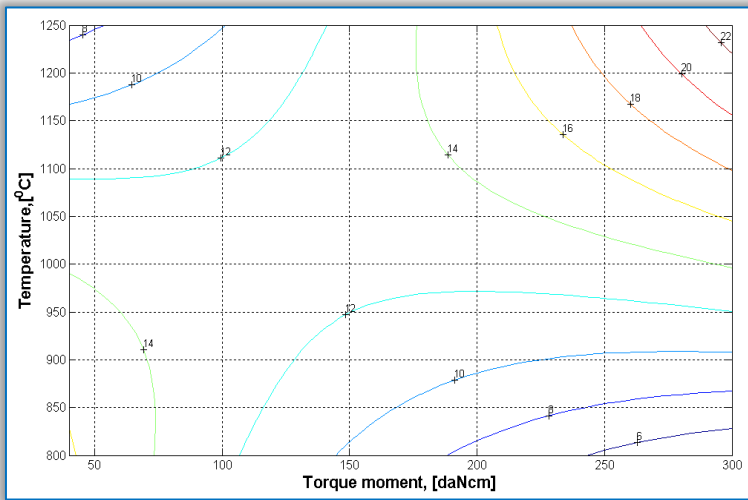
The austenitic stainless steel (nonmagnetic stainless steel), grade X5CrNi18–10, is the standard for the austenitic grades of stainless steel due to its good corrosion resistance, ease of formability and fabrication.

The ferritic stainless steel (non-hardenable stainless steel), grade X5CrNi18–10, is resistant to corrosion in most environments. Although the corrosion resistance of X6Cr17 is inferior to the austenitic grades of stainless steels, its ferritic microstructure makes it resistant to the effects of stress corrosion cracking, a form of corrosion to which most of the conventional austenitic stainless steels are susceptible to. The X6Cr17 is characterized by its good corrosion resistance is displayed in moderately corrosive media/environments.

The martensitic stainless steel (hardenable stainless steel), grade X46Cr13, is characterized by its good corrosion resistance in moderately corrosive environments. Stainless heat-resistant steels are always in demand when extreme technical requirements are imposed on the material, due of their outstanding chemical corrosion and mechanical properties.



a.



b.

Figure 3

Deformability diagram for the martensitic stainless steels (grade X46Cr13), at the experimental heating temperature values (800–1250°C)

a. the regression surface of plasticity and deformability characteristics, described by the number of torsions before failure [equation type: $z = a_{(1)}x^2 + a_{(2)}y^2 + a_{(3)}xy + a_{(4)}x + a_{(5)}y + a_{(6)}$, standard deviation: $r^2 = 0.8298$]

b. the level curves and the technological domains area

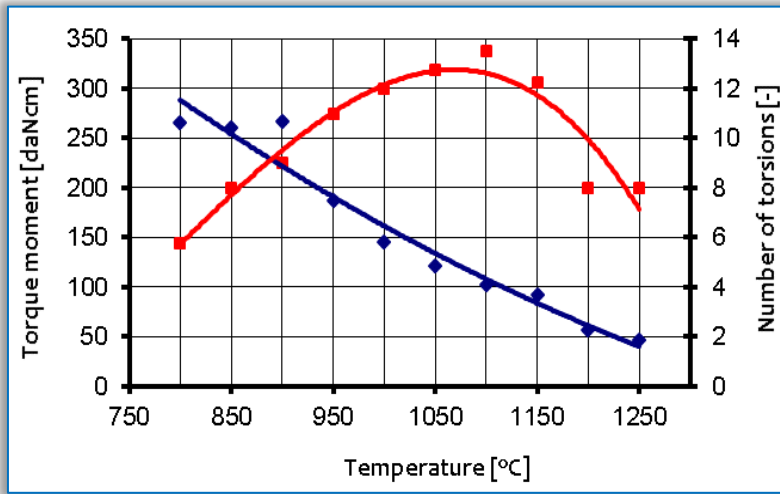


Figure 4

The variations of plasticity (number of torsions to failure) and deformation resistance (maximum torque) in case of the martensitic stainless steels (grade X46Cr13), at the experimental heating temperature values (800–1250°C)

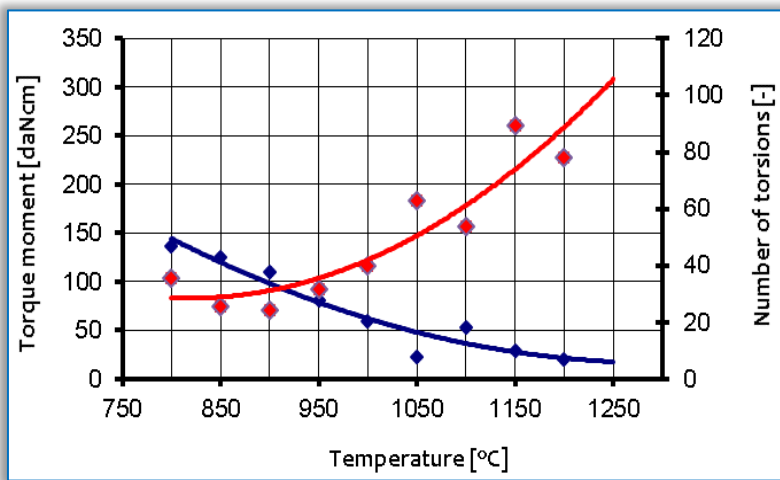
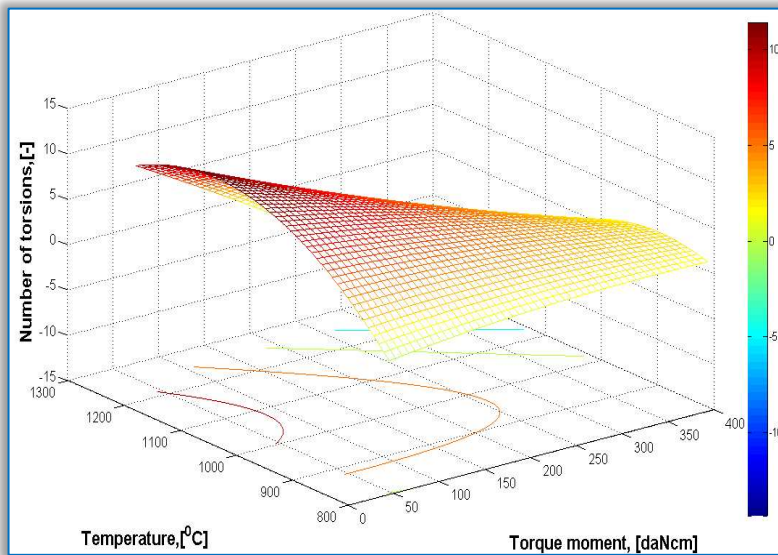
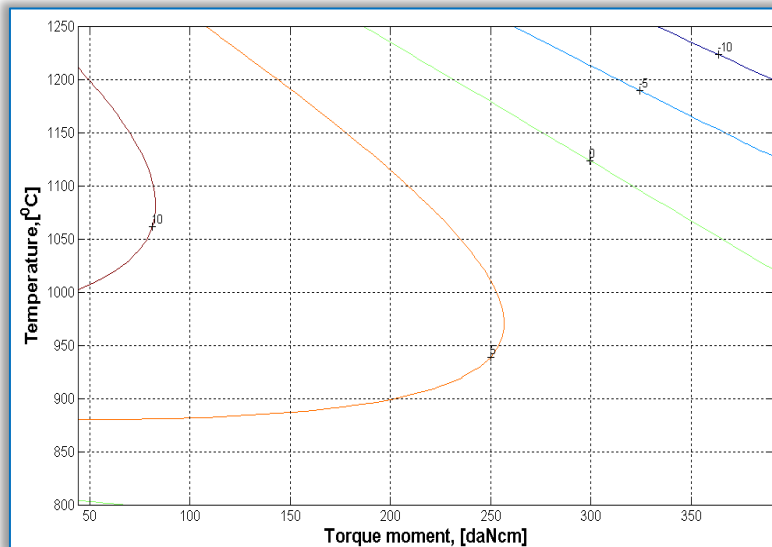


Figure 5

The variations of plasticity (number of torsions to failure) and deformation resistance (maximum torque) in case of the ferritic stainless steel (grade X6Cr17), at the experimental heating temperature values (800–1250°C)



a.



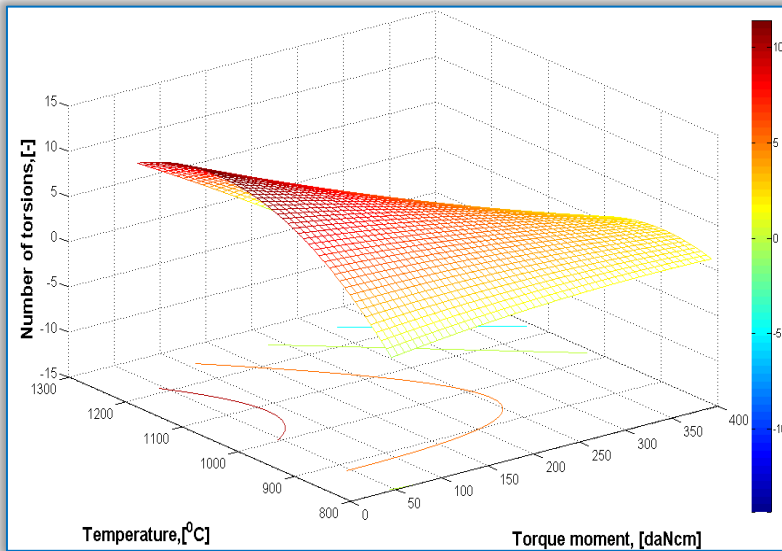
b.

Figure 6

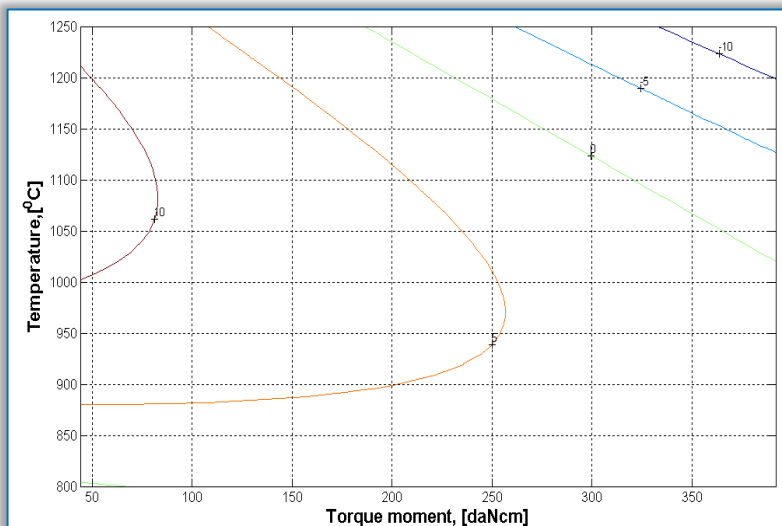
Deformability diagram for the ferritic stainless steel (grade X6Cr17), at the experimental heating temperature values (800–1250°C)

a. the regression surface of plasticity and deformability characteristics, described by the number of torsions before failure [equation type: $z = a_{(1)}x^2 + a_{(2)}y^2 + a_{(3)}xy + a_{(4)}x + a_{(5)}y + a_{(6)}$, standard deviation: $r^2 = 0.8056$]

b. the level curves and the technological domains area



b.



c.

Figure 7

Deformability diagram for the austenitic stainless steel (grade X5CrNi18–10), at the experimental heating temperature values (800–1250°C)

a. the regression surface of plasticity and deformability characteristics, described by the number of torsions before failure [equation type: $z = a_{(1)}x^2 + a_{(2)}y^2 + a_{(3)}xy + a_{(4)}x + a_{(5)}y + a_{(6)}$, standard deviation: $r^2 = 0.8056$]

b. the level curves and the technological domains area

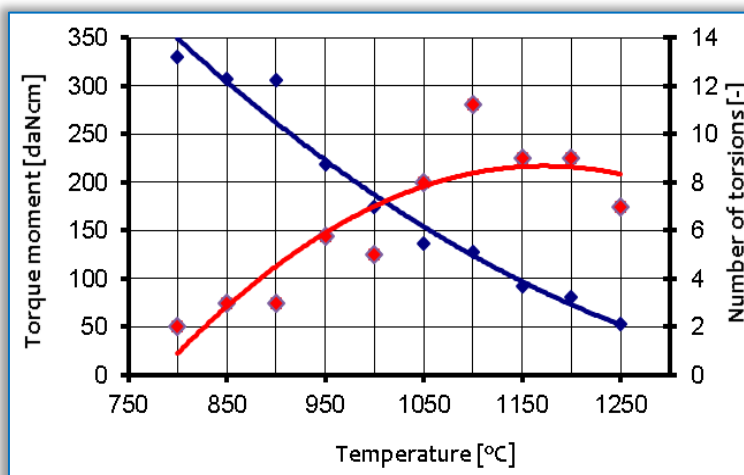


Figure 8

The variations of plasticity (number of torsions to failure) and deformation resistance (maximum torque) in case of the austenitic stainless steel (grade X5CrNi18-10), at the experimental heating temperature values (800–1250°C)

For the hot torsion test, we prepared 36 samples from each steel grade. They were subjected to torsional deformation by maintaining the deformation temperature in the experimental facility, from 50 to 50°C, within the range 800-1250°C.

The magnitude of the torque required to the specimen's torsion expresses the resistance to deformation, and the number of torsions to failure expresses the plasticity limit of that steel. The plasticity limit is expressed by the number of torsions to failure at a given temperature and deformation rate. Each point within the temperature range studied in the two diagrams (Figures 3-8) represents the arithmetic mean of four determinations.

In the graphical representation of the experimental tests results, presented above in the Figures 3-8, we have the following comments and remarks:

- ☐ the variations of plasticity (number of torsions to failure) and deformation resistance (maximum torque) are plotted in the Figure 4, Figure 5 and Figure 8. The variations, as shown in the above mentioned figures, indicate that the deformation resistance of a stainless steel (regardless of the steel grade) decreases with increasing the heating temperature; for the martensitic and austenitic steel grades (Figure 4 and Figure 8), due to the measurement error at high values of the torsion moment at low temperatures, the value of the maximum torque is lower than the technological requirement;
- ☐ the regression surface of plasticity and deformability characteristics of the martensitic stainless steels (grade X46Cr13), described by the number of torsions before failure, is shown in Figure 3(a); This can be interpreted as

deformability diagram, plotted in shown in Figure 3(b), which is typical for the martensitic stainless steels, X46Cr13 being such a steel grade;

- ☐ the regression surface of plasticity and deformability characteristics of the ferritic stainless steel (grade X6Cr17), described by the number of torsions before failure, is shown in Figure 6(a); This can be interpreted, plotted in shown in Figure 6(b), as deformability diagram, which is typical for the ferritic stainless steel, grade X6Cr17;
- ☐ the regression surface of plasticity and deformability characteristics of the austenitic stainless steel (grade X5CrNi18–10), described by the number of torsions before failure, is shown in Figure 7(a); This can be interpreted, plotted in shown in Figure 7(b), as deformability diagram for the grade X5CrNi18–10, which illustrate a much better the deformation resistance than the ferritic stainless steel grade X6Cr17, being an austenitic stainless steel grade;
- ☐ the upper limit of the optimum range of heating temperatures applied for deforming the studied steels, results clearly from the plasticity – temperature diagrams, as follows:
 - 1100°C, for the martensitic stainless steel, grade X46Cr13 (Figure 4);
 - 1050°C, for the ferritic stainless steel, grade X6Cr17 (Figure 5);
 - 1150°C, for the austenitic stainless steel, grade X5CrNi18–10 (Figure 8);
- ☐ the temperature may be limited due to the risk of excessive grain growth during heating under industrial conditions (phenomenon that does not occur during heating at the torsion machine – and, therefore, the values given for plasticity at high temperatures); [5, 10, 11, 15]
- ☐ regarding the end heating temperature, for the hot deformation of the studied stainless steel grades, we have the following experimental values (or ranges):
 - 900-950°C, for the martensitic stainless steel; it has a lower limit due to the high deformation resistance and the cracking hazard;
 - 800°C, for the ferritic stainless steel; sometimes it is recommended that the last two passes (processing) to be carried out at temperatures below 800°C, for completion of granulation;
 - 950°C, for the austenitic stainless steel.

5 Conclusions

This study includes the results of the experimental tests conducted to find the plasticity and deformability characteristics of several stainless steel grades: one martensitic stainless steel (grade X46Cr13), one ferritic stainless steel (grade X6Cr17) and one austenitic stainless steel (grade X5CrNi18–10).

The indications regarding the variation of plasticity with the temperature, using the hot torsion method, allowed for establishing the temperature range within which the steel plasticity is optimal and in which, in general, it is recommended to perform the entire hot plastic deformation. Also, depending on the plasticity variation with temperature, we can achieve a more rational distribution of the reduction coefficients per passes, so that the plasticity property of the steel to be used as much as possible.

Starting from the temperature of 900°C, all steel grades have a sufficient plasticity, but the value of the deformation resistance is still high up to the temperature of 950°C. The growth dynamic of the plasticity characteristics is continuous, reaching the maximum value at the temperature of 1250°C, while reducing the resistance to deformation. Thus, from the tests carried out to determine the hot deformability, it results that the optimal plasticity of the analyzed steels is found within the temperature range 950-1250°C.

Acknowledgement

This facility is subject to a patent registered with the State Office for Inventions and Trademarks (OSIM) under number 439/17.05.2010, entitled "Facility adapted for experimental determination of the resistance to thermal fatigue of samples placed tangentially on the generator of support discs", No. 54/2011.

References

- [1] A. J. Sedriks, Corrosion of stainless steel, John Wiley and Sons, New York, 2nd Edition (1996)
- [2] E. Cazimirovici, M. V. Suci, Rolling the special metallic materials (in Romanian), BREN Press, Bucuresti (2000)
- [3] D. Peckner, I. M. Bernstein, Handbook of stainless steels, McGraw-Hill, New York (1977)
- [4] G. Krauss, Steels: Processing, structure and performance, ASM International (2005)
- [5] J. Magaone, Studies and research on the thermostable stainless steel behavior in the plastic processing (in Romanian), Politehnica Press, Timișoara, (2010)
- [6] M. Trusculescu, A. Ieremia, Stainless and refractory steels (in Romanian), Facla Press, Timisoara, 1983
- [7] A. Momeni, K. Dehghani, Characterization of hot deformation behavior of 410 martensitic stainless steel using constitutive equations and processing maps, Materials Science and Engineering A, 527, 21-22, (2010), 5467-5473
- [8] D. Kuc, G. Niewielski, E. Hadasik, K. Radwanski, Structure and mechanical properties of hot deformed ferritic steel, Archives of Civil and Mechanical Engineering, 3 (2004), 85-92

-
- [9] D. Kuc, G. Niewielski, Technological plasticity and structure in stainless steels during hot-working, *Journal of Achievements in Materials and Manufacturing Engineering*, 32, 2 (2009), 154-161
- [10] J. Magaone, I. Ilca, V. Alexa, Hot deformability of austenitic stainless steel, *Știința și Inginerie*, 22, AGIR Press, București (2012), 241-248
- [11] J. Magaone, I. Ilca, Influence factors analysis and steels behavior within the high temperature operating, *Știința și Inginerie*, 16, AGIR Press, București (2009), 629-636
- [12] G. Niewielski, D. Kuc, Structure and properties of high-alloy steels, *Plasticity of Metallic Materials*, Silesian University of Technology Press, Gliwice (2004)
- [13] T. S. Byun, N. Hashimoto, K. Farrell, Temperature dependence of strain hardening and plastic instability behaviors in austenitic stainless steels, *Acta Materialia*, 52, 13 (2004), 3889-3899
- [14] Y. C. Lin, X.-M. Chen, A critical review of experimental results and constitutive descriptions for metals and alloys in hot working, *Materials and Design*, 32, 4 (2011), 1733-1759
- [15] I. Ilca, Optimization of hot-metal working of austenitic stainless steels, *Journal of Engineering Sciences and Innovation*, 2, 3 (2017), 103-117
- [16] EN 10088-2: 2005 Stainless steels – Technical delivery conditions for sheet/plate and strip of corrosion resisting steels for general purposes

Supportive Robotic Welding System for Heavy, Small Series Production with Non-Uniform Welding Grooves

Csongor Márk Horváth, Péter Korondi

Department of Mechatronics, Optics and Mechanical Engineering Informatics,
Faculty of Mechanical Engineering, Budapest University of Technology and Economics,
Műegyetem rakpart 1, 1111 Budapest, Hungary
hcsongorm@mogi.bme.hu, korondi@mogi.bme.hu

Abstract: Heavy welding is a demanding task with high robotization potential. This applies especially for the runners of Francis hydropower turbines, due to the high working costs and EHS requirements in Europe. However, heavy welding is often related to small-series production with long processing time. This sets high demands on the planning and monitoring functionality of the robot system. The research in this field is gaining momentum, yet very few articles suggest suitable solutions. This paper presents a robotic welding control system design and application that facilitates the planning, control, and monitoring of the welding process of non-uniform grooves of large-dimension joints. Its primary and unique characteristic is the simplified operator assisted programming method, where the three-dimensional path modification problem is translated into consecutive two-dimensional modifications. Therefore, reference cross-sections are created along the welding groove, where the sequence planning task of multi-pass weld bead placement is performed, and to the online modifications together with the adjustments are referred. The planning, changes and process supervision are supported by the robot system to handle uncertainties along the welding groove and adaptively utilize the robot operator experience. The activities are tracked and organized to supply information for later performance enhancement and reusability between similar processes. The supportive system design is particularly suitable for advanced, large-dimension, heavy robotic welding applications. A use case is presented on a welding a runner of Francis hydropower turbine.

Keywords: robotic welding; multi-pass welding; non-uniform groove; small series production

1 Introduction

1.1 Welding Robots for Small and Medium Sized Companies

The industry is facing major challenges increasing efficiency and productivity to stay competitive. The small and medium sized enterprises (SMEs) are an essential part of the countries' economy, as they represent 99 per cent of all enterprises. The domain of industrial robot usage and integration has been dominated by the large-scale automotive and electronics industries [19, 30]. With 27%, the automotive sector is the largest in the welding industry. From all application fields, the most common is welding and soldering (30%), which typically implemented for large volume production that requires high product mix and short production cycle. Although recent trends show an expansion of robot adoption outside of these areas, the progression into new fields is moderate.

Even though the SMEs are showing increasing demand for robotization, their demands differ from the traditional robot applications because their business models are more likely to involve wide range and small series production [26]. The tasks are often not well defined, heavy, fatiguing and hazardous, with substantial environmental load and stress level for the workers [35]. The limited proof of performance of the technologies are the technical barriers that limit the adoption of robotic systems by SMEs even in the most desired application areas.

Despite the quality and efficiency that a today's robotic welding systems can provide for the general welding industry, skilled human welders cannot yet be replaced in welding of joints in complex structures due to various reasons: high initial costs, tedious teaching procedure and long commissioning time. Thus, most of the welding is done manually or semi-automatically in fields such as the off-shore industry, ship manufacturing or hydropower turbine production [17].

1.2 Challenges in Heavy Multi-Pass Welding

Several challenges arise when the application comes to robotic heavy welding despite the convenience of using robotic welding systems. Typical challenges related to the small series production are the following:

1. Cost and personnel: SMEs have limited resources; The high initial costs of installations with the lack of dedicated and specialized personnel restrict the possibilities to deploy robotized solutions as well as the use of complex off-line programming systems at SMEs facilities.
2. Task complexity: Large-dimension welding joints typically have thick and non-uniform welding grooves. Therefore, significant amount of time and so-

phisticated approach are necessary to handle the multi-pass welding process. Such grooves are often still welded manually due to their complex shape.

3. Environment: The heavy, fatiguing and hazardous manual welding, with substantial environmental load and stress level for the workers effects directly the production and need to be conformed with the Environmental and health (EHS) regulations
4. Programming time: Small series production requires significant effort spent on the programming of the welding robot for the new part. The currently available robotic solutions are lacking a detailed model based multi-pass welding planning. An accurate multi-pass welding plan can shorten the preparation and welding time.
5. Handle uncertainties: Robots cannot make corrective decisions autonomously. Thus, decision making support is required, either by the sensors and the control system or through intuitive user interaction. Detailed and accurate knowledge about the process increase the applicability range of the planning, but, additional online handling of the arising uncertainties is inevitable.

The welding groove complexity of large-dimension joints originates mainly from the geometry and the varied thickness of the base materials. Regardless of the careful edge preparation and the standard conformity, weld joints have thick non-uniform grooves. Such examples are the tubular joints in pipeline manufacturing or the grooves on the hydropower turbine runners at the blades.

This paper presents a robotic welding control system design and application that facilitates the planning, control, and monitoring of the welding process of non-uniform grooves. Its primary and unique characteristic is the simplified operator assisted programming method. It contains an offline programming module with dedicated consideration of the non-uniformities of the welding groove and the simplified online programming module, supporting the welding path adjustment and process supervision to handle uncertainties. The supportive system design is presented on a use case with a runner of Francis hydropower turbine.

2 Background

2.1 Welding Robot Systems

Welding robots represent the largest fraction of applications deploying industrial manipulators. The most common techniques apply Metal Inert/Active Gas welding (MIG/MAG), the Tungsten Inert Gas (TIG), and Laser Beam Welding. Currently, automated robotic welding is gaining momentum due to the high wage

levels and the dropping installation and operation cost of a robot system. This offers new opportunities to automate small series production, although these are the result of several stages of development in the welding robot systems.

In the earliest, first generation robotic welding applications, the welding was performed in two runs; the first run was dedicated to learning the seam geometry and the second run was the actual tracking and welding. The second generation of robotic welding systems' development reduced the number of necessary runs by performing seam learning and tracking simultaneously, in real-time. The latest, third-generation welding robot systems are not only operated in real-time but within unstructured environments and learning the rapidly changing geometry of the seam during operation [36].

According to Pires [36], an automated robotic welding system design can be implemented in three different phases with the final goal to achieve decent performance and a high-quality weld. The first phase is the preparation, where the welding scene is set up and the offline programming is executed. The second phase is the welding phase, when the welding process is performed based on the continuous decisions made by the operator or the robot system to achieve the required weld quality. The last phase is the analysis phase, in which the welds are examined, and a decision made about the acceptance. The considered changes are collected and evaluated.

2.2 Hardware Components

Modern welding robot systems contain an integration of the robot manipulator, robot controller, welding equipment, work-piece positioner, supportive sensor system, and welding safety devices [12,31]. Those multiple units require coordinated or synchronized motion to access the entire work-piece, minimize idle time and maximize the arc/welding time. It often connected to a sensor system supporting the welding process and a computer for process control and data collection. In advanced operations, the standard computer peripheries are extended by additional Human-Machine Interfaces (HMI). A schematic of a general robot system is shown in Figure 1. Similar equipment used for the realization of the robotic welding system presented in this paper.

Sensors in robotic welding are used to detect and measure process features along with geometrical parameters, or monitor and control welding process parameters by technological sensors [13, 21, 48]. The first can be achieved in several ways applying most often optical sensors to detect and measure the joint geometry (seam finding, seam tracking) [49, 50], as well as the weld pool geometry and location [9, 37]. Research on robot systems for small series production has been conducted to determine the main factors for the users. Besides the flexibility, user-friendliness, shorter programming time, and robustness of operation, the possibility to integrate sensors both for simulation and during runtime was listed as signif-

icant. In this context, sensors used for seam tracking or to control the welding process are considered equally important [3]. Weld quality monitoring in robotic welding provides automatic detection of weld defects by analysing the process parameters and by comparing these with the nominal values [38]. It also could include non-destructive inspection methods such as radiography, ultrasonic, vision, magnetic detection, eddy current, acoustic measurements [55] or electromagnetic sensor [1].

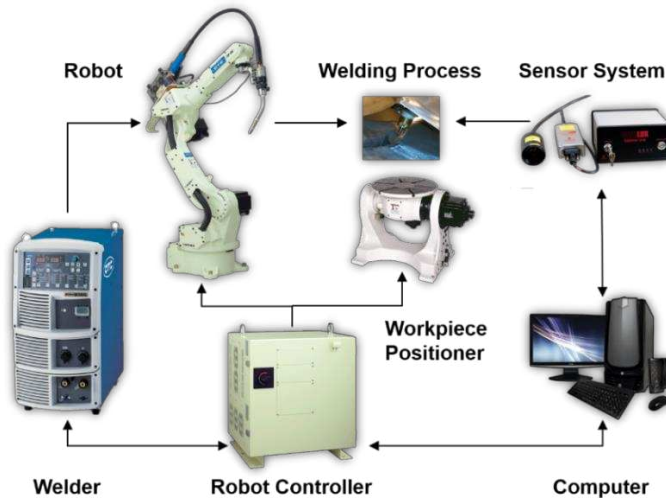


Figure 1
General robotic welding system

Due to the challenging formation of the high temperature welding environment (high current, spatter, liquid metal, high temperature), it is difficult to apply sensors to measure the welding parameters directly. These problems cause that the parameters that can be observed are not concurring with the parameters needed to be controlled. Furthermore, it is not trivial to carry out a simple feedback control. The complexity can be solved by developing models to map the observable parameters to appropriate actions on issues within the relevancy of the welding specification procedure. In this, the productivity and quality measures are defined together with the nominal welding process control parameters and geometry information to produce the desired weld. A model based control should, therefore, unify the data from the sensors, the welding procedure specifications and the robotic welding system specific restrictions [36].

2.3 Programming Methods

Two main categories of programming methods exist in practical industrial applications: online programming, including the lead-through and walk-through, and offline programming (OLP). Conventional online programming allows for precise control of the straightforward process with simple path definitions and work-piece geometry. Due to the low initial cost and low programming skills required, it is widely used. However, the entire production line is disrupted during teaching due to the downtime of the robot. Moreover, the taught program has limited flexibility and is unable to adapt to the current welding scenario and problems encountered in the welding operation without additional control [34].

