# The Modular Robots Kinematics

**Claudiu Pozna**

Department of Product Design and Robotics, University Transilvania of Brasov
Bd. Eroilor 28, 500036 Brasov, Romania
E-mail: cp@unitbv.ro

## 1 Introduction

The present paper intention is to develop a kinematical foundation for our next works in industrial robots (IR) modular design. The goal of this works is to develop cheap and improved robots which are adapted to the costumer needs. In order to achieve the mentioned goal, in [43], we have started a bibliographical research of the main modular design aspects. The mentioned analyze of the actual results in modular robots design gives us the possibility to establish our research program. The idea of this paper is to develop a kinematical formalism which will be use in the next dedicated to this subject.

The structure of the paper contains a presentation of our ideas about modular robots design, which will be followed by the presentation of our researches direction. From these directions we will focus on the implication of the modularity on robots kinematics and we will propose a new formalism.

## 2 Research Direction in Modular Robotics

The previous description of the actual knowledge, concerning the robots modular design, point out the following problems [1-42]:

- The scientific papers point out the importance of the modular design as a complementary direction of the integral design. The main benefits of this design method are: minimizing the time of design, increasing the number of configurations, an easy maintenance, a fall in prices, etc;

- There have been proposed principals of the modular design (generally) which has conducted to methodologies of design used in industry; this methodologies follow a certain type of modulating, which concerns the producer. Here do not raise the problem that a user can modify the product;

- In the case of modular robots, modulating refers at the user possibility to reconfigure the robot;

- Modularity of the Industrial Robots refers at the same time to the hardware and software aspects. We speck her about the possibility of the mechanical structure modification by combining certain hard modules as well about the possibility of redefining the architecture of the control program by using some programs modules;

- The impact of the modularization in the industrial robotic field is very small. We explain this fact by the complexity of the reconfiguration steps which must made by the user;

- In purpose of raising the configurations performance, the specialty writings mention the need of imagining some methods which incorporate optimization;

- The virtual prototyping represents a base which allows the using of optimization methods.

If we join all this considerations we notice the existence of tow main research directions:

- The implementation of the optimizing procedure in the modular design of industrial robots, the possibility of obtaining an optimal configuration for a certain type of tasks;

- The development of some *friendly* methodologies of reconfiguration: reducing the computing time, usage of some interface that can allow a natural language of programming, reducing the task of the user relishing him from the low level interfacing problems, etc.

If we associate the two directions of research we can observe the following: the implementation of the optimizing techniques supposes a rise of the complexity of design methodologies, implicitly of the reconfiguration process, while the development of some friendly methodologies of reconfiguration involves the simplification of this process.

We underlined that by robots modularity we understand a modularity which is taken upon oneself by the user. This idea belongs to the following scenario: the user buys a particular platform composed by several modules; chooses the appropriate configuration of the robot; constructs the robot from the modules. The user is a task specialist and not a robotic specialist, for these reasons the whole idea is based on the possibility to transfer knowledge from the robots manufacturer to the robots user. This means that, in order to create an ataractic concept of industrial modular robots, we must provide friendly interfaces. These interfaces are dedicated to obtain the robot configuration to assemble this configuration into a robot and to use this robot. For these reasons we have imposed the following design functions:

- The user interface must allow the robot construction:

    o   Obtain an optimal configuration related to this task;

    o   Configuration self recognition;

    o   Model building (kinematics and dynamics);

    o   Translate the user task into a robotic task;

    o   Control law building;

    o   Structure and sensors calibration;

- The user interface must allow the robot employment:

    o   Program the robot;

    o   Allow the robot maintenance.

In conclusion, the idea to use the user modular concept is possible only if appropriate interfaces are designed. We have considered that the first step on this direction is to imagine a kinematical tool which is able to describe the mentioned modularity. More precisely we intend to construct a formalism which will describe the kinematics of all particular construction which can be obtained from the main platform.

# 3   Modular Robots Kinematics

The kinematics researches are important because they offer the possibility to solve problems like: direct kinematics where we impose the desired movements in the robot joints and we obtain the effector's movements; inverse kinematics where we impose the effector's movement and we compute the joint movements; the working volume, where we can obtain the space where the robot task can be accomplish etc. We will mention here that the kinematics is a staring point for the dynamic analyze and the control system design. Our results are based on homogenous transformations described in [44]. Because we focus on the direct kinematical problem, our goal is to obtain a formalism which allows the kinematical description of the robot effectors (gripper, tools etc.) for each possible combination between the links and the joints. In order to do this we will construct the mathematical representation of the links and joints connections. The second step will be the construction of a graph which describes the links and joints connection possibilities. The third step will be to describe from mathematical point of view the previous graphical construction. In the end we will systematize our results in to an algorithm.

**The Connection between Joints and Links**

From the beginning we will mention that our study focuses only in robots with rotation joints which are reciprocally perpendicularly or parallel. The generality of our results is based on the robotic links (brackets) forms and on the various possibilities to attach joints to these brackets.

In Figure 1 we present the imagined general form of the mentioned links. Each link allows the connection with the previous joint at the referential $Ox_{1B}y_{1B}z_{1B}$ and with the follower joint on faces $F_{1...6B}$ at the referential $Ox_{2B}y_{2B}z_{2B}$. Using this form we can describe all the possible reciprocally orientation between the two joint which are connected to the bracket. More precisely, if the first connection (between the joint $j$ and the bracket) is limited to one face, the second connection can be one of the combination between the bracket faces ($F_{1...6B}$) and the joint $j+1$ faces ($F_{1...5A}$).
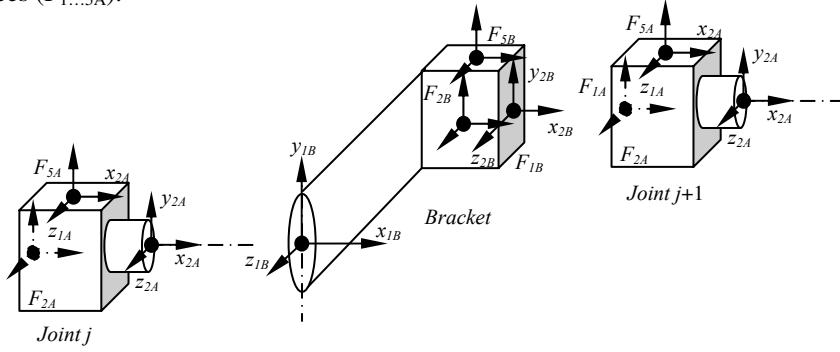


Figure 1

The link form

The link geometry, the positions and orientations of the connections faces, is defined relative to the first referential $Ox_{1B}y_{1B}z_{1B}$. Because of the initial assumptions (the joint are reciprocally perpendicular or parallel) the faces conserve the first referential orientation. That is the reason that from kinematical point of view the relation between the first referential and the faces referential are translations. If we use homogenous operators [44] we obtain the following equations:

$$P_{FiB}^{Bk} = P^{Bk} \cdot \begin{bmatrix} 1 & 0 & 0 & 0 \\ l_{FiB,x}^{Bk} & 1 & 0 & 0 \\ l_{FiB,y}^{Bk} & 0 & 1 & 0 \\ l_{FiB,z}^{Bk} & 0 & 0 & 1 \end{bmatrix} \tag{1}$$

where: $P_{FiB}^{Bk}$ is the position, orientation of the face $i$, which belong to the link $k$;

$P^{Bk}$ is the position, orientation of k link referential (relative to the main referential system);

$l_{FiB,x,y,z}^{Bk}$ are the coordinate of the face $F_{iB}$ center in the $Ox_{1B}y_{1B}z_{1B}$ referential system

$k = 1...n$ is the links type (there are several types of links);

$i = 1...6$ is the face number

According to figure 1 each bracket has two connections: the first with joint $j$, and the second with joint $j+1$. For the first connection we have identified six possibilities which are presented in Figure 2.



a),                                              b),                                              c)



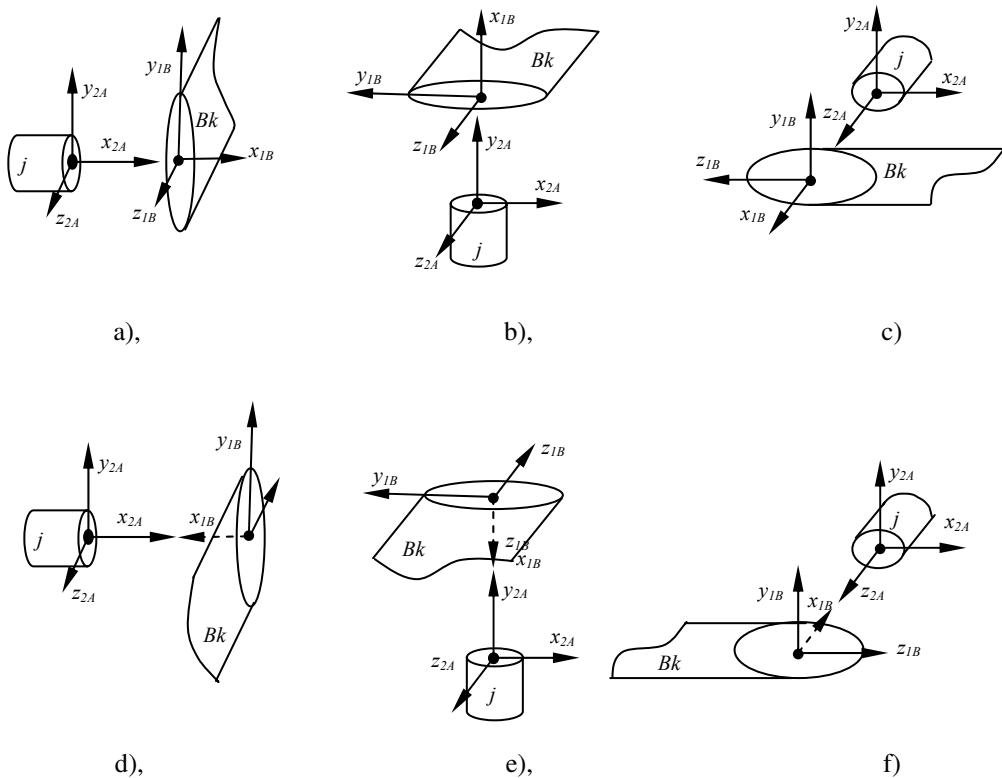d),                                              e),                                              f)

Figure 2
The six possibilities of the connections between joint $j$ and link $B_k$

We intend to *measure* the bracket dimensions (the position and orientation of the $F_{1...6B}$ faces) in the $Ox_{2A}y_{2A}z_{2A}$ referential system, which belongs to the joint $j$. Because the geometry of the link is defined in the $Ox_{1B}y_{1B}z_{1B}$ referential system

we must transform these geometrical data in conformity with the orientation of the joint connection. For this reason we can use the following transformations:

$$\begin{bmatrix} X_{FiB}^{Bk} & Y_{FiB}^{Bk} & Z_{FiB}^{Bk} \end{bmatrix}^T = {}^{j}S_{\beta_X,\beta_Y,\beta_Z} \begin{bmatrix} l_{FiB,x}^{Bk} & l_{FiB,y}^{Bk} & l_{FiB,z}^{Bk} \end{bmatrix}^T \qquad (2)$$

where: $l_{FiB,x,y,z}^{Bk}$ are the coordinate of the face $F_{iB}$ center in the $Ox_{1B}y_{1B}z_{1B}$ referential system;

$X, Y, Z_{FiB}^{Bk}$ are the coordinate of the face $F_{iB}$ center in the $Ox_{2A}y_{2A}z_{2A}$ referential system;

${}^{j}S_{\beta_X,\beta_Y,\beta_Z}$ is the rotation matrix (applied at joint $j$); $\beta_{X,Y,Z} \in \{-1,0,1\}$ :

- for the case presented in Figure 2a

$$ {}^{j}S_{1,0,0} = I_3 \qquad (3)$$

- for the case presented in Figure 2b

$$ {}^{j}S_{0,1,0} = \begin{bmatrix} 0 & 1 & 0 \\ -1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \qquad (4)$$

- for the case presented in Figure 2c

$$ {}^{j}S_{0,0,1} = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ -1 & 0 & 0 \end{bmatrix} \qquad (5)$$

- for the case presented in Figure 2d

$$ {}^{j}S_{-1,0,0} = \begin{bmatrix} -1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & -1 \end{bmatrix} \qquad (6)$$

- for the case presented in Figure 2e

$$ {}^{j}S_{0,-1,0} = \begin{bmatrix} 0 & -1 & 0 \\ -1 & 1 & 0 \\ 0 & 0 & -1 \end{bmatrix} \qquad (7)$$

- for the case presented in Figure 2f

$$
{}^{j}S_{0,0,-1} = \begin{bmatrix} 0 & 0 & -1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix} \tag{8}
$$

In the second extremity of the link we have 30 connection types with joint $j+1$. In Figure 3 we have presented two of these connections. More precisely in Figure 3a the named connection is between the link face $F_{5B}$ and the joint $(j+1)$ face $F_{1A}$; in Figure 3b the named connection is between the link face $F_{6B}$ and the joint $(j+1)$ face $F_{1A.}$

From kinematical point of view to describe these contacts means to use rotations operators. For example in Figure 3a we must apply a rotation right round $z$ axes in order to superpose $Ox_{2B}y_{2B}z_{2B}$ on $Ox_{1A}y_{1A}z_{1A}$:

$$
{}^{j+1}R^{Bk}_{F5B,F1A} = \begin{bmatrix} 0 & -1 & 0 \\ 1 & 0 & 0. \\ 0 & 0 & 1 \end{bmatrix} \tag{9}
$$

For the connection presented in Figure 3b we must apply a rotation right round $y$ axes in order to superpose $Ox_{2B}y_{2B}z_{2B}$ on $Ox_{1A}y_{1A}z_{1A}$:

$$
{}^{j+1}R^{Bk}_{F6B,F1A} = \begin{bmatrix} -1 & 0 & 0 \\ 0 & 1 & 0. \\ 0 & 0 & -1 \end{bmatrix} \tag{10}
$$

The conclusion is that for each connection type we must know the rotation operator which describes the contact. We will define this operator by $R^{Bk}_{FpB,FqA}$:

$$
{}^{j+1}R^{Bk}_{FpB,FqA} = \begin{bmatrix} r_{11} & \cdots & r_{13} \\ . & . & . \\ r_{31} & \cdots & r_{33} \end{bmatrix}; \tag{11}
$$

where: ${}^{j+1}R^{Bk}_{FiB,FiA}$ is the rotation matrix which describe the contact between the face $F_{pB}$ and the face $F_{qA}$; $r_{1\ldots3,1\ldots3}$ are the element of this matrix; $j+1$ is the joint number; $p,q = 1\ldots6$.

It is important to underline that these kinds of matrixes are known for each bracket and for each connection type.
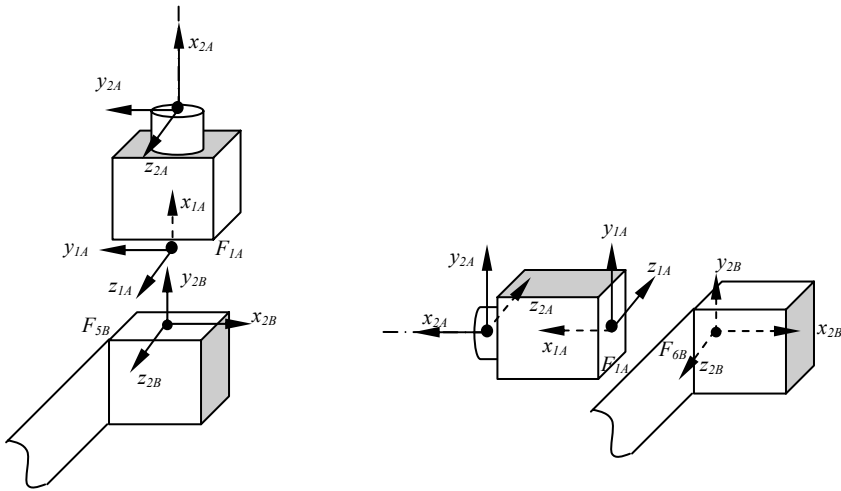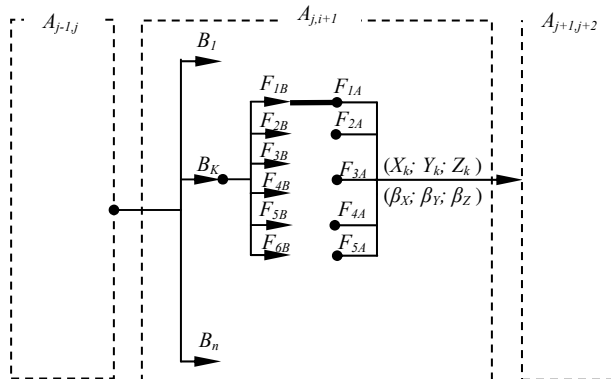
Figure 3

Two of the thirty possible connections between the link $B_k$ and joint $j+1$
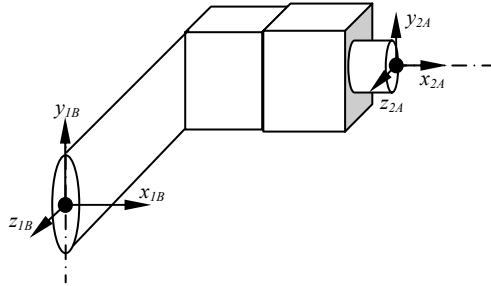
## The Connection Graph

The next step of our analyze focuses on a graphical description of the modularity. More precisely we intend to offer a picture of the modular robot construction from the previous discussed connection point of view. This graphical representation must contain all the possible connection and must bring out the chosen connection. Never the less the graphical construction is a graph which allow a future mathematical representation.

We have presented this graph in Figure 4a and for a better understanding in Figure 4b we have presented the picture of the chosen connection.



a)

b)

Figure 4

The connection graph

The graph (see Figure 4a) shows that we can choose one of the $n$ available brackets and one of the thirty connections between this bracket and the follower joint. The goal is to find a mathematical form which contain implicitly all these possibilities.

### The Homogenous Transformation between Joint $j$ and Joint $j+1$

Using the graph from Figure 4a we propose the following homogenous transformation between joint $j$ and joint $j+1$:

$$A_{j,j+1} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ X_k & |\beta_X| + (|\beta_Y| + |\beta_Z|)\cos(q_{i+1}) & -\beta_Z \sin(q_{i+1}) & \beta_y \sin(q_{i+1}) \\ Y_k & \beta_Z \sin(q_{i+1}) & |\beta_Y| + (|\beta_X| + |\beta_Z|)\cos(q_{i+1}) & -\beta_X \sin(q_{i+1}) \\ Z_k & -\beta_y \sin(q_{i+1}) & \beta_X \sin(q_{i+1}) & (|\beta_X| + |\beta_Y|)\cos(q_{i+1}) + |\beta_Z| \end{bmatrix} \quad (12)$$

Some comments are necessary:

- The homogenous transformation $A_{j,j+1}$ give us the position and the orientation of referential $Ox_{2A}y_{2A}z_{2A}$ relative to the referential $Ox_{1B}y_{1B}z_{1B}$ (see also Figure 4b);

- The $k$ index means that we have choused the link $B_k$ in order to lie joint $j$ to the joint $j+1$;

- $\beta_{X,Y,Z}$ are coefficients which define the direction of joint $j+1$ relative to the main referential system of the robot. More precisely $\beta_{X,Y,Z} \in \{-1,0,1\}$, $\beta_I = -1$ if joint $j+1$ has the direction $- I$, $\beta_I = 0$ if joint $j+1$ has the not the direction $- I$ or $I$, $\beta_I = 1$ if joint $j+1$ has the direction $I$. The mathematical expression of these coefficients can be obtained from the following equation:

$$\begin{bmatrix} \beta_X & \beta_Y & \beta_Z \end{bmatrix}^T = \prod_{p=0}^{j} {}^{p+1}R_{FiB,FiA}^{Bk} \cdot \begin{bmatrix} 1 & 0 & 0 \end{bmatrix}^T \tag{13}$$

- The position is defined by the coordinate $X_k, Y_k, Z_k$ which are computed with equations (2), were we use ${}^{j}S_{\beta_X,\beta_Y,\beta_Z}$ matrix from equation (3-8) according to the value of $\beta_{X,Y,Z}$ coefficients;

- $q_{j+1}$ is the rotation angle in joint $j+1$;

Using these transformations (12) we can compute the position, orientation of the robot end point:

$$P^E = \prod_{j=0}^{m-1} A_{j,j+1} \tag{14}$$

**The Algorithm**

If we summaries the previous results we can propose the following algorithm for the kinematical description of the modular robot:

- choosing a configuration means to choose a succession of brackets which are connected in a desired way to the joint;

- choosing a particular bracket means to know his dimensions $l_{FiB,x,y,z}^{Bk}$;

- a desired connection between the bracket and the joint allows us to know the rotation matrix ${}^{j+1}R_{FpB,FqA}^{Bk}$, which describes this connection (11);

- knowing this matrix we can compute $\beta_{X,Y,Z}$ coefficients (13);

- with these coefficients we can choose ${}^{j}S_{\beta_X,\beta_Y,\beta_Z}$ matrix (3-8) and compute $X_k, Y_k, Z_k$ dimensions (2);

- in the meantime these coefficients give us the possibility to compute the homogenous transformation between two successive joints (12);

- after we have defined our robot configuration we will obtain an equation which lies the joint rotation with position, orientation of the robot end (14);

- this equation can be used to solve kinematics problems (direct, indirect etc).

**Conclusions**

Present paper develops the research on modular robots. If in [1] we have made a bibliographical research, and we have presented our work strategy, in this paper we have started the kinematical analysis of the modular robots. This research focuses only on robots with rotation joints which are reciprocally perpendicularly or parallel. The generality of our study have been ensured by the general form of the link which lies two successive joint and the generality of the connection type between the link and the joint.

The main result that we have achieved is the algorithm which allows the mathematical construction of the homogenous transformation between the modular robots joint. This formalism gives us the possibility to solve the direct kinematics problem: to obtain the position and orientation of the modular robot end point when we impose desired trajectories in the robot joints.

We will continue this study by focusing in the inverse kinematics problem, in the working volume etc.

**Acknowledgements**

**References**

[1]　Bi, M., Zhang, W., Modularity Technology in Manufacturing: Taxonomy and Issuses, International Journal of Advance Manufacturing and Technology, 381-390,2001

[2]　Huang, C.-C., Kusiak, A. Modularity in Design of Products and Systems. IEEE Transactions on Systems, Man and Cybernetics, Part-A: Systems and Humans, 28(1), 66-77, 1998

[3]　Ulrich, T. K., Eppinger, S. D. Product Design and Development, $2^{nd}$ edition, New York:McGraw-Hill, Inc. 2001

[4]　Thyssen, J., Hansen, P. K., Impacts for Modularization. International Conference on Engineering Design (ICED) 01 Glasgow, 21-23 August 2001

[5]　Stake, R., A Hierarchical Classification of the Reasons for Dividing Products into Modules: a Theoretical Analysis of Module Drivers, Licentiate thesis, Royal Institute of Technology, Stockholm, Sweden, 2000

[6]　Siddique, Z., Rosen, D., Product Platform Design: a Graph Grammar Approach, Proceeding of DETC'99: 1999 ASME Design Engineering Technical Conferences, Las Vegas, Nevada, DETC99/DTM-8762, 12-16 September 1999

[7]     Dobrescu, G., Reich, Y., Design of a Gradual Modular Platform and Variants for a Layout Product Family. International Conference on Engineering Design (ICED) 01 Glasgow, 21-23 August 2001

[8]     Jiao, J., Tseng, M., A Methodology of Developing Product Family Architecture for Mass Customization, Journal of Intelligent Manufacturing, 10, pp. 3-20, 2001

[9]     Otto, K., Identyfing Product Family Architecture Modularity Using Function and Variety Heristics. PaperWork Center for Innovation on Product Development, Massachusetts Institute of Technology, Cambridge, MA, 02139, USA, 2001

[10]    Jose, A., Tollenaere, M., Using Modules and Platforms for Product Family Development: Design and Organizational Implications. Fifth International Conference in Integrated Design and Manufacturing in Mechnaical Engineering (IDMME), University of Bath, Bath United Kingdom, 2004

[11]    Hata, T., Kimura, F., Design of Product Modularity for Life Cycle Management, Department of precition Engineering, University of Tokio, Hongo, pp. 7-3-1, 2001

[12]    Chakrabarti, A., Sharing in Design—Categories, Importance, and Issues. International Conference on Engineering Design (ICED) 01 Glasgow, 21-23 August 2001

[13]    Agard, B. (2002) Contribution a` une me´thodologie de conception de produits a` forte diversite´. Ph.D. Thesis, INPG, Grenoble France. 2002

[14]    Venkat, A., Rahul, R., Module-based Multiple Product Design. IIE Proceedings, IERC, Sustainable Design Lab Engineering Management Department, University of Missouri-Rolla Rolla, MO-65401, 2002

[15]    Pedersen, K., Allen, V., Mistree, F., Numerical Taxonomy—a Systematic Approach to Identifyngpotential Product Platforms. International Conference on Engineering Design (ICED) 01 Glasgow, 21-23 August 2001

[16]    Mikkola J., Modularization in New Product Development: Implications for product Architectures, Supply Chain Management, and Industry Structures. Ph.D. Thesis, School of Technologies of Managing, Copenhagen Business School, Denmark, 2003

[17]    Dahmus, J. B., Gonzalez-Zugasti, J. P., Otto, K. N., Modular Product Architecture. Design Studies, 22(5), 409-424, 2001

[18]    Chidambaram, B., Agogino, A., Catalog-based Customization, Proceeding of DETC'99: 1999 ASME Design Engineering Technical Conferences, Las Vegas, Nevada, DETC99/DAC-8675, 12-16 September 1999

[19]    Sand, J. C. House of Modular Enhancement (Home): a Design Tool for Product Modularization, University of Saskachetwan, Master thesis, 1999

[20]   Jose, A., Tollenaere M., Modular and Platform Methods for Product Family Design: Literature Analysis, Journal of Intelligent Manufacturing 16, 371-390, 2005

[21]   Xu, H., Brussel, H., A Behaviour-based Arhitecture with Attention Control, Journal of Inteligent Manufacturing 1998,9,97-106

[22]   Young K, Hybrid Control for Authonomus Mobile robot Navigation Using Neuronal Network-based Behaviour Modules and Environement Clasification, 2003 Authonomus Robots 15, 193-206, Kluwer Academic Publishers, 2003

[23]   Mehmet O., Adaptive Fuzzy Sliding Mode Control for a Class of Bipartite Modular Robotic Systems Journal of Electrical and Electronice Engineering, Vol. 3 645-661, 2003

[24]   R. Fitch, Z. Butler, Z., 3D Rectilinear Motion Planning with Minimum Bend Paths. In Proc. of the Int'l Conf. on Intelligent Robots and Systems, 2001

[25]   Murata, S., Hardware Design of Modular Robotic System, IEEE-RSJ Int. Conf. on Intelligent Robots and System 2000

[26]   Yoshida, E Miniaturization of Self–Reconfigurable Robotic System using Shape Memory Alloy Actuators, J. of Robotics and Mechatronics, Vol. 12 No. 2, 96-102, 2002

[27]   Fromherz, M., Hogg, T., Shang, Y., Jackson, W., Modular Robot Control and Continuous Constraint Satisfaction IJCAI-01 Workshop on Modelling and Solving Problems with Constrains, 2001

[28]   www.lego.com

[29]   Chen, I., Yang, G., Automatic Generation of Dynamics for Modular Robots with Hybrid Geometry. In IEEE ICRA, pp. 2288-2293, 1999

[30]   Farritor, S., On Modular Design and Planning for Field Robotic Systems, Ph.D. Thesis, Massachusetts Institute of Technology, May 1998

[31]   R. Sinha, V. C. Liang, Modeling and Simulation Methods for Design of Engineering Systems, Journal of Computing and Information Science in Engineering. Vol. 1, pp. 84-91, 2001

[32]   Farritor, S., On Modular Design of Field Robotic System, Autonomous Robots 10, 57-65, 2001 Kluwer Academic Publishers. Manufactured in The Netherlands 2001

[33]   Chen, I., Burdick, J, Determining Task Optimal Robot Assembly Configurations. In IEEE Intl. Conf. on Robotics and Automation, pp. 132-137; 2002

[34]  Soshi, I., Paredis, J., Khosla, K., Interactive Multi-Modal Robot Programming, IEEE International Conference on Robotics and Automation Washington DC 2002

[35]  Kang, S-H., Pryor, M., Tesar, D., Kinematic Model and Metrology System for Modular Robot Calibration, 2005

[36]  S. Shin, D. Tesar, "Analytical Integration of Tolerances in Designing Precision Interfaces for Modular Robotics," Ph.D. Dissertation, The University of Texas at Austin, 2004

[37]  Cameron, T., Legault, J., Cox, D., Tesar, D., Configuration Management for Modular Flexible Small Automation System: A Case Study; 2004

[38]  Cox, D., Legault, J., Turner, C., Tesar, D., Automated Plutonium Processing Work Cell Technology, UAmerican Nuclear Society Proceedings of the 7UPUthUPU Topical meeting on Robotics and Remote SystemsU, April 1999

[39]  Legault, J., Tesar, D. Reducing Complexity of Automated Systems Through Configuration Management, Masters Thesis, The University of Texas at Austin, 2000

[40]  Slotine, J. J., Aplied Nonlinear Control Prentice Hall 2000

[41]  Dorf, R., Bishop, R., Modern Control Systems Prentice Hall 2004

[42]  Hung Vu, Dynamic Systems: Modeling and Analysis, Ed. McGraw Hill, 1997

[43]  Pozna, C., Modular Robots Design Concepts and Research Directions. In Proceedings of 5[th] International Symphosium on Inteligent Systems and Informatics. IEEE Catalog Number 07EX1865C, ISBN 1-4244-1443-1

[44]  Gogu, G., Representation du mouvement des corps solides, Hermes, Paris 1996

# Identification of Energy Distribution for Crash Deformational Processes of Road Vehicles

## István Harmati, Péter Várlaki

Department of Chassis and Lightweight Structures
Budapest University of Technology and Economics
Bertalan L. u. 2, H-1111 Budapest, Hungary
harmati@sze.hu, varlaki@kme.bme.hu

*Abstract: Car body deformation modelling plays a very important role in crash accident analyses, as well as in safe car body design. The determination of the energy absorbed by the deformation and the corresponding Energy Equivalent Speed can be of key importance; however their precise determination is a very difficult task. Although, using the results of crash tests, intelligent and soft methods offer an automatic way to model the crash process itself, as well as to determine the absorbed energy, the before-crash speed of the car, etc. In this paper a model is introduced which is able to describe the changing of the energy distribution during the whole deformational process and to analyze the strength of the different parts without any human intervention thus significantly can contribute to the improvement of the modelling, (automatic) design, and safety of car bodies.*

*Keywords: car-crash, deformational energy*

## 1 Introduction

Accident analysis and design of safety car become rather important task. In safety vehicle design one of the most important problems is the reconstruction of the deformation. The base of this method usually a classical model of a simpler dynamical system [4] [7], which can be completed with computer based numerical solutions [8] [11]. Unfortunately these models in most of the cases become unsolvable difficult, if we apply them in a very difficult real system such as a vehicle. A possible way to solve this problem is application of soft computing methods which are very useful for example in the determining the before crash speed of the crashed car [1] [2] [3].

For construction of the optimal car-body structures, a lot of data, so a lot of crash-tests were needed. These tests are quite expensive, thus only some hundred tests per factory are realized annually, which is not a sufficient amount for accident safety. Therefore, real-life tests are supplemented by computer based simulations

which increase the number of analyzed cases to few thousands, so there is an ever-increasing need for more correct techniques, which need less computational time and can more widely be used. Thus, new modelling and calculating methods are highly welcome in deformation analysis.

Through the analysis of traffic accidents we can obtain information concerning the vehicle which can be of help in modifying the structure or the parameters to improve its future safety. The energy absorbed by the deformed car body is one of the most important factors affecting the accidents thus it plays a very important role in car crash tests and accident analysis [5] [6]. The method of the finite elements can be usefully used for simulating the deformation process, but this kind of simulation requires very detailed knowledge about the parameters and energy absorbing properties, besides this it has high complexity and computational cost. The main aim of our experiment is to develop such method, which is able to simulate the deformation process more quickly as the recently used ones. For this purpose rough estimated parameters are used which enable to decrease the complexity and the computational cost.

## 2   The EES Method

In the following we briefly summarize how one can deduce the before crash speed of the vehicle from the absorbed (deformational) energy. From this it is clear that more accurate estimation of the absorbed energy yields more exact estimation of the speed. This method is based on the conservation of energy. It examines the equality of the sum of different kind of energies (kinetic, deformational, heat, etc.) before and after the crash.

Let we use the following usual notations:

$m$: mass of the vehicle      $v$: speed of the vehicle

$\Theta$: moment of inertia      $\omega$: angular velocity

*EES*: Energy Equivalent Speed

Kinetic energy: $E_K = \dfrac{1}{2}mv^2$     Rotational energy: $E_R = \dfrac{1}{2}\Theta\omega^2$

Deformational energy. $E_{DEF} = \dfrac{1}{2}mEES^2$

The energy conservation law for the crash is the following:

$$E_{K1} + E_{R1} + E_{K2} + E_{R2} = E_{K1}^{'} + E_{K2}^{'} + E_{R1}^{'} + E_{R2}^{1} + E_{DEF1} + E_{DEF2} + E_H$$

From this general equation the before crash speed of the vehicle can be estimated. Of course in special cases (for example crashing into a rigid wall) the equation becomes simpler, so the estimating process becomes simpler, too.

# 3   Preliminaries

Analysing the accident statistics it is clear that the most frequently occurring accident is a kind of head-on collision. According to these data the deformational models concentrate usually on the energy absorbed by the deformation of the frontal part of the car-body. The simplest method is to compare the damaged car body being under examination with other similar type of cars damaged under known conditions. In this case the examination is based on the experiences of the human experts. For this approach the required information about the car damaged in known circumstances: standard environments, detailed data of the car, photos, geometrical parameters, etc. Necessary data form the accident's photos, circumstances and environments of the accident, etc.

In the following we give a short historical overview on the energy net approach.

## 3.1   Campbell Model

The first model using energy net is related to K. Campbell (1974). It was based on the assumption that the absorbed deformational energy uniformly distributed on the whole width and height of the vehicle. Further assumption was that under low speed (4-5 km/h) the vehicle is able to tolerate the impact without deformation, and above this limit the deformation is almost linear.
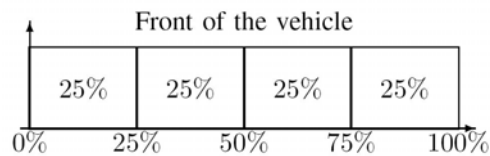
Figure 1
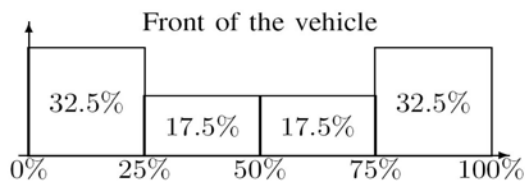The Campbell model

## 3.2   Röhlich Model

Figure 2
The Röhlich model

This model is a further development of Campbell's model. The model doesn't assume uniform distribution about the deformational energy along the width of the vehicle, but deal with a mathematical model based on crash-tests results. This method yields more accurate estimations, for example for the speed of the vehicle before the crash.

## 3.3    Schaper Model

The method worked out by Schaper based on the assumption that vehicles satisfying the prescriptions of FMVSS (Federal Motor Vehicle Safety Standards) have similar deformational characteristics. Using numerous crash-tests results few characteristic fields were defined in the car-body and the distribution of the absorbed deformational energy was examined in these regions.

Figure 3
The Schaper model

# 4    The Basic Concept

In the following we introduce the most important points of the suggested model [9] [10]. We assume having a 2 or 3 dimensional rectangular grid on the car body. The most important thing which we have to know by a detailed analysis of the car body from the point of view of its deformation energy absorbing properties, is the inside structure of the car body. In most of the cases such detailed data about the car-body structure are not available, therefore it is proposed to approximate them from known discrete data. The introduced method is starting from the so-called energy cells. These energy cells are got by discretization of the car-body, to each of which a value for describing its energy absorbing properties is assigned. Other characteristic of these cells is, that each of them can absorb and transmit energy, which depends on the amount of the absorbed energy, e.g. as the stiffness property of a cell is changing during the deformation process proportionately to the absorbed energy. As higher the stiffness of the cells is, as more amount of energy is necessary to achieve the same rate of deformation.

## 4.1   Absorbing Functions

The energy absorbing property of a cell is not a constant value, but changes during the deformational process. Characteristically cells become saturated in a certain sense, so they are able to absorb less and less from the deformational energy. Because of this the energy absorbing property is described by a monotone non–increasing function, instead of a constant value. The simplest is the following piecewise linear function:

$$f(x) = \begin{cases} 1 & x < a \\ \dfrac{b-x}{b-a} & a \leq x \leq b \\ 0 & x > b \end{cases}$$
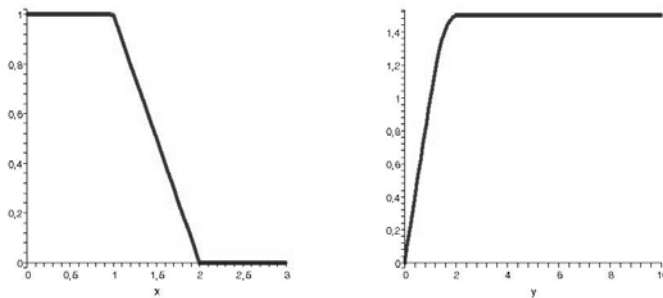


Figure 4

Energy absorbing function and the absorbed energy: simple case

For better approximation one can take a convex linear combination of these kinds of functions:

$$f(x) = \sum_{i=1}^{n} \lambda_i f_i(x) \text{ where } \lambda_i \geq 0 \text{ and } \sum_{i=1}^{n} \lambda_i = 1$$

If the input energy is $E$, then the absorbed energy is:

$$E_{abs} = \int_0^E f(x)dx$$

and the transmitted energy is

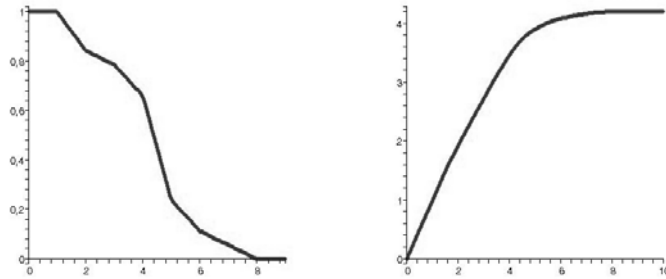$$E_{trans} = \int_0^E 1 - f(x)dx = E - \int_0^E f(x)dx$$

Figure 5

Energy absorbing function and the absorbed energy: convex combination

Principle of conservation of energy is fulfilled, so for absorbed and transmitted energy we have: $E_{abs} + E_{trans} = E$. Subsequently distribution of the transmitted energy between neighbours of the cell will play an important role.

## 4.2    Decomposition of the Impact

The deformational energy appears as kinetic energy (EES). According to this fact decomposition comes from the decomposition of the speed, according to a coordinate system with axes parallel with the grid applied on the car body. The orthogonal components of the velocity are $v \sin \varphi$ and $v \cos \varphi$ in two dimensional case, and $v \sin \vartheta \cos \varphi$, $v \sin \vartheta \sin \varphi$ and $v \cos \vartheta$ in three dimensional case. So the decomposition of the deformational energy is:

• In case of 2D: $E \cos^2 \varphi$ and $E \sin^2 \varphi$

• In case of 3D: $E \sin^2 \vartheta \cos^2 \varphi$, $E \sin^2 \vartheta \sin^2 \varphi$ and $E \cos^2 \vartheta$

## 4.3    Direction Depending Properties

Deformational characteristics of the parts are depending on the direction of the outer impact: a certain part is easy to deform in a certain direction and sometimes very difficult (demands more energy) to deform in other direction (esp. orthogonal direction). Because of this it seems practical to assign different absorbing functions to different orthogonal directions.

