# COGNITIVE 2010

The Second International Conference on Advanced Cognitive
Technologies and Applications

November 21-26, 2010 - Lisbon, Portugal

**ComputationWorld 2010 Editors**

Ali Beklen, IBM Turkey, Turkey

Jorge Ejarque, Barcelona Supercomputing Center, Spain

Wolfgang Gentzsch, EU Project DEISA, Board of Directors of OGF, Germany

Teemu Kanstren, VTT, Finland

Arne Koschel, Fachhochschule Hannover, Germany

Yong Woo Lee, University of Seoul, Korea

Li Li, Avaya Labs Research - Basking Ridge, USA

Michal Zemlicka, Charles University - Prague, Czech Republic

# COGNITIVE 2010

## Foreword

The Second International Conference on Advanced Cognitive Technologies and Applications [COGNITIVE 2010], held between November 21 and 26 in Lisbon, Portugal, targeted advanced concepts, solutions and applications of artificial intelligence, knowledge processing, agents, as key-players, and autonomy as manifestation of self-organized entities and systems. The advances in applying ontology and semantics concepts, web-oriented agents, ambient intelligence, and coordination between autonomous entities led to different solutions on knowledge discovery, learning, and social solutions.

We take here the opportunity to warmly thank all the members of the COGNITIVE 2010 Technical Program Committee, as well as the numerous reviewers. The creation of such a broad and high quality conference program would not have been possible without their involvement. We also kindly thank all the authors who dedicated much of their time and efforts to contribute to COGNITIVE 2010. We truly believe that, thanks to all these efforts, the final conference program consisted of top quality contributions.

Also, this event could not have been a reality without the support of many individuals, organizations, and sponsors. We are grateful to the members of the COGNITIVE 2010 organizing committee for their help in handling the logistics and for their work to make this professional meeting a success.

We hope that COGNITIVE 2010 was a successful international forum for the exchange of ideas and results between academia and industry and for the promotion of progress in the area of advanced cognitive technologies and applications.

We are convinced that the participants found the event useful and communications very open. We also hope the attendees enjoyed the beautiful surroundings of Lisbon, Portugal.


COGNITIVE 2010 Chairs:

Petre Dini, Concordia University, Canada/ IARIA, USA
Qiang Fu, Broadcom Corporation, USA
Hermann Kaindl, TU-Wien, Austria
JianQiang Li, NEC Labs, China
Sugata Sanyal, Tata Institute of Fundamental Research - Mumbai, India
Qin Xin, Simula Research Laboratory, Norway

# COGNITIVE 2010

## Committee

**COGNITIVE Advisory Chairs**

**Academia**
Hermann Kaindl, TU-Wien, Austria
Petre Dini, Concordia University, Canada/ IARIA, USA
Sugata Sanyal, Tata Institute of Fundamental Research - Mumbai, India

**Industry**
Qin Xin, Simula Research Laboratory, Norway
JianQiang Li, NEC Labs, China
Qiang Fu, Broadcom Corporation, USA

**COGNITIVE 2010 Technical Program Committee**

Taufik Abrão, State University of Londrina (UEL), Brazil
Emmanuel Adam, University Lille Nord de France, France
Rajendra Akerkar, Western Norway Research Institute, Norway
Giner Alor Hernández, Instituto Tecnológico de Orizaba - Veracruz, México
Ateet Bhalla, Technocrats Institute of Technology - Bhopal, India
Cagatay Catal, TUBITAK - Gebze, Turkey
Yaser Chaaban, Leibniz University of Hanover, Germany
François Charpillet, INRIA - Vandoeuvre lès Nancy, France
Po-Hsun Cheng, National Kaohsiung Normal University, Taiwan
Liviu Ciortuz, "Al. I. Cuza" University of Iasi, Romania
Sadie Creese, University of Warwick - Coventry, UK
Darryl N. Davis, University of Hull, UK
Leonardo Dagui de Oliveira , Escola Politécnica of University of São Paulo, Brazil
Lars Fredrik Høimyr Edvardsen, Intelligent Communication AS/Norwegian University of Science and
Technology, Norway
Simon Fong, University of Macau, Macao
Qiang Fu, Broadcom Corporation, USA
Joao Luis Garcia Rosa, University of Sao Paulo (USP) at Sao Carlos, Brazil
Alain Giboin, INRIA Sophia-Antipolis, France
Jürgen Graf, Karlsruher Institut für Technologie (KIT), Germany
Abdelhakim Hafid, University of Montreal, Canada
Ioannis Hatzilygeroudis, University of Patras, Greece (Hellas)
Tzung-Pei Hong, National University of Kaohsiung, Taiwan
Jianmin Jiang, University of Bradford, UK
Hermann Kaindl, TU-Wien, Austria
Kavi Kumar Khedo, University of Mauritius, Mauritius
Abderrafiaa Koukam, Université de Technologie de Belfort Montbéliard, France

Jana Krivec, Jožef Stefan Institute - Ljubljana, Slovenia
Narayanan Kulathuramaiyer, Universiti Malaysia Sarawak, Malaysia
JianQiang Li, NEC Labs, China
René Mandiau, LAMIH- Univ. de Valenciennes, France
Antonio Manzalini, Telecom Italia, Italy
Yiannis Papadopoulos, University of Hull, UK
Mengyu Qiao, New Mexico Institute of Mining and Technology, USA
Alejandro Rodríguez González, Universidad Carlos III de Madrid, Spain
Inès Saad, MIS Laboratory-University of Picardie Jules Vernes, France
Fariba Sadri, Imperial College London, UK
Abdel-Badeeh M. Salem, Ain Shams University, Egypt
Sugata Sanyal, Tata Institute of Fundamental Research - Mumbai, India
Paulo Jorge Sequeira Gonçalves, Instituto Politéchnico de Castelo Branco, Portugal
Sofia Stamou, Patras University, Greece
Vladimir Stantchev, Berlin Institute of Technology, Germany
Kenji Suzuki, The University of Chicago, USA
Antonio J. Tallón-Ballesteros, University of Seville, Spain
Uma Shanker Tiwary, Indian Institute of Information Technology, Allahabad, India
Shirshu Varma, Indian Institute of Information Technology - Allahabad, India
Qin Xin, Simula Research Laboratory, Norway
Kai Xing, University of Science and Technology of China, China
Anastasiya Yurchyshyna, University of Geneva, Switzerland

**Copyright Information**

# Table of Contents

# Gestalt and Computational Perceptual Approach

Brain responses tendencies given by visual and auditory basic stimuli

Bruno Giesteira, João Travassos, Diamantino Freitas

Electrotecnic and Computer Engineering Department
University of Porto
Porto, Portugal
bgiesteira@fba.up.pt, jpctravassos@gmail.com,
dfreitas@fe.up.pt

Diana Tavares
Neurophysiology Course
IPP: ESTSP-CEMAH
Porto, Portugal
tavares.diana@gmail.com

**GUI interfaces require considerable visual attention for their operation excluding the access to important information coded only in the layout. In an Era of mobile devices, we must enhance the auditory designs, to facilitate the interactive contents access to the blind, people with low vision, and/or in any use context. This essay is part of an experimental approach at the human perception based on the theories of form - Gestalt - and the Computational in order to process and implement the brain acquisition signal, obtaining relations between the visual and sound stimuli. We present a computational approach that underlay the electrical signal acquisition of the brain to stimuli response – "Event-Related Potentials" (P300) – based on a fundamental visual syntax that assumes the Gestalt phenomenology with new statistical interim results to the modeling multi-perceptive of information processing (visual and auditory), with the ultimate goal of framing a lexicon and/or basic patterns common that can be applied directly to a well-grounded development of GUI – "Graphic User Interfaces" and AUI – "Auditory User Interfaces".**

*Keywords - Perception; Event-Related Potentials; Gestalt; Computational Theory; GUI; AUI*

## I. INTRODUCTION

Since 2008 the Signals and Systems Laboratory, is leading a new approach in the perceptual field in order to recognize correlations or tendencies between brain responses elicited by two different stimuli modalities, namely visual and auditory [1] whose the main goal is to enhance the interaction multimodality in GUI – "Graphic User Interfaces" and AUI – "Auditory User Interfaces". Statistical data are presented in order to guide future development of auditory icons [2] and "hearcons" [3].

In the essay, first we present the "Research Fundamentals" where it explicit our main motivations and the state of the art regarding the perceptual theories and brain signal acquisition (Event-Related Potentials). Then we explained our laboratorial "Methodology" particularly regarding to the stimuli used and the brain acquisition signal, and then we organized the statistical interim results in the topic "Conclusions and Future Work" that correlate visual and auditory stimuli regarding the velocity of brain

recognition (m/s) as well its energy/resources to process the task. (m/v).

## II. RESEARCH FUNDAMENTALS

### A. Pleas and Motivations

Beyond the neurophysiologic and computational approach we primarily faced a perceptual issue. In terms of perception, we excluded any narrow approach based on only one line/school dogma of theoretical thinking. However, we identify ourselves with the Gestalt phenomenology [4] and Marr's computational theory [5]. Completely different conceptions about visual perception but, in our view, do not render and even complement each other in a Top-Down perspective. Perhaps because the first rests on to descriptive generalizations that make sense and definitely contribute to the understanding and discussion sustained on the phenomenology of visual perception in the XX and XXI centuries [6], but are difficult to reproduce in scientific terms, and the second because it triggers for the first time procedures and scientific methodologies to explain and replicate the way the human mind processes visual stimuli, falling nevertheless in computational reductionism (possibly suitable to the area of Artificial Intelligence) that, putting aside the individuals' phenomenological consciousness negatively conditioned Marr's theory. Nevertheless, revolutionized the way we currently investigate the areas of perception and cognition. In an increasable operative and neurophysiologic perspective [7].

### B. Theory of Form – Gestalt

Our research assumes, contrary to Marr's theory, the subjective nature of the stimuli by the direct influence of the individual conscience [4] [8] e.g., color or even dots, although isolated from a whole context, have subjective phenomenological dimensions inherent to the educational and cultural factors.

As well as the "Feature-Integration Theory of attention" it is assumed that the visual scene is initially encoded in a number of separable dimensions, such as color, orientation, spatial frequency, brightness and motion direction [9]. Any features presented in the same central of "fixation" of attention are combined to form a single object (Gestalt:

"Pragnanz"). This idea was inspired in part by Hubel and Wiesel studies [10] and others who provided evidence about "features" separated in the visual cortex in which each represents a different perceptual dimension such as color, orientation and movement. The organization of the elements is presumably carried out based on factors such as similarity, proximity, contiguity, direction and similarity [11]. These laws appear, however, to operate at a very early stage, presumably before the attention function and before the process responsible for the constancy of the properties of objects such as shape, size, brightness, and so on. It is believed that this is because the mechanisms of attention and constancy presuppose the prior existence of separate competitors' entities or objects about which they operate. The same authors [11] of the essay "Grouping based on phenomenal similarity of achromatic color" suggest that the organization is at an early stage, based on a new principle which they called uniform connectedness. Any features presented in the same central of "fixation" of attention are combined to form a single object (Gestalt: "Pragnanz"), suggesting that regions of uniform stimuli that are interrelated, such as dots, lines or large areas, are interpreted by the perceptual system as a unit. On the essay "Detection Signal Theory – STD" [12] Tanner and Swets, whose main concern involved the measurement of the relations between quality and intensity of a physical stimulus (e.g., light intensity or frequency of a tone) and the perceptual experience caused by that stimulus, attested that perceptual experiences have a continuum of magnitudes that are produced either by noise or by events.

There is thus a remarkable set of neurophysiological evidence indicating that the grouping of objects/stimuli according to Gestalt exists at an early stage, as well as the recognition of certain fundamental characteristics as color, texture, movement, etc., which by their similarity or difference, are distinguished by a perceptual level (joining or separating into different perceptual organization) [11].

### C. Computational Approach

As a neurobiologist and computer scientist, David Marr's works [5] led to a theoretical analysis of vision as a scientific problem proposing vision theories in several areas, including edge detection and perception of depth and shape. The distinction he made between algorithmic / representational / computational and achievement levels of analysis, guided the thinking of vision scientists since then. The levels of computational analysis relate to the objectives and purposes of the system under research. This analysis attempts to characterize, in the abstract, what the system is designed to do.

The algorithmic level is to specify an algorithm or procedure for carrying out the purpose specified by the computational level. Take vision as a computational problem has improved communication between disciplines such as psychophysics, neuroscience and computer science contributing to progress in these areas.

In our laboratory approach, for the first questions of Marr's computational model we found out about the importance of selecting, as visual stimuli to be tested, some of the key elements and basic concepts of visual communication [13] as dot, line, texture, color, scale, depth, movement, not only because the brain processes them differently (which still is the case, particularly for color, motion and depth) but because, in addition, also incorporate a basic visual syntax emphasizing precisely the Gestalt phenomenology that suggests the instinctive demand of the human being in perceive a whole with meaning – "Pragnanz" – in the most consistent, regular and simple way as possible, helping us to obtain correlations and/or trends between visual and sound stimuli conceptually similar and in the same context.

The fact that these are the minimum units perceived in any visual composition that by "Pragnanz" gives it a shape and Uno meaning is extremely relevant in the future correlation between visual and audio settings, constituting the basic units of visual communication capable of structuring more complex image and sound scenarios. Moreover, also being the first stage of the neurophysiologic journey, disparate and orthogonal in the activation of specialized cells of cerebral cortex (e.g., colors; silhouettes; movements; depth) responds to subsequent questions of Marr's computational theory. Namely:

*a) A set of preliminary questions to be asked, e.g., Why is it important to perceive dots, lines and colors?; What is its importance to the individual and his relationship with the world?; Why should the system work to make these visual stimuli explicit?; How can these be represented symbolically in the brain?;*

*b) Developing an algorithm capable of structuring the phenomenon in a neurophysiologic way at the cognitive evoked potentials (P300) level;*

*c) Testing and implementing the efficiency and robustness of the algorithm developed and the understanding of the data acquired at the neurophysiologic and perceptual level.*

### D. ERP-Event-Related Potentials

The recording of ERP – Event-Related Potentials – is a non-invasive electrophysiological investigation method whose goal is to evaluate some of the high level characteristics of information processing in the central nervous system. Each psychological operation in turn involves a temporal activation/inhibition pattern of neurons in a certain brain area. The sum of synchronously generated and event-locked postsynaptic potentials is recorded at the scalp in a form of an ERP component – a potential deflection that is spatially localized and temporally confined [14].

The analysis of the ERPs has been reported as a significant contribute to the knowledge of neural processes that underlie highly specific skills in humans such as language processing, comprehension, visual analysis of faces, processing of emotional stimuli (affective processing) [15] [16], affective picture processing [17], attention, auditory discrimination [18], visual selective attention [19] and mere recognition of stimuli [20].

## III. METHODOLOGY

### A. Basic Stimuli

All visual stimuli used were created to be capable to translate objectively the fundamentals of the visual syntax: The three primary light colors - Red, Green, Blue - plus White, and the basic visual elements and concepts - Dot; Line; Texture; Depth/Dimension; Movement [13]. For the auditory stimuli we selected an audiological grammar used on clinical exams since that is scientifically accepted [21], translating to sound, as far as possible, the basic visual syntax with the sound parameters used to clinical purpose. For the translation of colors we used a fundamental note with 60dB in Si6 (B6) that has a frequency of 1.975 KHz. The Si6 was the note that stayed closer to the 2 KHz, a frequency usually used in clinical context [21]. The musical instrument used to provide the tones varied according to the color we wanted to translate and with the previous volunteers' correspondence. Therefore we choose the sounds of a classic guitar with nylon strings; a piano, a synthesizer; and a glockenspiel, all of them in Si6 (B6 - 1.975 KHz) with 60dB. To translate the visual concepts we used pure tones also with 60dB. A "beep" at 2KHz for the dot, a 2s sound at 2KHz for the line, a 2s tone burst for texture; a 2s sound of 2KHz of frequency varying in intensity progressively from 60dB to 34dB to translate depth/dimension; and a fundamental sound composed by an octave that initiate with the note Si6 of 2KHz to 20 KHz for the movement concept.

Please, consult the "reference" address to full data access, "in press" [22].



Figure 1. Acquisition Signal Scheme

### B. Signal Acquisition

We divided the acquisition into three major parts and each part was preceded by a small interview. In the acquisition process the testing subject was instructed to discriminate a certain randomly appearing stimuli among other different randomly appearing stimuli by clicking on a button in his possession, while his brain activity was being recorded. The stimuli appeared one at a time, and we didn't cross stimuli modalities, i.e., we always separated visual stimuli from audiological stimuli (Figure 1). We always recorded a minimal of two times the brain activity for each stimulus the subject had to discriminate, to demonstrate reproducibility.



Figure 2. Algorithm Chart

All signals were processed and average techniques were applied to it, using "MatLab" software. We used "MatLab" software to develop the algorithm due to the fact that it is a more flexible tool and it allows us a deeper degree of analysis (Figure 2).

## IV. CONCLUSIONS AND FUTURE WORK

Our laboratory study assumes the phenomenology of perception streamlined by Gestalt, founded upon a procedural computational methodology in which it developed an algorithm capable of calculating the Cognitive Evoked Potential (P300) through the acquisition of electroencephalographic signal of some audible and visual stimuli cited that with the same electrodes position, enables us to obtain the response time of the brain against the recognition of a specific stimulus - latency - as well as brain energy resources necessary for this purpose - amplitude. This is only possible because the sensory information from different modalities converge (between 200 and 300 ms after the stimulus) to areas of the cerebral cortex that integrate all information on poly-sensory events, i.e., all the visual, auditory, somatosensitive and olfactory information converge in associative multimodal areas located in the prefrontal, parietotemporal and limbic cortex.

Please, consult the "reference" "in press" [23] address to full data access. The images presented on the previous address represent the average (full line), and the standard deviation (dashed lines) of all acquisitions for each stimulus used. For now, with this approach we can sustain some preliminary results in a Statistical Report.

### A. Statistical Report

The sample that here we presented consisted of 36 individuals, 10 (27.8%) of whom were male and the remaining 26 (72.2%) females. The average age of the sample was 22.75 years with a standard deviation of 5.699. The youngest person was 18 years old and the oldest 36 years. Statistical methods used:

a) *Arithmetic Mean for the associations made by the volunteers as well as for ratings of the stimuli;*

b) *Confidence intervals to "catalog" each stimulus in terms of latency and amplitude;*

c) *Correlations to see which of the pairs have stronger correlations.*

The complete statistical report could be consult in the following "reference" address [24].

## B. *Implications to GUI and AUI*

On the bottom of this essay we present one example of correlation between one pair of stimuli which, like all the report data, will guide our team in future work to develop more sustainable and efficient auditory icons [2] and "hearcons" [3] thereby improving the speed of recognition - latency (m/s) - and the brain resources/energy (m/v) required to interact with a system in a multimodal way.

GUI interfaces require considerable visual attention for their operation. Providing to blind users only the textual contents of the web pages, excluding the access to important information coded in the layout of web pages, the same happened on mobile devices. If interfaces move also to the realm of auditory designs – AUI – these problems are mitigated.

### REFERENCES

[1] S. Handel, "Perceptual Coherence: hearing and seeing", 2006, 1st., Oxford University Press

[2] W. Gaver, "The SonicFinder: An interface that uses auditory icons", Human-Computer Iinteraction, 1989, vol. 4, pp. 67-94, Lawrence Erlbaurn Associates, Inc.

[3] S. Brewster, P. Wright, and A. Edwards, "An evaluation of earcons for use in auditory human-computer interfaces", Interchi'93, 24-29 April 1993, pp. 222-227

[4] W. Kohler, "Gestalt Psychology: The definitive statement of the gestalt theory", 1992, Liveright

[5] D. Marr, "The Philosophy and the Approach", in Steven Yantis - Visual Perception, 2001, pp. 104-123, New York: Psychology Press

[6] W. Metzger, "Laws of Seeing", 2006, MIT Press

[7] S. Zeki, G. Watson, J. Lueck, J. Friston, C. Kennard, and J. Frackowiak, "A Direct Demonstration of Functional Specialization in Human Visual Cortex", in Steven Yantis - Visual Perception, 2001, pp. 193-202, New York: Psychology Press

[8] L. Kaufman and I. Rock, "The Moon Illusion", in Steven Yantis - Visual Perception, 2001, pp. 233-242, New York: Psychology Press

[9] A. Treisman and G. Gelade, "A Feature-Integration Theory of Attention", in Steven Yantis - Visual Perception, 2001, pp. 343-247, New York: Psychology Press

[10] D. Hubel and T. Wiesel, "Receptive Fields and Functional Architecture of Monkey Striate Cortex", in Steven Yantis - Visual Perception, 2001, pp. 147-167, New York: Psychology Press

[11] I. Rock, R. Nijhawan, S. Palmer, and L. Tudor "Grouping Based On Phenomenal Similarity of Achromatic Color", in Steven Yantis - Visual Perception, 2001, pp. 256-265, New York: Psychology Press

[12] P. Tanner and A. Swets, "A Decision-Making Theory of Visual Detection", in Steven Yantis - Visual Perception, 2001, pp. 48-55, New York: Psychology Press

[13] D. Dondis, "La Sintaxis de la Imagen", 1985, Gustavo Gili

[14] J. Kropotov, "Quantitative EEG, Event-Related Potentials and Neurotherapy", 2008, pp. 253-291, Academic Press

[15] M. Eimer, "Event-related brain potential correlates of emotional face processing", 2007, pp. 15-31, Neuropsychologia

[16] J. Petrek, "Pictorial cognitive task resolution and dynamics of event-related potentials", 2008, pp. 223-230, Biomed Pap Med Fac Univ Palacky Olomouc Czech Repub

[17] J. Olofsson, "Affective picture processing: An integrative review of ERP findings", 2008, pp. 247-265, Biol Psychol

[18] J. Duarte, "P300- long-latency auditory evoked potential in normal hearing subjects: simultaneous recording value in Fz and Cz", 2009, pp. 231-236, Braz J Otorhinolaryngol

[19] S. Hillard, "Event-related brain potentials in the study of visual selective attention", 1998, pp. 781-787, Proc. Natl. Acad. Sci.

[20] E. Meijer, "The Contribution of Mere Recognition to the P300 Effect in a Concealed Information Test", 2009, pp. 221-226, Appl Psychophysiol Biofeedback

[21] ASHA, "American Speech-Language-Hearing Association", http://www.asha.org/default.htm (15.08.2010)

[22] www.giesteira.net/Stimuli, B. Giesteira, J. Travassos, D. Tavares and D. Freitas, "Brain's electrical response to visual and auditory stimuli: Relations between the two stimuli modalities", 2010, BMEI 2010, in press (15.08.2010)

[23] http://www.giesteira.net/ERP_Charts/, B. Giesteira, J. Travassos, D. Tavares and D. Freitas, "Brain's electrical response to visual and auditory stimuli: Relations between the two stimuli modalities", 2010, BMEI 2010, in press (15.08.2010)

[24] http://www.giesteira.net/Statistical_Report.pdf (15.08.2010)

Figure 3. Average (full line), and the standard deviation of green color (left) and guitar (right) acquisition: http://www.giesteira.net/ERP_Charts/

TABLE I. P300 CORRELATION VISUAL AND AUDITORY STIMULI. SIGNIFICANCE LEVEL 0,05. PLEASE, CONSULT THE FOLLOWING ADDRESS TO FULL DATA ACCESS (STATISTICAL REPORT): HTTP://WWW.GIESTEIRA.NET/STATISTICAL_REPORT.PDF

| Pairs | FZ | | | | CZ | | | |
|---|---|---|---|---|---|---|---|---|
| | Latency (r) | Rate | Amplitude (r) | Rate | Latency (r) | Rate | Amplitude (r) | Rate |
| Green / Guitar | 0,445 | Moderate | 0,116 | Weak | 0,302 | Moderate | -0,002 | Weak |

TABLE II. "E.G.," AVERAGE & STANDARD DEVIATION OF THE VISUAL AND AUDITORY STIMULUS

| | N | FZ | | CZ | |
|---|---|---|---|---|---|
| | | *Latency* | *Amplitude* | *Latency* | *Amplitude* |
| Green | 34 | ,3348971 (,02670193) | -,0104118 (,01911638) | ,3404926 (,04178523) | -,0147663 (,01957021) |
| Guitar | 34 | ,3235221 (,03180304) | -,0276206 (,01850645) | ,3249632 (,03138339) | -,0280335 (,02009297) |

# Functional Segregation of Semantic Memory and Processing in Ventral Inferior Frontal Gyrus

Mi Li[1, 2], Shengfu Lu[1]*, Xiaofei Xue[1] and Ning Zhong[1,3]

1 The International WIC Institute
Beijing University of Technology
Beijing, China
lusf@bjut.edu.cn, xfx@emails.bjut.edu.cn
2 Liaoning ShiHua University
Liaoning, China
limi135@gmail.com
3 The Department of Life Science and Informatics
Maebashi Institute of Technology
Maebashi-City, Japan
zhong@maebashi-it.ac.jp

*Abstract*—**Although many studies in neuroimaging showed that semantic tasks activated the left ventral inferior frontal gyrus (LvIFG), whether there is the functional segregation of LvIFG in the semantic memory and semantic processing remains unclear. In order to determine neural differences of semantic memory and processing in LvIFG with functional MRI, thirty-six subjects performed reading tasks on triplets of either text or figure or text-figure. The text/figure/figure+text tasks activated two common areas located in an anterior portion of LvIFG and a posterior portion of LvIFG. The BOLD signal change of the posterior portion of LvIFG has semantic working memory characteristics because semantic repetition priming and the BOLD signal change of anterior portion of LvIFG has information processing characteristics. The results suggest that the posterior portion of LvIFG unrelated to information forms is for semantic memory, whereas the anterior portion of LvIFG related to information forms is for semantic processing.**

*Keywords-left ventral inferior frontal gyrus(LvIFG); function MRI; semantic memory; semantic processing; BOLD signal change*

## I. INTRODUCTION

Since the findings of Petersen et al. with Positron emission tomographic (PET) suggested that the left inferior frontal gyrus (LIFG) was identified in processing for semantic association [1], many researchers have focused on the study of the LIFG in neuroimaging. These functional neuroimaging studies have implied that the LIFG is involved in semantic processing, such as semantic judgement [1-6], the control of semantic retrieval [4,7], and the selection of semantic information [8]. Moreover, this region is also related to the processing nonverbal tasks such as object naming [2,9] and unfamiliar faces recognition [10]. Further studies showed that the LIFG was separated into two functional areas, including the dorsal (near the inferior

frontal sulcus involving BA44/45) and ventral parts (BA45/47), and the posterior and dorsal aspect of the left IFG related to phonological processing and the anterior and ventral aspect involved selectively in semantic processing [11-13]. For example, phonetic relative to pitch judgments for auditorily presented syllables activates BA44/6 and BA44/45 [14,15]. Similarly, BA45 is more significantly activated for phonemic orthographic decisions in the visual domain [16]. By contrast, the anterior and ventral LIFG appears to be involved in semantic processing. Kapur et al. [17] demonstrated that semantic-related activity in BA45/47, while others found such activity in the mid-ventrolateral frontal cortex (BA47) [18]. Additionally, most previous studies reported that the LIFG was also activated during the sentence comprehension [19-23]. The cognitive process at the sentence level is not only the semantic processing of words alone, but also refers to syntactic processing [19, 20], semantic working memory [21], and the integration of world knowledge [22,23]. Peter et al using ERP and FMRI demonstrated that the LIFG (BA45/47) is involved in the integration of both word meaning and world knowledge during reading a sentence, and the brain retrieves and integrates them at the same time [22].

Besides the LIFG (BA45/47) is directly related to semantic processing of words, the similar area is decreased activation during repeated semantic processing of those same words (namely semantic priming) [4][24-26]. Gabrieli et al. found that, when making semantic decisions about words, the repeated semantic processing is decreased activation in LIFG relative to initial semantic processing. This decrease in activation represents a semantic repetition priming effect that occurs under implicit test instruction [25]. Further, such repetition-induced decreases in LIFG activation appear specific to semantic processing: Repeated nonsemantic processing of words does not reduce LIFG activation [4]. Another study about the semantic repetition priming examined the stimulus generality of LIFG function during repeated relative to initial semantic processing of

words and of pictures. Their results suggested that the LIFG area (approximately to BA45/47 posteriorly) is decreased activation with repetition regardless of perceptual form [26].

Taken together, the left ventral inferior frontal gyrus (LvIFG) is a crucial area for semantic processing regardless of the verbal and nonverbal stimuli, and this region is also activated in the retrieval from semantic memory that terms semantic memory for short. Semantic processing and semantic memory are two different cognitive processes. It seems that these different processes are subserved by the individual subregions of LvIFG, however, whether the functional segregation of LvIFG in the semantic memory and semantic processing remains unclear. In order to determine neural differences of semantic memory and processing in LvIFG, the experimental materials were designed as the reading tasks on triplets of either text or figure or text-figure, which can describe the same information or content. Figure 1 gives an example of textual and figure tasks used in the experiment. These complex reading tasks involve many cognitive processes such as semantic processing of words, phrases or graph and using world knowledge to construct the whole meaning, and also requires to repeated retrieval of semantic knowledge to comprehend the whole meaning. Based on the above previous studies, we hypothesized that (1) the LvIFG related to semantic memory and processing would be commonly activated by the three present form (text, figure and figure+text) tasks; (2)there would be distinct subregions activated in semantic processing and semantic memory: the subregion of the LvIFG related to semantic processing would be increased activated, which showed the semantic processing characteristics; whereas the subregion of the LvIFG related to semantic memory would be decreased activated, which showed the semantic executive function characteristics.