More advanced programming methods are the operator-assisted online programming, such as the lead- and walk-through methods or the sensor guided programming. By walk-through programming, the robot arm itself is configured to be able to be moved by the operator, to teach the robot path based on the built-in [2, 7, 42] or external sensors [41]. Furthermore, experiments and research have been conducted to develop admittance controller driven teaching methods, deploying external tools [27, 44] and vision systems [33, 43, 45]. Besides the progress achieved on the online programming to make it more intuitive and fitting of the operator skills, most of the research outcomes are still not commercially available [34].

Using OLP methods, data based on CAD/CAM is a common practice in many areas of the industry, especially automation systems with large product volumes. Figure 2 illustrates the workflow of OLP. Many software and simulation tools are available to provide direct robot trajectories from CAD data of the work-pieces, robots, and fixtures used in the cell [20]. Some of the most advanced techniques apply the recent results of research in the field of Cyber-Physical Systems [10, 29, 39] and the Digital Twin [32, 46] related developments. The main advantages are that the generated code is reusable, flexible for modifications, and complex paths can be produced with reduced production downtime [18]. However, the OLP systems utilization in SMEs is limited due to the economic disadvantages for small volume production caused by the high cost of the OLP packages and the programming overhead for customization [34].

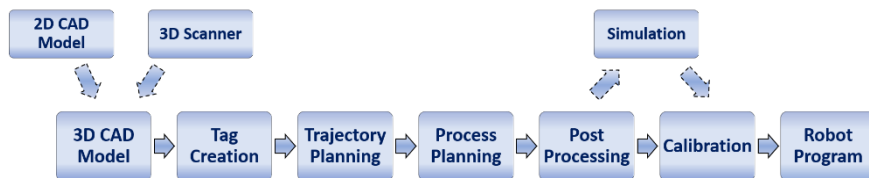


Figure 2

Key steps of offline programming. Reprinted from [34]

In welding, most of the available OLP software is considering the welding seam as a well-defined, uniform groove. The existing planning methods of multi-pass welding [25, 52, 53] based on a generally constant groove cross-sectional area where the differences in the geometry are results of errors. Only a few studies [6, 51] analysed how to handle the non-uniformity of the welding groove geometry systematically. These address the groove representation with straight edges, where the measured profile showed different shape, without consideration of the curvature of the edge preparation. The layer height calculation was based on trigonometrical principles. The introduced welding groove segmentation based on the weld bead placement strategy and the welding position difference. The groove geometry changes affected the weld bead numbers in the layers and the number of the layer number. One of the main conclusions was that the weld bead number in the layers should be constant, but the layer number would vary from segment to segment concerning the welding quality.

2.5 Human Behaviour Models and Human-Machine Interfaces

The mainstream trend in modern welding industry is mechanization and automation. However, human welders may be preferred over mechanized welding control systems in applications where experience-based behaviour in response to the received information is required [54]. Studies have been conducted to develop models of the mechanism of welders' experience-based behaviour to create a controller in automated welding. It has been found, that the welder makes decisions primarily based on past learned experiences and the humanistic approach of the acquired sensory information is imprecise. It only reflects partial truth about the instant status of the welding process [5, 23].

Another approach is to create HMI to overcome the barriers between the process and the operator, by improving the maintenance and support activities through remote communication [4]. This can be exploited by cyber-physical devices [8], cognitive info-communication methods [16, 22], or multi-modal man-machine communication (4MC) [28, 47]. Those latter methods utilize multiple senses of the human and create sensor bridging to transfer the otherwise naturally acquired data (NAD) [22]. Information from one sensor must be translated into another and transferred through non-conventional communication channels (Figure 3). Therefore, the goal of multi-modal human-machine communication is to realize natural, intuitive and efficient information flow between the remote operator and the local system [47] as well as create a virtual environment that makes the remote operator feeling next to the system [11, 14, 15, 24].

Based on the overviewed literature, the guidelines can be identified for the development of a heavy multi-pass welding robot system for SMEs. Cost and time efficient programming method is required to provide alternatives to the expensive and general OLP methods along with the slow but flexible online programming.

The development of a simplified OLP system is defined to achieve the necessary complexity level by automating the auxiliary, non-welding tasks; simplify online programming by developing HMI for the execution of the essential modifications integrated it into the control system. The sensor system must be integrated to support the operator's modification activity.

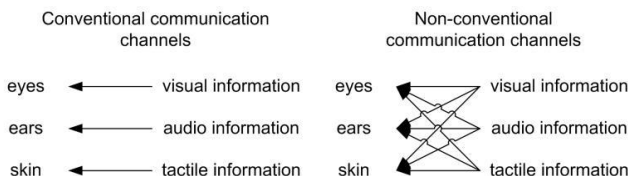


Figure 3

Differences between conventional and non-conventional information channels [16]

The simplified and process-oriented environment could balance out the missing skill set of the robot operator, and the supportive sensory system provides the necessary information to utilize the operator experience in welding.

3 Control System Structure for Heavy, Multi-Pass Robotic welding

This section provides a general description of the system design principles for welding tasks with large-dimension joints and non-uniform grooves. The system design is intended to replace the manual welding procedure directly, but it also needs to be able to compete with the online and offline programming methods. Figure 4 provides the schematic for such a system that can be considered as a cascade control system design. This contains three different control loops with different speed and functions, furthermore divided into the phases discussed in Section 2.1. The process consists of the preparation, offline planning and programming, the welding process control, and finally the observation and analysis.

3.1 Preparation and Offline Programming

The process starts with the welding scene setup, where the preparation includes the work-piece positioning, the welding method and the additional physical components definition (shielding gas, feed wire, preheating). The outermost loop of the cascade control system is offline programming and analysis loop, which performed between the different welding setups. Its forward section contains the offline programming, where the CAD/CAM models are handled. Based on the planning strategy of multi-pass welding and the weld bead models, the weld seam

is filled, and the robot trajectory is generated according to the calibration procedure. The feedback section includes the post weld analysis and the learning to update the planner algorithms for further applications.

The welding joint defined in the CAD model of the work-piece with the given groove geometry and the root weld path. Along this path, two-dimensional cross-sections can be extracted from the model that followed by multi-pass weld bead placement planning applied for each cross-section individually. A sectioning algorithm creates sections along the groove to create a unified weld bead pattern for the segment. The planning phase is closed by the trajectory generation in the model space that translated into robot trajectories after work-piece calibration. The direct paths transferred to the robot controller, where they become executable.

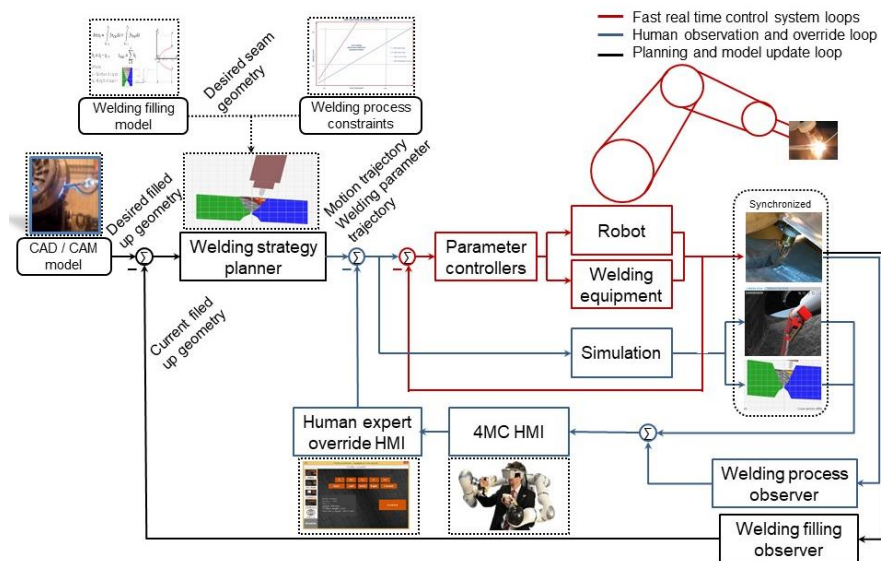


Figure 4
Scheme of the welding robot system

3.2 Welding Process Control

The inner part of the process structure is covering the welding phase with two overlapping control loops. The most inner loop represents the real-time control system of the welding process and the robot motion controlled by the robot controller. The feedback contains the robot system and welding process variables, such as the recent tool position for motion control and the measured values for the welding parameters.

The middle loop is the human interaction loop where the adaption is performed to the immediate situation during the welding process or to the desired path during the path setup and verification. Here, the feedback loop includes the observations of the welding process and the correction actions from the operator. On the given user interface, the cell operator could give commands to the system to perform the predefined sub-tasks that includes the path verifications and the welding executing. Furthermore, it offers path adjustments both during the dry-run and the weld-run.

3.3 Observation and Post-Weld Analysis

The post-weld analysis and observation are performed to validate the welding process goodness and decide about the acceptance or detect the defects of the welding. The proposed system is intended to handle all the available information collected during the preparation and the welding process, including the synchronized data gathered from the robot controller (speed, position and orientation information, input and output values, internal variable values), from the welding power source (variable welding parameters, pre-set welding parameters), and from the cameras and sensors. The data collection extended with the weld qualification measurements (visual inspection, destructive and non-destructive examination methods) can provide the information needed for a well-supported decision to adjust the reference parameters for the future welding processes.

4 Offline Programming and Path Verification

The programming of the robot and the verification of the welding path are linked together, and the proposed system supports this process with minimal user interaction. Figure 5 shows how the same path is represented in the different scenarios: first in the path definition phase, then in simulation, finally the path verification. This section provides descriptions about the offline programming system, including the transformation chain from the predefined machining path definition in model space to executable robot trajectory.

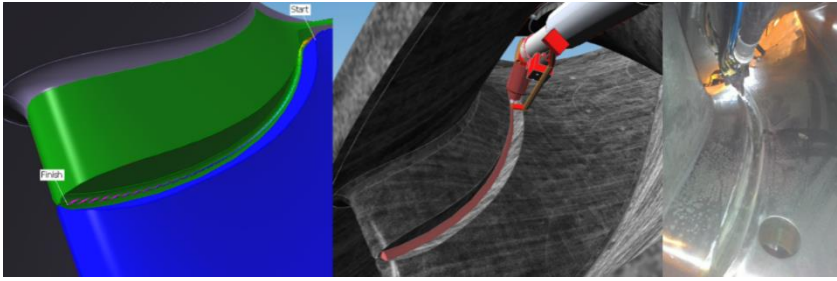


Figure 5

Root welding path verification utilizing a digital twin

4.1 Root Weld Path Definition and Reference Cross-Sections

The root weld path is defined during the offline programming and preparation phase and serves later as a reference trajectory of the multi-pass welding planning. The offline programming tool reads the CAD file of the work-piece then the groove definition is given including the reference cross-sections and the root weld path. The root weld path is built up from task points and normal vectors where the distribution and density of the points define the resolution of the path on the necessary level (straight grooves requires fewer control points compared to curvy grooves) and the normal vectors determining the initial welding torch orientation as shown in Figure 6. The schematic representation of the coordinate system and vector definitions are given in Figure 7. The reference coordinate system for the CAD/CAM data is defined as r , the robot's base coordinate system is defined as b .

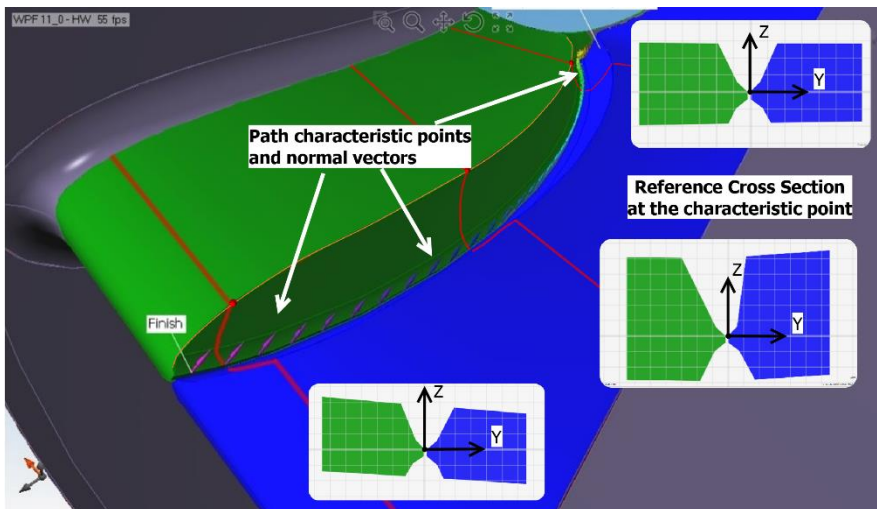


Figure 6

Root weld path trajectory definition in the model space as the digital representation of the work-piece

The task point coordinates \underline{C} are defined in the model \mathbf{r} coordinates and given in the path definition description with the path normal vector \mathbf{a} , which is a physical reference for the initial welding torch orientation. The tangent vector of the path \mathbf{n} is targeting the next task point respecting the predefined task direction. The third vector at the task point s is the cross product of the \mathbf{a} and \mathbf{n} . The task path description in the reference \mathbf{r} model space coordinated system is denoted as $\{T_c\}^r$ that includes each task points and their local coordinate system definition and provides the basis for the robot trajectory planning.

Reference cross-sections are generated from the CAD model along the root weld path to reduce the complexity of the path adjustments and to be used later during the multi-pass welding planning phase. The cross-sections are perpendicular to the path trajectory and defined for each task point on the plane of the local coordinate system \mathbf{t} , represented by the two vectors \mathbf{a} and s , where vector \mathbf{a} defines the z-axis and vector s defines the y-axis. The process of the transformation steps and matrixes is shown in Figure 8. and described in detail in the following.

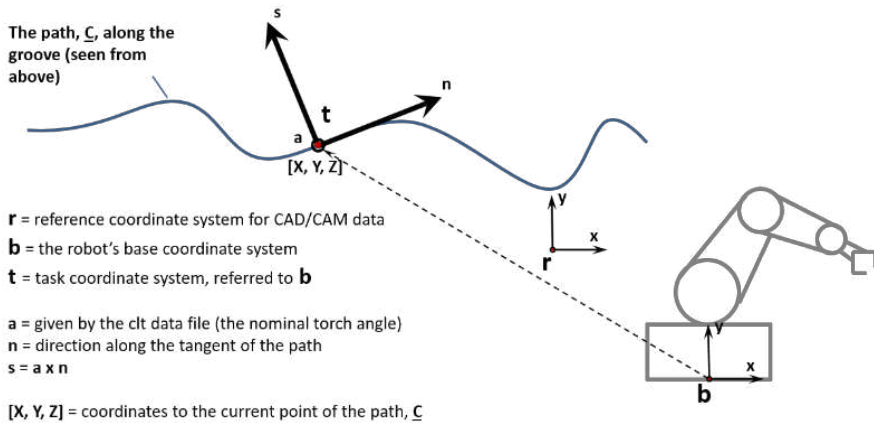


Figure 7

Definition of the reference coordinate systems

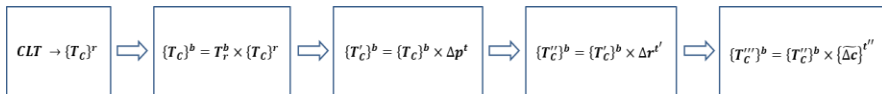


Figure 8

Structure of coordinate transformation – from CAD to executable motion trajectory

4.2 Welding Process Planning

The central part of the process planning in offline programming is the definition of the multi-pass weld bead pattern and the corresponding robot trajectory definition. During the multi-pass welding planning, the main controllable online variable settings collected for each weld bead that influences the welding process, namely arc voltage, arc current, torch travel speed, and wire feed rate. The welding parameters range is defined in the Welding Procedure Specifications as constraints for all weld bead related planning and modelling.

The commercially available welding systems do not contain model-based planning capability considering the weld bead profile properties. Such modules often only generate a symmetric and simplified weld bead layout, which usually requires major adjustments during the operation. In this proposed method, the positions of the weld beads are defined based on certain placement strategies and based on consideration of the groove geometry and the model of the weld bead profile function. Further plan-specific parameters are also included, such as the length of the seams, the welding torch orientation and collision avoidance modifications. The block diagram of the planning process is presented in Figure 9.

The planning process starts with the groove modelling (Block A1), when the groove's mathematical description made for each characteristic cross-section from the digital representation of the work-piece and the weld groove (CAD/CAM or profile scan data as I1-I3). The next step is the generation of the initial weld bead placement sequence in each given groove cross-sections handled by the Sequence Planner (Block A2). The weld bead sizes, shapes and welding parameters are defined by the Welding Filling Model (C1).

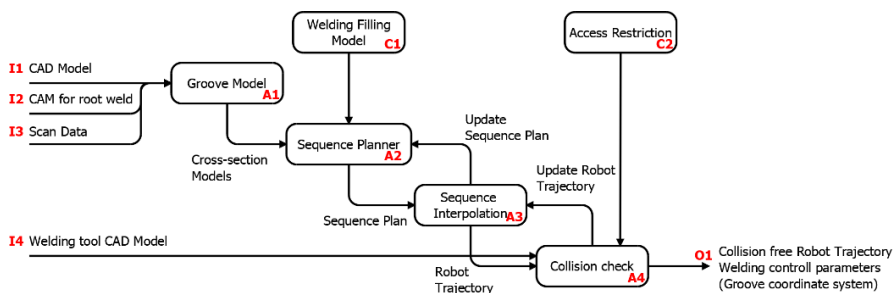


Figure 9

Planning process of multi-pass welding

The model uniqueness lays in the realistic representation of the weld bead profile function in the layer-by-layer deposition, instead of the conventional quadrilateral approximation, described by Yan, et al. [51]. The weld bead shapes are described as symmetric curve functions, and the edge preparation of the grooves defined as continuous convex functions. The produced ripple top surface of the layers is

better suited to reality than the flat surface approximations, therefore, the cumulating error is significantly reduced during the deposition. The exact implementation of the Welding Filling Model and the representation model of the weld bead profile is not synergic part of this paper. When the pattern is generated, the Sequence Interpolation section (Block A3) is activated to assure the pattern smoothness, creating sections for a consistent plan and starting the new iteration process to apply a generally accepted plan. This generates the initial robot trajectories with connected welding parameter settings. The last step (Block A4) is to adjust the recently created robot trajectories concerning the confined space access restriction, to avoid collisions.

4.3 Calibration and Path Definition

The trajectories generated by the multi-pass welding planner are referred to the local coordinate system in the model space but need to be transferred to the robot coordinate system before executions by coinciding with the location of the physical work-piece and the CAD model. This is done by performing a calibration procedure, through determining the position of the same reference coordinate system on the physical work-piece as being used in the virtual world where the CAD / CAM model is defined. During the calibration procedure, the T_r^b transformation matrix determines the translation and rotation from the model space r coordinate system to the robot's base coordinate system b , resulting the new coordinate definition as $\{T_c\}^b$, according to Equation 1.

$$\{T_c\}^b = T_r^b \times \{T_c\}^r \quad (1)$$

5 Online Process Control

By the end of the offline programming and process planning, the input parameters are available for the online process control that is the primary process in the welding phase [40]. The input parameters are the motion trajectory and the welding parameter trajectory. In this section, the block of the online process control is discussed (Figure 4). It includes the control of the physical robot system with the connected devices, the digital twin which is running parallel to the welding process, the welding process observer, which is acquiring the information about the process, and the human-in-the-loop.

The process flow can be described as the following: The reference motion trajectory and welding parameter trajectory are transferred to the parameter controllers. Those reference values translate into executable parameter sets and sections communicated to the physical devices (robot controller and welding power source). The physical signals feedback to the parameter controllers providing stable signals

to the welding process. The control loop implementation is distributed between the physical devices including the factory designed parameter controls. This parameter control with the devices is the most inner loop of the cascade control system. The digital twin is running parallel to this loop including the digital representations of the devices and the work-piece.

The welding process observer is the feedback of the welding process, including the supportive sensor system and overlapped with the information gained from the digital twin. Practically, the latter provides information about the hardly observable parameters, such as the current cross-sections, the already and the future deposited weld beads' reference torch position, as well as collision alerts. The feedback loop includes the human operator, for whom the information is translated through 4MC devices and if necessary overwrites the process references.

5.1 Applying Path Modifications: Translation and Rotation

Our approach to applying path modifications in the welding process is to separate the translation modifications from the rotations. Thus, the three-dimensional path modification problem is translated into consecutive two-dimensional modifications, where the reference cross-sections are serving as a modification plane. The reference cross-sections remain constant during the process regardless of the applied path modifications. The multi-pass welding planning becomes trackable for the operator. The user translation modifications are given along the reference cross-sections main axes as Δy and Δz , relative to \mathbf{t} task point. The new point \mathbf{t}' is the result of the translation $\Delta \mathbf{p}^t$ defined by $\{\mathbf{T}'_c\}^b$ as shown in Equations 2.

$$\{\mathbf{T}'_c\}^b = \{\mathbf{T}_c\}^b \times \Delta \mathbf{p}^t \quad (2)$$

The user rotations $\Delta \mathbf{R}_{x, \varphi}$ and $\Delta \mathbf{R}_{y, \theta}$ are applied to the translated point \mathbf{t}' , the transformation is combined as $\Delta \mathbf{r}^{t'}$ and is applied to resulting the new orientation transformation $\{\mathbf{T}''_c\}^b$ at the task point, according to Equation 3 and 4. The physical meaning of those transformations is that the $\Delta \mathbf{R}_{x, \varphi}$ defines the rotation of working angle of the torch by φ , the $\Delta \mathbf{R}_{y, \theta}$ defines the rotation of the travel angle by θ . Rotation around the path tangent vector (x -axis) is applied when the penetration on the groove face needs to be increased by asymmetrical heat distribution.

$$\Delta \mathbf{r}^{t'} = \Delta \mathbf{R}_{x, \varphi} \times \Delta \mathbf{R}_{y, \theta} \quad (3)$$

$$\{\mathbf{T}''_c\}^b = \{\mathbf{T}'_c\}^b \times \Delta \mathbf{r}^{t'} \quad (4)$$

Both, the translation and rotation modifications made by the operator can be applied to refine the predefined paths on the multi-pass welding plan to increase its accuracy and provide processed data for further analysis to enhance the planning.

5.2 Collision Avoidance and the Final Combined Transformation

In the confined working area, the final path transformation should be made to avoid collisions. The rotations are applied in the reference coordinate system r along the vectors a and n . The resulting transformation matrix is denoted by Δc^r as the cross product of rotation $\Delta R_{y,\theta}$ (around a) and $\Delta R_{z,\psi}$ (around n) (Equation 5). However, the Δc^r transformation should first change its base from r to t'' , therefore, Equation 6 should be applied to calculate $\Delta c_i^{t''}$. Introducing maximum limit for angle change in the collision avoidance $\Delta R_{z,\psi_{max}}$ and $\Delta R_{y,\theta_{max}}$ and performing the examination test in Equation 7, the limited rotation transformation would be $\widetilde{\Delta c}_i^{t''}$ and the final combined path description would be $\{T_C'''\}^b$ (Equation 8).

$$\Delta c^r = \Delta R_{y,\theta} \times \Delta R_{z,\psi} \quad (5)$$

$$\Delta c_i^{t''} = (T_C''^b)^{-1} \times T_r^b \times \Delta c_i^r \Rightarrow \Delta c_i^{t''} = T_r^b \times \Delta c_i^r \times T_C''^b \quad (6)$$

$$\text{Test: } \left\{ \begin{array}{l} \text{IF } |\Delta R_{z,\psi,i}| > \Delta R_{z,\psi_{max}} \Rightarrow \Delta R_{z,\psi,i} = (\pm) \Delta R_{z,\psi_{max}} \\ \text{IF } |\Delta R_{y,\theta,i}| > \Delta R_{y,\theta_{max}} \Rightarrow \Delta R_{y,\theta,i} = (\pm) \Delta R_{y,\theta_{max}} \end{array} \right\} \Rightarrow \Delta c_i^{t''} \rightarrow \widetilde{\Delta c}_i^{t''} \quad (7)$$

$$\{T_C'''\}^b = \{T_C''^b\}^b \times \{\widetilde{\Delta c}\}^{t''} \quad (8)$$

As shown above, several transformations need to be applied to achieve the collision-free trajectory in the complex groove geometry including planned multiple and related path definition, the operator modification during the online process and the continuous collision avoidance.

6 Experimental Verification

The proposed welding robot system is intended to replace manual welding methods by offering OLP and system wise process support. The performance of the system is compared to the manual metal arc welding procedure (which is the currently applied welding method for the examined manufacturing facility) and to online programming method. Each test case repeated for each of the three methods. The main properties for comparison of the different welding methods are 1) the total time spent on between the work-piece installation and final welding inspection, 2) time spent on the different tasks and their added value to the process, and 3) quality of the produced weld.

6.1 Experimental Setup

The robotic welding system design was implemented in a test robot cell for manufacturing Francis hydropower turbine runners. The robot cell was built up from an *OTC FD-V20A* high precision welding robot arm with 0.01 mm repetition accuracy, *Fronius MagicWave4000* welding power source including wire feeder unit and TIG welding torch, *PEMA 35 0000 FAS* manipulator unit, *Cavilux* welding camera system, together with additional safety and interfacing subsystems.

For the test setup, the base material of the runner was 1.4313 X3CrNiMo13-4 martensitic stainless steel. Argon 4.6 gas (purity over 99.996%) was used as the shielding gas with a constant 14 l/min flow rate. For deposition, 1.2 mm diameter CN 13/4-IG filler wire was used, continuously fed to the base material that was preheated to 80 °C temperature. The working angle of the welding torch is fixed at 90 degrees to the work-piece.

The range of the welding parameters was defined during the pre-welding procedure qualification, where wider limits were established. The sets were selected to produce heat input between 0.8 kJ/mm and 1.2 kJ/mm using direct current electrode negative (DCEN) current flow. The weld beads were placed in three 30-degree bevel angle V grooves of two 20 mm thick plates on 400 mm length with a gap of 2 mm and a root face of 2 mm. The plan consisted of 37 weld beads of each three test grooves. Their distribution is shown in Figure 10a. The welded structure went through heat treatment to improve the base material's mechanical properties by quenching and tempering.

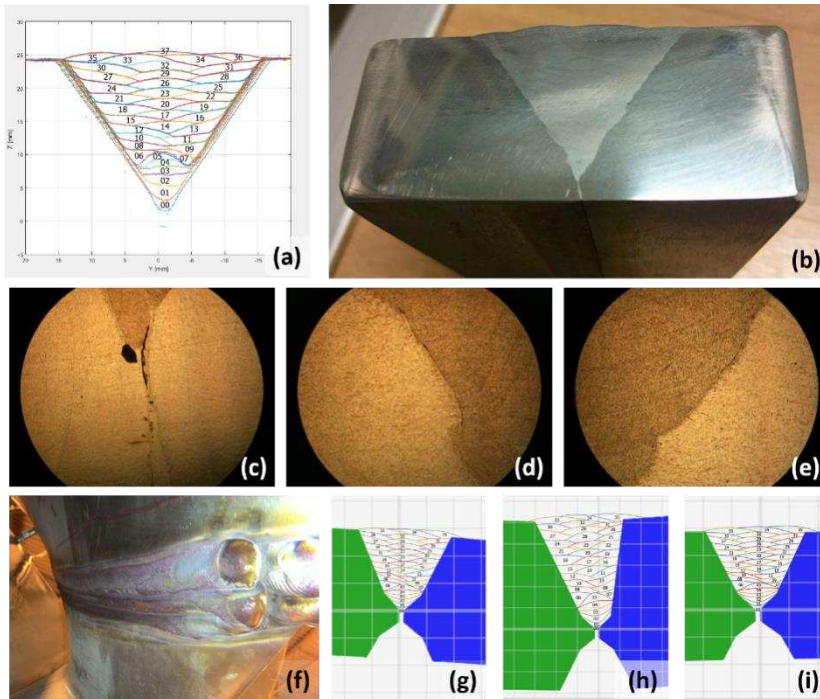


Figure 10

(a) Multi-pass weld bead placement pattern of the test work-piece, (b) prepared cross-section for macro etching, (c) impurity in the root of the weld, (d, e) merge on the height 6.4 mm on the left and right side, (f) filled seam on the runner of Francis hydropower turbine, (g, h, i) multi-pass weld bead placement pattern on the runner

The procedure of the experiments and their analysis followed the NS-EN ISO 15614-1 standard. The seams were examined by non-destructive methods such as penetrant testing, visual and ultrasonic inspection. After the stress relieving heat treatment, the test pieces were cut for destructive mechanical property testing for tensile, hardness and bend test. One of the cross-sections of the robot welded test work-piece is prepared for the macro etching, and the polished surface is shown in Figure 10b. The quality of the weld was examined under a microscope, where the root of the weld showed some impurity (Figure 10c), but the overall fusion found sufficient (Figure 10d-e). The exact test results of the mechanical property tests are not discussed due to industrial partner's restriction on data publication, but they were within the required range for each mechanical property and matched the base material's corresponding nominal values but outstanding excellent impact energy results. The welding parameter ranges defined during the welding procedure qualification test were the followings: arc voltage varies between 11 and 14 V, arc current is DCEN and ranges between 200 and

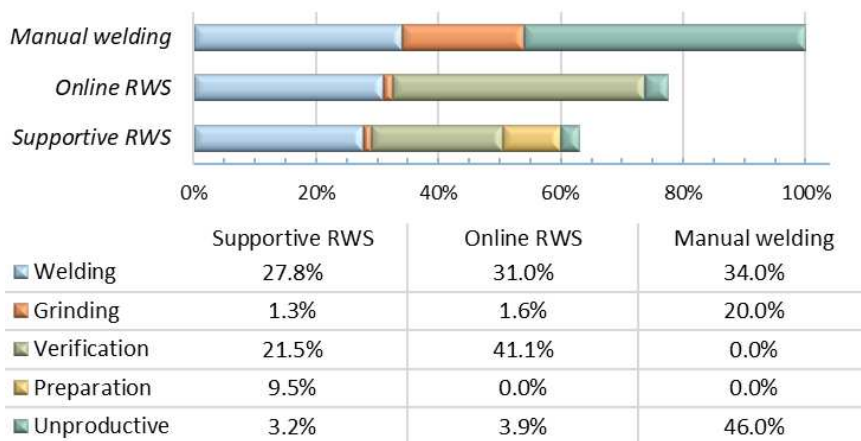
350 A, wire feed rate up to 200 cm/min, and welding speed ranges between 1.5 and 3.5 mm/s.