Using decomposition in the following we deal with outer impacts parallel with one of the axis of the grid. The final energy distribution will be the sum of the partial energy distributions got by computing with components of the outer impact.

# 5    Deformational Models

We could distinguish at least two models, according to what happens with the transmitted energy. If the transmitted energy is the input of the following cell in the 'tube' then we get a simple model (1$^{st}$ order model). In some cases it gives a right approximation, but in other cases we should determine weighted interconnections between the neighbouring cells (2$^{nd}$ order model).
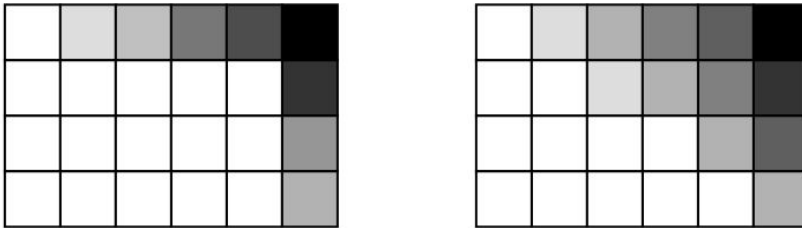


Figure 6

Schematic map of the deformation for the first and second order models

## 5.1    First Order Model

In this case the energy transmitted by the cells is assumed spreading in the direction of the original impact (as we mentioned above, we think of a component of the decomposed impact). It means that forces (screwing forces) between neighbouring cells are ignored; a cell can't deform or drag its neighbours. One can imagine this case as the deformational energy spreads in parallel, but independent 'tubes' along the car-body. This model can be applied cases in which the crash takes place in the whole width of the vehicle (frontal crash), but doesn't give satisfying result for partially overlapping crashes (for instance bumping against a tree).
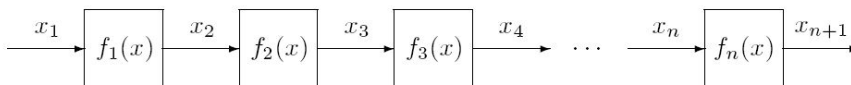


Figure 7

Energy transmission: first order model

## 5.2    Second Order Model

For applications, modelling of the energy spreading process between neighbouring cells not only in the direction of the original impact, but in orthogonal directions, is very important. A possible solution for this task is if we assume that the

transmitted energy is distributed in a certain proportion between all of the neighbours (but usually the next in the 'tube' get the most of the energy).

Let's see the following example. The crash meets the edge of the vehicle (for instance bumping against a pole), and then the deformational energy spreads from here into each directions (see Fig. 8). In the upper row every cell give the $s_{i1}$ times of the transmitted energy to the next cell, and $s_{i2}$ times to the cell bellow (in real beside). The latter case models the energy spreading caused by internal forces between the cells. The energy from this process is interpreted as result of an impact parallel with the original, so we deal with absorbing functions assigned to this direction. In the upper row the inputs are computed similar to the first order model, but in this case we calculate with the $s_{i1}$ times of the transmitted energy. There is a difference in the lower row: the input of cell $(2, 1)$ is $s_{12}$ times of the output of cell $(1, 1)$, and the output is distributed between different directions ($p_{11}$, $p_{12}$). Here we assume that there is no reaction: there is no energy transmission from the lower row to the upper. The input of cell $(2, 2)$ is the sum of $p_{11}$ times of the output of cell $(2, 1)$ and $s_{22}$ times of the output of cell $(1, 2)$. In general, the input of cell $(2, i)$ is the sum of $p_{i-1,1}$ times of the output of cell $(2, i-1)$ and $s_{i,2}$ times of the output of cell $(1, i)$. The $3_{rd}$, $4_{th}$ . . . row can be computed in a similar way.

In general, the computational method is the following:

1    Determine the absorbed and transmitted energy in the first row.

2    Determine the absorbed and transmitted energy at the first element of the second, third ... row applying the input energy assigned to the first element of the second row.

3    Determine the energy transmitted by cell $(2, 2)$ with the sum of its inputs.

4    Using this result determine the same for the third, fourth ... row.

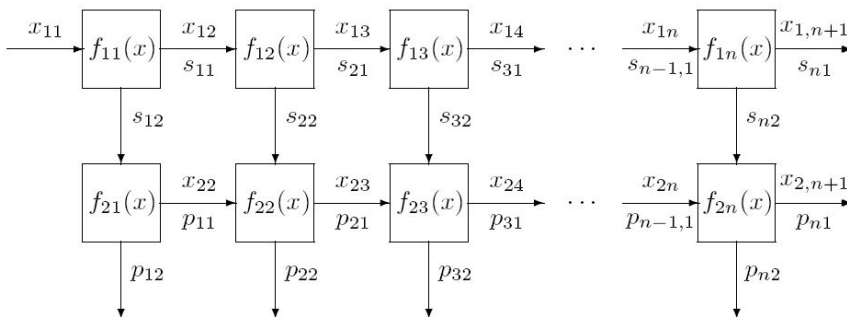5    Execute the procedure on the whole cell matrix.



Figure 8

Energy transmission: second order model

# 6   Illustrative Examples

In this section we show some examples for models introduced in the previous section. There are 6×10 cells in the model, but the parameters of the absorbing functions and the parameters describing the interactions between the cells are not real, just for the demonstration. In both of the cases we show two snapshots, with lower and higher input energy.
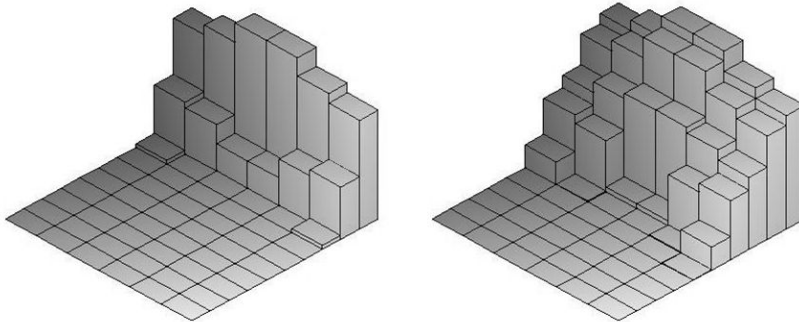


Figure 9
1st order model



Figure 10
2nd order model

**Conclusions**

A new theoretical approach was introduced for describing the changing of the energy distribution during the whole deformational process of the vehicle body. As can be seen from the simulation results, applying this treatment more sophisticated estimations can be made than with using the preliminary methods. The crucial point of the suggested method is the right determination or estimation of the parameters of the absorbing functions. This problem can be solved by analysing crash test results, crash data, digital photos and digital sequences with human experts and intelligent expert systems.

## Acknowledgement

## References

[1]     A. Rövid, A. R. Várkonyi-Kóczy, P. Várlaki, P. Michelberger: Soft Computing Based Car Body Deformation and EES Determination for Car Crash Analysis, In Proceedings of the Instrumentation and Measurement Technology Conference IMTC/2004, Como, Italy, 18-20 May, 2004

[2]     A. Rövid, A. R. Várkonyi-Kóczy, P. Várlaki: Intelligent Methods for Car Deformation Modeling and Crash Speed Estimation, In Proceedings of the 1st Romanian-Hungarian Joint Symposium on Applied Computational Intelligence, Timisoara, Romania, May 25-26, 2004

[3]     A. R. Várkonyi-Kóczy, P. Baranyi, Soft Computing Based Modeling of Nonlinear Systems, 1st Hungarian-Korean Symposium on Soft Computing and Computational Intelligence, Budapest, Hungary, Oct. 2-4, 2002, pp. 139-148

[4]     G. Molnárka: On the Buckling of a Viscoelastic Rood, Numerical Analysis and Mathematical Modeling Banach Center Publications,Warsaw, Poland, 1990, Vol. 24, pp. 551-555

[5]     H. Steffan, B. C. Geigl, A. Moser, H. Hoschopf: Comparison of 10 to 100 km/h Rigid Barrier Impacts, Paper No. 98-S3-P-12

[6]     P. Griškevičius, A. Žiliukas: The Crash Energy Absorption of Vehicles Front Structures, Transport Vol. 18, No. 2, 2003, pp. 97-101

[7]     T. Péter: Reduction of Parameters of Spatial Nonlinear Vehicle Swinging Systems, for Identification and Purposes, Periodica Polytechnica, Vol. 36, No. 1, 1992, pp. 131-141

[8]     T. Péter, E. Zibolen: Model Analysis in Vehicle Dynamics in Computer Algebraic Environment, Symposium on Euroconform Complex Pretraining of Specialists in Road Transport, Budapest, June 9-15, 2001, pp. 319- 331

[9]     I. Harmati, P. Várlaki: Estimation of Energy Distribution for Car-Body Deformation, In. Proc. of the 3rd International Symposium on Computational Intelligence and Intelligent Informatics, March 28-30, 2007, Agadir, Morocco

[10]    I. Harmati, A. Rövid, P. Várlaki: Energy Absorption Modelling Technique for Car Body Deformation, In Proc. of the 4th International Symposium on Applied Computational Intelligence and Informatics, Timisoara, Romania, May 17-18, 2007, pp. 269-272

[11]    E. Halbritter, G. Molnárka: Simulation of Upsetting. Proc. of the 10th International DAAAM Symposium, Vienna University of Technology, October 21-22, 1999. L, R.

# High-Frequency Soft-Switching DC-DC Converters for Voltage and Current DC Power Sources

## Jaroslav Dudrik, Juraj Oetter

Department of Electrical, Mechatronic and Industrial Engineering
Technical University of Košice
Letná 9, 04200 Košice, Slovak Republic
E-mail: jaroslav.dudrik@tuke.sk, tel.: +421 55 6022276, fax: +421 55 6330115
E-mail: juraj.oetter@tuke.sk, tel.: +421 55 6022271, fax: +421 55 6330115

*Abstract: The paper presents soft switching PWM DC-DC converters using power MOSFETs and IGBTs. The attention is focused mainly on the full-bridge converters suitable for high power applications. The properties of the PWM converters are described in comparison to other categories of soft switching converters. An overview of the switching techniques using in the DC-DC converters is included. Considerations are also given to the control methods. The principles of the switching and conduction losses reduction in the PWM converters are illustrated. Various types of snubber circuits are mentioned and their operation and limitations are discussed.*

*Keywords: DC-DC converter, PWM converter, soft switching, snubber*

## 1    Introduction

One of the major trends in power electronics is increasing the switching frequencies. The advances in semiconductor fabrication technology have made it possible to significantly improve not only voltage – and current capabilities but also the switching speed. The faster semiconductors working at high frequencies result in the passive components of the converters – capacitors, inductors and transformers – becoming smaller thereby reducing the total size and weight of the equipment and hence to increase the power density. The dynamic performance is also improved [1], [2], [5], [20].

This frequency elevation is responsible for the growing importance of pulse-width modulation on the one hand and for the use of resonance on the other hand. Another important trend resides in reduction of voltage and current stresses on the semiconductors and limitation of the conducted and radiated noise generated by the converters due large di/dt and du/dt [1], [2], [5], [19], [21].

Both these requirements, size and noise, are minimised if each switch in a converter utilises soft switching technique to change its status. The converter topologies and the switching strategies, which result in soft switching, are discussed in this paper.

# 2   PWM DC-DC Converters

## 2.1   Power Stages of the PWM DC-DC Converters

The full-bridge and half-bridge converters shown in Fig. 1 and Fig. 2 are mostly used in high power applications.
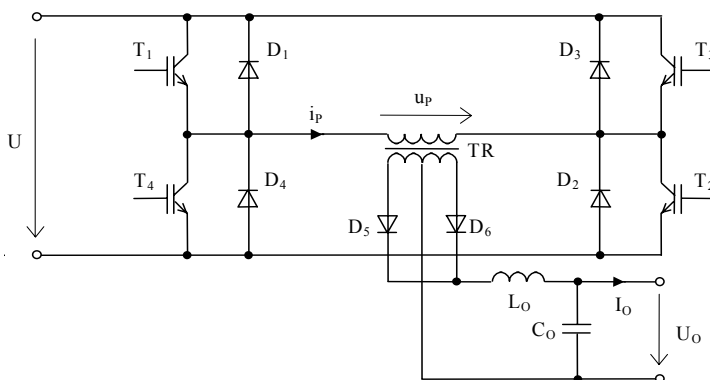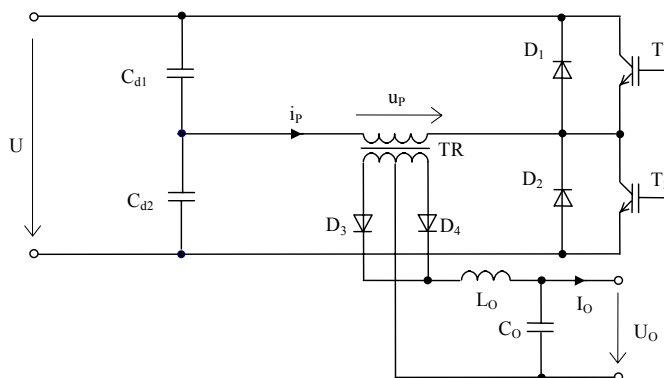
Figure 1

Full–bridge converter

Figure 2

Half–bridge converter

Pulses of opposite polarity are produced on the primary and secondary windings of the transformer by switching of the transistor.

In a connection with the half–bridge inverter, the capacitors $C_{d1}$ and $C_{d2}$ establish a voltage midpoint between zero and the input dc voltage.

The input voltage is equally divided between the capacitors. The relationship between the input and output voltage for the half–bridge is

$$\frac{V_0}{V_d} = \frac{N_2}{N_1} . D \tag{1}$$

and for the full–bridge is

$$\frac{V_0}{V_d} = 2 \frac{N_2}{N_1} . D \tag{2}$$

where duty cycle D= $t_{on}$/T and 0<D<0,5.

Comparison of the full–bridge (FB) converter with the half–bridge (HB) converter for identical input and output voltages and power ratings requires the following turn's ratio:

$$\left(\frac{N_2}{N_1}\right)_{HB} = 2\left(\frac{N_2}{N_1}\right)_{FB} \tag{3}$$

Neglecting the ripple in the current through the filter inductor at the output and assuming the transformer magnetizing current to be negligible in both circuits, the transistor currents $I_C$ are given by

$$(I_C)_{HB} = 2 (I_C)_{FB} \tag{4}$$

In both converters, the input voltage U appears across the switching transistors. However, they are required to carry twice as much current in the half–bridge converter. Therefore, in high power applications, it may be advantageous to use a full bridge over a half bridge to reduce the number of paralleled transistors in the switch.

## 2.2 PWM Strategies for Full–Bridge Converter

The conventional control diagram used for hard driven converters is shown in Fig. 3. The transistors ($T_1$, $T_2$) and ($T_3$, $T_4$) are switched as pairs alternatively at the selected switching frequency, which alternately places the transformer primary across the input supply U for same interval $t_{on}$. The maximum duty cycle is 50% (D=0.5).

A disadvantage of this switching mode is that when all four switches are turned off, the energy stored in the leakage inductance of the power transformer causes severe ringing with junction capacitance of switching devices.

Figure 3
Waveforms of hard–switching converter with conventional PWM

The control scheme in Fig. 4 is almost the same as previous except that the duty cycle in one leg (transistors $T_1$, $T_4$) is constant (D=0.5) and in second leg (transistors $T_2$, $T_3$) is variable in a range between zero and 50%.



Figure 4
Waveforms of converter with modified PWM control

The output voltage can be controlled also via phase control as shown in Fig. 5, [8], [9], [13], [17].

Both legs (transistors $T_1$, $T_4$ and $T_2$, $T_3$) of the bridge operate with a 50% duty cycle, and the phase shift between the legs is controlled. When the two legs operate in phase, the differential voltage applied to the transformer is zero, and zero DC output voltage is obtained. When the two legs of the bridge are in

opposite phase, the differential voltage applied to the transformer, and also the output voltage is maximal.

There are also other derivations of the switching strategies mentioned above [2], [10], [12], [15], [16], [21].



Figure 5
Waveforms of converter with phase–shifted PWM control

The phase shift pulse width modulation (PS-PWM) leads to asymmetrical switching waveforms. The **leading leg** consists of transistors $T_1$, $T_4$ and the **lagging leg** consists of transistors $T_2$, $T_3$. The transistor currents in these legs are not symmetrical. The PS-PWM control strategy leads to zero-voltage turn-on of the transistors in both of legs, as it is evident from oscillograms in Figs. 6 and 7.



Figure 6
Transistor voltage $u_{CE1}$ and current $i_{C1}$ in the leading leg at turn–on and turn–off

Figure 7

Transistor voltage $u_{CE2}$ and current $i_{C2}$ in the lagging leg at turn–on and turn–off

The turn-off losses occur as a result of hard turn-off of the transistors in both legs. The circulating current appears after turn-off of the transistors in the leading leg during freewheeling interval and consequently the conduction losses are increased in the lagging leg transistors.

# 3    Soft Switching PWM Converters

The soft switching PWM converter is defined here as the combination of converter topologies and switching strategies that result in zero–voltage and/or zero–current switching. This type of soft switching converter has been referred to as different names in the literature [1], [2], [3], [6], [9], [11], [14], [15], [18]. They are called also pseudo–resonant, quasi-resonant, resonant transition, clamped voltage topologies and other. In these converters the resonant transition is employed only during a short switching interval. The output voltage is usually controlled by PWM with constant switching frequency.

Soft switching PWM converters can be classified as follows:

1   ZVS PWM converters

2   ZCS PWM converters

3   ZVS ZCS PWM converters

This classification is explained further.

## 3.1    ZVS PWM Converters

The simplest ZVS PWM full bridge converter is shown in Fig. 8. The converter snubbers consist of capacitances $C_1$–$C_4$ and inductance $L_R$, which are represented by transistor and diode output capacitances and transformer leakage inductance respectively.

Figure 8
Full–bridge ZVS PWM converter

The converter is controlled by the phase–shifted PWM technique, which is shown in Fig. 5.

The transistors (MOSFETs or IGBTs) in leading or lagging leg are turned–on while their respective anti-parallel diodes conduct. Since the transistor voltage is zero during the entire turn-on transition, switching loss does not occur at turn–on (Figs. 9 and 10).



Figure 9
Switch (transistor MOSFET $T_1$ and its body diode $D_1$) voltage $u_{DS1}$ and current $i_{DS1}$ during turn–on and turn–off (leading leg)

Figure 10

Switch (transistor MOSFET $T_2$ and its body diode $D_2$) voltage $u_{DS2}$ and current $i_{DS2}$ during turn–on and turn–off (lagging leg)

By utilising small snubber capacitors $C_1 - C_4$ the turn–off losses are sufficiently reduced. If the transistor turn–off time is sufficiently fast, then the transistor is switched fully off before the collector voltage rises significantly above zero, and thus negligible turn–off switching loss is incurred (Figs. 9 and 10).

The detail of the turn-off process of the IGBT transistor in the leading leg of the converter is shown in Fig. 11. Using snubber capacitors in parallel to transistors, the turn-off losses are remarkably decreased.

The ZVS converter exhibits low primary–side switching loss and generated EMI.

However, conduction losses are increased with respect to an ideal hard switching PWM full bridge topology.

At light load, the leakage inductance energy is not sufficient to ensure zero–voltage switching in the lagging leg of the converter. This critical load condition is also a function of the line condition. The worst case is high input voltage when more capacitive energy is required.

Another consideration is the delay time from the turn–off $T_4$ until the turn–on of $T_1$ and visa versa. If the delay time $t_d$ is too short, then the device capacitance may not be fully discharged. However, if the delay time $t_{d2}$ (Fig. 5, Fig. 10) is too long, the capacitor voltage will peak, continue to resonate and drop. Fortunately, the time of peak charge is relative independent of the input voltage and load condition and is equal to one quarter of $L_R C$ time constant [13].

Figure 11

Transistor voltage $u_{CE1}$ and current $i_{C1}$ at turn–off –detail

The secondary–side diodes switch at zero current. This leads to switching losses and ringing as a result of interaction of diodes capacitance with the leakage inductance of the transformer. Additional snubber circuitry is usually required, for prevention of excessive diode voltage stress.

To remove the above-mentioned disadvantages a lot of derivations of the ZVS PWM converters were developed. The penalty for the improvement is usually higher complexity of the converter topology.

## 3.2   ZCS PWM Converters

The ZCS PWM converters can be derived from the ZVS PWM converters by applying the duality principle [14].

The scheme of the full–bridge ZCS PWM (FB–ZCS–PWM) converter is shown in Fig. 12, where $L_R$ is the resonant inductor and $C_R$ is resonant capacitor.



Figure 12

Basic circuit diagram of the FB ZCS-PWM converter

The transformer leakage inductance, the rectifier's junction capacitances, and the transformer winding capacitances can be utilised in this circuit.

Similar to the FB–ZVS–PWM converter, the FB–ZCS–PWM converter also uses phase–shift control at constant switching frequency to achieve required converter operation (Fig. 13).



Figure 13
Idealised waveforms of the FB ZCS PWM converter

The switches must have reverse-voltage blocking capability. The switch can be implemented by an IGBT or a MOSFET in series with a reverse blocking diode, an IGBT with reverse–voltage blocking capability, a MCT, or a GTO. An important advantage of the circuit is that the rectifier diodes do not suffer from reverse recovery problem since they commutate with zero–voltage switching.

This feature makes the converter attractive for applications with high output voltage e. g. power factor correction circuits, where the rectifiers suffer from severe reverse–recovery problems when conventional PWM, ZVS–QRC, or ZVS–PWM converter techniques are used [14].

The efficiency of the converter drops significantly at low line and heavy load since the switches begin to lose zero current switching.

The turn-on and turn-off process of the switches are shown in Figs. 13 and 14 respectively.

Some dual characteristics of the FB ZVS PWM converter and FB ZCS PWM converter are summarised in Table 2.

Figure 14

Switch $S_1$ voltage and current during turn-on and turn-off



Figure 15

Switch $S_3$ voltage and current during turn-on and turn-off

Table 2

Some dual characteristics of the FB ZVS PWM converter and FB ZCS PWM converter

|  | **FB-ZVS-PWM** | **FB-ZCS-PWM** |
|---|---|---|
| Topology type | Buck type | Boost type |
| Switching conditions for active switches | Zero-voltage switching | Zero current switching |
| Switching conditions for rectifiers | Zero current switching | Zero-voltage switching |
| Soft switching easy to achieve at | Heavy load | Light load |
| Implementation of active switches | Diode in parallel with transistor | Diode in series with transistor |

## 3.3    Zero-Voltage Zero-Current Switching PWM Converters

This type of converter is very attractive for high voltage, high power (>10 kW) applications where IGBTs are predominantly used as a power switches [1], [3], [4], [10], [12], [22].

The operating frequency of IGBTs is normally limited to 20-30 kHz because of their current tailing problem. To operate IGBTs at higher switching frequencies, it is required to reduce the turn-off switching losses. ZVS with substantial external capacitor or ZCS can be a solution. The ZCS, however, is deemed more effective since the minority carriers are swept out before turning off.

The zero-voltage zero-current switching (ZVZCS) PWM converters are derived from the full-bridge phase-shifted zero-voltage (FB-PS-ZVS) PWM converters. The PS-ZVS PWM converter is often used in many applications because this topology permits all switching devices to operate under zero-voltage switching by using circuit parasitics such a transformer leakage inductance and devices junction capacitance.

However, because of phase-shifted PWM control, the converter has a disadvantage that circulating current flows through a transformer and switching devices during freewheeling intervals (Fig. 4, Fig. 7).

The circulating current is a sum of the reflected output current and transformer primary magnetizing current. Due to circulating current, RMS current stresses of the transformer and switching devices are still high compared with that of the conventional hard-switching PWM full-bridge converter (Fig. 1). To decrease the circulating current to zero and thus to achieve zero-current switching for lagging leg, various snubbers and/or clamps connected mostly at secondary side of transformer are applied.

The principle by using all of the snubbbers and/or clamps is to secure disconnection of the transformer secondary side, as it is very simplified shown in

Fig. 16. It is usually realised by application of the reverse bias for the secondary side rectifier when transformer secondary voltage in the freewheeling interval becomes zero. The output rectifier ($D_5$, $D_6$) is then reverse biased and the secondary windings of the transformer are opened.



Figure 16
Principle of the ZVZCS converter operation

Therefore, both primary and secondary currents of the transformer become zero. Only a low magnetising current circulates during freewheeling interval as shown in Fig. 17. Thus, the RMS current of the transformer and switches are considerably reduced in the freewheeling interval.



Figure 17
Operation waveforms of ZVZCS PWM converter

Hence, the converter achieves nearly zero-current switching for the lagging leg (transistors $T_2$, $T_3$) due to minimised circulating current during interval of lagging leg transition and achieves zero-voltage switching for leading leg (transistors $T_1$,

$T_4$) due to reflected output current ($I_O/p=I_P$, $p=N_P/N_S$) during the interval of leading leg transition.

Several passive and active snubber and clamp circuits were developed to resolve the problem concerning the resetting the primary current of the transformer to achieve zero-current switching of the switches in the lagging leg of the converter [21].

An example of ZVZCS PWM converter is shown in Fig. 18. ZVS of the leading leg is achieved by the same manner as that of the ZVS full-bridge PWM converter, while the ZCS of the lagging leg is achieved by resetting the primary current during freewheeling period by using active clamp in the secondary side, which needs an additional active switch. Oscillogram of the collector-emitter voltage $u_{CE2}$ and collector current $i_{C2}$ in the lagging leg at turn–on and turn–off is shown in Fig. 19. The transistor is turned-on at zero voltage and turned-off at zero current. The circulating current does not occur, only negligible magnetizing current flows during freewheeling interval through primary winding of transformer. This combination of switching is very effective for IGBT transistors, which have problems at turn-off due to tail current effect.

The converter in Fig. 18 is operating very well at nominal load, but it is not capable operating over wide load range (from no-load conditions to short circuit) with zero-voltage or zero-current switching for all power switches.



Figure 18
ZVZCS DC-DC PWM converter

Figure 19

Transistor voltage $u_{CE2}$ and current $i_{C2}$ in the lagging leg at turn–on and turn–off



Figure 20

Improved ZVZCS DC-DC PWM converter

In order to achieve soft switching at no-load conditions and at short circuit the auxiliary circuits are needed. The example of such patented arrangement is shown in Fig. 20 [7], [22].

The auxiliary transformer $TR_2$ is the main part of the auxiliary circuits in this converter.

The transformer $TR_2$ should have considerably large air-gap to ensure sufficiently high magnetizing current $i_{m2}$ and at the same time to prevent core saturation. The saw-tooth magnetizing current $i_{m2}$ ensures the zero-voltage turn-off of the transistors $T_1$, $T_4$ not only at light load but also at no-load conditions.

Simultaneously, charging or discharging of the capacitors $C_1$, $C_4$ by magnetizing current $i_{m2}$ avoids high current spikes at transistors $T_1$, $T_4$ turn-on at light load and no-load.

In order not to loose the zero-current turn-off of the transistors $T_2$, $T_3$ at short circuit, it is necessary to charge up the capacitor $C_s$ to the rated value of the voltage. The capacitor $C_s$ can be charged from the rectifier GB1, which is connected to the secondary winding of the auxiliary transformer $TR_2$. Soft switching and reduction of circulating currents for full load range are achieved in this converter. The converter is especially suited for application where short circuit and no-load are normal states of the converter operation, e.g. arc welding.

### Conclusion

Principles of the zero-voltage and zero-current switching in PWM full-bridge high-frequency converters are described. The overview and division of the prospective soft-switching PWM converters for high power application is presented.

### Acknowledgement

### References

[1]     Choo, B., H., Lee, D., Y., Yoo, S., B., Hyun, D., S.: A Novel Full-Bridge ZVZCS PWM DC/DC Converter with a Secondary Clamping Circuit. PESC´98, Fukuoka, Japan, pp. 936-941

[2]     Lee, D. Y., Lee, B., K., Hyun, D. S.: A Novel Full-Bridge Zero-Voltage-Transition PWM DC/DC Converter with Zero-Voltage/Zero-Current Switching of Auxiliary Switches. PESC´98, Fukuoka, Japan, pp. 961-968

[3]     Dudrik, J., Dzurko, P.: An Improved Soft-Switching Phase-Shifted PWM Full-Bridge DC-DC Converter. PEMC 2000, Košice, 2000, pp. 2/65-69

[4]     Dudrik, J., Dzurko, P.: Arc-Welding Using Soft-Switching Phase-Shifted PWM Full-Bridge DC-DC Converter. Proc. of the Int. Conf. on Electrical Drives and Power Electronics, 1999, High Tatras, pp. 392-396

[5]     Dudrik, J., Dzurko, P.: Modern Voltage and Current Power Supplies. Proc, of the Int. Conf. EDPE´99, Industry Day, 1999, High Tatras, pp. 46-51 (In Slovak)

[6]     Dudrik, J.: Current–Mode Controlled DC Source for Arc Welding, EPE-PEMC 2004, Riga, Latvia, 2004, pp. 5-203-5-207 – CD Rom

[7]     Dudrik, J.: Circuits for Decreasing of Switching Losses in Extreme Conditions of the Converter. Slovak patent No. 283721, 2003

[8]	Tereň, A., Feňo, I., Špánik, P: DC/DC Converters with Soft (ZVS) Switching. In Conf. Proc. ELEKTRO 2001, section -Electrical Engineering. Zilina 2001, pp. 82-90

[9]	Feňo, I. Jadroň, E. Špánik, P.: Control Circuit for Partial Series Resonant Converter. In: proceedings "TRANSCOM 2001, section 2 – Electrotechnics. Zilina, June 25-27, 2001, pp. 33-36

[10]	Cho, J., G., Rim, G., H., Lee, F., C.: Zero Voltage and Zero Current Switching Full Bridge PWM Converter Using Secondary Active Clamp. Proc. IEEE PESC'96, pp. 657-663

[11]	Hamar, J., Nagy, I.: New Topologies of a Resonant DC-DC Converter Family. In: ELECTROMOTION'2001, Bologna, Italy, June19-20, Vol. 1, pp. 109-114

[12]	Michibira, M., Funaki, T., Matsura, K., Nakaoka, M.: Novel Quasi-Resonant DC-DC Converter Using Phase-Shift Modulation in Secondary Side of High-Frequency Transformer. Proc. IEEE PESC'96, pp. 670-675

[13]	Sabaté, J., A., Vlatković, V., Ridley, R., B., Lee, F., C., Cho, B. H.: Design Consideration for High-Voltage, High-Power, Full-Bridge, Zero-Voltage-Switched PWM Converter. Proc. VPEC, Vol. IV,1991, pp. 231-240

[14]	Hua, G., Lee, F.,C.: A Novel Full - Bridge Zero – Current -Switched PWM Converter. Proc. VPEC, Vol. IV,1991, pp. 215-224

[15]	Lee, D., Y., Lee, B., K., Hyun, D., S.: A Novel Full-Bridge Zero-Voltage-Transition PWM DC/DC Converter with Zero-Voltage/Zero-Current Switching of Auxiliary Switches, PESC´98, Fukuoka, Japan, pp. 961-968

[16]	Morimoto, T., Saitoh, K., Ogura, K., Mamun, A., A, Moiseyev, S., Nakamura, M., Nakaoka, M.: Transformer Parasitic Parameter - Assisted ZVS DC-DC Converter with Synchronous PWM Controlled Active ZCS Rectifier with Choke Input Smoothing Filter. EPE-PEMC 2000 Košice, 2000, Košice, Vol. 2, pp. 18-22

[17]	Trip, D., N., Popescu, V.: Small Signal Model for Phase Shift Control Zero Voltage Switching dc-dc Power Converters. Proc. of the Symposium on Electronics and Telecommunications, Timisoara, Romania, 2002, pp. 6-9

[18]	Rieux, O., Ladoux, P., Meynard, T.: Insulated DC to DC ZVS Converter with Wide Input Voltage Range. EPE´99, Lausanne, Switzerland, 1999, CD, p. 11

[19]	Carriero, C., Rains, F., Volpi, G., F.: Comparison Between Hard and Soft Switching Topologies for Low Voltage Low Power DC-DC Converter in Space Application. EPE´99, Lausanne, Switzerland, 1999, CD, p. 10

[20]	Bauer, P., Bauer, K.: Modern Power Electronics. ISBN 90-9010243-4, 1996

[21]   Liu R.: Comparative Study of Snubber Circuits for DC-DC Converters Utilized in High Power Off-line Power Supply Applications. Proc. IEEE PESC'99, pp. 821-826

[22]   Dudrik, J., Špánik, P., Trip, N.-D.: Zero Voltage and Zero Current Switching Full-Bridge DC-DC Converter with Auxiliary Transformer. IEEE Trans. on Power Electronics, Vol. 21, No. 5, 2006, pp. 1328-1335

# Improvement of the Power Transmission of Distribution Feeders by Fixed Capacitor Banks

## Abdellatif Hamouda

Department of Mechanics
a_hamouda1@yahoo.fr
Tel., fax: +21336925134
QU.E.R.E Laboratory, University Ferhat Abbas, 19000 Setif, Algeria


## Khaled Zehar

Department of Electrical Engineering
khaledzehar@yahoo.fr
Tel., fax: +21336925134
QU.E.R.E Laboratory, University Ferhat Abbas, 19000 Setif, Algeria

*Abstract: The aim of this paper is the presentation of a new analytical formulation of the reactive energy compensation on distribution lines, which are characterised by their radial configuration. It will be devoted to the determination of the sizes and locations of a given number of fixed capacitor banks placed on a non-homogeneous radial line with a non-constant voltage. For the network solution, it is required to know the voltage at each node and at the capacitors location. On the basis of what has been just said and due to the radial type of the line, an iterative method called voltage drop method will be applied. The voltage rms values and phase-angles at all the nodes and on the capacitor banks will be calculated. The mathematical models of the current distributions are made considering the line active and reactive power losses. In the reactive energy optimisation process and for the power and energy loss reductions, we have used new models. The latter take into account the effect of all the capacitors in the calculation of the loss reductions due to a particular one. The results obtained then, are compared to those of authors having previously worked on the subject.*

*Keywords: shunt capacitors, voltage drop, radial line, optimal capacitor sizes, optimal capacitor locations*

# 1   Introduction

The transit of a strong reactive component of the current in an electrical line causes power losses, voltage drop and thus a reduction of the line power transmission. Compared to transmission lines, the distribution ones have a low voltage and high current. The $RI^2$ loss (up of 13%) in distribution systems is than significantly high. To improve the line power transmission and to avoid turning to investments in new distribution lines, the power utilities are firstly forced to reduce the losses in distribution systems. To achieve power and energy loss reductions as well as voltage correction, shunt capacitor banks are widely used.

The capacitor installation techniques can be classified in several groups among of which we note; the analytical methods [1]-[7] which are easy to understand and gives detailed mathematical models, the numerical methods [8]-[11] which are iterative and with or without constraints, the heuristic methods [12]-[14] developed through intuition and experience. They are a fast practical way which leads to a near optimal solution and reduce the search space. Artificial intelligence based methods [15]-[20]. Based on the natural evolution and the annealing of solids, these methods can also be a combination of a set of methods (hybrid methods). In this class of methods we note: genetic algorithms [15]-[17]; fuzzy logic and genetic algorithms [18]; simulated annealing [19]-[20].

Our interest in this paper relates to the analytical methods. Several methods have been used to conduct the reactive energy compensation in an optimal way i.e., to have less power and energy losses. If the early analytical works [1]-[4] constituted an important step in the modelling of the optimisation of the reactive energy compensation, they remain however non-realistic because of the number of assumptions that they considered. Some of these assumptions are uniform line, uniform load and constant voltage along the line. Later works based on heuristic search [12]-14], although they do not make any simplifying assumptions, tackle the problem by initially identifying the possible nodes candidate to carry capacitor banks, then determining their optimal sizes. The problem is thus reduced to capacitor sizes optimisation.  The advantage of this method lies in the fact of being easy and very practical especially for the feeders with laterals. But, the capacitor locations are the more indicated but not the optimal ones. Moreover, for more than one capacitor, the effect of the batteries the ones on the others does not appear clearly in the objective function for a possible mathematical analysis of the problem. It appears implicitly in the reactive current updating when we perform the load flow. References [5]-[7], in our point of view, present best and complete mathematical models for leading the reactive energy optimisation. Indeed, this method allows the electrical energy suppliers to choose either to optimise only one parameter (capacitors location or size) or both of them at the same time.

Nevertheless, owing to the fact that the authors in [5]-[7] have defined the power and energy loss reductions due to a given capacitor as dependent only on the

powers of those located at its downstream and owing to fact that the formulation of the capacitor current [7] is not homogenous, we present in this paper new expressions for the power and energy loss reductions as well as the equations whose results give the capacitor optimal sizes and locations. However, we firstly define the objective function as well as a simple manner to calculate the nodes voltage and the voltage at the capacitor locations.