## II. METHODS

### A. Subjects

Thirty-six volunteers (eighteen female and eighteen male; mean age ± standard deviation (*S.D.*) = 22.5 ± 1.7) participated in this study. All of the subjects were right-handed and native-Chinese speaking. The subjects had no history of neurological or psychiatric illness, and no developmental disorders, including reading disablities. All of the participants gave their written informed consent, and the protocol was approved by the Ethical Committee of Xuanwu Hospital of Capital Medical University and. the institutional Review Board of the Beijing University of Technology.

### B. Materials and Procedure

In the experiment, 20 text, figure and figure+text stimuli, as well as 8 text-baseline, figure-baseline and (figure+text)-baseline stimuli were used. Each text stimulus was presented for a period of 16 seconds, the figure was presented for 14s, and the figure+text was presented for 18s; all the baseline tasks were presented for 8s. The

presentation time was set according to the behavioral experiment, in which participants can fully understand the information of text or figure presented to them. The text, figure and figure+text tasks describing the same event were counterbalanced across subjects; no individual read the same event twice [27].

The experiment consists of 4 sessions. The order of the text, figure and figure+text stimuli was pseudo-randomized in each session. All stimuli were presented on a blank background screen. The participants were instructed to read text, figure or figure+text information attentively. After a stimulus task disappeared, a question including two options was presented, then the subjects could press the selected buttons (left button refers to the first option; right button refers to the second option). The subjects were limited to answer the question during a period of 8s, and then the following rest task was presented for 6s. Four sessions were collected per each participant. The images for the initial 10s were discarded because of unsteady magnetization; the remaining images in the session were used in the analysis.



(a)



(b)

Figure 1 Examples of two types of tasks used in the experiment. (a) A figure task is an example of bar statistical graphs. (b) A text task is a paragraph ranging between 20 and 30 (mean 25) Chinese characters in length (here translated into English).

### C. Image acquisition

In each subject, functional ($T_2$* weighted) images, followed by an anatomical ($T_1$ weighted) image, were acquired with a Siemens 3-T Trio scanner (Trio system; Siemens Magnetom scanner, Erlangen, Germany). Functional images consisted of echo-planar image volumes which were sensitive to BOLD contrast in axial orientation

(TR = 2000 ms, TE = 31ms, flip angle = 90°). Prior to each run, the first two (10 s) discarded volumes were acquired to allow stabilization of magnetization. The volume covered the whole brain with a $64 \times 64$ matrix and 30 slices (voxel size = 4 mm× 4 mm× 4 mm, slice thickness = 4 mm, gap = 0.8 mm).

### D. Data analysis

Functional data was analyzed with statistical parametric mapping (SPM 2, Welcome Trust Centre for Neuroimaging, London, UK) implemented in Matlab 7.0 (Mathworks, Sherborne, MA, USA). The functional images of each participant were corrected for slice timing, and all volumes were spatially realigned to the first volume (the head movement was < 2 millimeters (mm) in all cases). A mean image created from the realigned volumes was coregistered with the structural T1 volume and the structural volumes spatially normalized to the Montreal Neurological Institute (MNI) EPI temple using nonlinear basis functions. Images were resampled into 2-mm cubic voxels and then spatially smoothed with a Gaussian kernel of 8 mm full-width at half-maximum (FWHM). The stimulus onsets of the trials for each condition were convolved with the canonical form of the hemodynamic response function (HRF) as defined in SPM 2. Statistical inferences were drawn on the basis of the general linear model as it is implemented in SPM 2. Linear contrasts were calculated for the comparisons between conditions. The contrast images were then entered into a second level analysis (random effects model) to extend statistical inference about activity differences to the population from which the participants were drawn. Activations are reported for clusters of 10 contiguous voxels (80 mm$^3$) that surpassed a corrected threshold of p < .05 on cluster level. The coordinates given by SPM 2 were corrected to correspond to the atlas of Talairach and Tournoux (1988).

### III. RESULTS

As shown in Figure 2 (a)-(d) and Table 1, we did the conjunction analysis of text, figure and figure+text, and the conjunction between them. All the results consistently showed that two distinct segregated areas were more significantly activated in the left ventral inferior frontal cortex involving an anterior portion (Talairach: -30, 30, -13, BA47/11) and a posterior portion (Talairach: -28, 13, -16, BA47) by text, figure and figure+text. The BOLD signal change percentages at the anterior portion of LvIFG showed the increased activation that has semantic processing characteristics, which suggests that the anterior portion of LvIFG is more related to semantic processing. In contrast, the BOLD signal change percentages at the posterior portion of LvIFG showed the decreased activation that has the semantic executive function characteristics because semantic repetition priming, which suggests that the posterior portion of LvIFG is more involved in semantic memory.

In addition, we also did the conjunction of tasks and rest

(text and rest, figure and rest, figure+text and rest, and text, figure, figure+text and rest), as shown in Figure 2 (e)-(h) and table1. The results showed that only the posterior portion of LvIFG (Talairach: -30, 10, -14, BA47) was more significantly activated during resting state, whereas the anterior portion of LvIFG was not activated. The BOLD signal change percentages at the posterior portion of LvIFG also showed the decreased activation.

Therefore, these results suggest that the semantic processing and semantic memory are dissociated in LvIFG, which means that the anterior portion of LvIFG is more involved in semantic processing, whereas the posterior of LvIFG is more involved in semantic memory.

TABLE I.  BRAIN ACTIVATIONS WITHIN THE LEFT VENTRAL FRONTAL CORTEX WITH CONJUNCTION ANALYSIS

| Region (BA) | Coordinate [a] | | | t | Cluster size (mm³) |
|---|---|---|---|---|---|
| | x | y | z | | |
| *T conj. F conj. FT* | | | | | |
| Posterior LvIFG (47) | -28 | 13 | -16 | 8.87 | 984 |
| Anterior LvIFG (47/11) | -30 | 30 | -13 | 8.28 | 400 |
| *T conj. F conj. FT conj. Rest* | | | | | |
| Posterior LvIFG (47) | -30 | 10 | -14 | 8.35 | 97 |
| | | | | | |
| *Task  conj. Task* | | | | | |
| *T conj. F* | | | | | |
| Posterior LvIFG (47) | -28 | 12 | -16 | 9.11 | 1752 |
| Anterior LvIFG (47/11) | -30 | 30 | -13 | 8.08 | 784 |
| *T conj. FT* | | | | | |
| Posterior LvIFG (47) | -28 | 12 | -16 | 8.84 | 1656 |
| Anterior LvIFG (47/11) | -30 | 30 | -13 | 8.19 | 1046 |
| *F conj. FT* | | | | | |
| Posterior LvIFG (47) | -30 | 10 | -14 | 9.28 | 1752 |
| Anterior LvIFG (47/11) | -30 | 30 | -13 | 8.58 | 800 |
| | | | | | |
| *Task  conj. Rest* | | | | | |
| *T conj. Rest* | | | | | |
| Posterior LvIFG (47) | -30 | 10 | -14 | 8.04 | 552 |
| *F conj. Rest* | | | | | |
| Posterior LvIFG (47) | -30 | 10 | -14 | 10.12 | 1216 |
| *FT conj. Rest* | | | | | |
| Posterior LvIFG (47) | -30 | 10 | -14 | 8.41 | 712 |

BA, Brodmann area; T, text meaning comprehension; F, figure meaning comprehension; FT, figure+text meaning comprehension; LvIFG, left ventral inferior frontal gyrus;

[a] The talairach coordinates of the centroid and associated maximum(Peak) *T* value within contiguous regions are reported.

Figure 2 Regions of significant activation in the left ventral inferior frontal cortex. (a-d) Results with respect to the conjunction analysis between task and task: (a) T conj. F, (b) T conj. FT, and (c) F conj. FT, (d) T conj. F conj. FT. All of (a-d) showed that two distinct segregated areas activated in the left ventral inferior frontal gyrus (LvIFG) involving an anterior portion (BA47/11) and a posterior portion (BA47). (e-h) Results with respect to the conjunction analysis between task and rest: (e) T conj. Rest, (f) F conj. Rest, (g) FT conj. Rest, and (h) T conj F conj. FT conj. Rest. All of (e-h) show that the consistent activation in the posterior portion (BA47) with (a-d) during resting state, whereas the anterior portion (BA47/11) of (a-d) has not activated in resting state. The results have implicated the posterior portion (BA47) related to semantic memory and the anterior portion (BA47/11) is not. The bar graph right shows the BOLD signal change percentages at the activated clusters in the anterior portion and posterior portion of the LvIFG. The statistical parametric map $T$ of all were presented a threshold of 5.05 ($P < 0.05$, corrected for multiple comparisons) and a 400 mm$^3$ cluster size. arLvIFG: anterior portion of LvIFG；prLvIFG: posterior portion of LvIFG.

## IV. DISSCUSSION

The goal of the present study was to examine the neural mechanism of the LvIFG related to semantic processing and semantic memory during reading comprehension of text, figure and figure+text. The three type tasks commonly activated the LvIFG involving an anterior portion and a posterior portion. The increase in BOLD signal of anterior portion showed the semantic processing characteristics, whereas the decrease in BOLD signal of posterior portion showed the semantic memory characteristics. In order to further verify this result, we also do the further analysis about the conjunction between task and rest. If the posterior portion of LvIFG was related to the semantic memory, there would be only this region activated during rest state, and the BOLD signal of this region would be negative. Consistent with this hypothesis, the results of conjunction between task and rest showed that only the posterior portion of LvIFG was decreased activated and the BOLD signal was indeed decreased, whereas the anterior portion was not activated. Thus, the present findings reveal dissociation between semantic processing and semantic memory within the LvIFG.

## A. *Anterior portion of the LvIFG and semantic processing*

This study showed that the activation of the anterior portion of the LvIFG was more related to the semantic processing regardless of text, figure and figure+text. Many previous neuroimaging studies using Positron emission tomographic (PET) and fMRI have consistently demonstrated that the anterior extent of the LIFG, corresponding to BA45/47, plays a crucial role in semantic processing of verbal and nonverbal tasks [1-10]. Petersen et al. (1988) with PET firstly reported that the left inferior frontal gyrus (LIFG) was identified in processing for semantic association [1]. Other functional neuroimaging studies have implicated that the LIFG is involved in semantic processing, such as semantic judgement [1-6], the control of semantic retrieval [4,7], and the selection of semantic information [8]. Moreover, this region is also related to the processing the nonverbal [9, 10]. Further studies that the LIFG was separated into two functional areas, including the dorsal (near the inferior frontal sulcus involving BA44/45) and ventral parts (BA45/47), and the posterior and dorsal aspect of the left IFG related to phonological processing and the anterior and ventral aspect involved selectively in semantic processing [11-13]. Additionally, many prior studies about sentence comprehension were also found the activation in LIPG [19-24]. The cognitive process at the sentence level is not only the semantic processing of words alone, but also refers to syntactic processing [19, 20], semantic working memory [21], and the integration of world knowledge [22,23]. Furthermore, the higher level about the discourse comprehension also found that the LIPG is activated [27,28].

In our study, subjects were instructed to read and comprehend the information presented by the text, figure or figure+text, which involves many cognitive processes such as semantic processing of words, phrases or graph and using world knowledge to construct the whole meaning. Thus, our finding suggests that the activation in the anterior portion of LvIFG contributes to the semantic processing and the integration of semantic and world knowledge.

## B. *Posterior portion of the LvIFG and semantic memory*

In the present study, the posterior portion of the LvIFG (BA47) was decreased activated by the conjunction of task and task, and the conjunction of task and rest, suggesting this region might be more closely related to the semantic memory. This is consistent with previous studies [4][24-26]. Gabrieli et al. found that, when making semantic decisions about words, the repeated semantic processing is decreased activation in LIPG relative to initial semantic processing. This decrease in activation represents a semantic repetition priming effect that occurs under implicit test instruction [25]. Further, such repetition-induced decreases in LIPG activation appear specific to semantic processing: Repeated nonsemantic processing of words does not reduce LIPG activation [4]. Another study about the semantic repetition priming examined the stimulus generality of LIPG function during repeated relative to initial semantic processing of words and of pictures. Their results suggested that the LIPG area (approximately to BA45/47 posteriorly) is decreased activation with repetition regardless of perceptual form [26]. Semantic memory refers to persons' general world knowledge [29,30], involving a wide range of information including facts, concepts and vocabulary [31]. Retrieval from semantic memory occurs during performance many cognitive tasks such as reading and making semantic decision. Other studies about the sentence comprehension have also reported that the LvIFG is activated and this region might be involved in verbal working memory during on-line sentence comprehension [20, 32]. In our study, the complex reading tasks require to repeated retrieval of semantic knowledge to comprehend the whole meaning. Together with the results during rest state, therefore, this study suggests that the decreased activation of the posterior portion of the LvIFG was more involved in semantic memory.

## V. CONCLUSION AND FUTURE WORK

In summary, this study investigated whether the functional segregation of LvIFG in the semantic processing and semantic memory. The present findings indicated that distinct subregions in the LvIFG support the functional segregation of semantic processing and semantic memory. Our results suggest that the anterior portion of the LvIFG is more related to semantic processing, whereas the posterior portion of the LvIFG is more related to semantic memory. We look forward to future studies that further detail these complementary functions and the cooperative work between the anterior and posterior LvIFG play in the reading comprehension.

## REFERENCES

[1] S. E. Petersen, P. T. Fox, M. I. Posner, M. Mintun, and M. E. Raichle, "Positron emission tomographic studies of the cortical anatomy of single-word processing," Nature, vol. 331, pp. 585-589, 1988.

[2] R. Vandenberghe, C. Price, R. Wise, O. Josephs, and R. Frackowiak, "Functional anatomy of a common semantic system for words and pictures," Nature, vol. 383, pp. 254-256, 1996.

[3] X. Wu, J. Lu, K. W. Chen, Z. Y. Long, X. Y. Wang, H. Shu, K. C. Li, Y. J. Liu, and L. Yao, "Multiple neural networks supporting a semantic task: An fMRI study using independent component analysis," Neuroimage, vol. 45, pp. 1347-1358, 2009.

[4] J. B. Demb, J. E. Desmond, A. D. Wagner, C. J. Vaidya, G. H. Glover, and J. Gabrieli, "Semantic encoding and retrieval in the left inferior prefrontal cortex- a functional MRI sutdy of task-difficulty and process specificity," Journal of neuroscience, vol. 15, pp. 5870-5878, 1995.

[5] J. X. Zhang, J. Zhuang, L. F. Ma, W. Yu, D. L. Peng, G. S. Ding, Z. Q. Zhang, and X. C. Weng, "Semantic processing of Chinese in left inferior prefrontal cortex studied with reversible words," Neuroimage, vol. 23, pp. 975-982, 2004.

[6] J. L. Wu, C. Cai, T. Kochiyama, and K. Osaka, "Function segregation in the left inferior frontal gyrus: a listening functional magnetic resonance imaging study," Neuroreport, vol. 18, pp. 127-131, 2007.

[7] A. D. Wagner, E. J. Pare-Blagoev, J. Clark, and R. A. Poldrack, "Recovering meaning: Left prefrontal cortex guides controlled semantic retrieval," Neuron, vol. 31, pp. 329-338, 2001.

[8] S. L. Thompson-Schill, M. D'Esposito, G. K. Aguirre, and M. J. Farah, "Role of left inferior prefrontal cortex in retrieval of semantic knowledge: A reevaluation," Proc. Natl. Acad. Sci. USA, vol. 94, pp. 14792-14797, 1997.

[9] M. Chee, B. Weekes, K. M. Lee, C. S. Soon, A. Schreiber, J. J. Hoon, and M. Chee, "Overlap and dissociation of semantic processing of Chinese characters, English words, and pictures: Evidence from fMRI," Neuroimage, vol. 12, pp. 392-403, 2000.

[10] J. V. Haxby, L. G. Ungerleider, B. Horwitz, J. M. Maisog, S. I. Rapoport, and C. L. Grady, "Face encoding and recognition in the human brain," Proc. Natl. Acad. Sci. U S A, vol. 93, pp. 922-927, 1996.

[11] R. A. Poldrack, A. D. Wagner, M. W. Prull, J. E. Desmond, G. H. Glover, and J. Gabrieli, "Functional specialization for semantic and phonological processing in the left inferior prefrontal cortex," Neuroimage, vol. 10, pp. 15-35, 1999.

[12] J. A. Fiez, "Phonology, semantics, and the role of the left inferior prefrontal cortex," Human brain mapping, vol. 5, pp. 79-83, 1997.

[13] K. B. McDermott, S. E. Petersen, J. M. Watson, and J. G. Ojemann, "A procedure for identifying regions preferentially activated by attention to semantic and phonological relations using functional magnetic resonance imaging," Neuropsychologia, vol. 41, pp. 293-303, 2003.

[14] J. F. Demonet, F. Chollet, S. Ramsay, D. Cardebat, J. L. Nespoulous, R. Wise, A. Rascol, and R. Frackowiak, "The anatomy of phonological and semantic processing in normal subjects," Brain, vol. 115, pp. 1753-1768, 1992.

[15] R. J. Zatorre, A. C. Evans, E. Meyer, and A. Gjedde, "Lateralization of phonetic and pitch discrimination in speech processing ," Science, vol. 256, pp. 846-849, 1992.

[16] J. A. Fiez, M. E. Raichle, F. M. Miezin, S. E. Petersen, P. Tallal, and W. F. Kaztz, "Pet studies of auditory and phonological processing - effects of stimulus characteristics and task demonds," Journal of cognitive neuroscience, vol. 7, pp. 357-375, 1995.

[17] S. Kapur, F. Craik, E. Tulving, A. A. Wilson, S. Houle, and G. M. Brown, "Neuroanatomical correlates of encoding in episodic memory - levels of processing effect," Proc. Natl. Acad. Sci. USA, vol. 91, pp. 2008-2011, 1994.

[18] S. E. Petersen, P. T. Fox, M. I. Posner, M. Mintun, and M. E. Raichle, "Positron emission tomographic studies of the processing of single words," Journal of cognitive neuroscience, vol. 1, pp. 153-170, 1989.

[19] M. Dapretto and S. Y. Bookheimer, "Form and content: Dissociating syntax and semantics in sentence comprehension," Neuron, vol. 24, pp. 427-432, 1999.

[20] Y. Uchiyama, H. Toyoda, M. Honda, H. Yoshida, T. Kochiyama, K. Ebe, and N. Sadato, "Functional segregation of the inferior frontal gyrus for syntactic processes: A functional magnetic-resonance imaging study," Neuroscience research, vol. 61, pp. 309-318, 2008.

[21] D. Caplan and G. S. Waters, "Verbal working memory and sentence comprehension," Behavioral and brain sciences, vol. 22, p. 77-94, 1999.

[22] P. Hagoort, L. Hald, M. Bastiaansen, and K. M. Petersson, "Integration of word meaning and world knowledge in language comprehension," Science, vol. 304, pp. 438-441, 2004.

[23] G. R. Kuperberg, T. Sitnikova and B. M. Lakshmanan, "Neuroanatomical distinctions within the semantic system during sentence comprehension: Evidence from functional magnetic resonance imaging," Neuroimage, vol. 40, pp. 367-388, 2008.

[24] J. Gabrieli, R. A. Poldrack and J. E. Desmond, "The role of left prefrontal cortex in language and memory," Proc. Natl. Acad. Sci. U S A, vol. 95, pp. 906-913, 1998.

[25] J. Gabrieli, J. E. Desmond, J. B. Demb, A. D. Wagner, M. V. Stone, C. J. Vaidya, and G. H. Glover, "Functional magnetic resonance imaging of semantic memory processes in the frontal lobes, " Psychological science, vol. 7, pp. 278-283, 1996

[26] A. D. Wagner, J. E. Desmond, J. B. Demb, G. H. Glover, and J. Gabrieli, "Semantic repetition priming for verbal and pictorial knowledge: A functional MRI study of left inferior prefrontal cortex," Journal of cognitive neuroscience, vol. 9, pp. 714-726, 1997.

[27] M. St George, M. Kutas, A. Martinez, and M. I. Sereno, "Semantic integration in reading: engagement of the right hemisphere during discourse processing," Brain, vol. 122, pp. 1317-1325, 1999.

[28] J. Xu, S. Kemeny, G. Park, C. Frattali, and A. Braun, "Language in context: emergent features of word, sentence, and narrative comprehension," Neuroimage, vol. 25, pp. 1002-1015, 2005.

[29] E. Tulving, "Episodic and semantic memory", In E. Tulving & W. Donaldson (Eds.), Organization of memory, New York: Academic Press, pp. 381- 403, 1972.

[30] E. Tulving, "Elements of Episodic Memory", Oxford: Clarendon Press, 1983.

[31] L. R. Squire, "Memroy and brain", Oxford Eng.: Oxford University Press, 1987.

[32] F. Dick, E. Bates, B. Wulfeck, J. A. Utman, N. Dronkers, and M. A. Gernsbacher, "Language deficits, localization, and grammar: Evidence for a distributive model of language breakdown in aphasic patients and neurologically intact individuals," Psychological review, vol. 108, pp. 759-788, 2001.

# Towards a Cognitive Handoff for the Future Internet:

## Model-driven Methodology and Taxonomy of Scenarios

Francisco A. González-Horta, Rogerio A. Enríquez-Caldera, Juan M. Ramírez-Cortés, Jorge Martínez-Carballido

Department of Electronics, INAOE
Tonantzintla, Puebla, México
{fglez, rogerio, jmram, jmc}@inaoep.mx

Eldamira Buenfil-Alpuche

Faculty of Engineering
Polytechnic University of the Guerrero State, UPEG
Taxco, Guerrero, México
eldamira@gmail.com

*Abstract*— **A cognitive handoff is a multipurpose handoff that achieves many desirable features simultaneously; e.g., seamlessness, autonomy, security, correctness, adaptability, etc. But, the development of cognitive handoffs is a challenging task that has not been properly addressed in the literature. In this paper, we discuss the difficulties of developing cognitive handoffs and propose a new model-driven methodology for their systematic development. The theoretical framework of this methodology is the holistic approach, the functional decomposition method, the model-based design paradigm, and the theory of design as scientific problem-solving. We applied the proposed methodology and obtained the following results: (i) a correspondence between handoff purposes and quantitative environment information, (ii) a novel taxonomy of handoff mobility scenarios, and (iii) an original state-based model representing the functional behavior of the handoff process.**

*Keywords- Cognitive handoff; handoff methodology; handoff scenarios*

## I. INTRODUCTION

A handoff is essential to support the mobility and quality of communications on wireless networks. Its main purpose is to preserve the user communications (continuity of services or seamlessness) while different kinds of transitions occur in the network connection. The resulting handoffs pursuing such purpose are obviously single-purpose handoffs, which we claim they are not enough to face the challenges of the future Internet [1], [2], [3], and [4].

The rationale for this claim is as follows: a seamless handoff provides service continuity, but it is worthless since it works only for the specific scenario to which was stated. Therefore, a handoff should also be adaptive to any possible scenario. Now, a seamless-adaptive handoff is useless if it demands online user interventions. Consequently, a handoff should also be autonomous. Even so, a seamless-adaptive-autonomous handoff is fruitless if new security risks appear during such handoff. Thus, a handoff should also be secure. Furthermore, seamless-adaptive-autonomous-secure handoff is still unproductive if it does not perform correctly, i.e., if it does not maximize the connection time to the best available network and minimize the handoff rate. Such rationale will

lead to a multipurpose handoff: seamless-adaptive-autonomous-secure-correct and thus a valuable handoff.

The development of handoffs achieving multiple desirable features has been "delayed" by the research community itself, despite it was advised since 1997 by Tripathi [1], because many authors preferred to focus on understanding and controlling very specific handoff scenarios (reductionist approach) instead of managing complex and generic handoff scenarios (holistic approach). However, recent handoff schemes, like the ones proposed by Altaf in 2008 [2] for secure-seamless-soft handovers, Cardenas in 2008 [3] for fast-seamless handoffs, and Singhrova in 2009 [4] for seamless-adaptive handoffs, show a tendency towards cognitive handoffs.

This paper presents a model-driven methodology for developing cognitive handoffs. This methodology represents the first attempt to systematically develop cognitive handoffs using a comprehensive model-based framework. The proposed methodology is founded on a synthesis of holism, reductionism, functional decomposition, model-based design, and scientific problem-solving theory.

As a result of deploying our methodology, we present a clear correspondence among cognitive handoff purposes and handoff environment information.

Besides, in order to test the resulting cognitive handoff when applying such methodology with the parameters associated to, and for a given scenario, we develop two things: i) A taxonomy of handoff mobility scenarios which gives a classification of handoff scenarios by considering all feasible combinations of several communication dimensions involved in, and ii) An original state-based model of the handoff process represented by five-state diagram which describes a general control handoff process coordinating the stages before, during, and after the handoff.

The rest of the paper is organized as follows. Section II presents the model-driven methodology we are using for developing cognitive handoffs. This section discusses the difficulties for developing cognitive handoffs and provides an overview of theoretical framework setting the basis of our methodology. Section III shows the first results we obtained from applying the methodology. These results include: (a) the correlation between context data and desirable handoff features through the definition of handoff purposes, objectives, and goals; (b) the taxonomy of handoff

scenarios derived from combining all the possible transition elements involved in handoffs; and, (c) a cognitive handoff state-based model that describes a general behavior of the control handoff process. Section IV presents a basic discussion on the applicability of preliminary results. Finally, Section V concludes the paper with a summary of contributions and future work.

## II. MODEL-DRIVEN METHODOLOGY FOR DEVELOPING COGNITIVE HANDOFFS

### A. Difficulties for Developing Cognitive Handoff

The simple idea of achieving multiple purposes simultaneously is challenging even for humans. Moreover, if the intended purposes represent opposing situations which all of them are desired, then even humans need a way to balance the different purposes in conflict; e.g., the conflict between doing the job accurately and doing it quickly. In optimization theory, multi-objective optimization states that improvements to a single purpose can be made as long as the change that made that purpose better off does not make any other purpose worse off. This is called a Pareto improvement. When no further Pareto improvements can be made, then the solution is called Pareto optimal [5].

Typically, a decision-maker chooses one optimal solution according to his preference. Therefore, the first difficulty in developing cognitive handoffs arises because there are many purposes, objectives, and goals all of them in conflict that need to be tradeoff.

A second significant difficulty emerges when numerous sources of environment information need to be considered to achieve the desired multiple purposes. Six sources of context we consider include: user, terminal, network, provider, application, and handoff process. Such sources produce context data that need to be collected, transformed, and distributed at the different handoff control entities (HCEs). The challenge is how to manage large amounts of unsorted high-dimensional data that have very complicated structures and at the same time reducing the signaling traffic overload produced by this task.

The last significant difficulty is originated by the different transition elements involved in the handoff process. These elements include radio channels, base stations, IP networks, service providers, user terminals, and all the feasible combinations. This variety of elements produces a large amount of scenarios that need to be considered for an adaptive handoff scheme.

### B. Theoretical Framework

First, we state the basis for establishing our methodology.

*1) Holism and Reductionism:* Holism and reductionism are two complementary and opposing approaches for analyzing complex systems [6]. They represent different views of the relationship between the whole and the parts. Holism states that parts cannot explain the whole, the whole states the behavior of parts; i.e., it is necessary to understand how the entire handoff system determines the behavior of its components. Conversely, reductionism states that parts can explain the whole, then the behavior of parts determine the behavior of the whole. We have seen how reductionist handoff schemes achieve its goals in specific scenarios but they quickly become special cases of more general models. Holistic models are more complex models that pretend to consider all the individual parts and to understand the purposes of the whole.

*2) Model-based Design:* The model-driven paradigm has emerged as one of the best ways to confront complex systems. As it was clearly expressed by Dr. Hoffman [7], models can capture both the structure of the system (architecture) and behavior (dynamism). Model-based systems engineering [8] helps to address complexity by raising the level of abstraction, enabling developers to view system models from many perspectives and different levels of detail while ensuring that the system is consistent. The Systems Modeling Language (SysML) [7, 8] is becoming an accepted standard for modeling in the systems engineering domain. Using SysML for modeling helps to reduce ambiguity in models. In fact, models can now show the dynamic behavior of systems, including how they transition between states and how the system behaves overall.

*3) Functional Decomposition*: refers to the process of resolving a functional relationship into its constituent parts in such a way that the original function can be reconstructed from those parts by function composition. The process of decomposition [9] is undertaken for the purpose of gaining insight into the constituent components.