The runner of the Francis hydropower turbine assembled from 17 blades; the grooves are with double U edge preparation in the middle section of the blades on a 560 mm length. Base material thickness is between 10 and 40 mm and changing gradually along the groove. The predefined welding parameter windows were used in all the three welding test cases, and the weld quality was examined by the previously mentioned non-destructive methods. The filled seam of blades is shown in Figure 10f, and the planned cross-sectional weld bead patterns, in the positions, marked earlier in Figure 6 are presented in Figure 10g-i.

6.2 Evaluation of the Experiments

The baseline for the comparison defined by the total time spent on the manual metal arc welding, where the processing time divided between the welding (34%), grinding (20%) and resting time (46%), later due to the EHS requirements. The performances of the robotized methods are presented in Figure 11.

The online programmed robotic welding robot system (Online RWS) program introduced TIG welding and resulted in significant improvement in most parameters compared to the manual welding. The lead time reduced with 22.4% and the proportion of the welding and grinding tasks improved to 40% and 2%, respectively. The remaining time is utilized as online programming time instead of non-productive resting time.



* values are relative to the total time spent on Manual welding

Figure 11

Distribution of activity time spent on subtasks*

The supportive robotic welding system (Supportive RWS) further reduced the total process time by 18.7% compared to the Online RWS, requiring only 63.1% of the manual process total time. The proportion of welding (44%) and grinding (2%) time is similar the Online RWS, but the introduction of OLP reduced the online programming time significantly, being the main factor of process time improvement.

In manual welding, the time of the process is directly translated into the work-piece, and the gained experience during execution cannot be transmitted to the following work-pieces. Thus, the lead time and the welding quality highly depends on the welder's skills. More consistent quality is achieved by the Online RWS and the Supportive RWS, where the set of welding parameters were defined more precisely, but increased amount of welding defects was detected during the online programming method. Those defects were traced back to the misjudged positioning due to the work-piece limited accessibility and the curvature of the groove. With the Supportive RWS decreased number of welding defects was detected.

Conclusions

In this paper, a supportive robot system design for multi-pass welding was introduced, that can handle non-uniform grooves in small series production. The proposed system design is based on a welding process modelling method as simplified offline programming (OLP), and process execution to support interfacing. The key component of the welding process modelling method is the multi-pass welding planning complexity reduction from a three-dimensional into consecutive two-dimensional with dedicated consideration of the non-uniformities of the welding groove. The modelling is applying a mathematical description approach, executed on each reference cross-section. It feeds the multi-pass welding planning module, where the weld beads are planned to be deposited layer by layer and their shapes are also given in mathematical models to keeping their and the groove's curvatures as accurate as possible.

The online system segment of the proposed system design includes simulation synchronization with the welding process and a human-in-the-loop control method with supportive adjustment functions; where the first provides non-observable information to the operator. The reference cross-sections generated during OLP serves as a modification plane that remains constant to ensure the trackability of the modifications during the operation and to provide information to the later refinement of the multi-pass welding plan. Involving the human operator in the loop enables online quality control and process modification to ensure high final quality of the welding. The system design was implemented for a use case of a Francis hydropower turbine runner. The welding experiments showed that it could support the robot operator during the welding process and to handle the non-uniform grooves.

Acknowledgement

The research reported in this paper was partially supported by the Norwegian Research Council through the project 245691 “Cognitive robot welding system (CoRoWeld)” and the Industrial PhD project 244972/O30 “Virtual presence in remote operation of industrial robot”. Further partial support was given by the Higher Education Excellence Program of the Ministry of Human Capacities in the frame of TOPIC research area of Budapest University of Technology and Economics (BME FIKP-MI).

References

- [1] B. M. Abdullah, A. Mason, A. Al-Shamma’a, Defect detection of the weld bead based on electromagnetic sensing, *Journal of Physics: Conference Series*. 450, pp. 12-39, 2013
- [2] M. H. Ang, L. Wei, L. S. Yong, An industrial application of control of dynamic behavior of robots-a walk-through programmed welding robot, in: *International Conference on Robotics and Automation*, 2000
- [3] G. Bolmsjö, M. Olsson, Sensors in robotic arc welding to support small series production, *Industrial Robot: An International Journal*. 32, 2005
- [4] B. Daniel, P. Korondi, G. Sziebig, T. Thomessen, Evaluation of Flexible Graphical User Interface for Intuitive Human Robot Interactions, *Acta Polytechnica Hungarica*. 11, 2014
- [5] D. R. Delapp, Observations of solidification and surface flow on autogenous gas tungsten arc weld pools., Ph.D. dissertation, Vanderbilt Univ., 2005
- [6] H. C. Fang, S. K. Ong, A. Y. C. Nee, Adaptive pass planning and optimization for robotic welding of complex joints, *Advances in Manufacturing*. 5, 2017
- [7] F. Ferraguti, C. T. Landi, C. Secchi, C. Fantuzzi, M. Nolli, M. Pesamosca, Walk-through Programming for Industrial Applications, *Procedia Manufacturing*. 11, pp. 31-38, 2017
- [8] M. Franke, B.-C. Pirvu, D. Lappe, B.-C. Zamfirescu, M. Veigt, K. Klein, K. Hribernik, K.-D. Thoben, M. Loskyll, Interaction Mechanism of Humans in a Cyber-Physical Environment, in: *Dynamics in Logistics*, Springer, 2016
- [9] M. Fridenfalk, Development of intelligent robot systems based on sensor control, Lund Univ., 2003
- [10] T. Gabor, L. Belzner, M. Kiermeier, M. T. Beck, A. Neitz, A Simulation-Based Architecture for Smart Cyber-Physical Systems, in: *International Conference on Autonomic Computing*, IEEE, pp. 374-379, 2016
- [11] P. Galambos, Vibrotactile feedback for haptics and telemanipulation: Survey, concept and experiment, *Acta Polytechnica Hungarica*. 9, 2012

-
- [12] G. Gökmen, Y. Karatepe, T. Ç. Ak, M. Kurtulmu, Spectrum Analysis of GMA Welter in Various Working Modes, *Acta Polytechnica Hungarica*. 9, pp. 5-16, 2012
- [13] W. P. Gu, Z. Y. Xiong, W. Wan, Autonomous seam acquisition and tracking system for multi-pass welding based on vision sensor, *Int J Adv Manuf Technol*. 69, pp. 451-460, 2013
- [14] T. Haidegger, L. Kovács, R.-E. Precup, B. Benyó, Z. Benyó, S. Preitl, Simulation and control for telerobots in space medicine, *Acta Astronautica*. 81, pp. 390-402, 2012
- [15] T. Hatano, C. M. Horvath, T. Thomessen, M. Niitsuma, A vibrotactile navigation aid for remote operation of an industrial robot, in: *International Symposium on System Integration, IEEE/SICE*, pp. 700-705, 2016
- [16] C. M. Horváth, S. Kovács, New cognitive info-communication channels for human-machine interaction, *Recent Innovations in Mechatronics*. 4, 2017
- [17] C. M. Horvath, T. Thomessen, P. Korondi, Robotized Multi-Pass Tungsten Inner Gas Welding of Francis Hydro Power Turbines, in: *IEEE IES, Edinburgh, Scotland, United Kingdom*, pp. 1759-1765, 2017
- [18] E. Horváth, C. Pozna, R.-E. Precup, Robot Coverage Path Planning Based on Iterative Structured Orientation, *Acta Polytechnica Hungarica*. 15, pp. 231-249, 2018
- [19] IFR, Executive Summary World Robotics 2017 Industrial Robots, 2017
- [20] P. Kah, M. Shrestha, E. Hiltunen, J. Martikainen, Robotic arc welding sensors and programming in industrial applications, *Int J Mech Mater Eng*. 10, 2015
- [21] H. C. E. Kjeldsen, Sensor Based Welding Automation Modelling System: Including a Specially Developed Low-cost Temp. Imaging System, 2007
- [22] T. Kosicki, T. Thomessen, Cognitive Human-Machine Interface Applied in Remote Support for Industrial Robot Systems, *International Journal of Advanced Robotic Systems*. 10, pp. 342, 2013
- [23] Y. Liu, W. Zhang, Y. Zhang, Dynamic Neuro-Fuzzy-Based Human Intelligence Modeling and Control in GTAW, *IEEE Transactions on Automation Science and Engineering*. 12, pp. 324-335, 2015
- [24] Y. Liu, Y. Zhang, Toward Welding Robot With Human Knowledge: A Remotely-Controlled Approach, *IEEE Transactions on Automation Science and Engineering*. 12, pp. 769-774, 2015
- [25] O. Madsen, C. Bro, M. B. Madsen, A Software module for planning of robotized multi-pass welding, in: *10th International Conference on Computer Technology in Welding and Manufacturing, Denmark, 2000*

- [26] J.A. Marvel, Collaborative Robots: A Gateway Into Factory Automation, ThomasNet News. 2014
- [27] D. Massa, M. Callegari, C. Cristalli, Manual guidance for industrial robot programming, *Industrial Robot: An International Journal*. 42, 2015
- [28] P. McGuire, J. Fritsch, J. J. Steil, F. Rothling, G. A. Fink, S. Wachsmuth, G. Sagerer, H. Ritter, Multi-modal human-machine communication for instructing robot grasping tasks, in: *International Conference on Intelligent Robots and Systems, IEEE/RSJ*, pp. 1082-1088, Vol. 2, 2002
- [29] L. Monostori, Cyber-physical Production Systems: Roots, Expectations and R&D Challenges, *Procedia CIRP*. 17, pp. 9-13, 2014
- [30] P. Muller, J. Julius, D. Herr, L. Koch, V. Peycheva, McKiernan, Sean, *Ann. Report on European SMEs 2016/2017: Focus on self employment 2017*
- [31] A. Y. Nee, *Handbook of Manufacturing Engineering and Technology* Springer London, 2015
- [32] E. Negri, L. Fumagalli, M. Macchi, A Review of the Roles of Digital Twin in CPS-based Production Systems, *Procedia Manufacturing*. 11, 2017
- [33] A. D. Nicholson, Rapid adaptive programming using image data, Ph.D. dissertation, University of Wollongong, 2005
- [34] Z. Pan, J. Polden, N. Larkin, S. Van Duin, J. Norrish, Recent progress on programming methods for industrial robots, *Robotics and Computer-Integrated Manufacturing*. 28, pp. 87-94, 2012
- [35] S. Pieskä, J. Kaarela, O. Saukko, Towards easier human-robot interaction to help inexperienced operators in SMEs, in: *3rd International Conference on Cognitive Infocommunications*, pp. 333-338, 2012
- [36] J. N. Pires, A. Loureiro, G. Bolmsjö, *Welding Robots: Technology, System Issues and Application* Springer Science & Business Media, 2006
- [37] J. N. Pires, A. Loureiro, T. Godinho, P. Ferreira, B. Fernando, J. Morgado, *Welding robots*, *IEEE Robotics Automation Magazine*. 10, 2003
- [38] T. P. Quinn, C. Smith, C. McCowan, E. Blachowiak, R. Madigan, Arc sensing for defects in constant-voltage gas metal arc welding, *Welding Journal*. 78, pp. 322-328, 1999
- [39] M. Riedl, H. Zipper, M. Meier, C. Diedrich, Cyber-physical systems alter automation architectures, *Annual Reviews in Control*. 38, 2014
- [40] K. Samu, B. Thamó, Internet based light quality measurement, *Recent Innovations in Mechatronics*. 4, p. 5, 2017
- [41] R. D. Schraft, C. Meyer, The Need For An Intuitive Teaching Method For Small And Medium Enterprises, in: *Proceedings of the ISR-Robotik., Munich, Germany*, p. 10, 2006

-
- [42] B. Siciliano, L. Villani, *Robot Force Control* Springer Science & Business Media, 2012
- [43] B. Solvang, G. Sziebig, P. Korondi, *Robot Programming in Machining Operations*, in: *Robot Manipulators*, Intechweb, pp. 479-496, 2008
- [44] S. Sugita, T. Itaya, Y. Takeuchi, Development of robot teaching support devices to automate deburring and finishing works in casting, *Int J Adv Manuf Technol.* 23, pp. 183-189, 2004
- [45] B. Takarics, P. T. Szemes, G. Nemeth, P. Korondi, Welding trajectory reconstruction based on the Intelligent Space concept, in: *2008 Conference on Human System Interactions*, pp. 791-796, 2008
- [46] F. Tao, J. Cheng, Q. Qi, M. Zhang, H. Zhang, F. Sui, Digital twin-driven product design, manufacturing and service with big data, *The International Journal of Advanced Manufacturing Technology.* 94, pp. 3563-3576, 2018
- [47] T. Thomessen, M. Niitsuma, K. Suzuki, T. Hatano, H. Hashimoto, Towards Virtual Presence Based on Multimodal Man-Machine Communication: A Remote Operation Support System for Industrial Robots, *IFAC-PapersOnLine.* 48, pp. 172-177, 2015
- [48] J. Wu, J.S. Smith, J. Lucas, Weld bead placement system for multipass welding, *IEE Proc.-Science, Measurement and Technology.* 143, 1996
- [49] Y. Xu, N. Lv, G. Fang, S. Du, W. Zhao, Z. Ye, S. Chen, Welding seam tracking in robotic gas metal arc welding, *J of Mat. Proc. Tech.* 248, 2017
- [50] Y. Xu, H. Yu, J. Zhong, T. Lin, S. Chen, Real-time seam tracking control technology during welding robot GTAW process based on passive vision sensor, *Journal of Materials Processing Technology.* 212, 2012
- [51] S. Yan, H. Fang, S. Ong, A. Nee, Optimal pass planning for robotic welding of large-dimension joints with nonuniform grooves, *Proceedings of the Institution of Mechanical Engineers, Part B: J. of Eng. Manufacture.* 2017
- [52] S. J. Yan, S .K. Ong, A. Y. C. Nee, Optimal Pass Planning for Robotic Welding of Large-dimension Joints with Deep Grooves, *Procedia CIRP.* 56, pp. 188-192, 2016
- [53] H. Zhang, H. Lu, C. Cai, S. Chen, Robot Path Planning in Multi-pass Weaving Welding for Thick Plates, in: T.-J. Tarn, S.-B. Chen, G. Fang (Eds.), *Robotic Welding, Intelligence and Automation*, Springer Berlin Heidelberg, pp. 351-359, 2011
- [54] W. J. Zhang, Y. M. Zhang, Dynamic Control of the GTAW Process Using a Human Welder Response Model, *Welding Journal.* 92, pp. 154-166, 2013
- [55] Z. Zhang, H. Chen, Y. Xu, J. Zhong, N. Lv, S. Chen, Multisensor-based real-time quality monitoring by means of feature extraction, selection and

modeling for Al alloy in arc welding, Mechanical Systems and Signal Processing. 60, pp. 151-165, 2015

List of Symbols

\mathbf{a}	path normal vector at the task point	$(\mathbf{s} = \mathbf{a} \times \mathbf{n})$
\mathbf{b}	robot's base coordinate system	\mathbf{t} local coordinate system at the task point
$\underline{\mathbf{C}}$	task point coordinates	\mathbf{t}' modified task point's coordinate system, only user translation, relative to \mathbf{t}
$\Delta \mathbf{c}^r$	transformation for collision avoidance, relative to \mathbf{r}	\mathbf{t}'' modified task point's coordinate system, only user rotation, relative to \mathbf{t}'
$\Delta \mathbf{c}_i^{t''}$	transformation for collision avoidance, relative to \mathbf{t}''	$\{\mathbf{T}_C\}^r$ task path description in \mathbf{r}
$\widetilde{\Delta \mathbf{c}_i^{t''}}$	range limited transformation for collision avoidance, relative to \mathbf{t}''	$\{\mathbf{T}_C\}^b$ task path description in \mathbf{b}
\mathbf{n}	Path tangent vector at the task point	$\{\mathbf{T}'_C\}^b$ task path description in \mathbf{b} after user translation modification
$\Delta \mathbf{p}^t$	Translation modification at \mathbf{t}	$\{\mathbf{T}''_C\}^b$ task path description in \mathbf{b} after user rotation modification
\mathbf{r}	CAD/CAM model coordinate system	$\{\mathbf{T}'''_C\}^b$ Final task path description in \mathbf{b} , including all modification combined
$\Delta \mathbf{r}^{t'}$	rotation transformation applied on \mathbf{t}'	\mathbf{T}_r^b transformation matrix from \mathbf{r} to \mathbf{b}
$\Delta \mathbf{R}_{x,\varphi}$	user rotation around the x -axis of the task point	φ welding torch working angle,
$\Delta \mathbf{R}_{y,\theta}$	user rotation around the y -axis of the task point	θ welding torch travel angle
$\Delta \mathbf{R}_{z,\psi}$	user rotation around the z -axis of the task point	ψ rotation angle around the electrode main axis
\mathbf{s}	third vector at the task point	

Study of the Drive Characteristics Affecting the Power Loss of V-Belt Drives

Péter Gárdonyi, István Szabó, László Székely and László Kátai

Faculty of Mechanical Engineering, Szent István University
Páter Károly út 1, 2100 Gödöllő, Hungary
gardonyi.peter@gek.szie.hu, szabo.istvan@gek.szie.hu,
szekely.laszlo@gek.szie.hu, katali.laszlo@gek.szie.hu

Abstract: The V-belt drive heats up during power transmission, i.e. a significant part of power loss turns into heat and is transferred to the environment. In this paper, the temperature rise in the V-belt was studied as loss intensity as a function of the drives parameters. It was justified in the scope of the major characteristics affecting power loss that by ideally selecting the parameters of the V-belt drive power loss can be measurably reduced. Based on earlier results as well, a regression model was used to examine the temperature rise in the V-belt. With the help of analysis of variance (ANOVA), the magnitude of the effect of each drive parameter was determined on the warming of the V-belt. According to our results, the most relevant drive parameter is the pulley diameter.

Keywords: ANOVA; Efficiency; Power loss; Temperature; V -belt

1 Introduction

Belt drives are widely used power transmission solutions, both in the field of industry and agriculture [1, 2]. Various types of belts are used, e.g. flat belts, V-belts, poly-V-belts. Flat belts have small bending rigidity, V-belts have large power transmission capability and poly-V belts combine these properties [3]. In regular industrial and agricultural use V-belts are the most wide-spread belt drive solutions.

The V-belt drive, like any machine structure, operates with certain efficiency, which is the ratio of useful and input power. The difference between these gives power loss, most of which turns into heat.

During the power transmission of belt drives, the heat build-up of the belt is the result of basically two impacts. The heat build-up due to the macroscopic friction of the contacting surfaces, and the proportion of hysteresis loss occurring because of the repeated bending of the belt (internal friction (slip among the molecules)) which turns into heat, Gerbert [4].

Childs and Cowburn [5] performed a series of measurements on the power loss of flat and V-belts associated with a very small pulley. They found that the reduced efficiency of belts that do not match their pulley groove angles may be due to greater radial compliance of these belts. In the companion paper [6], the experimental study quantified the effects of pulley radius, belt tension and belt deformation properties on speed and torque losses during power transmission.

Gerbert [7, 8, 9] from 1972 studied the mechanical behavior of V-belt drives in details. A unified slip theory [9] was proposed for the V-belt drives considering four factors that may affect belt slip, i.e. elastic creep along the belt, compliance in the radial direction, shear deflection that varies both radially and axially, and flexural rigidity during engagement/disengagement of the belt. This paper summarizes the reasons for the speed loss of rubber V-belt drives. On the other hand, Gervas and Pronin [10, 11] proposed that the torque loss of rubber V-belt drives resulted mainly from: (1) hysteresis losses in the materials when the belt was bent on and off the pulleys, and the belt was compressed into the groove of the pulleys; (2) the radial sliding losses as the belt is continuously wedged into and out of the pulleys. These loss components were studied by several researchers e.g. in case of CVT and flat belt drives [12, 13, 14].

Several researchers deal with V-ribbed drives. Song *et al.* [15] performed the thermal–mechanical finite element analysis of a two pulley V-ribbed belt drive system. In addition to material nonlinearity, large deformation, and frictional contact, the analysis took into account the thermal degradations and thermal expansions of belt rubber compounds. The temperature effects on stresses, strains, and belt-pulley contact slip rates were studied in detail. Yu *et al.* [16] presented a comprehensive study of the belt moving in pulley grooves. A three-dimensional finite element model represented a complete half rib operating between two pulleys.

Previous experimental studies on power loss behavior of belt drives usually considered V-belt and continuously variable transmission (CVT) belt drives. Childs and Cowburn [17] experimentally investigated the effects of mismatch between the wedge angles of pulley grooves and belt ribs on the power loss behavior of V-belt drives. During the tests, they kept the other parameters constant. They [5] also studied the effects of small pulleys on the power loss both theoretically and experimentally.

Lubarda [18] analytically formulated the variation in the belt force over the arc of contact of flat and V-belts before gross slip occurs. He separated the arc of contact into active and non-active regions, similar to the approach of Gerbert [6] and Johnson [19].

The V-belt drive, like any machine structure, operates with certain efficiency, which is the ratio of useful and input power. The difference between these gives power loss, most of which turns into heat. If the stabilized temperature of the V-belt is examined as loss intensity, the efficiency of the belt drive can be concluded

from it. Higher belt temperatures lead to the degradation of molecular chains, to the aging of the rubber, which significantly influences belt life.

In our earlier measurements we have investigated the heating generated inside V-belts due to repeated bending and engagement/disengagement and to define the damping coefficient during idle running, causing the increase in temperature [1]. The effect of the test parameters (pulley diameter, bending frequency, pre-tensioning) was studied separately. On the basis of our previous studies it can be concluded that bending frequency and belt temperature are linearly related [20].

Further tests were also carried out to find out to what extent drive set-up faults contribute to heat build-up in the V-belt and thus how they affect the efficiency of the drive [20].

In the present study, the effects of the belt-drive parameters on the temperature increase of V-belt are experimentally investigated. Although the temperature in the belt is affected by the contact between the belt and the pulley, the aim of our research is to determine the effect of the most important operating parameters in some fixed conditions (e.g. laboratory temperature, humidity, $i=1$ speed ratio, steel pulley, etc.).

For this purpose, test equipment was constructed with two equal-sized pulleys to measure the temperature of the V-belt. In comparison to the previous studies, a much larger number of parameters and their interactions are taken into consideration; they include belt tension, belt bending frequency, braking torque, pulley diameter. The effects and the interaction of these drive parameters to the belt temperature are investigated in the first time.

After the measurement and data collection process we constructed a nonlinear multiple regression model to see the relation between the response parameter and the independent parameters. The model strongly relies on the one-factor-at-a-time measurement results achieved in our previous papers [1, 20, 21].

2 Method

The objective of the test is to experimentally determine the drive characteristics affecting V-belt drive loss, and their relationship. The majority of the losses turn into heat, so the temperature rise in the V-belt was chosen as the test parameter, which means the power loss between the two steady states - between the workshop and the stabilized state of the operating temperature.

2.1 Experimental Plan

Table 1 summarizes the factors of the experimental settings on the basis of the effects determining heat buildup and the drive parameters involved in them. The

V-belt heats up due to the macroscopic friction of the contacting surfaces and the hysteresis loss occurring as a result of the repeated use of the belt, i.e. power loss turns into heat. These two fundamental effects are affected by the relationship of V-belt slippage and bending. The drive parameters of V-belt drive are not independent of each other, certain parameters can be combined, so our experiment is reduced to four independent variables.

Table 1
Defining the experimental factors

Nature of heat buildup	V-belt strain	Drive characteristics		Experiment variables (Factors)
friction of contacting surfaces	relative movement of V-belt in the pulley groove (flexible slippage, wedging effect)	pre-tensioning F_H		pre-tensioning F_H
		braking torque M		braking torque M
		grip angle β		pulley diameter d
hysteresis loss due to belt bending, inner friction	amplitude of bending stress in a load-unload cycle of the V-belt	belt profile		
		pulley diameter d		
		material properties of V-belt		
	frequency of V-belt use	belt speed v_b	revolution n	belt bending frequency f
pulley diameter d				
V-belt length L_d				

The test was conducted with a Z/10 profile V-belt, 5 pulley diameters ($d_l = 60, 90, 118, 150, 180$ mm), 1:1 ratio. With the fixed value of the free belt section length (345 ± 10 mm) the same flexibility of the different drive settings was ensured, which meant a different V-belt length per pulley diameter. With the revolution of the drive shaft in line with the interval generally occurring in practice 10 and 20 s⁻¹ belt bending frequency values were set. The V-belt drive was loaded with the largest transferable torque calculated on the basis of the given pulley diameter and bending frequency, and it was studied without power transmission. With the different experimental settings the value of pre-tensioning was determined based on the drive sizing, it was set at 50, 100, 150% of that. In practice pre-tensioning is the most difficult to set precisely from the drive belt parameters, therefore the setting of a wide range is justified. Table 2 shows the levels of experimental factors and factor space.

The V-belt drive was studied under normal operating conditions, within the limits of flexible slippage, for this the value of belt slippage was continuously measured. Overload occurred in the case of decreased pre-tensioning, with two experimental settings, they were ignored during the evaluation. In the case of settings without load the effect of pre-tensioning force could.

Table 2
Values of experimental factors and factor levels

Factor	Unit	Factor levels	High	Low
Pulley diameter d	Mm	5	60	180
Bending frequency f	s^{-1}	2	10	20
Braking torque M	N	2	0	M_N
Pretension F_H	N	3	$0.5 \cdot F_{HN}$	$1.5 \cdot F_{HN}$

2.2 Measurement of Temperature Rise in V-Belt

The temperature rise in the V-belt was determined from the function describing temperature change (Baule-Mitscherlich), in which the measured parameters change along a decreasing gradient towards the maximum of saturation. With the help of the mathematical model it was possible to calculate the steady-state operating temperatures, so they did not have to be reached during the experiments. Thereby the duration of the measurement was uniformly determined with each experimental setting regardless of whether the belt temperature was steady under the given circumstances or not.

The general equation of the function of saturation is:

$$Y = A \cdot (1 - e^{z+c \cdot X}) \quad (1)$$

Among the function parameters A indicates the stabilized temperature of the V-belt, c indicates the speed of heat buildup and z indicates the temperature of the belt at the beginning of the measurement.

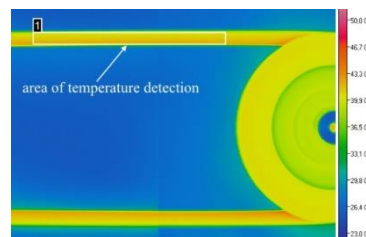


Figure 1
Thermal camera image and its evaluation

The temperature was measured with an infra camera type NEC H2640. During our experiments the side surfaces of the V-belt in contact with the groove were tested, which contained more information about the operation of the drive. Temperature data were obtained from the heat camera images taken of the active surface of the V-belt after image processing (Fig. 1).

Table 3
Design matrix

	Runs	1	2	3	4	5	6	7	8	9	10	11	12	13
Factor 1	d_1 [mm]	180	180	180	180	180	180	150	150	150	150	150	150	118
Factor 2	f [s^{-1}]	20	20	20	10	10	10	20	20	20	10	10	10	20
Factor 3	M [Nm]	1	1	1	1	1	1	1	1	1	1	1	1	1
Factor 4	F_H [N]	1	1.5	0.5	1	1.5	0.5	1	1.5	0.5	1	1.5	0.5	1
Response	ΔT [$^{\circ}C$]	6.2	6.4	6.2	4.1	4.6	4.1	8.5	8.7	8.7	6.1	6.1	6.1	16.2

Table 3 (continued)
Design matrix

14	15	16	17	18	19	20	21	22	23	24	25	26	27	28
118	118	118	118	118	90	90	90	90	90	90	60	60	60	60
20	20	10	10	10	20	20	20	10	10	10	20	20	20	10
1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
1.5	0.5	1	1.5	0.5	1	1.5	0.5	1	1.5	0.5	1	1.5	0.5	1
15.8	17.2	10.2	11.2	11.1	23.9	22.9	24.6	16.3	16.4	18.3	38.1	38.6	40.2*	28.2

Table 3 (continued)
Design matrix

29	30	31	32	33	34	35	36	37	38	39	40
60	60	180	180	150	150	118	118	90	90	60	60
10	10	20	10	20	10	20	10	20	10	20	10
1	1	0	0	0	0	0	0	0	0	0	0
1.5	0.5	1	1	1	1	1	1	1	1	1	1
28.5	32.6*	5.3	4.2	7.3	5.4	9.7	6.7	16.5	10.9	32.9	25.7

* *Outlier*

2.3 The Test Equipment and its Tools

The experiments were conducted on the universal test bench developed at the Department of Machine Design, Szent István University. Due to its design, it is suitable for testing a wide range of mechanical drives, clutches and rotating elements. The grooved table of the test bench offers numerous possibilities of placing the driving and driven units. For the testing of the belt drives, the drive unit was fixed to a tensioning mechanism guided by a linear bearing. The pre-tensioning of the belt can be set with an adjusting spindle and a load cell connected in series with it, whose line of action coincides with that of the shaft pulling force (F_H). In this way the pre-tensioning force can be measured directly. The structure of the universal test bench is shown in Figures 2-3. During the measurements it is possible to fix all the drive parameters through a measuring-data collector and to precisely define them with the help of a programmable logic controller (PLC).