## 2   Objective Function

To optimise the reactive energy compensation, the definition of an objective function is essential. This objective function called also 'economic savings' depends on the power and energy loss reductions due to all the installed capacitors and their costs. This function noted '$S$', is defined by [5]-[7]:

$$S \; = \; k_p \Delta P \; + \; k_e \Delta E \; - \; k_{cf} \sum_{k=1}^{n} Q_{ck} \tag{1}$$

By using the capacitor voltage and courant, the expression (1) becomes:

$$S \; = \; k_p \Delta P \; + \; k_e \Delta E \; - \; k_{cf} \sum_{k=1}^{n} \frac{V_{ck} \, I_{cqk}}{\cos \varphi_{ck}} \tag{2}$$

Expression (2) shows that '$S$' is dependent on the capacitors current and voltage. It depends also on the locations of the capacitors which appear explicitly in the expression of the power and energy loss reductions given below. These expressions are for a balanced radial main feeder on which '$n$' capacitors numbered, for the calculation suitability, from the end of the line to the substation end (*Figure 1*), as it follows:



Figure 1

Notation of the capacitors current and location

## 2.1   Power Loss Reduction

For a balanced three phase radial feeder and using the power loss reduction due to each capacitor, the total loss reduction can be written as it follows:

$$\Delta P = 3 \sum_{i=1}^{n} \Delta P_i \tag{3}$$

$\Delta P_i$ : is then the power loss reduction due to the $i^{th}$ capacitor. It is given by:

$$\Delta P_i = 2 R I_{cqi} \int_{0}^{h_i} \left( I_s(t) F_q(x) - \sum_{k=1}^{i-1} I_{cqk} \right) dx - 2 R I_{cqi} \sum_{k=i+1}^{n} I_{cqk} h_k - R h_i I_{cqi}^2 \tag{4}$$

## 2.2   Energy Loss Reduction

As for the power loss reduction, the energy loss reduction for this radial line is given by:

$$\Delta E = 3 \sum_{i=1}^{n} \Delta E_i \tag{5}$$

The capacitors being of fixed type and their in service duration is $T$ then, $\Delta E_i$ is equal to the integral between 0 and T of $\Delta P_i$ and it admits the following expression:

$$\Delta E_i = 2 R I_{cqi} T L_f \int_{0}^{h_i} I_s F_q(x) dx - 2 R I_{cqi} h_i T \sum_{k=1}^{i-1} I_{cqk} - 2 R T I_{cqi} \sum_{k=i+1}^{n} I_{cqk} h_k - R T h_i I_{cqi}^2 \tag{6}$$

# 3   Voltages Calculation

The load and battery currents are voltage-dependent thus, the calculation of the complex voltages is necessary. However, distribution networks being characterized by a high ratio *R/X* and a radial configuration, it is not recommended to use conventional load flow methods such as Gauss-Seidel or Newton decoupled which are essentially developed for transmission or strongly meshed networks. Applied for distribution networks, they can encounter convergence problems. In this case and for the voltage calculation, we suggest a method called voltage drop method or a backward and forward method. In the line backward sweep, the branch currents, the active and reactive powers and the power losses are calculated. After which, the forward sweep is carried out to determine the nodes

voltage and phase-angle. For a branch '$i$', the complex voltage of its receiving end is defined equal to that of its sending end '$i-1$' decreased by the branch voltage drop located between the two considered nodes. For the node '$i$' (*Figure 2*), the complex voltage is written as it follows:

$$\bar{V}_i = \bar{V}_{i-1} - (r_i + jx_i)\left[F_{di} - j(F_{qi} - F_{ci})\right]$$



Figure 2
Descriptive diagram of the line

Where the d and q components of the $i^{th}$ branch current and the current due to the capacitors located downstream this branch, are given by:

$$\begin{cases} F_{di} = \dfrac{P_i \cos \varphi_i + Q_i \sin \varphi_i}{V_i} \\[2ex] F_{qi} = \dfrac{Q_i \cos \varphi_i - P_i \sin \varphi_i}{V_i} \\[2ex] F_{ci} = \displaystyle\sum_{k=1}^{i} Icqk \end{cases} \tag{7}$$

The active and reactive powers injected into the node '$i$' are given by:

$$\begin{cases} P_i = P_{i+1} + P_{Li} + P_{loss\,i+1} \\[1ex] Q_i = Q_{i+1} + Q_{Li} + Q_{loss\,i+1} \end{cases} \tag{8}$$

The line active and reactive power losses are:

$$\begin{cases} P_{loss\,i+1} = r_{i+1} \dfrac{P_{i+1}^2 + Q_{i+1}^2}{V_{i+1}^2} \\[3ex] Q_{loss\,i+1} = x_{i+1} \dfrac{P_{i+1}^2 + Q_{i+1}^2}{V_{i+1}^2} \end{cases} \tag{9}$$

The $d$ and $q$ voltage-components, using the uniform normalised line concept, are:

$$\begin{cases} V_{di} = V_{di-1} - RL_{uni} \, F_{di} - x_{ni} L_{uni} \, F_{qi} + x_{ni} \, L_{uni} \, F_{ci} \\ V_{qi} = V_{qi-1} - x_{ni} L_{uni} \, F_{di} + RL_{uni} \, F_{qi} - RL_{uni} \, F_{ci} \end{cases} \tag{10}$$

Once the $d$ and $q$ components calculated, the voltage rms value and phase-angle of node '$i$' are obtained by:

$$\begin{aligned} V_i &= \sqrt{V_{di}^2 + V_{qi}^2} \\ \varphi_i &= ar \tan g \frac{V_{qi}}{V_{di}} \end{aligned} \tag{11}$$

The complex voltage of the $k^{th}$ capacitor is equal to that of bus '$i$' if it is located on it. If the capacitor location is between the buses '$i-1$' and '$i$', its $d$ and $q$ components are given by:

$$\begin{cases} V_{cdk} = V_{di-1} - R \, (L_{nni} - h_k) \, F_{di} + X_{ni} (L_{nni} - h_k) \, F_{qi} + X_{ni} (L_{nni} - h_k) \, F_{ci} \\ V_{cqk} = V_{qi-1} - X_{ni} \, (L_{nni} - h_k) \, F_{di} + R \, (L_{nni} - h_k) \, F_{qi} + R \, (L_{nni} - h_k) \, F_{ci} \end{cases} \tag{12}$$

From d and q components we get:

$$\begin{aligned} V_{ck} &= \sqrt{V_{cdk}^2 + V_{cqk}^2} \\ \varphi_{ck} &= ar \tan g \frac{V_{cqk}}{V_{cdk}} \end{aligned} \tag{13}$$

To determine both voltage magnitude and phase-angle, one will initialise all the complex voltages to that existing at the substation end (reference bus) and calculate initially the current distributions according to (7) by going up the line (backward sweep). Then $d$ and $q$ voltage components, in agreement with the expressions (10) for nodes and (12) for capacitors, are calculated by going down the line (forward sweep). The method being iterative, the computing process will be stopped only if the results converge. As a convergence test, we have adopted a per-unit difference of voltages of two successive iterations equal or less than 0.0001.

# 4   Optimisation of the Reactive Energy

Making the objective function maximum is equivalent to finding the batteries size and location which satisfy the following system:

$$\begin{cases} \partial S / \partial I_{cqi} = 0 \\ \partial S / \partial h_i = 0 \end{cases} \tag{14}$$

The solution strategy suggested is an iterative procedure being based on the solution of each equation of the system (14) for the two following major reasons:

- The solution facility which the iterative method offers in this case.

- Each equation of the system (14) taken separately, constitutes alone, a problem the importance of which is proven. Indeed, it happens that the interest of the electrical energy suppliers, for considerations which are peculiar to them, relates only to one of the two parameters independently of the other. Owing to the fact that one solves each equation separately, the access to the solution of only one of the two problems is than possible.

## 4.1    Optimisation of the Sizes

The substitution of «S» by its expression (2) in the first equation of the system (14) and after reorganizing the equation, we end up with the following contracted matrix expression (see the appendix for more details).

$$\widetilde{H}\,\widetilde{I}_{cq} = \widetilde{B} \tag{15}$$

Where:

- $\widetilde{H}$ : is an $n.n$ matrix called matrix locations the elements of which are such that:

$$h_{ij} = \begin{cases} h_i & if \quad i = j \\ 2h_j & if \ i \prec j \\ 2h_i & if \ i \succ j \end{cases} \tag{16}$$

- $\widetilde{I}_{cq}$ : is a $1.n$ matrix called capacitors size matrix the transposed of which is:

$$\widetilde{I}_{cq}^{\,t} = \big[ I_{cq1}, I_{cq2}, \dots, I_{cqn} \big]$$

- $\widetilde{B}$ : is a $1.n$ matrix the elements of which are such that:

$$B_i = \frac{k_p + k_e TL_f}{k_p + k_e T} \int_0^{h_i} I_s F_q(x)dx - \frac{k_{cf}V_{ci}}{2R(k_p + k_e T)\cos \varphi_{ci}}$$

Obtaining the optimal sizes passes by the resolution of the matrix equation (15) which gives the capacitor current $I_{cqk}$. The reduced optimal sizes of the capacitors ($Q_{ck}$) are then deduced from:

$$Q_{ck} = \frac{V_{ck}I_{cqk}}{\cos \varphi_{ck}} \tag{17}$$

## 4.2    Optimisation of the Locations

Just like for the sizes of the batteries, the resolution of the second equation of the system (14) and after the equation reorganisation, leads to (see appendix for details):

$$F_q(h_i) = \frac{2(k_p + k_e T)}{I_s(k_p + k_e TL_f)} \sum_{k=1}^{i-1} I_{cqk} + \frac{k_p + k_e T}{2 I_s(k_p + k_e TL_f)} I_{cqi} \quad (18)$$

Knowing $F_q(h_i)$, the per-unit optimal locations and consequently their real values are then deduced from the graph of the reactive current distribution function $F_q(x)$ shown in *Figure 3*. However, and to make the locations determination automatic, a program was envisaged for this purpose.



Figure 3
d and q branch current distributions before and after compensation

## 4.3    Optimisation of Two Parameters

To optimise the locations and the sizes of the batteries at the same time, an iterative procedure is planned. It calls upon the programs developed to determine each of the two parameters separately. The execution of this iterative method requires the knowledge of an initial solution. Arbitrary values are then assigned to the two required parameters. The determination of the optimal locations and sizes will be done, according to the following algorithm:

*Step 1    Read the line data.*

*Step 2    Read the arbitrary capacitors size and location.*

*Step 3    Initialise the tensions of the various bus bars and on the capacitors.*

*Step 4    Uniform and normalise the line and the loads.*

*Step 5    Calculate the normalised currents $I_{cqk}$ due to the capacitors from (17).*

*Step 6    While the convergence is not reached perform the following steps:*

   *a) Calculate the current distribution functions according to first and second expressions of (7).*

   *b) Calculate the bus voltages and the capacitors voltages bus from (11) to (13).*

   *c) While the capacitors location and size are not identical to the precedent ones, carry out the following steps:*

      *i)    Calculate the optimal locations according to (18).*

      *ii)   Calculate the optimal currents due to the capacitors according to (15).*

      *iii)  Calculate the relative powers of the capacitors according to (17).*

      *iv)   Calculate the cost reductions according to (1).*

   *d) Else continue.*

*Step 7    Else continue.*

*Step 8    Return to the real dimensions.*

*Step 9    Write the results.*

# 5    Application

As an application example and in order to be able to undertake a comparative study, we have considered the physically existing distribution line given by [5-7]. It's a non-homogeneous distribution line of medium voltage having nine sections of five wire-sizes. The line loads are non-uniform and are concentrated at the end of each section. As a base voltage we have adopted the voltage at the sub-station end (23 kV) which is also regarded as the angles origin. As a base power we have considered an apparent power of 4186 kVA. This value is equal to the sum of all the reactive loads.

As in [5-7], the three capacitors size and location optimisation problem is considered where: the load factor $L_f$ is equal to 0.45; the annual cost of the kW is

$k_p$ =168\$/kW; the annual cost of the kWh is $k_e$ =0.015 \$/kWh and the annual cost of the installed kVAr is $k_{cf}$ =4.9 \$/3 phasekvar. A 14.3% annual fixed charge rate is applied for capacitor cost. The obtained results for the optimal sizes and locations are consigned in *Table 1*. Some others interesting results are also given. The voltage rms values and phase-angle after the optimisation of the reactive energy are consigned in *Table 2* and those of the effect of the shunt battery current definition in *Table 3*.

Table 1

Results of the capacitors size and location optimisation

| N° of the capacitor | | 1 | 2 | 3 | $\Delta P$ (kW) | $\Delta P$ in[7] | $\Delta E$ (kWh) | S (\$) | number iterations |
|---|---|---|---|---|---|---|---|---|---|
| Initial sizes | (kVAr) | 300 | 600 | 800 | | | | | |
| Initial locations | (mile) | 13.27 | 07.32 | 3.02 | | | | | |
| Optimal Locations | (p.u) | 1.0000 | 0.2248 | 0.1074 | | | | | |
| Locations in [7] | (p.u) | 1.0000 | 0.2248 | 0.1074 | | | | | |
| Optimal Locations | (mile) | 16.27 | 06.32 | 4.02 | | | | | |
| Optimal Sizes | (p.u) | 0.1628 | 0.0964 | 0.0637 | 152.71 | 111.00 | 83818 | 19328 | 3 |
| Optimal Sizes | (kVAr) | 681.48 | 403.62 | 266.73 | | | | | |
| capacitors Currents | (p.u) | 0.1877 | 0.1040 | 0.0669 | | | | | |
| capacitors Sizes in [7] (kvar) | | 464 | 1070 | 2961 | | | | | |
| $Q_{cq[7]}/Q_{cq}$ Ratios | | 0.68 | 2.65 | 11.10 | | | | | |

Table 2

Bus voltage rms values and phase-angles after capacitors installation

| Bus | V (p.u) | $\varphi$ (rd) | $V_{[7]}$ | $\varphi_{[7]}$ | Bus | V (p.u) | $\varphi$ (rd) | $V_{[7]}$ | $\varphi_{[7]}$ |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1.0000 | 0.0000 | 1.0000 | 0.0000 | 5 | 0.9249 | -0.0726 | 0.9441 | -0.0831 |
| 1 | 0.9941 | -0.0095 | 0.9967 | -0.0102 | 6 | 0.9163 | -0.0811 | 0.9350 | -0.0906 |
| 2 | 0.9853 | -0.0212 | 0.9917 | -0.0232 | 7 | 0.9003 | -0.0926 | 0.9183 | -0.1001 |
| 3 | 0.9648 | -0.0419 | 0.9790 | -0.0481 | 8 | 0.8753 | -0.1135 | 0.8912 | -0.1169 |
| 4 | 0.9513 | -0.0495 | 0.9695 | -0.0600 | 9 | 0.8600 | -0.1319 | 0.8732 | -0.1311 |

Table 3

Effect of the current battery definition

| Battery | $I_{cq}$ (p.u) | $Q_{cq}$ (p.u) | $Q_{cq\ [7]}$ (p.u) | $Q_{cq[7]}$ / $Q_{cq}$ | Deviation (%) |
|---------|---------|---------|---------|---------|---------|
| 1 | 0.1877 | 0.1628 | 0.2168 | 1.33 | 33 |
| 2 | 0.1040 | 0.0964 | 0.1127 | 1.17 | 17 |
| 3 | 0.0669 | 0.0637 | 0.0704 | 1.11 | 11 |

## Conclusion

As consequences of the changes in the expressions (3) for $\Delta P_i$ and (4) for $\Delta E_i$, we note that:

The non-diagonal terms of the matrix locations are multiplied by two, if the capacitors size optimisation is of interest see expression (16). The first term of the right-hand side of the equation (18) giving $F_q(h_i)$ is multiplied by two, during the capacitors location optimisation.

From the results point of view, if the optimal locations (nodes: 9, 5 and 4) of the capacitors are the same to those given in [7], their sizes are completely different, see *Table 1*. The ratios of the capacitors size given in [7] to those which we have obtained, vary from 0.68 to 11.10 (see *Table 1*). The optimal choice of capacitors location and size, taking into account the line power losses, conduct to a power loss reduction equal to 152.71 kW, an improvement in the voltage profile (see *Table 2*) and a decreasing in the reactive current distribution $F_q(x)$ (see *Figure 3*). The cost reduction is then equal to 19328 \$ and would be better in our case. Indeed, if we count the total number of kVAr installed, it is equal to 1351.83 kVAr in our case and 4495 kVAr in [7]. Reported to the total reactive power (4186 kVAr), the ratio is of 32.29% in our case and 107.4% in [7]. This last value i.e. 107.4% means that the total requested reactive energy is satisfied by an external contribution and violates the maximum limit of the reactive energy to be compensated. In addition to the power loss reduction due the optimal capacitors placement, the improvement of the voltage, reduce the d component of the branch current (see Figure 3) and thus an additional reduction in the power losses. The power loss reduction due to this component of the branch current is equal to 20.28 kW and consequently, a supplementary cost reduction of 3407 \$. The total reductions of the power losses and the cost are respectively of 172.99 kW and 22735 \$.

Note that the optimal sizes of the capacitors obtained in our case are not standard ones. To overcome this, we suggest moving each non-standard size to that of smaller standard size or larger standard one and then choose those whose economic saving is the best.

If we consider that the optimal reactive currents are the same in both, reference [7] and our proposed method (*Table 3*), the difference in the capacitor power definition (Equation 17) leads to a deviation between 11% and 33%. This deviation increases as one move away from the sub-station end or as voltage magnitude decreases.

**Nomenclature**

- $k_p$ ( $k_e$ ): is the annual unit price of the kW (kWh).

- $k_{cf}$ : is the unit price of the three phase kVAr installed.

- $V_i$ ( $\varphi_i$ ): is the voltage rms value (phase-angle) at bus *i*.

- $V_{ck}$( $\varphi_{ck}$ ): is the $k^{th}$ battery voltage rms value (phase-angle).

- $X_{ni}$ : is the per-unit normalised reactance of the $i^{th}$ line section.

- $R$ : is the per unit resistance of the normalised uniform line.

- $h_k$ : is the normalised uniform location of the $k^{th}$ battery.

- $Q_{ck}$: is the $k^{th}$ battery size.

- $I_s(t)$ : is the time-dependent reactive current at the substation end.

- $T$ : is the in service duration of the batteries which are of fixed type. Its per-unit value is equal to 1.

- $L_{uni}$ : is the normalised uniform length of the branch '*i*'.

- $L_{nni}$ : is the total normalised uniform length (from the reference node to the node '*i*').

**Appendix**

*1   Optimisation of the sizes*

The optimal sizes of the batteries are obtained by the resolution of the system of equations $\partial S / \partial I_{qcj} = 0$. To do it, the derivatives of the power and energy loss reductions are required. We then obtain:

$$\partial \Delta P / \partial I_{cqj} = \sum_{i=1}^{n} \partial \Delta P_i / \partial I_{cqj}$$

Where $\partial \Delta P_i / \partial I_{qcj}$ is equal to:

$$\frac{\partial \Delta P_i}{\partial I_{cqj}} = \begin{cases} -2RI_{cqi}h_i & if \ i \succ j \\ 2R\int_0^{h_i} I_s F_q(x)dx - 2Rh_i \sum_{k=1}^{i-1} I_{cqk} - 2R\sum_{k=i}^{n} h_k I_{cqk} & if \ i=j \\ -2RI_{cqi}h_i & if \ i \prec j \end{cases} \quad (A1)$$

From where the expression of $\partial \Delta P / \partial I_{cqj}$ :

$$\frac{\partial \Delta P}{\partial I_{cqj}} = 3(\ 2\ R\int_0^{h_j} I_s\ F_q(x)\,dx - 4\,R\,h_j \sum_{k=1}^{j-1} I_{cqk} - 4\,R\ \sum_{k=j+1}^{n} h_k\ I_{cqk} - 2R h_j I_{cqj})$$

Making the same reasoning we obtain for $\partial \Delta E / \partial I_{cqj}$ :

$$\frac{\partial \Delta E}{\partial\ I_{cqj}} = 3(2\ R\ T L_f \int_0^{h_j} I_s\ F_q(x)\ dx\ -4\ R\ T\ h_j \sum_{k=1}^{j-1} I_{cqk} - 4\ R\ T\ \sum_{k=j+1}^{n} h_k\ I_{cqk} - 2\ R\ T h_j I_{cqj}) \quad (A2)$$

Finally we obtain for $\partial S / \partial I_{cqj} = 0$ , once a certain number of arrangements operated:

$$2h_j \sum_{k=1}^{j-1} I_{cqk} + 2\sum_{k=j+1}^{n} h_k I_{cqk} + h_j I_{cqj} = \frac{k_p + k_e TL_f}{k_p + k_e T} \int_0^{h_j} I_s F_q(x)dx - \frac{k_{cf}}{2R(k_p + k_e T)} \quad (A3)$$

Where: $\sum_{k=1}^{j-1} I_{cqk} = 0$ for $j=1$ and $\sum_{k=j+1}^{n} I_{cqk} = 0$ for $j=n$.

*(A3)* can be written in the matrix form as it follows:

$$\widetilde{H}\ \widetilde{I}_{cq} = \widetilde{B} \qquad\qquad (A4)$$

Where, $\widetilde{H}$ , $\widetilde{I}_{cq}$ and $\widetilde{B}$ are as defined in the section 4.1.

## 2    *Optimisation of the locations*

For the different derivatives to the locations, we obtain:

$$\frac{\partial \Delta P_i}{\partial h_j} = \begin{cases} 2R\,I_{cqi}\,I_s F_q(h_i) - 2\,R\,I_{cqi}\sum_{k=1}^{i-1} I_{cqk} - R\,I_{cqi}^2 & if \ \ j = i \\[2mm] 0 & if \ \ i \rangle j \\[2mm] -\,2\,R\,I_{cqi}\,I_{cqj} & if \ \ i \langle j \end{cases}$$

$$\frac{\partial \Delta E_i}{\partial h_j} = \begin{cases} 2RI_{cqi}TL_f I_s F_q(h_i) - 2RTI_{cqi}\sum_{k=1}^{i-1} I_{cqk} - RTI_{cqci}^2 & if \ \ j = i \\[2mm] 0 & if \ \ i \rangle j \\[2mm] -\,2RTI_{cqi}I_{cqj} & if \ \ i \langle j \end{cases}$$

And then:

$$\frac{\partial \Delta P}{\partial h_j} = 3(2\,R\,I_{cqj}\,I_s\,F_q(h_j) \; - \; 4\,R\,I_{cqj}\sum_{k=1}^{j-1} I_{cqc} \; - \; R\,I_{cqj}^2) \tag{A5}$$

$$\frac{\partial \Delta E}{\partial h_j} = 3(2RTL_f I_{cqj}I_s F_q(h_j) - 4RTI_{cqj}\sum_{k=1}^{j-1} I_{cqk} - RTI_{cqj}^2) \tag{A6}$$

At last $\partial S / \partial h_j = 0$ gives after having ordered it:

$$F_q(h_j) = \frac{2(k_p + k_e T)}{(k_p + k_e TL_f)I_s}\sum_{k=1}^{j-1} I_{cqk} \; + \; \frac{k_p + k_e T}{2(k_p + k_e TL_f)I_s}I_{cqj} \tag{A7}$$

## References

[1]    R. F. Cook: Calculating Loss Reduction Afforded by Shunt Capacitor Application, in IEEE Trans. Power App. and Syst.; PAS-83 (1964); 1227-1230

[2]    Nelson E. Chang: Determination of Primary-Feeder Losses, in IEEE Trans. Power App. and Syst.; PAS-87, n°12 (1968); 1991-1994

[3]    Nelson E. Chang: Location Shunt Capacitors on Primary Feeder for Voltage Control and Loss Reduction, in IEEE Trans. Power App. and Syst.; PAS-88, n°10 (1969); 71-77

[4]    Nelson E. Chang: Generalized Equations on Loss Reduction with Shunt Capacitor, in IEEE winter meeting New York, Jan. 30-Feb. 4 (1972); 2189-2195

[5]    J. J Grainger, S. H. Lee, A. M. Byrd, K. N. Clinard: Proper Placement of Capacitors for Losses Reduction on Distribution Primary Feeders, in Proc. Amer. Power Conf., n°42 (1980); 593-603

[6]    J. J Grainger, S. H Lee: Optimum Size and Location of Shunt Capacitors for Reduction of Losses on Distribution Feeders, PAS-100, n°3 (1981); 1105-1118

[7]    J. J Grainger, S. H Lee: Capacity Release by Shunt Capacitor Placement on Distribution Feeders: A New Voltage-dependent Model, in IEEE Trans Power App. Syst., PAS-101, n°5 (1982); 1236-1244

[8]    H. Duran: Optimum Number, Location, and Size of Shunt Capacitors in Radial Distribution Feeders, A Dynamic Programming Approach, in IEEE Trans Power App. and Syst., Vol. 87, n°8 (1982); 10-13

[9]    Tharwat H. Fawzi & Al: New Approach for the Application of Shunt Capacitors to the Primary Distribution Feeders, in IEEE Trans. Power App. and Syst., Vol. 102, n°1 (1983);10-13

[10]   M. Ponnavaiko, K. S. Prakasa Rao: Optimal Choice of Fixed and Switched Shunt Capacitors on Radial Distributors by the Method of Local Variations, in IEEE Trans. Power App. and Syst., Vol. 102, n°6 (1983), 1607-1615

[11]   M. E. Baran, F. F. Wu: Optimal Capacitor Placement on Radial Distribution Systems, in IEEE Trans. Pow. Deliv., Vol.4, n°1(1989); 725-734

[12]   M. H. Haque: Capacitor Placement in Radial Distribution Systems for Loss Reduction, in IEE Proc. Gener. Transm. Distrib. Sept., Vol. 146, n°5 (1999), 501-505

[13]   S. F. Mekhamer, M. E. El-Hawary, S. A. Soliman, M. A. Moustafa, M. M. Mansour: Reactive Power Compensation of Radial Distribution Feeders: A New Approach, in Trans. Dist. Conf. Eexhib 2002, Asia Pacific IEEE/PES 6-10 Oct, Vol. 1 (2002), 285-290

[14]   S. F. Mekhamer, M. E. El-Hawary, S. A. Soliman, M. A. Moustafa, M. M. Mansour: New Heuristic Strategies for Reactive Power Compensation of Radial Distribution Feeders, in IEEE Trans. Pow. Deliv. Oct., Vol. 17, n°4 (2002); 1128-1135

[15]   Kenji Iba: Reactive Power Optimisation by Genetic Algorithms, in IEEE Trans. Pow. Sys., May 4-7, Vol. 9, n°2 (1994); 685-692

[16]   Srinivasan Sundhararajan, Anil Pahva: Optimal Selection of Capacitors for Radial Distribution Systems Using a Genetic Algorithm, in IEEE Trans. Pow. Sys. August, Vol. 9, n°3 (1994); 1499-1507

[17]   Mohamed A. S. Masoum, Marjan Ladjevardi, Akbar Jafarian and Ewalds F. Fuchs: Optimal Placement, Replacement and Sizing of Capacitor Banks in Distorted Distribution Networks by Genetic Algorithms, in IEEE Trans. Pow. Deliv. October, Vol. 19, n°4 (2004); 1128-1135

[18]   Benemar Alencar de Souza, Helton do Nascimento Alves, B. A. de Souza, H. N. Alves, H. A. Ferreira: Microgenetic Algorithms and Fuzzy Logic

Applied to Optimal Placement of Capacitor Banks in Distribution Networks, in IEEE Trans. Pow. Sys. May, Vol. 19, n°2 (2004); 942-947

[19] H. D. Chiang, J. C. Wang, O. Cockings, H. D. Shin: Optimal Capacitor Placements in Distribution Systems; Part 1: New Formulation and the Overall Problem, in IEEE Trans. Pow. Deliv., Vol. 5, n°2 (1990); 634-642

[20] H. D. Chiang, J. C. Wang, O. Cockings, H. D. Shin : Optimal Capacitor Placements in Distribution Systems; Part 2: Solution Algorithms, in IEEE Trans. Pow. Deliv., Vol. 5, n°2 (1990); 643-649

# Investigation of the Possibilities for Interdisciplinary Co-operation by the Use of Knowledge-based System

**Ágnes Szeghegyi, Ulrich H. Langanke**

Keleti Károly Faculty of Economics, Budapest Tech Polytechnical Institution
szeghegyi.agnes@kgk.bmf.hu, langanke.ulrich@kgk.bmf.hu

*Abstract: Each problem always originates from constraints. The decision is a response to the challenges by the environment. In order to chose appropriate decision support techniques the structural complexity of the problem has to be determined. The aim of the application of knowledge based systems is to obtain decision support. In this paper the application of the system 'Doctus' is illustrated and exemplified in connection with processing the problem of the potential co-operation between the industrial companies and institutes of higher education. The analysis was carried out by the application of inductive and deductive inference procedures taking into account the requirements of the companies and the abilities and skills of the higher educational institutions. The assessment of the results obtained may generate further dilemmas for the solution of which appropriate knowledge bases can be brought about or the already existing ones have to be refined.*

*Keywords: Knowledge-based Systems, Knowledge-based Technologies, Deductive Graph, Deductive Reasoning, Model Graph, Inductive Reasoning, Interdisciplinary Co-operation*

## 1   Introduction

To what extent can human thinking be substituted by computers?

By seriously establishing the idea of automating abstract mathematical proofs rather than merely arithmetic, Turing greatly stimulated the development of general purpose information processing only in 1936. Previously, Hilbert had emphasized between the 1890s and 1930s the importance of asking fundamental questions about the nature of mathematics. Instead of asking 'is this mathematical proposition true?' Hilbert wanted to ask 'is it the case that every mathematical proposition can in principle be proved or disproved?' This was unknown, but Hilbert's feeling, and that of most mathematicians, was that mathematics was indeed complete. Gödel destroyed this hope by establishing the existence of mathematical propositions which were undecidable, meaning that they could be neither proved nor disproved.

The next interesting question was whether it would be easy to identify such propositions. After Gödel, Hilbert's problem was re-phrased into that of establishing decidability rather than truth, and this is what Turing sought to address. In the search for an automatic process by which mathematical questions could be decided, Turing envisaged a thoroughly mechanical device, the Turing machine (1936), in fact a kind of 'glorified typewriter'. Its significance arises from the fact that it is sufficiently complicated to address highly sophisticated mathematical questions, but sufficiently simple to be subject to detailed analysis [1]. Turing's universal computer can simulate the action of any other, in certain sense. This is the fundamental result of computer science. Indeed, the power of the Turing machine and its cousins is so great that Church [2] and Turing [3] framed the 'Church-Turing thesis,' to the effect that "*Every function 'which would naturally be regarded as computable' can be computed by the universal Turing machine.*" This thesis is unproven, but has survived many attempts to find a counterexample, making it a very powerful result.



Figure 1
Turing's 'glorified typewriter', the Turing Machine

In the possession of the idea of Turing Machine the computational complexity of a problem can be determined by the number of steps a Turing machine must make in order to complete any algorithmic method to solve the problem. If an algorithm exists with the number of steps given by any polynomial function of the amount of information given to the computer in order to specify the problem then the problem is deemed tractable and is placed in the complexity class 'P'. If the number of the necessary steps in solving a task rises exponentially with this information, then the problem is hard and is in another complexity class. There is an even stronger way in which a task may be impossible for a computer. Such problems are termed uncomputable. The most important example is the 'halting problem'. A feature of computers familiar to programmers is that they may sometimes be thrown into a never-ending loop.

According to the above outlined situation it can be stated that recent development of Information Technology made it possible the partial or in certain very special cases even the full automation of the process of human thinking, more formally speaking, at least in the realm of '*P class problems*'. In our days the decision supporting systems embodied in expert systems as software running on common computers are indispensable requisites for managers. Normally these decision supporting systems are used for managing decision dilemmas at the strategic level of decision making [4]. Their application can be extended for the rather 'soft' fields of application as social sciences as well as for the rather 'hard' subject areas of natural and technical sciences. In this paper, as an application example of the knowledge based technologies, analysis of the possible co-operation between the industrial companies and the higher educational institutions is presented.

The role and position of the higher educations seems considerably vary worldwide in these years. The situation of the Hungarian educational institutions is even more dubious: its staff frequently has to cope with often controversial demands that also vary in time. The basic knowledge obtained at the educational institutions has to be further developed in the practice. By our days this knowledge became of strongly technological nature, and it is expanded and enriched mainly by the research carried out by private companies. As a consequence its key elements ceased to be 'public properties', and became 'private' ones. The great majority of this technological knowledge does not even reach the educational institutions, and it is hopeless to acquire them by separate research conducted in these institutions. The employees of the industrial companies are provided with up-to-date technological knowledge in special education carried out within the companies. This situation automatically devalues the output of the traditional educational institutions that is the student who is provided with certain particular lexical knowledge that cannot be fresh and really up-to-date. The costs and time requirements of the post-education within the companies that is needed for obtaining really useful labor force also mean some burden to the industry that naturally wish to reduce them.

The education that is specific to the needs of the various companies can be brought about via the co-operation of the industrialists and the educational institutes by harmonizing the appropriate goals and conditions of the education. This harmonization has to result in the reduction of the financial burden of the postgraduate training, in the production of more marketable graduate students, and more popular institutional organizations of widely recognized reputation. It can also be expected that increasing proportions of the financial resources of the higher education will be covered directly by the students in Hungary.

The burden of the post education mainly can be reduced by deliberately providing the students rather with skills than with particular lexical knowledge in the education. On the basis of these skills the graduated students become able to learn in fast and efficient manner so the duration and the costs of the post-education can be reduced. For the development of these skills the higher educational institutions

have to operate as intellectual training camps the actual program of which also corresponds to the daily activities conducted in the industrial companies. In our view this practical demand is the key element on the basis of which the working connection between the educational institutes and the industry can be revitalized.

The up-to-date lexical part of the common, public knowledge to be acquired by the students during the education obtains special emphasis in this phase. The industrial development is accompanied by intensive standardization. The crystallization of the industrial standards is a long process in which various, often concurrent, competing companies take part simultaneously. It is much more expedient and advantageous for the educational institutions to take part in this process than only compiling the already existing, well crystallized standards. In this manner the graduate students can obtain knowledge and skills that are immediately marketable in the industry.

# 2    The Role of Computers in Decision Making

| Management decision making level | Information required for | Type of CIS support |
|---|---|---|
| Top Strategic | Planning long-term policy decision and planning | Decision Support System (DSS) |
| Middle Tactical | Controlling comparing results of operations with plans and adjusting plans or operations accordingly | Management Information System (MIS) |
| Lower Operational | Operating maintaining business records and facilitating the flow of work in a project | Data Processing System |

Table 1

The role of computers in decision making [8]

## 2.1    Data Processing Systems

A Data Processing System contains a series of procedures that deal with one or more types of relevant business transactions. Examples for the processed transactions are the payrolls, making out bills, acceptable bills, payable bills, stock control, purchase and others. Each data processing system is typically named as an application. Within the organization these systems support a wide scale of transactions. The required decision contains an individual transaction.

The data processing systems can be described as managers of many transactions, in turn they provide plenty of data and it is possible to have access to the required information. It is a fact that the data can increase as an avalanche. Here is arising the problem in the decision making of the management concerning the possible application of the systems data files namely the data volumes are too large. It is impossible to reach a general interpretation as the details are too much. Accordingly the capacity of the computers is applied to data summarizing and interpretation of reports and information [8].

## 2.2    Management Information Systems

A management information system summarizes and selects data from a large volume of data file and in such a way it makes a report to the system outputs. This process is often quoted as 'lock out reporting' or 'managing with lock out'. In the system the managers insert data selection criteria within the computer program on the basis of the content. Management information systems help especially in handling of problems with known structures soluble on the basis of the past experiences. Problem solving of such type is often associated with the tactical level of the management. The 'tactical' term refers to such starting operations within which the decisions can be made on the basis of known relations, rules, laws, that is as in the case of structured problems. There are reliable precedents applicable in decision making. In effect the knowledge and experience of the manager qualifies the problem as structured and the managing area as tactical. The nature of problems and the suitable solving method often depends on the situation, on the experience and knowledge on the given area of the problem solver [8].

## 2.3    Knowledge-based Systems (Decision Support Systems)

The purpose of the application of knowledge based systems is rather decision making than information processing. For solving problem the degree of an existent structure represents the dividing criterion for the application of decision support techniques. A knowledge based systems is most effective in the managing of semi-structured problems. The abilities of such systems are usually applied on the managing level of strategic planning. The tactical and strategic planning is unfortunately often defined with the term time range instead of its relative duties. The tactical planning deals with the starting situation including the application of know principles, rules, and laws. The strategic planning deals with situations including certain elements of the undoubtedly not predictable unknown. Strategic planning includes the future but not necessarily the long range future. The identifying characteristic can be catched rather in the semi-structured nature of the problem to be solved and in the uncertainty of the future [8].

The structurality of a problem is determined by that person who perceives and solves the problem. Therefore it is impossible to define the structurality of a problem in an absolutely correct manner. The degree of structurality of a problem is function of the knowledge, experience of the problem solver. The principle of the definition of problem structure is especially important in understanding of the knowledge based systems.

A knowledge based system is such a computer tool that is used by the manager in connection with his/her problem solving and decision making activity, helping him/her in decision making. A person has to define the problem structure and the criteria in connection with the problem evaluation. The manager makes certainly the decisions and solves the problems. Creation of alternative possible solutions is the duty of human creativity. The knowledge based system can be regarded as a problem solving tool kit that can be used in the 'valuation of alternatives' phase of the problem solving process [9].

Application of knowledge based systems requires 'new technology' in the sense that knowledge-based technology is qualitatively different to, and does not fall under the technological development trends of programming. It came about during the researches of the artificial intelligence in connection with human problem solving [10]. The aim of the artificial intelligence researches is the development of intelligent computer system. This is an artificial intelligence program that solves the problems such a way, and behaves so, as persons. On this basis it could be named as 'intelligent behavior'. It tries to create systems that imitate thinking and acting habits of persons. An artificial intelligence program must have attributes characteristic of the human problem solving behavior, that is in case of complicated problems having alternative possibilities with effective problem solving ability, with communication ability, with ability to handle uncertain situations, with ability to handle exceptions and with learning talent.

On the basis of the above ideas it can be stated that the knowledge based systems are artificial intelligence programs with new program structure appropriate to processing symbolical information. The information are processed with reasoning, with application of heuristics. The quality of problem solution depends on the quantity and quality of information available on the relevant area. The programming style is declarative [11]. Accordingly, the task of a knowledge based system operating on the knowledge based technology is not the realization of some mapping between the input and output domains and obtaining data adequate to the given conditions. Instead of that, its task is generation of statements on expert level on basis of the given data proposal [12]. That is the decision support system does not want replace the decision maker. It is satisfied by accelerating and supporting the process of human thinking. Further it can be said that it alloys the advantageous attributes of the human person (human knowledge) and the computer (artificial knowledge) [4].

Its advantages are listed as follows:

- Integration of knowledge of more experts, so it can make more perfect decision, can give better-founded proposal than any individual expert.

- Quick operation, the more hours for human decision can be shortened to a few minutes.

- In problem solving it comes always to identical solutions in contrast to the human decisions that normally are motivated by external conditions, atmosphere, and other effects.

- It can switch over flexibly from a problem to other one.

- Its application has no limits neither in time nor in space. Human expert can be found at definite date at definite place. He/she is not always available. He can step out taking away his knowledge.

- Its amortization is a gradual process while the costs in connection with human experts does not decrease.

- The expertise is summarized and stored on highest level.

- Knowledge Based Systems (KBS) are continuously developable open systems.

- They have modeling ability.

- They are especially efficient in the area of education/teaching.

Limits of the application of KBSs:

- Their unjustified application, forcing can lead to faults.

- They are always only able to solve problems of a narrow, mainly special field. The human person is all-round.

- They follow only given rules, cannot think with common sense.

- The human expert is ready to admit his/her incapability to give a real answer to a problem. In contrast, the system does not sense the limits of its applicability, so in situations not defined by rules it can come to incorrect solutions.

- In case of a system including too much rules it is difficult to check the success of relevant rules, the process of the program can slow down. In case of too few rules it can come to unreasonable conclusions.

The structure of the knowledge based system is custom-built. The knowledge basis includes in explicit form, separated from the other system components, the knowledge and terms describing the special field. The separated knowledge is easily available also for others.

The knowledge collection and arranging is a process describing the expert knowledge [10].

First step of building the knowledge basis is the determination of the decision dilemma [4].

# 3 Simulation Examination and Simulation Results

Problems always arise from difficult situations. The decision is an answer to the challenge of the environment. A decision maker can be able to make decisions only when he has recognized the difficult situation, has looked for the solution and has power and authorization to make decisions. For the decision maker the real difficulty arises when the decision attributes, their values and the relations between the decision attributes are uncertain [4]. Situations in which the object of decision is a great stake do not mean such difficulties.

In this particular case the decision dilemma is following:

What chances have the higher educational institutions to cooperate with companies? In this paper this problem field and the relevant decision possibilities, alternatives are examined via applying the 'Doctus' knowledge based shell system by deductive and inductive reasoning.