*4) Design as Scientific Problem-Solving:* In his inspiring paper, Braha [10] showed the similitude between the systems design process and the solving-problem process. Therefore, we developed his foundation and proposed a methodology establishing a general procedure that starts with a problem statement and ends up with the solution deployment. This theory views the problem statement as the initial state and then, by searching through a state-space, reaches a goal state representing the solution.

### C. Design and Development Procedure

Steps involved in a form of top-down procedure are:

*1) Stating the problem:* Develop a handoff procedure that can optimally achieve multiple desirable features simultaneously. The handoff procedure should be implemented for operating in real scenarios with multiple dimensions of heterogeneity. Then, as part of the problem: a) Identify and analyze the required system functions: Study the desirable handoff features that need to be implemented and determine the purpose, objectives, and goals associated to every feature. Associate a clear and single purpose to every desirable feature. Decompose each purpose into one or more objectives by identifying the performance parameters that help to quantify the achievement of every purpose. In the same way, divide every objective into one or more specific handoff goals, using optimization values and handoff context data and b) Determine the needed handoff context information: Establish what handoff criteria, handoff metrics, performance measures, handoff policies, handoff constraints, and handoff scenarios are needed to achieve every desired purpose. Study the availability,

locality, dynamicity, structure, and complexity of the variables, policies, and constraints to use.

*2) Design a subsystem structure or model-based framework:* State a cognitive handoff conceptual model, i.e. identify all external context information as well as all internal context information with the highest abstraction level. Whilst internal data constitutes self-awareness, external data constitutes context-awareness of the handoff process. Then, using functional decomposition divide up the conceptual model into a number of sub-models. Every sub-model corresponds to a particular sub-problem that functionally is part of the whole handoff problem. The structure of the system may be represented with a hierarchy of models or framework enclosing the parts of the whole system organized through functional relations. Models in this framework describe the system behavior in an accurate and unambiguous way if one uses a finite set of states and a set of transition functions, thus to ease this part: Identify the associated system states and phases. These dynamic models can be formally represented using finite automata, Petri nets, timed automata, etc. [11]. The states or phases of the handoff process should describe a general behaviour rather than specific details of particular sub-models.

*3) Execute the models:* Execution of models allows verification and validation of such models. This is the difference between just drawing pictures and making pictures "live" as it was pointed out by Hoffmann in [7]. However, verification and validation should not be confused. Model verification means to test if the model satisfies its intended purposes or specifications. Model validation tests if the model provides consistent outcomes that are accurate representations of the real world. We use three strategies for these tasks: simulation, prototyping, and analysis. Whatever the strategy we choose, model testing or model checking [12] requires the use of a formal notation; e.g., modelling languages for simulation, mathematic and logic for analysis, and programming languages or middleware for model prototype implementation. If a model cannot be properly validated or verified, then it must be redesigned within the framework.

*4) Implementation stages:* Once all the models in the framework have been individually tested, the design problem now reflects a well-structured solution. A detailed design can now be generated considering the entire framework of models. This whole system design should be implemented in a whole system prototype. The final prototype is ready to be tested in-situ; should any failure occur during testing, then a review of the conceptual model or any sub-model in the framework should be performed.

*5) Solution deployment:* The cognitive handoff solution is ready to operate on a real handoff environment. The solution system (cognitive handoff) provides a simultaneous acomplishment of the multiple purposes defined by the handoff problem. Each purpose should be associated to quantitative objective functions to measure the degree in which every handoff purpose was achieved.

## III. APPLYING THE MODEL-DRIVEN METHODOLOGY

### A. Purposes, Objectives, Goals, and Context Data

The handoff context information is extensive, heterogeneous, distributed, and dynamic. It supports the whole operation of the handoff process and the achievement of multiple desirable features. From the external and internal vision of the handoff environment, we have identified five external sources of context information (creating context-awareness) and one internal source which is the handoff process itself (creating self-awareness):

*1) User context:* This context includes the user preferences, user priorities, user profiles, and user history and it is used to respond to user needs, habits, and preferences.

*2) Terminal context:* This context domain includes the following evaluating parameters: (i) Link quality: Received Signal Strength (RSS), Signal-to-Noise Ratio (SNR), Signal-to-Noise-and-Interference Ratio (SNIR), Bit Error Rate (BER), Block Error Rate (BLER), Signal-to-Interference Ratio (SIR), Co-Channel Interference (CCI), Carrier-to-Interference Ratio (CIR), etc.; (ii) Power management: Battery Type (BT), Battery Load (BL), Energy-Consumption Rate (ECR), Transmit Power in Current (TPC), Transmit Power in Target (TPT), and Power Budget (PB); (iii) Geographic mobility: Terminal Velocity (Vel), Distance from a Base Station (Dist), Geographic Location (Loc), Moving Direction (MDir), and Geographic Coverage Area (GCA). All these evaluating parameters allow the deployment of QoS-aware handoffs, power-based handoffs, and location-aided handoffs.

*3) Application context:* It includes the QoS requirements of running applications; Lost Packets (LP), Delayed Packets (DP), Corrupted Packets (CP), Duplicated Packets (DuP), Data Transfer Rate (DTR- goodput), Packet Jitter (PJ), Out-of-Order Delivery (OOD), Application Type (AppT).

*4) Network context:* This information is necessary to select among networks (before handoff), to monitor service continuity (during handoff), and to measure network conditions (after handoff) thus they are: Network Bandwidth (NBW), Network Load (NL), Network Delay (ND), Network Jitter (NJ), Network Throughput (NT), Network Maximum Transmission Unit (NMTU).

*5) Provider context*: Information about connection fees, billing models, roaming agreements, coverage area maps, security management (AAA), types of services (data, voice, video), provider preferences, and provider priorities.

*6) Handoff performance context:* This information forms the self-aware part of our cognitive model and allowing evaluation of its performance. Call Blocking (CB), Call Dropping (CD), Handoff Blocking (HOB), Handoff Rate (HOR), Handoff Latency (HOL), Decisions Latency (DLat), Execution Latency (ExLat), Evaluation Latency (EvLat), Handoff Type (HOType), Elapsed Time Since Last Handoff (ETSLH), Interruptions Rate (IR), Interruption Latency (IL), Degradations Rate (DR), Degradations Latency (DL), Degradations Intensity (DI), Utility Function (UF), Signaling Overload (SO), Security Signaling Overload (SSO), Improvement Rate (ImpR), Application

Improvement Rate (AppImpR), User Improvement Rate (UsrImpR), Terminal Improvement Rate (TermImpR), Successful Handoff Rate (SHOR), Imperative Handoff Rate (IHOR), Opportunist Handoff Rate (OHOR), Dwell Time In the Best (DTIB), Authentication Latency (AL), Detected Attacks Rate (DAR), Online User Interventions Rate (OUIR), Tardy Handoff Rate (THOR), and Premature Handoff Rate (PHOR).

Once we have identified the context data from all the context sources and the desired handoff features that we wish to implement, then, we assign a qualitative purpose to every desired feature and, a set of quantitative objectives and goals to every handoff purpose. Tables I and II summarize such previous description.

TABLE I.     DESIRED FEATURES, PURPOSES, OBJECTIVES, AND GOALS

| Desired Handoff Features | Qualitative | Quantitative | |
|---|---|---|---|
| | Purposes | Objectives | Goals |
| Seamlessness | Maintain continuity of services or preserve user communications | Reduce DR, DL, DI, IR, IL | Minimize (BER, BLER, CCI, NL,ND, NJ, LP, DP, CP, DuP, PJ, TPC, TPT, ECR, CB, CD, HOB, HOL) Maximize (RSS, SNR, SNIR, SIR, CIR, NBW, NT, NMTU, DTR, BL, ETSLH) |
| Autonomy | Preserve handoff operation independent of users | Reduce OUIR | Maintain (IL < app.Timeout) |
| Security | Maintain a constant level of security along the handoff | Reduce SSO, DAR | Minimize (AL, SO, HOL) Maintain (High Encryption) |
| Correctness | Keep user always connected to the best network with minimal handoffs | Reduce HOR Increase DTIB | Minimize (HOR) Maximize (DTIB) |
| Adaptability | Keep success of all handoff objectives across any scenario | Multi-objective optimal balance Increase SHOR | Keep every desirable feature within its success range. Maximize (SHOR) |

TABLE II.     OTHER DESIRED PROPERTIES OF COGNITIVE HANDOFFS

| Desired Handoff Features | Qualitative | Quantitative | |
|---|---|---|---|
| | Purposes | Objectives | Goals |
| Necessary | Prevent unnecessary handoffs | Start HO only if it is imperative or opportunist Maint. HOR = IHOR + OHOR | Imperative if (UFcurr<Thinf) Opportunist if (UFcurr>Thsup) UFtarget is SuffB & ConB |
| Selective | Avoid selecting the wrong target | Verify target is consistently better (ConB) and sufficiently better | SuffB: UFtarget > (UFcurr + Δ) ConB: SuffB is maintained for |

| Desired Handoff Features | Qualitative | Quantitative | |
|---|---|---|---|
| | Purposes | Objectives | Goals |
| | | (SuffB) | SP time |
| Efficient | Operate quickly and well-organized to decide how to perform the handoff (HO) | Select the best method, protocol, or strategy according to the HOType, AppType, and Mobility state. Reduce DLat, ExLat, EvLat | Define HO policies or conditions for choosing MIP, SIP, MAHO, NAHO, or other protocols |
| Beneficial | Augment benefits to applications, users, and terminals after handoff | Have a better UF after HO or a maximum improvement rate (UFnew/UFold) | ImpR >> 1 Maximize (AppImpR, UsrImpR, TermImpR) |
| Timely | Initiate a HO not tardy and not prematurely | Reduce THOR and PHOR | Maintain (DLat within its tolerance range) |

These tables represent a relevant preliminary result of the applicability of cognitive handoff methodology. On one hand, they help to reduce the ambiguity and confusion on the usability of similar handoff features because every desirable handoff feature is defined in qualitative terms (purpose) and quantitative terms (objectives and goals). On the other hand, they help to correlate context data with desirable features. For instance, from Table I, we observe that RSS is correlated with seamlessness, IL with autonomy, AL with security, etc. This correlation is intended to select the context data that is needed to support every handoff purpose.

### B. Taxonomy of Handoff Mobility Scenarios

A second significant result obtained from the proposed model-driven methodology is a new taxonomy of handoff mobility scenarios derived from combining all the possible transition elements involved in handoffs; i.e., channels, cells, networks, providers, and terminals. This taxonomy depicts all different kinds of handoffs that are possible in real networks.

Nowadays, no handoff solution exists which comprehensively addresses the entire scale of heterogeneity. Multidimensional heterogeneity [13] is the reason for the large number of handoff scenarios. If we define a handoff scenario as an array $(d_1, d_2 \ldots, d_n)$ where $d_i$ is an instance of $D_i$ the $i$th dimension of heterogeneity and there are $|D_i|$ different ways to instantiate the $i$th dimension, then by the multiplication principle there will be $|D_1|\times|D_2|\times\ldots\times|D_n|$ possible handoff scenarios. However, for the user mobility dimension, the array (location, velocity, direction) may have distinct values at any instant along the path with infinite paths crossing the network; therefore, the number of possible mobility scenarios is infinite. Despite of such infinite scenarios, it is important to make a classification of handoffs according to the elements involved during the transition.

The complexity and treatment for a handoff depend on the type of transition that is occurring. A handoff will

require of services from distinct OSI model layers depending on the elements involved in the transition. For example, a handoff between channels of the same cell is a layer 1 handoff; a handoff between cells (base stations) is a layer 2 handoff, it is homogeneous if cells use the same wireless technology, otherwise is heterogeneous; a handoff between IP networks is a layer 3 handoff; a handoff from one provider to another or between user terminals will demand the services of layers 4-7. Fig. 1 depicts the hierarchical structure of a mobile Internet in a four-layer design (core, distribution, access, and mobile). We will use this figure to explain a handoff hierarchy that involves channels, cells, networks, providers, and terminals.
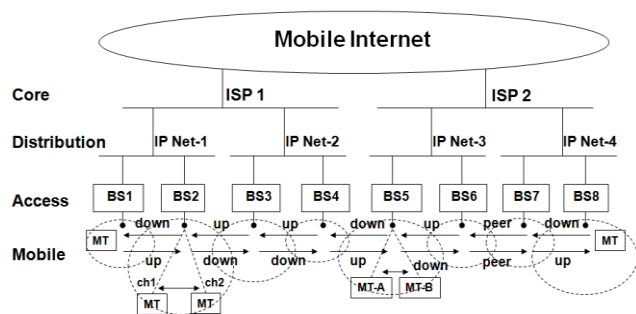


Figure 1.   Hierarchy of handoff mobility scenarios. Different overlay sizes for macro, micro, pico, and femto cells.

The mobile Internet is divided into independent administrative units called Autonomous Systems (AS). An AS is a network administrated by a single organization or person. The Internet is a network of autonomous systems. Fig. 1 depicts two autonomous systems called ISP1 and ISP2 for two distinct service providers. Every ISP uses a very high-speed core network where main servers are located. Providers divide their distribution networks, physically and logically, into a number of IP networks, subnets, or VLANs (Virtual LANs), where the types of services and users are separated. Each IP Net includes a group of base stations (BS) or access points with the same or different wireless access technology. Base stations get distributed across a geographic area to offer mobile communication services. Each base station controls a cell that may have a group of channels to distribute among the associated terminals or a single channel that is shared among several associated terminals.

In Fig. 1, BS2 illustrates a layer 1 handoff when the mobile terminal (MT) changes its connection between channels ch1 and ch2 without changing of BS, IP Net, ISP, or MT. A layer 2 handoff is illustrated between BS1-BS2, BS3-BS4, BS5-BS6, and BS7-BS8. A layer 2 handoff changes from one channel to another and from one base station to another, but keeps the same IP Net, ISP, and MT; however, if the cells involved are heterogeneous, then the handoff is *vertical*, otherwise is *horizontal*. A layer 3 handoff is depicted in BS2-BS3 and BS6-BS7. A layer 3 handoff changes from one channel to another, from one cell to another, and from one IP network to another, but

preserves the same provider and the same terminal; the layer 3 handoff may be heterogeneous, like in BS2-BS3, or homogeneous, like in BS6-BS7. We represent a layer 4-7 handoff, in BS4-BS5, when MT changes its communications from on channel to another, from one cell to another, from one IP Net to another, and from one ISP to another, but the user keeps the same terminal. The encryption schemes and data representation formats change from one provider to another, thus higher layer services are required. Inside the cell for BS5 we depict a handoff between terminals where the user transfers the whole session (current state of running applications) from terminal MT-A to terminal MT-B. Handoffs between terminals can be done for terminals within the same cell or different cells, within the same IP network or different IP networks, within the same provider or different providers. The terminal handoff depicted in BS5 keeps the same cell, same IP Net, and same ISP.

Fig. 2 presents a process diagram that generates the complete taxonomy of handoffs by following the different paths from the upper node to the lower nodes.



Figure 2.   Generation process for handoff taxonomy. There are 15 types of feasible handoffs that can be implemented in real wireless overlay networks. The 1Fh is not a handoff.

Every handoff type in this taxonomy should be complemented or further classified according to many other criteria by using the handoff classification tree of Nasser et al. in [14].

*C.  Cognitive Handoff State-Based Model*

By applying the second step of the model-driven methodology, design a subsystem structure, we created a cognitive handoff conceptual model and its first decomposition model both illustrated and discussed in [13]. Following the reductionist approach, we now focus on a major component of the handoff system, the cognitive handoff control system. At this stage, we designed a state-based model whose purpose is to understand the general behavior that should have the handoff control system. Thus, this model represents our third main result obtained from following the methodology.

Fig. 3 shows a five-state diagram modeling a general control handoff process. The states are: (1) Disconnection, (2) Initiation, (3) Preparation, (4) Execution, and (5) Evaluation. This model describes a generic control handoff system coordinating the stages before, during, and after the handoff. We describe each state briefly:



Figure 3. A handoff control model. This state diagram shows a reactive and deterministic behavior of cognitive handoffs.

*1) Disconnection:* is the initial state and one of the two final states. Here, the terminal is disconnected but discovering available networks. The process will stay here while there are no available networks.

*2) Initiation:* in this state the terminal is connected to the best available network and communications flow normally. This is another final state. The process stays here while there are no reasons (imperative or opportunistic [15]) to prepare for a handoff. If current connection breaks and no other network is available, then the process goes back to the disconnection state.

*3) Preparation:* as soon as a better network appears, the process changes to the preparation state. Here is where properly the handoff begins. This state decides why, where, how, who, and when to trigger the handoff. The handoff in progress can be rolled back to initiation if current link becomes again the best one.

*4) Execution:* once a control entity decides to trigger a handoff, there is no way to rollback; the handoff will be performed. This state knows the current and destination networks, the active application to be affected, and the strategy or method to use.

*5) Evaluation:* once the link switch is made, the control entity enters the evaluation state. This state recombines the measures for every objective function taken before and during the handoff, with new samples taken after the handoff to determine its successfulness. The evaluation latency is adjusted to a stabilization period [16].

## IV. RESULTS DISCUSSION

In this research, we have shown a new methodology to systematically develop cognitive handoffs, which are

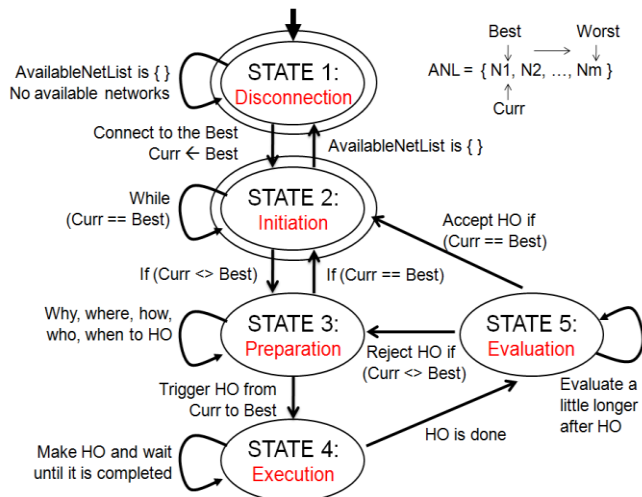expected to be in operation in the mobility scenarios of the future Internet. Such methodology is based on a sound theoretical framework including: methods for analyzing complex systems, the model-based systems engineering, the functional decomposition approach, and the scientific problem-solving theory. There are five stages in the proposed methodology: 1) state the problem, 2) design a model-based framework, 3) execute the models, 4) implement a prototype, and 5) deploy the solution. Thus, we have presented three main results obtained from applying the first two stages of the methodology: i) a cascade relationship of desired features, purposes, objectives, goals, and context data; ii) a taxonomy of handoff mobility scenarios; and iii) a generic state-based model for a cognitive handoff control system.

Furthermore, there are some other issues that require detailed discussion: (a) the complexity of a cognitive handoff system, (b) the evaluation of cognitive handoff models, and (c) the implementation of cognitive handoffs.

### A. Cognitive Handoff Complexity

In [13] we showed two main properties of complex systems that are also present in cognitive handoffs: the hierarchic structure of systems and the property of emergence. In this section we provide other reasons of why cognitive handoffs are complex software systems: (1) Cognitive handoffs exhibit a rich set of behaviors: reactive, proactive, deterministic, non-deterministic, context-aware, self-aware, etc.; behavior is determined by the particular desirable features associated to handoffs. (2) Cognitive handoffs can be stated as multi-objective optimization problems. (3) Cognitive handoffs are driven by events in the physical world; e.g., the user mobility, the user preferences, the provider services, the coverage areas, etc. (4) Cognitive handoffs maintain the integrity of hundreds or thousands of records of information while allowing concurrent updates and queries. (5) Context information is extensive, heterogeneous, dynamic, and distributed. (6) Cognitive handoffs control real-world entities, such as the switching of data flows through a large set of available networks, providers, and terminals. (7) Handoff management has a long-life span; handoffs will exist in all future wireless networks. (8) Handoff management is a key issue for wireless industry and standardization bodies. Grady Booch in [17] provides further discussion on the attributes of complex software systems.

### B. Evaluation of Cognitive Handoff Methodology and Models

Now, as a result of applying our proposed methodology, one gets a set of models that are different in purpose (intentions), usability (applicability), notation (language), and abstraction (hierarchy).

Methodology and each model must be evaluated, either by quantitative evaluation, which comprises the definition of criteria and metrics intended to measure one specific property or, conversely by a qualitative evaluation which is related to credibility that comes from the way in which the

cognitive maps are built and the clarity it represents the opinion's of most experts [18].

In relation to a qualitative evaluation of the methodology, one requires to think on the stages proposed by the development process, the kind of activities to accomplish in each stage, the strength of its theoretical basis, the kind of lifecycle in the development process, etc. Meanwhile, corresponding quantitative evaluation, metrics should be applied to all asociated parametres in the stages of the process.

With respect to evaluate models, we made a clear distinction in Section II.C between verification and validation. The verification tests if the model satisfies its purpose, whilst validation tests if the model outcomes are representations of reality. During the development process of a new system, special purpose models are built to support the understanding that goes on during the development and no hard data emerge from such models, thus, they can only be verified, but not validated.

It is worth to notice that in this paper, we deal with a specific kind of model belonging to those known as soft models [18]. Soft models are intended to understand rather than to predict and therefore verification is the way to qualitatively evaluate such models. Specifically, the theoretical framework in Section IIB has solid and proven bases.

### C. Cognitive Handoff Implementation

We envision the implementation of cognitive handoffs as a network of distributed agents cooperating and competing to take any type of handoff to success. We distinguish between agents for controlling the handoff process (HCEs) and agents for managing the handoff context data (CMAs). The CMAs are responsible for recollecting the context data and updating the handoff information base at the HCEs. CMAs are located in user terminals and distributed in different layers of the network infrastructure. HCEs are located also in every user terminal and at the network access layer; HCEs perform a handoff control process like the one depicted in Fig. 3. Thus, let us develop the state-based model as follows.

A dynamic ordered list of available networks (ANL) is organized from best to worst, according to the value of desirability calculated for every network. The desirability metric is a utility function combining a broad set of network selection criteria. The best network is the one with highest desirability. The value of desirability for the $n$th network, named $D_n(\mathbf{v})$, may have a geometric or stochastic distribution depending on the dynamic nature of context variables used as selection criteria, and arranged in a criteria vector $\mathbf{v} = (V1, V2, …, Vm)$. We use Equation (1) to represent a general mathematical model for the desirability function:

$$D_n(\mathbf{v}) = \sum(K + Wi)\log(Vi^+) - \sum(K + Wj)\log(Vj^-) \qquad (1)$$

The set of decision variables $(V1, V2, …, Vm)$ fetched for the $n$th available network is partitioned in two subsets: $Vi^+$ and $Vj^-$; where $Vi^+$ is the set of criteria that contribute to the desirability (e.g. NBW and NT) and $Vj^-$ is the set of variables that contribute to the undesirability (e.g. NL and

ND). Wi and Wj are weights corresponding to each variable such that Wi and Wj $\in \mathfrak{R}[0,1]$, $\sum Wi = 1 = \sum Wj$ and K is a scaling factor so that small changes in the context variables reflect big changes in $D_n(\mathbf{v})$.

For geometric distributions, a proactive handoff strategy may anticipate handoff decisions and for stochastic distributions a reactive handoff strategy with thresholds, hysteresis margins, and dwell-timers may prevent unnecessary handoffs. The control handoff process illustrated in Fig. 3 shows a reactive and deterministic procedure; reactive, because the process starts the preparation for a handoff until another network with higher desirability is present and, deterministic, because it is always possible to determine the current state of the process within one of five states.

Fig. 4 and Fig. 5 depict geometric distributions of desirability with different handoff strategies. Fig. 4 shows a proactive strategy where the handoff preparation starts before the target network improves the current connection. Fig. 5 shows a reactive strategy where handoff preparation starts after the target network has improved the current connection.
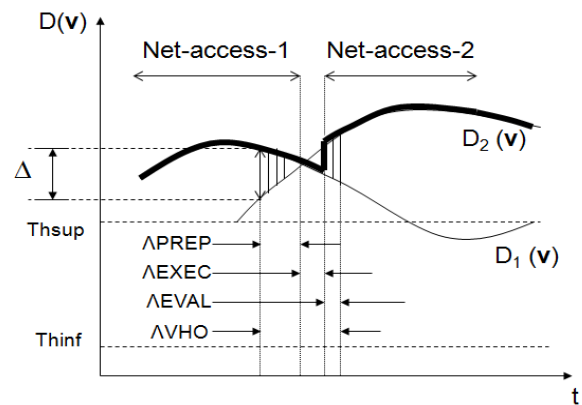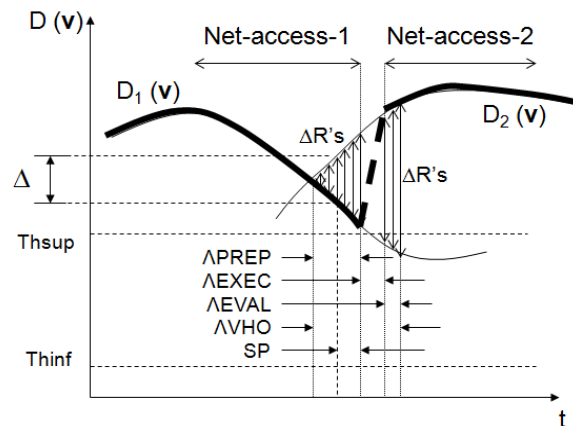


Figure 4.    A proactive handoff strategy.



Figure 5.    A reactive handoff strategy.

The darken line over the desirability functions illustrate the current connection. The performance parameters ΛPREP,

ΛEXEC, ΛEVAL, and ΛVHO depict the latencies for the different stages: preparation, execution, and evaluation. Configuration parameters include Δ (hysteresis margin), desirability threshold (Thsup, Thinf), and dwell-timer (SP). Relative Desirability measures are (ΔRs) which are equal to |Dcurr − Dbest|.

The available network list (ANL) is a data structure located at the HCEs, but continuously updated by the CMAs. When the ANL is empty, the terminal goes to the disconnection state (State 1) and stays there while such list is empty. CMAs are continuously discovering new networks and ordering the list from the highest desirable networks to the lowest desired networks.

The change from disconnection state to initiation state (State 2) occurs as soon as new networks are available. The HCE selects the best available network from the list and connects the terminal to it. The State 2 is the Always Best Connected state because the terminal will stay connected to the best network as long as no other available network improves the current connection.

The change from initiation to preparation (State 3) occurs when a new network is improving or has improved the current network. Handoff decisions, in State 3, start by identifying a reason to begin the preparation for a handoff (why). Next, selecting the target network (where). Then, deciding what strategy, method, or protocol to choose (how). Then, deciding what HCE will be responsible to trigger the handoff (who), and finally, deciding the best moment to trigger the handoff (when). The chosen handoff strategy, method, or protocol depends on the current handoff scenario (as those depicted in Fig. 1) and the type of handoff in progress (as those illustrated in Fig 2).

The decision to trigger a handoff in one terminal changes the control process from preparation to execution (State 4). The trigger handoff decision activates a procedure to change the data flows of an application from one access network to another, within specific handoff and time constraints. The switching mechanism takes a time ΛEXEC to complete.

Once the switching process is completed, the HCE enters to the evaluation state (State 5). This is an important stage of feedback to the handoff control process. At this stage, the HCE has a constrained period of time to decide to accept or reject the recently executed handoff. One condition for handoff success occurs if the new current connection is the best available connection, but others include measuring the objective functions, associated to every handoff purpose, and if all these measures are within a boundary region of acceptable quality, then the cognitive handoff is successful, otherwise it is defective and outliers should be corrected.

## V.   CONCLUSION AND FUTURE WORK

Cognitive handoffs are multipurpose handoffs achieving many desirable features simultaneously; e.g., seamlessness, autonomy, security, correctness, and adaptability. The development of cognitive handoffs is a challenging task that has not been properly addressed in the literature. Therefore, we proposed a new model-driven methodology for developing cognitive handoffs. We applied the proposed

methodology and obtained a clear relationship between handoff purposes and handoff context information, a new taxonomy of handoff scenarios, and an original state-based model of a generic control handoff process.

We continue developing and integrating the models generated by the cognitive handoff methodology. A future work is to organize such models in a comprehensive framework of models representing the functional issues for the whole cognitive handoff process. Further work is needed to study the availability, locality, dynamicity, structure, and complexity of variables, metrics, polices, and constraints involved in cognitive handoffs. The evaluation of the cognitive handoff methodology by quantitative techniques demands more work. We are preparing a manuscript to analyze the cognitive handoff problem as a multi-objective optimization problem using the cellular automata approach to simulate complex handoff scenarios.

## REFERENCES

[1]   N. D. Tripathi, "Generic adaptive handoff algorithms using fuzzy logic and neural networks," Ph.D. dissertation, Virginia Polytechnic Institute and State University, August 21, 1997.