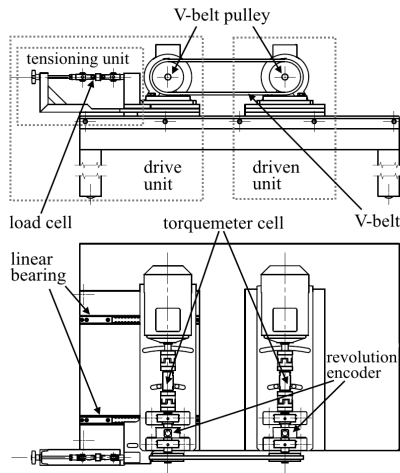


Figure 2
Schematic of the experimental setup

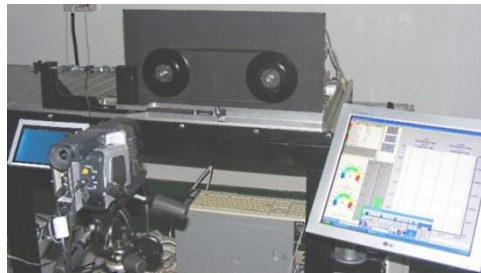


Figure 3
The experimental setup

3 Results and Discussion

3.1 Nonlinear Regression Model for Temperature Change

In earlier studies the individual effect of pulley diameter [17] and frequency [18] of belt bending on temperature change of the belt has already been investigated with one-factor-at-a-time (OFAT) method. The following experimental relation of the temperature and the diameter was obtained in [17]:

$$\Delta T = a_0 + \frac{a_1}{d} \quad (2)$$

where a_0 and a_1 are constants.

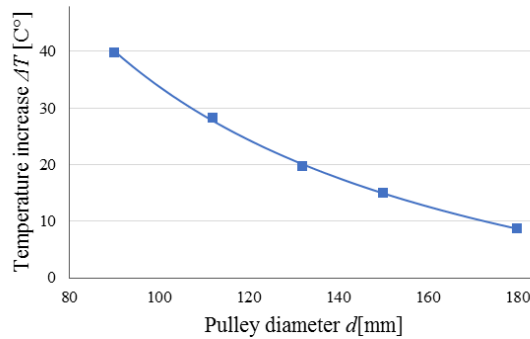


Figure 4

Temperature increase of the belt sides as a function of the pulley diameter
(profile: SPA; $d_i = 90, 112, 132, 150, 180 \text{ mm}$; $i = 1$; $L_d = 1207$; $f = 20 \text{ s}^{-1}$; $M_l = 0 \text{ Nm}$) [17]

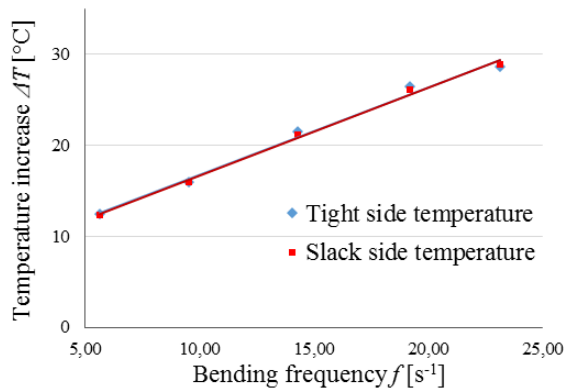


Figure 5

Temperature increase of the belt sides as a function of the bending frequency
(profile: SPA; $d_i = 112 \text{ mm}$; $i = 1$; $L_d = 1207$; $f = 5,6 - 23,1 \text{ s}^{-1}$; $M_l = 0 \text{ Nm}$) [17]

Furthermore, it turned out [17] that ΔT is approximately a linear function of f . In order to see the relationship between the pre-tensioning and the change in temperature, similarly to previous experiments, the values of the other independent variables were fixed and only the extent of the studied parameter was changed (see Figs. 6-7). It shows that ΔT is also a linear function of F_H (R^2 varies between 0.9964 - 0.9991). Based on these results we build a nonlinear regression model to investigate the qualitative relationship between the control parameters and the response parameter. In our models the dependence of ΔT on d , f and M is considered in the previously discussed characteristic way and for simplicity we assume a linear dependence of ΔT on M . Our model is

$$\Delta T = a_0 + \frac{a_1}{d} + a_2 f + a_3 M + a_4 F_H \quad (3)$$

where a_0, \dots, a_4 are constants.

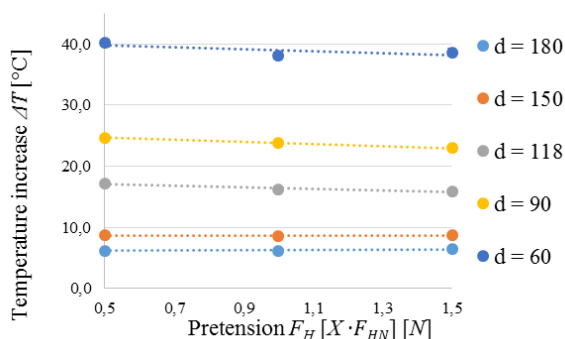


Figure 6

Temperature increase of the belt sides as a function of the pretension on frequency level $f = 20 \text{ s}^{-1}$
 ((profile: Z10; $d_1 = 60, 90, 118, 150, 180 \text{ mm}$; $i = 1$; $a = 345 \pm 10 \text{ mm}$; $M = M_N$)

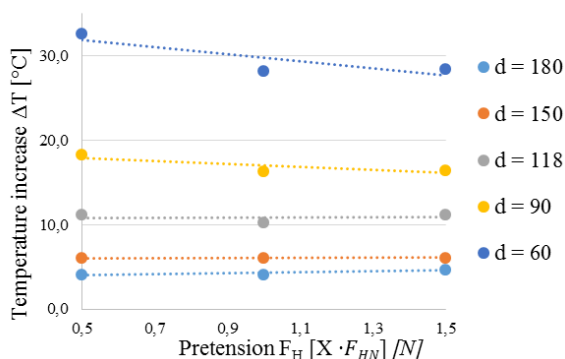


Figure 7

Temperature increase of the belt sides as a function of the pretension on frequency level $f = 10 \text{ s}^{-1}$
 ((profile: Z10; $d_1 = 60, 90, 118, 150, 180 \text{ mm}$; $i = 1$; $a = 345 \pm 10 \text{ mm}$; $M = M_N$)

3.2 Analysis of Variance (ANOVA)

For a regression model analysis of variance (ANOVA) is the tool to check its relevance. If the resulting significance level is high ($p > 0.05$) then the model is not applicable for the quantitative description of the phenomena. Furthermore, using ANOVA, the total variance of the response variable is decomposed into the sum of variance explained by the control parameters plus the residual, which is the difference between the measured and the predicted value. A T-test is used to evaluate the significance of each coefficient in the regression equation. If the significance level of the T-test for a coefficient is high ($p > 0.05$) it means that the coefficient of the corresponding control variable is statistically 0. In Table 4 we summarize the results of ANOVA.

Table 4
ANOVA for temperature change

Source	Sum of squares	Degree of freedom	Mean square	F/T value	P
Model	3574.22	37	1158.16	394.75	<0.001
<i>d</i>	3152.84	4	788.21	29.824	<0.001
<i>f</i>	211.43	1	211.43	8.489	<0.001
<i>M</i>	129.39	2	64.67	6.228	<0.001
<i>F_H</i>	2.43	2	1.22	-0.423	0.675
Residual	99.75	34	3.99		

Table 5
Coefficient of the parameters in the regression model

Model	coefficient	Beta coefficient	t	p
constant	-17.100		-14.973	<0.001
<i>1/d</i>	2317.476	0.888	29.824	<0.001
<i>f</i>	0.472	0.243	8.489	<0.001
<i>M</i>	4.430	0.185	6.228	<0.001

From the table one can see that the F-value of the model is significant, that is our model is relevant. Except the coefficient of pre-tension the coefficients of all other control variables significantly differ from 0. With the obtained values of the coefficients our model is of the form

$$\Delta T = -17.1 + \frac{2317.476}{d} + 0.472 f + 4.43 M \quad (4)$$

The value for goodness of fit of the model is $R^2 = 0.970$. According to the variance components of the control variables the reciprocal of the diameter affects the temperature change in the largest extent, the effect of frequency and pretension is almost the same and a magnitude smaller than of the diameter.

3.3 Verification of the Model

Using the same experimental setup with arbitrarily chosen control parameter values we performed additional tests to verify our model. These experiments reveal that the difference between the predicted values for the temperature change and the experimental values are smaller than 5% and is our model is reasonably accurate (see Table 6).

Table 6
Model verification control parameters

	Pulley diameter d [mm]	Bending frequency f [s^{-1}]	Braking torque M [Nm]	Pretension F_H [N]	Measured value ΔT [$^{\circ}C$]	Calculated value ΔT [$^{\circ}C$]	Relative difference [%]
1	118	20	$0,5 M_N$	$1 F_{HN}$	13.7	14.19	3.6
2	180	20	$1/3 M_N$	$1 F_{HN}$	6.4	6.69	4.6
3	118	15	$0 M_N$	$1 F_{HN}$	9.5	9.62	1.3
4	150	20	$1/3 M_N$	$1 F_{HN}$	8.8	9.27	4.9
5	60	15	$0 M_N$	$1 F_{HN}$	28.5	28.60	0.2

Conclusions

In this paper, the drive parameters affecting the power losses of V-belt drives were studied. The majority of the losses turn into heat, so the temperature rise in the V-belt was chosen as the test parameter, which means the power loss between the two steady states - between the workshop and the stabilized state of the operating temperature.

It can be seen from the experiments conducted, that the selection of drive parameters plays an important role in the temperature increase in V-belts. The loss of V-belt drive is directly proportional to bending frequency, pre-tensioning and load, whereas, inversely proportional to the diameter of the pulley. With the help of the analysis of variance, it was concluded that from the drive parameters, the diameter of the pulley affects the temperature buildup in the V-belt to the largest extent; the effects of the bending frequency, pre-tensioning and load are nearly identical, but at the same time, this effect is more than an order of magnitude less than that of the pulley diameter.

A significant part of the temperature rise is due to internal friction caused by bending, due to the viscoelastic behavior of the V-belt from its material properties. The bending stress is inversely proportional to the diameter of the pulley, this relationship is also justified by the model of temperature rise from bending.

When designing a highly efficient belt drive it is necessary to select the above mentioned characteristics optimally.

Nomenclature

F_H : Pre-tensioning of the V-belt drive

M : Braking torque

b : Arc of contact

d : Pulley diameter

f : V-belt bending frequency

- L_d : V-belt length
 ΔT : V-belt temperature increase
 i : speed ratio
 a : drive center distance
 p : significance level

References

- [1] L. Kátai and I. Szabó, Identification of V-belt power losses with temperature measurement, *Journal of Mechanical Science and Technology*, 29 (8) (2015) 3195-3203
- [2] F. H. Schafer, *Antriebsriemen*. Arntz Optibelt Gruppe Höxter, 2007, ISBN 978-3-00-0217113-5
- [3] A. Grinčová, M. Andrejiová, D. Marasová: Analysis of Conveyor Belt Impact Resistance Data Using a Software Application, *Acta Polytechnica Hungarica*, Vol. 14, No. 2 (2017) 113-130
- [4] B. G. Gerbert, Power loss and optimum tensioning of V-belt drives, *Journal of Engineering for Industry*, 96 (1974) 877-885
- [5] T. H. Childs, D. Cowburn, Power Transmission Losses in V-Belt Drives Part 2: Effects of Small Pulley Radii, *Proceedings of the Institution of Mechanical Engineers, Part D: Journal of Automobile Engineering*, 201 (1) (1987) 41-53
- [6] B. G. Gerbert, *On flat belt slip*, *Vehicle Tribology, Tribology Series*, 16 (1991) 333-339
- [7] B. G. Gerbert, Force and slip behavior in V-belt drives, *Acta Polytechnica Scandinavica, Mechanical engineering series*, 67 (1972)
- [8] B. G. Gerbert, Pressure Distribution and Belt Deformation in V-Belt Drives, *Journal of Engineering for Industry*, 97 (1975) 976-982
- [9] B. G. Gerbert, G. Belt slip — a unified approach, *Journal of Mechanical Design*, 118 (3) (1996) 432-438
- [10] K. J. Gervas and B. A. Pronin, Calculation of Power Losses in Belt Drives, *Russian Engineering Journal*, 47 (3) (1967) 26-29
- [11] K. J. Gervas, Determining the Power Losses in V-belt Drives During Flexure, *Soviet Rubber Technology*, 28 (2) (1969) 42
- [12] L. Bertini, L. Carmignani, F. Frendo, Analytical model for the power losses in rubber V-belt continuously variable transmission (CVT), *Mechanism and Machine Theory*, 78 (2014) pp. 289-306

-
- [13] L. D. Pietra, F. Timpone, Tension in a flat belt transmission: Experimental investigation, *Mechanism and Machine Theory*, 70 (2013) pp. 129-156
- [14] C. Zhu, H. Kiu, J. Tian, Q. Xiao, X. Du, Experimental investigation on the efficiency of the pulley-drive CVT, *International Journal of Automotive Technology*, 11 (2010) pp. 257-261
- [15] G. Song, K. Chandrashekhara, W. F. Breig, D. L. Klein and L. R. Oliver, Analysis of cord reinforced poly-rib serpentine drive with thermal effect, *Journal of Mechanical Design*, 127 (2005) 1198-1206
- [16] D. Yu, T. Childs and K. Dalgarno, Experimental and finite element studies of the running of V-ribbed belts in pulley grooves, *Proceedings of the Institution of Mechanical Engineers C, Journal of Mechanical Engineering Science*, 212 (1998) 343-354
- [17] T. H. C. Childs and D. Cowburn, Power transmission losses in V-belt drives, Part 1: Mismatched belt and pulley groove wedge angle effects, *Proceedings of the Institution of Mechanical Engineers D, Transport Engineering*, 201 (1987) 33-40
- [18] V. A. Lubarda, Determination of the belt force before the gross slip, *Mechanism and Machine Theory*, 83 (2015) 31-37
- [19] K. L. Johnson, *Contact Mechanics*, first ed. Cambridge University Press, London, UK (1987)
- [20] P. Gárdonyi, L. Kátai and I. Szabó, Examination of drive misalignment and v-belt temperature conditions, *International Journal of Science, Technics and Innovations for the Industry*, 12 (2015) 56-59
- [21] P. Gárdonyi, L. Kátai and I. Szabó, Relationship between the pulley diameter and V-belt temperature conditions, *Scientific Conference of Young Engineers, XX*, Cluj-Napoca, Romania (2015) 151-154

DQNET: Assessment of Quality Regulation System as Complex Information Network

Tamás Csiszér

Óbuda University, Bécsi út 96/b, 1034 Budapest, Hungary
csiszer.tamas@rkk.uni-obuda.hu

Abstract: This article introduces a set of indicators and their interpretation called DQNET for the assessment of information structures in documents of quality regulatory systems. This complex system is considered a network with a piece of information in documents nodes; and links between them arcs. Like in citation network of scientific publications there are several network indicators in such information networks, which can reflect the 'positions' and 'roles' of elements in this system. By in- and out-degrees and other matrices documentations can be identified with, e.g. 'high importance' or with 'high sensitivity', requiring different ways of handling. By the indicators of structure functional suitability of regulation can be analyzed and predicted too.

Keywords: Documentation analyses; importance and sensitivity; networkscience

1 Introduction

Grouping or clustering documentation by calculating the similarity or the distance of documents or of their parts as the entities of regulation systems are one of the most important fields of the research of complex information networks.

Many scientists proposed indicators of document similarities, focusing on different elements of documents like words and phrases. One of these researchers, Wang proposes a method to represent a document as a typed Heterogeneous Information Network (HIN), where the entities and relations are annotated with types [5]. He and his colleagues underline that most of researches in the field of documentation networks are focusing on similarities between documents and do not put enough efforts on links sourced in heterophily, i.e. the difference between documents [7]. Yang proposes hierarchical attention network for classifying documents according to its hierarchy and the importance of content (word, sentence and document vectors) [6]. Tan presents the latent quality model (LQM). LQM associates each document with a latent quality score, which provides a measure of the impact or popularity of a document [3]. Wan proposes Cluster-based Conditional Markov Random Walk Model (ClusterCMRW) and the

Cluster-based HITS Model (ClusterHITS) to find parts of different documentations related to the same content to summarize information [4]. Cao developed a Ranking framework upon Recursive Neural Networks to rank sentences for multi-document summarization [1]. Carley applies Dynamic Network Analysis (DNA) approach to create and analyze multi-mode and multi-link networks [2].

All of these approaches were involved into the development process of DQNET. Documentation systems consist of many elements such as manuals, descriptions of procedures and products, forms, templates and others, published on paper or in electronic format. There is a huge number of links among their parts indicating the connections of regulations. One can find regulation holes and redundancy too. Due to this complexity, these systems are difficult to create, maintain, assess, upgrade and improve, so these activities should be supported by analytical and development methods based on qualitative and quantitative measurements. The purpose of these methods is to give evidence of proper or improper structure of documentation or – in general – information systems. In the following chapters, we introduce a set of indicators of DQNET that can represent the internal and external properties of the elements of such kind of systems.

2 Theory

2.1 Structure

A network-like representation of a documentation system can be seen in Figure 1. General nodes (black dots) represent the elements of the system. Links between documents are represented by grey arrows, indicating the direction of links too. There can be seen some special types or groups of nodes as well, represented by colored dots and circles as follows:

- Blue dots: reversed regulation. It may show the problem of regulating something with a link to another document, which has a link the other way around. It may be useful if these links belong to the same parts of documents and indicate the two directions of the same connection, or if these links belong to different parts of documents and indicate different connections, but it may indicate that these documents are linked to each other with a regulation hole. In these cases, the higher the number of links between two documents is, the stronger the connection of them can be detected.
- Purple dots: regulation loop. It may show the similar property of regulation described above, with involving more than two documents.

Such kinds of grouping of documents should be assessed according to the same approach.

- Orange dots: regulation chain. It shows how a rule defined with a set of documents with one-direction links. Obviously, it should be considered as chain if their links belong to the same parts of documentations.
- Red circles: hubs. These documents may have important roles in the system due to their links to other elements.
- Green circles: regulation trees. One document links to more than one. Changes of this document can have huge influence on others or vice-versa.
- Blue circles: regulation islands. These have noconnection to the other parts of the regulation network.

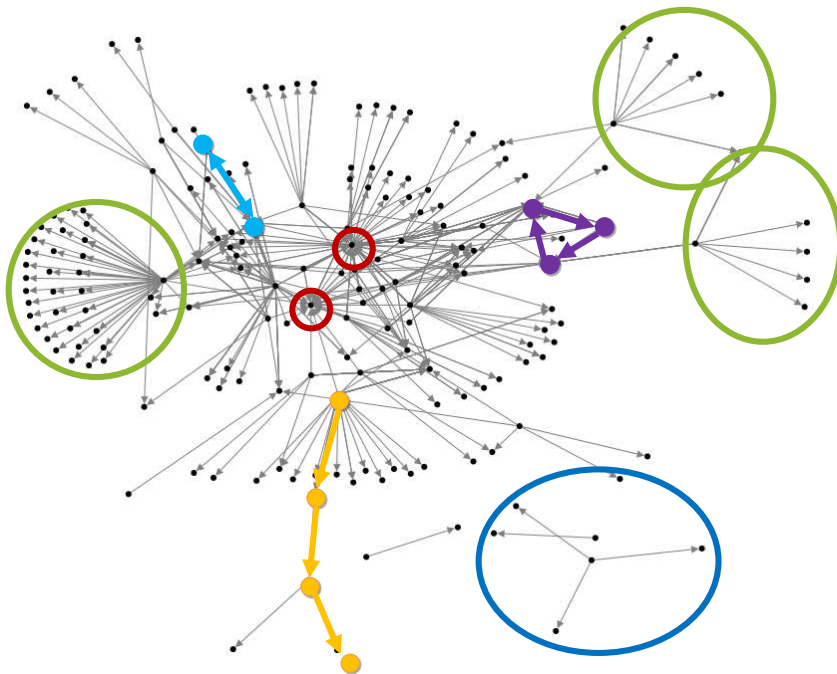


Figure 1

An example for the network-like implementation of documentation systems

Networks can be described with the number and structure of these special types and groups of nodes.

2.2 Internal Properties

Internal properties of documents determine how easy it is to understand, memorize and apply regulations described in documentation system. Most important related indicators:

- Size related indicators: number of pages, words, sentences and lines.
- Sentence structure related indicators: length of sentence and words, rate of number of words and sentences, rate of number of commas and sentences, rate of long sentences.
- Text structure related indicators: rate of number of sentences and paragraphs, rate of number of paragraphs and pages.
- Document structure related indicators: number of appendices and chapters.

With these indicators documents can be qualified from different perspectives as follows:

- Understandability: how easy it is to understand regulations.
- Notability: how easy it is to memorize regulations.
- Accountability: how easy it is to identify responsibilities.
- Searchability: how easy it is to find user or case relevant information.
- Applicability: how easy it is to apply the regulations during operation.

2.3 External Properties

External properties can be described by well-known network indicators as follows:

- In-degree - Importance: number of incoming links of a nod. The higher the in-degree of a nod is, the more important the document represented by the nod is.
- Out-degree – Sensitivity: number of outgoing links of a nod. The higher the out-degree of a nod is, the more sensitive the document represented by the nod is.
- Degree distribution – Evenness: how links are distributed to nodes. It shows how evenly nodes are connected to each other. Some questions that can be answered by this indicator: 1) Can a chain of links among all documents be found?; 2) Are there isolated elements or groups in network?; 3) Are there big differences among the degrees of nodes?

- **Betweenness – Criticality:** in what part of the shortest paths between any pair of nodes the associated node takes part. The higher the betweenness is, the more critical role the document has in the network.
- **Closeness – Simplicity:** how close nodes are to each other. The shortest the average distance (number of links on the path) among the nodes is, the simpler the network is. It may help us to make the system of connections simpler.
- **Clustering coefficient – Looping:** rate of realized and possible numbers of triangles of nodes. It shows how many connected circles of 3 documents have been created.
- **Reciprocated vertex pair ratio – Reciprocity:** rate of two-directional to one-directional links of nodes.

Knowing the values of network indicators, documentation network can be qualified. Some examples of qualification:

- **Clarity:** documents are connected to each other precisely; sender and receiver documents of the links can be identified exactly.
- **Relevance:** links connect the proper parts of proper documents.
- **Redundancy:** two-directional links between two documents are not redundant, i.e. indicate two different connections.
- **Contradiction:** rules defined in connected documents are consistent.
- **Completeness:** regulation hole cannot be found.

2.4 Network-based Optimization of Documentation System

There are several ways to optimize a documentation system. It depends on the goals, organization structure and culture, skills of users, technical environment, level of automatization, etc. Due to this complexity there is no single ideal solution, but some important features can be defined based on the properties described above.

One of the fundamental goals of creating documents is to define regulation for operation, which must be easy to find, understand, memorize and apply. It can be ensured if rules for conducting a particular activity are handled as individual information package represented by only one node in the regulation network. This information package consists of short sentences and graphical elements. Two or more nodes are connected by links if activities represented by these nodes 1) form a predecessor-successor pair of process steps, 2) are allocated to the same equipment, 3) need the same human skills to be done, etc. Nodes and different types of links among them form different networks of operation rules. Different subgraphs of these networks belong to processes, products, resources, organization

groups and localizations. According to the grouping principle, different types of documentation (e.g. process manuals, product descriptions, etc.) can be created too. An example of this rule-based network can be seen in Figure 2.

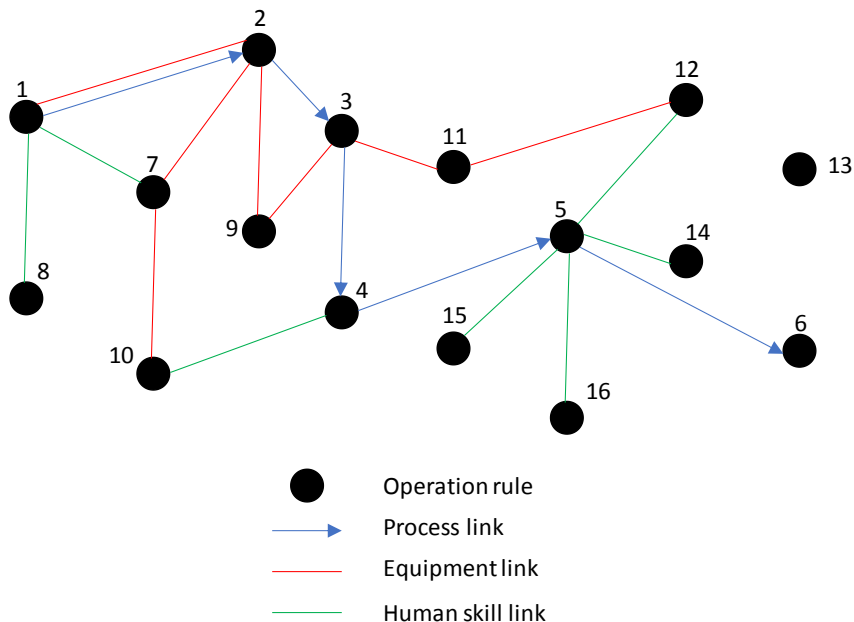


Figure 2

Example for rule-based network (part of whole network)

3 Case Study

An international financial organization has a complex system of documentations. Due to the order of Central Regulatory Office process documentations have to be modified to meet new requirements. The management decided to analyze the documentation structure with DQNET network indicators. The associated graphs and calculations are generated by NodeXL application.

3.1 Overall Metrics

The whole system can be seen in Figure 3. The Overall metrics are presented in Table 1.

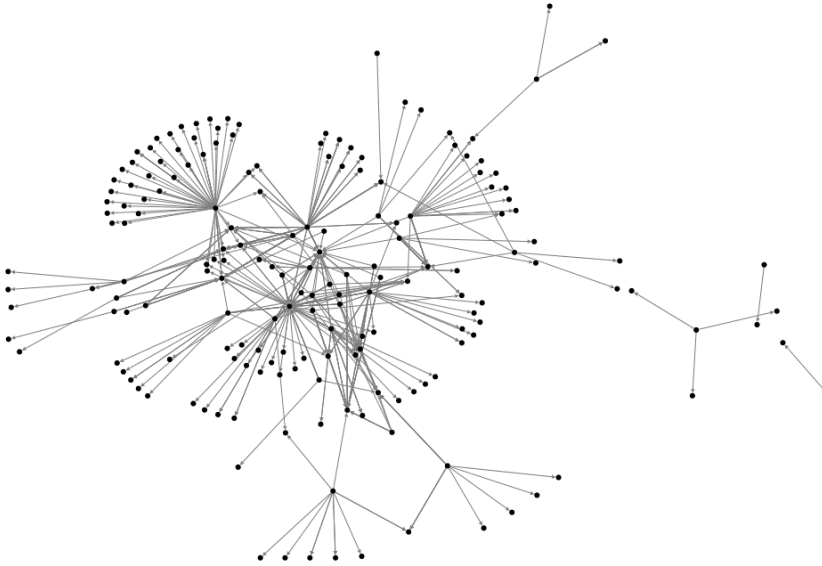


Figure 3

The whole structure of documentation system

Table 1

Overall metrics of documentation system

Metrics	Value
No. of Vertices	181
No. of Unique Edges	185
No. of Edges With Duplicates	300
No. of Total Edges	485
No. of Connected Components	4
Maximum Vertices in a Connected Component	173
Maximum Edges in a Connected Component	480
Maximum Geodesic Distance (Diameter)	8
Average Geodesic Distance	3,621674
Graph Density	0,008133824

Overall metrics mostly reflect the structural properties of the documentation network.

The rate of numbers of unique and duplicated arcs shows that elements of this system are complex documents and not small and task-related regulations units introduced in 2.4. In general, the bigger the rate of Numbers of Unique Edges to Numbers of Total Edges, the smaller part of the operation is regulated by nodes of the documentation network. Obviously, it is true when only one type of connection is applied in the network.

The Number of Connected Components represents the connectivity of the network. Having 4 such components here means that this documentation structure is highly connected. This conclusion is supported by the fact that the far biggest part (173 nodes) of all units (181 nodes) belongs to the same subgroup.

The relatively high Maximum and Average Geodesic Distances ¹ reflect that regulation chain is the typical structural element (see in 2.1).

The very small Graph Density ² ($8.1 \cdot 10^{-3}$) denotes that references among documents are seldom.

According to the overall indicators mentioned above we can conclude that this documentation system consists of complex documents forming a highly connected network with regulation chains as a typical structural element and with relatively few links among the nodes.

3.2 Individual Metrics

To identify the different roles of nodes, individual metrics are calculated too. Their values can be seen in Table 2, Table 3 and Table 4. The associated subgraphs are highlighted in the following Figures.

The value of In-Degree represents the importance of a node. The most important document – from this point of view – has 22 individual incoming links, – due to its 22 connected neighbors (see in Table 2 and Figure 4). We can realize its important role in the graph too, since it is located in the middle of the network. There are two more nodes with more than 10 in-degree (14 and 11). For further reference, we call them ID1, ID2 and ID3.

Table 2
In-Degree metrics

Metrics	Value
Minimum In-Degree	0
Maximum In-Degree	22
Average In-Degree	1,464

¹ Geodesic distance is the distance between two vertices along the shortest path between them.

² Number of unique edges per maximum number of edges the graph would have if all the vertices were connected to each other.

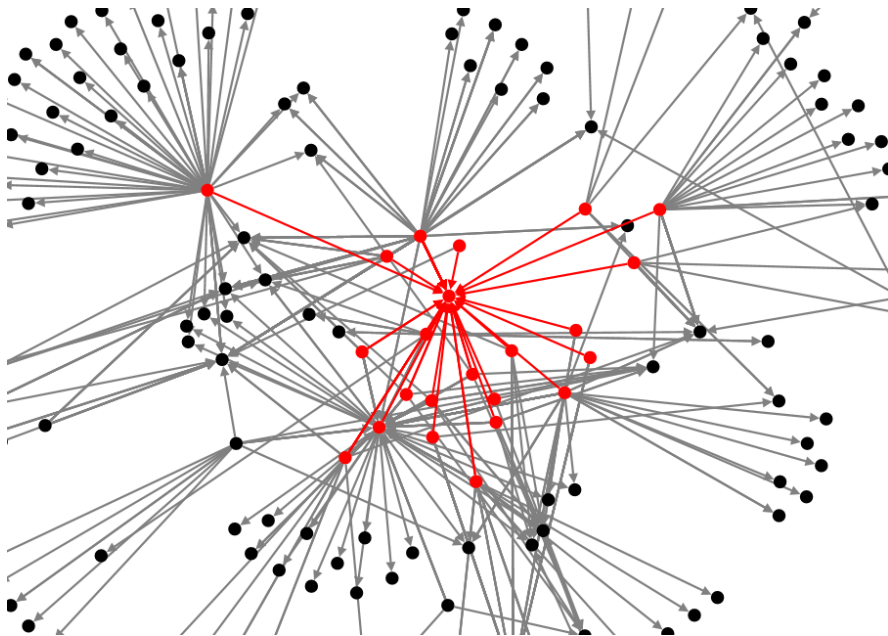


Figure 4

Subgraph of the nod with the biggest in-degree and its connected neighbors are marked in red

If we look at the list of nodes with out-degrees, we can see that the most sensitive document (call it OD1) seems to have 44 outgoing links (Table 3). If we see the graph (Figure 5), we can realize that these are concurrent links, which means all of them connect only two nodes. The conclusion is that despite the high value of out-degree, this document does not play a significant role in this network. It is interesting that out-degrees of D1 and D3 are zero, while D2 has the second biggest out-degree (24). However, D2 has only two neighbors, so its role is not significant either. Instead, the node (call it OD3) with the third biggest (19) out-degree has 19 neighbors (Figure 6). It means that it is the most sensitive document in this system.