## 3.1 Reasoning on Basis of Rules (Deductive Reasoning)

### 3.1.1 Knowledge Base

The knowledge base building denotes the decision preparing and the decision support denotes the decision proposal [4]. The decision attributes are shown in the first column of Table 2. The values pertaining to the decision attributes are given on nominal or ordinal scale. They are discrete terms. Their formulation and quantity is a delicate task, as at determining the rules of the depending attribute not suitable formulation and excessive refining of the values can cause problems. In such case it is necessary a subsequent refinement, and correction [4]. The situation could be made easier by presenting on interval scale, while the application of fuzzy sets should eliminate all difficulties as the fuzzy logics deals with the mathematical handling of the uncertainty of the linguistic terms. Fuzzy systems can successfully be used for various purposes e.g. in motion control of wheeled systems [5], identification of various physical systems [6], and control of test devices [7]. Handling of the uncertainty will be realized most effectively in the fuzzy expert systems [13].

| Result | No | promising | Yes |
|---|---|---|---|
| Competence | Confuse | clear | |
| Partner relationship | Subordinate | equal | |
| References | of dubious origin | missing | convincing |
| Communication ability | Weak | convincing | |
| Communication possibility | Clumsy | just acceptable | good |
| Personal relationships | Problematic | acceptable | good |
| Cultural differences | Great | soluble | small |
| Capacity | Weak | to be extended | satisfactory |
| Infrastructure | Problematic | acceptable | |
| Technical background | to be developed | acceptable | |
| Human resources | to be developed | acceptable | |
| Professional teachers | weak | acceptable | excellent |
| Motivation of teachers | Burden | accepts | challenge |
| Flexibility | Inflexible | flexible | |
| Confidence | Missing | partial | exists |
| Financial possibilities | Critical | promising tendencies | good |
| Separable resources | Hopeless | possible | already solved |
| Sponsors | Missing | little | considerable |
| Manager's attitudes | Burden | passive | active support |
| Manager's motivation | Prestige | self-realization | |
| Manager's skills | Burden | passive | active |
| Manager's social commitments | Little | much | |
| Investment costs | Little | just acceptable | much |
| Image | Bad | acceptable | good |
| Controllability | Casually | posteriorly | continuous |
| Number of competitors | Missing | little | much |
| Strength of competitors | Small | middle | great |
| Economic advantage | Little | uncertain | considerable |
| Available knowledge | Negligible | useful | |
| Experience | Beginner | experienced | |
| Attitude | Constrained | accepted | readily |
| Activation of results | in a single field | in various fields | |
| Problem identification | not real | suspicious | real |
| Way of realization | False | acceptable | |
| Solution | Disquieting | just acceptable | excellent |
| Position in competition | Difficult | advantageous | |

Table 2

Decision criteria and typical 'values'

### 3.1.2   Deductive Graph

It is an empirical fact that the expert is not able draw together more than 3-4 decision attributes with 'if..then' rules. Therefore it is practical to classify the attributes in graph form, arranging them in a hierarchy. The deductive graph presents the dependence conditions of the attributes [4].
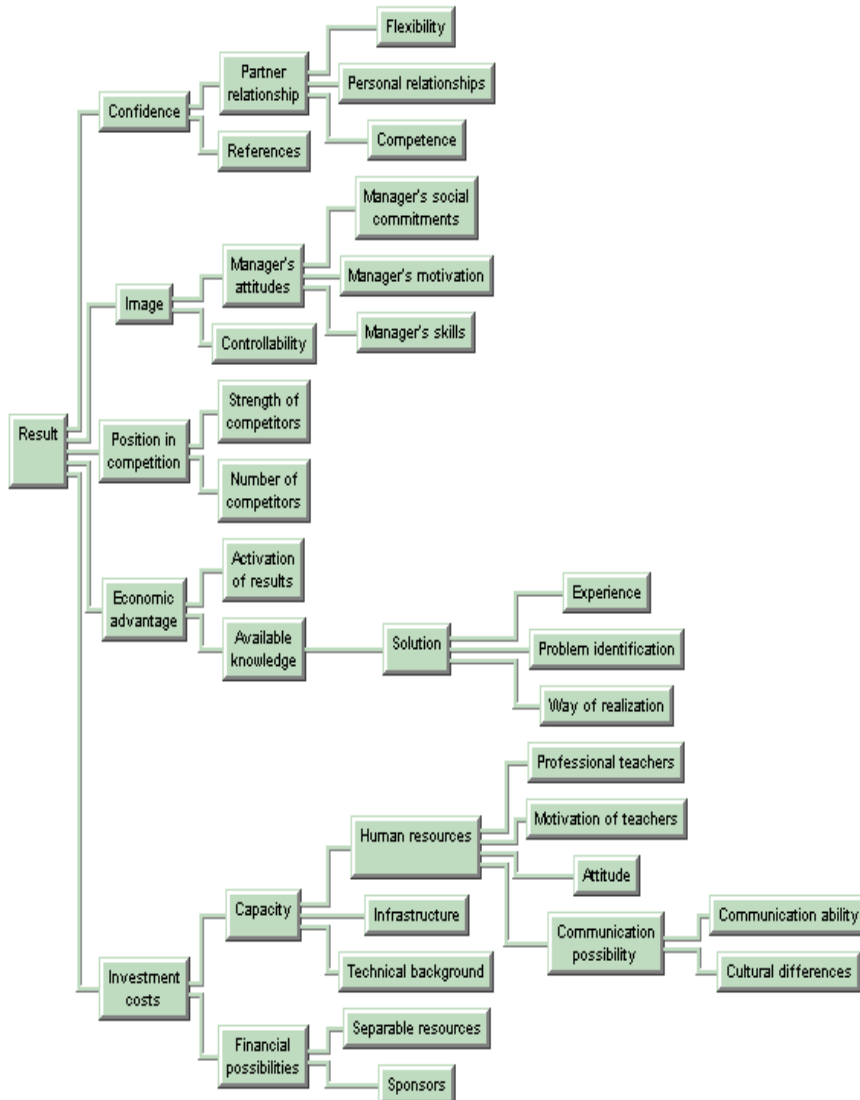


Figure 2
Deductive graph

### 3.1.3    Rules

The conclusion runs along the graph upwards from below, from the input attributes through the depending medians to the top of the decision tree, to the final decision [4].

### 3.1.4    Cases

Next step is presenting the cases.

| | Partner relationships | References | Communication ability | Communication possibility | Personal relationships | ➔ |
|---|---|---|---|---|---|---|
| **A** | subordinate | missing | convincing | just acceptable | problematic | ➔ |
| **B** | equal | dubious origin | weak | clumsy | acceptable | ➔ |
| **C** | equal | missing | convincing | just acceptable | good | ➔ |
| **D** | equal | convincing | weak | just acceptable | problematic | ➔ |
| **E** | subordinate | dubious origin | convincing | just acceptable | acceptable | ➔ |
| **F** | equal | missing | weak | clumsy | good | ➔ |
| **G** | subordinate | convincing | convincing | good | problematic | ➔ |
| **H** | equal | dubious origin | weak | clumsy | acceptable | ➔ |
| **I** | equal | missing | convincing | just acceptable | good | ➔ |
| **BMF** | subordinate | convincing | convincing | good | acceptable | ➔ |

Table 3

Cases and their values (excerpt)

As the aim of this study is to present the applicability and operation of the knowledge based systems, and since no authentic input data for the real higher educational institutions are available, the institutions denoted with are fictitious institutions. The values of input attributes are casually chosen. The final BMF is included as a real institution.

Reasoning consists in the activation of rules. The shell is reasoning on basis of the input rules through the depending attributes on the output attributes, in the present case on the 'result'.

| | Result |
|---|---|
| **A** | promising |
| **B** | promising |
| **C** | no |
| **D** | yes |
| **E** | no |
| **F** | no |

| G | promising |
|---|---|
| **H** | yes |
| **I** | yes |
| **BMF** | promising |

Table 4

Results of the deductive reasoning

## 3.2   Reasoning on the Basis of Cases (Inductive Reasoning)

On the basis of known cases it is possible to come to general conclusion regarding the decision rules.

The knowledge base includes the decision attributes and their values (Table 2).

### 3.2.1   Model Graph

The system generates the model graph, using the inputs of the cases.



Figure 3

Model graph

The 'if..then' rules constituting the base of decision making are readable from the model graph. It also is visible that the qualifying attributes influence the decision.

Through the application of the inductive method it is visible that in reality the weight of decision attributes is different. With regarded to this the informativy degree gives a numerical value that reflects the gain originating from the given information. The informativy degree can have a value from '0' to '1'.

**Conclusions**

With regard to the possibilities of interdisciplinary cooperation the conclusion is more reliable if more experts take part in building up the knowledge base, in the creation of the rules.

**Deductive Reasoning**

On the basis of the performed analysis regarding the BMF, the 'Doctus' knowledge based system gives the 'promising' result in the case of the rule based reasoning. This means the possibility of cooperation for the realization of which there are still tasks to be done. With full knowledge of these, not striving for completeness, further decision dilemmas can arise as follows (it is possible to repeatedly build up the knowledge based on them):

- What steps are to be made for the necessary realization?

- Are these capable to be realized?

- Is it worth to realize them?

- To what extent are they compatible with the strategic plans of BMF?

- Etc…

**Inductive Reasoning**

Parameters of particular cases were not available. Supposing that the conclusions determined by the deductive method are acceptable, they were considered as cases. On the basis of the model graph the determining attribute is the 'Strength of competitors'.

The companies can choose from the supply of the market, so the market position of the higher educational institutes effects mostly the realization of the cooperation.

The 'strength of competitors' compared to the given educational institution influences the conclusions as follows:

- If small, the cooperation will be realized.

- If middle, a further attribute will be the quality of 'personal relationships'.

- If great, a further attribute the decision will be influenced by the 'controllability' of the higher educational institutions.

**References**

[1]     Andrew Steane: Quantum computing, Department of Atomic and Laser Physics, University of Oxford, Clarendon Laboratory, Parks Road, Oxford, OX1 3PU, England, July 1997

[2]     A. Church: An Unsolvable Problem of Elementary Ónumber Theory, Amer. J. Math. 58 345-363, 1936

[3]     A. M Turing: On Computable Numbers, with an Application to the Entschneidungsproblem, Proc. Lond. Math. Soc. Ser. 2 42, 230); see also Proc. Lond. Math. Soc. Ser. 2 43, 544), 1936

[4]    Szeghegyi-Velencei: Üzleti döntéshozásra alkalmas tudásalapú döntéstámogató rendszerek BMF-KGK-4007, Budapest, 2003 (in Hungarian)

[5]    Gy. Schuster: Simulation of Fuzzy Motion Controlled Four-wheel Steered Mobile robot, IEEE Intenational Conference on Intelligent Engineering Systems (INES '97), September 15-17, 1997, pp. 89-94, ISBN 0-7803-3627-5

[6]    Gy. Schuster: Fuzzy Approach of Backward Identification of Quasi-linear and Quasi-time-invariant Systems, in Proc. of the 11[th] International Workshop on Robotics in Alpe-Adria-Danube Region (RAAD 2002), June 30-July 2 2002, Balatonfüred, Hungary, pp. 43-50,ISBN 963 7154 09 04

[7]    Gy. Schuster: Adaptive Fuzzy Control of Thread Testing Furnace, in Proc. of the ICCC 2003 IEEE International Conference on Computational Cybernetics, August 29-31, Gold Coast, Lake Balaton, Club Siófok, Siófok, Hungary, pp. 299-304, ISBN 963 7154 17 5

[8]    W. E. Leigh, M. E. Doherty: Decision Support SystemsSouth-Western Publishing Co., Cincinnati, Ohio, 1996

[9]    J. Raggett, W. Bains: Mesterséges intelligencia A-Z, Akadémiai kiadó, Budapest, 1992 (in Hungarian)

[10]   Starkné Werner Ágnes: Mesterséges intelligencia-szakértői rendszerek, Veszprémi Egyetemi Könyvkiadó, Veszprém, 1997 (in Hungarian)

[11]   J. Kelemen, S. Nagy: Bevezetés a mesterséges intelligencia elméletébe, tankönyvkiadó, Budapest, 1991 (in Hungarian)

[12]   Sántáné Tóth Edit: Tudásalapú technológia, szakértő rendszerek, Miskolci Egyetem Dunaújvárosi Főiskolai Kar, Dunaújváros, 1997 (in Hungarian)

[13]   Á. Szeghegyi: Bizonytalanságok kezelésére is alkalmas tudásbázisú döntéstámogató rendszerek (manuscript in Hungarian)

# Notes on the Components for Intelligent Tutoring Systems

## Ladislav Samuelis

Department of Computers and Informatics, Technical University of Košice
Letná 9, 04200 Košice, Slovakia
ladislav.samuelis@tuke.sk

*Abstract: The aim of this contribution is to present briefly several views on the components for the intelligent tutoring systems (ITSs) and to analyze the concepts behind them. We will briefly analyze the concepts provided by the classical machine learning (ML) approaches and then we discuss the role of the components in distributed environment. Finally we synthesize the converging trends of the ITSs technologies and try to outline further possible research directions in building ITS from components.*

*Keywords: intelligent tutoring systems, components, e-learning, inductive and deductive inference*

## 1    Background and Motivation

What is the idea behind building ITSs? Human beings built tools for centuries for remembering and manipulating with data. The components of those systems were based always on the contemporarily available tools and technologies. In other words, human beings are trying continuously to automize cognitive activities, which are understood and could be described by algorithm. The advent of personal computers enabled widening the automation of further cognitive activities.

Learning is a very complex process and in fact how people learn, is not fully understood till the recent time. We know unfailingly that the learning efficiency varies significantly from individual to individual. For more than 6000 years the human civilization learned in some groups. Personal computers do not only support automation of cognitive processes but also enable computer-supported group personalized and learning. ITSs became gradually complex software systems and the reusability of their components plays crucial role in their sustainability and further evolution. This idea motivated the development of the contemporary Learning Management Systems (LMSs) with the aim to accept and run distributed courses. In order to ensure the compatibility between the LMSs and the distributed courses, the ADL initiative elaborated the implementation of

the Shareable Content Object Reference Model (SCORM [1]). It is obvious that the more reliable components are available, the more sustainable the education is. This metalevel of abstraction consists of two components: LMS and the SCORM compatible courses. The next chapter emphasizes several remarks on the evolution of ITSs.

## 2 Overview of the Evolution of Intelligent Tutoring Systems

### 2.1 A Short Historical Perspective

As illustrated in the following Figure 1, the advent of ITSs is around the late sixties.



Figure 1
The evolution of e-learning [2]

The differentiation of the two main streams is based on the one hand side by modeling the human cognition by the application of the deduction and the induction concepts. Deduction based applications evolved later, around the late eighties, into the application of the PROLOG [3] and similar programming languages in building expert systems. It was also the era of the fifth generation computers initiative [4] launched by the Japans. On the other hand-side were the dedicated applications, based on procedural languages, for supporting pedagogical

purposes, like tracking the assignments and their evaluation. There are many monolithic tutoring (or e-learning) systems, which aim is to conduct and track the students' progress in special areas.

We observe in Figure 1 that computers have been used in education for over 30 years. Computer Based Instruction (CBI) and Computer Based Training (CBT) are the predecessors of ITSs. The field of *automating tutoring* (e.g. Computer Aided Instruction – CAI or Intelligent Computer Aided Instruction – ICAI) was one of the first most fruitful fields, which served aims of AI researchers. Application of machine learning algorithms in tutoring has a long history. The field of automating tutoring was one of the first most fruitful fields, which served aims of *artificial intelligence* (AI) researchers. These were more or less monolithic systems with low level of flexibility.

Splitting of the practice into two streams has begun since the beginning of the late sixties mainly due to the different research orientation. In further decades we observe a gradual convergence and synthesis mainly due to the penetration of the Internet.

The upper stream is characterized by the application of the artificial intelligence algorithms and computational architectures, such as production rules systems, generative grammars, Bayesian networks, Markov models, neural networks, higher-order semantic spaces, fuzzy control systems, and non-linear dynamical systems. With these methods ITSs use tutorial dialogue to adaptively respond to the learner's motivational states, and efficiently bridge the man-machine interface, adapting to a learner's profile. One of the aspects of the profile of the learner or e-learning user is a user model. The user model can be used for solving the problem of e-learning users' cognitive load.

The bottom stream is characterized by the application of actual software technologies to support functionalities which a human teacher might do: select (or generate) appropriate material, set up the exercise, monitor student activity, give hints during exercises and feedback afterwards, understand why students make mistakes, customize presentation style to the student's style, ask and answer questions.

## 2.2    Some Selected Definitions of the ITSs

The most general way to describe ITS, is to say that it is the application of AI to the education. During the several last decades the penetration of computers essentially influenced the architectures of the so-called 'intelligent tutoring' systems. It is fashionable recently to mark sophisticated software systems with this attribute. The definition of *intelligence* is context dependent and we will not deal with the phenomenon of intelligence. The aim of this paragraph is to present several definitions of the ITSs.

- 'ITSs are computer software systems that seek to mimic the methods and dialog of natural human tutors, to generate instructional interactions in real time and on demand, as required by individual students. Implementations of ITSs incorporate computational mechanisms and knowledge representations in the fields of artificial intelligence, computational linguistics, and cognitive science.' [5]

- 'Broadly defined, an intelligent tutoring system is educational software containing an artificial intelligence component. The software tracks students' work, tailoring feedback and hints along the way. By collecting information on a particular student's performance, the software can make inferences about strengths and weaknesses, and can suggest additional work.' [6]

- 'In particular, ITSs are computer-based learning systems which attempt to adapt to the needs of learners and are therefore the only such systems which attempt to 'care' about learners in that sense. Also, ITS research is the only part of the general IT and education field which has as its scientific goal to make computationally precise and explicit forms of educational, psychological and social knowledge which are often left implicit.' [7]

We may distil from these definitions, that they stress the expert systems' point of view and they consider ITS as a monolithic system.

# 3    Notes on the Theoretical Foundations for Intelligent Tutoring Systems

It is generally accepted (as it is seen also from the above mentioned definitions) to refer to an ITS if the system is able to:

- build a more or less sophisticated model of cognitive processes,

- adapt these processes consecutively and

- control a question-answer-interaction.

Conventionally, ITS provides individualized tutoring or instruction and has the following 4 models or software components [8]

- knowledge of the *domain* (i.e. knowledge of the domain expert, refers to the topic or curriculum being taught)

- knowledge of the *learner* (e.g., what he/she knows, what he/she's done, how he/she learns, …)

- knowledge of *teacher strategies* (or pedagogical issues i.e. how to teach, in what order, typical mistakes and remediation, typical questions a student might ask, hints one might offer a student who is stuck)

- *User interface* (i.e. interactive environment interface)

We observe in general the following two approaches in building components of the ITS:

*Deductive inference*: In this case we are asking, whether particular statement(s) coincide with a general theory. Deductive systems provide inference engines to conduct the deduction. The task is to prove that the input predicates fulfill the general theory. E.g. applications based on Prolog programming language represent such systems.

*Inductive inference*: In this case we are building, or inferring, a general rule from the special input examples. In other words, we are seeking a rule, which covers the particular set of examples. We have to note, that in general, it is impossible to infer a rule, which will satisfy the intended algorithm. There may exist always one more example, which could modify the recently obtained rule. Inductive inference methods are used for inferring and maintaining the learner model. Application areas are as follows:

- Programming by examples – PBE [9] (e.g. programming of robots)
- Inference of grammars [10]
- Pattern recognition [11]

Recently we are witnesses of the research in components of ITSs, which are distributed over the computer network. As mentioned above, we distinguish between the learning management systems and the courses, which are accepted and embedded into these systems. The market is already matured with the LMS products and active research is conducted in the construction of reusable distributed courseware repositories. We will talk about this issue in the next paragraph.

# 4 Intelligent Tutoring Systems and the E-Learning Standards

## 4.1 The Importance of the Standards

Why are learning technology standards important? Standards in Learning Technologies (LT) will have a powerful impact on the way education will work in the near future. Whether learning takes place in a classroom or over the Internet, the relationships between educators, learners, and study materials will be greatly influenced by the development of standards for learning technology. In engineering intelligent tutoring systems important role play the reusability concept and the standards [12].

For educators, LT standards may make it easier to share course materials with colleagues, and to use materials produced by a much wider range of publishers without worrying about those materials being incompatible with their existing course management software. On the other side, the applications that are developed based on LT standards will also influence the way in which they teach.

For students, it may provide the ability to move between institutions, anywhere in the world, with far greater ease than it is currently possible, taking their academic record with them. The key issue is what this record contains, and who has access to it, and this will largely depend on how the standards for learner profiles are defined.

For institutions, there are clear benefits from connecting up systems for academic records, course delivery, and assessment. Provided that those standards are adopted by the key vendors of those systems, and the standards are framed in such a way that the activities and values of the institution are supported. If the majority of vendors adopt the standards and those standards don't suit education institutions, then the lack of alternative systems available on the market will either force institutions into changing their practices or purchasing expensive proprietary solutions.

For vendors the adoption of open standards widens the playing field, allowing small and medium sized companies to create education solutions that are compatible with other compliant systems. For example, any vendor could create a Virtual Learning Environment (VLE) that can deliver the same course materials as the market-leading products. This in turn leads to greater choice for institutions, educators, and students. For publishers standards mean reduced costs and time to market, as content does not need to be developed for multiple VLE platforms.

What happens if there aren't any standards? A lack of open standards results in a fragmented market for education products, reducing choice and locking users into proprietary systems. Having standards broadens the choice for end users, by allowing small and medium sized vendors to compete and to increase the range of materials available to educators and students. Rather than being forced into purchasing expensive total solutions, institutions will have the option to mix-and-match components that have the features they want, without having to worry about integration and data format issues.

## 4.2   The Role of the Metadata

Description of the knowledge by metadata is a way of storing and retrieving information from knowledge pools. This is the main mechanism of the knowledge description. The danger is that the metadata describe data in a certain context and that is why they are inherently context dependent. In other words, data and information are distinct because pieces of information exist only within the cognition of human beings.

If digital courses have to be interchanged, some metadata standards have to be adopted (like SCORM); also, the metadata has to be generated. While this is often only possible to do this manually it would be better to try to extract the metadata at least partially from the corpus of the material at issue. This is a formidable complex task, where only partial successes have been reported.

The primary developmental goals of the ITSs community are aligned with *advanced distributed learning's* (ADL) long-term vision: 'To generate, assemble, and sequence content that dynamically adapts to the learner to optimize learning. Specifically, ADL is actively engaging in research and implementation of the digital knowledge environment of the future in the areas of standards and authoring tools that give instructors the ability to create ITS functionality within a virtual training environment.'

ADL Technologies help to create new markets for training materials, reduce the cost of development and increase the potential return on investment. Platform neutrality and software reusability are considered essential for the sustained investments necessary to create the dynamic ADL environment. The ADL initiative is currently pursuing: the SCORM, which is 'a collection of specifications adapted from multiple sources to provide a comprehensive suite of e-learning capabilities that enable interoperability, accessibility and reusability of Web-based learning content'.

# 5 Course Repositories as Reusable Components

Finding the suitable component is not a trivial task. Repository systems provide key infrastructure for the development, storage, management, discovery and delivery of all types of electronic content. Repositories must provide a basic set of functions in order to provide access to learning objects and other assets in a secure environment.

There exist initiatives for the development of efficient solutions, e.g.:

- Content Object Repository Discovery and Registration/Resolution Architecture – CORDRA[13]. It is an 'open, standards-based model for how to design and implement software systems for the purposes of discovery, sharing and reuse of learning content through the establishment of interoperable federations of learning content repositories.'

- Another European initiative is the Ariadne [14] project 'A European Association open to the World, for Knowledge Sharing and Reuse. The core of the Ariadne infrastructure is a distributed network of learning repositories.'

**Discussion and Conclusions**

To sum up, the component approach supports solving the following issues of recent ITS

1  *reusability* on the courses (learning objects and learning objects metadata),

2  *design patterns for ITS* (ITS prototyping, service-oriented frameworks),

3  *adaptable software* for Personalized E-Learning Services (PELS),

4  *adaptive hypermedia*[15], and software agents in ITS.

The subject domain of engineering adaptive learning systems is considered today one of the 'hottest' in the area of the use of artificial intelligence in e-learning. This domain gathers specialists from many disciplines, like artificial intelligence, software engineering, hypermedia and web engineering. We may conclude that the trend is to assemble ITSs with components (e.g. patterns, models, frameworks) in 'software factories'.

There is no doubt that the machine learning methods will be utilized in the emerging e-learning industry. The aim of this contribution was to clarify the terminology and to summarize prevailing concepts. As the software tools and paradigms evolve in rapid manner, it is necessary to follow the existing results in the field of the artificial intelligence and software engineering in parallel in order to synthesize reasonable methods for the effective and efficient implementation purposes.

**References**

(URLs retrieved on Sept. 10, 2006)


[1]  SCORM 2004 3[rd] Edition Documentation Suite Public Draft, http://www.adlnet.gov/scorm/index.cfm

[2]  Gibbons, A. S. & Fairweather, P. G. Computer-based Instruction. (2000) In, S. Tobias and J. D. Fletcher (Eds.), Training and Retraining: A Handbook for Business, Industry, Government, and the Military. New York: Macmillan Gale Group

[3]     Ivan Bratko, PROLOG Programming for Artificial Intelligence, 2000

[4]     Edward A.Feigenbaum and Pamela McCorduck, The Fifth Generation: Artificial Intelligence and Japan's Computer Challenge to the World, Michael Joseph, 1983

[5]     Intelligent tutoring, http://www.adlnet.gov/technologies/tutoring/index.cfm

[6]     Intelligent    Tutoring    Systems    (a    subtopic    of    education), http://www.aaai.org/AITopics/html/tutor.html

[7]     Self J., The Defining Characteristics of Intelligent Tutoring Systems, Research: ITSs Care, Precisely, International Journal of Artificial Intelligence in Education (1999), 10, 350-364

[8]     Kearsley, G. 1987. Artificial Intelligence and Instruction, Reading, MA: Addison Wesley

[9]     H. Leiberman, Your Wish is My Command, Programming by example, Media Lab., MIT, Morgan Kaufrmann Publishers, 2001, ISBN 1-55860-688-2

[10]    Fu. K., Booth T. L.: Grammatical Inference; Introduction and Survey-Part I, IEEE Trans. on Systems, Man and Cybernetics, Vol. SMC-15, No. 1, Jan/Feb. 1985

[11]    Bhagat P.: Pattern Recognition in Industry, Elsevier, 2005

[12]    The Learning Object Metadata Standard, http://lttf.ieee.org/techstds.htm

[13]    An Introduction to CORDRA, http://cordra.net/introduction/

[14]    Strategy Paper, Ariadne Strategy Status, http://www.ariadne-eu.org/

[15]    Cristea A. Editorial: Authoring of Adaptive Hypermedia, On-line Journal "Education Technology and Society", Special Issue on Authoring of Adaptive Educational Hypermedia, July 2005 (Volume 8, Issue 3)

# Microgeometry Tests of 'Contradictory' Surfaces with Various Evaluation Techniques

**Gábor Fekete**

Budapest University of Technology and Economics, H-1111 Budapest, Hungary
E-mail: fucso01@gmail.com

**Sándor Horváth, Árpád Czifra**

Budapest Tech, H-1081 Budapest, Hungary, E-mail: czifra.arpad@bgk.bmf.hu

*Abstract: The article deals with a special problem area of the filtering and evaluation techniques of surface roughness measurement and calls attention to the contradictions of parameters to characterize surface microgeometry. Differences in the value of parameters are presented through tests of the operating surfaces of three disparate component types. Each chapter discusses the evaluation of measurement results according to various methods, on the basis of which it can be established that value differences are caused by the use of different filters, on the one hand, and by discrepancies between the profiles detected. From the practical point of view, it is most expedient to characterize surfaces with unfiltered parameters.*

*Keywords: surface roughness, filtering, motif, height difference correlation*

## 1 Introduction

A number of measurement and evaluation techniques are known to test surface microgeometry [1]. At the same time, parameter-based characterization is used nearly exclusively in industrial practice: it has provided general roughness ($R_a$) and other parameters since the 1930s. This list of parameters has considerably expanded by today and the applicable standard [2] accurately defines the method and parameters of evaluation.

Filtering the measured profile, that is, the separation of roughness and waviness forms an important part of microgeometry evaluation. Its significance was formerly discussed in detail [3]. This study calls attention to a special probleam area of filtering and parameter based evaluation techniques. It was already stated by [4] that parameter-based characterization is uncertain and limited in many

respects because results depend, to a great extent, on scanning length and resolution. As recommended by the standard, scanning and evaluation length are closely correlated with filtering, which also exerts a significant impact on results [5].

Our study characterizes the profiles measured by not only parameters but also by a motif-based technique and a height-difference correlation function, thereby seeking an explanation and solution for the problems arisen.

# 2    Test Methods, Problem Description

## 2.1    Measurement Technique

The tests presented in this article were performed using a Mahr Perthometer Concept type stylus instrument according to the measurement arrangement shown in Figure 1.
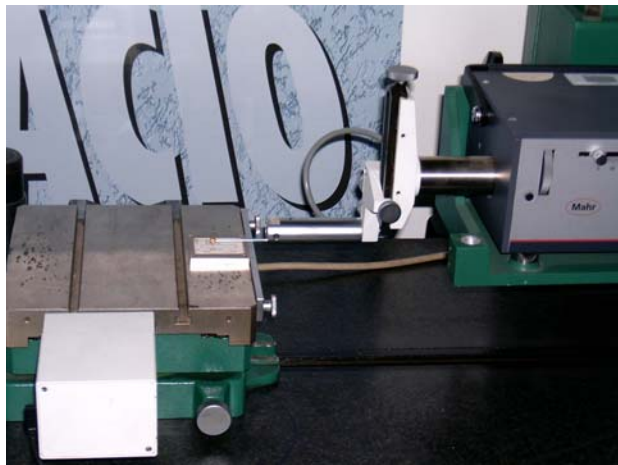


Figure 1
Measurement arrangement

The unit on the left side of the picture is the object table, on which various fitting and fixing devices can be fastened. On the right side of the picture, the unit holding and moving the stylus instrument is shown, whose main function is to drag the stylus instrument at the appropriate speed, to position it vertically and to hold it fixed. The signals detected by the stylus instrument are transmitted to a PC through the control unit, thereby data can be recorded and evaluated promptly after the measurement. Passive vibration proofing is provided by the granite table constituting the machinery unit base.

The tests to be presented in later chapters were performed using an RHT 3-50e type stylus instrument. Three different components were used as test pieces for the measurements: a turbo loading blade, a carburettor nozzle and a fuel feeder. In order to detect the problem at an appropriate level, several measurements were performed with the same setup of the measuring equipment; however, our article only presents the evaluation of each representative measurement.

## 2.2    Evaluation Technique

The results yielded by measurements were analyzed using three evaluation methods. The reason for this is that the more or less advantageous and disadvantageous applicability of each method is demonstrated and the measured surfaces are characterized from various aspects. The basic principle of evaluation methods is discussed below.

### 2.2.1    Parameter-based Evaluation Method

This procedure is the most frequently applied evaluation technique in both scientific research and industries. The reason for this is that surface topography data reveal statistical features, so the measurement points specified in an x-y-z coordinate system can be used for specifying parameters and functions by statistical means. [6]

Parameter evaluation is started from the so-called median line, which divides the profile detected by the measuring instrument in a way that the sum of the square of the profile ordinates is the smallest. The profile yielded this way is broken down into a waviness and roughness profile using a filter prescribed by the standard. Each parameter can be specified for both unfiltered and filters waviness and roughness profiles. Parameters pertaining to the roughness profile are used in practice nearly exclusively.

### 2.2.2    Motif Method

One of the basic ideas of the profile analysis described in the title is that the tribological behaviour of the operating surfaces of components is significantly affected only by surface asperities greatly protruding from the median plane. Thus, in the course of evaluation, some details of the profile detected from the surface can be disregarded under certain circumstances.

The profile section between two adjacent peaks of the profile detected is termed a motif. In the course of profile analysis, motif combinations are used to study whether the common peak of two adjacent motifs can be disregarded when calculating roughness and waviness parameters, that is, whether two adjacent motifs can be substituted by a common motif. The rules of combining and quitting motifs are described in the literature. [7] After establishing the motif combination

of the profile detected, roughness parameters will result from the evaluation of the transformed profile. [7]

### 2.2.3    Height-Difference Correlation Function

A number of people in scientific circles urge the introduction of a function-based approach to evaluate surface microtopography. Several correlation functions have been defined in this respect. The height-difference correlation function (HDCF) is defined as follows [8]:

$$C_z(\lambda) = \left\langle \left( z(x + \lambda) - z(x) \right)^2 \right\rangle, \text{ where:}$$

λ              wavelength,

z(x)          height coordinate of measured point located in x,

z(x+ λ)     height coordinate of measured point located in (x+ λ),

$\langle \cdot \rangle$          average value over the x range,

By representing the function $C_z(\lambda)$ in a logarithmic coordinate system, some characteristic parameters can be read (Figure 2). Correlation length $\xi_\perp$ is to characterize the profile in the vertical direction, while correlation length $\xi_{\parallel}$ is to characterize the profile in the longitudinal direction. Furthermore, D fractal dimension can also be read from the curve.



Figure 2
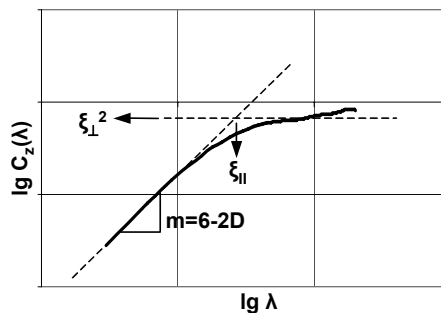Representation of function $C_z(\lambda)$

## 2.3    Problems of Evaluation Techniques

In the course of parameter-based evaluations, geometrical discrepancies arising from different sources are separated, roughness and waviness are examined separately. Separation is performed on the basis of the wavelength of the two profile components, which is characterized by λc (Cut-Off). λc is identical with

the basic measurement length lr, which can be used to interpret a considerable part of roughness parameters. The entire evaluation section is a multiple of the basic length (in general: $ln = 5 \cdot \lambda c$). The basic measurement and evaluation settings of the instruments are also specified by standards. The prescribed values of settings are selected depending on the fineness of the surface, according to the following table.

Table 1
Cut-Off and sampling distance in function of roughness parameters

| Cut-Off values according to ISO 4288-1996 | | | | |
|---|---|---|---|---|
| **Periodical profile** | **Non-periodical profile** | | **Cut-Off** | **Basic length of roughness / evaluation section** |
| Spacing Distance $S_m$ (mm) | $R_z$ (μm) | $R_a$ (μm) | $\lambda c$ (mm) | lr/ln (mm) |
| >0.013 to 0.04 | to 0.1 | to 0.02 | 0.08 | 0.08/0.4 |
| >0.04 to 0.13 | >0.1 to 0.5 | >0.02 to 0.1 | 0.25 | 0.25/1.25 |
| >0.13 to 0.4 | >0.5 to 10 | >0.1 to 2 | 0.8 | 0.8/4 |
| >0.4 to 1.3 | >10 to 50 | >2 to 10 | 2.5 | 2.5/12.5 |
| >1.3 to 4 | >50 | >10 | 8 | 8/40 |

When performing 2D measurements – as the characteristic parameters of the surface are not known – the required measurement length is specified on the basis of the expected surface roughness values of the surface examined. If the results yielded by the measurements are not within the intervals pertaining to standard measurement lengths, then the measurement must be repeated at another standard test length.

This article intends to present contradictory cases experienced in the course of surface roughness measurements where it was a problem to specify the parameters 'really' characterizing the surfaces. Controversial results were yielded by the examination, in the course of the measurements, of the operating surfaces – with the required surface roughness as designed – of turbo loading blades, carburettor nozzles and fuel feeders applied in diesel engines. We experienced that in case of different standard measurement lengths, different roughness values were yielded, all of which were within the value ranges prescribed by the standard for the respective measurement lengths. In the course of tests, the same part of surfaces was scanned by the stylus instrument after the adjustment of measurement lengths. The characteristic surface of each test piece was measured several times under identical measurement conditions; dominant measurement results are summarized in Table 2.

Table 2

Values set and measured in the course of tests

| Parameters | Turbo loading blade | | Carburettor nozzle | | Fuel feeder | |
|---|---|---|---|---|---|---|
| Cut-Off (mm) | 0.25 | 0.8 | 0.08 | 0.25 | 0.08 | 0.25 |
| Measurement length (mm) | 1.75 | 5.6 | 0.56 | 1.75 | 0.56 | 1.75 |
| Ra (µm) | 0.088 | 0.217 | 0.013 | 0.023 | 0.015 | 0.022 |
| Rz (µm) | 0.463 | 1.379 | 0.083 | 0.275 | 0.086 | 0.151 |

When examining measurement results, conspicuous discrepancies can be seen in the case of the most frequently used indices. Average and maximum surface roughness values sometimes present discrepancies of orders of magnitude at the same test piece, while each measurement can be deemed as a standard test.

# 3    Results

## 3.1    Parameter-based Characterization

Table 3 shows filtered roughness profile parameters while Table 4 shows unfiltered profile parameters for the three different components, where the signage is the following:

$$Turboloadingblade \quad \Rightarrow \quad a$$
$$Carburettornozzle \qquad \Rightarrow \quad b$$
$$Fuelfeeder \quad \Rightarrow \quad c$$

Table 3

Filtered profile parameters

| Compo nent sign | Measurem ent length (mm) | Measure ment length | Ra (µm) | Rz (µm) | RSk (-) | Rku (-) | RS (µm) |
|---|---|---|---|---|---|---|---|
| a | 1.75 | Short | 0.088 | 0.463 | 0.035 | 2,303 | 8,254 |
| | 5.6 | Long | 0.217 | 1.379 | -0.914 | 4,623 | 12,18 |
| b | 0.56 | Short | 0.013 | 0.083 | -0.858 | 4,1 | 5,841 |
| | 1.75 | Long | 0.023 | 0.275 | -3.663 | 24,43 | 6,958 |
| c | 0.56 | Short | 0.015 | 0.086 | -0.433 | 2,989 | 5,842 |
| | 1.75 | Long | 0.022 | 0.151 | -1.25 | 6,446 | 6,524 |
| Parameter average | | | **1.9** | **2.683** | **-6.319** | **3.374** | **1.261** |

Table 4
Unfiltered profile parameters

| Compon ent sign | Measurem ent length (mm) | Measure ment length | Pa (µm) | Pt (µm) | PSk (-) | Pku (-) | PS (µm) |
|---|---|---|---|---|---|---|---|
| a | 1.75 | rövid | 0.308 | 2.266 | -1.339 | 5,306 | 12,667 |
| | 5.6 | hosszú | 0.327 | 2.835 | -0.925 | 4,32 | 14,758 |
| b | 0.56 | rövid | 0.016 | 0.146 | -1.319 | 6,676 | 5,963 |
| | 1.75 | hosszú | 0.023 | 0.444 | -3.528 | 25,32 | 7,083 |
| c | 0.56 | rövid | 0.018 | 0.124 | -0.607 | 3,136 | 6,059 |
| | 1.75 | hosszú | 0.022 | 0.252 | -1.357 | 7,030 | 6,599 |
| Parameter average | | | 1.24 | 2.1 | -1.867 | 2.283 | 1.147 |

Average values serving as a basis for evaluation were provided by the following formula:

$$Parameter average = \frac{\frac{a_h}{a_r} + \frac{b_h}{b_r} + \frac{c_h}{c_r}}{3}$$

In the formula, the indices refer to measurement lengths. Out of the average values calculated, those parameters can be considered as appropriate surface characteristics where the values are in the range of 1 to 1.5. The reason for this is that such degrees of discrepancy may result from other measurement errors as well.

It can be established on the basis of the results that filtered parameters present significant differences in case of dissimilar measurement lengths, in spite of the fact that they refer to the same surface. Contradictory results can arise for two reasons: on the one hand, the difference in measurement lengths; and on the other hand, the different filters applied in the course of measurements.

As regards roughness indices, height-type characteristics presented (Ra, Rz) 2 to 2.5-fold differences. Even greater discrepancies arose in distortion characteristics (RSk, RKu). Measurement results fell in the acceptable range only for the width parameter RS.

On the whole, it can be stated that filtered parameters reflect surfaces measured in various ways to a lesser degree, therefore it is reasonable to apply unfiltered parameters.

As regards the indices in Table 4, it can be established that the value differences of measurements pertaining to identical surfaces are smaller than in the case of filtered profile parameters. These characteristics refer to identical test surfaces more specifically. Discrepancies can be explained by the various profiles detected, from which it directly follows that roughness indices react sensitively to the filters set.

In the case of the height parameters derived from unfiltered profile evaluation, discrepancies account for up to 75% of the differences experienced with filtered profiles. The smallest difference can be detected in longitudinal characteristics (RS, PS), which could also result from a measurement errors.

## 3.2 Motif Characterization

Table 5

Motif evaluation parameters

| Component sign | Measurement length (mm) | Measurement length | R (µm) | Rx (µm) | AR |
|---|---|---|---|---|---|
| a | 1.75 | Short | 0.6 | 2,26 | 141,9 |
| | 5.6 | Long | 0.898 | 2,348 | 205,6 |
| b | 0.56 | Short | 0.068 | 0,136 | 29,2 |
| | 1.75 | Long | 0.143 | 0,443 | 143,2 |
| c | 0.56 | Short | 0.069 | 0,102 | 26,0 |
| | 1.75 | Long | 0.108 | 0,252 | 53,8 |
| **Parameter average** | | | **1.72** | **2.26** | **2.81** |

It is interesting to observe that in the Motif evaluation, maximum motif roughness values (Rx) nearly entirely correspond to the maximum roughness indices of unfiltered profile parameters (Pt). As regards the other indices, it can be observed that there are 1.5 to 3-fold value differences between the figures of the two measurement lengths, which are lower than the discrepancies in filtered profile parameters, but graeter than the differences in unfiltered parameters P. Therefore, motif-based 'filtering' affects results to a greater extent than unfiltered evaluation does; moreover, this procedure is also sensitive to measurement length. It should be noted that parameters A and B to define motifs were identical in each measurement.