[2]   A. Altaf, F. Iqbal, and M. Y. Javed, "S3H: A secure, seamless and soft handover between WiMAX and 3G networks," *Intl. Conf. on Convergence and Hybrid Information Technology*, pp. 530-534, 2008.

[3]   L. R. Cardenas, M. Boutabia, and H. Afifi, "An infrastructure-based approach for fast and seamless handover," *The Third International Conference on Digital Telecommunications*, pp. 105-109, 2008.

[4]   A. Singhrova and N. Prakash, "Adaptive vertical handoff decision algorithm for wireless heterogeneous networks," *11th IEEE Intl. Conf. on High Performance Computing and Communications*, pp. 476-481, 2009.

[5]   J. Branke, et al. (Eds.), Multi-objective Optimization: Interactive and Evolutionary Approaches, Springer, Germany, 2008.

[6]   V. V. Raman, "Reductionism and holism: two sides of the perception of reality," The Global Spiral, an e-publication of Metanexus Institute, published on July 15, 2005. URL: http://www.metanexus.net/magazine/tabid/68/id/9338/Default.aspx, retrieved on September 4, 2010.

[7]   H. P. Hoffman, C. Sibbald, and J. Chard, "Systems engineering: the foundation for success in complex systems development," IBM Corporation (white paper), Software Group, December 2009, pp. 1-11.

[8]   D. M. Buede, The Engineering Design of Systems: Models and Methods, 2nd Edition, John Wiley & Sons, USA, 2009.

[9]   O. Maimon, and L. Rokach, Decomposition Methodology for Knowledge Discovery and Data Mining: Theory and Applications, Series in Machine Perception Artificial Intelligence, Vol. 61, World Scientific Publishing Co., London, 2005.

[10]   D. Braha, and O. Maimon, "The design process: properties, paradigms, and structure," *IEEE Trans. On System, Man, and Cybernetics – Part A: Systems and Humans*, 27(2), March 1997, pp. 146-166.

[11]   P. A. Fishwick (Ed.), Handbook of Dynamic System Modeling, Chapman & Hall /CRC, USA, 2007.

[12] C. Baier and J.P. Katoen, Principles of Model Checking, The MIT Press, USA, 2008.

[13] F. A. González-Horta, R. A. Enríquez-Caldera, J. M. Ramírez-Cortés, J. Martínez-Carballido, and E. Buenfil-Alpuche, "Towards a cognitive handoff for the future Internet: a holistic vision," *2nd International Conference on Advanced Cognitive Technologies and Applications,* COGNITIVE 2010, Lisbon, Portugal, Nov. 2010.

[14] N. Nasser, A. Hasswa, and H. Hassanein, "Handoffs in fourth generation heterogeneous networks" *IEEE Communications Magazine*, pp. 96-103, October 2006.

[15] W. Zhang, J. Jaehnert, and K. Dolzer, "Design and evaluation of of a handover decision strategy for 4th generation mobile networks," The 57th Semiannual Vehicular Technology Conference, VTC 2003, Jeju, Korea, 2003.

[16] H. J. Wang, R. H. Katz, and J. Giese, "Policy-enabled handoffs across heterogeneous wireless networks," *WMCSA 99*, New Orleans, Louisiana, USA, 1999.

[17] G. Booch, et al. Object-Oriented Analysis and Design with Applications, Third Edition, Addison-Wesley, Chap. 1, 2007.

[18] M. Pidd (Ed.), Systems Modelling: Theory and Practice, John Wiley & Sons, England, pp. 1-42, 2004.

# Agent and Swarm Views of Cognition in Swarm-Array Computing

Blesson Varghese and Gerard McKee
*School of Systems Engineering, University of Reading, Whiteknights Campus*
*Reading, Berkshire, United Kingdom, RG6 6AY*
*Email: b.varghese@student.reading.ac.uk, g.t.mckee@reading.ac.uk*

*Abstract*—The current state of work in 'Swarm-array computing' requires the theoretical concepts proposed in the framework to be formalised. As a preliminary effort to this end, a recently proposed computational intelligence based hierarchical layered architecture for cognitive agents is mapped onto the intelligent agent based approach of swarm-array computing. The cognitive capabilities of two components of the intelligent agent based approach, namely an agent and a swarm, are considered and the components view of perception, reasoning, judging, response and learning are presented. The layered cognitive architecture maps well onto the microscopic or agent level than on the macroscopic or swarm level of the intelligent agent approach.

*Keywords*-cognitive layered architecture, intelligent agents, swarm-array computing.

## I. Introduction

Research pursued in the field of cognitive agent architectures has explored the development and deployment of agents with cognitive capabilities in a computing environment. The interest in cognitive agents has increased rapidly in recent times since cognitive agents are not merely reflexive agents and if employed in parallel computing systems can provide solutions to a wide variety of scientific problems requiring memory, reasoning, and problem solving capabilities.

A recent effort to enhance the fault tolerance of large scale parallel computing systems incorporated under the swarm-array computing framework, referred to as the intelligent agent approach is one such application that can benefit from cognitive agent architectures. The current state of work requires the theoretical concepts proposed in the swarm-array computing framework to be mapped onto cognitive architectures, hence formalising the framework as layers. To this end, a layered cognitive agent architecture is required.

After an extensive survey of literature it was noted that research in cognitive architectures focused on component based architectures with minimal effort towards layered cognitive architectures. However, one recent research reported a hierarchical and layered architecture for cognitive agents [1]. The work reported in this paper is motivated towards mapping the layered architecture for cognitive agents onto the intelligent agent based approach of swarm-array computing.

The remainder of this paper is organised as follows. Section II presents the related work in the area of cognitive agent architectures. Section III considers the swarm-array computing framework, particularly the intelligent agent based approach. Section IV considers the layered cognitive agent architecture. Section V presents the formalised framework of swarm-array computing by mapping the layered cognitive agent architecture onto the intelligent agent based approach. Section VI performs a qualitative evaluation of the cognitive architecture. Section VII concludes this paper with a discussion of and consideration of future work.

## II. Related Work

Among the wide variety of cognitive agent architectures presented in literature, a few relevant architectures, namely the ECLAIR, LIDA and ICARUS which are relatively recent cognitive architectures and ACT-R and Soar which are architectures that have undergone a greater development cycle are briefly reviewed in this section.

ECLAIR, otherwise known as the Engine for Composable Logical Agents with intuitive Reorganization is a recent architecture with emphasis on adaptation and learning [2]. One feature of the ECLAIR model is that it handles unknown inputs to the model by processing them as if they were known to the system. The modules of the architecture enable stimulus or perception of the world, agent awareness, two types of agent behaviours, namely reflex and plan-based, adaptation and decision making.

LIDA, the abbreviation for Learning Intelligent Distribution Agent, is a more recent cognitive agent architecture based on deriving a working model of machine consciousness [3]. An associative memory which provides a perceptual knowledge base, episodic memory for long term storage of autobiographical and semantic memory, functional consciousness that plays the role of a daemon watching for an appropriate condition for acting, procedural memory which is a graph based memory for representing an action and its context and result, and a module for high level action selection of feelings and emotions form the major components of the LIDA architecture. Multiple learning mechanisms is another feature of the LIDA architecture.

ICARUS another recent cognitive agent architecture is based on conceptual inference and skill execution, which are two approaches of handling knowledge in the architecture [4]. By conceptual inference an agent understands its state and situation by inferring from percepts and beliefs while by skill execution an agent achieves goals by decomposing them

into ordered subgoals. Other features of the architecture include goal selection, means-end problem solving and skill learning.

ACT-R is a cognitive agent architecture that aims to model human behaviour. The architecture comprises six sensory modules and each responsible for vision based processing, executing actions, achieving goals, long-term declarative knowledge, relational declarative knowledge and short-term memory [5]. The model is also capable of learning and updating its knowledge.

Soar is a goal oriented cognitive agent based architecture and represents long-term knowledge in the form of production rules, episodic memory and semantic memory [6]. Multiple learning mechanisms are implemented in the Soar architecture. For example, procedural long term knowledge is acquired through reinforcement learning while declarative knowledge is acquired by episodic and semantic learning.

The agent architectures presented above do not follow a layered architecture. Due to the modular or component based structure of sub systems in the above architectures, the modelling of communication and interaction between the agents tend to be taxing when compared to agent interaction and communication that could be modelled in a layered approach. Further, layered approaches enable the incorporation of additional sandwich layers for extending existing architectures by minimal modifications of the supporting layers.

Clearly there is a need for agent architectures to be developed such that they are layered. Research in the direction of developing layered cognitive architectures are sparse. However, one recent research reported in [1] has proposed a hierarchical five layered architecture for cognitive agents. The layers of the architecture are based on the sequence of activities that contribute to the cognitive capabilities of an agent.

In the next sections, we explore how the layered architecture for cognitive agents considered above can be mapped onto an approach in swarm-array computing, namely intelligent agents. The purpose of mapping the layered architecture onto the swarm-array computing approach is a part of the effort made towards formalising the theoretical concepts of the swarm-array computing framework, which is considered in the next section.

### III. SWARM-ARRAY COMPUTING FRAMEWORK

Research in swarm-array computing has progressed in the direction of applying autonomic computing concepts to large-scale distributed parallel computing systems, thereby improving the fault tolerance of parallel computing systems. The framework for parallel computing deals with constituents, namely the computing platform, the problem/task to be executed, the landscape and the swarm. Moreover, three approaches that bind these constituents, namely the



Figure 1. Illustration of the Intelligent Agent Approach in Swarm-Array Computing

intelligent agent based, intelligent core based and a hybrid approach are proposed.

In this section, the current state of work of the intelligent agent based approach of swarm-array computing is presented. Parallel reduction [7], a class of algorithms which are of importance and employed in a variety of applications in the high performance computing domain has been considered in the intelligent agent based approach. These algorithms are based on tree structures and figure 1, left, shows an example. The data flows from the leaves of the tree towards the root and at each intermediate node the converging data input is transformed into a result that is passed forward to the next intermediate node. The interconnection of a node in the tree represents its dependencies and the complexity of communication and coordination between the nodes also increases with its dependencies.

In the intelligent agent based approach, when a parallel reduction algorithm needs to be run on a high performance computing platform each node of the tree is scheduled onto a separate computing node. Since these computing nodes are susceptible to failures, there is a need to deal with the isolation of faults. Traditional methods such as checkpointing are challenged by drawbacks and reduce the efficiency of high performance computing systems [8][9].

However, the efficiency can be improved if the algorithm is self-managing such that if a node is about to fail the component of the algorithm can be moved off the node and the input and output dependencies re-established on another node. To incorporate this level of intelligence in the algorithm it would be appropriate to implement agent-like intelligence whereby a computing node can be monitored and a component moved if a failure is anticipated.

In the intelligent agent based approach, the parallel components of the reduction algorithm are mapped onto agents, such that the algorithm in effect is a payload to the set of agents. Figure 1, middle, shows the swarm with its payload. The set of agents are intelligent due to a few cognitive capabilities that they possess. Further, the set of agents

which carry the payload onto the computing nodes can be viewed as a robot swarm and the array of computing nodes can be viewed as a landscape. The robot swarm can then move over this terrain to find a suitable area to execute the algorithm that is mapped onto them. Moreover, if one of the nodes on which the swarm is located fails, a local adjustment can be made by the swarm agent relocating to a nearby part of the landscape and re-instantiating its dependencies, hence offers the potential to improve fault tolerance by minimizing human intervention as in traditional fault isolation methods, and therefore increase the efficiency of high performance computing systems.

These concepts have been investigated practically through both a simulation and implementation [10][11]. The implementation employed a computer cluster with thirty three compute nodes and one head node. A Message Passing Interface (MPI) [12] implementation, namely Open Message passing interface (OpenMPI) [13] was used as the middleware for the implementations. A parallel summation algorithm with fifteen nodes was implemented using both the classical approach and the swarm-array computing approach. The implementation of the approach followed a layered structure such that an abstraction layer was implemented over the actual hardware layer. The agents traversed on the abstraction layer as shown in figure 1, right. The results obtained from the implementations proved that the swarm-array computing approach improved fault tolerance as measured by the mean time taken for reinstatement of the algorithm if a node failed.

The current state of work of the intelligent agent based approach in swarm-array computing is as described above. The approach meets the aims for which it has been proposed yet lacks the formalisation of an agent architecture. Hence, a direction for progressing research would be to investigate an appropriate agent architecture that can be mapped onto the intelligent agents in swarm-array computing.

The requirements for the agents in the intelligent agent approach of swarm-array computing need to be considered before looking at cognitive agent architectures. There are atleast four requirements that agents in the intelligent agent approach need to meet so as to demonstrate cognitive capabilities, namely perception, reasoning, judging, responding and learning.

Firstly, an agent needs to be aware of its environment which includes both the computing cores on which it can carry a task onto and other agents in its vicinity and agents with which it interacts or shares information. Secondly, an agent needs to situate itself on a computing core that may not fail soon and can provide necessary and sufficient consistency in executing the task. Thirdly, an agent needs to predict core failures by consistent monitoring (for example, heat dissipation of the cores can be used to predict failures). Fourthly, an agent needs to be capable of shifting gracefully from one computing core to another, without causing interruption to the state of execution, and notifying other interacting agents in the system when a core on which a sub-task being executed is predicted to fail.

## IV. LAYERED COGNITIVE AGENT ARCHITECTURE

The need to formalise the existing body of work in swarm-array computing requires the mapping of a generic cognitive agent architecture onto the approaches in swarm-array computing. The layered structure followed in the swarm-array computing approaches including the intelligent agent based approach of interest in this paper hence requires a layered cognitive architecture so that the approach can be formalised.

Recently a Computational Intelligence based Architecture for Cognitive Agents, referred to as CIACA in this paper, has been proposed and reported in [1]. Though the architecture cites examples of activities from highway traffic modelling, the architecture aims to formalise the sequence of activities, namely perceiving, reasoning, judging, responding and learning an agent needs to be capable of performing so that it is cognitive. One notable feature of the architecture is that it follows a hierarchical and layered structure, therefore of interest in the context of swarm-array computing. A perceptual layer, a reasoning layer, a judging layer, a response layer and a learning layer are the five layers constituting the architecture and is shown in figure 2.

The perceptual layer aims to perceive information from the environment by sensing and from information provided by other interacting agents [1]. For example, in a traffic highway model, if a car on the highway is assumed to be an agent and the highway its environment, then information such as other cars ahead and behind and their approximate distances from the agent is acquired information by perception.

The reasoning layer supports coherent and logical thinking by obtaining information from the perceptual layer and processing the information using existing knowledge [1]. Fuzzy logic is suggested as a useful tool for reasoning which is a component of the natural thinking process performed unconsciously by humans. For example, humans can successfully park a car though many approximations are made.

The judging layer receives information from the reasoning layer and may send it back to the reasoning layer after processing or refining information to make decisions at the reasoning layer [1]. For example, if a vision system is used for parking a car between two parked vehicles, then information irrelevant for successful parking obtained through the vision system needs to be eliminated. Feature extraction through edge and corner detection algorithms may prove useful.

The response layer instructs response commands to the perceptual layer after applying rules to the information obtained from the judging layer [1]. For example, move ahead, accelerate or decelerate and switch lane are response commands. Additional reasoning and judging capabilities

| CIACA Architecture | Traffic Model | Intelligent Agent Approach - Swarm-Array Computing | |
|---|---|---|---|
| | | Agent View | Swarm View |
| Learning | driving on a new lane system | update core dependencies | update knowledge of whats happening on the landscape |
| Responding | Response commands - move ahead, accelerate, decelerate, switch lane | move to a core | transform pattern shape |
| Judging | Eliminate information not necessary for parking | decision-making for which core to move onto | decision-making for which area can the swarm move onto |
| Reasoning | Parking a car by approximations | the core an agent can move onto if a core fails | the area a swarm can move onto if a set of cores fail |
| Perceiving | Cars ahead and behind | agents and cores in vicinity, will a core fail | identify an area for execution |

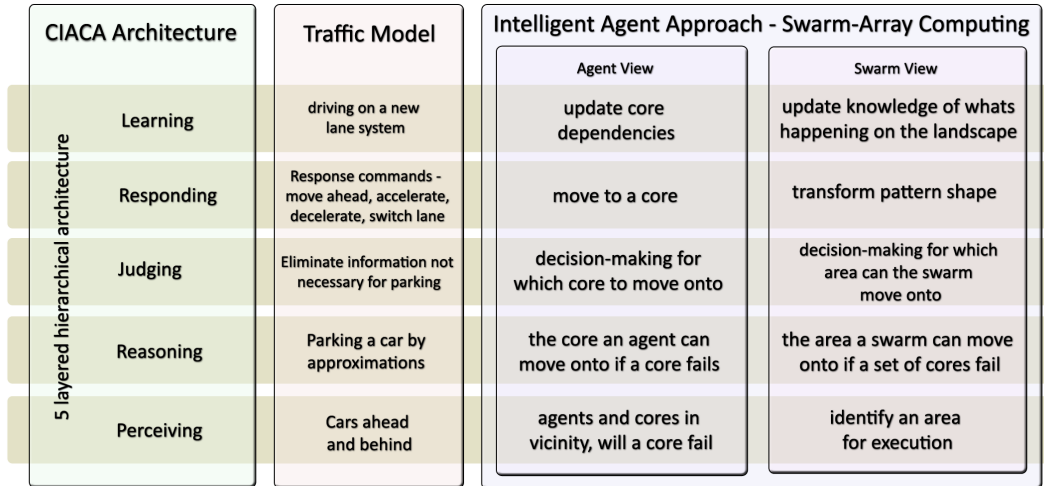(left vertical label: 5 layered hierarchical architecture)

Figure 2.   Illustration of the CIACA architecture and the traffic model, and the Agent and Swarm views in the Intelligent Agent Approach of Swarm-Array Computing

can also be implemented in this layer for exhibiting higher level of intelligence.

The learning layer modifies existing knowledge by using information gathered by the agent. The ability to learn is implemented within algorithms such that an agent is able to derive new knowledge. More sophisticated learning such as extended learning by different generation of agents leading towards evolution of knowledge can also be implemented in this layer.

The above layers meet the requirements for achieving an almost fully functional cognitive architecture. Further, the layers sufficiently formalise the sequence of activities of an agent in a hierarchical layered sequence such that additional computational strategies can be implemented in these layers. Moreover, the layered approach can meet an agent's requirement for demonstrating cognitive capabilities. Therefore, the concepts of the CIACA architecture are chosen to be mapped onto the intelligent agent based approach of swarm-array computing approach.

## V.  Mapping the Layered Architecture for Intelligent Agent Approach

An efficient mapping of a cognitive architecture onto the intelligent agent approach can be obtained only by considering different views or perspectives of the components, namely an agent, the swarm of agents, a landscape, multiple landscapes, that contribute to the approach. The perception, reasoning, judging, response and learning capabilities of an agent and the swarm will be explored in this section. The landscape and multiple landscape view, more applicable for the intelligent core based approach of swarm-array computing, will be reported elsewhere.

### A.  Agent View

In the first instance an agent's view is considered. In the intelligent agent approach, the environment in which an agent is situated comprises agents with which it can interact and computing resources. Perception in this context would mean to acquire information concerning the environment. An agent needs to answer questions such as 'are there other agents in my vicinity?' and 'which computing cores are functional in my vicinity?' To achieve this, an agent can probe the environment, i.e., by sending 'are you alive' signals to the agent and the computing resources. Perception for an agent also includes gathering information for answering the question will the core that I am situated on fail? Sensory information for predicting such failures can be gathered through sensors that consistently probe the hardware core. For example, the rise in temperature of a core can be sensed to predict a computing core failure.

Reasoning becomes necessary once an agent predicts the computing core it is situated on to fail. An agent needs to answer questions such as which cores in the computing environment would it be possible to move onto? Since an agent has options to move onto a few cores in its vicinity, an agent needs to make an appropriate choice. Hence, the agent needs to also think will the core that I will move onto fail? For this an agent should have perceived sensory information of the cores in its vicinity.

Judging for an agent is necessary for decision making. For example, an agent may think about which core do I move to, but a decision has to be made concerning the core to which an agent can move. As suggested above, the sensory information perceived by an agent aids decision making. However, in the context of the intelligent agent based approach in swarm-array computing, the judging layer

necessarily need not be implemented as a separate layer; it may complement the reasoning layer.

After an agent makes a decision as to which core it can move onto, a response needs to be initiated. If the response is instructed explicitly through an external controller then the agent's cognitive capability is challenged. On the other hand, a response initiated by the agent itself is appropriate in the context of achieving intelligence by demonstrating cognitive capabilities. An instruction such as 'move to' or 'move to core x' initiated so that an agents move onto a core other than which it is situated on is an example.

A sophisticated agent also requires mechanisms whereby it learns about its environment from the perceived sensory information and uses it for decision making. For example, in the context of the intelligent agent algorithm, an agent updates its information on the cores it is dependent on. The core dependencies known to the agent and the knowledge gained from 'are you alive' signals contribute to the knowledge of an agent about the landscape.

From an agent view, the CIACA architecture maps well onto the intelligent agent based approach and is shown in figure 2. This is so since the mapping on the agent level provides a microscopic view of the intelligent agent approach. The perception layer provides functionalities for acquiring information from the agent's environment. The reasoning layer enables an agent to think logically while the judging layer assists an agent in narrowing down an agent's choice. The response layer enables an agent to initiate a response while the learning layer updates existing knowledge of an agent.

With the current state of work in the intelligent agent approach it is to be noted that reasoning and judging need not be implemented as separate layers since decision making does not involve many choices. The CIACA architecture does not address issues such as security, resource management and providing generic services. However, since the intelligent agent approach implemented in parallel reduction algorithms considered low-level aspects, such issues did not have to be considered.

### B. Swarm View

Having considered the capabilities of an agent, namely perceiving, reasoning, judging, responding and learning, above, there is also a need to consider how a group of agents or a swarm comprising agents can demonstrate intelligence in the swarm-array computing approach. Hence, the second view considered in this section is the swarm's view. The cognitive capabilities of the swarm are emergent since individual agents contribute to the swarm's behaviour. In other words, an agent demonstrates intelligence on a microscopic or individual behaviour level, whereas the swarm demonstrates its intelligence on a macroscopic or abstract behaviour level. For example, an agent might need to move onto many cores

while executing a task and update its dependencies resulting in the displacement of the swarm on the landscape.

The emergent behaviour of the swarm due to the perception of individual agents is the capability to identify an area comprising a set of cores on the computing space where the swarm can situate itself to execute a task. In other words, the swarm perceives which area it can situate to execute a task and whether the set of cores can provide sufficient resources and access to resources for successful execution of the task mapped onto the swarm agents.

Reasoning on the swarm level includes questions such as 'where can I move onto if a set of cores in the computing space fail?' and 'will the new set of cores that I move onto fail?' The swarm's decision making or judging capabilities are demonstrated by the definitive and unequivocal decision it makes when a question such as 'which area on the computing space do I move to' arises after reasoning when a set of cores has been predicted to fail.

Response at the swarm level includes the capability of a swarm to displace itself from one location to another in the computing space if a set of core is predicted to fail. The displacement can occur by transforming the shape of the pattern formed by the agents whereby individual agents reposition to other cores such that the task mapped onto an agent is seamlessly executed.

The knowledge of what is happening on the computing landscape aids decision making when the swarm has to move about on the landscape. This knowledge is learnt by the swarm from sensory information that is perceived and local interactions such as sending and receiving 'are you alive' signals by the agents comprising the swarm. On an implementation level, for the sake of convenience, knowledge can be acquired, maintained and updated centrally. However, a decentralized strategy for acquiring, maintaining and updating the knowledge-base is closer to the swarm concept.

In general, though the CIACA architecture mapped onto the swarm view, as shown in figure 2, all layers of the architecture did not prove useful in the mapping. Perception of the swarm is an emergent behaviour and therefore is of less importance in this context. However, reasoning, judging and response can be seen on both microscopic (agent) and macroscopic (swarm) levels. Learning involves representation and storage of knowledge, which needs to be decentralized to be in similar lines of a swarm, and therefore operates on the microscopic level.

The mapping of the CIACA architecture on the swarm view is unlike the mapping of the CIACA architecture on the agent view for two reasons. Firstly, the level of abstraction for the agent and swarm view is different, since in the swarm view, macroscopic properties providing an abstract view is considered. However, in the agent view microscopic properties were considered. The CIACA architecture maps well on the microscopic level. Secondly, the interactions

in the swarm level are more complex since they comprise both inter-agent and agent-environment interactions of all agents comprising the swarm. Hence, on an abstract level, the CIACA architecture proves less effective as seen in the swarm view than on the microscopic level as seen in the agent view.

## VI. Evaluating the CIACA Architecture

This section evaluates the mapping of the CIACA architecture onto the swarm-array computing approach considered in previous sections. The set of six evaluation criteria presented in [14] are used to perform the evaluation. The evaluation criteria are: (a) Generality, versatility, and taskability, (b) Rationality and optimality, (c) Efficiency and scalability, (d) Reactivity and persistence, (e) Improvability, (f) Autonomy and extended operation. This set of criteria is a general set of principles relevant to cognitive agent architectures and are broad in its scope of evaluation and hence adopted for the qualitative evaluation of the architecture.

*1) Generality, Versatility and Taskability:* Generality evaluates how well the architecture can support intelligent behaviour in a broad range of environments. The CIACA architecture is proposed as a general cognitive agent architecture and illustrated for traffic models. The applicability of the CIACA architecture for the intelligent agent based approach of swarm-array computing is illustrated in this paper. Though the CIACA architecture has been recently proposed yet has illustrated two applications on which the architecture can be mapped onto. To exemplify and evaluate the generality of the CIACA architecture more applications need to be investigated so that the CIACA architecture can be mapped onto.

Versatility evaluates how taxing is the process of constructing intelligent systems across a given set of tasks and environments. The architecture maps well on the microscopic level, i.e., agent view in swarm-array computing, but does not map well on the macroscopic level, i.e., swarm view in swarm-array computing. Therefore, the CIACA architecture is not necessarily versatile on an abstract view.

Taskability evaluates how an agent can carry out tasks not only by knowledge it has acquired but also by explicit communication with humans or other agents. In the intelligent agent based approach an agent's response to move off from a failing core is not only based on the knowledge it has of its environment but also from commands it may obtain as signals from other agents in its vicinity. The aim of the intelligent agent based approach is to create self-managing systems to execute a task by minimising human administrator intervention; therefore the approach does not consider receiving explicit commands from humans. However, in traffic models, humans need to make explicit commands to an agent representing a car to reach a goal.

*2) Rationality and Optimality:* Rationality evaluates how an agent's knowledge and action will lead towards its goal or in other words the relationship between an agent's goal, its knowledge and actions. In the agent view of the intelligent agent based approach, the primary goal of an agent is to execute a sub-task that is mapped onto it. To achieve this goal an agent may have to relocate on different computing cores. In this context, the degree of rationality of an agent will be based on how well an agent utilises its knowledge of the computing environment to execute a task.

Optimality evaluates whether an agent's selected behaviour yields an optimal solution. The degree of optimality will be high if an agent can successfully complete its task by being rational. In the case of parallel summation algorithms, every agent on whom the task of summation is mapped receives information from and yields information to other agents in the environment. The states that an agent can enter into, thereby showing different behaviours, is limited in this case as against agent behaviours that demonstrate different behaviours as shown in traffic models.

*3) Efficiency and Scalability:* Efficiency evaluates the amount of time and space required by the computing system. The experimental studies on the cluster proved that the time for reinstating the execution of an algorithm once a failure occurred was significantly reduced when compared to the time taken by other existing traditional approaches. Clearly the efficiency of the approach increases with the adoption of cognitive agent architectures.

Scalability evaluates the architecture's performance in varying conditions including task difficulty, environment uncertainty and time of operation. The intelligent agent approach was implemented on a computer cluster focusing on space applications and simulated uncertainty in the environment which was sensed as a hazard by the agent. Scalability studies on other experimental platforms have not yet been explored for the approach. Moreover, the swarm-array computing approach has been implemented for parallel reduction algorithms, an important class of algorithms in parallel computing, but has not yet moved towards implementations for more complex algorithms.

*4) Reactivity and Persistence:* Reactivity evaluates how well an agent can respond to unexpected situations or events. The unexpected situation considered in the intelligent agent approach is a failure of the computing core. The failure is anticipated by an agent and the agent responds to this situation by making a decision to which core it must relocate. It is noted from experimental results that mean times taken for reinstating the execution of an algorithm if a core fails is in the order of milliseconds, and therefore confirms that agents respond and react quickly in the swarm-array computing framework.

Persistence evaluates how the architecture pursues its goals despite changes in the environment. An agent is not only affected by the failure of a computing core in its environment, but also by another agent in its vicinity. If an agent situated on a core predicted to fail is dependent on one

or more agents in the environment, then dependency information needs to be circulated such that agents can continue to pursue goals despite changes in their environment.