Table 3
Out-Degree metrics

Metrics	Value
Minimum Out-Degree	0
Maximum Out-Degree	44
Average Out-Degree	1,464

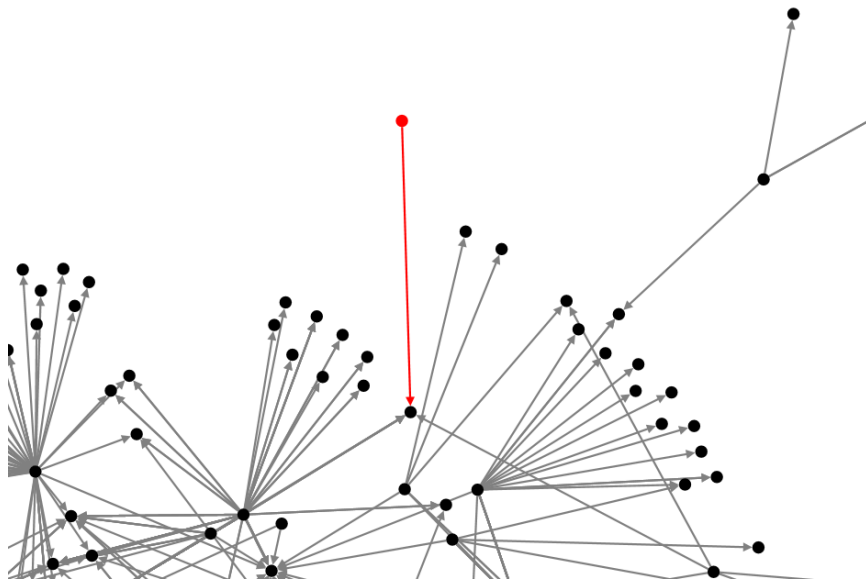


Figure 5

Subgraph of the nod with the biggest out-degree and its connected neighbors in red

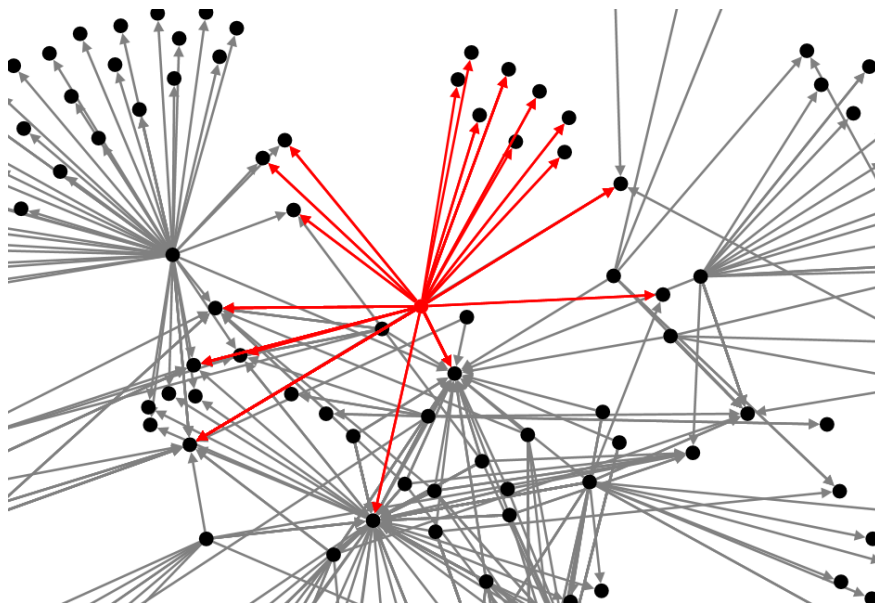


Figure 6

Subgraph of the nod with the third biggest out-degree and its connected neighbors in red

The averages of in- and out-degrees are equal. We could come to the conclusion that the sensitivity and the importance of the elements of this network are the same. If we check the degree distributions (Figure 7), we can see that there are very few nodes with high degrees while most of the nodes have much less connections. But we should not forget the fact that leaders of these lists have different numbers of neighbors, which influences the evaluation of the roles of the nodes.

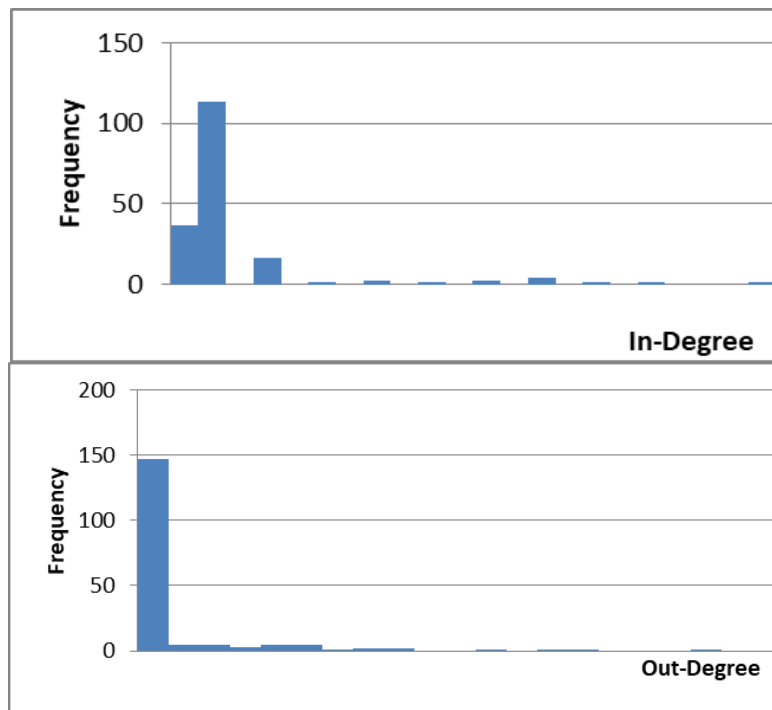


Figure 7
Degree distribution

As we wrote in chapter 2.3, Betweenness Centrality shows how critical the role of the node is in connecting network parts. 5 of the top 6 nodes in this list are the top 3 nodes of in- and out-degree lists. The exception is the node with the fourth biggest betweenness centrality (call it BC4) that has 14 out-degree and 0 in-degree (see in Figure 8). BC4 is a typical representative of network bridges, but in documentation system it is not so obvious. Here the different rates of in- and out-degrees show different types of bridges. If it has few in-degrees and a big number of out-degrees, the node is a so-called fork-bridge, which means it is a rather sensitive document. On the other hand, a node with few out-degrees and a big number of in-degrees is a join-bridge, which means it is an important document. If any of the degrees is null, the node is not a real bridge.

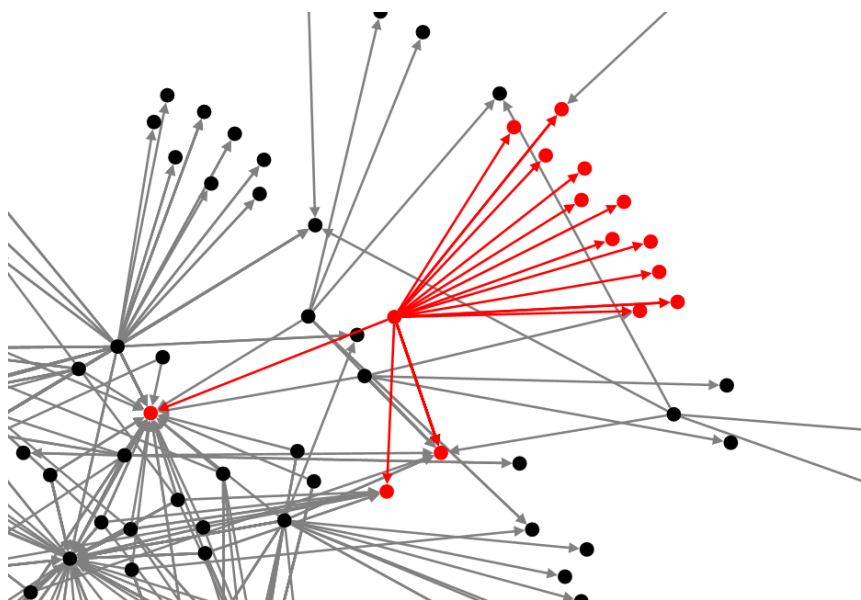


Figure 8

Subgraph of nod with fourth biggest betweenness centrality and its connected neighbors in red

The Closeness Centrality is the reciprocal of farness, which reflects how far a node is from other nodes it is connected with, i.e. in network terms how long the paths are between the connected nodes. In documentation networks the smallest closeness centrality belongs to the documents that take place in long regulation chains. In our sample network 173 nodes have 0.001, 0.002 or 0.003 value, and only 8 nodes have more (0.2 for 3 nodes, 0.333 for 1 node, 1 for 4 nodes). It means that there are typically long regulation chains between documents and most of the documents take part in these paths. Documents with high values are the part of regulation islands, i.e. isolated groups of nodes. Such distribution of closeness centrality shows that this documentation system is uniform but has a very long and complex set of connections that makes it difficult to easily overview and understand it.

Another type of centrality related indicator is calculated for our sample network to highlight the importance of documents more sophisticatedly. This is the Eigenvector Centrality, which takes into account not only the number of connected nodes but the degree of connected nodes as well. It means that in documentation networks the bigger degrees the neighbors of a selected document have, the more important or sensitive it is. The distribution of values (Figure 9) may be more interesting than its nominal value (Table 4). It demonstrates that there are noticeable differences among the eigenvector centrality values of nodes. In our sample network it fine tunes the description of importance of documents.

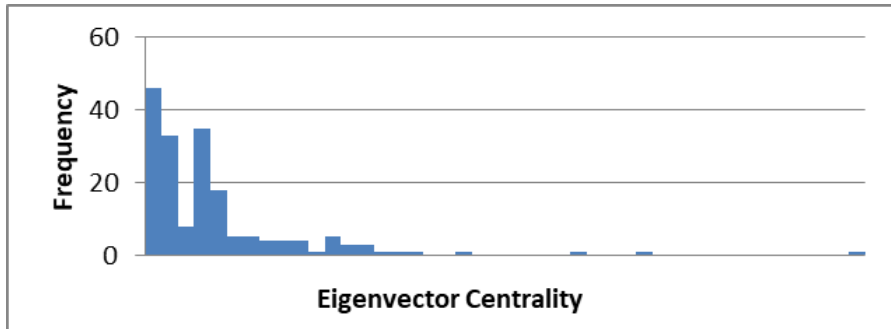


Figure 9
Distribution of eigenvector centrality

Table 4
Eigenvector Centrality metrics

Metrics	Value
Minimum Eigenvector Centrality	0,000
Maximum Eigenvector Centrality	0,057
Average Eigenvector Centrality	0,006

Conclusions

DQNET can be applied to identify and map quality regulation networks, to describe the properties of documents and to identify optimization opportunities. To conduct these activities properly individual and group network indicators have to be reinterpreted according to the specific characteristics of the documentation system. As an example, hubs can be important or sensitive documents, bridges can be fork- or join-bridges, subgroups of nodes can be regulation-chains of regulation loops.

References

- [1] Cao, Z., Wei, F., Dong, L., Li, S., Zhou, M. (2015): Ranking with Recursive Neural Networks and Its Application to Multi-Document Summarization, pp. 2153-2159, Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence
- [2] Carley, K. M. (2015): Crisis Mapping: Big Data from a Dynamic Network Analytic Perspective, Journal of Organization Design
- [3] Tan, L. S. L., Chan, A. H., Zheng, T. (2015): Latent quality models for document networks, arXiv:1502.07190v1, Annals of Applied Statistics
- [4] Wan, X., Yang, J. (2008): Multi-Document Summarization Using Cluster-Based Link Analysis, pp. 299-306, Proceedings of the 31st annual international ACM SIGIR conference on Research and development in information retrieval

- [5] Wang, C. (2015): KnowSim: A Document Similarity Measure on Structured Heterogeneous Information Networks, Data Mining (ICDM), 2015 IEEE International Conference on
- [6] Yang, Z., Yang, D., Dyer, C., He, X., Smola, A., Hovy, E. (2016): Hierarchical Attention Networks for Document Classification, pp. 1480-1489, Proceedings of NAACL-HLT
- [7] He, Y., Wang, C., Jian, C. (2017): Modeling Document Networks with Tree-Averaged Copula Regularization, pp. 691-699, Proceedings of the Tenth ACM International Conference on Web Search and Data Mining

Risk Assessment of the Human Factor in the Field of Building and Infrastructure Defense

Tamás Berek¹, Judit Kovács²

¹ Institute of Military Leadership Training, Faculty of Military Sciences and Officer Training, National University of Public Service, Hungária krt. 9-11, H-1101 Budapest, Hungary, Berek.Tamas@uni-nke.hu

² Institute of Microelectronics and Technology, Kandó Kálmán Faculty of Electrical Engineering, Óbuda University, Tavaszmező u. 15-17, H-1084 Budapest, Hungary, kovacs.judit@kvk.uni-obuda.hu

Abstract: In order to establish the concept of building and infrastructure defense, a complex security system must be created by making, analyzing and interpreting an appropriate plan. This task is especially difficult and complex for defending buildings of unknown functions. Industrial projects usually differ from what was planned both in space and in time. The authors of the article introduce the general aspects of security personnel and the characteristics of risk assessment. The basic points of configuring the labor force components of building and infrastructure defense are also introduced.

Keywords: complex security; defense concepts; risk assessment; human factor

1 Introduction

The threat level for any building and its respective infrastructure is determined by several factors. Some of these factors are the security degree of operation, the demand and the value of the used materials, technical equipment and information, and the criminal infection of the area. The time of the day, the reliability of the applied security system, the speed of action and troubleshooting, and the features and territorial impact of undesirable acts are also of great importance [1].

Analyzing the question from a distant approach, the aim is to maintain a safe state of the building and its respective infrastructure. This state, providing the ideal status that the operation of the security system is fault-free, may seem steady in time, though this steady state is only an outward seeming. All acts that are performed inside the area of the building, the equipment, the quantity and the risk of the materials used are relatively easy to be determined. From certain points of view, the changes in the safe state may be prevised, knowing – among others – the

feature of the building, the acts performed inside, the applied technology and the materials used.

Security could be defined as the safe state of somebody or something. However, this safe state does not literally exist, because security is the complex outcome of some specific existence or actions and the endangering factors of them. It means that security may be interpreted only together with endangering factors. It is the very moment when an endangering factor appears that the expression of security gets its deeper meaning. The higher level the endangering factors of existence or normal operation are, the lower level of security is [2].

It follows that the status of security is fundamentally determined by the endangering factors and by the protection applied against them [3]. Creating the complex security system of a building and its respective infrastructure, one must be aware of and recognize the nature of outer and inner endangering factors that may affect security. After the evaluation of these factors, the acts in response and the whole structure of defense must be laid out.

Simplifying it, in the case of any building and its respective infrastructure, the subject of defense and the sources of dangers must be specified by recognizing and analyzing the endangering factors arising from the environment. The security system must be planned and carried out by knowing these factors.

2 The Role of Security Personnel in Property Protection Complexes

Complex property protection is made up of components based on one another. The aim is to reduce the probability of certain risks, as well as to moderate the adverse consequences of possible incidents [4].

To identify the rate of the components of the complex security system is an inevitable task during the design process, since this act will grant the effective and fault-free operation as well as the phenomenon of synergy: the interaction of subsystems will produce a combined effect greater than the sum of their separate effects.

Electronic and electrotechnic subsystems are among the most important elements of security systems and their reliable operation is essential. The development of low-voltage subsystems that are optimal from the aspect of reliability is greatly facilitated by the test method of the principle of determining disturbance states. As part of that process, the analysis of disturbances of both technological and human error origins must be worked out [5]. A disturbance state is a state of the system when it cannot perform its function, due to the effects of well-determined technological or human disturbances.

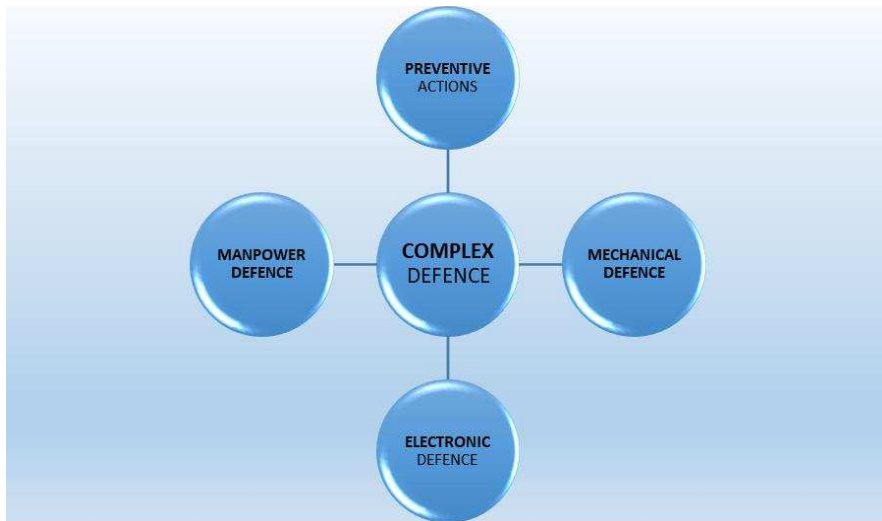


Figure 1

Components of complex property protection

Source: edited by Berek

Any type of errors cause a disturbance state when it is concerned with processes. Disturbances may also result in unacceptable consequences when critical failure is due to a disturbance state. Two main groups of the factors that lead to human errors causing disturbances may be distinguished as internal and external factors. The next basic categories of internal factors may be physical, emotional, cognitive and social effects that also include more categories like personality, intelligence, motivation and ability. External factors can be divided into organizational and environmental factors. In each of the categories, separate analysis is needed to determine to what extent certain factors cause a specific error. [6]

Nevertheless, in certain areas, special features apply. The mechanical, electronic and electrotechnic elements of security systems for construction-industrial projects, are generally inefficient and sometimes are even absent, especially at the early stages. In this case, due to the occurring variance, security personnel has the main role in defense.

Construction-industrial projects are especially the ones where rates of security components have to be changed at the different investment phases. These changes may only be handled properly with security systems that were designed to involve a possible option for flexibility.

This flexibility is provided mainly by security guard. It is also this flexibility that guarantees adherence to specific features of the security systems: at the loss of any elements – that often happens during constructions – active elements may cover all parts of security (though perhaps, at a lower level of efficiency) [7].

Another special issue is the physical security of facilities that store dangerous industrial materials. In order to prevent unauthorized access to and/or theft of radioactive and infectious materials and toxins, providing the defense of laboratories and other facilities is of extreme importance. At the same time, the special staff of these institutions, in certain positions, is exposed to physical, chemical, biological and radiological risks. However, these risks can significantly be reduced by the development of a well-designed defense program. In order to avoid direct threat by everyday working conditions to users of dangerous work areas, devices or materials, the careful construction of the physical protection of the affected areas and equipment is extremely important. The same protection is vital so that the potential hazards of costly devices and dangerous materials or devices containing dangerous materials do not leave the controlled working areas or property of the institution, by unlawful appropriation.

There may be serious risks of abuse by competent persons with access rights to hazardous substances, against which, not only certain components of security or protection should be reinforced, but also, protection elements against intentional personal abuse should be developed. Laboratory accidents and the release of dangerous biological materials may not be due solely to deliberate illegal activities or sabotage, but the accidental release of hazardous substances may also result from the improper use of infectious substances in laboratories or by their inappropriate packaging or transportation [8].

The protection of controlled work areas and working processes, in the lab complex, including the personnel involved and the protection of hazardous materials and waste storage facilities are also of great importance. The same degree of protection is required for the lab areas not considered to be working places and for the external environment of the laboratories.

In order to ensure a continuous and comprehensive protection, when designing the security system, there is a great need to coordinate the efficiency and harmonization of each independent autonomous subsystem and to ensure the conditions of supervision. The effectiveness of physical guarding is provided by the effective combination of mechanical and electric devices and security personnel and not overlooking the role of preventive measures.

When ensuring the protection of the hazardous areas of the lab facilities, there is little chance to use labor force, therefore, the rate of electronic protection devices should be increased. At the same time, the efforts to reinforce internal control are also in the forefront. It is the task of the lab staff to control the regulations and procedures, and to operate and maintain the security systems. It is well-known that the efficiency of the entire property defense system is determined by the efficiency of its weakest element. In improperly built systems, the living component is quite often the weakest link. So, it cannot be emphasized enough, how important it is that the human factor is taken into account, when structuring security systems.

There are situations where the presence of temporarily or operationally stored substances is a threat in itself. During the planning and development of protection in the controlled area, one of the main aspects, is that the safety engineering subsystems are designed to meet the intended function of the laboratory and to provide the highest level of technical, mechanical, electronic and personal security.

In an emergency event, the controlling system is able to perform several actions simultaneously; nevertheless, its basic task, is the prevention of emergency situations. So, in case of a possible occurrence, the personnel support of the operation of the system is needed. Monitoring the controlled areas, it has to alert the operator immediately so that they can intervene in time. Security guard has a very important role, in this case as well.

In contrast to technical systems with average parameters, security personnel is capable of managing compound or unforeseen situations [3]. However, the subjectivity of the human factor might as well be its own vulnerability, since personnel may base the sources of inside hazards that are difficult to detect and identify, and it is also hard to provide protection against them. The occurrence of events like damages due to any improper execution of tasks, sabotage, theft or participation in them, or releasing important information that would provide the strength of protection may threaten the whole defense system. The prevention of these events is extremely problematic.

Building and infrastructure defense is a very compound task. The lack or the weakness of any part of the property protection complex will affect the overall efficiency of the security system. The components of complex systems (access control systems, security monitoring systems, etc.) are also involute security subsystems. It is essential to meet the requirements of fault-free operation. In addition to the high-level integration of electronic, electrotechnic and mechanical security subsystems, the activities and preparedness of the operating crew are of great significance [9].

Since the human factor has a key role, the analysis of it may have a considerable influence on the establishment of risk assessment and risk management. Human performance has a fundamental impact on the reliability and security level of various systems. Generally, the role of the human factor, in connection with the occurrence of any events, may be divided into three main groups. People may cause, prevent or be the victims of particular events, thereby giving an improved approach to risk assessment and risk management.

From the aspect of security, the human factor appears in two, rather contradictory ways, of the previously mentioned, as follows. Within manufacturing security systems, design is a highly challenging activity. Since the environment is constantly changing, the proper designs do not perform as expected, with the frequent overestimation of how efficiently people will work [10]. On the other hand, people are able to cope with unforeseen situations, to analyze and to create

solutions. Without human actions many incidents could lead to accidents. Safe behavior does not mean the absence of errors, but the positive human contributions to safety, even in the form of prevention [11].

As the human factor is always present among the main reasons of accidents and disasters, human contribution has priority in any analysis of risk assessment. According to different surveys, 45-80% of errors are due to the human factor, varying with ways of approach. The special role of the human factor was recognized decades ago, and research on human factors has been present since then. Human errors have been categorized and the broad use and development of human reliability assessment has been urged. Initially, it was discovered that specific systems must be developed to analyze the events related to human factors. Later, it was shown, that human factor-associated common cause failures may appear in any kind of security systems.

3 Analyzing the Human Factor in Risk Assessment

Among the reasons that may turn incidents to accidents, as well as, among the main reasons of industrial accidents, the human factor is always present. Consistent explorations of consequences will recognize human errors even in the depths of technical reasons. Based on the research of Rankin and Krichbaum, the role of the human factor in the occurrence of accidents shows a dramatic rise, reaching up to a 70-80% level, regardless of the technological conditions [12].

This significant increase has two main reasons. One of them is the sophistication and the high-level reliability of the mechanical and electrical equipment, while the other one is the greater human involvement in the controlling processes that is due to the complexity of systems. Not only do the sophistication and the high-level reliability of the mechanical and electrical equipment greatly reduce the number of technical errors, but they also give opportunities to manage critical processes, even at the events of system failure and breakdown. The greater human involvement in the controlling processes, as a consequence of the complexity of systems, means that humans primarily become the supervisor of automated processes.

Human contribution has a place of utmost importance in any analysis of risk assessments. The first progressive development of risk-based approaches occurred in relation to the analysis of electronic and electrotechnic systems, in the fields of space technology, nuclear power and the chemical industry. However, due to the diversity of physical and chemical processes, as well as, control strategies and procedures, special techniques have been developed for the specific needs of each area. As an example, regard the method of hazard and operability study (HAZOP). It was first introduced in the chemical industry and has since been considered

necessary for the preliminary assessment of any complex system that consists of several processes of either serial or parallel structures that involve subsystems of dangerous chemical or thermodynamic reactions.

The estimation of the impact and consequences of the risk events on people, property and environment is realized in the risk assessment process. The calculation of the probability of these risk events actually happening, as well as, determining their potential impact are important parts of the risk assessment process.

By its nature, the process of quantitative risk assessment is based on probabilities. It recognizes that accidents are rare, and that the potential risks and events may not be completely avoided. As serious incidents occur or not, over the lifetime of a given process or building and its infrastructure, it is not appropriate to base the evaluation process on the consequences of isolated events alone. However, the probability of the cases that have actually happened should also be considered. These probabilities and the levels of risks derived from them must have an impact on both the design level and the operational and organizational controls and revisions.

In the process of integrated risk assessment, “risk identification” as the first step involves describing the system, determining the possible events and the responses to them, as well as the classification and the filtering of events. The second step, i.e. “modeling event scenarios”, is based on event tree analysis, and its objective is to place the sequences of events among the states of losses. The main parts of the next step, “analysis of consequences”, are the assessment of the consequences and the analysis of the moderating effects. The following “evaluation of the frequencies of events” is one of the most complex tasks. In addition to the actual evaluation of such frequencies and system analysis, the analysis of the human factor is usually performed at this step. Finally, in “risk assessment”, determining the risks is brought off by means of the consequences and frequencies.

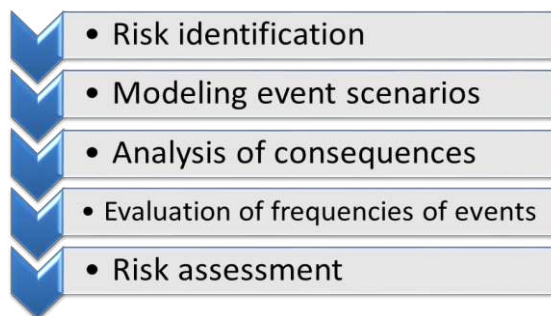


Figure 2

A detailed plan for the integrated risk assessment process

There is no doubt that integrated risk assessment is the best way for meeting the requirements of risk assessment today. In the comprehensive plan of the risk assessment process, human factor analysis is interpreted as a subtask of the evaluation of event frequencies. It means that analyzing the human factor does not get a role in the preceding or simultaneous steps of the process: neither in risk identification, nor in modeling event scenarios, nor in the analysis of consequences. However, it is very important to be aware of the impact of human factors even at the beginning of the processes, as in this case the corresponding details of plans may be modified easily and at a low cost. It is therefore recommended to take the human factor into account already from the first step in any integrated risk assessment.

The probability that certain events occur due to the human factors involved may be determined by an appropriate human error analysis method. Human error analysis covers the systematic specification of the factors affecting human performance and the exploration of situations that are likely to give rise to errors; this is the way leading to incidents. This analysis may involve the identification of interfaces that are influenced by the errors. Based on the frequency of occurrence or on the severity of the consequences, a relative ranking of the errors may be established. The results may be qualitative or quantitative as to their nature. They also involve the systematic listing of errors that are likely to occur during normal or emergency operation. The error rate depends on many factors, ranging from stress over experience to the complexity of the task or to the right skills, including situation-specific specialties.

Human error is a general concept, which includes every situation when the planned sequence of mental or physical actions fails to achieve its planned and desired aim. This failure is not due to any kind of stochastic circumstances [13]. Human error may also be considered a failure of a human action, due to internal human failure mechanisms, which is to loosely describe any sub-optimal human performance. Two main groups of human errors are errors of commission (wrong human actions) and errors of omission (missing human actions). A human error as the consequence of the difference between the planned and the realized action or performance, may be categorized as a slip, a lapse or a mistake. A separate group of errors is violation, when the action is not allowed, prohibited or not appropriate. Latent errors may also play an important role, although this type of errors is usually difficult to identify because of their distance from the occurring event both in time and in space [14]. Human failure is the failure of a defined human action in any Human Reliability Analysis (HRA) model. There may be more reasons leading to failures attributed to human errors. A human failure may affect components – that is called a fault, and it processes when disturbances occur. A failure that results in unacceptable consequences such as, unavailability or malfunction leading to personal or property damage is called a critical failure. Another possible classification of human errors, usually taking place in Probabilistic Safety Assessment (PSA) models, depends on the chronology of the

human error and the occurring event. Three types of errors may be distinguished by this chronology: an error of human performance type A is an error that is committed during a human action before the initial event, mainly in connection with the availability of the system (for example in connection with the actions of maintenance), an error of human performance type B is an error which causes a direct initial event, while an error of human performance type C is an error that is committed during the human actions made for averting breakdowns or accidents. In the case of errors of human performance type C, the following groups may be differentiated: the lack of a needed action, an action made by mistake and the error of an action made for compensating the lack of a needed action [15].