## 3.3 Height-difference Correlation Function

Height-difference correlation functions were also studied in case of the three different test pieces already presented in the previous chapters. Due to the extra-fine surfaces of the carburettor nozzle and the fuel feeder test pieces, the breakpoint of the curve only came about after 3 to 4 calculated points, which made it difficult to calculate the matching line. This made it particularly uncertain to specify the fractal dimension (D), therefore these results are not published.
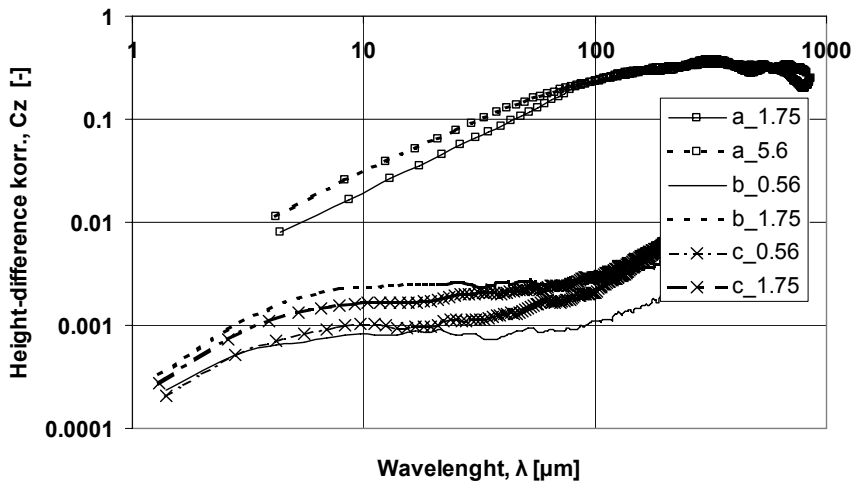
Figure 3
Height-difference correlation functions of 2D measured profiles

Figure 3 shows the HDCF curves, while Table 6 shows the correlation lengths to be read from each curve. It can be observed that the run-offs of the curve pairs pertaining to the two evaluation lengths of each surface are similar. In case of the turbo loading blade (a) correlation length also show good agreement, while in case of the other two components the difference is more significant, particularly as regards factor $\xi_\perp$.

Table 6
Parameters to be read from Figure 3

| workpiece | Measurement length [mm] | $\xi_\perp$ [μm] | $\xi_{\parallel}$ [μm] |
|---|---|---|---|
| a) | 1.75 | 0.565 | 117 |
| | 5.6 | 0.574 | 108 |
| b) | 0.56 | 0.028 | 9,7 |
| | 1.75 | 0.049 | 7,8 |
| c) | 0.56 | 0.031 | 8,33 |
| | 1.75 | 0.041 | 5,21 |
| **Parameter average [-]** | | **1.356** | **0.784** |

The lowermost line in Table 6 shows the average of the ratio of corresponding value pairs. So, it can be generally stated that the dissimilar measurement lengths produced from identical surfaces present 30 to 40% differences in case of height-difference corrrelation, a technique not applying filtering, which discrepancy can be explained by surface inhomogeneity, on the one hand, and different measurement lengths, on the other hand. Assuming that the surface is

homogeneous, the difference arising from dissimilar measurement lengths can result in up to 50% discrepancy on the basis of the tests above.

## Conclusions

In summary, it can be stated that the parameter-based evaluation system is very sensitive to the measurement length and the filter set. Differences between standard measurements performed at various measurement lengths can reach or even exceed 200 to 300% in case of some parameters.

In case of the parameters of unfiltered profiles, motif-based evaluation, and height-difference correlation, lesser differences of 50 to 100% arise for the different measurement lengths.

Further tests are required to find out whether the difference to be observed is characteristic of the surface in case of the same component or a consequence of the measurement length difference.

Contradictory parameter values were detected mostly at nominal roughness values near or at the borderline of a range of the filtering technology table. So it would be expedient to formulate extended rules for such limit values. These conditions can be based on unfiltered profile parameters or even on the characteristic parameters of the height-difference correlation function.

It can also be established that the differences detected are the smallest in the case of profile direction parameters.

## Acknowledgements

## References

[1]     Tom R. Thomas: Rough Surface, Imperial Collage Press, London (1998)

[2]     ISO 4288-1996

[3]     Horváth, S., Czifra, Á.: The Importance of Waviness in Study of Microtopography of Cutting Surface, DMC 2005, The 5[th] International Scientific Conference, Development of Metal Cutting, Kosice, Slovakia, September 12-13, 2005, pp. H 1-4

[4]     Thomas, T. R, Rosén, B. G.: Determination of the Sampling Interval for Rough Contact Mechanics, Tribology International 33, pp. 601-610 (2000)

[5]     Farkas, G., Czifra, Á.; Palásti K., B. Horváth, S.: Műszaki felületek mikrogeometriai vizsgálatában alkalmazott 2D-s és 3D-s paraméterek összevetése, információtartalmuk elemzése, Gép, 2005/2-3, pp. 51-59

[6]     Kovács K., Palásti Kovács B.: Műszaki felületek mikro-topográfiájának jellemzése háromdimenziós paraméterekkel. I. A háromdimenziós

topográfiai paraméterek áttekintése. Gépgyártástechnológia, 1999/8, pp. 19-24

[7]     Horváth, S.: Forgácsolt felületek hullámosságának vizsgálata, Doktori értekezés, Budapest, 1990

[8]     Klüppel, M., Müller, A., Le Gal, A., Heinrich, G.: Dynamic Contac Of Tires with Road Tracks, Meeting of the Rubber Division, American Chemical Society, San Francisco, April 28-30, 2003

# The Analysis of the Principal Eigenvector of Pairwise Comparison Matrices

**András Farkas**

Budapest Tech, Faculty of Economics
1084 Budapest, Tavaszmező út 17, Hungary
e-mail: farkas.andras@kgk.bmf.hu

*Abstract*. This paper develops the spectral properties of pairwise comparison matrices (PCM) used in the multicriteria decision making method called analytic hierarchy process (AHP). Perturbed PCMs are introduced which may result in a reversal of the rank order of the decision alternatives. The analysis utilizes matrix theory to derive the principal eigenvector components of perturbed PCMs in explicit form. Proofs are presented for the existence of rank reversals. Intervals over which such rank reversals occur are also established as function of a continuous perturbation parameter. It is proven that this phenomenon is inherent in AHP even in the case of the slightest departure from consistency. The results are demonstrated through a sample illustration.

*Keywords*: multiple criteria decision making, algebraic eigenvalue-eigenvector problem, rank reversal issue

## 1 Introduction

The analytic hierarchy process (AHP) is a multicriteria decision making method that employs a procedure of multiple comparisons to rank order alternative solutions to a multiobjective decision problem. Ever since the development of the AHP in the late 1970's by Saaty [14], a great number of criticisms of this approach have appeared in the literature. One of its most controversial aspects is the phenomenon of rank reversal of the decision alternatives. Both proponents and opponents of the AHP agree that rank reversal may occur, but disagree on its legitimacy. The problem has been considered by many authors and a persistent debate has followed; see Watson and Freeling [22], Saaty and Vargas [18], Belton and Gear [3], Vargas [21], Harker and Vargas [10], Dyer [5], Saaty [17], Harker and Vargas [11], Salo and Hämäläinen [19] and Pérez [13].

Despite the amount of work done on the subject, there are virtually no papers presenting a formal study of the algebraic eigenvalue-eigenvector problem of AHP's pairwise comparison matrix (PCM). This paper provides a rigorous mathematical presentation of this problem and gives proofs for the existence of rank reversal for a certain case. The foregoing research has been shown that a rank reversal may occur in AHP, (*i*) by introducing continuous perturbation(s) at one or more pairs of elements of a consistent PCM (see e.g., Watson and Freeling [22], Dyer and Wendell

[6]), or, (*ii*) by adding a new alternative to a perturbed PCM that is a replica (copy) of any of the old alternatives (see e.g., Belton and Gear [2], Dyer and Wendell [6]) and (*iii*) due to the normalization when aggregating the weights of the alternatives from the data even if the PCMs are each consistent to determine the overall priorities of the alternatives (see e.g., Barzilai and Golany [1]). In this paper, intervals are also established for the case (*i*) over which such rank reversals occur for situations when a PCM departs from perfect consistency even in only an arbitrarily small degree. The paper considers PCM's with a single criterion only.

**Definition 1** A square matrix $\mathbf{A}$ of order $n$ is called a *symmetrically reciprocal* (SR) matrix if its elements $a_{ij}$ are nonzero complex numbers and

$$a_{ij}a_{ji}=1, \quad \text{for } i\neq j; \quad i,j=1,2,...,n,$$
$$a_{ii}=1, \quad \text{for } i=1,2,...,n. \tag{1}$$

**Definition 2** A square matrix $\mathbf{A}$ of order $n$ is called a *transitive* matrix if its elements $a_{ij}$ are nonzero complex numbers and

$$a_{ij}a_{jk}=a_{ik}, \quad \text{for all } i,j,k. \tag{2}$$

**Definition 3** A square matrix $\mathbf{A}$ of order $n$ is a *one-rank* matrix if its elements $a_{ij}$ can be expressed as

$$a_{ij}=p_iq_j, \quad \text{for all } i,j. \tag{3}$$

**Theorem 1** *Let* $\mathbf{A}=[a_{ij}]$ *be a square matrix of order n, $n\geq3$. (i) If $\mathbf{A}$ is transitive, then $\mathbf{A}$ is a one-rank* SR, *as well. (ii) If $\mathbf{A}$ is a* SR *matrix, then $\mathbf{A}$ is transitive if and only if it is a one-rank matrix.*
(The proof of Theorem 1 is given in Farkas, Rózsa and Stubnya [8].)

The concept of a SR matrix defined by relation (1) was introduced by Saaty [14], who used the term reciprocal matrix. We prefer to designate this property according to Definition 1 since reciprocal matrices are the equivalent terms for the inverse matrices. In the framework of AHP, Saaty [14] developed such a SR matrix, $\mathbf{A}=[a_{ij}]$, called a *pairwise (paired) comparison* matrix, entries of which represent the relative importance ratios of the alternative $A_i$ over the alternative $A_j$, $i,j=1,2,...,n$, with respect to a common criterion. Elements of $\mathbf{A}$ are *positive, real* numbers. Saaty [14] called $\mathbf{A}$ a *consistent* matrix if the transitivity property (2) holds for $\mathbf{A}$ as well (cardinal consistency). In the AHP, every decision maker should provide ratio estimates for each possible pair of the alternatives [$n(n-1)/2$].

Using an eigenvalue-eigenvector approach, for a finite set of alternatives the AHP develops weights (and thus the priority ranking) of the alternatives on a *ratio scale*. Due to the properties of most of the decision problems occurring in practice the *rank order* of the alternatives, however, is usually generated on an *ordinal scale*. As it is well known, an ordinal ranking is said to be complete (it contains no ties) if the ordinal transitivity condition (ordinal consistency) holds, i.e.,

$$A_i \rightarrow A_j \text{ and } A_j \rightarrow A_k \text{ imply } A_i \rightarrow A_k \quad \text{for all } i,j,k, \tag{4}$$

where, the symbol $A_i \rightarrow A_j$ is interpreted as $A_i$ is preferred to $A_j$.

Saaty [15, p.848] proved that the weight (priority score) of an alternative, what he called the *relative dominance* of the $i$th alternative $A_i$, is the $i$th component of the principal right eigenvector of $\mathbf{A}$, $u_i$, provided that $\mathbf{A}$ is consistent, i.e., $\mathbf{A}$ is a transitive PCM. The *principal* right eigenvector belongs to the eigenvalue of largest modulus. The eigenvalue of largest modulus will be called *maximal* eigenvalue. By Perron's theorem, for matrices with positive elements, the maximal eigenvalue is always positive, simple and the components of its associated eigenvector are positive (see e.g., in Horn and Johnson [12]). Saaty [15, p.853] claimed to prove that this result also holds for a SR matrix that is *not* necessarily consistent, i.e., if it is *not* transitive. At this point the question can be raised, whether or not the components of the principal right eigenvector produce the true relative dominance of the alternatives, if the PCM is perturbed. Therefore, in this paper we shall study the behavior of the components of the principal eigenvectors of perturbed PCMs.

In Section 2, PCMs of specific form are defined and their spectral properties are described. In Section 3, PCMs of general form are introduced and their spectral properties are developed. In Section 4, the issue of rank reversal is examined for the specific versus the simple perturbed case. Proofs are given for the intervals of rank reversals. A sample illustration is provided in Section 5. The characteristic polynomials of PCMs of general form and the development of their principal eigenvector components are presented in Appendices A and B.

## 2 Pairwise Comparison Matrices of Specific Form

**Definition 4** A square matrix with positive entries is called a *specific* PCM denoted by $\mathbf{A}$, if it is transitive.

According to Theorem 1, any transitive matrix can be expressed as the product of a column vector $\mathbf{u}$ and a row vector $\mathbf{v}^{\mathrm{T}}$ as :

$$\mathbf{A} = \mathbf{u}\mathbf{v}^{\mathrm{T}}, \tag{5}$$

where

$$\mathbf{v}^{\mathrm{T}} = \left[1, x_1, x_2, \ldots, x_{n-1}\right], \text{ and } \mathbf{u}^{\mathrm{T}} = \left[1, \frac{1}{x_1}, \frac{1}{x_2}, \ldots, \frac{1}{x_{n-1}}\right].$$

Introducing the diagonal matrix $\mathbf{D} = \mathrm{diag}\langle 1, 1/x_1, 1/x_2, \ldots, 1/x_{n-1}\rangle$ and the vector $\mathbf{e}^{\mathrm{T}} = [1,1,\ldots,1]$, obviously $\mathbf{D}^{-1}\mathbf{A}\mathbf{D} = \mathbf{e}\mathbf{e}^{\mathrm{T}}$. It is easy to see that the characteristic polynomial of $\mathbf{A}$, $p_n(\lambda)$, can be obtained in the following from

$$p_n(\lambda) \equiv \det[\lambda \mathbf{I}_n - \mathbf{A}] = \det[\lambda \mathbf{I}_n - \mathbf{e}\mathbf{e}^{\mathrm{T}}] = \lambda^{n-1}(\lambda - n), \tag{6}$$

where $\mathbf{I}_n$ is the identity matrix of order $n$. From (6), it is apparent that $\mathbf{A}$ has a zero eigenvalue with multiplicity $n-1$ and one simple positive eigenvalue, $\lambda=n$, with the corresponding right and left eigenvectors, $\mathbf{u}$ and $\mathbf{v}^{\mathrm{T}}$, respectively. The relative dominance of the alternatives are given by the components of $\mathbf{u}$. Conventionally, this solution is normalized so that its components sum to unity.

## 3  Pairwise Comparison Matrices of General Form

In applied problems, decision makers give subjective judgements on the relative importance ratios. As a common consequence, usually, a failure of the relation (2) to hold is manifested in their PCMs. Hence, it seems to be apparent to explore how the maximal eigenvalue and its associated eigenvector vary when matrix $\mathbf{A}$ is perturbed such that it remains in SR, however, its transitivity is lost.

**Definition 5** A square matrix with positive entries is called a *perturbed* PCM and denoted by $\mathbf{A}_{\mathrm{p}}$, if the matrix is symmetrically reciprocal and not transitive.

Consider now the transitive matrix $\mathbf{A}=\mathbf{D}\mathbf{e}\mathbf{e}^{\mathrm{T}}\,\mathbf{D}^{-1}$ with the elements $a_{ij}=1/a_{ji}=x_j/x_i$, $i,j=0,1,2,...,n-1$. Let the elements of matrix $\mathbf{A}$ be perturbed in its first row and in its first column in a multiplicative way. This perturbed SR matrix $\mathbf{A}_{\mathrm{P}}$ can be written as

$$
\mathbf{A}_{\mathrm{P}} =
\begin{bmatrix}
1 & x_1\delta_1 & x_2\delta_2 & \cdot & \cdot & \cdot & x_{n-1}\delta_{n-1} \\[2mm]
\dfrac{1}{x_1\delta_1} & 1 & \dfrac{x_2}{x_1} & \cdot & \cdot & \cdot & \dfrac{x_{n-1}}{x_1} \\[2mm]
\dfrac{1}{x_2\delta_2} & \dfrac{x_1}{x_2} & 1 & \cdot & \cdot & \cdot & \dfrac{x_{n-1}}{x_2} \\[2mm]
\cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\
\cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\
\cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\[2mm]
\dfrac{1}{x_{n-1}\delta_{n-1}} & \dfrac{x_1}{x_{n-1}} & \dfrac{x_2}{x_{n-1}} & \cdot & \cdot & \cdot & 1
\end{bmatrix},
\tag{7}
$$

where $\delta_1,\delta_2,...,\delta_{n-1}$ are arbitrary positive numbers with $\delta_i \neq 1$, $i=1,2,...,n-1$. Performing a similarity transformation, the characteristic polynomial of $\mathbf{A}_{\mathrm{P}}$, $p_{\mathrm{n}}^{\mathrm{P}}(\lambda)$, is obtained as

$$
p_n^{\mathrm{P}}(\lambda) \equiv \det[\lambda\mathbf{I}_n - \mathbf{A}_{\mathrm{P}}] = \det[\lambda\mathbf{I}_n - \mathbf{D}^{-1}\mathbf{A}_{\mathrm{P}}\mathbf{D}] = \det\mathbf{K}_{\mathrm{P}}(\lambda),
\tag{8}
$$

where

$$\det \mathbf{K}_P(\lambda) = \begin{vmatrix} \lambda-1 & -\delta_1 & -\delta_2 & . & . & . & -\delta_{n-1} \\ -\dfrac{1}{\delta_1} & \lambda-1 & -1 & . & . & . & -1 \\ -\dfrac{1}{\delta_2} & -1 & \lambda-1 & . & . & . & -1 \\ . & . & . & . & . & . & . \\ . & . & . & . & . & . & . \\ . & . & . & . & . & . & . \\ -\dfrac{1}{\delta_{n-1}} & -1 & -1 & . & . & . & \lambda-1 \end{vmatrix}.$$

The matrix $\mathbf{K}_P(\lambda) = \lambda \mathbf{I}_n - \mathbf{D}^{-1}\mathbf{A}_P\mathbf{D}$ may be interpreted in the form of a modified matrix: with the notation

$$\mathbf{U}_P = \begin{bmatrix} 0 & 1 \\ 1-\dfrac{1}{\delta_1} & 0 \\ . & . \\ . & . \\ . & . \\ 1-\dfrac{1}{\delta_{n-1}} & 0 \end{bmatrix} \text{ and } \mathbf{V}_P^T = \begin{bmatrix} 1 & 0 & . & . & . & 0 \\ 0 & 1-\delta_1 & . & . & . & 1-\delta_{n-1} \end{bmatrix}.$$

To find the inverse of a matrix that is modified by a one-rank matrix [see the determinant (A1) in Appendix A] through applying the Sherman-Morrison formula [20, p.126] let us introduce the matrix $\mathbf{T}_P(\lambda)$ as

$$\mathbf{T}_P(\lambda) = \lambda \mathbf{I}_n + \mathbf{U}_P\mathbf{V}_P^T. \tag{10}$$

Thus, the modified matrix $\mathbf{K}_P(\lambda)$ can now be described as

$$\mathbf{K}_P(\lambda) = \mathbf{T}_P(\lambda) - \mathbf{e}\mathbf{e}^T. \tag{11}$$

It is shown in Appendix A that the determinant of matrix $\mathbf{K}_P(\lambda)$, i.e. the characteristic polynomial $p_n^P(\lambda)$, yields

$$p_n^{\mathrm{P}}(\lambda) \equiv \det \mathbf{K}_{\mathrm{P}}(\lambda) = \lambda^{n-3} \left\{ \lambda^3 - n\lambda^2 + (n-1)\sum_{i=1}^{n-1}(1-\delta_i)(1-\frac{1}{\delta_i}) - \sum_{i=1}^{n-1}(1-\frac{1}{\delta_i})\sum_{i=1}^{n-1}(1-\delta_i) \right\}. \quad (12)$$

Expression (12) shows clearly, that even if all elements are perturbed in one (arbitrary) row and in its corresponding column of matrix $\mathbf{A}$, then, $\mathbf{A}_{\mathrm{P}}$ has a zero eigenvalue with multiplicity$\geq n-3$, if $n>2$ and a trinomial equation is obtained for the nonzero eigenvalues. From (12), the characteristic polynomial, $p_n^{\mathrm{P}}(\lambda)$, can be rewritten in the simplified form

$$p_n^{\mathrm{P}}(\lambda) \equiv \lambda^{n-3}\left\{\lambda^3 - n\lambda^2 - C\right\}, \quad (13)$$

where the constant term $C$ contains the perturbation factors $\delta_i \neq 1$, $i=1,2,...,n-1$, as

$$C = -(n-1)\sum_{i=1}^{n-1}(1-\delta_i)(1-\frac{1}{\delta_i}) + \sum_{i=1}^{n-1}(1-\frac{1}{\delta_i})\sum_{i=1}^{n-1}(1-\delta_i).$$

From now on, we restrict our investigations to PCMs with *one* SR perturbed pair of elements only, say $\delta_1 = \delta \neq 1$, while $\delta_i = 1$ for $i \neq 1$.

**Definition 6** If one pair of elements, $a_{12}$ and $a_{21}$ of a specific PCM has the form $a_{12} = x_1\delta$, $a_{21} = 1/x_1\delta$, and $\delta > 0$, then it is called a *simple perturbed* PCM, denoted by $\mathbf{A}_{\mathrm{S}}$.

In this special case, we have the simple perturbed matrix $\mathbf{A}_{\mathrm{S}}$ as

$$\mathbf{A}_{\mathrm{S}} = \begin{bmatrix} 1 & x_1\delta & x_2 & . & . & . & x_{n-1} \\ \dfrac{1}{x_1\delta} & 1 & \dfrac{x_2}{x_1} & . & . & . & \dfrac{x_{n-1}}{x_1} \\ \dfrac{1}{x_2} & \dfrac{x_1}{x_2} & 1 & . & . & . & \dfrac{x_{n-1}}{x_2} \\ . & . & . & . & . & . & . \\ . & . & . & . & . & . & . \\ . & . & . & . & . & . & . \\ \dfrac{1}{x_{n-1}} & \dfrac{x_1}{x_{n-1}} & \dfrac{x_2}{x_{n-1}} & . & . & . & 1 \end{bmatrix}. \quad (14)$$

Performing a similarity transformation [see (6) and (8)], the characteristic polynomial of $\mathbf{A}_S$, $p_n^P(\lambda)$, can be written as

$$p_n^S(\lambda) \equiv \det[\lambda \mathbf{I}_n - \mathbf{A}_S] = \det[\lambda \mathbf{I}_n - \mathbf{D}^{-1}\mathbf{A}_S\mathbf{D}] = \det\mathbf{K}_S(\lambda), \tag{15}$$

where

$$\det\mathbf{K}_S(\lambda) = \begin{vmatrix} \lambda-1 & 1-\delta & -1 & . & . & . & -1 \\ 1-\dfrac{1}{\delta} & \lambda-1 & -1 & . & . & . & -1 \\ -1 & -1 & \lambda-1 & . & . & . & -1 \\ . & . & . & . & . & . & . \\ . & . & . & . & . & . & . \\ . & . & . & . & . & . & . \\ -1 & -1 & -1 & . & . & . & \lambda-1 \end{vmatrix}.$$

Similarly to (9), the matrix $\mathbf{K}_S(\lambda)$ in (15)

$$\mathbf{K}_S(\lambda) = \lambda \mathbf{I}_n - \mathbf{D}^{-1}\mathbf{A}_S\mathbf{D}, \tag{16}$$

may be interpreted as a modified matrix

$$\mathbf{K}_S(\lambda) = \lambda \mathbf{I}_n + \mathbf{U}_S\mathbf{V}_S^T - \mathbf{e}\mathbf{e}^T, \tag{17}$$

with the notation

$$\mathbf{U}_S = \begin{bmatrix} 0 & 1 \\ 1-\dfrac{1}{\delta} & 0 \\ . & . \\ . & . \\ . & . \\ 0 & 0 \end{bmatrix} \text{ and } \mathbf{V}_S^T = \begin{bmatrix} 1 & 0 & . & . & . & 0 \\ 0 & 1-\delta & . & . & . & 0 \end{bmatrix}.$$

Introducing

$$\mathbf{T}_{S}(\lambda) = \lambda \mathbf{I}_{n} + \mathbf{U}_{S} \mathbf{V}_{S}^{\mathrm{T}} = \begin{bmatrix} \lambda & 1-\delta & 0 & . & 0 \\ 1-\dfrac{1}{\delta} & \lambda & 0 & . & 0 \\ 0 & 0 & \lambda & . & 0 \\ . & . & . & . & . \\ 0 & 0 & 0 & . & \lambda \end{bmatrix}, \tag{18}$$

the modified matrix $\mathbf{K}_{S}(\lambda)$ can be written as

$$\mathbf{K}_{S}(\lambda) = \mathbf{T}_{S}(\lambda) - \mathbf{e}\mathbf{e}^{\mathrm{T}}. \tag{19}$$

In this special case, the characteristic polynomial (12) has the form

$$p_{n}^{S}(\lambda) \equiv \lambda^{n-3}[\lambda^{3} - n\lambda^{2} - C_{S}], \tag{20}$$

where, the constant term, $C_{S}$, now becomes

$$C_{S} = -(n-2)(1-\delta)(1-\frac{1}{\delta}) = (n-2)Q,$$

and $Q$ is expressed as a function of the perturbation factor $\delta$ as

$$Q = \delta + \frac{1}{\delta} - 2, \quad \delta > 0 \ (\delta \neq 1). \tag{21}$$

Let $r$ denote the maximal eigenvalue of a simple perturbed PCM, $\mathbf{A}_{S}$. Then, $r$ can be obtained from the equation [cf.(20)]:

$$r^{3} - nr^{2} - (n-2)Q = 0, \tag{22}$$

where $Q$ is given by (21). Since $Q>0$, from (22) it is easy to see that $r>n$. The proof can be found in Farkas, Rózsa and Stubnya [8]. The components of the principal eigenvector can be obtained from the one-rank matrix

$$\mathrm{adj}(r\mathbf{I}_{n} - \mathbf{A}_{S}) = [u_{ij}^{S}(r)], \tag{23}$$

since any column of the adjoint gives the elements of the principal eigenvector. In Appendix B, we show that the elements, $u^{S}_{ij}(r)$, of the principal eigenvector for the simple perturbed case are:

$$\begin{bmatrix} u_{11}^{S}(r) \\ u_{21}^{S}(r) \\ \dots \\ u_{i1}^{S}(r) \\ \dots \end{bmatrix} = \begin{bmatrix} r^{n-2}[r-(n-1)] \\ \dfrac{1}{x_1}r^{n-3}\left\{r-\left(1-\dfrac{1}{\delta}\right)[r-(n-2)]\right\} \\ \dots \\ \dfrac{1}{x_{i-1}}r^{n-3}\left\{r-\left(1-\dfrac{1}{\delta}\right)\right\} \\ \dots \end{bmatrix}; \quad i=3,4,\dots,n, \qquad (24)$$

$$\begin{bmatrix} u_{12}^{S}(r) \\ u_{22}^{S}(r) \\ \dots \\ u_{i2}^{S}(r) \\ \dots \end{bmatrix} = \begin{bmatrix} r^{n-3}\{r+(\delta-1)[r-(n-2)]\} \\ \dfrac{1}{x_1}r^{n-2}[r-(n-1)] \\ \dots \\ \dfrac{1}{x_{i-1}}r^{n-3}\{r+(\delta-1)\} \\ \dots \end{bmatrix} x_1; \quad i=3,4,\dots,n, \qquad (25)$$

and

$$\begin{bmatrix} u_{1j}^{S}(r) \\ u_{2j}^{S}(r) \\ \dots \\ u_{ij}^{S}(r) \\ \dots \end{bmatrix} = \begin{bmatrix} r^{n-3}[r+(\delta-1)] \\ \dfrac{1}{x_1}r^{n-3}\left\{r-\left(1-\dfrac{1}{\delta}\right)\right\} \\ \dots \\ \dfrac{1}{x_{i-1}}r^{n-2}\left\{\dfrac{r-2}{n-2}\right\} \\ \dots \end{bmatrix} x_{j-1}; \quad i,j=3,4,\dots,n. \qquad (26)$$

## 4 The Issue of Rank Reversal

The concept of rank reversal is now introduced. Consider the simple perturbed matrix $\mathbf{A}_S$ defined by (14). In the specific versus the simple perturbed case, the maximal eigenvalue $r$ of matrix $\mathbf{A}_S$ can be determined from (22), where $r>n$ ($n\geq3$) always holds [8]. The components of the principal eigenvector can be obtained from the one-rank matrix (23). Since any column of this matrix gives the elements of the

principal eigenvector, let us choose the $n$th column, hence, let $j=n$. Suppose that for two consecutive elements, $u_i$ and $u_{i+1}$ of the principal eigenvector of a *specific* PCM

$$u_i < u_{i+1} \tag{27}$$

holds. Furthermore, suppose that for the corresponding two elements, $u_{in}^{S}(r)$ and $u_{i+1,n}^{S}(r)$, of the adjoint matrix (23), i.e., for those of the principal eigenvector of a *simple perturbed* PCM

$$u_{in}^{S}(r) > u_{i+1,n}^{S}(r) \tag{28}$$

holds. If this case occurs, then, the rank order of the alternatives $A_i$ and $A_{i+1}$ has been reversed. This phenomenon is called the *rank reversal* of the alternatives which are in question.

It is well known in the cardinal theory of decision making that an opposite order of the corresponding components of the principal eigenvector *cannot* be yielded. In contrary to this, in the sequel, we give proofs for the occurrence of such rank reversals in the AHP between the alternatives $A_1$ and $A_2$. For this purpose, it will be sufficient to compare the order of the first two components of the principal eigenvectors.

For the specific case, the maximal eigenvalue of **A** equals $n$. The first two components of the principal eigenvector of **A** are as follows [cf. (5)]

$$1 \quad ; \quad \frac{1}{x_1} \tag{29}$$

i.e., the components of the principal eigenvector are monotonously increasing for $x_1 < 1$, whereas they are monotonously decreasing for $x_1 > 1$. In Theorem 2, necessary and sufficient condition is given for the occurrence of a rank reversal in the specific versus the simple perturbed case.

**Theorem 2** *Let* $\mathbf{A}=[a_{ij}]$ *be a transitive (consistent) pairwise comparison matrix of order n, n $\geq$3. Between the alternatives $A_1$ and $A_2$ when the elements $a_{12}$ and $a_{21}$ of are perturbed, a rank reversal occurs if and only if*

$$1 > x_1 > \frac{r-1+\frac{1}{\delta}}{r-1+\delta} = 1 - \frac{\delta-\frac{1}{\delta}}{\delta+(r-1)}, \quad for \ \delta > 1, \tag{30}$$

*or*

$$1 < x_1 < \frac{r-1+\frac{1}{\delta}}{r-1+\delta} = 1 + \frac{\frac{1}{\delta}-\delta}{\delta+(r-1)}, \quad for \ 0 < \delta < 1, \tag{31}$$

*Proof.* Using (B5) given in Appendix B, after performing the necessary algebraic manipulations the first two elements of the $n$th column of $\text{adj}(r\mathbf{I}_n - \mathbf{D}^{-1}\mathbf{A}_S\mathbf{D})$, i.e., the cofactors corresponding to the first two elements of the $n$th row $(r\mathbf{I}_n - \mathbf{D}^{-1}\mathbf{A}_S\mathbf{D})$ are obtained as [cf. (26)]

$$\left\{\text{adj}(r\mathbf{I}_n - \mathbf{D}^{-1}\mathbf{A}_S\mathbf{D})\right\}_{1n} = r^{n-3}[r-(1-\delta)] \; ; \; \left\{\text{adj}(r\mathbf{I}_n - \mathbf{D}^{-1}\mathbf{A}_S\mathbf{D})\right\}_{2n} = r^{n-3}[r-(1-\frac{1}{\delta})]. \quad (32)$$

Taking into account (B6) given in Appendix B, the first two components of the principal right eigenvector of the simple perturbed PCM, $\mathbf{A}_S$, are proportional to A rank reversal occurs if the elements in (33) are monotonously decreasing for $A_i$ and $A_{i+1} < 1$, or they are monotonously increasing for $x_1 > 1$ [cf. (29)].

$$r - 1 + \delta \; ; \; \frac{1}{x_1}(r - 1 + \frac{1}{\delta}). \quad (33)$$

Depending on than whether $\delta$ is greater unity, or $\delta$ is less than unity, two cases are distinguished:

(*i*) if $\delta > 1$ and $x_1 < 1$, then the elements in (33) are monotonously *decreasing* if $x_1$ resides in the interval given by (30), and

(*ii*) if $0 < \delta < 1$ and $x_1 > 1$, then the elements in (33) are monotonously *increasing* if $x_1$ resides in the interval given by (31).

This means that the condition is *necessary*. Furthermore, since all operations in the proof can be performed in the opposite direction, the condition is *sufficient* as well.

We note that according to (21) and (22), $r$ is dependent on the value of $\delta$. This fact, however, has no impact on the *existence* of the intervals (30) and (31), over which a rank reversal occurs. □

Concerning the other elements of the principal eigenvector, they can be obtained by making similar considerations. As a result, for these elements we have

$$u_{in}^P = \frac{1}{x_{i-1}} r \frac{r-2}{n-2}, \quad i = 3,4,...,n. \quad (34)$$

From (34), it is obvious that rank reversal cannot occur between any pair of the alternatives $A_3, A_4,...,A_n$. The occurrence of a rank reversal between alternatives $A_1$ and $A_i$, $i=3,4,...,n$, or between $A_2$ and $A_i$, $i=3,4,...,n$, could be analyzed in a similar way as was shown above. This investigation, however, is left to the reader.

# 5 A Sample Illustration

A widely used concept of measuring the degree of consistency of a perturbed PCM in the AHP framework is to calculate the *consistency index*, *CI*. Saaty (1977) introduced the following formula

$$CI = \frac{r-n}{n-1}. \tag{35}$$

The *consistency ratio*, *CR*, can be computed by comparing the *CI* with the corresponding random consistency index, *RI*, derived from a sample of 500, of randomly generated PCMs using the scale of [1/9,1/8,...,1,...,8,9] (see in Saaty [16]). He proposed that if this consistency ratio *CR=CI/RI* is less than or equal to 0.10, then the results be accepted. Otherwise, the problem should be studied again and its corresponding PCM revised. He also stated that such a small error does not affect the order of magnitude of the alternatives and hence, the relative dominance would be about the same.

Given *n*, and specifying a value for *CI*, from (35), the maximal eigenvalue *r*, of a simple perturbed PCM, $\mathbf{A}_S$, given by (14) can be obtained as

$$r = n + CI(n-1), \tag{36}$$

then, from (22), for the term *Q* we have

$$Q = \frac{n-1}{n-2} r^2 CI. \tag{37}$$

Next, using (21), the roots of the following equation can be calculated from

$$\delta^2 - (2+Q)\delta + 1 = 0. \tag{38}$$

Finally, using (30) and (31), the intervals for the values of $x_1$ over which a rank reversal occurs are

$$1 > x_1 > \frac{(n-1)(1+CI)+\dfrac{1}{\delta}}{(n-1)(1+CI)+\delta} = \frac{r-1+\dfrac{1}{\delta}}{r-1+\delta}, \quad for \ \delta > 1, \tag{39}$$

and

$$1 < x_1 < \frac{(n-1)(1+CI)+\dfrac{1}{\delta}}{(n-1)(1+CI)+\delta} = \frac{r-1+\dfrac{1}{\delta}}{r-1+\delta}, \quad for \ 0 < \delta < 1. \tag{40}$$

Consider a simple perturbed PCM of order *n*=3, that departs from consistency arbitrarily small. Let *CI*=0.01. Using the appropriate table in Saaty [16], the corresponding *RI*=0.58. Thus, *CR*=0.017. From (36), (37), and (38), the computed parameters are, *r*=3.02, *Q*=0.1824, $\delta$=1.5279, $1/\delta$=0.6545, respectively. Using (39) and (40), the values of $x_1$, with any of which a rank reversal occurs lie in the interval 0.7538 to 1.3266. This result demonstrates that due to the fact that $\mathbf{A}_S$ is an inconsistent PCM even in the slightest degree, yet there exists a relatively large

interval, over which rank reversal occurs between alternatives $A_1$ and $A_2$. That is, the fundamental ordinal transitivity relation given by (4) is being violated by this phenomenon. The occurrence of such a rank reversal might be serious in practice when an undesired alternative is chosen by the decision maker as the best.

At this point the question might be raised as to whether it would be meaningful to revise a given perturbed PCM and then, to make attempts to reduce its *CR* measure. It is remarkable, that meanwhile in the literature, several other more promising ways have been proposed for improving the measure of inconsistency of a general PCM (see Salo and Hämäläinen [19], Genest and Zhang [9] and Bozóki and Rapcsák [4]). The study of these approaches is, however, left to the reader.

**Conclusions**

This paper presented a matrix theory based analysis for the eigenvalue-eigenvector approach of the AHP. It was shown that this approach produces a perfect solution to the decision making problem if the PCM is consistent. However, the method cannot give the true ranking of the alternatives if the PCM is inconsistent, i.e., if it is *not* transitive. Therefore, if a PCM is inconsistent, even in the slightest degree, then the principal eigenvector components do not give the true relative dominance of the alternatives. Obviously, this result can be extended to PCMs with arbitrary number of perturbed pairs of elements, since, in the practical applications of the AHP, neither the cardinal consistency, nor the ordinal consistency of the expert's judgements can be ensured a'priori.

**Appendix A**

In order to obtain the characteristic polynomial, $p_n^P(\lambda)$, of the perturbed PCM, $\mathbf{A}_p$, [see (8)], let us write the determinant of the modified matrix $\mathbf{K}_p(\lambda)$ given by (9), in the following form

$$\det\mathbf{K}_P(\lambda) = \det\left[(\lambda\mathbf{I}_n + \mathbf{U}_P\mathbf{V}_P^T) - \mathbf{e}\mathbf{e}^T\right] = \det(\lambda\mathbf{I}_n + \mathbf{U}_P\mathbf{V}_P^T)\det\left[\mathbf{I}_n - (\lambda\mathbf{I}_n + \mathbf{U}_P\mathbf{V}_P^T)^{-1}\mathbf{e}\mathbf{e}^T\right]. \quad (A1)$$

It is easy to show that

$$\det[\mathbf{I}_n + \mathbf{W}\mathbf{Z}^T] = \det[\mathbf{I}_m + \mathbf{Z}^T\mathbf{W}],$$

where $\mathbf{W}$ is an ($n\times m$) matrix and $\mathbf{Z}^T$ is an ($m\times n$) matrix. Rewriting (A1), then using (10) and (11) we get

$$\det\mathbf{K}_P(\lambda) = \det(\lambda\mathbf{I}_n + \mathbf{U}_P\mathbf{V}_P^T)\left[1 - \mathbf{e}^T(\lambda\mathbf{I}_n + \mathbf{U}_P\mathbf{V}_P^T)^{-1}\mathbf{e}\right] = \det\mathbf{T}_P(\lambda)\left[1 - \mathbf{e}^T\mathbf{T}_P^{-1}(\lambda)\mathbf{e}\right]. \quad (A2)$$

The inverse of a matrix modified by a low-rank matrix may be written in the following form (see in Woodbury, [23])

$$\mathbf{T}_P^{-1}(\lambda)=(\lambda\mathbf{I}_n+\mathbf{U}_P\mathbf{V}_P^T)^{-1}=\frac{1}{\lambda}\mathbf{I}_n-\frac{1}{\lambda}\mathbf{U}_P(\lambda\mathbf{I}_2+\mathbf{V}_P^T\mathbf{U}_P)^{-1}\mathbf{V}_P^T. \tag{A3}$$

Using (A2) and by performing the necessary operations in (A3) (see in [8, p.426]), the characteristic polynomial of the perturbed PCM is obtained in the form

$$p_n^P(\lambda)\equiv\lambda^{n-3}\left\{\lambda^3-n\lambda^2+(n-1)\sum_{i=1}^{n-1}(1-\delta_i)(1-\frac{1}{\delta_i})-\sum_{i=1}^{n-1}(1-\frac{1}{\delta_i})\sum_{i=1}^{n-1}(1-\delta_i)\right\}. \tag{A4}$$

*Remark.* It is easy to show that, if the number of the rows (and their corresponding columns) which contain at least one perturbed pair of elements in the specific PCM, **A** [see matrix (7)], is $m\le(n-1)/2$, then, the rank of matrix **A** increases by $2m$, i.e., the multiplicity of the zero eigenvalues becomes $n-2m-1$, and we obtain an equation of degree $2m+1$ for the nonzero eigenvalues.