*5) Improvability & Autonomy and Extended Operation:* Improvability evaluates the ability of an agent to perform a task with addition of knowledge when compared to the state it did not possess knowledge. Clearly, in the intelligent agent approach considered in this paper, an agent makes its decision as to which core it should move onto in the case of a predicted failure is based on its knowledge of which cores in its vicinity are not likely to fail. If the agent did not possess knowledge of its environment, an agent would make a decision that would not be necessarily optimal, i.e., moving off to a core that is likely to failure thereby requiring a further relocation at the expense of time and slowing the execution of the task.

Autonomy evaluates the personal independence of an agent. The degree of autonomy in the architecture can be evaluated based on the cognitive ability of the agents seen in the approach rather than merely being reflexive agents.

Extended operation evaluates whether an agent can operate on its own for prolonged periods of time. To start off the intelligent approach was proposed to isolate faults when single nodes failed and continue seamless execution of a task for a prolonged period. Additional work will be required to enable the approach to handle multiple node failures, thereby extending the operation of agents for prolonged time frames in more realistic scenarios.

## VII. CONCLUSIONS

The work reported in this paper has aimed to formalise the intelligent agent based approach in swarm-array computing by mapping a layered cognitive architecture, namely the CIACA architecture, onto the intelligent agent approach. Primarily, the conceptual aspects of such a mapping has been presented in this paper. An agent view and swarm view of perception, reasoning, judging, response and learning in the swarm-array computing framework has been presented. A comparative evaluation of the mapping using the cognitive agent architecture against a set of general criteria is performed.

Future work will aim to formalise the intelligent core based swarm-array computing approach using the cognitive layered architecture. Immediate efforts will be also made to consider the intelligent core based approach for exploring the landscape and multiple landscape views.

## REFERENCES

[1] A. T. Lawniczak and B. N. Di Stefano, "Computational Intelligence Based Architecture for Cognitive Agents", in the Proceedings of he International Conference on Computational Science, 2010.

[2] A. L. Buczak, K. Greene, D. G. Cooper, M. Czajowski, M. O. Hofmann, "A Cognitive Agent Architecture Optimized for Adaptivity" in the Proceedings of the Conference on Artificial Neural Networks in Engineering, 2005.

[3] U. Ramamurthy, B. Baars, S. K. D'Mello and S. Franklin, "LIDA: A Working Model of Cognition" in the Proceedings of the 7th International Conference on Cognitive Modeling, 2006, pp 244–249.

[4] P. Langley, K. Cummings, and D. Shapiro, "Hierarchical Skills and Cognitive Architecture" in the Proceedings of the 26th Annual Conference of the Cognitive Science Society, 2004, pp. 779–784.

[5] J. R. Anderson, D. Bothell, M. D. Byrne, S. Douglass, C. Lebiere and Y. Qin, "An Integrated Theory of the Mind", Psychological Review, 111, (4). 2004, pp. 1036–1060.

[6] J. E. Laird, "Extending the Soar Cognitive Architecture" in the Proceedings of the Artificial General Intelligence Conference, 2008, pp. 224–235.

[7] J. Dongarra, I. Foster, G. C. Fox, W. Gropp, K. Kennedy, L. Torczon and A. White, "The Sourcebook of Parallel Computing", Morgan Kaufmann Publishers, 2003.

[8] J. P. Walters and V. Chaudhary, "Replication-Based Fault Tolerance for MPI Applications" in the IEEE Transactions on Parallel and Distributed Systems, Vol. 20, No. 7, July 2009, pp. 997–1010.

[9] X. Yang, Y. Du, P. Wang, H. Fu and J. Jia, "FTPA: Supporting Fault-Tolerant Parallel Computing through Parallel Recomputing" in the IEEE Transactions on Parallel and Distributed Systems, Vol. 20, Issue 10, October 2009, pp. 1471–1486.

[10] B. Varghese and G. T. McKee, "Can Space Applications Benefit from Intelligent Agents?" in the Proceedings of the 3rd International ICST Conference on Autonomic Computing and Communication Systems, Limassol, Cyprus, 2009 and in the Proceedings of AUTONOMICS 2009, LNICST 23, 2009, pp. 197–207.

[11] B. Varghese, G.T. McKee and V. N. Alexandrov, "A Cluster-Based Implementation of a Fault-Tolerant Parallel Reduction Algorithm Using Swarm-Array Computing" in the Proceedings of the 6th International Conference on Autonomic and Autonomous Systems, Cancun, Mexico, 2010, pp. 30–36.

[12] W. Gropp, E. Lusk and A. Skjullum, "Using MPI-2: Advanced Features of the Message Passing Interface", MIT Press, 1999.

[13] E. Gabriel, G. E. Fagg, G. Bosilca, T. Angskun, J. Dongarra, J. M. Squyres, V. Sahay, P. Kambadur, B. Barrett, A. Lumsdaine, R. H. Castain, D. J. Daniel, R. L. Graham, T. S. Woodall, "Open MPI: Goals, Concept, and Design of a Next Generation MPI Implementation" in the Proceedings of the 11th European PVM/MPI Users' Group Meeting, Budapest, Hungary, 2004, pp. 97–104.

[14] P. Langley, J. E. Lair and S. Rogers, "Cognitive Architectures: Research Issues and Challenges" in Cognitive Systems Research, Volume 10, Issue 2, 2009, pp. 141–160.

# The Significance of Verbal and Spatial Attentional Resources on Mental Workload and Performance

Abdulrahman Basahel, Mark Young and Marco Ajovalasit

Human-Centered Design Institute, School of Engineering and Design

Brunel University

Uxbridge, Middlesex UB8 3PH, UK

Abdulrahman.Basahel@brunel.ac.uk          M.young@brunel.ac.uk          Marco.ajovalasit@brunel.ac.uk

*Abstract*— For decades, scientists have been studying the impact of task workload on individual performance. The purpose of this study was to examine and validate the difficulty levels of two visual tasks (verbal and spatial) to be used in a later experiment studying the interaction of physical and mental workload on attentional resources and performance. Additionally, the study was conducted to determine if a significance difference exists between how men and women perform these types of tasks. The verbal and spatial task workloads satisfied the difficulties levels.

*Keywords-mental; workload; attentional resources; verbal; spatial*

## I. INTRODUCTION

Workload has been identified as one of the main variables that impacts human performance [1][2]. Task workload can be divided into two types, i.e., mental demands and physical demands [2][3]. In fact, mental workload has increased more than physical workload in many jobs due to the rapid increase in technology in recent years [4]. However, most tasks still include physical and mental demands on the operators [1][ 2][5][ 6].

The mental workload concept does not have a conventional definition; however, mental workload can be defined as the resources (i.e., human capacity and skills) that are needed to complete the demanded tasks [2][7]. On the other hand, the concept of physical demand can be defined as the demand associated with tasks that require physical work from the operators, thereby utilizing the musculoskeletal system, the cardiorespiratory system, and the nervous system of the human body [8].

The description of the interference of multi-task demands in terms of shrinkage of some energy for information processing is the goal of resource theory [9]. The attentional resources model includes three orthogonal components. These components are: processing stage (perception, working memory, and response), processing codes (spatial and verbal), and processing modalities (visual and auditory) [9].

The purpose of this study was to examine and validate the difficulty levels of two tasks visually resource (verbal and spatial) that will be used in a subsequent series of studies. Also, it aimed to identify whether there is a difference between male and female performance in both tasks. The gender factor is important since the second part of

the future research programme will be the investigation of the effect of verbal and spatial through auditory resources on attentional resources and performance for both males and females. Therefore, the gender difference is necessary for both studies (visual and auditory tasks) to determine if any significant differences occur in completing the tasks before this study is implemented in the real domain. In addition, the last part of this programme of research will be conducted in the industrial field and will include tasks performed by male operators.

## II. MENTAL AND PHYSICAL WORKLOAD AGAINST PERFORMANCE

In terms of the effects of workload, reference [10] reported that workload can affect and reduce the ability level of the user. In addition, an unexpected rise in the task demand level may lead to a decrease in performance accuracy and an increase in response time of the operator and the operating system [10]. In addition, most researchers agree that the optimum levels of mental load produce acceptable job performance and response; in other words, the job demands should not be too low or too high [1][2][11][12]. Moreover, mental workload includes two major parts; these parts are "stress (task demand) and strain (the resulting impact upon individual)" [7]. Furthermore, a balance must be reached between the physical load of the assigned tasks and user's physical functional capacity in order to produce acceptable performance and reduce injuries and errors [13]. Accordingly, the interaction between the impacts from physical and mental workload on performance is not uniform.

## III. VERBAL AND SPATIAL TASKS AGAINST GENDER DIFFERENCES

Variations in the ability of males and females to perform verbal or spatial tasks have been reported by researchers in several studies [14][15][16]. Research indicates that gender could affect operator performance in tasks that require spatial ability, such as the mental rotation test [15][16]. Researchers report that males do better than females in spatial ability tasks. According to [14], the differences between males and females in their ability to perform mathematical tasks are small and decreased in severity over the course of a year-long study. Furthermore, the difference in the ability of males and females to perform well in the spatial ability task

is dependent upon the type of spatial task [17]. Other researchers reported that women participants received a lower score than men when performing several tests of spatial ability, including Money's Road Map Test, the Geometric Forms test, and the Mental Rotation Test [17]. However, in some types of cognitive tasks, females have faster response rates than males [17]. Therefore, differences of males and females in ability to perform cognitive tasks or tests are related with the type of the job.

## IV. WORKLOAD MEASUREMENTS TECHNIQUES

Various techniques are used to measure mental workload, including performance measures, subjective assessment tools e.g., "NASA-TLX, SWAT" and physiological incidents [2][4][7][18]. Furthermore, subjective tools have been frequently implemented for measuring mental and physical workload in separate situations [18]. Although all of these techniques have advantages and disadvantages and have been widely examined in numerous studies [2][7], many researchers believe that physiological indicators, such as heart rate (HR), heart rate variability (HRV), blood pressure and eye blink, are more sensitive to mental demand and thus more reliable [2][4]. Reference [19], found that when the HRV of individuals was high the appropriate responses increased and errors declined, whereas when the HRV was low the inappropriate responses increased. Furthermore, they believed that there is a relationship between HRV and mental task performance.

## V. METHODS

### A. Design

This experiment was conducted to evaluate and validate the impact of three levels of mental workload for two tasks upon attentional resource performance of operators: an arithmetic task (verbal) and a spatial figures task (spatial); the experiment was a full factorial repeated measures design. Each study was implemented separately with a separate group of participants, so the counterbalancing between the types of tasks was not necessary in this research study. The arithmetic mental task included the following three different levels:

1- The low level involves addition/subtraction problems with numbers between 1 and 10.
2- The intermediate level involves addition/subtraction problems with two numbers between 3 and 35 for the subtraction operation and two numbers between 6 and 35 for the addition operation.
3- In the difficult level, participants are asked to complete high level addition/subtraction problems with two numbers between 20 and 150 for the subtraction operation and between 20 and 105 for the addition operation.

The second task, i.e., the spatial figures task, also included the following three various levels:

1- For the low level, participants are asked to compare three figures with an original figure.
2- In the intermediate level, participants are asked to select two identical figures from six figures.
3- In the difficult level, participants are presented with nine figures and are asked to choose two identical figures from the group.

These types of tasks were chosen to represent the effect of typical workload levels upon the attentional resources performance of users [20]. The arithmetic task was used on three different levels to emulate the demand of a verbal task, according to some [21], who used the mathematical mental task as a verbal task that placed a load on attentional resource capacity and information processing. In addition, the arithmetic mental task is considered to be a verbal task that places a load on working memory capacity since individuals memorize the numbers as words in short-term memory [17]. Also, according to reference [17], mental rotation is considered to be a spatial task that relies on spatial resources; numerous studies have measured the spatial reasoning abilities of individuals. Moreover, several studies have employed mental rotation tasks (the figures were published by [22], in order to evaluate the load on the spatial ability resources of individuals (see, [16][23][ 24]).

Overall, this research included two independent variables: the types of tasks (an arithmetic task and a spatial figures task), and the difficulty of each task. Furthermore, it contained three different dependent variables: performance (number of correct responses and time of correct responses); physiological indices (obtained by measuring the heart rate and heart rate variability); and subjective assessments of mental workload (observed by using the NASA-TLX scores) [25].

### B. Materials

All performance trials were conducted using an Acer-compatible PC with a Pentium II 300 processor operating at 266 MHz and a Tangent 17-inch monitor; it involved the MathsNet Mental Tests 1.5 Software [26], "see Figure 1", and the Mental Rotation Test Software [27], "see Figure 2". All participants were comfortable with and clearly understood how to present the task on the PC screen.
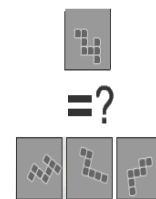


Figure 1. Screenshot of the mental math software.

Figure 2. Screenshot of the spatial figures program.

The testing was conducted using Web-based software for two reasons: first, the programs met the experimental requirements in terms of recording the correct responses and

trial time. Secondly, there is other software (i.e., E-prime) that can do these same tasks, but it is expensive, whereas the software chosen is available online and is easy to use. Also, Polar CS600 chest-electrode was used in order to measure heart rate (HR) and heart rate variability (HRV). Polar ProTrainer 5 software (5.35.164) was used to analyze the heart rate and heart rate variability. In addition, the NASA-TLX [25] was used to evaluate the mental workload of each task.

### C. Participants

Twelve participants (ages 25–35) were selected randomly from Brunel University. This sample size included 6 males and 6 females who were chosen in order to find one standard deviation for the independent variables as well as normality. Also, the same sample was utilized for both studies. The task procedures were explained to all participants. Participants were invited by an announcement that was issued through the university Web site. In addition, the experiments met the ethical rules of the School of Engineering and Design at Brunel University.

### D. Procedures

At the beginning, participants were given a brief introduction about the experiment in order to familiarize them with the steps. Also, the participants were provided instructions and advice on how to perform an arithmetic mental task and a spatial figures task. Then, the participants were asked to affix the chest electrodes for the heart rate monitor on their chest such that we could record the HR and the HRV for each participant as they completed the assigned tasks. In addition, the height, weight, age, and gender of each participant was recorded and used to set-up the heart rate monitor tool.

Subsequently, the first experiment began with the presentation of the arithmetic mental tasks. In addition, the participants were presented the levels of difficulty randomly in order to reduce the potential carryover effects and fatigue. Each volunteer completed 25 questions within each level as accurately and quickly as possible in the allotted 5 minutes. The number of correct responses and the actual time required to complete the correct responses and the section were recorded directly by the software. The HR and HRV were recorded at rest level and continuously throughout the completion of each condition using chest electrodes made by Polar. Also, immediately after completing each trial, the participants were asked to complete the NASA-TLX scale in the 2 to 3 minute interim between each level. After completing the first experiment, i.e., three levels of arithmetic problems, the subjects were given 5 minutes to complete the NASA-TLX and rest.

Then, the second experiment (i.e., the spatial figures test) began. The participants were asked to continue wearing the chest electrodes for the HR monitor such that we could continue measuring HR and HRV. The Mental Rotation program generated different figures with various angles of rotation at three different levels (i.e., low, intermediate, and high). The program also recorded the number of correct

choices and the time required to complete the task. Each condition level included 25 problems, and participants were given 5 minutes to complete each level. In addition, they took 2 to 3 minutes to rest and complete the NASA-TLX between each condition.

## VI. DATA ANALYSIS

Analysis of variance (ANOVA) was used to investigate the impact of the independent factors (i.e., the arithmetic and spatial figures mental tasks) on the dependent variables (i.e., performance, mental workload, and physiological arousal). Also, repeated contrast comparisons were used to determine whether or not homogeneous levels of arithmetic tasks were significantly different from that of the spatial figures tasks. Furthermore, independent t-tests were implemented in order to examine the significance of the mean differences of each type of task and their interaction. A 95 % confidence level (i.e., $\alpha = 0.05$) was applied in these studies.

## VII. RESULTS

### A. Performance

The descriptive statistics for the participants are illustrated in Table I. The participants' performance was measured by recording the number of correct responses for the arithmetic tasks and the spatial figures tasks (i.e., the mental rotation test). In addition, the responses were related to the task workload levels for each task arithmetic and spatial figures tasks; "see Figure 3" and "see Figure 4", respectively.

TABLE I.  DESCRIPTIVE STATISTICS FOR SAMPLE SIZE

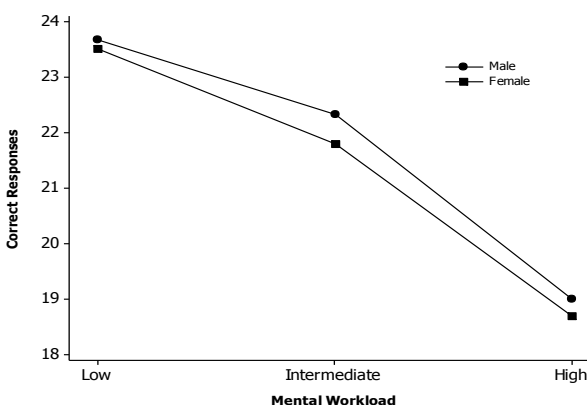| Variable | Male(n=6) | | Female (n=6) | |
|---|---|---|---|---|
| | *Mean* | *SD* | *Mean* | *SD* |
| Age | 28 | 2.7 | 27.8 | 2.9 |
| Height (cm) | 165.7 | 8.8 | 161.7 | 7.6 |
| Weight (kg) | 65.2 | 11.3 | 58.2 | 4.4 |



Figure 3. Arithmetic mental workload levels against correct responsesforboth male and female.
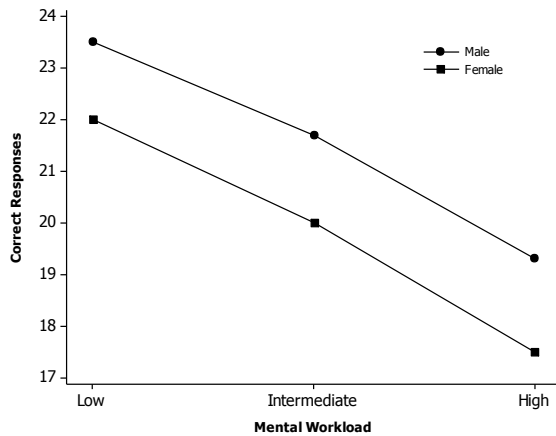
Figure 4. Spatial mental workload against correct responses for both male and female

The ANOVA technique showed that the level of difficulty of the arithmetic mental task significantly influenced the participants' ability to provide correct responses ($F(2,20) = 23.50$, p=0.000). The ANOVA findings revealed that the effect of the interaction between gender and arithmetic difficulty level on the number of correct responses was not significant ($F(2,20) = 0.118$, p=0.889). Moreover, a significant difference was observed between the intermediate level and the high level of the arithmetic test (p<0.001), whereas the difference between the low and intermediate levels was not significant (p=0.136). According to an independent t-test, no significant differences were observed between the performance of the male and female participants at all levels (i.e., low, intermediate, and high) (p>0.05). Table II illustrated the descriptive statistics of participants' correct responses at all tasks levels.

TABLE II.    CORRECT RESPONSES MEAN AND STANDARD DEVIATIONS ACROSS ALL PARTICIPANTS IN ARITHMETIC AND SPATIAL FIGURES TASKS.

| Task | Low level | | Intermediate level | | High level | |
|---|---|---|---|---|---|---|
| **Arithmetic Task** | Mean | SD | Mean | SD | Mean | SD |
| Male | 23.7 | 1.40 | 22.3 | 1.80 | 19.0 | 2.80 |
| Female | 23.5 | 1.52 | 21.8 | 2.86 | 18.7 | 2.50 |
| **Spatial Figures Task** | Mean | SD | Mean | SD | Mean | SD |
| Male | 23.5 | 1.5 | 21.7 | 1.6 | 19.3 | 2.4 |
| Female | 22.0 | 0.9 | 20.0 | 2.4 | 17.5 | 2.4 |

On the other hand, the ANOVA analysis illustrated that the impact of spatial figures mental task conditions was significant ($F(2,20) = 15.85$, p= 0.001). In contrast, no significant effect of the interaction between gender and spatial figures levels on the participants' performance ($F(2,20) = 0.023$, p=0.946). However, according to contrast tests, there was a significant difference between the low and intermediate levels and between the intermediate and high levels (p= 0.016 and p= 0.034, respectively).

### B. Participants' Accuracy and Time of Correct Responses

The percentage of correct responses (accuracy) of participants and the average time of correct responses for both tasks (arithmetic and spatial figures tasks) were generally related with the task difficulty levels.

The ANOVA showed that the levels of difficulty of the arithmetic mental task and spatial task significantly impacted the participants' accuracy ($F(2,20) = 40.909$, p=0.000). The ANOVA findings revealed that the effect of the interaction between gender and task type on response accuracy was not significant ($F(2,20) = 0.70$, p=0.480). Moreover, a significant difference was observed between the intermediate level and the high level of the arithmetic test (p=0.001), as well as the difference between the low and intermediate levels (p=0.005). In addition, when the task level (arithmetic and spatial) increased the accuracy decreased "see Figure 5".



Figure 5. Response accuracy associated with the three levels of mental workload for both arithmetic and spatial figures tasks.

In terms of response time, the repeated-measures ANOVA indicated that a significant impact was made by the task levels on participants' average correct responses ($F(2,20) = 606.46$, p<0.001), and when the task difficulty increased, the speed significantly decreased as shown in "see Figure 6". On the other hand, the ANOVA findings revealed that the effect of the interaction between gender and tasks types on response accuracy was not significant ($F(2,20) = 0.25$, p=0.778).. Moreover, a significant difference was observed between the intermediate level and the high level of the arithmetic test (p<0.001), also the difference between the low and intermediate levels was significant (p<0.001).

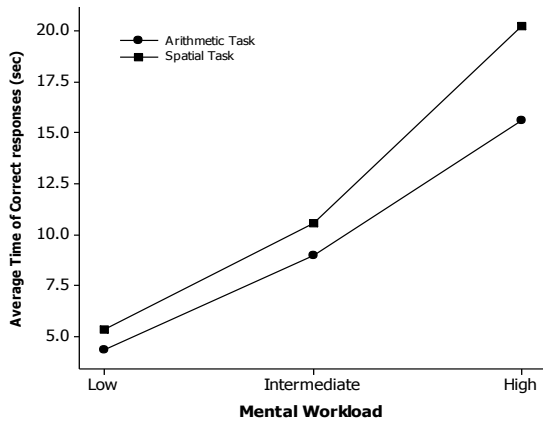Figure 6.    Average of correct response time (time in sec) associated with the three levels of mental workload for both arithmetic and spatial figures tasks.

### C. Physiological Parameters

- Heart Rate (HR)

The HR and HRV parameters were measured in order to determine the impact of the mental workload on the arousal level of the participants. As in previous research, a correlation was observed between these parameters and the difficulty level of the task. Table III presents the mean HR for participants as they completed the low level, intermediate level, and high level mental workload (i.e., both the arithmetic and spatial figures tasks). The table reveals that on average, the participants' HRs raised as the mental workload increased. However, the ANOVA indicated a significant effect of task type (i.e., arithmetic and spatial figures) on HR, $(F(1,11) = 30.28, p<0.001)$. Also, the data analysis indicated that a significant impact was made by the task levels on participants' HRs $(F(2,22) = 50.07, p<0.001)$, and when the task difficulty increased, the HR significantly increased as shown in "see Figure 7".

TABLE III. HAERT RATE OBSERVATION(BEATS/MIN) MEAN AND STANDARD DEVIATION ACROSS ALL PARTICIPANTS

| Task | Low level | | Intermediate level | | High level | |
|---|---|---|---|---|---|---|
| | Mean | SD | Mean | SD | Mean | SD |
| Arithmetic Task | 72.8 | 8.6 | 78.3 | 9.6 | 86.2 | 9.6 |
| Spatial Figures Task | 81.8 | 8.7 | 88.6 | 10 | 95.8 | 12 |

On the other hand, no significant impact from task type interaction and their levels on HR was observed $(F(2,22) = 0.224, p=0.775)$. According to the ANOVA analysis, there was a significant difference between the tasks $(p<0.001)$. The ANOVA results for the arithmetic task demonstrated that no significant influence was observed by gender and arithmetic levels interaction on participants' HRs $(F(2,20) = 0.531, p=0.596)$ "see Figure 8".
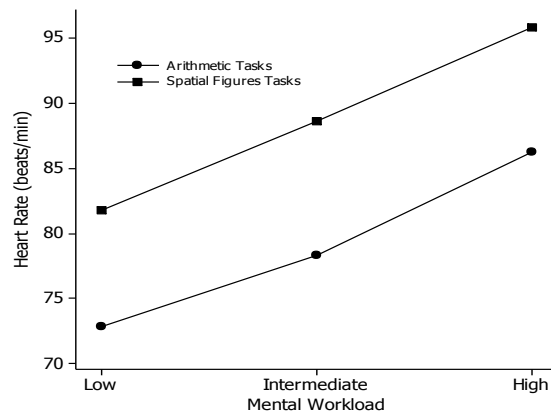


Figure 7.   Heart rate (HR) variable associated with the three levels of mental workload for  both arithmetic and spatial figures mental tasks.

In contrast, the difficulty level of the arithmetic task significantly affected the HRs of the participants $(F(2,20) = 31.15, p<0.001)$. According to repeated contrast comparisons, the mean HR increased significantly during the participants' completion of the high level arithmetic task $(p= 0.001)$ as compared to that of the intermediate level. Also, the HRs of participants rose significantly $(p= 0.007)$ during their completion of the intermediate level arithmetic task versus that of the low level arithmetic task. The independent t-test showed that no significant difference was observed between the mean HRs of both genders during completion of the low, intermediate, and high level task $(p>0.05)$.



Figure 8. Heart rate (HR) variable associated with the three levels of mental workload of arithmetic mental task for male and female.

Additionally, repeated-measures ANOVA results reported that the the spatial figures task workload affected the HRs of the subjects $(F(2,10) = 28.54, p<0.001)$. Conversely, the difference of the impact of gender and spatial task level interaction on the participants' HR was not significant $(F(2,20) = 2.51, p=0.107)$ "see Figure 9". According to repeated contrast comparisons, the HR significantly increased when participants completed the high level spatial figures tasks as compared to that of the intermediate and low levels $(p< 0.05)$. Also, the HR of participants rose significantly $(p< 0.05)$ when participants completed intermediate level spatial figures tasks versus low level spatial figures tasks. An independent t-test indicated

that no significant difference was observed between the mean HRs for both genders when completing the low, intermediate, and high levels (p>0.05) of spatial figures tasks.



Figure 9. Heart rate (HR) variable associated with the three levels of mental workload of spatial figures task for male and female.

- Heart Rate Variability (HRV)

ANOVA analysis indicated that task type (i.e., arithmetic versus spatial figures tasks) significantly affected the mean HRV $(F_{(1,11)} = 8.93$, p=0.012). Also, the data analysis indicated that the task level significantly impacted the participants' HRV $(F_{(2,22)} = 38.14$, p<0.001); i.e., when the arithmetic and spatial figures task difficulty level increased, the HRV significantly decreased as shown in "see Figure 10". Contrast comparisons were used to explore the differences that occurred in HRV when participants completed tasks in different difficulty levels. The HRV during completion of the difficult level was lower than that of the intermediate level (p<0.001), and the mean HRV was lower during completion of the intermediate level than that of the low level condition (p< 0.001).



Figure 10. Heart rate variability (HRV) variable associated with the three levels of mental workload for both arithmetic and spatial figures mental tasks.

In contrast, no significant difference was observed between the interaction of both tasks and their levels $(F_{(2,22)} = 0.884$, p=0.386). According to the ANOVA analysis, there was a significant difference between the tasks (p = 0.002).

However, according to the ANOVA analysis, the difficulty level of the arithmetic tasks significantly influenced the mean HRV $(F_{(2,20)} =20.13$, p<0.001). In contrast, the effect of the interaction between gender and arithmetic levels on the HRV of participants was not significant $(F_{(2,20)} =3.85$, p=0.065) "see Figure 11". A significant difference was observed between the low and intermediate levels (p=0.003), as well as between the intermediate and high levels (p=0.004). An independent t-test analysis indicated that a significant difference was observed between the HRV of males and that of female during the completion of low, intermediate, and high level of arithmetic tasks (p<0.05).