In any of the above categorizations, the role of violation is not handled as being as important as it is in reality. In Reason's categories, the violation of intentionally causing harm is not even regarded as being a human error. However, if human error is considered to be the consequence of the expected and the realized actions or behavior, violation is an error, as it differs from the expected human behavior. In reality, deliberate actions of this type may have serious conclusions, their number is significantly increasing and they cannot be regarded as "low probability", isolated events any more. The violation of intentionally causing harm being treated as a human error makes it possible to be a part of human factor research, which is the cornerstone of prevention.

In addition to the previous ideas on the violation of intentional causing harm, violation in the traditional sense is also a serious problem. Violation may be motivated by the search for simpler and faster solutions. For example when a worker crosses a conveyor belt because it is simpler than bypassing it, or by risk-seeking behavior as when the worker crosses the conveyor belt, because he wants to show his courage to peers. The offenders are often unaware of the risks; thereby violation may become a habit [6]. According to Skriver [16], after a specified period, the prescribed processes are no longer evident in these cases. They consider the main causes of violation in organizational factors: in the lack of adequate equipment, working environment and supervision and due to the fact that there is no consequence of the committed violation.

Prevention here, as in all other areas, has an important role. Moreover, the energy spent on prevention will be compensated. Based on the connection between violation and the number of rules to be kept, there is a number K given by the number of rules to be kept weighted by the difficulty of keeping them, which should be analyzed in the given situation. There exists a number K_0 such that if $K < K_0$, then no violation will happen (or only with negligible probability), and if $K_0 < K$, then violation will happen (with a considerable probability). The aim is to make a system of sufficient rules for the given task where $K < K_0$ [6].

In building and infrastructure defense, the implementation of the daily operational work is specified by the operational rules of the system and the service instructions of the security service. Among these rules, the daily tasks, protocols and procedures must be defined precisely for different situations. The daily work of a security guard's life is carried out according to an appropriately assembled policing scheme, with the tasks and procedures being unchanged. Though, there may be differences in its implementation. Each error resulting from time management may be declared a typical example that may cause a disturbance state based on human factors. One of the basic principles of an operational process, is the assumption, that workers are highly predictable and standardized in their behavior, regarding their schedule. They always start work on time, operate at a constant rate throughout the day, take breaks at planned times, rotate properly, etc. Nevertheless, such regular behavior of workers rarely occurs in practice. According to a test made in the UK [13], the analysis of the data suggested that up to one third of the potential time for production is lost due to stoppages, extended breaks and disruptions to the flow of the line. Not only does the loss of time cause a recession of production, but it may also cause disturbance states. [10]

Knowledge and awareness of the human factors are basically important in preventing the development of disturbance states, and therefore, they have a distinguished place in design processes. When the human factor is taken into account, the reliability of complex systems may indirectly be increased. There is no doubt that human performance has a fundamental impact on the reliability and safety level of complex technical systems, such as security systems. Among the major contributing factors of disturbance states, the human factor can be found in each case. As a consequence, the human factor must always be taken into account when analyzing disturbance states in building and infrastructure defense.

Conclusions

The maintenance of the labor force component of a complex property protection requires constant control, and security service managers must monitor the system. It is well-known that the effectiveness of a security system is characterized by that of its weakest component. That weakest link often happens to be the security personnel. Apparently, as far as crew becomes unreliable, the entire security complex is threatened. In the cases when it is recognized, the negative effects may often be outweighed by technical upgrades. Control systems that can be implemented are nowadays indispensable. Such systems may include the camera surveillance of workers, the establishment of patrol monitoring systems, etc. [17].

Generally, despite the fact that the vast majority of errors, including technical reasons, are due to the human factor, people are able to maintain safe and economical operations, and are also capable of providing a responsive action to disturbance states at the same time. In this way, human performance affects the probability of all unexpected situations and their consequences. Today, the well-established industrial security applications and the procedures of design and

operation make the basis of risk management. The wide-spread awareness of the possible dangers has implied the development and use of systematic approaches, methods and tools of risk assessment procedures. These are often referred to as hazard analysis or quantitative risk assessment [6].

A risk analysis is required in different areas of production: in business, industrial production and environmental protection in the field of work safety. Although laws and standards regulate its implementation, they do not include specific execution [18].

Risk assessment and risk management are two of the most important jobs done today to achieve maximum security levels. Analyzing human factor can have a major influence on the risk assessment and risk management process, because the role of the human factor is crucial. Human performance has a fundamental impact on the reliability and security level of various systems.

The maintenance and operation of the labor force component of a complex property protection requires an active presence, as the risk of this component is continuously assessed. It was concluded that human performance affects property protection – as a complex system, which is based on technological and human factors – on the whole. Also, it was shown that human performance has a basic impact on the safety levels and reliability of complex technical systems in building and infrastructure defense.

Beyond the exploration of errors, it is also vital that the mapping of the reasons of errors are done, which is proven to be suitable, by using methods based on cognitive theories.

Within manufacturing security systems in building and infrastructure defense, the design and redesign activities are both challenging. The competitive environment is constantly changing and there is a demand to make products cheaper, better and faster. In this kind of environment, people who carry out repetitive, manual tasks seem to remain critical to the success of the system. Designers of security systems often have little appreciation of the wide range of factors that influence human performance. This can lead to “proper” designs not performing as expected, with engineers frequently overestimating how efficiently and effectively people will work. [10]

The key to a successful solution is to improve the awareness of engineers, concerning the impact the human factor has on the design. It is especially important to improve this awareness at beginning of the design process; at this stage most negative factors can be more easily and inexpensively altered.

References

- [1] Lukács György: Új vagyónvédelmi nagykönyv, CEDIT Kft., Budapest, 2002

- [2] Berek Lajos: Biztonságtechnika ÁROP – 2.2.21 Tudásalapú közszolgálati előmenetel jegyzete NKE 2014
- [3] Báthori B.- Bodrogi F. – Szili L.: Őrzés védelem, jegyzet, Pro Lex Oktató és Szolgáltató KKT, Budapest, 1995
- [4] Utassy Sándor: Komplex villamos rendszerek biztonságtechnikai kérdései, Doktori (PhD) értekezés, 2009, ZMNE
- [5] Zsigmond Gyula: Biztonságtechnikai rendszerek hibamentességéről Bolyai Szemle 2010:(4) pp. 207-213 (2010) <http://uni-nke.hu/downloads/bsz/bszemle2010/4/15.pdf>
- [6] Kovács Judit: Az emberi tényező matematikai modellezésének lehetőségei a katasztrófavédelmi kockázatértékelés és kockázatkezelés területén Doktori (PhD) értekezés, 2011, ZMNE
- [7] Teke András: Az őrzés mint rendészeti alaptevékenység VI., in: Rendvédelmi Füzetek 2000/45, a Rendőrtiszti Főiskola kiadványa, Budapest, 2000
- [8] Berek Tamás - Pellérdi Rezső: ABV (CBRN) kihívásokra adott válaszlépések az EU-ban 2011, Bolyai Szemle XX. évf. 2. szám, ISSN: 1416-1443
- [9] Bodrácskó Gyula – Berek Tamás: Megelőző intézkedések szerepe a komplex vagyónvédelem területén, építőipari beruházások során, 2010. Hadmérnök, www.hadmernok.hu/2010_1_bodracska_berek.php
- [10] T. S. Baines, R. Asch, L. Hadfield, J. P. Mason, S. Fletcher, J. M. Kay: Towards a theoretical framework for human performance modelling within manufacturing systems design (Simulation Modelling Practise and Theory 13, 2005, 486-504)
- [11] NEA (2003): Nuclear Regulatory Challenges Related to Human Performance, ISBN: 92-64-02089-6, OECD, Paris, 21 pages
- [12] W. Rankin, L. Krichbaum, Human Factors in Aircraft Maintenance, Integration of Recent HRA Developments with Applications to Maintenance in Aircraft and Nuclear Settings, June 8-10, 1998, Seattle, WA, USA
- [13] James Reason & Alan Hobbs: Managing Maintenance Error- A Practical Guide, Ashgate Publishing Company, 2003
- [14] James Reason: Managing the Risks of Organizational Accidents, Ashgate Publishing Company, 2004
- [15] Gyula Zsigmond-Judit Kovács: Determination of Disturbance States with a Special Focus on the Human Factor Bolyai Szemle 2007 :(3) pp. 187-193 (2007) http://uni-nke.hu/downloads/bsz/bszemle2007/3/16_kovacsjudit_new.pdf

- [16] Jan Skriver: The Human Factor Human and organisational aspects of RCA, part I and II, Resilience IAEA Workshop on Root Cause Analysis 9-13 November, 2009
- [17] Berek Tamás - Bodrácska Gyula: Az élőerős őrzés az objektumvédelem építőipari ágazatában Hadmérnök, V. Évfolyam 4. szám - 2010. december http://www.hadmernok.hu/2010_4_berek_bodracska.php
- [18] Berek Lajos- Tóth Georgina Nóra: Risk and chance of practices Hungarian Journal of Industry and Chemistry 38:(2) pp. 193-196 (2010), <http://mk.uni-pannon.hu/hjic/index.php/hjic/article/view/300/279>

The Role of Knowledge Management in Developing Quality Culture

Andrea Bencsik

Széchenyi István University, Egyetem tér 1, 9026 Győr, Hungary;
J. Selye University, Bratislavská cesta 3322, 94501 Komarno, Slovakia
bencsika@sze.hu

Gabriella Horváth-Csikos

Szent István University, Páter Károly utca 1, 2100 Gödöllő, Hungary
Horvath.Csikos.Gabriella@gtk.szie.hu

Abstract: The study reveals the connecting points, where the elements of knowledge management system as part of the corporate processes (studiously the quality systems) play an important role in the successful operation of the company. The characteristic features of quality culture are in focus, of which support is put in parallel with the characteristic features of the learning organisation's culture phrased as a precondition of knowledge management system. With this comparison, the author guarantees the understanding that by solving quality problems and developing quality culture with supporting the elements of knowledge management system will contribute to reaching corporate success and strategic goals.

Keywords: knowledge management; organizational culture; learning organization; quality culture; TQM

1 Introduction

Knowledge management, as a young research field is still fighting its battles nowadays to reach its reason for existence. However, several studies (in theory and in practice as well) justify its contribution to the successful operation of the corporation, but the presence of this way of thinking is not natural in the every day life of corporations. Developing a knowledge management system is a challenging task nowadays for company managers. However, more and more people realize its significance and its positive economic results – which is based on the proper evaluation and management of knowledge, the human capital –, the actual

development is still to come [6; 10]. Operating the knowledge management system is a question of approach and way of thinking, it does not require a separate corporate unit or responsible staff, and it only requires the development of a proper culture and the integration of knowledge into everyday processes. The demand (requirement) of knowledge sharing is present more stressfully among the goals of the organisation, but as we know, the pressure can have results only in the short-run. Taking this requirement, from a professional aspect, into consideration, we can state that knowledge sharing cannot be realized under such workplace conditions, where corporate culture is working against the development of a trustful atmosphere supporting knowledge sharing [22]. The research results – based on recent inland practice – show that in more cases the demand for integrating knowledge management’s way of thinking into every day practice appears on strategic level, but this can only be observed in isolated solutions in reality. We cannot mention any domestic organisation, where the base of its operation would be a well-developed system or would contribute totally to reach the strategic goals. This means that although the number of organisations increases, where their requirement (at least on strategic level) is phrased, but still there is a lot to be done by us, the knowledge “evangelists”. In order to accept, to get to know and to apply the logic of knowledge management, its way of thinking, its models and their application we need to do a lot on a daily base [21; 23].

This battle for acknowledging the reason for its existence makes us remember the period when quality management and the development of quality system was fighting their professional battles.

Similar to the hard times that the demand for developing quality systems had to live through some 10 years ago, made we think about what solutions supporting mutually each other can be mentioned as argument in connection with the role of the two systems, which make the connection obvious with regard to corporate success.

Although the emphasized status of quality questions, the role of quality in corporate operation cannot be questioned nowadays, in several cases we can still see formal solutions, which do not help in solving problems regarding quality. The results of a lot of former researches highlight connections between knowledge management and quality management from different viewpoints (in theory and in practice as well). These connections are confirmed in service sector, in libraries, in higher education and in companies as well [32; 27; 19; 17; 12; 8; 4; 2; 1]. Some of the most important thoughts will be presented in the followings, which support the solution of quality problems in organisations from the side of knowledge management and also the layman reader can notice the tight relationship and cooperation possibility between the two systems. In order to adapt this train of thoughts in practice, we should not forget that the top management of organisations play a key role in the success of developing the systems and operating them [20; 25].

The Association for Excellence Public Company, in Hungary – Szövetség Kiválóságért Közhasznú Egyesület – [31] in a 2005-study revealed several problems, which can be traced back to the deficiencies of knowledge management in connection with quality problems. The research, which involved 27 companies, identified the following problems:

- repeated similar mistakes (costs),
- duplicated carrying out of tasks (previous projects, their results are not known),
- lack of information (e.g.: customer service),
- lack of sharing inner good ideas, best practices,
- weak link (1-2 people have the knowledge/information only),
- integration of gained knowledge is slow or missing (slow product development, market competitor overtakes),
- information/knowledge sources cannot be reached easily (frustrated worker).

The results of the research show that among the problems influencing quality work in several cases there were problems referring to knowledge, knowledge share and knowledge management. In order to prove the above mentioned tight relationship between knowledge management and quality systems, some basic concepts and models have to be introduced.

2 Material and Methods

2.1 Briefly about Knowledge Management

Nowadays changes are increasing much faster than ever before. These changes go hand in hand with overvaluing knowledge, as a production factor. Parallel to overvaluing knowledge, the speed of its term of limit increases as well. Therefore, the task of management is to ensure substitution and care properly about the maintenance of values. The bigger value novelty-developing knowledge has, the more difficult it is to obtain, the sooner it can fall into disuse and the more hidden it is, the bigger its significance is in order to maintain competitive operation. To integrate, to manage and to develop this useful personal (hidden/tacit) knowledge can happen only with the tools and methods of knowledge management [5].

The history of knowledge management (KM) dates back to the years of 1980s. The top managers of companies have already been talking about it since 1990

when they were forced to rethink their knowledge about management and business operation. The concept of knowledge management entered common knowledge in 1991, due to an article, which was published in Fortune magazine. This article was written by Tom Stewart, and the title was Brainpower [30]. After some years it was already called the next big challenge following BPR and TQ.

Several definitions of knowledge management have appeared since the years of 1990, and its concept was phrased in different ways, which was formed on one hand by knowledge perception and on the other hand by the idea about how to manage knowledge. According to Sándori [28], Fehér [13], Davenport and Prusak, [9] knowledge management is the effective connection between those who know something and those who want to know something. Some people say it is quite a trendy phenomenon, while others think that with the help of it, organisations are able to react to processes happening in the present and they are able to make the necessary steps by using the knowledge material of the past and by relying on analyses made in connection with the future. Not only the effectiveness and the global competitiveness of organisations can improve by using knowledge management, but also the given country can make profit out of it by the organisation. The Work Committee of Knowledge Management at the Hungarian Academy Of Sciences defined the content of knowledge management by consensus the following way: knowledge management (KM) is a process (management subsystem) and culture, in which disclosing, gathering, creating, recording, keeping, transferring and continuously increasing knowledge capital is managed in an integrated way and supported by information technology. Its aim is to increase the production of added value of the organisation and to enlarge its innovation potentials; its key concept is synergy [24].

The significance and importance of KM is approved by more and more companies and the research of the years 2000 show that the four-fifth of the European organisations considers it as strategic tool. At present, managers do not think of knowledge management as a technique, but rather they concentrate on knowledge as a key-resource. Managing knowledge has become a tool for increasing corporate competitiveness nowadays by managing knowledge consciously and systematically.

As a general summary, we can conclude that knowledge management wants to find the answer for questions, who, when, where and in what form needs knowledge. The role of knowledge management is to ensure the competitiveness of the company in the new, knowledge-based economy [5].

2.1.1 The Focuses of Knowledge Management

From the various definitions of knowledge management, from a practical point of view there are two basic approaches known: the human-centred and the informatics-centred. From the two theories it is the human-based direction, which creates the base of my thinking as primarily this has reason for existence from the

point of view of the quality systems discussed in the introduction. Although in the followings, the previously mentioned two focuses will be presented briefly, which are dominant in the process of developing knowledge management, emphasizing that IT is an important precondition formed technical point of view, while it is impossible to operate successfully even a developed system without a trustful culture.

Figure 1 shows the position of knowledge management on the curve depicting management fashion-wave. From the figure it can be seen that from the aspect of maturity, the quality systems have reached the much more accepted and stabilized phase.

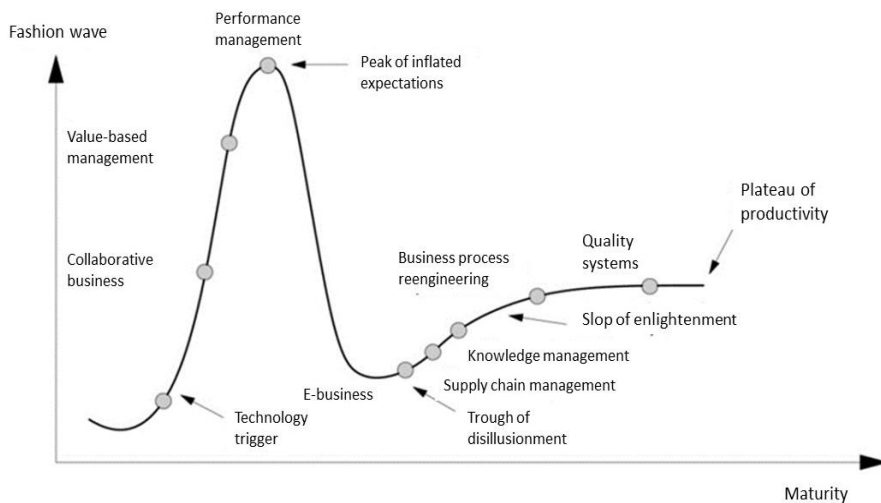


Figure 1
Management fashion [18]

2.1.2 IT as One Precondition for Building KMS

The debate between researchers and practical professionals about the role of technological solutions is the same age as knowledge management as a management tool's coming to the front. Even the most famous researchers have different opinions. Those who belong to the first group represent the opinion that states the almightiness of IT, while the other group does not argue with the necessity of its existence, but does not regard it to be the only reason for success. According to the opinion of Bögel [7] knowledge management could become popular because technological development made a more effective knowledge management possible. However, Dougherty [11] completely refused the role of technological solutions. The truth is somewhere between the two. According to our present judgment of values IT is necessary, but not an ample precondition [14]. The operation of a well-structured, logically built and reasoned IT system

can be phrased as the background support of knowledge management systems. The main aim of informational and communicational technologies in the field of knowledge sharing is to connect the concerned parties.

2.1.3 Cultural Characteristics Ensuring the Operation of KM

Neither the corporate structure, nor technological solutions provide value themselves or an effective knowledge management system. The organisation is made up of people, knowledge is inside people, they create it and use it; therefore, the role of human resources has to be handled as a high priority. The development of a proper corporate culture is necessary so that the workers can operate the organisation the way their management expects them too and the interest of the company requires. The proper corporate culture in this case means the willingness to knowledge share, the use of others' knowledge and also the cooperation developing from the common knowledge of the organisation.

The guru of corporate culture, Hofstede and his fellow workers [16] compared culture to the software of the brain. They developed the definition further, according to which, culture is the social programming of the brain. It differentiates the members of a group from each other.

If in an organisation the culture is not the proper one from the view of knowledge management, then the workers there can have a negative approach to certain processes of knowledge management. This affects mainly knowledge share, which has to be an activity of top priority. If the workers feel that they are not supported in knowledge share, then they keep their precious (tacit) knowledge, which cannot be documented and if they leave the company, they take this knowledge with them.

2.1.4 The Reason for Existence of a Learning Organisation - Senge Model

The aim of the operation of a learning organisation and a knowledge management system is to mobilize the divided or hidden knowledge in the company through organisational groups. With this it will be possible to react to market demands and to the steps of competitors faster and in a more flexible way. As a result of it better quality can be produced by better planning and by more effective work, and finally the innovation skills of the company will increase.

The implementation of criteria of a learning organization means the requirements of a knowledge management system. In this case, organizational members, individuals and groups, are open to acquiring new skills, to continuous renewal and learning (double-loop and the deuterio learning are of significance from the point of view of learning [3]). Such an organizational atmosphere supports the success of knowledge sharing, which means that everyone aims to transmit his/her knowledge, to share it with colleagues and the other members of the organization for the sake of the collective goals. This fact separately contributes to the fact that

people will be able to work and produce the expected results in the framework of a balanced organizational operation on a higher level of knowledge.

If the performance is higher, not only in quantity, but also in quality, it is to be seen in the efficiency of the enterprise, since it will make it possible to produce more modern, higher quality, more marketable products and services in competitive organizational conditions.

The establishment of conditions of knowledge within the company – first of all, innovative knowledge – is beyond the operation of a learning organization, such as the grounded internal knowledge base that is the condition for permanent development and renewal. As well as, the organizational atmosphere (culture) that establishes creativity, the conditions of continuous learning handle the requirement on a strategic level to the employee's satisfaction, thus a reliable quality and competitive performance of the work are ensured by putting the right person in the right place.

A learning organisation is an approach, a philosophy on one hand, and on the other hand it involves certain features connected to philosophy, which indirectly influence the success of the given organisation. The learning organisations therefore possess the emergence of the five principles, which are not characteristics of other organisations [29].

As a consequence of *thinking in a system* during the process of concentrating on the changing process, which appears as a constant demand, people concentrate on revealing the reason-cause correspondence hidden in the background of the problems and on the holistic examination of organisation involving its surrounding world instead of concentrating on the “here and now” solutions. The concept of *guiding ourselves* means the fact that people are able to learn independently, they possess a view of the future, which ensures them to be able to prioritize among the tasks. They are able to concentrate their creativity in order to reach their personal, individual goals and as a consequence of this to reach all the goals of the organisation (individual learning – corporate knowledge). The *samples of thoughts* primarily influence our attitude, often unconsciously. They influence our activities and our way to react to things. By making these samples conscious the members of the learning organisation can help – primarily with using one of the techniques applied in *group communities* – our ability and willingness to changes and our real activities. Common *prospects for the future* have to be established if we want that the goals set by certain individuals should contribute to the success of the organisation in the long-run. Its feature is that it involves the individual ideas and consequently the organisational groups and members will be able to identify themselves with it.

In case of the existence of preconditions, the application of a model is practical, which involves all the steps of knowledge management processes and can be applied easily in practical life. A possible model can be seen on the following Figure 2 [26].

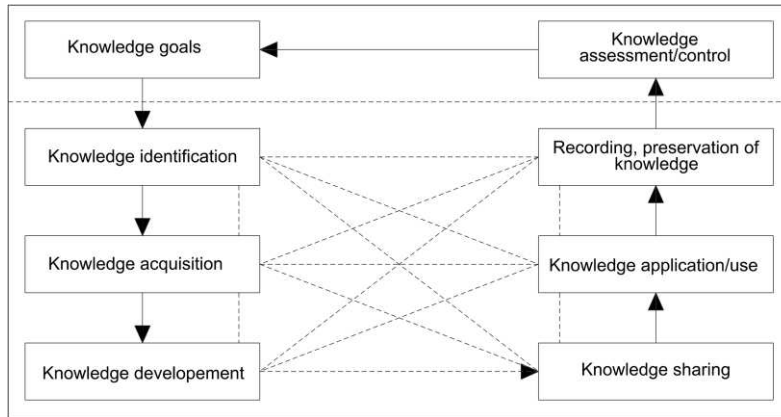


Figure 2
Probst's knowledge management model [26]

Leaders have to be closer to the way of thinking that a condition of competitiveness is the employee satisfaction and that the satisfied employees' willingness to learn and change can contribute to the possibility that the knowledge should be utilized and distributed in the organization at a higher level. To accept this thinking we have to turn to the eight elements of the Probst model.

These elements have a logical order and leaders have to know the steps to build a KM system. Theoretical knowledge and interactions among colleagues (communication) create real knowledge which can be utilized by an organization in order to accomplish its competitiveness. A chance to reach personal and organizational desires is the most important guarantee of the pursuit of knowledge, which is a precondition of learning and satisfaction at the same time. To sum up, the elements of the Probst model again:

Determination of Knowledge Goals:

- normative goals: to create an organizational culture, a joyful workplace and atmosphere which supports the competitiveness;
- strategic goals: the existing organizational knowledge helps to realize the strategic goals and with this knowledge we can work out a new, competitive strategy on a higher level;
- operational goals: they support to apply the knowledge management, if the normative and strategic goals are concrete enough (make sure that normative and strategic knowledge goals will be translated into action).

Knowledge Identification:

View of the internal abilities and possessed knowledge. The once applied and experienced things do not have to be discovered again and again. It is suitable to do a comparison with the environment, to use benchmark, then the organizational memory has a significant role to call out previous experience.

Knowledge Acquisition:

It means to learn from competitors, rivalry, stakeholders, buying know-how, unknown knowledge integration, to apply every cultured method of knowledge acquisition. (Dishonest tools which are important from the perspective of competition, that cannot be made consistent with employee satisfaction, do not belong to this step of KM system, for example, company acquisition, stealing knowledge, industrial espionage, etc.).

Knowledge Development:

New abilities, ideas, develop more effective technology, and new knowledge collection. (Obstacles and supporting elements which are in close correlation with employee satisfaction have to be identified.) It is a direct fact to influence competitiveness.

Knowledge Distribution:

Separated knowledge has to be integrated into the whole organization. Who, what, when, how, has to be known. (It is a prerequisite of satisfaction and a result at the same time. Its base is communication and a knowledge sharing organizational culture in each case.)

Knowledge Use:

Possessing knowledge does not bring a result, but to use it, does. This is in connection with an ability of knowledge acquisition, with willingness to learn and change, which is in interaction with satisfaction. To realize this step the previous phase must be accomplished. (This is also true in the case of internal and external knowledge.)

Knowledge Preservation:

Knowledge has not existed forever. In preservation and in forgetting knowledge the organizational memory has a significant role. This is very important as well from the point of view of creation of a balance between learning – forgetting, knowledge losing processes – and fluctuation. People who join an organization bring with them their knowledge. This can be anticipated from the development of knowledge bases and knowledge stores. With the increase of satisfaction the chance to preserve work forces and knowledge will grow. (Knowledge maps, knowledge catalogues)

Knowledge Measurement:

In the estimation of competitiveness the indexes have significant roles which are suitable for measuring normative, strategic and operational goals. In the stage of determination of knowledge goals the possibilities of success evaluation have to be fixed. (Market conditions, competitors, sales, etc.)

The elements have to be handled from a system view and reviewing connections with them is very important.

It can be stated, together with the heading 'self-management', which appears among the operating conditions of a learning organizational culture, that people, in such organizational conditions, who get a chance to accomplish their own goals, perform higher quality work by harmonizing their own imagination with the organizational goals, something that is a determining factor of employee satisfaction at the same time. On the other hand, qualitative products influence customer satisfaction as a consequence of competitive operations.

Examining the activities of enterprises with regard to quality, researchers face two kinds of problems:

- Failures adjustable by the employee, which arises if the employee has all the three criteria of self-control:
 - Knows what his/her task is.
 - Knows what he/she is doing now.
 - Is able to control his/her activity.

These criteria are only the conditions of self-control: they do not mean self-control has automatically been achieved. For self-control to be achieved an employee needs to possess the appropriate approach and responsibilities, and also know and want to use these tools.

- System failures, namely, the failures that are influenced by the management of the enterprise, occur if one or more of the criteria of self-control are not fulfilled. This indicates imperfection in creating the conditions for self-control by management, if the employees do not know the expectations, the mode of actions, and the possible tools, and if they do not get feedback on their performances and the appropriate support, either physically or in terms of human resources.

In the case of striving for quality work (the competitiveness criterion), it is about the change of management's point of view, achieving long-term thinking. This brings about a new 'lifestyle', a modified behaviour and a new emphasis on the life of the enterprises. This means that the increase in general effectiveness and value creation comes to the forefront instead of permanent cost minimization over a short-term period. Enterprises want to produce at a cheaper and cheaper price, they try to offer more and more to meet customer needs. But, these companies do not want to give more or something else that the customers want. This way of thinking is to be seen as a concept and a program wherein the implementation and the 'steps' of change have the same significance as the result itself. This is a permanent learning process, meaning the formation of a learning organization that supports continuous quality improvement.

3 Results

3.1 Quality Culture

In order to understand the relationship also from the side of quality conceived in the title of the study, first the meaning of quality culture has to be cleared. It also has more definitions in professional literature, among which we have chosen one. The definition describes the train of thoughts, which can be used as good example in order to show the relationship with the learning organisational culture establishing the base of the operation of knowledge management. According to this:

Quality culture is a corporate environment, where a certain approach, behaviour and attitude are prevailed, which is accepted by all the participants and which makes everybody be responsible for quality. It is an approach, which can be characterized by the ambition to be excellent, by continuous quality development and by taking other concerned parties' demand into account. For the sake of successfulness such a quality culture has to be made tangible by a proper quality management system, which is able to maintain the monitoring and evaluation processes and results of the organisation.

A proper cultural background is the precondition of developing quality systems and of the operation of quality management systems. The proper culture is also a requirement for developing knowledge management system, which is ensured by the above described learning organisation solution. To prove the similarities between learning organisational culture and quality culture, the following comparison should be observed. In both cases trust is the basic requirement and in both cases the human-centred approach prevails.

The base of quality management system is trust, its precondition is: quality culture

- Approach, behaviour
 - Process-centeredness
 - System-approached guidance
 - Continuous development
 - Involvement of workers - Teamwork
- } human-centred

The base of building knowledge management system is trust, its precondition is: learning organisational culture

- System-approach
 - Self-development, self-guidance
 - Common prospects for future
- } human-centred

- Inner persuasion, sample of thoughts
- Learning in group

The success of knowledge management systems depends on the way of thinking of the top management; it depends on how much they understand that knowledge management can help reach corporate goals. If quality is important, knowledge is overvalued, [15] sharing, preserving and developing knowledge is necessary, which means building knowledge management system is necessary. Thus, the result of the operation of the two systems is a fact presupposing each other mutually.