**Appendix B**

To develop the principal eigenvector of the simple perturbed PCM, $\mathbf{A}_S$, let us calculate the one-rank matrix

$$\text{adj}(r\mathbf{I}_n-\mathbf{A}_S)=\mathbf{ab}^T, \tag{B1}$$

any column of which produces the principal (right) eigenvector. First, the proof of the following lemma will be given that refers to the calculation of the adjoint of a modified matrix.

**Lemma** *If $\mathbf{T}_P$ is a nonsingular matrix of order n, furthermore, $\mathbf{a}$ and $\mathbf{b}$ are column vectors of order n, then the adjoint of the modified matrix $\mathbf{T}_P-\mathbf{ab}^T$ can be obtained in the form* (see in Elsner and Rózsa [7]):

$$\text{adj}[\mathbf{T}_P-\mathbf{ab}^T]=\text{adj}\mathbf{T}_P\left\{(1-\mathbf{b}^T\mathbf{T}_P^{-1}\mathbf{a})\mathbf{I}_n+\mathbf{ab}^T\mathbf{T}_P^{-1}\right\}. \tag{B2}$$

*Proof.* By the Sherman-Morrison formula [20, p.126)], the inverse of the modified nonsingular matrix $\mathbf{T}_P-\mathbf{ab}^T$ exists if

$$1-\mathbf{b}^T\mathbf{T}_P^{-1}\mathbf{a}\neq0,$$

and it can be written as

$$(\mathbf{T}_P-\mathbf{ab}^T)^{-1}=\mathbf{T}_P^{-1}+\frac{\mathbf{T}_P^{-1}\mathbf{ab}^T\mathbf{T}_P^{-1}}{1-\mathbf{b}^T\mathbf{T}_P^{-1}\mathbf{a}}. \tag{B3}$$

By (A2), the determinant of a nonsingular matrix $\mathbf{T}_P$ modified by a one-rank matrix $\mathbf{ab}^T$ is given as

$$\det[\mathbf{T}_P - \mathbf{ab}^T] = (1 - \mathbf{b}^T \mathbf{T}_P^{-1} \mathbf{a}) \det \mathbf{T}_P. \tag{B4}$$

Multiplying (B4) by the inverse (B3), the formula (B2) for the adjoint follows. &#9633;

**Corollary** *Since the determinant is a continuous function of its elements,* (B2) *is valid also in the case if* $1 - \mathbf{b}^T \mathbf{T}_P^{-1} \mathbf{a} = 0$, i.e.,

$$\text{adj}[\mathbf{T}_P - \mathbf{ab}^T] = (\text{adj} \mathbf{T}_P) \mathbf{ab}^T \mathbf{T}_P^{-1}, \quad \text{if} \quad 1 - \mathbf{b}^T \mathbf{T}_P^{-1} \mathbf{a} = 0. \tag{B5}$$

Let us now apply these results for the *simple perturbed* PCM, $\mathbf{A}_S$. Making use of (16), it is easy to show that

$$\text{adj}[\lambda \mathbf{I}_n - \mathbf{A}_S] = \mathbf{D}\{\text{adj}[\lambda \mathbf{I}_n - \mathbf{D}^{-1} \mathbf{A}_S \mathbf{D}]\} \mathbf{D}^{-1} = \mathbf{D}\{\text{adj}[\mathbf{K}_S(\lambda)]\} \mathbf{D}^{-1}. \tag{B6}$$

Let us introduce the notation $\mathbf{P}_S(\lambda) = \text{adj}[\mathbf{K}_S(\lambda)]$. According to (B2), by letting $\mathbf{a} = \mathbf{b} = \mathbf{e}$, and using (19) we can write that

$$\mathbf{P}_S(\lambda) \equiv \text{adj}[\mathbf{K}_S(\lambda)] = \text{adj}[\mathbf{T}_S(\lambda) - \mathbf{ee}^T]. \tag{B7}$$

Substituting $r$ for $\lambda$, by (22) and (21) it is obvious that $1 - \mathbf{e}^T \mathbf{T}_S^{-1}(r) \mathbf{e} = 0$. Thus, (B5) can be applied, and for the adjoint $\mathbf{P}_S(r)$ we have

$$\mathbf{P}_S(r) = [p_{ij}^S(r)] = \text{adj}[\mathbf{K}_S(r)] = \{\text{adj} \mathbf{T}_S(r)\} \mathbf{ee}^T \mathbf{T}_S^{-1}(r). \tag{B8}$$

Consequently, $\mathbf{P}_S(r)$ is a rank-one matrix, and therefore, any (column) vector of $\text{adj}[\mathbf{T}_S(r)]$ is the principal eigenvector corresponding to the maximal eigenvalue $r$ of the simple perturbed PCM $\mathbf{A}_S$. Hence, making use of (18), (B6), and (B8) the eigenvectors $u_{ij}^S(r)$, given by formulas (24), (25) and (26), can be obtained from

$$\text{adj}(r\mathbf{I}_n - \mathbf{A}_S) = \mathbf{D}\left\{ r^{n-3} \begin{bmatrix} r^2 + (\delta-1)r \\ r^2 - \left(1 - \dfrac{1}{\delta}\right)r \\ r^2 + Q \\ . \\ . \\ . \\ r^2 + Q \end{bmatrix} \left[ \dfrac{r - \left(1 - \dfrac{1}{\delta}\right)}{r^2 + Q}, \quad \dfrac{r + \delta - 1}{r^2 + Q}, \quad \dfrac{1}{r}, \quad . \quad . \quad ., \quad \dfrac{1}{r} \right] \right\} \mathbf{D}^{-1}, \tag{B9}$$

as the $k$th column of $\mathbf{P}_S(r)$ is premultiplied by $\mathbf{D}$ and is multiplied by $x_{k-1}$, $k=1,2,...,n$. In (B9), $Q$ is given by (21).

# References

[1]    Barzilai,J and B.Golany, "AHP Rank Reversal, Normalization and Aggregation Rules", *INFOR*, **32**, (1994), 57-64.

[2]    Belton,V. and T.Gear, "On a Short-coming on Saaty's Method of Analytic Hierarchies", *Omega*, **11**, (1983), 228-230.

[3]    Belton,V. and T.Gear, "The Legitimacy of Rank Reversal - A Comment", *Omega*, **13**, (1985), 143-144.

[4]    Bozóki,S. and T.Rapcsák, "On Saaty's and Koczkodaj's inconsistencies of pairwise comparison matrices," *Journal of Global Optimization*, (2007), [in press].

[5]    Dyer,J.S. "Remark on the Analytic Hierarchy Process", *Management Sci.*, **36**, (1990), 249-258.

[6]    Dyer,J.S., and R.E. Wendell, "A Critique of the Analytic Hierarchy Process", Working Paper. 84/85-424. Department of Management, The University of Texas at Austin, 1985.

[7]    Elsner,L. and P.Rózsa, "On Eigenvectors and Adjoints of Modified Matrices," *Linear and Multilinear Algebra*, **10**, (1981), 235-247.

[8]    Farkas,A., Rózsa,P. and E.Stubnya, "Transitive Matrices and Their Applications", *Linear Algebra and its Applications*. **302-303**, (1999), 423-433.

[9]    Genest,C. and S.S.Zhang, "A Graphical Analysis of Ratio-scaled Pairwise Comparison Data", *Management Sci.*, **42**, (1996), 335-349.

[10]   Harker,P.T. and L.G.Vargas, "The Theory of Ratio Scale Estimation: Saaty's Analytic Hierarchy Process", *Management Sci.*, **33**, (1987), 1383-1403.

[11]   Harker,P.T. and L.G.Vargas, "Reply to 'Remarks on the Analytic Hierarchy Process'", *Management Sci.*, **36**, (1990), 269-273.

[12]   Horn,R.A. and C.R.Johnson, Matrix Analysis. Cambridge University Press, Cambridge, 1985.

[13]   Pérez,J., "Some Comments on Saaty's AHP", *Management Sci.*, **41**, (1995), 1091-1095.

[14]   Saaty,T.L., "A Scaling Method for Priorities in Hierarchical Structures", *Journal of Math. Psychology*, **15**, (1977), 234-281.

[15]   Saaty,T.L., "Axiomatic Foundation of the Analytic Hierarchy Process", *Management Sci.*, **32**, (1986), 841-855.

[16]   Saaty,T.L.,"The Analytic Hierarchy Process-What It Is and How It Is Used.", *Math. Modelling*, **9**, (1987), 161-176.

[17] Saaty,T.L., "An Exposition of the AHP in Reply to the Paper 'Remarks on the Analytic Hierarchy Process'", *Management Sci.*, **36**, (1990), 259-268.

[18] Saaty,T.L. and L.G.Vargas, "The Legitimacy of Rank Reversal", *Omega*, **12**, (1984), 514-516.

[19] Salo,A.A. and R.P.Hämäläinen, "On the Measurement of Preferences in the Analytic Hierarchy Process", Research Report. A47, Helsinki University of Technology, Systems Analysis Laboratory, Espoo, Finland, 1992.

[20] Sherman,J. and W.J.Morrison,"Adjustment of an Inverse Matrix Corresponding to Changes in a given Column or a given Row of the Original Matrix", *Ann. Math. Stats.*, **21**, (1949), 124-127.

[21] Vargas,L.G., "A Rejoinder", *Omega*, **13**, (1985), p.249.

[22] Watson,S.R. and A.N.S.Freeling, "Comment on: Assessing Attribute Weights by Ratios", *Omega*, **11**, (1983), p.13.

[23] Woodbury,M.A.,"Inverting Modified Matrices", Memorandum Report. 42, Statistical Research Group. Institute for Advanced Study, Princeton, New Jersey, June 14, 1950 .

# Behaviour of RHEED Oscillation during LT-GaAs Growth

## Ákos Nemcsics

Institute for Microelectronics and Technology, Budapest Tech
Tavaszmező utca 15-17, H-1084 Budapest, Hungary
and
Research Institute for Technical Physics and Materials Science,
Hungarian Academy of Sciences, H-1525 Budapest, POB 49, Hungary
e-mail: nemcsics.akos@kvk.bmf.hu

*Abstract: The behaviour of decay constant of RHEED oscillation during MBE growth on GaAs (001) surface at low temperature growth conditions is studied in this work. The dependence of decay constant on As-to-Ga ratio, substrate temperature and the excess of As content in the layer are examined here.*

*Keywords: MBE, RHEED, LT-GaAs*

## 1    Introduction and Experimental Preliminaries

Nowadays, molecular-beam-epitaxial (MBE) growth of GaAs at low temperature (LT) – around 200 ºC – has become an expanding important method since it provides highly insulating films and contributes to the synthesis of magnetic semiconductors [1]. It was shown that LT growth leads to incorporation of excess As in the crystal up to 1.5% depending on the growth parameters [2, 3]. The high concentration of excess As in LT-GaAs results in many new properties. As-grown and annealed LT-GaAs layers exhibit extremely high electrical resistivity and very short lifetimes of photoexcited carriers [4]. Their electrical parameters can be interpreted using the combined band and hopping conduction model [5, 6]. The majority of excess As is in antisite position, while the remaining As excess originate from interstictial As or Ga vacancies [7, 8]. The uniqueness of LT-GaAs is its high density of midgap states resulting from excess As, while the structure of the matrix remains perfect [9].

The use of reflection high-energy electron diffraction (RHEED) to control the growth of LT-GaAs has been reported in [10-12]. It is not easy to observe RHEED oscillations at LT growth. The RHEED oscillations are very strongly influenced by the growth parameters, such as deposition temperature, As-to-Ga ratio, etc.

RHEED oscillations were observed in two regions of As-to-Ga ratio at LT growth. One of these regions is close and another is far from the unity of As-to-Ga ratio. The strongest oscillations were observed when this ratio was nearly unity [10, 11]. Oscillations were also found if the ratio was larger than hundred [12].

The RHEED and its intensity oscillations at LT-GaAs growth exhibit certain particular behaviours. The intensity, phase and decay constant of oscillations depend on the As-to-Ga ratio, excess As content and substrate temperature, too. We investigate here the decay constant of oscillation during the growth of LT-GaAs. The deposition temperature and the range of the As-to-Ga ratio are 200 ºC and 0.9-1.3, respectively. This investigation is based on the measurement and the observed intensity oscillations of RHEED which are described in the literature [9-11].

## 2    Results and Discussion

The temporal evaluations of RHEED specular intensity during the LT-GaAs growth – where the As-to-Ga ratio is close to unity – are shown in Figs. 1 and 2 of Refs. [10] and [11], respectively. It can be observed in these figures that when the As-to-Ga ratio moves off from unity, then the decay time of oscillations becomes small. If the ratio is 1.3 then the oscillation intensity becomes weak so its evaluation is difficult. The decay constant of the oscillations was determined as described in our previous work [13]. The amplitude decay of oscillations was investigated peak to peak. The peak to peak series are determined with the subtraction of the background. After the subtraction, an exponential function is fitted to determine the decay of intensity with the help of least-squares method. The exponential approximation of the time dependence of intensity is $I(t) = B_0 \exp(-t/\tau_d)$, where $\tau_d$ is the decay time constant and $B_0$ is a scaling factor. The extracted decay time constants vs. As-to-Ga ratio are shown in Fig. 1.
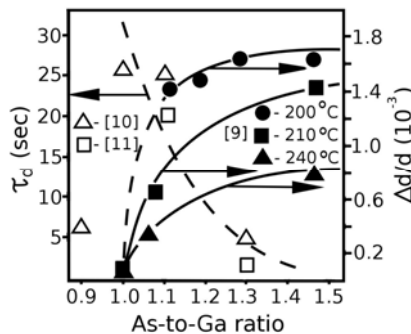


Figure 1

*left:* The decay constant vs As-to-Ga ratio at 200°C, *right:* The lattice spacing vs As-to-Ga ratio at 200, 210 and 240°C. The lines serves guide the eye only

It is known that the effect of the strain in the grown layer can be observed also with the help of RHEED oscillations. If the strain is large in the grown structure then the decay time constant is small. If the strain is small or absent in the structure then the decay time constant is large. This effect is demonstrated and described in the case of InGaAs/GaAs heterostructures – as a good model system – in our previous work [13, 14]. We can observe very strong change in the decay time constant depending on the As-to-Ga ratio at 200 °C (see Fig. 1). Depending on the growth parameters these LT-GaAs layers may contain large amount of excess As atoms. The majority of excess As is in the antisite position. The lattice spacing of LT-GaAs becomes greater than in the stoichiometric case. The relative increase of lattice spacing ($\Delta d/d$) of the non-stoichiometric LT-GaAs was determined in Ref [9]. The functions of $\Delta d/d$ vs As-to-Ga ratio are depicted also in Fig. 1.

We can observe (see Fig. 1) that the decay time constant of oscillations $\tau_d$ decreases rapidly during the LT-GaAs growth with increasing of As-to-Ga ratio, that is, also with increasing of the excess As content. The decay of oscillations can have several causes. The excess As gives rise to lattice mismatch, so also to strain, in the grown layer. This strain can influence the decay of intensity oscillations. At first, we investigate this effect because the strain effects are known and the influence of growth phenomena are unknown. From the given parameters (see Fig. 1) the mismatch dependence of the oscillation decay can be determined. The variation of decay time constant vs. $\Delta d/d$ is shown in the insert of Fig. 2. It is clear, that not only the mismatch is responsible for the decay but the growth conditions, too. Changes of the excess As modify not only the mismatch but the growth conditions, e.g. growth rate, too [15]. So, both the mismatch and the growth parameters influence the behaviour of the oscillation decay. We approximate this decay with an exponential function. Furthermore, we suppose that the both effects such as the mismatch and the growth influence can be separated from each other. In this way the decay phenomenon can be described by two time constants, as follows $I(t) = B_0 exp[(-t/\tau_G)+(-t/\tau_M)]= Bexp(-t/\tau_M)$, where $\tau_G$ and $\tau_M$ are the assumed time constans of the separated influences, such as growth and mismatch, respectively. $B$ and $B_0$ are the scaling factors which depend on the excess As and also on the $\Delta d/d$. The decay originating from the mismatch can be expessed as follows:

$$\frac{1}{\tau_M(\Delta d/d)} = \frac{1}{\tau_d(\Delta d/d)} - \ln\left(\frac{B_0}{B(\Delta d/d)}\right)\frac{1}{t} = \frac{1}{\tau_d(\Delta d/d)} - e(\Delta d/d)$$

where the factors are functions of $\Delta d/d$ and also of the As:Ga ratio and the term $e(\Delta d/d)$ serves as an operational aid only. In the case of stoichiometric LT-GaAs growth ($\Delta d/d$=0) there is no decay due to mismatch. This means, that for $\Delta d/d$=0 the reciprocal value of the decay time constant originates fully from the crystal growth phenomenon ($\tau_d(0)=\tau_G(0)=\tau_{Go}$), that is the value of $1/\tau_M$ (0) is zero. The

value of $\tau_{Go}$ is constant. The other component of $\tau_G$, $\tau_{G1}$ dependens on As-to-Ga ratio (or $\Delta d/d$), where the whole $\tau_G$ is $\tau_G = \tau_{G1}(As\text{-}to\text{-}Ga) + \tau_{Go}$. So the second term of the $1/\tau_M(\Delta d/d)$ expession $e(\Delta d/d)$ has also an independent and dependent part on As-to-Ga ratio (or $\Delta d/d$). The interchange between As-to-Ga ratio and $\Delta d/d$ may be only in the case of the narrow range of growth parameters where these ratios are proportional with each other.

We have separated the supposed strain effect from the effects due to growth which can be responsible for the decay of oscillation. In the case of InGaAs growth, we have supposed that the effect of growth remains constant in low In content region, because the one of the most important growth parameters, the deposition temperature, remained the same during the experiment. With this supposition we have obtained good agreement between the theoretical critical layer thickness and the threshold thickness, which is derived from the $\tau_M$ decay constant [14]. In the case of InGaAs In substitutes Ga in the lattice. Both elements estabilish similarly strong $sp^3$ type bonding in the lattice because the similar bonding structure. The situation in the case of LT-GaAs is quite different. The excess As which substitutes Ga in the lattice has different and weaker bonding than $sp^3$ hybrid. This fact modifies locally the probability of chemisorbtion of As atoms so also the probability of the incorporation of the further excess As atoms in the crystal [16]. The concentration of excess As can be determined from the chemisorbtion rate of As atoms. As atoms that are chemisorbed on the arsenic-terminated GaAs (001) surface serve as precursors of excess As, and hence, the concentration of excess As depends directly on the steady-state coverage of the chemisorbed As atoms [9]. The presented excess atoms As perturbs the bonding behaviour in the crystal, that is, the energy distribution along the surface. We use a simple description for the changing of the unperturbed surface layer by layer. At the first step, the unperturbed area $A^*$ can be written as follows: $A_1{}^* = Ab - Aa$, where $A$ is the whole area of the investigated surface. The factors $b$ and $a$, which are less than one, give the areas on the surface which can be covered by chemisorbed As and which can be incorporate excess As, respectively. The second step can be described as follows: $A_2{}^* = (Ab - Aa)b - Aa$. The $n^{th}$ layer we can get by follow-up the former given algorithm. The size of the perturbed area depends also on the number of the grown layers. This dependence can be neglible if the number of the layers is not large [16]. Among the surface reconstructions of the GaAs (001) surface, the c(4x4) reconstruction occurs at LTs under high As flux [17-20]. The value of $b$ can be estimated because the maximum coverage of chemisorbed As atoms is 0.75 monolayers like in the case of this reconstruction. The value of $a$ can be estimated by the maximum excess As content which is 0.015 [9]. It can be seen that the factor $b$ is larger than $a$, so we can get, after arrangement of the expression $A^*$ and neglectig small terms, the following simple power function for $n^{th}$ step: $A_n{}^* = Ab^n$. We suppose that the intensity of RHEED is proportional to the size of the smooth part of surface. A continuous description by replacing of $n$ by $rt$, yields

$I(t) = cA^*(t) = cb^{rt}A$, where $r$ is the growth rate, $t$ is the growth time and $c$ is a constant characterizing the diffraction power. This can be written in the following form: $I(t) = cA\exp(-t/\tau_{G1})$, where $\tau_{G1}$ is the decay time constant originated from growth phenomena, which depends on the As-to-Ga ratio, this is, $\tau_{G1} = -1/r\ln b$. The $\tau_{Go}$ and $\tau_{G1}$ dependence on $b$ are depicted in the insert of Fig. 2.
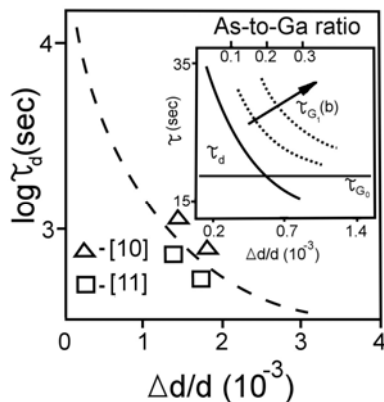


Figure 2

The function of $\tau_M$ vs $\Delta d/d$ which is derived from the high temperature InGaAs growth. The calculated data originated from the LT-GaAs growth. *insert:* The decay time constant vs lattice spacing derived from Fig. 1.

To justify our analysis we compare the values of $\tau_M$ extracted from the oscillation decay of LT-GaAs growth and the material independent decay constant, which is calculated from the effect of mismatch. The variation of $\tau_M(\Delta d/d)$ should be determined as follows: $1/\tau_M \propto 1/\tau_d - 1/\tau_{Go} - 1/\tau_{G1}$, similarly as described in Ref [14]. The strain dependent decay time constant vs composition in the case of InGaAs is given [14]. The composition independent variation of $\tau_M$ vs $\Delta d/d$ can be derived from the above mentioned dependence with the help of the modified Vegard`s law [21, 22]. The material independent variation is shown in Fig. 2. The calculated $\tau_M$ data from LT-GaAs are depicted in this figure where the fitting parameter of $b$ was 0.63. The value one of As-to-Ga ratio serves as a reference point for the calculation of $\tau_M$. In this calculation we have taken into consideration also the As-to-Ga ratio of 1.3. The $\tau_M$ determined from LT growth corresponds to the calculated dependence, but we have to note here that the ratio of 1.3 is very difficult to evaluate. We can estabilish that the separation of the growth and mismatch effect on the decay of RHEED oscillations can describe the LT growth only in a narrow range. The intensity oscillation at the As-to-Ga ratio of 1.3 is very uncertain to evaluate because the weak intensity. This drastical intensity damage can result not only from the mismatch joined with the reduction of unperturbed area but it can be also explained the change of the sticking coefficient of the deposited species.

**Conclusion**

The decay and absence of the RHEED intensity oscillations at LT-GaAs growth can origin from several effects e.g. change of sticking coefficients, change of the morphology of the grown surface and change of the mechanical strain in the layer. Here was found, that the separation of growth and strain influence on the RHEED oscillation decay in the case of LT-GaAs is possible in a narrow region of As-to-Ga ratio.

**Acknowledgements**

**References**

[1]    M. Missous: Fundamental Issues of Device-Relevant Low Temperature GaAs and Related Materials Properties; Material Science Engineering B 44, 304 (1997)

[2]    D. C. Look: Molecular Beam Epitaxial GaAs Grown at Low Temperatures; Thin Solid Films 231, 61 (1993)

[3]    G. J. Witt: LT MBE GaAs: Present Status and Perspectives; Material Science Engineering B 22, 9 (1993)

[4]    P. Kordos, A. Foster, J. Betko, M. Morvic, J. Novak: Semi-Insulating GaAs Layers Grown by Molecular-Beam Epitaxy; Applied Physics Letters 67, 983 (1995)

[5]    J. Novak, M. Kucera, M. Morvic, J. Betko, A. Foster, P. Kordos: Characterization of Low-Temperature GaAs by Galvanomagnetic and Photoluminescence Measurement; Material Science Engineering B 44, 341 (1997)

[6]    J. Betko, M. Morvic, J. Novak, A. Foster, P. Kordos: Magnetoresistance in Low-Temperature Grown Molecular-Beam Epitaxial GaAs; Journal Applied Physics 86, 6243 (1999)

[7]    X. Liu, A. Prasad, J. Nishino, E. R. Weber, Z. L. Weber, W. Walukiewicz: Native Point Defects in Low-Temperature-Grown GaAs; Applied Physics Letters 67, 279 (1995)

[8]    N. Otsuka, S. Fukushima, K. Mukai, F. Shimogishi, A. Shuda: Surface Atomic Process of Incorporation of Excess Arsenic in Low-Temperature-Grown GaAs; Proc. of 3[rd] Symp. on Non-Stoichiometric III-V Compounds (Eds. S. Malzer, T. Marek, P. Kiesel), Erlangen 127 (2001)

[9]    A. Shuda, N. Otsuka: Surface Atomic Process of Incorporation of Excess Arsenic in Molecular-Beam Epitaxy of GaAs; Surface Science 458, 162 (2000)

[10]   R. P. Mirin, J. P. Ibbetson, U. K. Mishra, A. Gossard: Low Temperature Limits to Molecular Beam Epitaxy of GaAs; Applied Physics Letters 65, 2335 (1994)

[11]   M. Missous: Stoichiometric Low-Temperature GaAs and AlGaAs: a Reflection High-Energy Electron-Diffraction Study; Journal Applied Physics 78, 4467 (1995)

[12]   A. Shen, H. Ohno, Y Horikoshi, S. P. Guo, Y. Ohno, F. Matsukura: Low-Temperature GaAs Grown by Molecular-Beam Epitaxy under High As Overpressure: a Reflection High-Energy Electron Diffraction Study; Applied Surface Science 130/132, 382 (1998)

[13]   Á. Nemcsics, J. Olde, M. Geyer, R. Schnurpfeil, R. Manzke, M. Skibowski: MBE Growth of Strained InxGa1-xAs on GaAs (001); Physica Status Solidi (a) 155, 427 (1996)

[14]   Á. Nemcsics: Correlation between the Critical Layer Thickness and the Decay Time Constant of RHEED Oscillations in Strained InxGa1-xAs/GaAs Structures; Thin Solid Films 367, 302 (2000)

[15]   A. Nagashima, M. Tazima, A. Nishimura, Y. Takagi, J. Yoshino: STM and RHEED Studies on Low-Temperature Growth of GaAs (001); Surface Science 514, 350 (2002)

[16]   Á. Nemcsics, to be published

[17]   P. K. Larsen, J. H. Neave, J. F. van der Veen, P. J. Dobson, B. A. Joyce: GaAs (001)-c(4x4): A Chemisorbed Structure; Physical Review B 27, 4966 (1983)

[18]   V. V. Preobrazhenskii, M. A. Putyato, O. P. Pchelyakov, B. R. Semyagin: Surface Structure Transition on (001) GaAs during MBE; Journal Crystal Growth 201/202, 166 (1999)

[19]   V. V. Preobrazhenskii, M. A. Putyato, B. R. Semyagin: Measurements of Parameters of the Low-Temperature Molecular-Beam Epitaxy of GaAs; Semiconductors 36, 837 (2002)

[20]   V. V. Preobrazhenskii, M. A. Putyato, O. P. Pchelyakov, B. R. Semyagin: Experimental Determination of the Incorporation Factor of As4 during Molecular Beam Epitaxy of GaAs; Journal Crystal Growth 201/202, 170 (1999)

[21]   A. Adachi: Material Parameters of In1-xGaxAsyP1-y and Related Binaries; Journal Applied Physics 53, 8775 (1982)

[22]   E. Villaggi, C. Bocchi, N. Armani, G. Carta, G. Rosetto, C. Ferrari: Deviation from Vegard Law in Lattice-Matched InGaAs/InP Epitaxial Structures; Japan Journal Applied Physics 41, 1000 (2002)

# Performance Enhancement of Air-cooled Condensers

## M. M. Awad [*], H. M. Mostafa [**], G. I. Sultan [*], A. Elbooz [*], A. M. K. El-ghonemy [**]

* Faculty of Engineering, Mansoura University, Egypt
** Higher Technological Institute, Tenth of Ramadan City, Egypt

*Abstract: Heat transfer by convection in air cooled condensers is studied and improved in this work. In order to enhance the performance of air cooled condensers, it is important to take into consideration both of condensation inside condenser tubes and convection outside, where the enhancement in convection side is the dominant one. Aluminum extruded micro-channel flat tubes improve the performance of condensation more than conventional circular tubes but still has potential for air side improving. So the enhancement of convective heat transfer in air side is achieved in this study by inclination of the flat tubes by a certain angle with respect to horizontal in two cases. The first proposed case is to make convergent and divergent channels for air flow (case 1), while the second case is tilting all tubes in parallel to each other (case 2). A parametric study is performed to investigate the optimum inclination angle (β) and aspect ratio (Ar). Mathematical modeling for air cooled condensers was applied to aluminum flat tubes to study and evaluate these proposed two cases. A computational fluid dynamic software (CFD) is used to solve the problem. Theoretical results show that the optimum angle for the proposed two cases is about 4 deg. With corresponding aspect ratio of 0.58. This leads to enhancement of heat transfer coefficient by factor (Kh) of 1.469 and 1.46 against increase in pressure drop factor (KP) of 2.12 and 1.95 for case 1 and case 2 respectively.*

## 1 Introduction

Air-cooled finned-tube condensers are widely used in refrigeration and air-conditioning applications. For the same amount of heat transfer, the operation of air cooled condensers is more economic as compared with water cooled condensers [1]. Typically air-cooled condensers are of the round tube and fin type. To improve the performance of air-cooled condensers multiple techniques can be achieved such as enhancements on inner pipe surface, changing the tube geometry from round to flat shape and external fins.

A micro-channel flat tubes heat exchanger is one of the potential alternatives for replacing the conventional finned tube heat exchangers. This kind of heat exchangers is made of a flat tube with several independent passages in the cross-section, and formed into a serpentine or a parallel flow arrangement. In these heat exchangers, a multitude of corrugated fins with louvers are inserted into the gaps between flat tubes. The flat tube design offers higher thermal performance and lower pressure drop than the finned-round tube heat exchangers [2]. Brazed aluminum heat exchanger is made from micro-channel flat tubes in parallel to each other which is called parallel flow heat exchanger (PFHE). As a result of its superior performance, some companies in heating, ventilating and air conditioning are considering the flat tube heat exchanger as a high efficiency alternative in order to save electricity when used in window and split type air conditioners which consume large amounts of electricity and contribute to the severe electricity shortage in the peak period. The key advantage of the brazed aluminum design is smaller size and lower weight than finned-round tube condensers. The heat capacity of a parallel-flow heat exchanger (PFHE) is 150-200% larger than that of the conventional heat exchanger [3]. This high heat capacity of the PFHE can meet the requirements of compactness and lightness. Oval and flat cross-sectional tube for finned tube heat exchangers provides a higher heat transfer performance as compared to those formed with round tube geometry as mentioned by Chang et. al. [1]. The effect of tube profile change from round to flat shape on condensation has been investigated experimentally by Wilson et. al [4].They used horizontal copper tubes were initially round with 9.52 mm outer diameter and 8.91 mm inner base diameter. The tubes were successfully flattened into an oblong shape with inside heights of 5.74, 4.15, 2.57, and 0.974 mm. Refrigerants R-134a and R-410A were investigated over a mass flux range from 75 to 400 kg/m$^2$. S, and quality range from approximately 10-80%. They summarized the following results:

1  For a given mass flow rate, there is a significant reduction in refrigerant charge due to flattened tubes.

2  The pressure drop increases as the tube profile is flattened at a given mass flux and quality.

3  There is enhancement of condensation heat transfer coefficient as the tube profile is flattened.

4  Heat transfer enhancement is dependent on the mass flux, quality, and tube profile.

The condensation of refrigerant in multi-port micro-channel extruded tubes has been investigated by many authors [5-7]. All of them concluded that the micro-channel flat tube enhance the inside heat transfer many times than conventional round one. So the present work is mainly concentrated on air side heat transfer from flat tube condensers which is the dominant one.

Although, the PFHE has the above mentioned good thermal performance, but there is still a lot of potentials for improving the air side convective heat transfer which is the dominant one. Therefore, the present study is directed to enhance the convection side heat transfer by inclination of its flat tubes, one is inclined towards clockwise and the next in counter clockwise direction by angles up to 16 deg with respect to horizontal to make convergent and divergent channels for air flow (case 1). Furthermore, without the need of replacing any equipment of production line that producing PFHE, another construction for inclination of all tubes by the same angle range (0:16 deg.) but all tubes are kept in parallel with each other (case 2) is also included in the present study. Finally the effect of aspect ratio (Ar) has been investigated at the optimum inclination angle (β). The best choice for correct range of inclination angles from 0:16 deg that leads to enhancement was obtained from the researches [8-9].

**Nomenclature**

**h**       air-side heat transfer coefficient, $W/m^2 . K$

**Dh**    hydraulic diameter, mm

**Vf**     Air face velocity, m/ s

**H**      transverse pitch of parallel tubes. mm

**L**       Width of flat tube cross section, mm

**Ar**     Aspect ratio =H/L

**P β**   pressure drop for case of inclined flat tubes by angle β, Pa.

**P0**    pressure drop for case of parallel flat tubes with β=0 deg.

**β**      inclination angle of flat tubes with respect to horizontal, deg.

**Re**    Reynolds number, dimensionless

**ΔP**    pressure drop, Pa

**PFHE** parallel flow heat exchanger (Aluminum Brazed heat exchanger = PFHE or serpentine flow heat exchanger)

**η**      overall performance =Kh/KP

**Kh**    Enhancement factor of h = hβ/h0

**KP**    pressure drop increase factor= Pβ/P0

*Subscripts*

*av*    average

# 2   Mathematical Model

Many industrial applications, such as air cooling in the coil of an air conditioner, can be modeled as two-dimensional heat flow. All pre-generated meshes for the studied cases were prepared first by GAMBIT software. Then modeled as bank of tubes in cross-flow, and the air outside flow is classified as turbulent and steady.

The model is used to predict the Flow and temperature fields that result from convective heat transfer. Due to symmetry of the tube bank, only a portion of the geometry was modeled in FLUENT. Domain is discretized into a finite set of control volumes or cells. General transport equations for mass, momentum and energy are applied to *each* cell and discretized. The governing equations are solved to the studied flow field. The numerical solution was conducted to investigate the influence of inclination angle (β) and aspect ratio (Ar) on the performance of air cooled condensers.

The following values which applicable to window and split air conditioning systems, are used as input data for solving the studied problem:

1   Air flow is steady, 2 dimensional and turbulent

2   Air face velocity (Vf)=2.5, 5 and 7.5 m/s

3   The condenser saturation temperature of refrigerant=323 K. Hence study is based on constant wall temperature=323 K

4   Ambient air temperature=308 K

5   The flat tube condenser configurations: tube height (b)=1.8 mm, tube width (L)=18 mm, tubes transverse pitch=10.4 mm.

## 2.1   Numerical Technique

Flow and heat transfer characteristics is obtained for forced convection of air flow across flat tubes at different operating parameters. By using CFD software, the flat tubes condensers shown in Fig. 2a has been studied first, which is called parallel flow heat exchanger (PFHE). Then the proposed modifications in the following sequence: Case 1: construction of convergent and divergent channels for air flow through inclination of flat tubes by angle up to 16 deg. With respect to horizontal, one is inclined towards clockwise and the next in counter clockwise as shown in Fig. 2b.

Case 2: Tilting of all tubes in parallel to each other by angle up to 16 deg with respect to horizontal) either forward or backward as illustrated in Fig. 2c.

# 3    Results and Discussions

In order to study the performance of the proposed two cases, the obtained results are presented relative to those of parallel flat horizontal tubes at the same operating conditions. Contour lines for temperature and velocity in axial direction are shown in Figs. 3 and 4 for flat horizontal tubes, convergent divergent (case 1), and tilted sections of tubes (case 2). Generally, it is observed from Fig. 3 that there is a decrease in fluid temperature towards the centre between pipes in flow direction as the flow is developing. Also, it is found from Fig. 4 for the case of convergent divergent passage, the velocity increases in convergent passes and decreases in divergent passage.

Figure 5 shows the contour of local surface heat transfer coefficient (h) for the same studied cases. It is clear from this figure that the local values of heat transfer coefficient (h) for the studied two cases are higher than those of horizontal tubes The effect of inclination angle β on the performance of flat tube air cooled condenser is illustrated in Figs. 6a and 6b for the studied two cases compared with flat horizontal tubes. As shown in Fig. 6a the increase in ΔP is small in the first part up to 8 deg. then ΔP increases sharply. Therefore it is preferable to operate in this first part. Also, it is found that there is a peak value at β =4 deg. Also, there is a higher values for both of hav, ΔP in the second part of the curve which is not preferable practically. In the other hand, for case 2 presented in Fig. 6b with increasing β up to 12 deg there is a continuous increase in both of hav, Δp. Then for β more than 12 deg. Leads to decreasing in hav and ΔP.

To compare between the two cases and to choose the optimum β, it is important to evaluate the enhancement process as a whole. Therefore the effectiveness of the process (η) (which is defined as η=(Kh/Kp) is plotted against β for the studied cases in Fig. 7b at different values of air face velocities (Vf). It is clear from Fig. 7a and Fig. 7b that for different velocities (Vf), the optimum β for both cases is 4 deg. Also to study the effect of the tested face velocities (Vf), it is clear from Fig. 7b that varying Vf from 2.5 up to 5 m/s leads to considerable change in the performance as a whole (η). But with increasing Vf from 5 to 7.5 m/s, there is a small change in the performance as the two curves are nearly coincident. So, the value of 5 m/s is considered as the max. limit for operation (practical value).

Then for the investigated optimum β, the effect of aspect ratio (Ar) on the performance is shown in Fig. 8. Easily the optimum (Ar) is chosen at value of 0.58 which corresponds to H=10.4 mm.

Finally, to verify the obtained results, a comparison with similar researches that investigated experimentally are shown in Figs. 9a and 9b.

For case 1: The convective heat transfer to air flow in converging-diverging tubes were studied experimentally by Ariad et. al. [8]. Their study was based on constant wall temperature at different values of β from 0 up to 16 deg., which is

similar to the proposed studied cases. They reported that the obtained enhancement comparing to equivalent straight tube at the same mean diameter is Kh=1.45 against KP of 2.2 value. The corresponding values (at 4 deg., Vf=5 m/s) for the present proposed cases are, Kh=1.469, 1.46 against KP=2.12, 1.9 for, convergent –divergent and tilting case respectively. Also the experimental optimum β was 5,30' is agreed with the present one obtained theoretically (4 deg). As shown in Fig. 9a both of present theoretical results and experimental results [4] for Nu are plotted against β. which demonstrates a good agreement.

For case 2: The effect of inclination angle on the performance of aluminum brazed heat exchanger was investigated experimentally by Kim M. H. [10]. From Fig. 10b the comparison of results is showing acceptable agreement. Also they reported that there is enhancement in hav. With increasing β up to 12 deg which agreed with the present results.

## Conclusion

1   Using the proposed convergent divergent construction of heat exchanger with optimum angle of 4 deg offers the best enhancement in heat transfer coefficient. For one row coil which is used in car air condition, the enhancement factor is about Kh= 1.467 with increase in pressure drop (KP factor=2.12).

2   To keep the production line that manufacturing the PFHE, the proposed construction of tilting the all tubes in parallel by 4 deg with respect to horizontal is recommended. This leads to enhancement factor of Kh=1.46 with increase in pressure drop of KP=1.9.