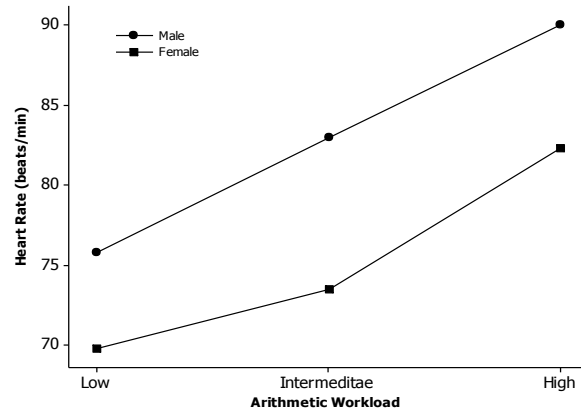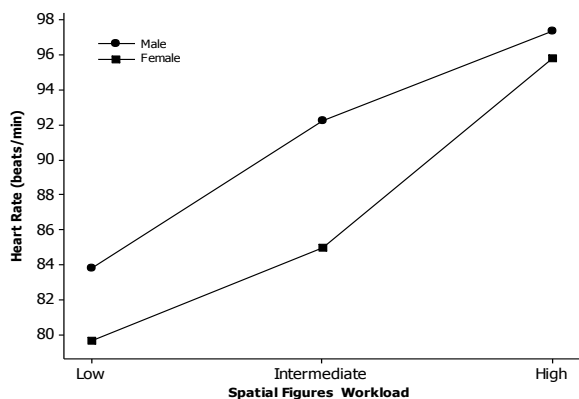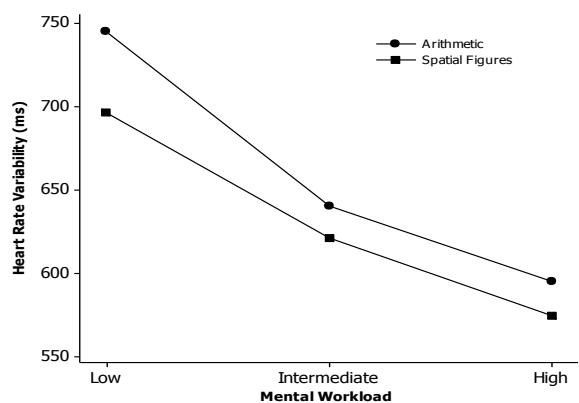Figure 11. Heart rate variability (HRV) variable associated with the three levels of mental workload of arithmetic mental task for male and female.

Additionally, the analysis findings remarked that the level of difficulty for the spatial figures tasks significantly impacted the participants' HRV (i.e., both male and female) $(F_{(2,20)} = 122.79$, p<0.001). However, the interaction between gender and spatial figures task level did not impact the participants' HRV $(F_{(2,20)} = 3.45$, p=0.082) "see Figure 12". Contrast comparisons indicated that the HRV decreased significantly when participants completed the high level spatial figures task as compared to that of the intermediate level (p<0.001). Also, the mean HRV dropped significantly (p<0.001) when participants completed the intermediate level spatial figures task versus that of the low level spatial figures task. Furthermore, the independent t-test presented that a significant difference was observed between the HRV of males and females when completing low, intermediate, and high level spatial figures tasks (p=0.008, p=0.01 and p=0.01, respectively).
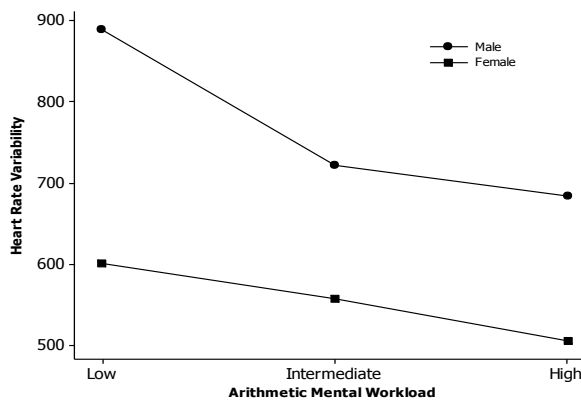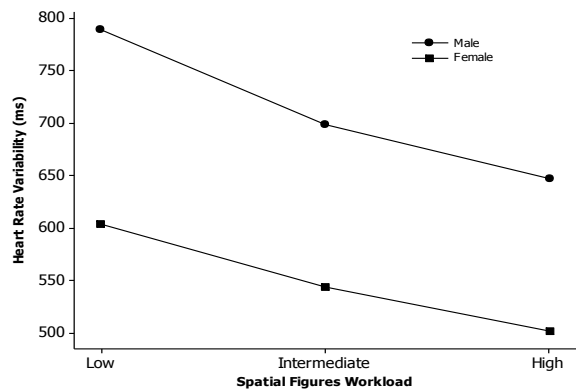
Figure 12. Heart rate variability (HRV) variable associated with the three levels of mental workload of spatial figures task for male and female.

- NASA-TLX Assessment Tool

The impact of the arithmetic task levels on overall NASA-TLX ratings was significant for both male and female participants $(F(2,20) =182.79, p<0.001)$. As arithmetic task difficulty increased, the overall NASA-TLX scores also increased $(p<0.001$ from low to intermediate; $p<0.001$ from intermediate to high). According to the independent t-test results, no significant difference was observed between the male and female overall workload scores for intermediate and high level arithmetic tasks $(p>0.05)$, whereas a difference in means was observed for the low level arithmetic tasks $(p= 0.026)$ "see Figure 13". However, the ANOVA results indicated that the impact from the interaction between gender and arithmetic levels on the NASA-TLX scores was not significant $(F(2,20) =0.320, p=0.730)$.
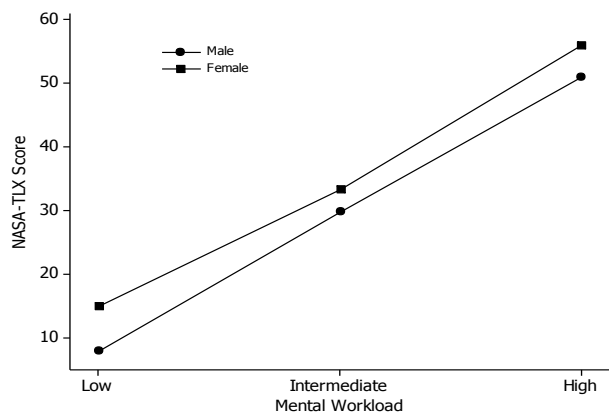


Figure 13.  Mean NASA-TLX ratings assessment tool associated with three levels of arithmetic task for both gender

Statistically significant differences were observed in the overall NASA-TLX ratings for the spatial figures tasks $(F(2,20) =379.85, p<0.001)$. In addition, the overall workload from the NASA-TLX increased when the task level became more difficult (low to intermediate, $p<0.001$; intermediate to high, $p<0.001$). On the other hand, the ANOVA results did not reveal any significant effect from the interaction between gender and spatial mental task levels on

overall NASA-TLX scores $(F(1,5) =3.21, p>0.05)$ "see Figure14". The independent t-test indicated that no significant difference was observed between the mean scores on the NASA-TLX for both genders in the intermediate and high levels $(p>0.05)$, whereas a significant difference was observed between the mean scores on the NASA-TLX males and that of females when completing the low level task $(p= 0.017)$. However, the overall NASA-TLX scores were significantly related to the scores on the mental demand dimension for both the arithmetic and spatial figures tasks $(r= 0.99, p<0.001$ and $r=0.99, p<0.001$, respectively).

In addition, all dimensions of the NASA-TLX increased except the physical dimension, which was not affected, since no physical effort was expended with either of these tasks.



Figure 14. Mean NASA-TLX ratings assessment tool associated with three levels of spatial figures task for both gender

## VIII.   DISCUSSION

The number of correct responses indicated that the arithmetic and spatial figures tasks levels were varied in difficulty. In addition, it presents that the compatibility of mental workload of the tasks was satisfied from low level to high workload level. Moreover, the participants' responses decreased in sequence as the mental demand from the arithmetic and spatial figures tasks increased. These findings are similar to that of previous research studies. Reference [18], found that performance decreased with increasing mathematical operation load. Regarding mental workload (i.e., monitoring and arithmetic tasks), Reference [2], pointed out that an increase in the monitoring and arithmetic process demand led to a degradation in the individual performance. Generally, the results showed that no significant difference has been observed between the genders when completing arithmetic task, whereas a considerable difference has been observed between the genders in terms of spatial ability. These findings are supported by [28], who reported no significant difference between the performance of men and women in arithmetic and language tasks. Furthermore, the findings revealed that a significant difference was observed between the HRV of each gender in terms of both tasks.

Also, the NASA-TLX analysis showed a significant gender difference in the mathematics and spatial tasks at low level, whereas there was no difference at intermediate or high levels in either task.

In addition, a significant difference in gender performance appeared only in the spatial ability task. This may be because the female participants took more time than the male participants in order to be as accurate as possible on the arithmetic task, which leads to more performance stress and perception of strain. This resulted in significant differences in the NASA-TLX ratings and the HRV variable between genders, yet resulted in no significance difference in their performances. References [15][16], mentioned that many previous studies determined the differences between males and females in terms of spatial capability on the mental rotation test; i.e., men tend to outperform women on spatial tasks. These results support the finding from this study; i.e., males provided more correct responses on the spatial figures task than females. Furthermore, the present results indicate that the interaction between gender and task levels did not significantly impact the participants' ability to provide correct responses. Finally, the participants' performance analysis of the arithmetic task showed that there was no significant correct response difference between the low and intermediate level, whereas there was a significant performance difference between intermediate and high level. This may create a potential problem for subsequent research, but it might be related to the sample size (in other words, a more significant effect may be found with a larger group of participants). However, the physiological and NASA-TLX scores analyses revealed a significant performance difference between the low vs. intermediate and intermediate vs. high levels of arithmetic tasks.

Physiological indices in this study were affected by the mental demand levels of the tasks (i.e., both arithmetic and spatial figures tasks). The physiological data analysis found a significant difference between the tasks, which appears to indicate that participants found the spatial task to be more difficult than the arithmetic task. That difference may be reduced by increasing the difficulty level of the arithmetic task, but this would produce a potential problem in the validation of these levels by pilot study. However, both tasks had a significant impact on physiological parameters, and these effects were parallel; therefore, this will not affect the study. However, the participants' (male and female) HR parameter was positively related to the arithmetic and spatial task difficulty levels. Specifically, as the arithmetic task level increased, the HR also increased, and the HR increased as the complexity of the spatial figures task increased. This was consistent with previous experimental studies [2][4] who found that the HR of subjects was affected by the complexity levels of the mental workload in the form of arithmetic and monitoring tasks. Their results indicated that when the difficulty level of the mathematical and monitoring tasks increased, the HR of the participants also increased. On the contrary, this study presented that HRV declined as the difficulty level of the arithmetic and spatial figures tasks increased. Reference [2], found that HRV declined as the

monitoring workload level increased. In addition, the findings showed that a significant difference was observed between the HRV of participants completing arithmetic tasks and the HRV of participants completing spatial figures tasks, whereas no significant difference was observed between the HR of subjects completing either task. According to [19], the increase in HRV leads to an increase in the acceptable responses; however, the decrease of HRV produces a decrease in the correct responses.

The study results indicate that a significant relationship exists between the subjective mental assessment tool (i.e., the NASA-TLX) ratings and the arithmetic and spatial figures task levels. Specifically, the NASA-TLX score increased with the increase in task level for both types of tasks. In general, the experimental data analysis indicated that the NASA-TLX scores were sensitive to changes in mental demand levels. This finding is similar to that of numerous papers. For example, [2][18], concluded that increases in NASA-TLX ratings were related to an increase in mental workload. For most participants, the highest NASA-TLX rating was recorded after completing the most complex arithmetic and spatial figures tasks. On the other hand, the lowest score was recorded after completing the simplest level of both tasks. However, the interaction between gender and task levels did not significantly impact the NASA-TLX scores. Furthermore, the study analysis indicated that no significant differences were observed between the male and female performance on intermediate and difficult levels for both tasks. In contrast, a significant difference was observed between the genders on the low level for both the arithmetic and spatial figures tasks.

## IX. CONCLUSION AND FUTURE WORK

In conclusion, the level of difficulty for arithmetic and spatial figures tasks were validated, which was the target of this experiment. Indeed, all of the variables (i.e., performance, physiological variables, and NASA-TLX scores) that were measured in this study indicated that the design of both tasks achieved three intensity levels (i.e., low, intermediate, and high) of mental effort. Furthermore, the participants' correct answers, HR, and NASA-TLX ratings increased when the arithmetic and spatial figures levels increased. In contrast, the HRV of the participants correlated negatively with the complexity level for both tasks; in other words, the HRV declined when the arithmetic and spatial task levels increased. Based on the findings of this study, each of these tasks appears to include three cognitive load conditions that are demanding enough to produce reliable differences between participants. Therefore, both tasks are seemingly suitable to use in the subsequent research programme. That future study aims to examine the effects of the combination of physical and mental demands on human attentional resources performance in dual occupations in a laboratory setting. The study will include visual resources (verbal and spatial resources), and the second study will include auditory resources (verbal and spatial), while future work is planned to investigate the overlaps between physical

and mental workload in a real domain such as factory production lines.

In summary, this paper studied the impact of mental workload on verbal and spatial attentional resources in order to validate levels of difficulties of two mental tasks. The first task was an arithmetic mental task used to show verbal resource. The second task used spatial figures to reflect spatial resource. In addition, gender difference was considered as a factor in this study because it will be looked at in a second, future experiment. Three parameters were used as dependent variables (performance, physiological parameters, and NASA-TLX score). The participants' responses decreased in sequence as the mental demand from the arithmetic and spatial figures tasks increased. Physiological indices in this study were affected by the mental demand levels of the tasks. The HR parameter of the participants (male and female) was positively related to the arithmetic and spatial task difficulty levels. On the other hand, this study showed that HRV declined as the difficulty level of the arithmetic and spatial figures tasks increased. In addition, the NASA-TLX score increased with the increase in task level for both types of tasks. The three difficult levels of arithmetic and spatial figure tasks were validated according to the results of the study.

REFERENCES

[1] Xie B. and Salvendy G, "Review and reappraisal of modeling and predicting mental workload in single- and multi-task environments,". Work & Stress,vol. 14, 2000, pp. 74-99.

[2] Hwang, S.L., Yau, Y.J., Lin, Y.T., Chen, J.H., Huang, T.H., Yenn, T.C. and Hsu, C.C, "Predicting work performance in nuclear power plants," Safety Science, vol. 46, 2008, pp. 1115-1124.

[3] Macdonald, W.A, "Workload, performance, health and well-being: a conceptual framework," In Karwowski, W. (Eds.). International Encyclopedia of Ergonomics and Human Factors, (Taylor and Francis Group, USA), 2001, pp2802-2807.

[4] Fredericks, T. K., Choi, S.D., Hart, J., Butt, S.E. and Mital A,"An investigation of myocardial aerobic capacity as a measure of both physical and cognitive workloads," International Journal of Industrial Ergonomics, vol. 35, 2005, pp. 1097-1107.

[5] Megaw, T. "The definition and measurement of mental workload,". In J. R. Wilson and N. Corlett (Eds.). Evaluation of Human Work, (Taylor and Francis Group, London), 2005, pp. 524-551.

[6] Perry, C.M., Sheik-Nainar, M.A., Segall, N., Ma R. and Kaber, D.B, "Effects of physical workload on cognitive task performance and situation awareness," Theoretical Issues in Ergonomics Science, vol. 9, 2008, pp. 95-113.

[7] Young, M.S. and Stanton, N.A,"Mental Workload". In N. A. Stanton , A. Hedge, K. Brookhuis and E. Salas (Eds.). Handbook of Human Factors and Ergonomics Methods, (Taylor and Francis Group, London), 2004, pp. 39.1-39.7.

[8] Louhevaara, V. and Kilbom, A, "Dynamic work assessment," In Wilson, J.R. and Corlett, N. (Eds.). Evaluation of human work, (Taylor and Francis Group, US), 2005, pp. 429-451.

[9] Matthews, G., Davies, D.R., Westerman, S.J. and Stammers, R.B,"Divided attention and workload". In. Human Performance Cogntion, Stress and Individual Differences, (Taylor and Francis Group, New York), 2000, pp. 87-106.

[10] Cox-Fuenzalida, L.E, "Effect of workload history on task performance," Human Factors vol..49, 2007, pp. 277-291.

[11] Neerincx, M.A. and Griffioen, E, "Cognitive task analysis: harmonizing tasks to human capacities," Ergonomics, vol. 39, 1996, pp. 543-561.

[12] Young, M.S. and Stanton, N.A,"Attention and automation: new prespectives on mental underload and performance," Therotical Issues in Ergonomics Science, vol. 3, 2002, pp. 178-194.

[13] De Zwart, B.C.H., Frings-Dresen, M.H.W. and Van Dijk, F.J.H, "Physical workload and the ageing worker: a review of the literature," Int. Arch Occupation Environment Health, vol. 68, 1995, pp. 1-12.

[14] Hyde, J.S., Fennema, E. and Lamon, S.J, "Gender differences in mathematics performance: A meta-analysis," Psychological Bulletin, vol. 107,1990, pp. 139-155.

[15] Voyer, D., Butler, T., Cordero, J., Brake, B., Silbersweing, D., Stern, E. and Imperato-McGinley, J, "The relation between computerized and paper-and-pencil mental rotation tasks: a validation study," Journal of Clinical and Experimental Neuropsychology, vol. 28, 2006, pp. 928-939.

[16] Peters, M. and Battista, C, "Application of mental figures of the Shepard and Metzler type and description of a mental rotation stimulus library," Brain and Cognition,vol. 66, 2008, pp. 260-264.

[17] Halpern, D.F, "Empirical evidence for cognitive sex differences," In. Sex Differences in Cognitive Abilities (3$^{rd}$ed), (Lawren Erlbaum Associates, USA), 2000, pp. 81-130.

[18] Didomenico A. and Nussbaum M. A, "Interactive effects of physical and mental workload on subjective workload assessment," International Journal of Industrial Ergonomics, vol. 38, 2008, pp. 977-983.

[19] Hansen, A.L., Johnsen, B.H. and Thayer, J.F, "Vagal influence on working memory and attention,"International Journal of Psychophysiological, vol. 48, 2003, pp. 263-274.

[20] Wickens, C. D, "Multiple resources and performance prediction," Theoretical Issues in Ergonomics Science, vol. 3, 2002, pp. 159-177.

[21] Wickens, C.D, "The structure of attentional resources,". In R. S. Nickerson (Ed.). Attention and performance VIII, (Lawrence Erlbaum Associates, Inc., New Jersey), 1980, pp. 239-257.

[22] Shepard, S. and Metzler, D,"Mental rotation of three dimentional objects," Science, vol. 171, 1971 , pp. 701-703.

[23] Hooven, C.K., Chabris, C.F., Ellison, P.T. and Kosslyn, S.M, "The relationship of male testosterone to components of mental rotation," Neuropsychologia, vol. 42, 2004, pp. 782-790.

[24] Sanz de Acedo Lizarraga, M.L. and Garcia Ganuza, J.M, "Improvement of mental rotation in girls and boys," Sex Roles, vol. 49, 2003, pp. 277-286.

[25] Hart, S. G. and Staveland, L. E, "Development of NASA-TLX (Task Load Index): results of empirical and theoretical research," In P. A. Hancock and N. Meshkati (Eds.), *Human Mental Workload*, (North-Holland: Amsterdam),1988, pp. 138-183.

[26] MathsNet, "Mental test,". Available from:www.mathsnet.net/form_mental.html [Accessed 18.06.2010], 2007.

[27] Bjornson, "Mental rotation training,". Available from: http://bjornson.inhb. de/?p=55 [Accessed 09.05. 2010], 2008.

[28] Skrandies, W., Reik, P. and Kunze, Ch, "Topography of evoked brain activity during mental arithmetic and language tasks: sex differences," Neuropsychologia, vol. 37, 1999, pp. 421-430.

# Guided Local Search in High Performance Detectors for MIMO Systems

Jaime L. Jacob, Yuri M. Mostagi, Taufik Abrão

DEEL - *Electrical Engineering Department*

*State University of Londrina* UEL

*Londrina, PR, Brazil*

*jaime.jacob@uel.br,   yuri.mostagi@gmail.com,   taufik@uel.br*

*Abstract*—**This work analyzes efficient non-spreading (a)synchronous MIMO detection topologies under realistic channels which results in high throughput and good performance × complexity trade-off. In this sense, we look for near-optimum efficient MIMO detections suitable for (un)coding schemes. Main system and channel parameters are analyzed, such as increasing number for transmitter and receiver antennas, number of iterations for convergence under AWGN and Rayleigh fading channels. Two heuristic local search MIMO detectors are compared with other near-optimum detectors, specifically SDR (semidefinite relaxation), expectation-maximization (EM), and linear multiuser detectors. Besides, the MIMO detectors performances under large MIMO systems (high number of transmitter and/or receiver antennas) are analyzed. The performance × complexity tradeoff results have indicated promising features for the guided local search (GLS) procedures in high capacity MIMO detectors.**

*Keywords*-**MIMO system, heuristic detectors, semidefinite relaxation.**

## I. Introduction

The capacity of a DS/CDMA system in multipath channel is limited mainly by the multiple access interference (MAI), self-interference (SI), near-far effect and fading. The Rake receiver explores the path diversity in order to reduce fading effect, but it is not able to mitigate the MAI [1], [2].

An alternative to solve this limitation is to apply the multiuser detection (MUD). The best performance is acquired by the optimum multiuser detection (OMUD), based on the log-likelihood function (LLF) [2]. However, this is achieved at costs of huge computational complexity which increases exponentially with the number of users. In the last decade, a variety of multiuser detectors with low complexity and sub-optimum performance were proposed, such as linear detectors, subtractive interference canceling [1], [2], semidefinite programming (SDP) approach [3]–[5] and heuristic methods [6]–[10].

In the near-optimal multiuser detection based on semidefinite relaxation (SDR-MuD), the optimal maximum likelihood (ML) detection problem is carried out by *relaxing* the associated combinatorial programming problem into an semidefinite programming (SDP) problem with both the objective function and the constraint functions being convex functions of continuous variables. SDR-MuD approach has been shown to yield near-optimal detection perfor-

mance in detecting binary/quadrature phase shift keying (BPSK/QPSK) signals [3]. On the other hand, there are few works dealing with high-order modulation heuristic detectors (HeurD) for MIMO systems. Particle swarm optimization (PSO) approach for MIMO detection with 16- and 64-QAM was considered in [11], [12]. A 16-QAM local search (LS) and hybrid PSO heuristic multiuser detectors suitable for DS/CDMA systems under SISO multipath channels has been considered in [13].

This work proposes a framework analysis for near-optimum detection suitable for non-spreading high-order modulation MIMO systems based on heuristic guided local search (GLS) approach, comparing with others well established detectors methods in the literature. For several MIMO detectors, the performance×complexity trade-off analysis is carried out, considering different systems and channels parameters in order to confirm the efficiency of the heuristic GLS-MIMO detectors approach. The non-spreading squared $M$-PSK MIMO system configurations in flat fading channels have been explored.

## II. Non-Spreading MIMO System Model

Figure 1 illustrates different four configurations possibilities for the channel that can be treated with the system model described herein. In this work we have explored the configuration a) single-user MIMO (non-spreading systems) non-selective fading channels with $M$-PSK or squared $M$-QAM modulation formats. Next we describe the adopted system model.

Consider a generic MIMO system with $K$ transmit antennas and $N$ receive antennas, and not necessarily $K \leq N$, where $K$ symbols are transmitted from $K$ transmit antennas simultaneously. Let $s_k$ be the symbol transmitted by the $k$th transmit antenna. Each transmitted symbol goes through the wireless channel to arrive at each of $N$ receive antennas. Denote the path gain from transmit antenna $k$ to receive antenna $n$ by $h_{nk}$. Considering a baseband discrete-time model for a AWGN or flat fading MIMO channel, the signal received at $n$th antenna is given by

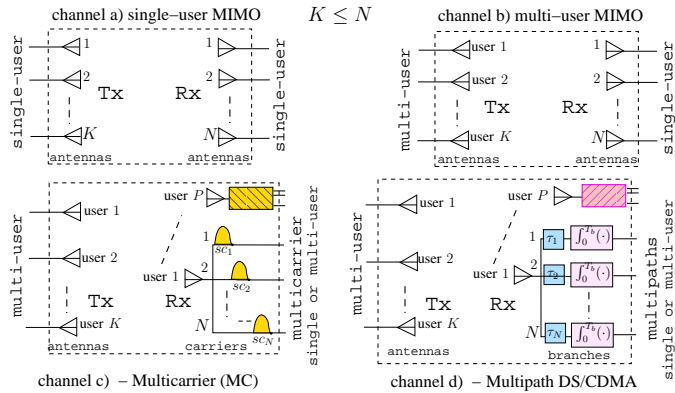$$y_n = \sum_{k=1}^{K} h_{nk} s_k + \eta_n \qquad (1)$$

Figure 1.    Channel configuration possibilities for the MIMO detection problem.

The $h_{nk}, \forall n \in \{1, 2, \dots, N\}, \forall k \in \{1, 2, \dots, K\}$ are assumed to be or i.i.d. complex Gaussian r.v's (fade amplitudes are Rayleigh distributed) with zero mean and $\mathbb{E}[(h_{nk}^I)^2] = \mathbb{E}[(h_{nk}^Q)^2] = 0.5$, where $h_{nk}^I$ and $h_{nk}^Q$ are the real and imaginary parts of $h_{nk}$, or the channel matrix is assumed unitary for the case of AWGN channel.

The noise sample at the $n$th receive antenna is assumed to be complex Gaussian with zero mean, and the samples $\{\eta_n\}, n = 1, \dots, N$, are assumed to be independent with:

$$\mathbb{E}[\eta_n^2] = N_0 = \frac{K E_s}{\gamma}$$

where $E_s$ is the average energy of the transmitted symbols, and $\gamma$ is the average received SNR per receive antenna [14].

The received signals are collected from all receive antennas, so (1) can be re writing in a vectorial form as:

$$\mathbf{y} = \mathbf{H}\mathbf{s} + \mathbf{n} \qquad (2)$$

where $\mathbf{y} = [y_1\, y_2\, \dots\, y_N]^T$ is the received signal vector, $\mathbf{s} = [s_1\, s_2\, \dots\, s_K]^T$ is the transmitted symbol vector, the $N \times K$ channel matrix $\mathbf{H}$, with channel coefficients $h_{nk}$, and $\mathbf{n} = [\eta_1\, \eta_2\, \dots\, \eta_N]^T$ is the noise vector. In a first analysis, $\mathbf{H}$ is assumed to be known perfectly at the receiver, and afterward errors in channel estimation matrix at received can be introduced. At the transmitter, the channel matrices are assumed completely unknown.

### III.   OPTIMUM DETECTION

The optimal maximum likelihood (ML) detector estimates the symbols for all $K$ users by choosing the symbol combination associated with the minimal distance metric among all possible symbol combinations in the $M = 2^m$ constellation points. So, ML detection in a memoryless non-spreading MIMO Gaussian channels ($K \times N$) can be formulated as:

$$\min_{\mathbf{s}} \qquad \|\mathbf{y} - \mathbf{H}\mathbf{s}\|_2^2 \qquad (3)$$
$$\text{s.t.} \qquad \{s_k\} \in \mathcal{A}, \quad k = 1, \dots, K$$

In order to avoid handle complex-valued variables, the separable squared QAM or $M-$PSK constellation is adopted. Hence, (3) can be redefined as the following decoupled optimization problem:

$$\min_{\mathbf{r}} \qquad \|\mathbf{z} - \mathbf{M}\mathbf{r}\|_2^2 \qquad (4)$$
$$\text{s.t.} \qquad r_i \in \mathcal{C} \subset \mathbb{Z}, \quad i = 1, \dots, 2K$$

with definitions:

$$\mathbf{z} := \left[ \begin{array}{c} \text{Re}\{\mathbf{y}\} \\ \text{Im}\{\mathbf{y}\} \end{array} \right] \in \mathbb{R}^{2N \times 1}; \quad \mathbf{r} := \left[ \begin{array}{c} \text{Re}\{\mathbf{s}\} \\ \text{Im}\{\mathbf{s}\} \end{array} \right] \in \mathbb{R}^{2K \times 1}$$
$$(5)$$
$$\mathbf{M} := \left[ \begin{array}{cc} \text{Re}\{\mathbf{H}\} & -\text{Im}\{\mathbf{H}\} \\ \text{Im}\{\mathbf{H}\} & \text{Re}\{\mathbf{H}\} \end{array} \right] \in \mathbb{R}^{2N \times 2K}; \qquad (6)$$

Clearly, (4) is a quadratic optimization problem with discrete variables in the set $\mathcal{A}$ and can be expressed as:

$$\min_{\mathbf{r}} \qquad \mathbf{r}^T \mathbf{Q} \mathbf{r} + \mathbf{q}^T \mathbf{r} \qquad (7)$$
$$\text{s.t.} \qquad r_i \in \mathcal{C} \subset \mathbb{Z}, \quad i = 1, \dots, 2K$$

where $\mathbf{Q} = \mathbf{M}^T \mathbf{M}$, $\mathbf{q} = -2\mathbf{H}^T \mathbf{z}$, and $\mathbf{r} = [r_1^I, r_2^I, \dots, r_K^I, r_1^Q, r_2^Q, \dots, r_K^Q]^T$, with $r_k^I$ and $r_k^Q$ the in-phase and quadrature component, respectively, for the $k$th user evaluated symbol. Note that the solution $\mathbf{r}^*$ in (7) represents the estimation symbol for all $K$ users, simply by composing the in-phase and quadrature components as: $r_k^* = r_k^{I*} + j r_k^{Q*}$. Note that if $K > N$, $\mathbf{Q}$ could become singular for some channel realization, once that $\mathbf{Q}$ is merely positive semidefinite. This difficulty can be removed by adding $\epsilon \mathbf{I}$ with small $\epsilon > 0$ to $\mathbf{Q}$.