Choosing some systems of the quality systems, the following Table 1 presents the features, which justify their human-centred approach.

Table 1
Human-centred features of quality systems

Quality management systems			
ISO	TQM	EFQM	LEAN
<i>Process-centred</i>	<i>Proper management</i>	<i>Management</i>	Effectiveness
<i>System-approach</i>	Target in focus	Strategy	Flexibility
Effectiveness	Active participation	Integration of outer partners	Elimination of waste
<i>Continuous development</i>	<i>Continuous quality development</i>	<i>Teamwork</i>	<i>Flexible, skilled workers</i>
Satisfaction of concerned parties	Outer partner-relationship	<i>Human-centred</i>	Operation without mistakes
Customer-centeredness	Effective resource management	<i>Resource management (knowledge database, knowledge capital, etc.)</i>	Kaizen
	<i>Human-centered</i>	Corporate self-evaluation	Added value
	<i>Cultivating and developing corporate culture</i>		Perfection
	<i>Teamwork</i>		Cost-management
	Social-level learning		<i>Teamwork</i>

Using the features, the list containing the features of quality management systems and knowledge management systems can be put side by side, where the same colours show the relationship conceived in the title and in the introduction (Table 2)

Table 2
Features characterizing the similarities between quality management systems and knowledge management systems

Quality management systems	KMS
Customer-centred	<i>System-approach</i>
<i>Determining role of management</i>	<i>Built on trust</i>
<i>Integrating workers</i>	<i>Human-centred</i>
Process-centred approach	Knowledge is basic requirement
<i>System-centred managerial approach</i>	<i>Own set of methods, tools</i>
<i>Continuous development</i>	<i>Strategic role</i>
Factual approach of decision-making	<i>Struggle for acceptance</i>
Suppliers' connections in favour of mutual profit	<i>Formal application</i>
<i>Struggle for acceptance</i>	<i>Determining role of management</i>
<i>Formal application</i>	<i>Continuous development</i>
<i>Trust is basic requirement</i>	
<i>Own set of methods, tools</i>	
<i>Strategic role</i>	

The operation of knowledge management system, the applied set of tools – can be accepted based on the above evidence –, how the logic of knowledge management supports the operation of quality systems and the success of developing quality culture. The further features strengthening cooperation, which characterize the operation of both the learning organisational culture and of the quality culture is summarized in the followings. The elements of similarities between the systems:

- Trust,
- Key-role of humans - workers, management (Who?),
- System-approach (technology, processes, supporting services) (How?),
- Knowledge is basic expectation (What?),
- Own set of methods, set of tools – (overlaps),
- Modelled way of seeing things,
- Strategic role,
- A run „walk of life“ – struggle for accepting reason for existence,
- Formal application – unexploited opportunities.

Based on the above shown logical frame, it can clearly be seen that the steps of knowledge management system can be used according to the following correspondence in order to solve the problems described in the beginning of the

study (Table 3). Consequently, the help of which the quality-centred way of thinking and behaviour, the development and maintenance of quality culture can be supported without any doubt.

Table 3
KM solution supporting solving quality problems

Quality problems	KM solutions
Repetition of similar mistakes (costs),	Preserving, recording, sharing knowledge
Duplicated fulfilling of tasks (projects form the past, their results are known)	Preserving, recording knowledge
Lack of information (e.g.: customer service),	Knowledge share – building on trust
Lack of sharing inner good ideas, best practices	Sharing, gaining and developing knowledge
Weak link (1-2 people have the key knowledge/information)	Key people – knowledge share – building trust
Slow/ no integration of gained knowledge (slow improvement of products/rival overtakes),	Utilization of knowledge
Information/knowledge sources are hardly available (frustrated workers).	Developing knowledge-database and corporate memory

4 Discussion and Conclusion

As a summary it can be stated that the solution of quality problems and the logic of knowledge management live in symbioses. The quality issues can hardly be solved without knowledge and operating knowledge management system in order to fulfill strategic goals is worthwhile only by taking quality efforts into consideration.

Quality culture is built on strengths, where the features are:

- challenging individual and common goals,
- real conditions,
- clear, clean-cut rules,
- relations built on trust,
- open communication.

In such a corporate environment, the followings can dominate:

- positive atmosphere to reach goals,
- optimism and desire to do something,
- knowledge share,
- mutual help,
- goal-orientation and self-confidence,
- focusing on the task.

We have to realize this situation in organizations. It is confirmed by the idea of Donna Denehy as well.

„The marriage of KM and QM is something that many organizations are now thinking about. Integrating the two can begin to:

- Improve the baseline knowledge of your representatives and understand where knowledge gaps exist,
- Positively impact employee morale and empowerment,
- Provide customers more ease of doing business,
- Help ensure consistent information is being provided to help improve compliance,
- Drive consistent improvement in your quality results,
- Improve customer experience and satisfaction.

Now is the time to start thinking about bridging the gaps.”

References

- [1] Akdere, M. (2009) The Role of Knowledge Management in Quality Management Practices: Achieving Performance Excellence in Organizations. *Advances in Developing Human Resources* 11(3) 349-361, doi:10.1177/1523422309338575
- [2] Alimohammadlou, M. & Eslamloo, F. (2016) Relationship between Total Quality Management, knowledge Transfer and knowledge Diffusion in the academic settings, in *Procedia - Social and Behavioral Sciences* 3rd International Conference on New Challenges in Management and Organization: Organization and Leadership. 230. 104-111, Dubai: Bernard McKenna, Farzad Sattari Ardabili, and Nezameddin Faghieh
- [3] Argyris, C. & Schön, D. A. (1978) *Organizational Learning, A Theory of Action Perspective*. Boston: Addison-Westley

-
- [4] Baltus, R. (2001) Integrating Knowledge Management and Quality Management. In *Software Quality*. Springer. ed. M. Wieczorek and D. Meyerhoff, 107-125, Berlin: Springer
- [5] Bencsik, A. (2015) *A tudásmenedzsment elméletben és gyakorlatban*, Budapest: Akadémiai Kiadó
- [6] Berber, N. & Slavić, A. (2016) Human Resource (HR) Outsourcing in European Compensation Management in the Light of CRANET Research, *Acta Polytechnica Hungarica*, 13(3), 207-225
- [7] Bögel, Gy. (1999) *Tudásmenedzsment – a láthatatlan hatalom*. Budapest: Magyar Távközlés
- [8] Choi S. L., Kowang, T. O., Fei, G. Ch & Mang, H. P. (2016) Importance of Knowledge Management on Total Quality Management. A Review. *World Applied Sciences Journal* 34 (12): 1829-1833, doi: 10.5829/idosi.wasj.2016.1829.1833
- [9] Davenport, T. H. & Prusak, L. (2001) *Tudásmenedzsment* Budapest: Kossuth Kiadó
- [10] Denehy, D. (2016) *The Marriage of Knowledge Management and Quality Management*, Accessed April 14, 2016, <http://blog.verint.com/customer-engagement/the-marriage-of-knowledge-management-and-quality-management>
- [11] Dougherty, V. (1999) Knowledge is about people, not databases *Industrial and Commercial Training* 31(7) 262-266
- [12] Duran, C., Çetindere, A. & Şahan, Ö. (2013) An analysis on the relationship between total quality management practices and knowledge management: The case of Eskişehir, In *Procedia - Social and Behavioral Sciences 2nd World Conference on Business, Economics and Management - WCBEM 2013*, Vol. 106, 65-77 Amsterdam: Elsevier
- [13] Fehér, P. (2002) *Tudásmenedzsment: Problémák és veszélyek*. *Vezetéstudomány* 33(4) 36-45
- [14] Fehér, P. (2005) A technológiák szerepe a tudásmenedzsment folyamatok támogatásában *Vezetéstudomány*, 36(4) 11-22
- [15] Gyulay, T. (2017) *Minőség és tudásmenedzsment*. Paper presented at “Tudásmenedzsment Műhely” Budapest November 04
- [16] Hofstede, G. & Hofstede, G. J. & Minkov, M. (2010) *Cultures and Organizations: Software of the Mind*. New York: McGraw-Hill
- [17] Honarpour, A., Jusoh, A., & Nor, K. M. (2017) Total quality management, knowledge management, and innovation: an empirical study in R&D units. *Total Quality Management & Business Excellence*. Accessed November 12, 2017, <https://doi.org/10.1080/14783363.2016.1238760>

- [18] Hype Cycle elemzési keret Gartner Group 1995 Accessed Oktober 12 (2017) <http://www.gartner.com/technology/research/methodologies/hype-cycle.jsp>
- [19] Kahreh, Z. S., Shirmohammadi, A., & Kahreh, M. S. (2014) Explanatory Study Towards Analysis the Relationship between Total Quality Management and Knowledge Management, In *Procedia - Social and Behavioral Sciences 2nd World Conference on Business, Economics and Management - WCBEM 2013*, Vol. 109, 600-604, Amsterdam: Elsevier
- [20] Kolnhofer-Derecskei, A., Reicher, R. Zs., & Szeghegyi, A. (2016) The X and Y Generations' Characteristics Comparison, *Acta Polytechnica Hungarica*, 14(8) 107-125
- [21] Lazányi, K. (2015) A családi vállalkozások és a tudásmenedzsment, Taylor: *Gazdálkodás- és Szervezéstudományi folyóirat: A virtuális intézet Közép-Európa kutatására közleményei* 7(3-4) 254-260
- [22] Lazányi, K. (2010) Tudásmenedzsmenttel a vidékfejlesztésért, In: *Lifelong Learning Magyarország Alapítvány (szerk.) Tanulás, tudás, gazdasági sikerek avagy a tudásmenedzsment szerepe a gazdaság eredményességében: tudományos konferencia kiadványa: Győr, 2010. április 14. Konferencia helye, ideje: Győr, Magyarország, Budapest: Lifelong Learning Magyarország Alapítvány*, 407-411
- [23] Mura, L., Ključnikov, A., Tvaronavičienė, M., & Androniceanu, A. (2016) Development Trends in Human Resource Management in Small and Medium Enterprises in the Visegrad Group, 14(7) 105-122
- [24] Noszkay, E. (2007) Tudás és menedzsment (Tendenciák és jelenségek a tudásmenedzsment hazai alkalmazásai alapján) Miskolci Egyetem, *Gazdaságtudományi Kar VI. Nemzetközi Konferencia „A közgazdász képzés megkezdésének 20. évfordulója alkalmából” Konferencia Kötet* 120 – 127. Miskolc – Lillafüred
- [25] Poor, J., Vinogradov, Sz., Gábríelné, Gy., Antalik, I., Horbulák, Zs., Juhász, T., Kovács, I. É., Némethy, K., & Machová, R. (2017) Atypical Forms of Employment on Hungarian-Slovakian Border Areas in Light of Empirical Researches, *Acta Polytechnica Hungarica*, 14(7), 123-141
- [26] Probst, G., S. Raub, S., & Romhardt, K. (2006) *Wissen Managen, Wie Unternehmen ihre wertvollste Ressource optimal nutzen.* Wiesbaden:Gabler GmbH
- [27] Reddy, T. R. (2012) Total quality management and knowledge management integrations in library and information centers: a study. *Journal of Research in International Business and Management* Vol. 2(11) 292-298, Accessed November 21, 2017, <http://www.interestjournals.org/JRIBM>

- [28] Sándori, Zs. (2001) Mi a tudásmenedzsment? Magyar Elektronikus Könyvtár Accessed December 11 2017, <http://mek.oszk.hu/03100/03145/html/km4.htm>
- [29] Senge, P. M. (1998) Az 5. alapelv. A tanuló szervezet kialakításának elmélete és gyakorlata. Budapest: HVG Kiadó
- [30] Stewart, T. A. (1991) Brainpower Intellectual capital is becoming corporate America's most valuable asset and can be its sharpest competitive weapon. Fortune Magazin, Accessed July 08 2017, http://archive.fortune.com/magazines/fortune/fortune_archive/1991/06/03/75096/index.htm
- [31] Szabó, K. (2017) Tudásmenedzsment és kiválóság, Paper presented at "Tudásmenedzsment Műhely" Budapest November 04
- [32] Zakeri, S., Goudarzi, H. T., Atamanesh, H. & Koochaki, H. (2014) Total Quality Management and Knowledge Management. A Researches Review, Technical Journal of Engineering and Applied Sciences. Accessed October 13, 2017, www.tjeas.com

Regional Disparities of Small and Medium Enterprises in Slovakia

Jarmila Lazíková¹, Anna Bandlerová¹, Oľga Roháčiková², Pavol Schwarcz³, Ľubica Rumanovská³

¹ Slovak University of Agriculture in Nitra, Department of Law, Tr. A. hlinku 2, 949 76 Nitra, Slovakia, jarmila.lazikova@uniag.sk, anna.bandlerova@uniag.sk

² Slovak University of Agriculture in Nitra, Department of Public Administration, Tr. A. hlinku 2, 949 76 Nitra, Slovakia, olga.rohacikova@uniag.sk

³ Slovak University of Agriculture in Nitra, Department of EU Policies, Tr. A. hlinku 2, 949 76 Nitra, Slovakia, pavol.schwarcz@uniag.sk, lubica.rumanovska@uniag.sk

Abstract: Small and medium enterprises are mostly considered as key elements of a market economy. Their share is 99% of a total number of all enterprises in Slovakia on average. The aim of the paper is to identify regional disparities of the SMEs development in Slovakia. As methods, we used chain indexes to compare the changes among numbers of small, medium and large enterprises, time series analysis, non-parametric method for investigation of the statistically significant differences and correlation analysis. According to the results we expect increasing numbers of SMEs in Slovakia in the next four years, mainly an increasing of small enterprises. The development of SMEs is very different in particular regions of Slovakia. The number of enterprises in remote rural areas grow less rapidly than the number of those in more accessible rural areas of Slovakia. The strong correlation between SMEs and large enterprises indicates suitable conditions for doing business in rural areas for SMEs as well as for large enterprises. The support policy of SMEs should be more intensive, especially in the rural areas that are not suitable for a large scale business.

Keywords: small and medium enterprises; regional disparities; forecasting; rural areas

1 Introduction

Micro, small and medium enterprises are the engine of the European economy [14, 32, 48]. Small and medium enterprises (SMEs) as a part of each market – oriented economy add flexibility, competitiveness and innovation activities into the market environment and contribute to the overall regional development [37].

They are able to fulfill the gaps in the market which cannot be covered by large enterprises due to their robustness [19]. SMEs are an important element of market economy, job opportunities, added value or foreign trade [17, 41] and essential presumption for the stable progress of a country [20]. They are an essential source of jobs, create entrepreneurial spirit and innovation in the EU and are thus crucial for fostering competitiveness and employment [28, 29, 36]. SMEs are still an issue that is interesting to study because it is recognized that small enterprises have a major role in the employment and contribution to the gross domestic product [35, 45]. Therefore, the support of SMEs development is one of the political priorities of countries and supranational organisations (such as the European Union). It is, probably, the reason why many scientific papers, studies and reports are oriented on the factors that improve the success of the SMEs in the market [13, 24, 46, 47]. There are usually identified external and internal factors with the regard to the success of SMEs. As for the external factors, there are studied macroeconomic, political, legal, social, technological and demographic factors and competitive environment, as well [3, 4, 5, 10, 21, 31]. As for the internal factors, there are considered business experiences and the motivation of the owners/managers will be able to manage the organization [35], knowledge management [15], firm size and age of the business [25, 27], marketing of their products, qualification of employees in the marketing department, finance to undertake marketing research [34], lack of infrastructure [8], lack of innovativeness [18, 35, 39, 42], location and human capital [11, 25]. Storey [37] suggests that the location of a small business is a factor, which influences its performance because the bulk of sales of small enterprises are too highly localised markets. Dahlqvist, et al. [12] added that the geographic area, where a firm is located, has implications for its access to markets and resources such as: finance, skilled labour, subcontractors, infrastructure, and other facilities. Enterprises located in urban and commercial areas were more likely to survive, during a given year, than those located in rural areas [25]. However, Keeble [23] suggests that, whilst on balance rural enterprises may grow more rapidly than their urban counterparts; enterprises in remote rural areas in the United Kingdom grow less rapidly than those in more accessible rural areas.

2 Material and Methods

The aim of the paper is to evaluate the impact of the location on the SMEs development in Slovakia and to confirm or refuse the above mentioned findings on location of SMEs in the Slovak environment. Slovakia is a rural country, only the Bratislava region is considered as urban area, all other regions are classified as rural or semi-rural counties. For this purpose we tried to identify the most attractive regions for SMEs. Moreover, we compare the trend development in the regions with the average trend given by the whole country as well. For this

purpose we need to identify the SMEs; unfortunately, there is no single definition and the criteria are different in various political and legal documents. So, firstly, we need to define SMEs for the purpose of this paper. Therefore, the paper is organised as follows: The first chapter introduces the various notions of SMEs and explains what is considered as a SME for our further analysis. The second chapter describes the current state and the development of SMEs during the period of 1996-2015 in Slovakia and provides a forecasting for the next four years. The third chapter identifies the most attractive regions for SMEs and regional disparities of the SME development in urban, rural and semi-rural areas using the time series and cross-sectional data as well. The last chapter provides discussion on the results of the above-mentioned issues.

The data about SMEs was received from the published data of the Statistical Office of the Slovak Republic for particular regions (NUTS III) and counties (LAU 1) for a period of 1996-2015 (the starting year is the year of the new administrative zoning of Slovakia that is used in the analysis, the last year is the year of the newest data at the time of the paper submission).

As methods, we used chain indexes to compare the changes among numbers of small, medium and large enterprises, time series analysis to provide forecasts by Statistical Analytical System (SAS), non-parametric method for investigation of the statistically significant differences and correlation analysis.

Chain index is an index number, in which, the value of any given period is related to the value of its immediately preceding period as described Pacáková [31].

For non-parametric testing, the Kruskal-Wallis test was used characterised as follows:

$$H = \left(\frac{12}{N(N+1)} \sum_{j=1}^k \frac{R_j^2}{n_j} \right) - 3(N+1) \quad (1)$$

H – Kruskal – Wallis test characteristics

N – total number of counties (all regions combined)

R_j – rank total for each region

n_j – number of counties in each above mentioned region

k – number of regions

For time series analysis, we used the linear trend with auto-regressive errors for forecasting of SMEs, the combination of 4 models (log linear trend with auto-regressive errors, Winters method additive and multiplicative, and linear (Holt) exponential smoothing) for forecasting of small enterprises; and combination of two models (linear trend with auto-regressive errors and log linear trend with auto-regressive errors) for forecasting of medium enterprises. The models are described in various publications [2, 9, 37].

Linear trend with auto-regressive errors calculated as follows:

$$y_t = b_0 + b_1 t + \varepsilon_t \quad (2)$$

$$\varepsilon_t = \delta \cdot x_{t-1} + u_t \quad (3)$$

where $b = (b_0, b_1)$ is a vector parameter and $\{\varepsilon_t\}$ represents the auto-regressive errors.

Log – linear trend with auto-regressive errors, in which the dependent variable changes at an exponential rate over time or constant growth at a particular rate calculated as follows:

$$\ln(y_t) = b_0 + b_1 t + \varepsilon_t \quad (4)$$

$$\varepsilon_t = \delta \cdot x_{t-1} + u_t \quad (5)$$

where $b = (b_0, b_1)$ is a vector parameter and $\{\varepsilon_t\}$ represents the auto-regressive errors.

Linear (Holt) exponential smoothing calculated as follows:

$$\hat{y}_{t+1} = u_t + v_t \quad (6)$$

$$u_t = \alpha \cdot y_t + (1 - \alpha)(u_{t-1} + v_{t-1}) \quad (7)$$

$$v_t = \beta \cdot (u_t - u_{t-1}) + (1 - \beta)v_{t-1} \quad (8)$$

$$u_1 = y_1 \quad (9)$$

$$v_1 = 0 \quad (10)$$

$$0 < \alpha \leq 1 \quad (11)$$

$$0 \leq \beta \leq 1 \quad (12)$$

Winters method additive calculated as follows:

$$y_t = (\beta_0 + \beta_1 t) + s_t + \varepsilon_t \quad (13)$$

Winters method multiplicative calculated as follows:

$$y_t = (\beta_0 + \beta_1 t) \times s_t \times \varepsilon_t \quad (14)$$

where s_t is seasonal pattern and ε_t is irregular component.

To quantify the association between the small, medium and large enterprises, we used the correlation analysis. We used the Pearson correlation coefficient calculated as follows:

$$r_{xy} = \frac{\text{COV}(x, y)}{\sqrt{s_x^2 s_y^2}} \quad (15)$$

where $\text{cov}(x, y)$ is covariance of two variables in a data set and s_x^2 , s_y^2 are variances of x and y as described Pacáková [31].

3 Results

3.1 Notion of SMEs

There is still no universally accepted definition what small and medium enterprises are. In scientific papers, international legal binding or non-binding documents or political documents, many definitions of SMEs are included, but they differ from one another. We can find the definitions of SMEs based on two approaches; qualitative and quantitative ones. Bolton report [7] defines three qualitative criteria of SMEs: management of firm by its owner(s) in a personalized manner, relatively small share of the market in economic terms, independence in the sense that it does not form a part of a larger enterprise is relatively free from outside control in its principal decisions. Marwede [27] regards legal form, the role of the firm owner, the firm's position on the market, organizational structure and economic and legal autonomy. Loecher [26] deals with the qualitative measures such as personal principle, unity of leadership and capital. Despite the volume of SME definitions, there is a tendency to accept quantitative criteria, first and foremost the headcount or employee number criterion as the main determinant in categorizing SMEs [6]. Ardic, Mylenko and Saltane [1] confirm in their cross-country analysis that the most common definitions used by regulators are based on the number of employees, sales and/or loan size. The most common among the three is the number-of-employees criterion. Within the World Bank Group, IFC and MIGA have official definitions but also define SMEs in other ways. IFC and MIGA formally define SMEs as fulfilling two of three criteria: (1) having more than 10 and fewer than 300 employees; (2) having between 100 000 and 15 million dollars in sales; (3) having between 100 000 and 15 million dollars in assets [40]. The enterprises under the above-mentioned minimum level are considered as micro enterprises. The SME definition has been developing also in the legal acts of the European Union. The first definition was incorporated in the article 11 of Fourth Council Directive 78/660/EEC on the annual accounts of certain types of companies.¹ It permits some exemptions from the detailed annual accounts for companies, which on their balance sheet dates do not exceed the limits of two of the three following criteria: (1) balance sheet total; (2) net

¹ Fourth Council Directive 78/660/EEC of 25 July 1978 based on Article 54 (3) (g) of the Treaty on the annual accounts of certain types of companies (OJ L 222, 14.8.1978, pp. 11-31)

turnover; (3) average number of employees during the financial year. The first two criteria were changed five times but the number of employees is stable over the time. The overview is presented in the Table 1.

Table 1
Overview of the changes in the financial criteria of SME definition in EU directives

	Directive 78/660/EC	Directive 94/8/EC ²	Directive 1999/60/EC ³	Directive 2003/38/EC ⁴	Directive 2006/46/EC ⁵ and Directive 2012/6/EU ⁶	Directive 2013/34/EC ⁷
Micro-enterprises						
Balance sheet total EUR	-	-	-	-	350 000	350 000
Net turnover EUR	-	-	-	-	700 000	700 000
Average number of employees per year	-	-	-	-	10	10

² Council Directive 94/8/EC of 21 March 1994 amending Directive 78/660/EEC as regards the revision of amounts expressed in ecus (OJ L 82, 25.3.1994, pp. 33-34)

³ Council Directive 1999/60/EC of 17 June 1999 amending Directive 78/660/EEC as regards to amounts expressed in ecus (OJ L 162, 26.6.1999, pp. 65-66)

⁴ Council Directive 2003/38/EC of 13 May 2003 amending Directive 78/660/EEC on the annual accounts of certain types of companies as regards to amounts expressed in euro (OJ L 120, 15.5.2003, pp. 22-23)

⁵ Directive 2006/46/EC of the European Parliament and of the Council of 14 June 2006 amending Council Directives 78/660/EEC on the annual accounts of certain types of companies, 83/349/EEC on consolidated accounts, 86/635/EEC on the annual accounts and consolidated accounts of banks and other financial institutions and 91/674/EEC on the annual accounts and consolidated accounts of insurance enterprises (Text with EEA relevance) (OJ L 224, 16.8.2006, pp. 1-7)

⁶ Directive 2012/6/EU of the European Parliament and of the Council of 14 March 2012 amending Council Directive 78/660/EEC on the annual accounts of certain types of companies as regards micro-entities (OJ L 81, 21.3.2012, pp. 3-6)

⁷ Directive 2013/34/EU of the European Parliament and of the Council of 26 June 2013 on the annual financial statements, consolidated financial statements and related reports of certain types of enterprises, amending Directive 2006/43/EC of the European Parliament and of the Council and repealing Council Directives 78/660/EEC and 83/349/EEC Text with EEA relevance (OJ L 182, 29.6.2013, pp. 19-76)

Small enterprises						
Balance sheet total EUR (EUA, ⁸ ECU ⁹)	1 000 000	2 500000	3 125 000	3 650 000	4 400 000	4 000 000
Net turnover EUR (EUA, ECU)	2 000 000	5000 000	6 250 000	7 300 000	8 800 000	8 000 000
Average number of employees per year	50	50	50	50	50	50
Medium enterprises						
Balance sheet total EUR (EUA, ECU)	4 000 000	10000000	12 500 000	14 600 000	17 500 000	20 000 000
Net turnover EUR (EUA, ECU)	8 000 000	20 000 000	25 000 000	29 200 000	35 000 000	40 000 000
Average number of employees per year	250	250	250	250	250	250

This definition is used only for the purpose of this directive on the annual accounts of certain types of companies and its amendments. Regardless of this limited use, the changes were very often. Moreover, the European Commission adopted two recommendations that define SMEs. The indicators are the same as in the above-mentioned directives but the highest levels were changed again. In 1996, the recommendation of EC¹⁰ established the first common SME definition mainly for the purposes of the implementation of various Community policies. The definition could be used in general for various purposes, but a recommendation compared to a directive is not a legally binding act. It is binding

⁸ according to the Fourth Council Directive 78/660/EEC of 25 July 1978 based on Article 54 (3) (g) of the Treaty on the annual accounts of certain types of companies (OJ L 222, 14.8.1978, pp. 11-31)

⁹ according to the Council Directive 94/8/EC of 21 March 1994 amending Directive 78/660/EEC as regards to the revision of amounts expressed in ecus (OJ L 82, 25.3.1994, pp. 33-34)

¹⁰ Commission Recommendation 96/280/EC of 3 April 1996 concerning the definition of small and medium enterprises (Text with EEA relevance) (Official Journal L 107, pp. 4-9)

for the European institutions, but it is only voluntary for individual Member States. According to the recommendation a small enterprise has fewer than 50 employees and has either, an annual turnover not exceeding ECU 7 million, or an annual balance-sheet total not exceeding ECU 5 million. A medium enterprise has fewer than 250 employees, and either, an annual turnover not exceeding ECU 40 million, or an annual balance-sheet total not exceeding ECU 27 million. In 2003, the European Commission adopted a new recommendation¹¹ because of the need to adapt it to economic developments. It entered into force on January 1, 2005 and applies to all EU policies, programmes and measures for SMEs. Article 2 of the Annex of this recommendation defines a microenterprise as an enterprise which employs fewer than 10 people and an annual turnover and/or annual balance sheet total of which does not exceed EUR 2 million; a small enterprise as an enterprise which employs fewer than 50 people and an annual turnover and/or annual balance sheet total of which does not exceed EUR 10 million; a medium enterprise is made up of enterprises which employ fewer than 250 people with an annual turnover not exceeding EUR 50 million, and/or an annual balance sheet total not exceeding EUR 43 million. According to the overview of the above-mentioned changes only within the EU legal acts we can state that the number of employees is the most stable criterion. However, Curran and Blackburn [11] point out that the definition of SMEs by number of employees has become difficult due to part-time work, casual work or temporary work becoming more widely used by employers. Gibson and Vaart [16] consider the criterion of turnover as the most consistent of the three quantitative criteria. On the other hand, the financial criteria are changed relatively often and then, it is impossible to use permanently changing statistical data for a mathematical-statistical analysis mainly for the time series analysis. Precisely, product of these definitions is the definition of SMEs legitimized by the European Union and, which is used by most of the researchers [6]. Therefore, we regard in further analysis the SMEs by number of employees as provided by the Statistical Office of Slovakia. Small enterprises are considered as enterprises with the number of employees within the range 0 to 49; medium enterprises within the range of employees 50 to 250 and large enterprises have 250 employees or more.

3.2 Development of SMEs in Slovakia

SMEs represent 99% of all enterprises in the European Union [14]. It is also the case of Slovakia, where the share of SMEs is on average 99.85% during the period of 1996-2015. Out of it, the share of small enterprises is on average 99.31% of a total number of SMEs and the share of medium enterprises in Slovakia is quite negligible (only 0.69% on average). The development of SMEs was more

¹¹Commission Recommendation of 6 May 2003 concerning the definition of micro, small and medium enterprises (Text with EEA relevance) (notified under document number C(2003) 1422) (Ú. v. EÚ L 124, 20.5.2003, s. 36-41)

intensive after the accession of Slovakia in the European Union in 2004. The number of SMEs was increasing when the economic crisis broke out. Since 2008, the number of SMEs is quite stable without important changes until today (Figure 1).