3   This proposed heat exchanger is the strong candidate for use in industrial applications, which is named 'convergent-divergent flow heat exchanger'.

## References

[1]   Y. P. Chang, R. Tsai, J. W. Hwang: Condensing Heat Transfer Characteristics of Aluminium Flat Tubes, in Applied Thermal Engineering, 1997, Vol. 17, No. 11, pp. 1055-1065

[2]   R. L. Webb, X. M. Wu: Thermal and Hydraulic Analysis of a Brazed Aluminum Evaporator, in Applied Thermal Engineering, 2002, 22: 1369-1390

[3]   K. Chung, K. S. Lee, W. S. Kim: Optimization of the Design Factors for Thermal Performance of a Parallel Flow Heat Exchanger, in Int. J. of H&M Transfer, 2002, 45: 4773-4780

[4]   M. J. Wilson, T. A. Newell, J. C. Chato, C. A. Ferreira: Refrigerant Charge, Pressure Drop, and Condensation Heat Transfer in Flattened Tubes, in Int. J. Refrigeration, 2003, 26: 442-451

[5]     J. W. Coleman, S. Garimella: Characterization of Two Phase Flow Patterns in Small Diameter Round and Rectangular Tubes, in International Journal of Heat and Mass Transfer, 1770, 31 "0888# 1758

[6]     E. Bari, J. Y. Noel, G. Comini, G. Cortella: Air-cooled Condensing Systems for Home and Industrial Appliances, in Applied Thermal Engineering, 2002, 25:1446-1458

[7]     H. K. Varma, C. P Gupta: Heat Transfer during Forced Convection Condensation inside Horizontal Tube, in Int. J. Refrigeration, 1995, 18:210-214

[8]     F. F. Araid, M. A. Shalaby, M. M. Awad: Convective Heat Transfer to Gas Flow in Converging Diverging Tubes, in Mansoura University Bulletin, June 1986, V, 11, No. 1

[9]     L. H. Rabie, A. A. Sultan, Y. E. Abdel Ghaffar: Convective Heat Transfer and Pressure Drop in Periodically Convergent Divergent Variable Area Annuli, 2001, IMPEC12

[10]    M. H. Kim, S. Song, C. W. Bullard: Effect of Inlet Humidity Condition on the Air Side Performance of an Inclined Brazed Aluminium Evaporator, in I. J. of Refrigeration, 2002, 25: 611-620

[11]    G. Lazza, U. Merlo: An Experimental Investigation of Heat Transfer and Friction Losses of Interrupted and Wavy Fins for Fin-and-Tube Heat Exchangers, in Int. J. Refrigeration, 2001, 24: 209-416

[12]    S. Sanitjai, R. J. Goldstein: Forced Convection Heat Transfer from a Circular Cylinder in Cross Flow to Air and Liquids, in International Journal of Heat and Mass Transfer, 2004, 47: 4795-4805

[13]    J. S. Jabardo, G. W. Manami, M. R. Lanella: Modeling and Experimental Evaluation of an Automotive Air Conditioning System with a Variable Capacity Compressor, in Int. J. of Refrigeration, 2002, 25: 1157-1172

[14]    X. M. Wu, R. L. Webb: Thermal and Hydraulic Analysis of a Brazed Aluminum Evaporator, in Applied Thermal Engineering, 2002, 22: 1369-1390

[15]    C. Y. Yang, R. L. Webb: Condensation of R-12 in Small Hydraulic Diameter Extruded Aluminum Tubes with and without Micro-Fins, in I. J. of H&M Transfer, 2005, 39:791-800

[16]    A. Bensafi, S. Borg, D. Parent: CYRANO: a Computational Model for the Detailed Design of Plate–Fin-and Tube Heat Exchangers using Pure and Mixed Refrigerants, in Int. J. Refrigeration, 1997, Vol. 20, No. 3:218-228

Figure 1

Micro-channel fat tube cross-sectional view (Lxb)



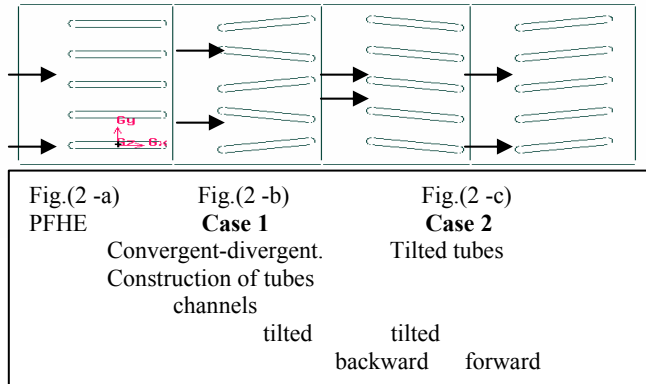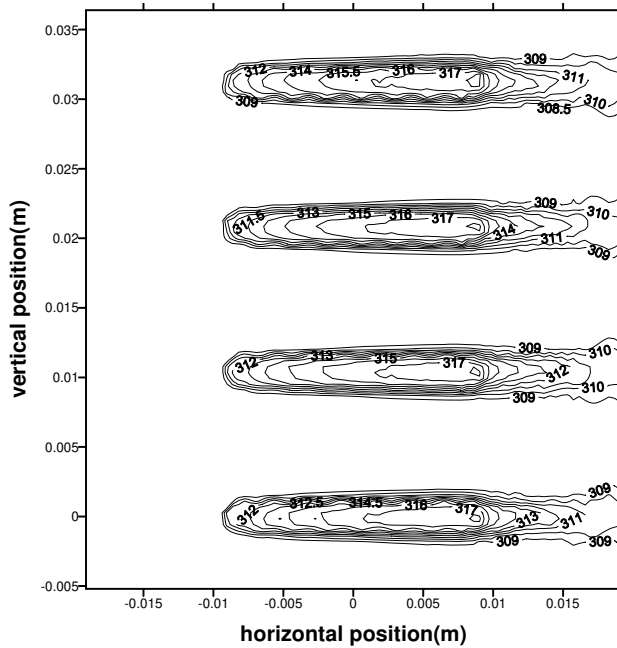| Fig.(2 -a) | Fig.(2 -b) | Fig.(2 -c) | |
|---|---|---|---|
| PFHE | **Case 1** | **Case 2** | |
| | Convergent-divergent. | Tilted tubes | |
| | Construction of tubes | | |
| | channels | | |
| | tilted | tilted | |
| | backward | forward | |

Figure 2

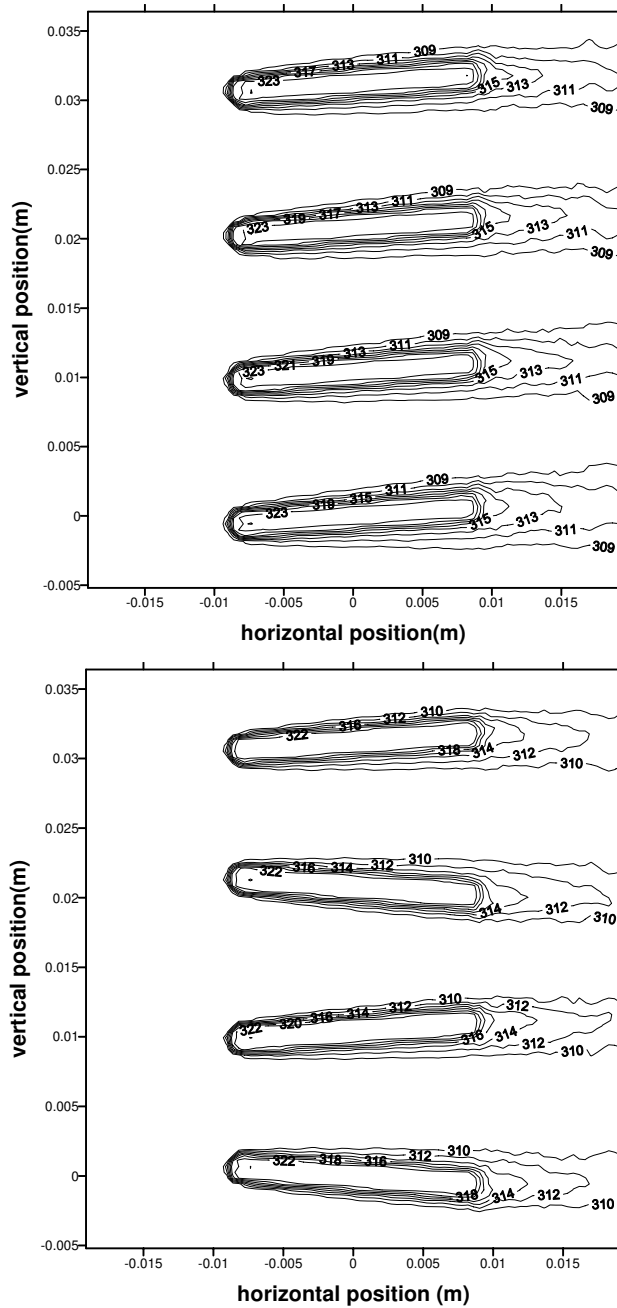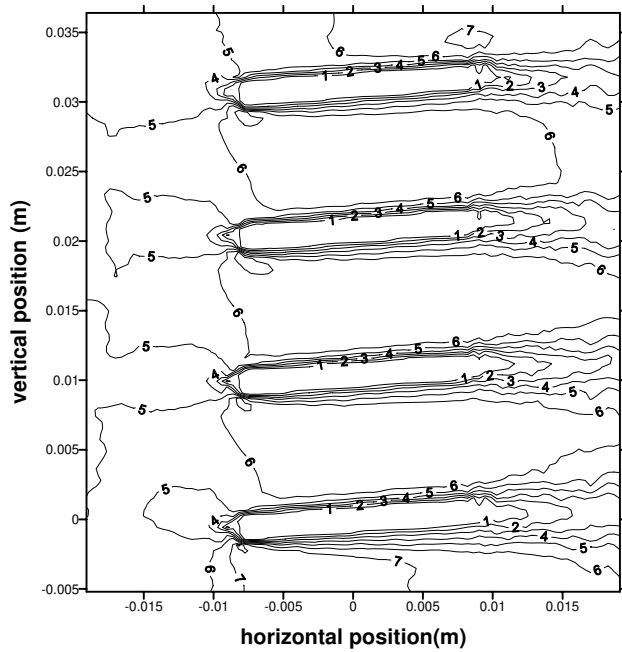Layout of flat horizontal tubes (PFHE) and the proposed two cases of modifications

Figure 3
Temperature contour for the studied two cases compared with horizontal flat tubes case
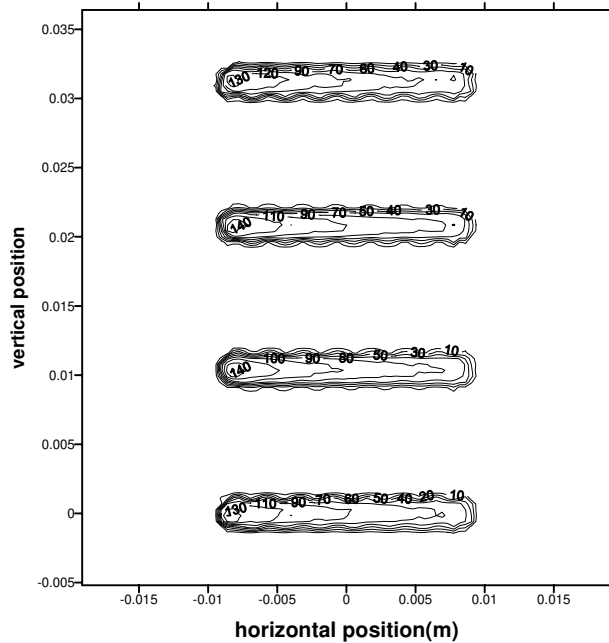
Figure 4

Velocity contour for the studied two cases compared with horizontal flat tubes case
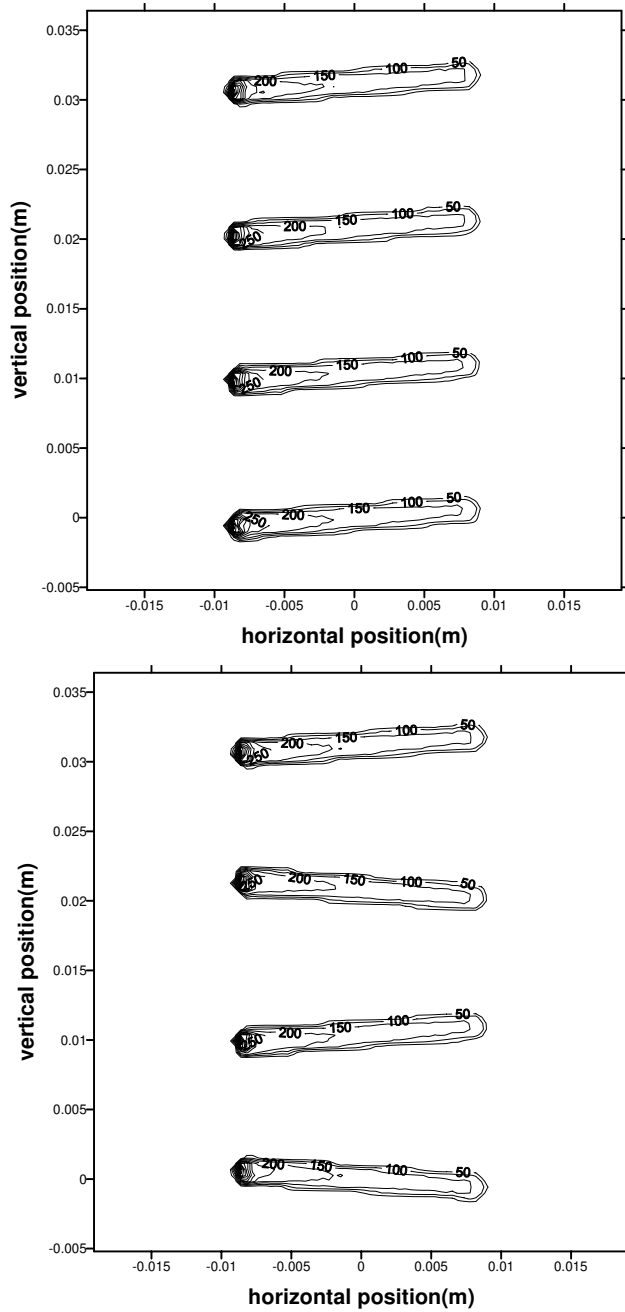
Figure 5

Local heat transfer coefficient contour for the studied two cases compared with horizontal flat tubes
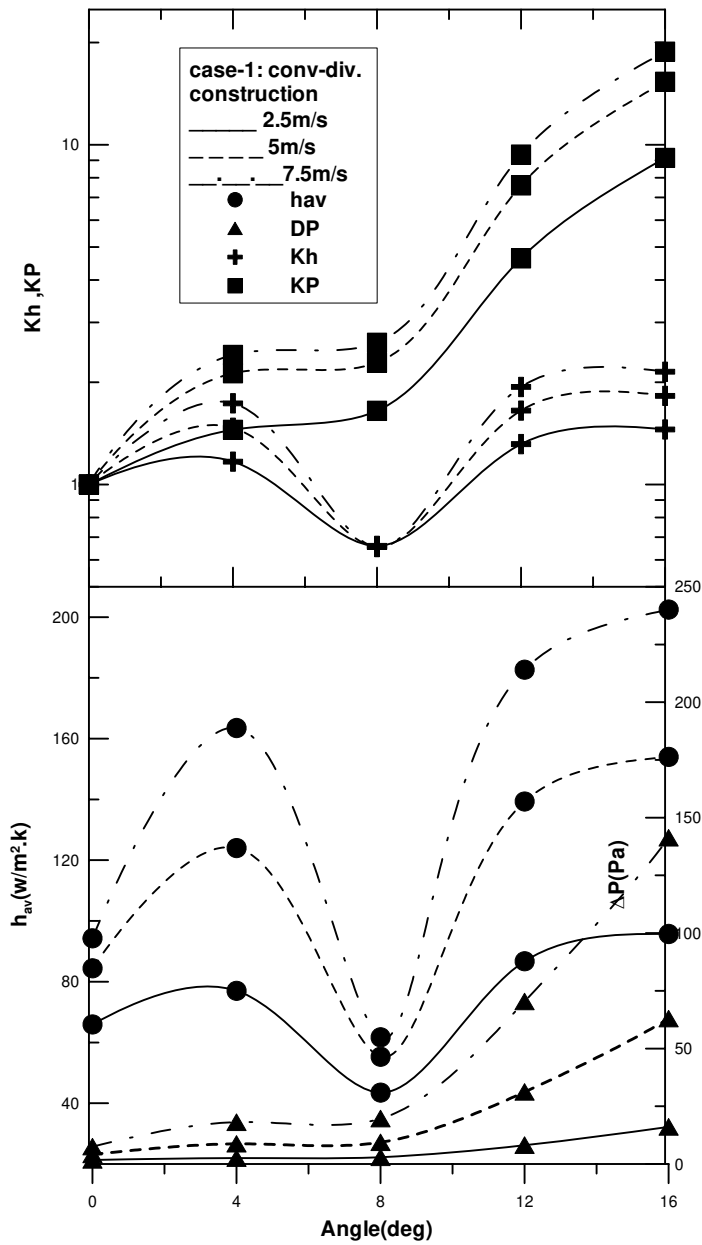
case

Figure 6a

Variation of performance against β for case 1 (convergent – divergent)

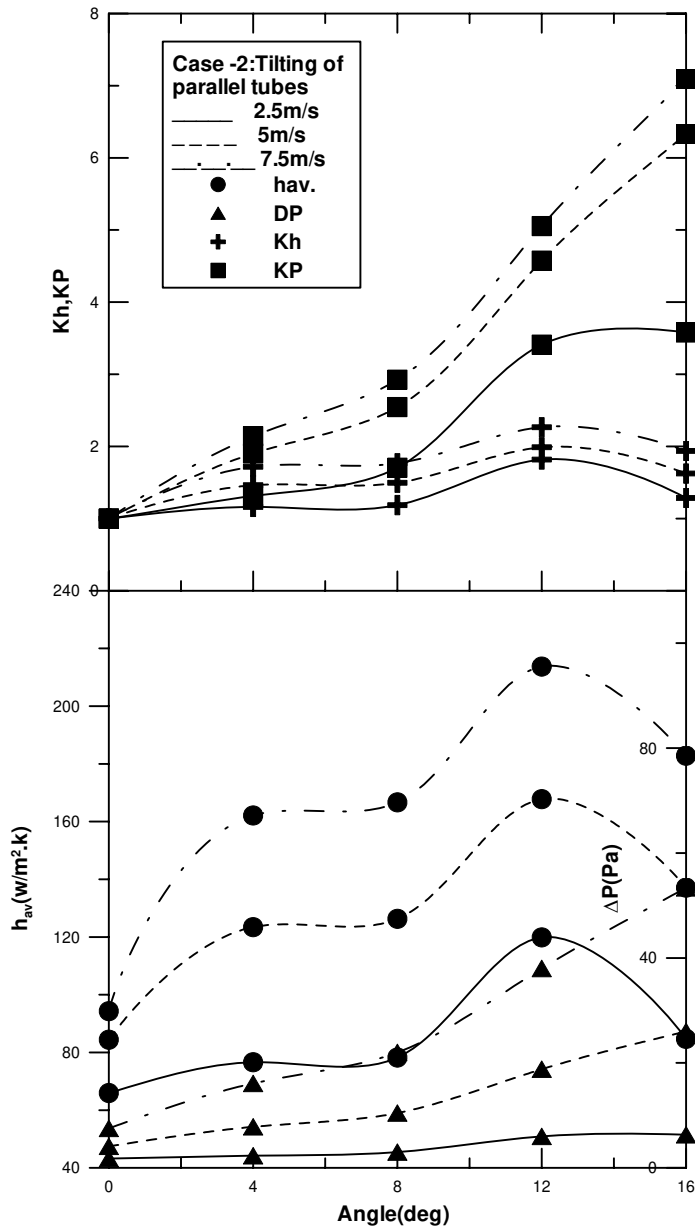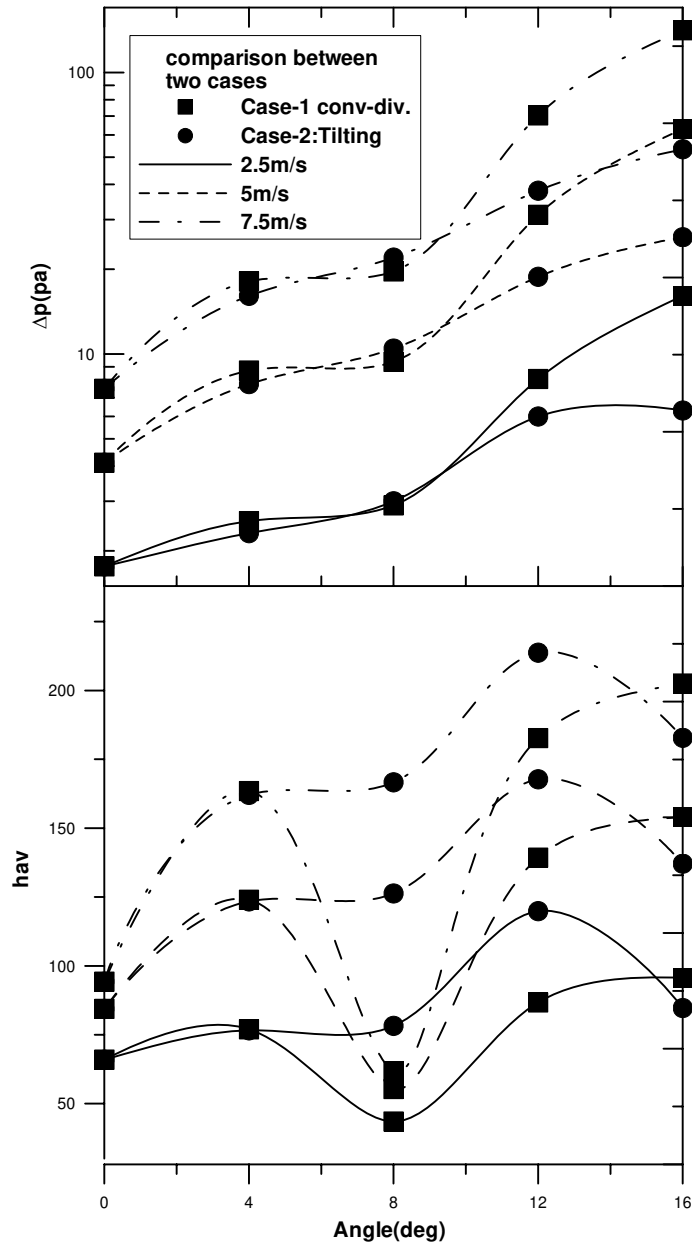Figure 6b

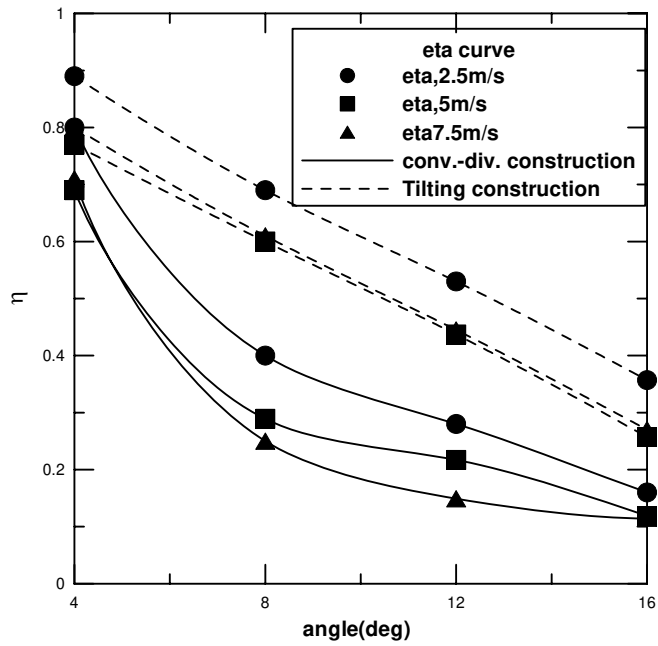Variation of performance against β for case 2 (tilting of tubes)

Figure 7

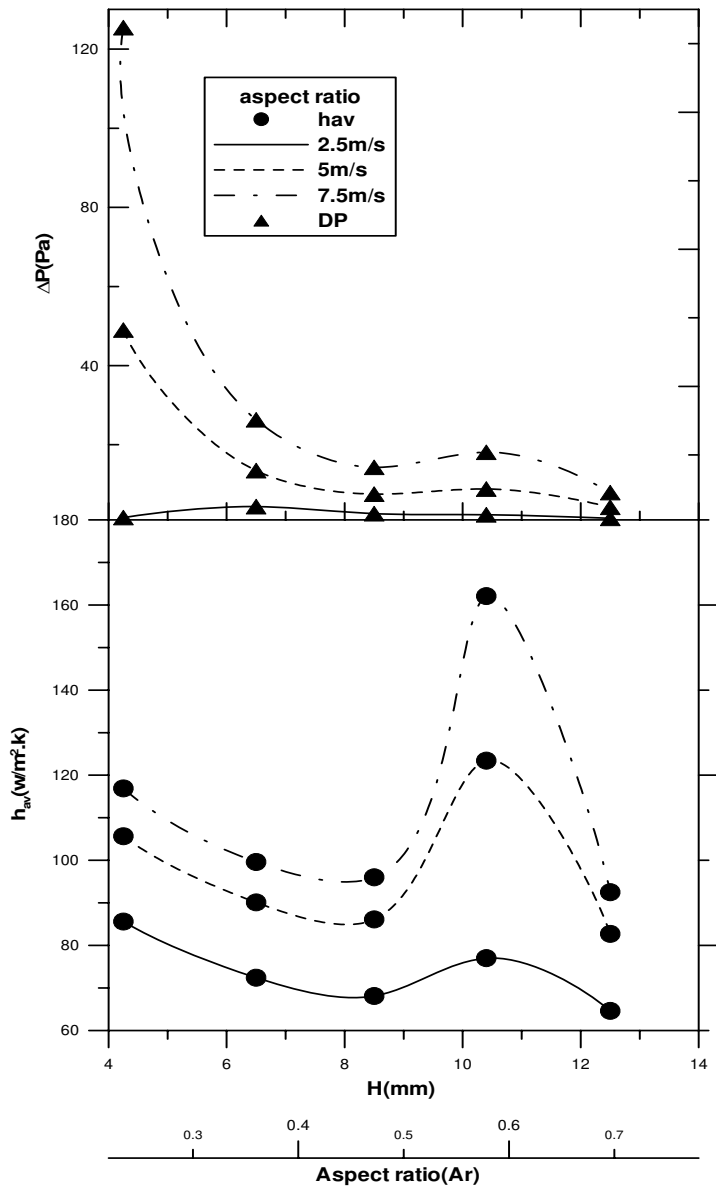Comparison between performance of the proposed two cases

Figure 8
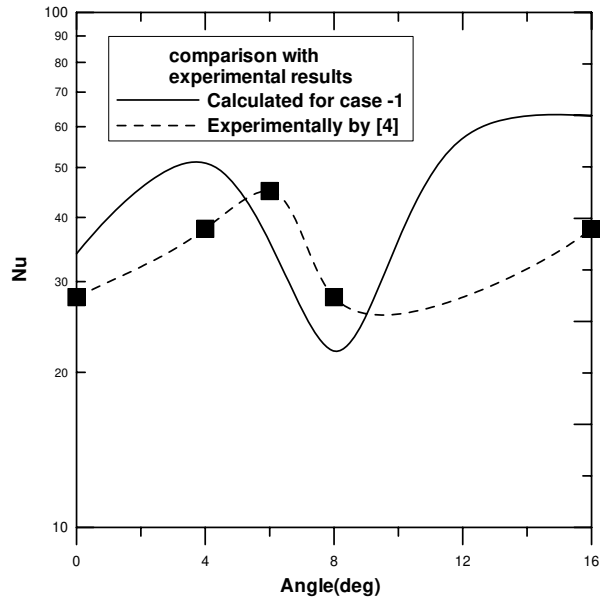Performance variation against aspect ratio

Figure 9a

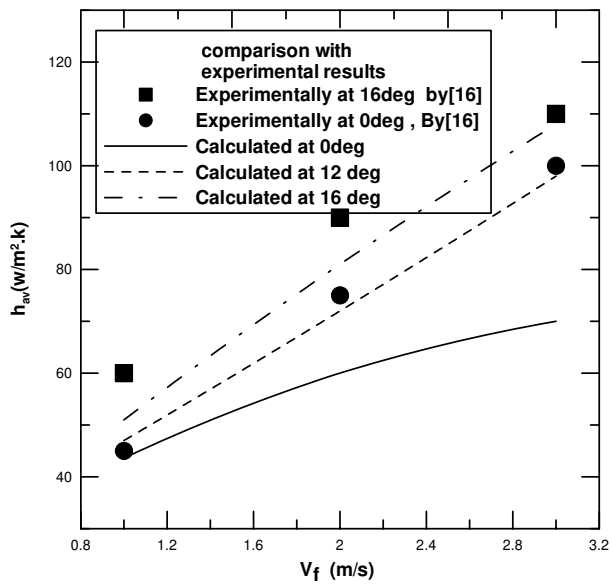Comparison with experimental results for case 1



Figure 9b

Comparison with experimental results for case 2

# Two-Phase Heuristic for the Vehicle Routing Problem with Time Windows

**Sándor Csiszár**

Institute of Microelectronics and Technology
Kálmán Kandó Faculty of Electrical Engineering, Budapest Tech
Tavaszmező u. 17, H-1084 Budapest, Hungary
E-mail: csiszar.sandor@kvk.bmf.hu

*Abstract: The subject of the paper is a complete solution for the vehicle routing problem with time windows, an industrial realization of an NP hard combinatorial optimization problem. The primary objective –the minimization of the number of routes- is aimed in the first phase, the secondary objective –the travel distance minimization- is going to be realized in the second phase by tabu search. The initial route construction applies a probability density function for seed selection. Guided Route Elimination procedure was also developed. The solution was tested on the Solomon Problem Set and seems to be very compeitive with the best heuristics published in the latest years (2003-2005).*

*Keywords: Metaheuristics; Vehicle routing problem; Time windows; Tabu search; Combinatorial optimization*

## 1    Introduction, Problem Definition

At transportation an important question is -next to the in time delivery- the cost of the service. The lowest number of routes is primarily important because it determines the number of vehicles applied. The second priority is the minima of the total travel distance. There are studies where the secondary objective is the minimum schedule time when quick and in time service is more important than the travel distance. The VRP has a large literature, so here only a short introduction is given. Let $G = \{V, E\}$ be a graph and $V = \{v_0, v_1, \ldots , v_n\}$ a set of vertices representing the customers around a depot, where vertex $v_0$ denotes the depot. A route is a closed tour, starting from and entering the depot: $\{v_0, \ldots, v_i, v_j, \ldots, v_0\}$. The cost of a tour is: $C_t = \sum_{(route)} d_{i,j}$, where $i$ and $j$ are consecutive customers on the route and $d_{i,j} \in E$. The objectives of the solution are to determine the lowest number of routes -or number of vehicles ($N_v$)- and the lowest cost (total travel distance, $C = \sum_{(s)} C_t$, where ($s$) is the actual solution), provided

that each customer can be visited only once. The specific problem instances are given by the the number of customers, their demand, delivery time windows, service time and customer coordinates. The vehicle capacities are identical and given for each problem type. The distances between the customers are the Euclidian distances. The service must be started within the given time window, the vehicle travel time constraint is determined by time window of the depot. (The violation of the constraints is not allowed in a feasible solution.)

Mathematical formulation of VRP with Time Windows (VRPTW) can be found in the relevant literature [14]. Although there are exact methods [10] [11], their application is limited because the computation time is excessively increasing with the number of customers. The solution of this problem has double objective, and it implies the real advantage of the two-phase solution. The chance of one phase heuristics for route elimination is quite low. The explanation is: neighbourhood solutions don't represent so significant changes that would be required for eliminating quite long routes. In addition the routes become ever longer as the search is progressing. The other reason is, the lowest cost is -sometimes- not at the lowest number of routes. If we analyse the performance of the tabu search (TS) we must admit that despite it is one of the most successful metaheuristics it has difficulties in special cases when the route elimination goes together with cost increment. The objective function of the TS is usually designed for finding cheaper solutions. We can change the objective function and the length of the tabu-list but it is difficult for a pure TS algorithm to get out from 'deep valleys' so the search is basically guided by the secondary objective. In the route elimination respect -although it is the primary objective- the pure TS loses to other –lately developed- hybrid metaheuristics [2].

The remaining part of the paper is structured as follows. Section 2 is about the initial route construction, Section 3 describes the guided route elimination, Section 4 is the second phase by tabu search, and finally Section 5 is about the conclusion, computational results and future plans.

**Notations used in the article:**

$C$    : Total cost of the solution or total travel distance (*TTD*),

$TC_i$ : Time constraint factor of customer $i$,

$\overline{TC}$  : Average time constraint of the problem,

$N_v$   : Number of vehicles (number of routes),

$N_r$   : Number of customers on the actual route,

*ItNo* : Actual iteration number,

$t_{li}$   : Latest time to start service at customer $i$.

$t_{ei}$   : Earliest time to start service at customer $i$.

$t_i$    : Actual arrival time at customer $i$,

$r_i$    : Customer distance from depot,

$r_{max}$   : Depot distance of the farthest customer,

$r_{ij}$    : Distance between customer $i$ and $j$,

$r_i'$      : Relative customer distance from depot,

$n_s$       : Initial seed number for parallel route construction,

$n$         : Total number of customers,

$n_0$       : Number of customers in the seed selection zone;

$t_{wb}$, $t_{wa}$ : waiting time before and after the insertion,

α, λ, ω:  Cost function parameters.

# 2    Initial Solution

At the design of the Initial Route Construction (IRC) the main objective was to obtain quick and good quality solutions (first of all in respect of the number of routes) within reasonable computation time. The intention was to design a heuristic that maps the most important considerations of a human being. These are as follows: (It is noted that nodes and customers are used as synonyms in the article.)

   a)   distance from the depot,

   b)   waiting time before and after the insertion,

   c)   distance between customers affected in the insertion,

   d)   savings (cost difference realized by the actual solution, compared to serving the single customer by another extra vehicle),

   e)   route building in both directions,

   f)   time window of the given customer,

   g)   demand of customers in order to fill the trucks optimally -especially at those problems, where the vehicle capacity represents a strict constraint,

   h)   taking care of not leaving a single customer alone unrouted (otherwise an expensive extra route should be devoted to serving that),

   i)   parallel route construction,

   j)   step back if the partial solution seems to be unfavourable.

## 2.1    The Objective Function

In the literature a known objective function is used for insertion heuristics (Eq. 1) calculating the cost if node $k$ is inserted between node $i$ and node $j$. Items (a,) (b,)

(c,) (d,) (e,) are also considered in Eq. 1. Item (f,) has an important role in the insertion sequence, because the narrower the time window of a certain customer the more difficult it is to insert this customer into a route, so we have to give preference to customers with narrow time window. The time constraint factor ($TC_i$) was introduced in order to realize this preference in the mathematical formulas: $TC_i = (t_{li} - t_{ei})/(t_{l0} - t_{e0})$. The average time constraint factor is also needed for dimensionless description: $\overline{TC} = 1/n \sum TC_i$. The new objective function is formulated in Eq. 2, where ω is used for getting a realistic weighting in order to moderate the drastic time constraint effect.

$$c_k = \alpha(r_{ik} + r_{kj} - r_{ij}) + (1-\alpha)(t_{wa} - t_{wb}) + \lambda r_k \tag{1}$$

$$c_k = (TC_k / \overline{TC})^\omega \left[ \alpha(r_{ik} + r_{kj} - r_{ij}) + (1-\alpha)(t_{wa} - t_{wb}) \right] + \lambda r_k \tag{2}$$

Computational tests show that $\omega = [0.35 \cdots 0.7]$ gives the best results. Customers - having narrower time windows- seems to be cheaper for the algorithm, this way we can stimulate them. This change in Eq. 1 made a perceptible improvement in the results: 1.57% improvement in the number of routes. It worth mentioning that Eq. 2 was tested with reduced time constraint factor: $TC_i = (t_{li} - t_{act})/(t_{l0} - t_{e0})$, where $t_{act}$ is the actual possible time of starting service at the selected customer. It seems to be rational to consider only that part of the time window that is available for service in the actual situation. Even if the customer has a wide time window and it is unrouted at a certain situation, it may be critical for the service if the remained available time for service is short. With the application of the idea no further improvement could be detected on the sample problems.

## 2.2    Seed Selection for Route Initialization

The essence of this concept is the following: seed points are selected from a cirkular-ring zone ( $r' > r_{sb}$ ), those customers that are far from the depot or have a low value of *TC* factor are favoured as seed points. The applied probability function for seed selection is Eq. 3. A detailed description of the probability based seed selection can be found in [6].

$$p_i = \begin{cases} 0 & \text{if } (r'_i < r_{sb}) \\ (r'_i / r_{sb})^a & \text{if } (r'_i \geq r_{sb}) \text{ and } (TC_i > 0.4) \\ (r'_i / r_{sb})^a \; (\overline{TC}/TC_i)^b & \text{if } (r'_i \geq r_{sb}) \text{ and } (TC_i \leq 0.4) \end{cases} \tag{3}$$

Let the relative depot distance be: $r_i' = r_i / r_{max}$, where $r_{max}$ is the distance of the farthest customer from the depot. It seems to be rational to define a minimum radius around the depot, within that radius seeds are not selected. Let's define this radius as a relative seed border: $r_{sb}$. Usually $r_{sb} = 0.5$ is sufficient for this but in

special cases it is necessary to reduce this value according to have sufficient number of customers -min$\{2N_v, 0.35n\}$- in the circular-ing zone. Reduce the actual value of $r_{sb}$ until the above condition is satisfied.

The parameters of Eq. 2 are $\alpha$, $\lambda$, and $\omega$. The following parameter combinations [1] resulted 42 initial solutions: $\alpha = [0.5 \cdots 1]$ in 0.1 steps, $\lambda = [0.5 \cdots 1.7]$ in 0.2 steps and $\omega = 0.5$.

## 2.3   IRC Algorithm

($TC_i$ and $\overline{TC}$ are calculated after data retrieve);

1. **Set** $\alpha$ from the interval $[0.5 \cdots 1]$ and $\lambda$ from $[0.5 \cdots 1.7]$ ;

2. **Set** $r_{sb} = 0.5$ (for the seed border);

3. Preliminary route construction;

    (to find out the preliminary number of routes: $N_{vp}$);

4. Count the number of customers in the seed selection zone ($n_o$);

5. **while** $n_0 < \min(2N_{vp}, 0.35n)$ **do**

    $r_{sb} := 0.975 \cdot r_{sb}$ ;   Count $n_0$;

6. **end while**;

7. **Set** $R := \varnothing$ (set of routed nodes);

8. **Set** $U := \{u_i\}$ $i = 1, 2, \ldots, n$ (set of unrouted nodes);

9. **Set** $S$ (set of unrouted nodes in the seed selection zone);

10. **Repeat**

11.   **Set** $p := -1$ ; $r := 3 \cdot r_{\max}$ ;

12.   **for** ( $\forall i \in R$ ) **and** ( $\forall j \in S$ ) **do**

13.     **If** $r_{ij} < r$ **then**

14.         $p := j$; $r := r_{ij}$ ;

15.     **end if;**

16.   **end for**;

17.   **if** $r < 1.5 \cdot \bar{r}$ **then** $v := p$ (it means that a customer remained unrouted

                    close to the full route in the seed zone, so it is selected)

18.   **else**

19.     **if** ( $p := -1$ ) **then**

20.       **Let** $v := i \mid r_{0i} > r_{0j}, \forall (i, j) \in U$   (there are no unrouted customers

                  in  seed zone, so the farthest one is selected)

21.     **else**

22.      $v$ is selected by Eq. 3

23.     **end if;**

24.   **end if;**

25.    $R \leftarrow \{v\}$; $U := U - \{v\}$;

26.   **If** $r_v \geq r_{sb}$ **then** $S := S - \{v\}$ **end if**;

27.   Initialize a route $k$: 0 - $v$ - 0;

28.   $v$:=Select node for insertion using Eq. 2;

29.   **while** ( $v \neq \varnothing$ ) **do**

30.    Insert node $v$ into route $k$;

31.     $R \leftarrow \{v\}$; $U := U - \{v\}$;

32.    **If** $r_v \geq r_{sb}$ **then** $S := S - \{v\}$ **end if**;

33.    $v$:=Select node for insertion using Eq. 2;

34.   **end while;**

35. **Until** ( $U \neq \varnothing$ )

The described algorithm was embedded in a cycle to compute the 42 initial solutions, then the best one (with the lowest number of routes) was selected for further processing. 10 computational runs were made on each of the 56 Solomon problems and compared to the best results found by metaheuristics. The average of the ever found best number of routes is 7.23 the same figure of IRC was 7.80. The result of the initial route construction algorithm is good also in the best one comparison.