The vector $\mathbf{r}$ in (7) is a discrete set with size dependence of $M$ and $K$ and can be solved directly using $m-$dimensional ($m = \log_2 M$) search method. Therefore, the associated combinatorial problem in an exhaustive search fashion has an exponential computational complexity that becomes prohibitive even for a moderate product $M K$ This ML detection problem can be solved efficiently by expanding the discrete feasible set into a continuous and convex feasible region [15]. Hence, manipulations, simplifications and relaxation over (4) or (7) is explored in the next section.

### IV.   SUB-OPTIMAL MIMO DETECTORS

Based on the recently proposed non-spreading MIMO detectors suitable either for coding as uncoding MIMO schemes, herein we investigate near-optimum detectors MIMO detectors under same system framework used in [16], but aiming to improve the performance $\times$ complexity trade-off. The goal is to obtain a structure which high overall throughput with good performance and relatively simple detection (or even decoding) using low complexity detectors topologies. Under uncoded-MIMO systems context, we have proposed LS, Hyb-opt LS and PSO (named GLS-MIMO) heuristic detectors. Hence, the performance $\times$ complexity trade-off of the proposed HeurD, specially GLS-MIMO

detectors, are compared with SDR, EM, linear parallel interference cancellation (PIC) and minimum mean squared error (MMSE) MIMO detectors. Below we discuss relaxations, simplifications, heuristic and classical criteria such as expectation-maximization and minimum mean squared error approaches suitable to non-spreading MIMO detectors.

### A. Semidefinite Programming Relaxations (SDR)

In a brute-force fashion, the conventional ML detector requires to examine all symbol combinations, i.e., $2^{mK}$ possibilities, or $M^{2K}$ for the equivalent decoupled optimization problem (4). Hence, the difficulty in adopting the OMUD is its high computational complexity, which is proportional to $\mathcal{O}(M^K)$. Therefore, when $K$ or/and $M$ increase, the computational complexity increases rapidly and this option becomes impractical. SDR and/or heuristic approaches are alternatives to deal with this problem, reducing complexity substantially, avoiding this huge complexity at an affordable performance loss in relation to optimum performance.

*1) Relaxation for Decoupled ML Uniform QAM MIMO Detection Problem:* Utilizing upper and lower bounds on the symbol energy in the relaxation step, a high-order QAM SDR MIMO detector with complexity that is independent of the constellation order for uniform QAM was proposed [17].

Under the hypothesis that $\mathcal{A}$ is an square alphabet and symmetric about the origin[1], the decoupled optimization problem posed by (4) can equivalently be rewritten as [17]:

$$\min_{\mathbf{X}} \quad \text{Trace}(\mathbf{Q}\mathbf{X}) \tag{8}$$
$$\text{s.t.} \quad \mathbf{X} \geq \mathbf{0}; \quad \text{rank}(\mathbf{X}) = 1;$$
$$\mathbf{X}_{2K+1,2K+1} = 1 \quad \mathbf{X}_{i,i} \in \mathcal{C}^2, \ i = 1, \ldots, 2K;$$

$$\text{with:} \quad \mathbf{x} := \begin{bmatrix} \mathbf{r}^T \\ t \end{bmatrix} \in \mathbb{R}^{2K+1}; \quad t \in \{\pm 1\} \tag{9}$$

$$\mathbf{Q} := \begin{bmatrix} \mathbf{M}^T\mathbf{M} & -\mathbf{M}^T\mathbf{z} \\ -\mathbf{z}^T\mathbf{M} & 0 \end{bmatrix}; \quad \text{and} \quad \mathbf{X} := \mathbf{x}\mathbf{x}^T \tag{10}$$

Since the optimization problem (8)–(10) has nonconvex constraints: a) rank constraint, $\text{rank}(\mathbf{X}) = 1$; b) squared finite alphabet constraints; $\mathbf{X}_{i,i} \in \mathcal{C}^2, \ i = 1, \ldots, 2K$, than, dropping the rank-one constraint a), and relaxing the constraints b) to the convex half-space (lower and upper) constraints:

$$\text{L} := \min_{a \in \mathcal{C}} a^2 \leq \mathbf{X}_{i,i} \leq \max_{a \in \mathcal{C}} a^2 =: \text{U}, \quad i = 1, \ldots, 2K,$$

we finally obtain the the SDR detector for non-spreading MIMO system:

$$\min_{\mathbf{X}} \quad \text{Trace}\left( \begin{bmatrix} \mathbf{M}^T\mathbf{M} & -\mathbf{M}^T\mathbf{z} \\ -\mathbf{z}^T\mathbf{M} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{r}^T \\ t \end{bmatrix} \begin{bmatrix} \mathbf{r}^T & t \end{bmatrix} \right)$$
$$\text{s.t.} \quad \mathbf{X} \geq \mathbf{0}; \quad \text{L} \leq \mathbf{X}_{i,i} \leq \text{U} \ i = 1, \ldots, 2K;$$
$$\mathbf{X}_{2K+1, 2K+1} = 1 \tag{11}$$

[1]Always valid for the QAM constellations.

with $\mathbf{x} \in \mathbb{R}^{2K+1}$, $t \in \{\pm 1\}$. As suggest in [17], the relaxed problem in (11) can be solved using any of the available modern SDP solvers, based on interior point (IP) methods, such as SeDuMi [18].

After this step, an approximate solution to the original problem can be generated using Gaussian randomization:

- drawing random vectors $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{X}_{\text{opt}})$, where $\mathbf{X}_{\text{opt}}$ denotes the solution of (11),
- quantizing each element of $\mathbf{x}$ to the nearest point in $\mathcal{C}$: $\check{\mathbf{x}} = \text{quantize}(\mathbf{x})$,
- reconstructing $\mathbf{s}$ from the quantized $\mathbf{x}$, i.e., $\mathbf{s}_i \leftarrow \check{\mathbf{x}}_i$, with e.g. $\mathbf{s}_i \in \{\pm 1, \pm 3\}$ for $16-$QAM.
- and finally picking the $\widehat{\mathbf{s}}$ that yields the smallest cost in original minimization problem (4).

Other strategies for approximating $\mathbf{s}$ in a squared-QAM NO-MUD problem can be obtained using a simple quantization or eigenvalue decomposition, as well as a simple randomization procedures [19], [20].

Previous results using different SDR higher-order QAM MIMO near-optimum detection strategies [17], [20]–[26]. In this sense, an efficient SDR detector was proposed in [26]. The focus is to achieve near-optimal BER performance with worst-case polynomial complexity. This is made by combining an optimized dual-scaling IP method for the relaxed SDP with a truncated version of the Sphere Decoder (SD) [27] and a dimension reduction strategy.

*2) SD × SDR Detectors:* The SD demonstrates impressive low running time for small systems operating in the high SNR regime; however for large systems or under low SNR regions, the running time grows exponentially. The core of the SD is based on exhaustive search. This feature is responsible for the increasing complexity under low SNR and/or huge number of users or antennas in MIMO systems.

On the other hand, the SDR detector is by nature insensitive to SNR changing, and its running time scales gradually with problem size. The insensitivity to SNR is a major ally in the low SNR regime, where the ML detection problem shows great difficulty. However, in the high SNR regime, the SDR algorithm fails to take advantage of the low noise property of the channel, when the SD takes advantage. The Sphere Decoding algorithm with adjustable radius search serves as a fast heuristic test of low noise channel realizations. The maximum number of sphere expansions is selected to ensure that complexity of the truncated Sphere Decoder does not dominate complexity of the dual-scaling algorithm.

It is worth to note that unlike the excellent performance obtained with BPSK modulation, the performance of SDR detector under higher-order QAM modulation formats is still considerably worse than achievable with the ML detector. This observation motivate us to propose heuristic alternatives for low-order modulation as well as high order squared-QAM formats.

## B. Guided Local Search Heuristic MIMO Detectors

Other approach to reach near-optimum performance consists in to apply heuristic procedures over (4). Comparisons among heuristic techniques performed for the MUD problem with BPSK modulation were carried out in [13], [28] and show that they are able to achieve performances close to the OMUD with low complexity. A local search $M-$QAM SISO MuD based on the BPSK s-LS-MUD has been analyzed in [28]. In a hybrid heuristic detection, the conventional Rake, or MMSE, or EM receiver output is adopted as the initial solution. Next, all unitary Hamming distance (from the initial solution) are evaluated individually through the equivalent quadratic minimization problem, eq. (7). The third step consists of switch the search to the simplified $k$-opt local search multiuser detector (s-LS-MUD), with $k = 1, 2$ [28], or even adopting PSO, genetic algorithm (GA), simulating annealing (SA), or other heuristic approach.

*1) 1-LS MIMO Detector:* The MIMO detector is based on the guided local search (GLS-MIMO detector) of all candidates with unitary Hamming distance regarding to the current solution.

In an exhaustive search approach, there are $M$ branches originating from each node in any tree search algorithm. So, for BPSK modulation ($M = 2$) always there are two branches, while in 16-QAM there are 16 branches originating from each node, regarding to a symbol of $k$th user or Tx antenna in a MIMO context. Therefore, the complexity is exponential with the number of users: $M^K$ candidates exist.

On the other hand, in a $k$-opt local search ($k-$LS) approach with binary modulation, at each iteration, the fitness function (LLF) is evaluated only $\binom{K}{k}$ times. Under strong interference environment, 1-LS MIMO degradation can be significant. However, the performance degradation can be mitigated combining $k$-opt and $p$-opt LS, with $p > 1$ and $p \neq k$, with marginal complexity increasing regarding to single $k$-opt LS. In general, the simplest strategy in $k$-$p$-opt LS consists in swapping the local search from $k$-opt to $p$-opt and vice-versa along the iterations every time there are no improvement in the fitness function value, eq. (4) or (7). Indeed, the performance-complexity of 1-2-opt LS (named "Hyb. 1-2-LS") MIMO detector are evaluated in Section V. The local search setting is swapped (1-opt to 2-opt) after 3 iterations with no fitness values improvement and a total number of iterations is set to $It = 3K$.

## C. Iterative Expectation-Maximization (EM) Detector

In order to compare the performances and complexities of the proposed GLS-MuD for MIMO channels, in this section an extension of the iterative expectation-maximization multiuser detector (EM-MuD) under BPSK DS/CDMA systems [29], [30] is provided. In [29], a EM-MuD for BPSK synchronous DS/CDMA under SISO AWGN channels was proposed, while in [30] an extension for MIMO flat-fading channels has been proposed.

In [30], the EM algorithm is applied to the maximum likelihood detection of BPSK synchronous DS/CDMA under MIMO (layered space-time codes) flat fading channels systems. The single data stream in the input is demultiplexed into $K$ substreams, and each substream is modulated independently; then transmitted over a rich-scattering wireless channel to $N$ received antennas The conditional log likelihood function (LLF) of a single layer is iteratively treated, rather than maximizing the intractable likelihood function of all layers. Computer simulations demonstrate some improvement of the EM-MIMO detection scheme with BPSK modulation in relation to the original V-BLAST scheme.

Herein, we analyze the performance-complexity of the BPSK EM-MIMO detector [30] under different scenarios. Since the symbols can only take the values $\{\pm 1\}$, the iterative decisions for the $k$th substream of a synchronous EM detector take the form:

$$\widehat{s}_k^{n+1} = \text{sign}\left\{\text{Re}\left[\widehat{s}_k^n \mathbf{h}_k^H \mathbf{h}_k + \beta_k \mathbf{h}_k^H \left(\mathbf{y} - \mathbf{H}\widehat{\mathbf{s}}^n\right)\right]\right\} \quad (12)$$

where $\widehat{\mathbf{s}}^n = [\widehat{s}_1^n \ \widehat{s}_2^n \ \dots \ \widehat{s}_K^n]^T$ is the symbol vector estimative at the $n$th iteration, and $\beta_k$'s are arbitrary real valued scalars satisfying $\sum_{k=1}^{K} \beta_k = 1, \ \beta_k \geq 0$.

For QPSK modulation, the $k$th symbol estimation at $n$th iteration above is given by $s_k^{n+1} = \text{csign}\{\cdot\}$, with the same argument of (12), where the complex decisor $\text{csign}\{a + jb\} = \text{csign}\{a\} + j\text{csign}\{b\}$.

Eq. (12) provides an iterative method to detect the symbols of all substream (or users). An appropriate initial value for symbol estimate is given by the output of the conventional maximum ratio combining (MRC) receiver $\mathbf{s}^0 = \mathbf{s}_{\text{mrc}}$. On the other hand, parameter $\beta_k$ has a critical role in the EM-based algorithm convergence. By setting $\beta_k = 0$, eq. (12) loses its iteration capability and reduces to the MRC receiver:

$$\widehat{s}_k = \text{sign}\left\{\text{Re}\left[\mathbf{h}_k^H \mathbf{h}_k \widehat{s}_k\right]\right\} = \text{sign}\left\{\text{Re}\left[\mathbf{h}_k^H \mathbf{y}\right]\right\}. \quad (13)$$

Assuming $\beta_k = 1$, eq. (12) becomes a linear PIC detector:

$$\widehat{s}_k^{n+1} = \text{sign}\left\{\text{Re}\left[\widehat{s}_k^n \mathbf{h}_k^H \mathbf{h}_k + \mathbf{h}_k^H \left(\mathbf{y} - \mathbf{H}\widehat{\mathbf{s}}^n\right)\right]\right\} \quad (14)$$

In [31], the $\beta_k$'s values were found experimentally, and found $\beta_k = 0.8$ for the best performance in a system with one and two receive antenna scenarios. In our simulation results, for $K$ and $N$ in the range of $[5; 20]$ antennas, the best $\beta_k = 0.8$ value was confirmed.

## D. Linear MMSE and Pseudo-Inverse **H** MIMO Detectors

The well known linear MMSE and channel pseudo-inverse (Pinv-H) based multiuser sub-optimal detectors are easily represented for $M-$PSK MIMO detection, respectively, as

$$\widehat{s}_k = \text{sign}\left\{\text{Re}\left[\mathbf{h}_k^H \left(\mathbf{h}_k^H \mathbf{h}_k + \sigma_k \mathbf{I}_N\right)^\dagger \mathbf{y}\right]\right\}, \quad \text{(MMSE)} \ (15)$$

and
$$\widehat{\mathbf{s}} = \text{sign}\left\{ \text{Re}\left[\mathbf{H}^{\dagger}\mathbf{y}\right] \right\}, \qquad \text{(Pinv-H)} \quad (16)$$

where $(\cdot)^{\dagger}$ represents the pseudo-inverse operator.

Numerical results for the analyzed and proposed NO-MUD MIMO detectors are discussed in the next section.

## V. NUMERICAL RESULTS

The performance of MIMO detectors were obtained by Monte-Carlo simulations, considering both AWGN and NLOS flat Rayleigh fading channels; the transmitted and received antennas were grouped into two categories: determined ($K \leq N$) and undetermined ($K > N$) MIMO channels. Figure 2 shows typical statistics for flat Rayleigh channel coefficients deployed in simulations. In order to facilitate the performance-throughput comparison analysis among the several MIMO detectors, low order modulation (BPSK) was assumed.



Figure 2. Typical statistics for the $h_{1,1}$ and $h_{10,10}$ Rayleigh channel coefficients.

### A. Performance under AWGN Channels

Figures 3 presents the MIMO detectors performance tendency under AWGN channels when number of transmitted antennas increases from $K = 5$ to 10, and to 20, while the number of receive antennas is held, $N = 10$. The GLS-MIMO detector (1-LS) advantage increases when the channel approaches to the determined limit condition, i.e. $K = N$ antennas. For the undetermined MIMO channel condition, Fig. 3.c indicates that a single guided local search (1-LS) is not enough to deal with the interference generated under degraded spatial eigen-mode ($K >> N$). Alternativaly, a low complexity GLS-MIMO is evaluated in the next subsection, namely hybrid 1-shift-2-LS.

Note that although the SDR approach result in high performance, the complexity is quite high when compared to the GLS-MIMO detectors, mainly when $K \geq N$, as discussed in Section V-C.



Figure 3. MIMO detectors performance for AWGN channels: a) $K = 5$ $N = 10$; b) $K = 10$ $N = 10$; c) $K = 20$ and $N = 10$ (undetermined).

### B. Performance under Flat Rayleigh Channels

Figures 4 and 5 illustrate the BER degradation reduction obtained with 1-2-opt LS ("Hyb. 1-2-LS" in legend) under BPSK modulation and flat Rayleigh channels. The swapping procedure between the two guided LS algorithms ($1 \rightleftarrows 2$-opt LS) occurs after 3 iterations with no fitness values improvement; the total number of iterations was set to

$It = 3K$. Figure 5 indicates the performance degradation *versus* SNR for the Conventional (MRC), Linear PIC, 1-opt LS, 1-2-opt LS, and SDR detectors under the increasing number of transmitted antennas (from $K = N = 3$ to $K = N = 12$).

The same MIMO detectors performance tendency in Figure 4 was observed with 16-QAM modulation. For the sake of space limitation, the MIMO detectors performance under higher $M-$QAM modulation orders were not shown herein.

It is worth to note that when the number of transmit and/or receive antennas increase, characterizing large MIMO systems, i.e., high number of $K$ and/or $N$ antennas (tens to hundreds) [32], Figs. 4.b and 5.c indicates a relative improvement performance of GLS MIMO detectors (1-LS and Shift 1-2-LS) regarding to the linear strategies (MMSE and Pinv-H) for low SNR. Additionally, the computational cost/complexity to obtain the channel matrix inverse is high in comparison with heuristic strategies, indicating a relative gain for heuristic approaches in these scenarios.



Figure 4. Performance results with $N = 10$, non-selective Rayleigh channels. a) $K = 5$; b) $K = 10$. $1 \rightleftarrows$ 2-opt LS occurs after 3 iterations with no improvement; $It = 3K$.



Figure 5. Performance for 1-LS and 1-2-LS with $K = N$ antennas under flat Rayleigh channel: a) $K = N = 3$; b) $K = N = 8$; c) $K = N = 12$. $1 \rightleftarrows$ 2-opt LS occurs after 3 iterations with no improvement; $It = 3K$.

### C. Complexity

Table I shows the complexity equations representing the number of complex multiplication/division and addition/subtraction operations for each analyzed MIMO detector. The four basic operations were considered with the same computational complexity. The complexity analysis was ex-

pressed by one iteration per substream/symbol (or antenna). The pseudo-inverse operator complexity was calculated with the Golub-Reinsch SVD [33].

Table I
NUMBER OF OPERATIONS COMPLEXITY FOR THE MIMO DETECTORS

| Detector | Eq. | Complexity |
|---|---|---|
| SDR | (11) | $16K^2N + 17KN + 4K^2 - 2N + 14K + 2$ |
| MRC | (13) | $2N - 1$ |
| Pinv-**H** | (16) | $11K^3 + 9K^2N + 5KN^2 - N$ |
| MMSE | (15) | $24N^3 + N^2 + 5N - 2$ |
| EM | (12) | $2NK + 6N - 1$ |
| Lin PIC | (14) | $2NK + 5N - 1$ |
| 1-LS | Ref. [28] | $8NK + N - 2K - 1$ |

Table II indicates the number of operations for each MIMO detector considering the antennas scenarios discussed before. Hence, in AWGN channel the operation are considered real values, while under flat Rayleigh channels, the operations are assumed complex values. In all cases it was assumed $N = 10$ received antennas. The complexity for the 1-LS and Shift 1-2-LS is almost the same, been omitted the Shift 1-2-LS complexity. One can see the MIMO detector with less number of operation for all $K = 5, 10$ or 20 transmitted antennas is the MRC, but with the worst performance. While the larger complexity is achieved by the MMSE, followed by the EM (with $It = 100$), and Pinv-**H**, respectively. Note that under $K = 10$ transmitted antennas scenario, the Pinv-**H**, MMSE and Lin-PIC with $It = 100$ detectors present approximately the same complexity of EM MIMO detector with $It = 100$.

Table II
NUMBER OF OPERATIONS COMPLEXITY FOR THREE MIMO CHANNELS SCENARIOS. $N = 10$ RX ANTENNAS.

| Detector (@$N = 10$) | Eq. | $K = 5$ | $K = 10$ | $K = 20$ |
|---|---|---|---|---|
| SDR | (11) | 5002 | 18218 | 69262 |
| MRC | (13) | 19 | 19 | 19 |
| Pinv-**H** | (16) | 6115 | 24990 | 133990 |
| MMSE | (15) | 24148 | 24148 | 24148 |
| EM $It = 18$ | (12) | 2862 | 4662 | 8262 |
| EM $It = 100$ | (12) | 15900 | 25900 | 45900 |
| Lin PIC $It = 18$ | (14) | 2682 | 4482 | 8982 |
| Lin PIC $It = 100$ | (14) | 14900 | 24900 | 49900 |
| 1-LS $It = 18$ | Ref. [28] | 7002 | 14202 | 28242 |
| 1-LS $It = 10$ | Ref. [28] | 3890 | 7890 | 15690 |

Analyzing the performance-complexity trade-off provided by Figs. 4, 5 and Table II one can see that the best choice is the GLS-MIMO Detectors (1-LS and Shift 1-2-LS).

## VI. CONCLUSIONS

The proposed simple guided local search MIMO detectors (1-LS and Shift 1-2-LS) have been shown attractive option regarding to the linear strategies (MMSE and Pinv-H), Expectation-Maximization approach, and conventional MIMO receivers (V-Blast and MRC topologies) as well, due to either high computational complexity in obtaining the channel matrix inverse (in case of linear strategies)

or the very poor performance although lower complexity of MRC and EM strategies. On the other hand the SDR performance is always better than GLS-MIMO detectors, but the computational demand is much more intensive than heuristic approaches.

Finally, under large MIMO systems scenarios (tens to hundreds $K$ and $N$), the GLS MIMO detectors presents a relative improvement performance under the conventional, linear, EM and inverted channel matrix approaches.

## REFERENCES

[1] S. Moshavi, "Multi-user detection for ds-cdma communications," *IEEE Communication Magazine*, vol. 34, pp. 132–136, Oct. 1996.

[2] S. Verdú, *Multiuser Detection*. New York: Cambridge University Press, 1998.

[3] P. H. Tan and L. K. Rasmussen, "The application of semidefinite programming for detection in cdma," *IEEE Journal on Selected Areas in Communication*, vol. 19, no. 8, pp. 1442–1449, Aug. 2001.

[4] W.-K. Ma, T. N. Davidson, K. M. Wong, Z.-Q. Luo, and P. C. Ching, "Quasi-maximum-likelihood multiuser detection using semi-definite relaxation with applications to synchronous cdma," *IEEE Trans. on Signal Processing*, vol. 50, no. 4, pp. 912–922, Apr. 2002.

[5] X. M. Wang, W.-S. Lu, and A. Antoniou, "A near-optimal multiuserdetector for ds-cdma using semidefinite programming relaxation," *IEEE Transactions on Signal Processing*, vol. 51, no. 9, pp. 2446–2450, Sept. 2003.

[6] C. Ergün and K. Hacioglu, "Multiuser detection using a genetic algorithm in cdma communications systems," *IEEE Transactions on Communications*, vol. 48, pp. 1374–1382, 2000.

[7] H. S. Lim and B. Venkatesh, "An efficient local search heuristics for asynchronous multiuser detection," *IEEE Communications Letters*, vol. 7, no. 6, pp. 299–301, June 2003.

[8] F. Ciriaco, T. Abrão, and P. J. E. Jeszensky, "Ds/cdma multiuser detection with evolutionary algorithms," *Journal Of Universal Computer Science*, vol. 12, no. 4, pp. 450–480, 2006.

[9] L. D. Oliveira, F. Ciriaco, T. Abrão, and P. J. E. Jeszensky, "Particle swarm and quantum particle swarm optimization applied to ds/cdma multiuser detection in flat rayleigh channels," in *ISSSTA'06 - IEEE International Symposium on Spread Spectrum Techniques and Applications*, Manaus, Brazil, 2006, pp. 133–137.

[10] ——, "Simplified local search algorithm for multiuser detection in multipath rayleigh channels," in *The 16th IST Mobile and Wireless Communications Summit*, Budapest, Hungary, July 2007, p. 5 pp.

[11] A. Khan, S. Bashir, M. Naeem, and S. Shah, "Heuristics assisted detection in high speed wireless communication systems," in *IEEE Multitopic Conference*, Dec. 2006, pp. 1–5.

[12] H. Zhao, H. Long, and W. Wang, "Pso selection of surviving nodes in qrm detection for mimo systems," in *GLOBECOM - IEEE Global Telecommunications Conference*, Nov. 2006, pp. 1–5.

[13] T. Abrão, L. D. de Oliveira, B. A. Angélico, and P. Jeszensky, *PSO Assisted Multiuser Detection for DS-CDMA Communication Systems*, ser. Particle Swarm Optimization: Theory, Techniques and Applications. Nova Science Publishers, Sept 2010, vol. 1, pp. 247–278.

[14] H. Jafarkhani, *Space-Time Coding: Theory and Practice*. Cambridge University Press, 2005.

[15] H. Peng, L. Rasmussen, and T. Lim, "Constrained maximum-likelihood detection in cdma," *IEEE Transactions on Communications*, vol. 49, no. 1, pp. 142–152, Jan. 2002.

[16] N. S. J. Pau, D. P. Taylor, and P. A. Martin, "Robust high throughput space time block codes using parallel interference cancellation," *IEEE Transactions on Wireless Communications*, vol. 7, no. 5, pp. 1603–1613, May 2008.

[17] N. Sidiropoulos and Z.-Q. Luo, "A semidefinite relaxation approach to mimo detection for high-order qam constellations," *IEEE Signal Processing Letters*, vol. 13, no. 9, pp. 525–528, Sept. 2006.

[18] J. F. Sturm, "Using sedumi 1.02, a matlab toolbox for optimization over symmetric cones," *Sofware Methods Optimization*, vol. 11-12, pp. 625–653, 1999, available online: http://sedumi.mcmaster.ca.

[19] W. K. Ma, P. C. Ching, and Z. Ding, "Semidefinite relaxation based multiuser detection for m-ary psk multiuser systems," *IEEE Transactions on Signal Processing*, vol. 52, no. 10, pp. 2862–2872, Oct. 2004.

[20] A. Wiesel, Y. C. Eldar, and S. S. (Shitz), "Semidefinite relaxation for detection of 16-qam signaling in mimo channels," *IEEE Signal Processing Letters*, vol. 12, no. 9, pp. 653–656, Sept. 2005.

[21] X. Wang and Z. Mao, "Multiuser detection for mc-cdma system with m-qam using semidefinite programming relaxation," in *PACRIM - IEEE Pacific Rim Conference on Communications, Computers and signal Processing*, Aug 2005, pp. 530–533.

[22] Y. Yang, C. Zhao, P. Zhou, and W. Xu, "Mimo detection of 16-qam signaling based on semidefinite relaxation," *IEEE Signal Processing Letters*, vol. 14, no. 11, pp. 797 – 800, Nov 2007.

[23] Z. Mao, X. Wang, and X. Wang, "Semidefinite programming relaxation approach for multiuser detection of qam signals," *IEEE Transactions on Wireless Communications*, vol. 6, no. 12, pp. 4275 – 4279, Dec. 2007.

[24] Y. Zhang, W.-S. Lu, and T. Gulliver, "Integer qp relaxation-based algorithms for intercarrier-interference reduction in ofdm systems," *Can. J. Elect. Comput. Eng*, vol. 32, no. 4, pp. 199–205, Fall 2007.

[25] S. Park, H. Zhang, H. Hongchao, D. Han, J. Kim, E. S. Kang, and W. W. Hager, "A fast suboptimal algorithm for detection of 16-qam signaling in mimo channels," in *MILCOM - IEEE Military Communications Conference*, Oct. 2007, pp. 1–7.

[26] M. Kisialiou and Z.-Q. Luo, "Efficient implementation of a quasi-maximum-likelihood detector based on semi-definite relaxation," in *ICASSP'07 - IEEE International Conference on Acoustics, Speech and Signal Processing*, vol. 4, April 2007, pp. 1329–1332.

[27] U. Fincke and M. Pohst, "Improved methods for calculating vectors of short length in a lattice, including a complexity analysis," *Mathematics of Computation*, vol. 44, pp. 463–471, 1985.

[28] L. de Oliveira, F.Ciriaco, T. Abrao, and P. Jeszensky, "Local search multiuser detection," *AEÜ International Journal of Electronics and Communications*, vol. 63, pp. 259–270, 2009.

[29] M. Borran and M. Nasiri-Kenari, "An efficient detection technique for synchronous cdma communication systems based on the expectation maximization algorithm," *IEEE Transactions on Vehicular Technology*, vol. 49, no. 5, pp. 1663–1668, Sept. 2000.

[30] K. R. Rad and M. Nasiri-Kenari, "Iterative detection for v-blast mimo communication systems based on expectation maximisation algorithm," *Electronics Letters*, vol. 40, no. 11, pp. 684–685, 27 May 2004.