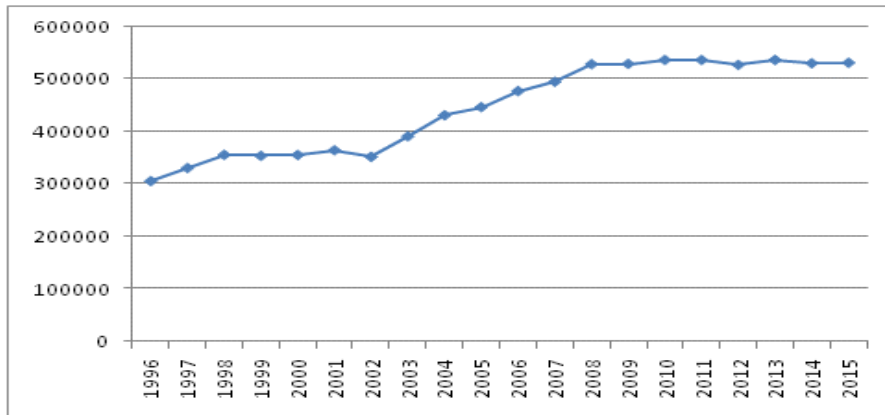


Figure 1

Development of SMEs 1996–2015 in Slovakia

The comparison of the development of small, medium-sized and large enterprises is only possible according to the chain indexes. The absolute numbers of enterprises are not comparable because of a high share (99%) of small enterprises. The chain indexes are documented in Figure 2. The impact of the economic crisis was reflected in 2009 and 2010 when the highest decreasing was recorded in the number of large enterprises. Small enterprises were the first to recover from the economic crisis. In 2010, they were increasing in number, but medium and large enterprises were still decreasing. The number of these categories of enterprises increased one year later. During the period of 1996–2015, the development of medium and large enterprises was very similar and the fluctuation was higher in numbers than in the number of small numbers. Small enterprises are more able to help in the stabilisation process during the economic recession.

A similar development of medium and large enterprises indicated long-term relation between them. However, no co-integration relations were confirmed (neither between the numbers of large and medium enterprises nor between the numbers of medium and small enterprises). There are no long-term balanced relations among the numbers of all three groups of enterprises. Based on the above mentioned results we can state that the macroeconomic, legal and political changes are able to influence the state and the development of the numbers of these three groups of enterprises in a very different way. The number of large enterprises is more sensitive to these changes than the number of SMEs, mainly the number of small enterprises.

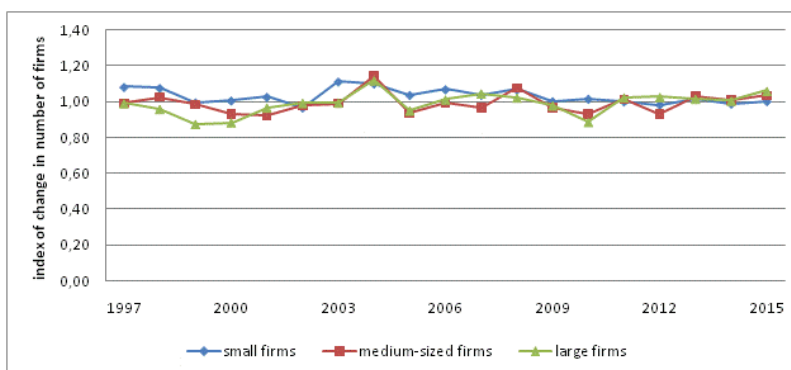


Figure 2

Chain indexes of change in number of enterprises according to their size

After the economic crisis the numbers of SMEs fluctuates around 530,000. The probability of a further trend is proved by the models of time series analysis that enable us to predict the development of the number of SMEs in the next four years. The forecast models were developed by the statistical analytical system (SAS) and the SAS Time Series Forecasting System was used to predict the development of SMEs in Slovakia, given the historical data of the absolute data of the number of SMEs in the period of 1996-2015. We chose three models that predict (1) development of number of SMEs together by the linear trend with auto-regressive errors; (2) development of small enterprises by the combination of 4 models (log linear trend with auto-regressive errors, winters method additive and multiplicative, and linear (Holt) exponential smoothing; (3) development of medium enterprises by the combination of two models (linear trend with auto-regressive errors and log linear trend with auto-regressive errors). The forecasting results are documented in Table 2.

Table 2

Forecast of development of the SMEs in Slovakia based on the historical data

	Year	2016	2017	2018	2019
1 st model	Predicted value	549377	573229	602475	631023
	Upper 95% confidence	580357	615491	653968	686162
	Lower 95% confidence	518396	530968	550983	575884
2 nd model	Predicted value	543790	560475	577394	594301
	Upper 95% confidence	562181	587218	610093	631652
	Lower 95% confidence	525400	533732	544695	556951
3 rd model	Predicted value	2684	2606	2515	2486
	Upper 95% confidence	2870	2799	2707	2686
	Lower 95% confidence	2498	2412	2324	2285

All three models were compared by Mean Absolute Percent Error (MAPE), R-Square, Akaike Information Criterion and Schwarz–Bayssian Information Criterion [9, 2, 37] and the best values of indicators were considered to choose particular model for forecasting of SMEs together and individually. The results are presented in Table 3.

Table 3
Selected indicators for evaluation of model's quality

Models	MAPE	R-Square	Akaike Criterion	Schwarz Bayssian Criterion
1 st model	2,661	0,975	391,912	397,886
2 nd model	3,066	0,966	393,594	397,576
3 rd model	3,495	0,730	195,191	197,182

MAPE criterion measures the size of the error in percentage terms. The model is acceptable if the MAPE criterion is less than 10. We chose the models with the smallest value of MAPE and all three models have MAPE of about 2-3%, which is acceptable for forecasting. The values of Akaike criterion and Schwarz–Bayssian criterion are useful when comparing more models. In this case, we chose the models with the lowest values for each forecasting. The first and second models have similar values because the small enterprises prevail in the number of SMEs. The R-square characteristic is more than 90% in the case of forecasting of number of SMEs together and forecasting of number of small enterprises. We prefer models according to the highest R-square and the smallest MAPE.

Based on the results of the first model (SMEs together) we can expect an increasing trend of numbers of SMEs in Slovakia. According to this model, the number of SMEs will increase by approx. 100,000 enterprises during the next four years. This model was selected as the best according to the above – mentioned criteria. We assume that it is very optimistic forecasting because of a relatively stable number of SMEs from 2008 until nowadays. Therefore, we separated the number of SMEs enterprises between the small and medium enterprises and did the forecasting again. The second model for small enterprises indicates an increasing trend of the number of small enterprises by approx. 60,000 enterprises. We assume that it is more realistic forecasting than the forecasting in the first model mainly when expecting a decreasing number of medium enterprises. The third model of forecasting of medium-sized enterprises will indicate a decrease by about 300 enterprises. Small enterprises are more adaptive when markets fail while the medium enterprises are more sensitive to political and economic changes, such as the actual migration crisis, preparation of the negotiation process between the EU and the Great Britain on the secession from the EU and the negotiation process between the EU and the USA on the trade agreement. We conclude that the number of SMEs together will indicate an increasing; however, this increasing will be probably a little bit smaller than the forecasting according to the first model.

3.3 Development of SMEs in Slovak Regions

The Slovak Republic has eight regions (NUTS III) with various levels of development and living standard. Therefore, we are interested in allocation of SMEs in particular regions of Slovakia. The most developed region of Slovakia is the Bratislava region. The number of SMEs is much higher in this region than in all other regions of Slovakia during the whole selected period of 2000-2015 (on average 47 SMEs per 1 km², minimum 33 per 1 km² in 2002 and maximum 60 per 1 km² in 2014). All other regions are more comparable considering the number of SMEs per 1 km². The view is provided by Figure 3.

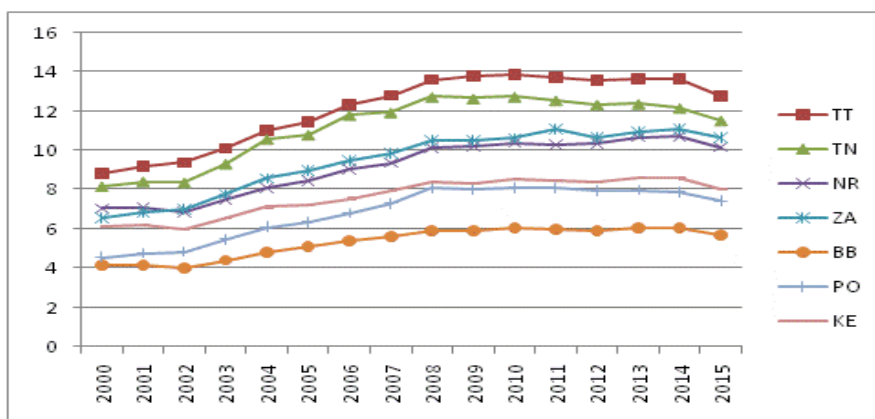


Figure 3

Number of SMEs per 1 km² during 2000-2015 in some regions of Slovakia¹²

The development of number of SMEs per 1 km² in particular regions has copied the development of number of SMEs in the whole country. The order of the regions has been retained during the period of 2000-2015. We can state that the SME enterprises prefer more developed regions when they decide on the location of their business. There are significant differences between the Bratislava region and all other regions. According to the Programme of rural development for the programming period of 2014-2020, only the Bratislava region is considered as urban region. Other regions of Slovakia are considered as rural (Nitra region, Banská Bystrica region, Prešov region and Trnava region) or semi-rural regions (Trenčín region, Žilina region, Košice region). It confirms the results of Liedholm (2002) that enterprises located in urban and commercial areas are more likely to survive during a given year than those located in rural areas or those being

¹² * TT – Trnava region, TN – Trenčín region, NR – Nitra region, ZA – Žilina region, BB – BanskáBystrica region, PO – Prešov region, KE – Košice region; BA – Bratislava region needs a separate figure due to data considered as outliers – its data are between 30-60 SMEs per 1 km²

operated out of home; urban and commercial location is also associated with faster growth, as measured by the number of employees hired in a given year. However, statistical differences among all other regions are not distinct. Therefore, we use the Kruskal-Wallis test to identify statistically significant differences among the Slovak regions. We used the data on the number of SMEs per 1 km² in particular counties (LAU1) of each region in 2015. The total number of observations is the number of counties (79) that are organized in 8 regions (NUTS 3). Due to small number of observations a non-parametric test (Kruskal-Wallis test) was used. Statistically significant differences were defined by the multiple range tests in Statgraphic. If we regard all 79 counties, the Kruskal-Wallis test confirms the statistical significance only between the Bratislava region and the rest of Slovakia; it is not possible to follow the potential statistical differences among other regions. Therefore, we left out the Bratislava region from the observation; the Kruskal-Wallis test confirms the statistical significance only between the Košice region and the rest of Slovakia. In spite of the fact that the Košice region was among the last three regions when comparing their development in the period of 2000-2015 (Figure 3), it was considered as the second best by the Kruskal-Wallis test. In the first case (development in 2000-2015), the number of SMEs per 1 km² was distributed on the whole area of the region and the best counties such as Košice I, Košice II, Košice III and Košice IV with the highest number of SMEs were down by the worst counties such as Rožňava and Sobrance with the smallest number of SMEs. In the second case, the number of SMEs per 1 km² was considered only for a particular county, so the best counties in the Košice region can use their impact on the results. If we want to consider the potential statistical significance in the regions of Slovakia, we need to leave out the outliers caused by the best counties from the Bratislava region (i.e. Bratislava I-V counties) and the Košice region (i.e. Košice I-IV counties). The number of observations was reduced to 70 counties.

Table 4
Differences of number of SMEs among regions of Slovakia

Region	Mean	Variance	p-value	K-W test statistic
Bratislava	17,00	68,65	0,00013	29,32
Trnava	12,79	12,19		
Trenčín	11,23	13,37		
Nitra	10,72	19,66		
Žilina	10,77	34,01		
Banská Bystrica	5,89	15,71		
Prešov	6,86	12,27		

Slovak regions; p-value is smaller than 0.05. In addition, according to the multiple range tests, there are statistically significant differences:

- between the Bratislava region and every other region (except the Trnava region and the Trenčín region);

- between the Trnava region and the Prešov, Košice and the Banská Bystrica regions;
- between the Trenčín region and the Prešov, Košice and the Banská Bystrica regions;
- between the Nitra region and the BanskáBystrica and the Košice regions;
- between the Žilina region and the Prešov, Košice and the Banská Bystrica regions.

Table 5
Multiple Range test results

Region	Count	Homogenous groups			
Košice	7	X			
Banská Bystrica	13	X			
Prešov	13	X	X		
Nitra	7		X	X	
Žilina	11			X	
Trenčín	9			X	X
Trnava	7			X	X

According to the above mentioned classification of rural, semi-rural and urban regions, we can state that there is no significant difference among semi-rural and rural regions regarding the number of SMEs. The counties of the Bratislava region which remain after excluding the most developed counties Bratislava I to V are not considered as urban area any more, only together with Bratislava I to V the data rank these counties as urban ones. We can state that the location is an important factor for SMEs development and urban areas are more appropriate for SMEs enterprises than rural and semi-rural regions. Finally, we regarded the numbers of small, medium and large enterprises in each of the 79 counties of Slovakia and found very strong correlation between the pairs of all three groups of enterprises (Table 6).

Table 6
Correlation matrix between particular group of enterprises

	Small enterprises	Medium enterprises	Large enterprises
Small enterprises	1		
Medium enterprises	0,895971011	1	
Large enterprises	0,834352581	0,949558998	1

If there are many small enterprises in a particular county, there is also higher number of medium or large enterprises. There is an extremely strong correlation between medium and large enterprises (0.95). We can suppose that the conditions for doing business in a county are suitable for SMEs as well as for large enterprises. The support policy of SMEs should be more intensive especially in the counties that are not very attractive for doing business either.

4 Discussion

There is still no universally accepted definition what small and medium enterprises are. The most usual criteria are the financial criteria of turnover, sales or assets and a number of employees. While the limits of the financial criteria are still being changed because of its adaptation to economic development, the number of employees is more stable during the period of time and so more suitable for statistical analysis of the SME development. Therefore, we were able to analyse the development of SMEs in Slovakia from 1996 to 2015. The share of SMEs is on average 99.85% during the period of 1996-2015. Out of it, the share of small enterprises is on average 99.31% of a total number of SMEs and the share of medium enterprises in Slovakia is quite negligible (only 0.69% on average). The development of SMEs during the economic crisis confirms the fact about a higher flexibility of small enterprises which increased in number in 2010 while the number of medium and large enterprises was still decreasing. During the period of 1996-2015, the situation of medium and large enterprises was developed in a very similar way, but it was different from the development of small enterprises. We assume that small enterprises are able to help the stabilisation process during the economic recession more than medium enterprises. It is a question if the support policy of SMEs should not be oriented only on the support of small enterprises because medium enterprises are more similar to large enterprises than to small ones. We had not found any long-term balanced relations among the numbers of all three groups of enterprises. Therefore, we suppose that the macroeconomic, legal and political changes are able to influence the state and development of the quantity of these three groups of enterprises in a very different way. The amount of large enterprises is more sensitive to these changes than the amount of SMEs, mainly the amount of small enterprises. In the future, we expect an increasing of numbers of SMEs in Slovakia (an increasing of small enterprises by app. 60 000 and a decreasing of medium enterprises by app. 300). The medium enterprises are more sensitive to the political and economic changes than the small ones.

The development of SMEs is very different in particular regions (NUTS III) of Slovakia. The most developed region of Slovakia is the Bratislava region. The number of SMEs is much higher in this region than in all other regions of Slovakia. We can state that the SMEs prefer more developed regions when deciding on the location of their business. There are significant differences between the Bratislava region (urban region) and all other regions (semi-rural and rural regions) of Slovakia. The result of Kruskal-Wallis test confirms statistically significant differences among the Slovak regions after excluding the urban areas. We expected the confirmation of the statistically significant differences between rural and semi-rural regions. However, it was not confirmed and the statistically significant differences were measured between the Nitra region (rural region) and the Banská Bystrica region (rural region) as well as between the Nitra region (rural region) and the Košice region (semi-rural region). We conclude that the best

conditions for the SMEs development are indicated in the counties of Bratislava I-V and Košice I-IV. After excluding these counties from the analysis, the best region is still the Bratislava region. The second best are the Trnava region and the Trenčín region. The third place is occupied by the Žilina region and the Nitra region. The Prešov region rank on the fourth place. The least attractive regions for SMEs are the Banská Bystrica region and the Košice region (excluding the counties of Košice I-IV). We can conclude that the number of enterprises in remote rural areas grows less rapidly than the number of those in more accessible rural areas.

Finally, we found a very strong relation between the pairs of the numbers of small, medium and large enterprises in each of 79 counties of Slovakia. We suppose that the conditions for doing business in a county are suitable for SMEs as well as for large enterprises. The support policy of SMEs should be more intensive especially in the counties that are not very attractive for large enterprises either. In spite of the effort to eliminate the divergences among the regions of the EU, there are still rather considerable differences among the regions of a given country, not to mention the differences in the whole EU.

5 Conclusions

The share of SMEs was 99.85% during the period of 1996-2015 on average. Out of it, the share of small enterprises was on average 99.31% of a total number of SMEs and the share of medium enterprises in Slovakia was quite negligible. In the future, we expect an increasing trend of numbers of SMEs in Slovakia, mainly an increasing of small enterprises. Medium enterprises are more sensitive to political and economic changes. The development of SMEs is very different in particular regions of Slovakia. The most developed region of Slovakia is the Bratislava region. The number of enterprises in remote rural areas grow less rapidly than the number of those in more accessible rural areas of Slovakia. The conditions for doing business in a county are suitable for SMEs as well as for large enterprises. The support policy of SMEs should be more intensive especially in the counties that are not very attractive for large enterprises. In spite of the efforts to eliminate the divergences among the regions of the EU, there are still rather considerable differences among the regions of Slovakia.

Acknowledgement

The paper is an output of project P7-INCO-2013-9 Reinforcing cooperation with European Neighbourhood Policy countries on bridging the gap between research and innovation (R2I-ENP), KEGA 001UCM-4/2016 Creating Innovative Study Materials for the newly accredited Programme Management in Public Administration and GA 7/2017 Impact of the supporting mechanisms of CAP on the agricultural land market in Slovakia.

References

- [1] Ardic, O. P., Mylenko, N., Saltane, V. Small and Medium Enterprises. A cross – Country Analysis with a New Data Set. Policy Research Working Paper. The World Bank, 2011, 32 p.
- [2] Arlt, J., Arltová, M. Economic time series. Praha: Professional Publishing, 2009, 290 p.
- [3] Attahir, Y. Critical success factors for small business: Perceptions of South Pacific entrepreneurs. In: Journal of Small Business Management, 1995, Vol. 33, No. 2, pp. 68-73
- [4] Beck, T. et al. The influence of financial and legal institutions on firm size. Journal of Banking & Finance, 2006, Vol. 30, No. 11, pp. 2995-3015
- [5] Benzing, C. et al. Entrepreneurs in Turkey: A Factor Analysis of Motivations, Success Factors, and Problems. Journal of Small Business Management, 2009, Vol. 47, No. 1, pp. 58-91
- [6] Berisha, G., Pula, J. S. Defining Small and Medium Enterprises: a critical review. Academic Journal of Business, Administration, Law and Social Sciences, 2015, Vol. 1, No. 1, pp. 17-28
- [7] Bolton, J. E. and Committee of Inquiry on Small Enterprises. Small enterprises: report of the Committee of Inquiry on Small Enterprises. London: Her Majesty's Stationary Office, 1971, 435 p.
- [8] Chowdhury, M. S., Alam, Z., Arif, I. Success Factors of Entrepreneurs of Small and Medium Sized Enterprises: Evidence from Bangladesh. Business and Economic Research, 2013, Vol. 3, No. 2, pp. 38-52
- [9] Cipra, T. Analysis of time series. Praha: SNTL, 1986. 248 p.
- [10] Clover, T. A., Darroch, M. A. G. Owners' perceptions of factors that constrain the survival and growth of small, medium and micro agribusiness in Kwazulu-Natal, South Africa. Agrekon, 2005, Vol. 44, No. 2, pp. 238-263
- [11] Cseh Papp, I., Varga, E., Schwarczová L., Hajós, L. Public work in an international and Hungarian context. Central European Journal of Labour Law and Personnel Management, 2018, Vol. 1, No. 1, pp. 6-15
- [12] Curran, J. Blackburn, R. A. Researching the Small Enterprise. London: SAGE, 2001, 216 p.
- [13] Dahlqvist, J., Davidsson, P., Wiklund, J. Initial conditions as predictors of new venture performance: A replication and extension of the Cooper et al. study. In: Enterprise and Innovation Management Studies, 2000, Vol. 1, No. 1, pp. 1-17

- [14] De Alwis, C. Owner family and business succession in family owned companies. *Acta Oeconomica Universitatis Selye*, 2016, Vol. 5, No. 1, pp. 40-54
- [15] European Commission. The new SME definition. User guide and model declaration. Brusel: EC, 2005, 52 p.
- [16] Gholami, M. H. et al. Investigating the influence of knowledge management practices on Organizational performance: An empirical study. *Acta Polytechnica Hungarica*, 2013, Vol. 10, No. 2, pp. 205-216
- [17] Gibson, T., van der Vaart H. J. Defining SMEs: A Less Imperfect Way of Defining Small and Medium Enterprises in Developing Countries. *Brookings Global Economy and Development*, 2008, 29 p. <<http://seaf.com/wp-content/uploads/2014/10/Defining-SMEs-September-20081.pdf>>
- [18] Gogolová, M. The doorstep sales and the selling events in Slovakia. *Acta Oeconomica Universitatis Selye*, 2014, Vol. 3, No. 2, pp. 33-39
- [19] Grancay, M., et al. Gravity model of trade of the Czech and Slovak Republics 1995-2012: How have determinants of trade changed. *Politická Ekonomie*, 2015, Vol. 63, No. 6, pp. 759-777
- [20] Harabi, N. Determinants of Firm Growth: An Empirical Analysis from Morocco. MPRA Paper, Switzerland: University of Applied Sciences, 2005, 33 p. <https://mpra.ub.uni-muenchen.de/4394/1/MPRA_paper_4394.pdf>
- [21] Hitka, M., et al. Load-carrying Capacity and the Size of Chair Joints Determined for Users with a Higher Body Weight. *Bioresources* 2018, Vol. 13, No. 3, pp. 6428-6443
- [22] Holešová, H. Small and medium business. Bratislava: Úrad vlády SR, 2003, 19 p.
- [23] Horvátová, L. et al. Aspects of effective support of small and medium sized enterprises. *Economics, Management, Innovation*, 2012, Vol. 4, No. 2, pp. 49-59
- [24] Ivanova, S., Latyshov, A. Sustainable entrepreneurship: agrarian policy in South Korea. *Entrepreneurship and Sustainability Issues*, 2018, Vol. 5, No. 4, pp. 748-760
- [25] Jasra, J. M. et al. Determinants of Business Success of Small and Medium Enterprises. *International Journal of Business and Social Science*, 2011, Vol. 2, No. 20, pp. 274-280
- [26] Jeck, T. Small and medium enterprises in Slovakia and in the European Union: Barriers, Financing and Innovations. Bratislava: Ekonomický ústav SAV, 2014, 26 p.

- [27] Keeble, D. Small firm creation, innovation and growth and the urban-rural shift. In: J. Curran and D. Storey (Eds) *Small Enterprises in Urban and Rural Locations*, London: Routledge, 1993, pp. 54-78
- [28] Korcsmáros, E. Factor affecting the development of SMEs (example from Slovakia based on primary research in Nitra region). *Acta Oeconomica Universitatis Selye*, 2018, Vol. 7, No. 1, pp. 70-78
- [29] Liedholm, A. Small Firm Dynamics: Evidence from Africa and Latin America. *Small Business Economics*, 2002, Vol. 18, No. 1, pp. 225-240
- [30] Loecher, U. Small and medium sized enterprises: delimitation and the European definition in the areas of industrial business. *European Business Review*, 2000, Vol. 12, No. 5, pp. 261-264
- [31] Marwede, E. Die Abgrenzungsproblematik mittelständischer Unternehmen - Eine Literaturanalyse (in German). Augsburg: Institute for Economics of University of Augsburg, 1983, 119 p.
- [32] Mura, L. Current situation in family businesses. Managerial trends in the development of enterprises in globalization era, *Conference Proceedings*, 2017, pp. 178-185
- [33] Mura, L. et al. Economic freedom – classification of its level and impact on the economic security. *AD ALTA-Journal of Interdisciplinary Research*, 2017, Vol. 7, No. 2, pp. 154-157
- [34] Nucci, A., Bates, T.. An Analysis of small business size and rate of discontinuance. *Journal of Small Business Management*, 1989, Vol. 27, No. 1, pp. 1-7
- [35] Orlova, L., Gagarinskaya, G., Gorbunova, Y., Kalmykova, O. Start-ups in the field of social and economic development of the region: a cognitive model. *Entrepreneurship and Sustainability Issues*, 2018, Vol. 5, No. 4, pp. 795-811
- [36] Pacáková, V. et al. *Statistics for Economicists*. Bratislava: Ekonómia, 2003, 358 p.
- [37] Rajić, T., Milošević, I. An empirical analysis of the determinants of SME's customer loyalty: evidence from Serbia. *Acta Oeconomica Universitatis Selye*, 2016, Vol. 5, No. 1, pp. 128-138
- [38] Ranjith, J. G. S., Banda, O. G. D. Determinants of Success of Small Business: A Survey-Based Study in Kuliyaipitiya Divisional Secretariat of Sri Lanka. *International Journal of Business and Social Research*, 2014, Vol. 4, No. 6, pp. 38-50
- [39] Rogerson, C. In Search of the African Miracle: Debates on Successful Small Enterprise Development in Africa. In: *Habitat International*, 2001, Vol. 25, No. 1, pp. 115-142

- [40] Sarwoko, E. Frisdiantara, Ch. Growth Determinants of Small Medium Enterprises (SMEs) *Universal Journal of Management*, 2016, Vol. 4, No. 1, pp. 36-41
- [41] Simo, D., Mura, L., Buleca, J. Assessment of milk production competitiveness of the Slovak Republic within the EU-27 countries. *Agricultural Economics-Zemedelska Ekonomika*, 2016, Vol. 62, No. 10, pp. 482-492
- [42] Spyros, G. et al. *Forecasting methods for management*. Michigan: Wiley, 1989, 470 p.
- [43] Storey, D. J. *Understanding the Small Business Sector*. London: Routledge, 2016, 280 p.
- [44] Szabo, K. Z, Herman E. Productive Entrepreneurship in the EU and Its Barriers in Transition Economies: A cluster Analysis. *Acta Polytechnica Hungarica*, 2014, Vol. 11, No. 6, pp. 73-94
- [45] Takáč, I. Competitiveness of small and medium enterprises in Slovakia. Competitiveness and innovations on agricultural land in SR. Nitra: Slovak University of Agriculture, 2011, pp. 61-66
- [46] Takáč, I. EU Policy to promote entrepreneurship of small and medium-sized enterprises in rural areas. *Entrepreneurship in rural areas. EU Business Law I*. Nitra: Slovak University of Agriculture, 2011, pp. 257-265
- [47] Takáč, I. Support of SME business in the EU. Lazíková *et al.* (2013): *EU Business Law*, Nitra: Slovak University of Agriculture, 2013, pp. 160-177
- [48] World Bank. *The big business of small enterprises*. Washington: World Bank, 2014, 262 p.

2018 Reviewers

Ales, Prochazka	Fahlbruch, Babette
Andoga, Rudolf	Földesi, Péter
Ballagi, Áron	Főző, Ladislav
Barabás, Péter	Fullér, Róbert
Barányi, István	Fustos, Janos
Baranyi, Péter	Gajzágó, Éva
Bareith, Attila	Galambos, Péter
Bekut, Dusko	Gašpar, Vladimír
Bencsik, Andrea	Gaul, Emil
Beneda, Károly	Gilányi, Attila
Bicsák, György	Gosztolya, Gábor
Blazic, Saso	Haidegger, Tamás
Bochicchio, Vincenzo	Halawa, Krzysztof
Bogárdi, István	Hegyesi, Franciska
Borbás, Lajos	Hercegfői, Károly
Brassai, Sándor Tihamér	Horváth, Ildikó
Cedomir, Milosavljevic	Horváth, Zoltán
Chmielewska, Katarzyna	Husi, Geza
Cohen, Eliahu	Illés, Tibor
Czampa, Miklós	Izsó, Lajos
Czifra, Árpád	Janjic, Aleksandar
Csapó, Ádám	Juhász, Tímea
Csendes, Tibor	Kárász, Péter
Csernoch, Mária	Károly, Dóra
Danaj, Adela	Kasanicky, Tomas
Deák, István	Kasik, László
Dineva, Adrienn	Kipyatkova, Irina
Doslic, Tomislav	Kiss, Orhidea
Dragan, Antic	Komjaty, Maros
Drégelyi-Kiss, Ágota	Kordic, Branislav
Drexler, Dániel András	Košinár, Peter
Eigner, György	Koszttyánné Mátrai, Rita
Elek, Renáta	Kovács, András
Fábián, Csaba	Kovács, László
Fábián, Enikő Réka	Kovács, Levente
Fábián, Réka	Kovács, Szilveszter

Kovacs, Tibor
Kovacs, Tunde
Kővári, Attila
Krész, Miklós
Kuczmann, Miklós
Laššák, Miroslav
László, Gábor
Lazányi, Kornélia
Machová, Kristína
Machova, Renata
Maros, István
Midner, Vesna
Mildner, Vesna
Milosavljevic, Cedomir
Molnár, Gyöngyvér
Molnár, György
Mucsi, András
Mura, Ladislav
Nagy, Dénes
Nagy, István
Németh, Huba
Oldal, István
Oniga, Stefan
Oran, Ahmet
Orosz, Tamás
Oroszvály, László
Páles, Zsolt
Picariello, Simona
Pintér, Ákos
Pokorádi, László
Prochazka, Ales
Puheim, Michal
Rafajłowicz, Wojciech
Roman, Raul-Cristian
Rövid, András
Sallai, Gyula
Saso, Blazic
Sik-Lányi, Cecília
Šimko, Marián

Simonak, Slavomir
Stevanovic, Dragan
Szabolcsi, Robert
Szakács, Tamás
Szántai, Tamás
Szénási, Sándor
Sziebig, Gábor
Szilágyi, László
Takács, Márta
Takarics, Béla
Tar, József
Tomasek, Martin
Tóth-Laufer, Edit
Török, Ágoston
Vadász-Bognár, Gabriella
Varga, Árpád
Varju, Attila
Várlaki, Péter
Vascak, Jan
Velencei, Jolán
Veress, Árpád
Vicsi, Klára
Zelenka, Jan
Zeleny, Milos
Zsigmond, Gyula