# 3   First Phase

An exercise was made in the study how a general attribute of the problem can guide the search and how successful it is. This attribute is the total cost of the solution.

## 3.1   The Route Elimination

The developed route elimination algorithm is a recursive procedure that applies the *in depth-first search*. Depth of the search tree ( $6 \cdots 8$ ) depends on the average time constraint ( $\overline{TC}$ ), because the wider the time windows are the higher the complexity of the search is.

The most promising route has to be selected first for elimination. Three types of route selection methods are used. The *first method* is based on the number of nodes on the route (the shorter routes are preferred). The *second* one takes also the insertion frequency of the nodes into account. It can not be used at the beginning of the search and 65-35% weighting is applied between the number of customers on the route and the insertion frequency by the following equation:

$$selCrit = \min[0.65 N_r \cdot N_v / n + 0.35(1 - N_v /(ItNo \cdot N_r)\sum ins)] \qquad (4)$$

The *third* route selection procedure based on the route selection frequency. This latest one prefers those routes that are selected rarely. The route selection is controlled by the block management unit. The successful insertion frequency and the number of route elimination trials per route are collected in a database. These data are used for two purposes: route selection for elimination, insertion sequence of nodes. If the remained nodes after trial have a favourable insertion statistics it means the chance for success is higher at the post search.

The *in depth-first search* is executed within given cost limits, otherwise expensive insertions and changes cause dramatic total travel distance increment. The only exception is the last node, provided that all the previous insertions were successful -in this case the cost limit is not considered. If a certain route is selected for elimination all its nodes are tried to be inserted onto other routes. The first task is to determine the insertion sequence of the customers on the selected route. As the computational trials show it is an important part of the route elimination algorithm.

Insertions are tried with reasonable cost limit: $[2 \cdots 2.6] \cdot \bar{r}$. The purpose of this limit is to avoid drastic cost increment that would hinder further insertions. If the insertion fails then 3-Opt insertions are tried provided the time windows are wide. A learning process is embedded in the route construction phase that counts the success of the 3-Opt intra route changes and these data are used for making decision about the 3-Opt insertions. If the success ratio reaches a certain percent the 3-Opt internal reordering is used before insertions otherwise not. If 3-Opt insertions fail intelligent reordering is tried [1].

Until no unsuccessful insertion occurs, in case of failure a repair procedure is activated after a couple of cost reduction steps. First the graph has to be modified by the TS algorithm in order to reduce the total travel distance. On both sides around the unsuccessfully inserted customer all the routes have to be identified in two times 40° sector (or in a user given angular domain). Try to combine these routes every feasible way (2-Opt) and at each route combination try the search again.

After the route elimination if a reasonable number of customers remains unrouted and their time constraint factors and depot distances satisfy certain criteria a Post Search is taking place. The evaluation of the remained nodes is made by an

algorithm that considers the depot distance, *TC* factors and also the insertion data provided the number of iterations is higher than a defined value ($5 \cdot N_v$).

If the route elimination was not successful and only a few customers –less then $(0.2 \cdots 0.3)n / N_v$ - remained unrouted it seems to be rational to fill the route up as much as possible in order to draw off customers from other routes and at the same time to increase diversification. The filling up is done by varying the parameters ($\alpha$, $\lambda$, $\omega$) in Eq. 2.

## 3.2    Concept of Guided Route Elimination (GRE)

GRE is the core procedure of the developed route elimination. During the route elimination a list is used to prevent evolution of circles in customer exchange. It was observed that a number of nodes don't move from their place during the route elimination. That means the search is limited to a certain neighbourhood. Obviously it comes from the neighbourhood graph definition that only a limited subset of solutions can be reached from a given solution. The main idea was to enforce the unmoved nodes to a certain extent. It is well known at the tabu search if no cheaper solution is found in the neighbourhood -provided the tabu list is not hurt, in this case those nodes are penalized which move frequently. The penalization is made by the following equation: $d_k^{'} = d_k - const \cdot f_m$, where $f_m$ is the move frequency registered in a database. The value of the constant $\begin{bmatrix} 5 \cdots 10 \end{bmatrix}$ has no effect on the search. In this case the modified cost ($d_k^{'}$) is used for selection instead of the real cost ($d_k$) and the Tabu Search finds the move of the least penalized customers cheaper and prefers their move to reveal new regions in the solution space.

The essence of the GRE concept is: not only the moves enforced by tabu search but those evoked by the '*in depth first search*' are counted and added to the move frequency. This way the unmoved nodes are forced. In order to moderate this effect the cost increment has to be controlled. The Search management checks regularly the total cost and compares it to the initial cost. If the relative cost increment is higher than the user defined value $\begin{bmatrix} 1.1 \cdots 1.2 \end{bmatrix}$ the customer move frequency data of TS are readjusted to the starting value (100). For this reason the algorithm is divided into search blocks. The starting value was defined to 100 because it gave a realistic block cycle number $(45 \cdots 50)$ and the move ratio between the TS and the GRE is about $80 / 20\%$.

## 3.3    A Route Elimination Cycle

1.  **Let** $p_1 = Random(100)$ be; (start of the route selection)

2.  **if** ($p_1$ < 30) **then**  Select one of the rarely tried routes (based on statistics)

3.  **else**

4.     **if** (no statistical data available) **then**

5.       Find the route with the lowest number of customers;

6.     **else**

7.       Selection based on statistical data;

8.     **end if** ;

9.  **end if** ; (end of the route selection)

10. *firstFail*:=False;

11. *costContrCount* := 0;

12. **Repeat**

13.    Select customer for insertion (using MSA);

14.    *InsOK*:=False;

15.    **if** Recursive insertion is successful then *InsOK*:=True **end if**;

16.    **if** Not *InsOK* **then**

17.      **if** 3-Opt Recursive insertion is successful **then** *InsOK*:=True **end if**;

18.    **end if**;

19.    **if** Not *InsOK* **then**

20.      **if** Intelligent Reordering recursion is successful **then**
        *InsOK*:=True **end if**;

21.    **end if**;

22.    **if** Not *InsOK* **and** Not *firstFail* **then**

23.      Cost reduction cycles;

24.      *costContrCount* := 0;

25.      **if** Not Repair algorithm (recursion) is successful **then**
        *firstFail*:=True **end if**;

26.    **else**

27.      Inc(*costContrCount*);

28.    **end if**;

29    **If** (*costContrCount* > *allowedNo*) **then**

30      Cost control cycles;

31      *costContrCount* := 0;

32    **end if**;
    (The *allowedNo* depends on the initial number of customers
    on the route.)

33.  **Until** No more nodes is found;

34.  **If** (No. of unrouted customers>0) **then** Post Search **end if**;

## 3.4   Multi Strategy Application

This is a successful tool used several time in the solution. Route selection for elimination, the route elimination cycle and the node insertion sequence combines several strategies. According to the experiences an intelligent strategy works well at a problem instance and fails at the other. Certain strength can be the weakness at another problem type. The question is: how can we maximize the exploitation of the good strategies? The following method is proposed:

-   Find out promising strategies.

-   Decide on the priority order (or probability values) and the application rules.

-   Use the application rules and probability functions for the application of strategies.

-   Change the probability values during the search according to the success ratio (optional learning).

## 4   Second Phase: Tabu Search

In the second phase of the developed metaheuristic the same tabu search was applied for finding cheaper solution as it was used for controlling the *TTD* at the guided route elimination. The basis of this TS solution was the reactive tabu search [3] with some differences, first of all in the used operator set and the tabu list. The extinct arcs are stored on the tabu list instead of route and node identifiers. This is less restrictive, although in case of arc (*i-j*) on tabu the re-creation of this arc is not allowed, but if the route identifier and the relevant node is stored on tabu it is possible on another route.

The used insertion and the interchange operators are always evaluated using global best (GB) strategy. If no cheaper solution is found the next operators is used in the following order: intra route 3-Opt (Fig. 1), interchange-21 (two nodes from one route and one from the other are swapped if it is feasible), interchange-22, interchange-32. These operators slow down the procedure but gives better results. The tabu list tenure is managed in the interval of $[6\cdots12]$. At the beginning of the second search phase the tabu list is initialized with an empty list and the best solution is taken from the first phase. The aspiration criteria is applied, if the found solution is better than the global best the tabu list is neglected. In case of a significant cost or waiting time reduction intensification is

made realizing that by an intensification list. A maximum of 3000 iteration steps was adjusted.
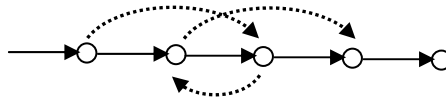


Figure 1
Intra route 3-Opt changes

# 5    Computational Results, Conclusion, and Future Plans

The presented two-phase algorithm (HGRE) is tailored to middle size VRP with Time Windows. A summary of its main contributions are as follows:

1   The objective function (Eq. 2) qualifies the unrouted nodes for insertion taking also the width of the time windows into account (Eq. 1). A similar concept was applied at the seed selection for route initialization [6].

2   It was found the continuous control of Total Travel Distance (*TTD*) aids route elimination. A Guided Route Elimination concept (GRE) was realized to explore the solution space as much as possible.

3   A learning process was built into the initial route construction phase (the success ratio of 3-Opt intra route changes is collected and calculated) to use that at the route elimination in order to save computation time. If the success ratio reaches a given limit the 3-Opt insertions are tried otherwise not.

The computer program for the developed two-phase algorithm was written on Delphi platform by dynamic memory programming and was tested on the Solomon Problem Set on 1.7 GHz computer. The user surface of the developed computer probram can be seen on Figure 2. A maximum pre-adjusted search time was 30 minutes. At the initial route construction and at the route filling up α, λ, ω were varied in order to generate initial solutions, the minimum tabu tenure was 6 the maximum was 12. The cost limit used at the *in depth-first search* was: $r_{ik} + r_{kj} \leq 2.6\bar{r}$. It is quite difficult to find perfect comparison, first because of the continuous improvement of the computer platforms, processors and RAM capacities, secondly because of the differences in the used parameters, the optimization criteria and the programming languages. The performance of the developed HGRE seems to be very competitive with the lately published best algorithms found in the literature (Table 1).

| Problem class | Mean Values | BBB (2003) | BH (2004) | HG (2005) | *HGRE (2006)* |
|---|---|---|---|---|---|
| C1 | MNV | 10.00 | 10.00 | 10.00 | *10.00* |
|    | MTD | 828.48 | 828.38 | 828.38 | *828.38* |
| C2 | MNV | 3.00 | 3.00 | 3.00 | *3.00* |
|    | MTD | 589.93 | 589.86 | 589.86 | *590.32* |
| R1 | MNV | 11.92 | 11.92 | 11.92 | *11.92* |
|    | MTD | 1221.10 | 1211.10 | 1212.73 | *1227.89* |
| R2 | MNV | 2.73 | 2.73 | 2.73 | *2.73* |
|    | MTD | 975.43 | 954.27 | 955.03 | *987.91* |
| RC1 | MNV | 11.50 | 11.50 | 11.50 | *11.50* |
|     | MTD | 1389.89 | 1384.17 | 1386.44 | *1400.78* |
| RC2 | MNV | 3.25 | 3.25 | 3.25 | *3.25* |
|     | MTD | 1159.37 | 1124.46 | 1108.52 | *1156.63* |
| All | CNV | 405 | 405 | 405 | *405* |
|     | CTD | 57952 | 57273 | 57192 | *58239* |
|     | Tot.Time | 1800 | 7200 | - | 2798 |

Table 1

Comparison on Solomon benchmark problems. BBB: Berger et al. [16], BH: Bent and Van Hentenryck [17], HG: Homberger and Gehring [9]
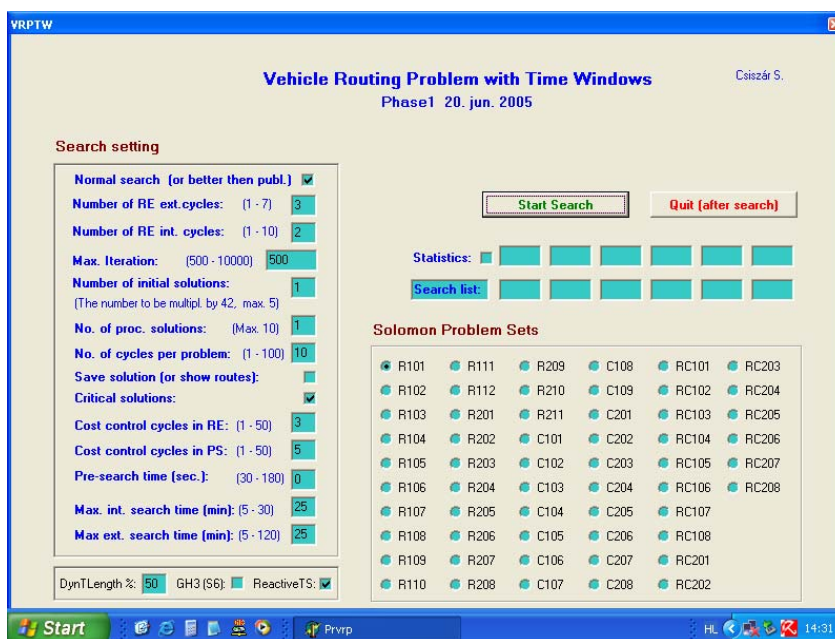


Figure 2
HGRE application

The best results were selected from a series of 10 independent runs per problem instance then the mean number of vehicles (*MNV*), mean travel distance (*MTD*) and the mean computation time per problem type (*MCT*) were calculated. Additionally the cumulated number of vehicles (*CNV*) and the cumulated travel distance (*CTD*) are reported. Bold letters are used if the found value is the best one or equals to the best known result.

It is important to analyse the time requirement of the HGRE algorithm. The time percentage of the initial route construction is 5.5%, that of the route elimination is 15.5% and the tabu search is the most significant time consuming part of the algorithm with 79%. The route elimination is quite quick, although there are several 'heavy' problems for the algorithm, especially where the best number of routes is less then 4, consequently 25-30 customer nodes have to be inserted.

For further development a possible way to use a simpler cost control algorithm in order to further reduce the computation time. For large problem instances more sophisticated intelligence is needed in order to decide about the application of the time consuming search. At the neighbourhood structure the Global Best strategy has to be replaced by First Best or investigate only a certain part of the neighbourhood to get a reasonable computational time at large problems.

### Acknowledgement

### References

[1]     O. Bräysy, A Reactive Variable Neighbourhood Search for the Vehicle-Routing Problem with Time Windows, Informs Journal on Computing 15 (2003) pp. 347-368

[2]     O. Bräysy, M. Gendreau, Tabu Search Heuristics for the Vehicle Routing Problem with Time Windows, Sociedad de Estadística e Investigación Operativa, Madrid, Spain (December. 2002)

[3]     W. C. Chiang, R. A. Russell, A Reactive Tabu Search Metaheuristic for the Vehicle Routing Problem with Time Windows, Journal on Computing 9 (1997) pp. 417-30

[4]     W. C. Chiang, R. A. Russell, Simulated Annealing Metaheuristic for the Vehicle Routing Problem with Time Windows, Annals of Operation Research 63 (1996) pp. 3-27

[5]     G. Clarke, J. W. Wright, Scheduling of Vehicles from a Central Depot to a Number of Delivery Point, Operation Research 12 (1964) pp. 568-581

[6]     S. Csiszár, Sequential and Parallel Route Construction based on Probability Functions for Vehicle Routing Problem with Time Windows, $5^{th}$ International Conference IN-TECH-ED (2005) pp. 379-388

[7]     Glover F. Tabu Search, Part I. ORSA J. Computing, 1 (1989) pp. 190-206

[8]     Glover F. Tabu Search, Part II. ORSA J. Computing, 2 (1990) pp. 4-32

[9]     J. Homberger, H. Gehring, A Two-Phase Hybrid Metaheuristic for the Vehicle Routing Problem with Time Windows, European Journal of Operation Research, 162 (2005). pp. 220-238

[10]    G. Laporte, The Vehicle Routing Problem: An Overview of Exact and Approximate Algorithms, European Journal of Operation Research 59 (1992) pp. 345-58

[11]    G. Laporte, Y. Nobert, Exact Algorithms for the Vehicle Routing Problem, In S.Martello, G. Laporte, M. Minoux, C. Riberio, editors, Surveys in Combinatorial Optimization, Amsterdam North-Holland (1987) pp. 147-84

[12]    M. M. Solomon, Algorithm for the Vehicle Routing and Scheduling Problem with Time Window Constraints, Operation Research 35 (1987) pp. 254-265

[13]    E. Taillard, Parallel Iterative Search Methods for Vehicle Routing Problems, Networks 23 (1993) pp. 661-72

[14]    S. R. Thangiah, I. H. Osmann T. Sun, Hybrid Genetic Algorithms, Simulated Annealing and Tabu Search Methods for Vehicle Routing Problem with Time Windows, Technical Report UKC/OR94/4, (1995) Institute of Mathematical Statistics, University of Kent, Canterbury, UK

[15]    A. Van Breedam, A parametric Analysis of Heuristics for the Vehicle Routing Problem with Side-Constraints, European Journal of Operation Reports 137 (2002) pp. 348-370

[16]    J. Berger, M. Barkaoui, O. Bräysy, A Route Hybrid Generic Approach for the Vehicle Routing Problem with Time Windows, INFORMS 2003, 41 (2)

[17]    R. Bent, PV. Hentenryck, A Two-Stage Hybrid Local Search for the Vehicle Routing Problem with Time Windows, Transportation Science 2004, 38 pp. 515-530

# Resource-oriented Programming Based on Linear Logic

**Valerie Novitzká, Daniel Mihályi**

Department of Computers and Informatics
Faculty of Electrical Engineering and Informatics
Technical University of Košice
Letná 9, 042 00 Košice, Slovakia
valerie.novitzka@tuke.sk, daniel.mihalyi@tuke.sk

*Abstract: In our research we consider programming as logical reasoning over types. Linear logic with its resource-oriented features yields a proper means for our approach because it enables to consider about resources as in real life: after their use they are exhausted. Computation then can be regarded as proof search. In our paper we present how space and time can be introduced into this logic and we discuss several programming languages based on linear logic.*

*Keywords: linear logic, programming languages, resource-oriented programming*

## 1 Introduction

The aim of our research is to describe solving of large scientific problems by computers as constructive logical reasoning in some logical formal system over type theory. Predicates we consider as predicament-forming functors with nomenclative arguments (symbols), where unary predicates express properties and predicates of higher arity express relations. Using logical reasoning in categorical logic over type theory we can get a mathematical solution of a given problem and the existence of a proof in the intuitionistic linear logic provides us a computable solution of a problem [16, 17].

The role of computer program is to execute instructions under whose the computer system is to operate, to perform some required computations [15]. A running program should provide us a desirable solution of a given problem. We consider programming as a logical reasoning over axiomatized mathematical theories needed for a given solved problem. A program is intuitively understood as data structures and algorithms. Data structures are always typed and operations between them can be regarded as algorithms. Results of computations are obtained

by evaluation of typed terms. Due to discovering a connection between linear logic and type theory - a phenomenon of Curry-Howard correspondence [5], we are able to consider types as propositions and proofs as programs. Then we can consider a program as a logical deduction within linear logical system. Precisely, reduction of terms corresponding to proofs in the intuitionistic linear logic can be regarded as computation of programs. Thus computation of any resource-oriented program is some form of goal-oriented proof search in linear logic. One of the main goals of this approach is to avoid eventual mistakes of correctness generated by implementation of a programming language. This approach also keeps us away from potential problems in verification of programs.

Linear logic (LL) was defined by Girard in [6,7] as resource-oriented logic, where formulae are actions. It enables reasoning similar as in real life, where resources are exhausted after their using.  The basic ideas about LL we present in section 2. There are several programming languages based on LL and Section 2.2 contains the short disscussion about them. The rest of this paper presents the most important features of  LL that enable to introduce resources as space and time into linear logic.

# 2    Linear Logic

LL is regarded as a continuation of the constructivisation that began with both classical and intuitionistic logic and it can serve as a suitable interface between logic and computer science [1]. Whereas classical logic treats sentences with stable values depended on Tarsky semantic tradition (i.e. interested in denotation of sentences), values of linear intuitionistic logic sentences depend on an internal state of a dynamic system according to Heyting semantic tradition (i.e. interested in constructing the proof of a given sentence).

## 2.1    Logical Connectives of LL

LL introduces new logical connectives. *Linear implication A* ⟶o *B* describes that an action *A* is a cause of action *B*; a formula *A* can be regarded as a resource that is exhausted by linear implication so as in real life. It is the most important feature of LL and also of programming based on LL. In clasical logic the truth value of formula *A* after the implication *A* → *B* remains the same. Classical implication can be rewritten in linear manner as *(!A)* ⟶o *B*, where '*!*' is *modal operator* (*exponential*) expressing that we can use a resource *A* repeatebly, as many times as we need. LL defines *multiplicative conjunction* (MC), *A* ⊗ *B,* expressing that both actions *A* and *B* will be done. *Additive conjunction* (AC), *A* & *B,* expresses that only one of these actions will be performed and we can choose which one. *Additive disjunction* (AD), *A* ⊕ *B,* also describes that only one action of *A* and *B*

will be performed but we do not know which one. And finally, *multiplicative disjunction* (MD), *A ℘ B,* expresses: if *A* is not performed then *B* is done, or if *B* is not performed then *A* is done. MD is similar to disjunction in classical logic. *Neutral element* for MC is *1*, for AC is ⊤, for MD is ⊥, and for AD is *0*. *Linear negation,* $A^\perp$, is obtained by analogy with the dual spaces in algebra [9], as it is expressed in (1):

$$A \multimap B = B^\perp \multimap A^\perp \tag{1}$$

Negation is *involutive*, i.e. $(A^\perp)^\perp = A$. An *entailment* in LL is a sequent of the form

$$\Gamma \vdash A \tag{2}$$

where $\Gamma = (A_1, ..., A_n)$ is a finite sequence of linear formulae and *A* is a linear formula deductible from the premises in *Γ*. If *Γ* is empty, then a formula is provable without assumptions. Inference rules of LL have a form

$$\frac{S_1 \ ... \ S_n}{S} \tag{3}$$

where $S_1, ..., S_n$, and *S* are sequents, $S_1, ..., S_n$ are assumptions and *S* is a conclusion of the rule. A proof in this calculus has the form of (proof-) tree, the direction of a proof is from-bottom-to-up, i.e. from the root to the leaves of the tree, where the root is the proved formula and leaves are axioms. In every proof step we can apply the inference rules of LL in Fig. 1 that introduce logical connectives and negation:

$$\frac{\Gamma \vdash A, \Sigma}{\Gamma, A^\perp \vdash \Sigma}(\perp - l) \qquad \frac{\Gamma, A \vdash \Sigma}{\Gamma \vdash A^\perp, \Sigma}(\perp - r)$$

$$\frac{\vdash \Gamma, A \qquad \vdash B, \Delta}{\vdash \Gamma, A \times B}(\otimes)$$

$$\frac{\vdash \Gamma, A}{\vdash \Gamma, A \oplus B}(\oplus - l) \qquad \frac{\vdash \Gamma, B}{\vdash \Gamma, A \oplus B}(\oplus - r)$$

$$\frac{\vdash \Gamma, A \qquad \vdash B, \Delta}{\vdash \Gamma, A \,\&\, B}(\&) \qquad \frac{\vdash \Gamma, A, B}{\vdash \Gamma, A \,⅋\, B}(⅋)$$

Figure 1
Inference rules for LL connectives and negations

In LL are valid the following De Morgan laws:

$$(A \otimes B)^\perp = A^\perp \,℘\, B^\perp \qquad\qquad (A \,℘\, B)^\perp = A^\perp \otimes B^\perp \tag{4}$$

$$(A \,\&\, B)^\perp = A^\perp \oplus B^\perp \qquad\qquad (A \oplus B)^\perp = A^\perp \,\&\, B^\perp \tag{5}$$

Due to meaning of linearity, LL refuses *weakening*, i.e. the constant function *F(x)=a* and *contraction*, i.e. the quadratic function *F(x)=G(x,x)*, but they can be reintroduced as logical rules by using modal operators. Just this restriction gives linear logic it's resource conscious nature. On the other side, due to reintroducing weakening and contraction as logical rules, proof symetry of sequents can be restored again.

## 2.2    Polarity in LL

Logical connectives of LL can be divided into two classes: *positive* and *negative* connectives. Girard [8] regarded positive connectives of LL as algebraic style and negative connectives of LL as logical style. Positive connectives and constants of LL are $\otimes, \oplus, 1, 0$ and negative ones are $\&, \wp, \top, \bot$.

A formula of LL is positive if its outermost logical connective is positive; dually it is negative if its outermost logical connective is negative. A rule of LL calculus is *invertible* if it introduces a negative connective. Negative connectives are introduced by just one rule and the decomposition in proof from bottom-to-up is deterministic. From this it follows a very important consequence: we can get together several proof steps as a *single* step if we have a cluster of negative formulae. Dual property to invertibility is *focalisation* [2], it says that if we have a cluster of positive connectives called *synthetic connective,* we can perform corresponding rules simultaneously. The dual properties of invertibility/focalisation express associativity of logic that is the LL analogue of the Church-Rosser property of λ-calculus.

The *polarity* between positive/negative properties is general in LL. A cluster of rules with the same polarity can be performed as a single rule, by invertibility in the negative case, by focalisation in the positive case. So, the change of polarity in a proof means a new step in this proof and it can be consider as a *time incrementation*. In this manner we can introduce *time* into LL. If we forget truth of formulae and their content and we consider only their locations in proofs, then we can introduce *space* into LL and explicitly handle resources in LL. Due to the important property of invertibility and focalisation linear connectives can be organized by polarities. Summary of significant properties of LL connectives is shown in the Table 1.

## 2.3    Programming Languages Based on Linear Logic

LL forms a base for several functional and logic programming languages. Between functional programming languages [12] based on LL we mention Lilac [13] designed in 1994, but it is not wide-spread language.

| Abbrev. | Symb. | Descr. | Neutral | Polarity | NonDeter. | Human influence |
|---------|-------|--------|---------|----------|-----------|-----------------|
| MC | $\otimes$ | parallel | $1$ | + | | we know how |
| AC | $\&$ | choice | $\top$ | - | internal | we know how |
| AD | $\oplus$ | choice | $0$ | + | external | we don't know how |
| MD | $\wp$ | until | $\bot$ | - | | we know how |

Table 1
Overview of linear connectives and constants

Logic programming may be viewed as the interpretation of logical formulae as programs and proof search as computations. Our view to logic programming is that computation is a kind of a proof search: let us have program $P$ and goal $G$. We are trying to find a proof of the sequent $P \vdash G$. Different goals correspond to different computation sequences. The result of a computation with logic program is a proof that the goal is the logical consequence of the facts and rules. The search strategy is determined by the structure of the goal and the program supplies the context of the proof. The goal is 'active' whereas program is 'passive' and provides a context in which the goal is executed. This principle is called goal-directed provability.

Whereas classical logic programming is based on the first order logic, the resource-oriented programming is based on LL. Resource-oriented program also consists of facts and rules and user-query starts a calculation that it answers the question whether query result belongs to a given program or not. But the difference is only that for performing any action in this logic a clause in a corresponding resource-oriented program must be used exactly once.

There are numerous programming languages that make use of all resource-oriented benefits of LL. In these languages, it is possible to create and consume resources dynamically as logical formulae. Lolli, Lygon, and Forum are implemented as interpreter systems; Lolli [10] is based on SML and Prolog, Lygon [18] is implemented over Prolog, and Forum [14] is based on SML, λProlog and Prolog. Nowadays we have good experiences with Lygon programming language, its implementation can be viewed as extended module over Prolog with full support of LL.

The programming languages ACL and HACL [11] introduce concurrent paradigm in LL programming. Every formula of LL is regarded as a process in some state, the proceses run concurrently with asynchronous communication. Minerva is a commercional language based on Java with own development environment. It enables using LL, but also clasical logic features. The programming language Jinni (Java Inference eNgine and Networked Interactor) enables only a fragment of LL with linear implication, AC and AD and it makes a connection between object-oriented programming and logic programming.

# 3   Resource Handling in LL

In this section we present how it is possible to introduce the resources of space and time into LL. Our approach follows the famous idea published in [9] and it represents a novel approach to proof theory, where proofs are written in locative structure of Gentzen's sequent calculus.

Objects of linear logic are called *designs* and they play the role of proofs, λ-proofs, etc. in usual syntax. Designs are located somewhere. Therefore we need a concept of *location*. A design represents a cut-free proof of a LL formula *A,* in which all information has been erased, only locations in sequents are kept.

Let *A* be a LL formula. Immediate subformulae (w. r. t. focalisation) are denoted by natural numbers called *biases* written as $i,j,k,...$ A finite set of *biases* is called *ramification*, denoted by $I,J,K,...$ An address, *locus,* is a finite sequence $<i_1,...,i_n>$ of biases. We denote addresses by $\sigma,\tau,\upsilon,\xi,...$ A locus denotes a place or spatial location of a formula. If we have a formula *A* and its proof, then *A* is in the root of the proof tree. Formulae in this tree are subformulae and they have precise locations (absolute or relative to the root). But if some subformula occurs in the proof tree more times, then its every occurrence need to receive a distinct location. If locus of *A* is σ and *B* is a subformula of *A* with bias *i* then locus of *B* is $\sigma*<i>$. Let *σ* be a locus. Then $\sigma*\tau$, where symbol *'*'* denotes concatenation, is a sublocus of *σ*. Sublocus $\sigma*\tau$ is called

- *strict* if *τ* is non-empty sequence, $\tau \neq <>$;
- *immediate* if *τ* consists of just one bias, $\tau = <i>$.

We can say that a locus *σ* is in ordering relation *'≤'* with all its sublocuses, i.e. $\sigma \leq \sigma*<i>$.

Because in sequent calculus we proceed in a proof from its conclusion (root of proof tree), a locus occurs before its subloci w. r. t. time relation. Two loci can be comparable and incomparable (disjunct) w.r.t. relation *'≤'*. When two loci are

- *incomparable,* their relation is spatial, i. e. they are completely independent;
- *comparable,* their relation is timeable.

Ramification is needed for multiplicative rules where are two subformulae (assumptions) at the same time. We assume one-sided sequents of the form $\vdash \Gamma$ if all formulae in *Γ* are positive. If a sequent $\vdash A,B,C$ consists of two positive formulae *B,C* and a negative formula *A*, then we replace this one-side sequent by two-side one

$$A^\perp \vdash B,C \tag{6}$$

that consists only of positive formulae. Focalisation enables restriction to sequents with at most one formula on the left side.

Every formula in a sequent has a locus, i. e. a finite sequence of biases. If we forget everything in a sequent except loci of formulae, we get an expression called *pitchfork:*

$$\{\xi\} \vdash \Lambda \tag{7}$$

where $\{\xi\}$ is singleton containing one locus $\{\xi\}$ that can be empty set and $\Lambda$ is a finite set of loci. A locus on the left side is incomparable with every locus in $\Lambda$. A pitchfork $\{\xi\} \vdash \Lambda$ consists of a *handle* $\xi$ and the *tines* (loci) in $\Lambda$. A pitchfork is *positive* if it has no handle and *negative* if it has a handle. A pitchfork is *atomic* if it has a form $\xi \vdash$ , or $\vdash \xi$.

These definitions enable to introduce space and time into LL. In the following example we show how it can be done for a LL formula.

**Example:** Let $A=((P^{\perp} \ \oplus \ Q^{\perp}) \otimes R^{\perp})$ be a LL formula with $P, Q, R$ positive formulae. We can construct the following three proofs of $A$. Proofs 1 and 2 differs only in using left- or right- rules for introducing AD, in the 3 proof we firstly rewrite $A$ using De Morgan rules and then we build its proof.

1.

$$\cfrac{\cfrac{\cfrac{\overline{P \vdash \Gamma}\,(id)}{\vdash P^{\perp},\Gamma}\,(\perp-r)}{\vdash P^{\perp}\oplus Q^{\perp},\Gamma}\,(\oplus-l) \quad \cfrac{\overline{R \vdash \Delta}\,(id)}{\vdash R^{\perp},\Delta}\,(\perp-r)}{\vdash ((P^{\perp}\oplus Q^{\perp})\otimes R^{\perp}),\Gamma,\Delta}\,(\otimes)$$

(-)     (+)

2.

$$\cfrac{\cfrac{\cfrac{\overline{Q \vdash \Gamma}\,(id)}{\vdash Q^{\perp},\Gamma}\,(\perp-r)}{\vdash P^{\perp}\oplus Q^{\perp},\Gamma}\,(\oplus-r) \quad \cfrac{\overline{R \vdash \Delta}\,(id)}{\vdash R^{\perp},\Delta}\,(\perp-r)}{\vdash ((P^{\perp}\oplus Q^{\perp})\otimes R^{\perp}),\Gamma,\Delta}\,(\otimes)$$

(-)     (+)

3. Using De Morgan rules we can write

$$A=(( P^{\perp} \ \oplus \ Q^{\perp}) \otimes R^{\perp}) = (( P \oplus Q)^{\perp} \otimes R^{\perp}) = (( P \oplus Q) \,\wp\, R)^{\perp} .$$

Then the proof tree is:

3.

$$\cfrac{\cfrac{\cfrac{\overline{\vdash P,R,\Lambda}\,(id) \quad \overline{\vdash Q,R,\Lambda}\,(id)}{\vdash (P^{\perp}\&Q^{\perp}),R,\Lambda}\,(\&)}{\vdash ((P^{\perp}\&Q^{\perp})\wp R^{\perp}),\Lambda}\,(\wp)}{A \vdash \Lambda}\,(\perp-l)$$

(+)     (-)

The left-side sign '*(-)*' in proofs 1 and 2 denotes the place where polarity changes from positive to negative. In the proof 3 the left-side sign '*(+)*' denotes the place where polarity changes from negative to positive. We can form clusters of the rules with same polarity for these proof trees and perform them as the following (single) rules that express *time incrementation* in proofs:

1.
$$\frac{P \vdash \Gamma \qquad R \vdash \Delta}{\vdash A, \Gamma, \Delta}(r1)$$

2.
$$\frac{Q \vdash \Gamma \qquad R \vdash \Delta}{\vdash A, \Gamma, \Delta}(r2)$$

3.
$$\frac{\vdash P, R, \Lambda \qquad \vdash Q, R, \Lambda}{A \vdash \Lambda}(r3)$$

Now we forget everything except locations of formulae above. We assume that the locus of *A* is *ξ*, and biases for *P, Q, R* are *3, 4, 7,* respectively. We note that *Γ, Δ, Λ* are no longer formulae but loci. Then we can rewrite previous rules and we get explicit information about space resources needed for proofs. Every sequent in these proofs is a pitchfork:

1.
$$\frac{\xi 3 \vdash \Gamma \qquad \xi 7 \vdash \Delta}{\vdash \xi, \Gamma, \Delta}(\xi r1)$$

2.
$$\frac{\xi 4 \vdash \Gamma \qquad \xi 7 \vdash \Delta}{\vdash \xi, \Gamma, \Delta}(\xi r2)$$

3.
$$\frac{\vdash \xi 3, \xi 7, \Lambda \qquad \vdash \xi 4, \xi 7, \Lambda}{\xi \vdash \Lambda}(\xi r3)$$

□

From this example we can see that we can represent any LL cut-free proof as pitchfork rules. In all last three rules we explicitly treat with space and time. This approach creates the new possibilities in logic programming and enables to see logic as a means how to precisely reason about real life problems and solved them by computers.

**Conclusions**

In our paper we tried to present foundations of the new approach of logic and logic programming. LL enables to handle resources in explicit way, what is the main disadvantage of classical logic used in logical programming. This idea is very recent and it enables to see logic, namely LL, as a logic for reasoning close to real life. In problem solving we are thinking in notions of space and time and there are the most important notions also in programming computers. These concepts form new criteria also for programming languages based on logic and we beliefe that this approach will create a new age of logical reasoning and logic programming.

**Acknowledgement**

**References**

[1]    Abramsky S.: Computational Interpretations of Linear Logic, in Theoretical Computer Science, 1992, pp. 1-53

[2]    Andreoli J. M.: Proposition pour une synthése des paradigmes de la programation logique et de la programmation par objects, Thése d'informatique de l'Université de Paris VI, 1990

[3]    Bierman G. M.: On Intuitionistic Linear Logic, Tech. Rep. No. 346, University of Cambridge, 1994

[4]    Faggian C., Hyland J.: Designs, Disputes and Strategies, In Proccedings of CSL '02, 2002, pp. 1-21

[5]    Girard J. Y.: Proofs and Types, Cambridge University Press, 2003, pp. 1-175

[6]    Girard J. Y.: Linear Logic, Theoretical Computer Science, 50, 1987, pp. 1-102

[7]    Girard J. Y.: Linear Logic: Its Syntax and Semantics, In: J. Y. Girard, Y. Lafont, and L. Regnier, editors, Advances in Linear Logic, Cambridge, 1995, pp. 1-42

[8]    Girard J. Y.: Locus Solum. Mathematical Structures in Computer Science, Vol. 11, 2001, pp. 301-506

[9]    Girard J. Y.: On the Meaning of Logical Rules I: Syntax vs. Semantics. In Computational Logic, U. Berger and H. Schwichtenberg, Eds., NATO ASI Series 165, Vol. 14. Springer-Verlag, New York, 215-272

[10]   Hodas J. S, Miller D.: Logic Programming in a Fragment of Intuitionistic Linear Logic. Information and Computation, Vol. 110, No. 2, 1994, pp. 327-365

[11]   Kobayashi N., Yonezawa A.: Typed Higher-Order Concurrent Linear Logic Programming. Technical Report 94-12, University of Tokyo, 1994

[12]   Kozsik Tamás: Tutorial on Subtype Marks. Z. Horváth (Ed.) Central European Functional Programming School (The First Central European Summer School, CEFP 2005, Budapest, Hungary, July 4-15, 2005), Rev. Selected Lectures. LNCS 4161. Springer-Verlag, 2006, pp. 191-222

[13]   Mackie I.: Lilac – a Functional Programming Language Based on Linear Logic, J. of Functional Programming, Vol. 4, No. 4, 1994, pp. 395-433

[14]   Miller D.: A multiple-Conclusion Specification Logic. Theoretical Computer Science, 165(1):201-232, 1996

[15]   Novitzká V.: Formal Foundations of Correct Programming, Advances in Linear Logic, Elfa, 1999, pp. 1-70

[16]   Novitzká V., Mihályi D., Slodičák V.: Linear Logical Reasoning on Programming, Acta Electrotechnica et Informatica, 6, 3, 2006, pp. 34-39, ISSN 1335-8243

[17]   Novitzká V., Mihályi D., Slodičák V.: How to Combine Church's and Linear Types, ECI'2006, Košice - Herľany, September 20-22, 2006, Košice, 2006, 6, pp. 128-133, ISBN 80-7099-879-2

[18]   Winikoff M. D.: Logic Programming with Linear Logic, PhD. Thesis, Univ. Melbourne, 1997