[31] S. Iraji and J. Lilleberg, "Interference cancellation for space-time block-coded mc-cdma systems over multipath fading channels," in *VTC 2003 - IEEE 58th Vehicular Technology Conference*, vol. 2, Oct. 2003, pp. 1104–1108.

[32] K. V. Vardhan, S. K. Mohammed, A. Chockalingam, and B. S. Rajan, "A low-complexity detector for large mimo systems and multicarrier cdma systems," *IEEE Journal on Selected Areas in Communication*, vol. 26, no. 3, pp. 473–485, April 2008.

[33] G. H. Golub and C. F. V. Loan, *Matrix Computations*, 3rd ed. The Johns Hopkins University Press, 1996.

# Towards a Cognitive Handoff for the Future Internet:

## A Holistic Vision

Francisco A. González-Horta, Rogerio A. Enríquez-Caldera, Juan M. Ramírez-Cortés, Jorge Martínez-Carballido

Department of Electronics, INAOE
Tonantzintla, Puebla, México
{fglez, rogerio, jmram, jmc}@inaoep.mx

Eldamira Buenfil-Alpuche
Faculty of Engineering
Polytechnic University of Guerrero State (UPEG)
Taxco, Guerrero, México
e-mail: eldamira@gmail.com

*Abstract*— **Current handoffs are not designed to achieve multiple desirable features simultaneously. This weakness has resulted in handoff schemes that are seamless but not adaptive, or adaptive but not secure, or secure but not autonomous, or autonomous but not correct, etc. To face this limitation, we initiated a research project to develop a new kind of handoff system which attains multiple purposes simultaneously by using context information from the external and internal handoff environment. We envision a cognitive handoff as a multipurpose, multi-criteria, environment-aware, and policy-based handoff that trades-off multiple objectives to reach its intended goals. This paper presents a conceptual (soft) model of cognitive handoffs using a holistic approach. We applied the proposed model to identify cognitive handoff performance parameters and tradeoffs between conflicting objectives. We argue that cognitive handoffs are the archetype of handoffs for the future Internet.**

*Keywords-Cognitive handoff; future Internet; holism*

## I. INTRODUCTION

A handoff is intended to preserve the user communications while different kinds of transitions occur in the network connection. Thus, a handoff is the process of transferring communications among radio channels, base stations, IP networks, service providers, mobile terminals, or any feasible combination of these elements.

Significant desirable handoff features mentioned in the literature are, e.g., seamless [1], adaptive [2], autonomous [3], secure [4], and correct [5]; however, many others can be found in the vast literature of handoffs: transparent, reliable, flexible, robust, balanced, immune, fast, soft, smooth, lossless, efficient, proactive, predictive, reactive, QoS-based, power-based, location-aided, time-adaptive, intelligent, generic, etc. Despite the rich variety of desirable handoff features, two important problems remain unsolved: (1) how can be combined different desirable features into a single handoff process so that it can achieve many purposes simultaneously? (2) how to define every desirable feature so that ambiguity and subjectivity can be reduced?

This gap in knowledge about handoffs has produced a number of single-purpose schemes that successfully achieve one attractive feature but completely ignore others; e.g., seamless handoffs with poorly or null adaptation to other

scenarios or technologies [6]; adaptive handoffs that do not consider any security goal [2]; secure handoffs that ignore user autonomy [4]; etc. Also, there is a growing confusion in literature about similar features; e.g., accurate-correct, fast-timely, smooth-seamless, robust-reliable, etc. In order to reduce misuse and ambiguity of these attributes is convenient to associate a qualitative property (purpose) and quantitative measures (objectives and goals) to each desirable feature. By doing so, we can qualify and quantify their performance individually or in comparison with others. Major contributions of this research paper include:

*1) A new holistic vision of handoffs.* Many handoff solutions follow a reductionist approach; i.e., they achieve one desirable feature, use a small amount of handoff criteria, and work only in very specific scenarios. Although these simplistic solutions provide understanding and control of particular situations, we have seen how they quickly become special cases of more general models. Thus, we claim that the handoff problem for the future Internet requires holistic solutions, achieving multiple desirable features, using a diversity of context information, and adapting to any handoff scenario.

*2) A new conceptual model for cognitive handoffs.* We propose a new kind of handoff that is multipurpose, multi-criteria, context-aware, self-aware, policy-based, and trades-off multiple conflicting objectives to reach its intended goals. This paper provides the conceptual model and its first level of functional decomposition.

## II. SINGLE-PURPOSE VS MULTI-PURPOSE HANDOFFS

Dr. Nishith D. Tripathi in his outstanding thesis work published in 1997 [7] probably was the first author in considering a handoff that can simultaneously achieve many desirable features. His inspiring work served for many years as a basis for developing high performance handoffs; however, the complexities of handoff scenarios from 1997 to present days have changed significantly. For instance, the handoff concept changed from simple lower-layer transitions between base stations and channels to more elaborated cross-layer transitions among networks, providers, and terminals. The limited scope of Tripathi's handoff concept has brought in consequence that his algorithms and models become today special cases of more

general models. Holism is relevant in this way to provide a long-term solution for the handoff problem. Another author who describes several desirable handoff features is Nasser et al. [8] in 2006. Both, Tripathy and Nasser, described various desirable features, but they did not make any difference among features, purposes, objectives, and goals. A handoff model needs a clear distinction to such former concepts.

The holistic vision to the handoff problem has also been studied by Dr. Mika Ylianttila in his exceptional thesis work [9] published in 2005. He presented a holistic system architecture based on issues involved in mobility management areas (e.g., mobility scenarios, handoff strategies, handoff control, handoff algorithms, handoff procedures, mobility protocols, mobility parameters, performance measures, and handoff metrics). The work of Ylianttila improved the architecture of handoff issues that Pahlavan [10] published in 2000. However, these architectures have some drawbacks: i) they did not include the context management problem in their models; ii) they did not mention the tradeoffs that handoffs should consider in a multi-objective scenario; and iii) their architectures are based on types of issues and not in the functionality aspects of the handoff process.

Besides the above related work, we use two criteria to classify handoff schemes that are approaching to cognitive handoffs: the number of desirable features they achieve and the amount of context information they use. Handoff schemes, like the ones proposed by So [6] and Zhang [11], achieve only one desirable feature using limited context information; they provide seamless handoffs between particular network technologies and specific mobility scenarios. The schemes proposed by Siddiqui [12] and Hasswa [13] use broad context information, but they are focused only in one feature (seamlessness).

Conversely, the solutions proposed by Sethom [4] and Tuladhar [14] provide seamless and secure handoffs on a variety of handoff scenarios because they use broad context. The schemes proposed by Singhrova [2] and Chen [15] achieve seamless and adaptive mobility, but they cannot adapt to any handoff scenario because they use limited context. Finally, the scheme proposed by Altaf [16] achieves seamless, secure, soft, and adaptive handoffs, but just between WiMAX and 3G networks because they use limited context.

Considering this tendency, it will be common to observe in the near future a new generation of handoffs that can achieve many desirable features using broad handoff context information. In current literature, none architecture, model, or algorithm is reported to have this property.

Regarding the related work of standardization bodies, like the IEEE 802.21 and the IETF MIPSHOP, we observed that they are focusing in seamless heterogeneous handoffs; they are not taking into account the vast diversity of desirable features that handoffs could have. The IEEE 802.21 workgroup has approved three task groups to face very particular handoff scenarios: the IEEE 802.21a for security extensions to media independent handovers, the IEEE 802.21b for handovers with downlink only technologies, and the IEEE 802.21c for optimized single radio handovers. We believe they are following a

reductionist approach, but they lack the holistic vision of cognitive handoffs. Emmelman, in [17], discusses ongoing activities and scopes of these standardization bodies.

### III. THE COGNITIVE HANDOFF HOLISTIC VISION

#### A. Origin of Single-Purpose Handoffs

The thoughtful study of handoffs started in the early 1990s with the first generation (1G) cellular networks (e.g. AMPS [18]). These networks provided seamless conversations while the mobile phone switched between channels and base stations. The decision to perform a handoff was made only on a signal strength basis, but the handoff execution should be imperceptible to users. For this reason, the AMPS system required that the handoff gap be no more than 100 ms to avoid the possibility of dropping a syllable of speech [18]. These traditional handoffs are single-purpose/single-criterion or seamless/signal strength.

#### B. Major Challenges in the Future Internet

*1) Multidimensional Heterogeneity:* A major trend in future communication systems is the coexistence of multiple dimensions of heterogeneity integrated into a seamless, universal, uniform, ubiquitous, and general-purpose network. This future Internet will be seamless if it hides heterogeneity to users, universal if it can be used by anyone with any terminal, uniform if it is an all-IP network, ubiquitous if it is available anywhere and anytime, and general-purpose if it can provide any service. We divide heterogeneity into five dimensions as illustrated in Fig. 1 and explained in the next paragraphs. The arrows going down from the service provider dimension to the user mobility dimension depict two different handoff scenarios created by instantiating objects in each dimension.



Figure 1.  Multidimensional heterogeneity in the future Internet.

*a) Diversity on service providers and operators:* Offer different classes of services, billing models, security policies, and connection prices. They deploy different wireless technologies around the world and make roaming agreements and alliances with other providers and operators.

*b) Variety of applications and services:* Intend to fulfill the distinct ways of human communication; e.g., voice, video, data, images, text, music, TV, telephony, etc.

*c) Several access network technologies:* Include wired and wireless access technologies [19]; e.g., Ethernet, Bluetooth, WiMAX, WiFi, UMTS, MBWA, IMT-2000, GPRS, GSM, EDGE, LTE/SAE, DVB-HS, etc. They differ in terms of electrical properties, signaling, coding, frequencies, coverage, bandwidth, QoS guarantees, mobility management, media access methods, packet formats, etc.

*d) Plethora of mobile user terminals:* Users can be humans, machines, or sensors. Terminals for machines are integrated parts of machines. Sensor terminals collect information from networked sensors [19]. Terminals for humans are mobile and multimode, equipped with telecommunication capabilities and different saving energy characteristics; they change its factor form from those looked like computers (laptops, netbooks) to those looked like cell phones (PDAs, smartphones).

*e) Numerous user mobility states:* Network terminals can be located anywhere – in space, on the ground, under the ground, above water, underwater, and they can be fixed in a geographic position or moving at any speed – pedestrian, vehicular, ultrasonic [19].

Nowadays, no handoff solution exists which comprehensively addresses the entire scale of heterogeneity. Moreover, multidimensional heterogeneity has three main attributes: is inevitable, is the source of great amounts of context information, and produces an infinite number of handoff scenarios.

*2) Ubiquitous Connectivity:* It enables connectivity for anyone or anything, at any time, from anywhere. A myriad of wireless access technologies are spread across the entire world overlapping one another but avoiding interferences among them. Two requirements for ubiquitous connectivity are: (i) to develop scalable architectures to integrate any number of wireless systems from different service providers [20] and (ii) to develop smart multimode mobile terminals able to access any wireless technology [21].

*3) Cognitive Mobility:* It allows roaming mechanisms where the user is always connected to the best available network, with the smaller number of handoffs, service disruptions, user interventions, security threats, and the greater number of handoff scenarios.

### C. External and Internal Handoff Environment

We envision a cognitive handoff as a process that is both context-aware and self-aware. This implicates to make the handoff process aware of its external and internal environment. We borrowed the term 'cognitive' from Dr. Dixit vision of cognitive networking [22]. He defines *cognitive networking* as an intelligent communication system that is aware of its environment, both external and internal, and acts adaptively and autonomously to attain its intended goals. We believe cognitive handoffs not only should behave adaptively or autonomously to attain its intended goals, but also seamlessly, securely, and correctly.

On one hand, the external environment is directly related with all the external entities that provide a source of context information to the handoff process. These entities are users, terminals, applications, networks, and providers; a cognitive handoff should adapt to any kind of these entities. These entities maintain a strong cyclic relationship as follows: users interact with terminals, terminals run applications, applications exchange data through networks, networks are managed by providers, and providers subscribe users. The cyclic relationship of external entities suggests that all external context information emanates just from these five basic entities and no more; hence, if we ignore information of any of these entities, the handoff process will not adapt properly to all the scenarios. Therefore, a cognitive handoff should consider all the five entities.

On the other hand, the internal environment is another source of context and it is directly related with the behavior or performance of handoffs. This behavior directly depends on the desirable features of handoff. Next, we identified and describe five major desirable features which are considered highly significant for the current and future scenarios.

### D. Multiple Desirable Features of Handoff

*1) Seamlessness:* It means to preserve the user communications before, during, and after the handoff thus reducing service degradation or interruption. Service degradation may be due to a continuous reduction in link quality, network quality, handoff quality, QoS guarantees, and energy savings. Service interruption may be due to excessive degradations or a "break before make" approach.

*2) Autonomy:* This desirable feature is closely related to seamlessness. A handoff is autonomous, automatic, or autonomic when no user interventions are required during a handoff in progress. However, this does not mean that user interventions are not required in handoffs. It is good that users participate in the handoff configuration process by defining their preferences, priorities, or necessities; but, it is convenient that users can perform this activity offline to prevent any distraction during online communications.

*3) Security:* We say a handoff is secure if not new threats appear along the handoff process and security signaling traffic does not overload the network and degrades the communication services. This is a very challenging task, but if optimization techniques are used together with our model it could be shown that by minimizing handoff latency, authentication latency, and signaling overload, the risk of new threats appearance may be reduced.

*4) Correctness:* A handoff is correct if it keeps the user always connected to the best available network with the smaller number of handoffs; this is similar to the Gustaffson's vision of ABC defined in [23]. We consider that the best network is the one that is sufficiently better and consistently better. Furthermore, correctness can bring other additional features to the handoff process:

- *Beneficial*: if quality of communications, user expectations, or terminal power conditions get improved after handoff.
- *Timely*: if handoff is executed just in time; i.e., right after target is properly selected and before degradations or interruptions occur.
- *Selective*: if it properly chooses the best network among all the available networks.
- *Necessary*: if it is initiated because of one imperative or opportunist reason.

- *Efficient*: if it selects the most appropriate method, protocol, or handoff strategy, according to the types of: handoff in progress, user mobility, and application.

These handoff attributes derived from correctness, take special relevance during the decision-making phase, where it must be decided why, where, how, who, and when to trigger a handoff.

*5) Adaptability:* An adaptable handoff should be successful across any handoff scenario. A handoff is successful if it achieves a balance of every desirable feature at a minimum level of user satisfaction.

### E. Structure of Handoff Context Information

The handoff context information is extensive, heterogeneous, distributed, and dynamic. It supports the whole operation of the handoff process and the achievement of multiple desirable features. Therefore, such context information should be arranged in a clear structure. Table I and Table II show the structure of handoff context information according to a pair of criteria: the source of context and the class of information respectively. The sources of context originated in the external handoff environment support context-awareness while the one originated in the internal environment (the handoff process itself) will provide self-awareness.

TABLE I.    STRUCTURE FOR SOURCE OF CONTEXT INFORMATION

| |
|---|
| *User context:* This context allows users to customize the handoff according to their own needs, habits, and preferences. It includes: user preferences, user priorities, user profiles, user history, etc. |
| *Terminal context:* Allows the deployment of QoS-aware handoffs, power-based handoffs, and location-aided handoffs:<br>(a) *Link quality:* Received signal strength (RSS), signal to noise ratio (SNR), signal to interference ratio (SIR), signal to noise and interference ratio (SNIR), bit error rate (BER), block error rate (BLER), co-channel interference (CCI), carrier to interference rario (CIR), etc.<br>(b) *Power management:* Battery type (BT), battery load (BL), energy-consumption rate (ECR), transmit power in current (TPC), transmit power in target (TPT), power budget (PB), etc.<br>(c) *Geographic mobility:* Velocity (Vel), distance to a base station (Dist), location (Loc), direction (MDir), coverage area (GCA), etc. |
| *Application context:* This context includes the QoS requirements of active applications: Lost packets (LP), delayed packets (DP), corrupted packets (CP), duplicated packets (DuP), data transfer rate (DTR-goodput), packet jitter (PJ), out-of-order delivery (OOD), application type (AppT), etc. The consideration of these QoS parameters makes provisions for application-aware handoffs. |
| *Network context:* This context is needed to avoid selecting congested networks (befor handoff), to monitor service continuity (during handoff), and to assess the handoff success by measuring network conditions (after handoff): Network bandwidth (NBW), network load (NL), network delay (ND), network jitter (NJ), network throughput (NT), network maximum transmission unit (NMTU), etc. |
| *Provider context:* Connection fees, billing models, roaming agreements, coverage area maps, security management (AAA), types of services (data, voice, video), provider preferences, and provider priorities. A negotiation model may be required to equate the differences between service providers, network operators, and mobile users. |
| *Handoff performance context:* Call blocking (CB), call dropping (CD), handoff blocking (HOB), handoff rate (HOR), handoff latency (HOL), decision latency (DLat), execution latency (ExLat), evaluation latency (EvLat), handoff type (HOType), elapsed time since last handoff |

(ETSLH), interruptions rate (IR), interruption latency (IL), degradations rate (DR), degradation latency (DL), degradation intensity (DI), utility function (UF), signaling overload (SO), security signaling overload (SSO), improvement rate (ImpR), application improvement rate (AppImpR), user improvement rate (UsrImpR), terminal improvement rate (TermImpR), successful handoff rate (SHOR), imperative handoff rate (IHOR), opportunist handoff rate (OHOR), dwell time in the best (DTIB), authentication latency (AL), detected attacks rate (DAR), online user interventions rate (OUIR), tardy handoff rate (THOR), premature handoff rate (PHOR), etc.his context allows users to customize the handoff according to their own needs, habits, and preferences. It includes: user preferences, user priorities, user profiles, user history, etc.

TABLE II.    STRUCTURE FOR CLASS OF INFORMATION

| |
|---|
| *Handoff criteria:* Network discovery, decision-making, and performance evaluation. Some examples of handoff criteria include variables or parameters from the external/internal environment such as RSS, NL, BL, LP, HOL, Vel, connection price, etc. |
| *Handoff metrics:* Mathematical models used to measure several significant tasks of the handoff process; for instance, the quality of links, the quality of communications, the quality of different networks, the quality and quantity of handoffs, the quality of different providers, the achievement of user preferences, the power budget of a mobile terminal, the geographic mobility of a user, etc. Handoff metrics may combine a variety of handoff criteria and help any specific handoff algorithm to make optimal decisions. |
| *Performance measures:* Set of handoff metrics that are used to quantify performance of communications, performance of networks, performance of handoffs, and to evaluate the degree of achieving a handoff objective. |
| *Handoff policies:* Users and providers define a series of policies to the handoff operation. Policies define and specify rules for making handoff decisions in any particular situation; for instance, what to do if the link quality drops below a level required for an acceptable service. User and provider may have different views of the handoff process; provider may be interested in QoS while user in connection charges. Both points of view must be consistently integrated into a single handoff policy management database. |
| *Handoff constraints:* Conditions that must be satisfied in a particular handoff scenario and used to control the handoff operation by keeping performance parameters within specific limits. For instance, for a seamless handoff process, the delay has to be kept within certain boundaries; for real-time applications a delay of 50 ms could be acceptable, whereas non-real-time applications might accept delays as long as 3-10 sec [9]. |
| *Handoff configuration:* Defines preferences, priorities, and other configuration parameters required to customize the handoff operation. Typically, the configuration information is organized in a handoff profile linked to a particular user, provider, and terminal and should be initially performed offline either by the user, the provider, both or an auto-configuration setup. But, depending on the type of handoff algorithm, different configuration parameters may be required to be initialized, e.g. thresholds, timers, hysteresis, weights, etc. |

### F. Cognitive Handoff Conceptual Model

Once we have established and justified the necessity for developing a new handoff system, we present our conceptual model based on the statement that "a cognitive handoff should intend to achieve multiple desirable features and be aware of its entire environment by using information coming from multiple context domains". Fig. 2 depicts this basic idea by interconnecting multiple desirable features

with multiple context domains that we already explained separately in III.D and III.E.

The purpose of this model is to help people debate and discuss about the complexity of cognitive handoffs. Thus, topics of discussion would be related to level of complexity, correlation among desired features and context data, and the possibility of establishing handoffs as a multi-objective optimization problem as well as to give specifications for practical implementations. Used in this way this model is not intended for predicting, designing, or implementing cognitive handoffs, but for understanding and explaining such difficult and complex process  All the above issues have not been addressed in the handoff literature; therefore , in effect, the purpose of this conceptual model is being achieved.  Models like the one we present here are validated by credibility, and credibility comes from the way in which the cognitive maps are built and the clarity it represents most of the opinion's experts [24]. In the next section we provide some advances towards the development of cognitive handoffs.
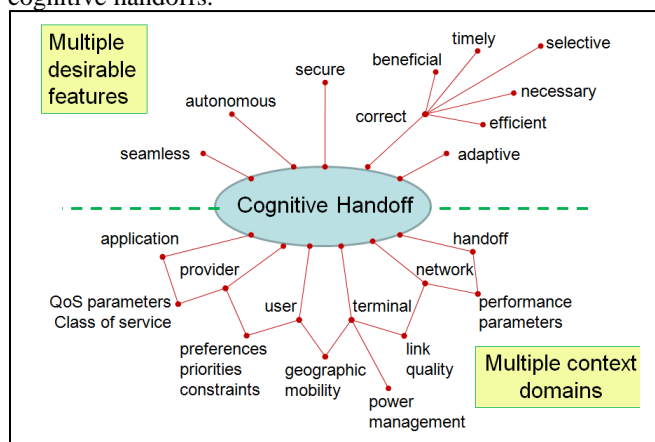


Figure 2.   Cognitive handoff conceptual model. The desired features to achieve determine the context data to use and vice versa.

## IV.   COGNITIVE HANDOFF MODEL AT WORK

### A.  Cognitive Handoff and Complex Systems

Cognitive handoffs are complex adaptive systems because: (1) they exhibit a complicated hierarchical structure (e.g., a power saving system is part of a network discovery system, which is part of a handoff system, which is part of a mobility system, which is part of a wireless communication system, and so on, but also a power saving system is part of the decision system, which is part of the handoff system, and so on); (2) the whole cognitive handoff system achieves purposes that are not purposes of the parts (e.g., a cognitive handoff purpose is to maintain the continuity of services, but this purpose is not defined in any of the parts or subsystems of the cognitive handoff system); and, (3) the handoff environment is dynamic and therefore adaptability is a desired handoff feature.

### B.  Correlating Desired Features and Context Data

With respect on whether all previously described context data are necessary to describe limitations on the model; one has to realize that the usage of certain context parameters

depends on the desirable features being implemented and the context data available in a moment will allow to accomplish or not a particular desired feature. Thus, we need to state a correct relationship or dependence between each desirable handoff feature and the subset of context data necessary to be accomplished. We made a correlation between desired features and context data by transforming desired features into purposes, purposes into objectives, objectives into goals, and goals into context data. For the sake of space in this manuscript the mentioned correlation will be shown in our next paper subtitled "Model-driven methodology and taxonomy of scenarios" already accepted for publication [25].

### C.  Advances for a Practical Implementation

The cognitive handoff system, represented in Fig. 2 by the oval in the middle, can be expanded into several sub-systems by using a functional decomposition approach [26]. Fig. 3 shows the main functional sub-systems for cognitive handoffs represented in ovals: handoff control algorithm, network discovery, handoff decisions, handoff execution, handoff evaluation, and handoff context information management. We briefly describe them:



Figure 3.          Functional decomposition model. The desired features provide purposes, objectives, and goals to achieve, while context domains provide the information needed to attain such goals.

- **Handoff Control Algorithm**: This is the main director of the handoff procedure. The entity which implements the control algorithm is called Handoff Control Entity (HCE). There should be one HCE in every user terminal and also there may be many others distributed across the network infrastructure. HCEs are agents that cooperate and compete to take a particular handoff to succeed.
- **Network Discovery**: This is the system for detecting and discovering available access networks. An available network is a reachable and authorized network considered for an eventual handoff.
- **Handoff Decisions**: The handoff decisions system is intended to answer the questions of why, when, where, how, and who should trigger the handoff. Typically this system has focused only in where and when to handoff [27]. The holistic vision extends the scope of handoff decisions.
- **Handoff Execution**: This system is intended to change the physical and logical connection from one network to

another, from one provider to another, or from one terminal to another. This change requires the most effective method, protocol, or strategy according to the current handoff scenario. The MIPSHOP group at IETF and the IEEE 802.21 standard are creating tools for implementing media-independent handoffs since 2003.

- *Handoff Evaluation*: This system measures the achievement of every desirable handoff feature and decides whether the executed handoff was successful or not. The evaluation results should be delivered after the handoff execution but within strict time constraints, thus this task is proactively distributed along the handoff process.
- *Handoff Context Information Management*: This system is intended to collect the distributed handoff context data, transform the data in information, and redistribute this information to the HCEs which are responsible for making handoff decisions and control.

Discovery, decisions, execution, and evaluation systems can be viewed as sequential stages of the handoff process; however, the context manager is a background process which permanently supplies the handoff control entities with fresh information about the handoff environment.

### D. Cognitive Handoff Performance Measures

The performance evaluation of cognitive handoffs requires a performance metric for each handoff purpose and a graphical representation to visualize multivariate data [28]. These metrics combine mathematically several performance measures that are associated to every handoff purpose. It is possible that metrics can normalize heterogeneous data into a single value representing the performance of each handoff purpose. Moreover, metrics can also be designed as utility functions so that greater values are better and all values are on the same scale. Fig. 4 exemplifies a radar graph comparing the performance of multiple handoff purposes simultaneously. We say that if all measures are within a boundary circle of acceptable quality, then the cognitive handoff is successful, otherwise the handoff is defective and outliers should be corrected.
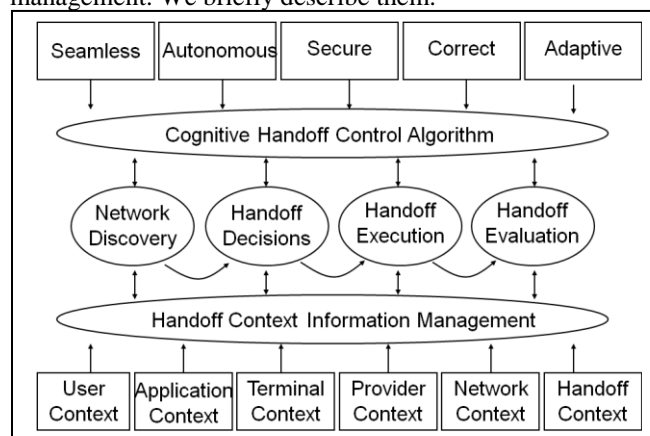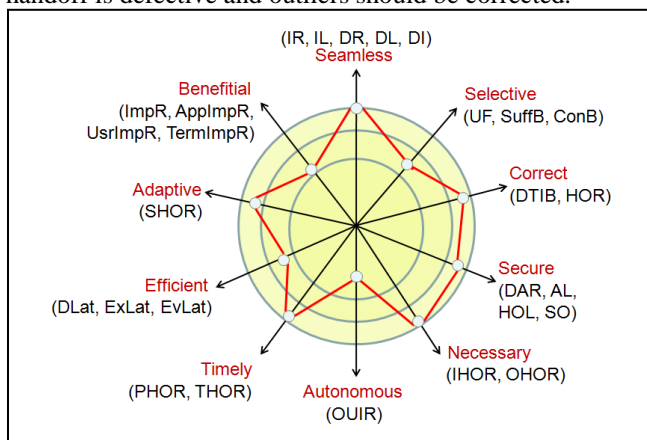


Figure 4. Functional decomposition model. The desired features provide purposes, objectives, and goals to achieve, while context domains provide the information needed to attain such goals.

### E. Formulating the Cognitive Handoff as a MOP

Let F be the set of desirable handoff features and C be the set of context data. We say that a context variable $Vi \in C$ is *correlated* with a desired feature $f \in F$ if and only if a change on the value of $Vi$ impacts on the purpose of f. For instance, some changes on the value of SNR may degrade or improve the link quality and impact on the purpose of seamlessness that is to maintain the continuity of services; thus, we say that SNR is correlated with seamlessness.

Let $Vf$ be the set of correlated variables with f, where $Vi \in Vf \subseteq C$. We say that $Vi$ is *positively correlated* with f if and only if increments on the value of $Vi$ produce improvements on the purpose of f and decrements on $Vi$ produce degradations on the purpose of f. For instance, increments on SNR improve the link quality, which improves the service continuity of seamlessness, and conversely, decrements on SNR degrade the link quality, which degrades the service continuity of seamlessness. Therefore, SNR is positively correlated with seamlessness.

$\uparrow$SNR $\rightarrow$ $\uparrow$LINKQUALITY $\rightarrow$ $\uparrow$SEAMLESSNESS
$\downarrow$SNR $\rightarrow$ $\downarrow$LINKQUALITY $\rightarrow$ $\downarrow$SEAMLESSNESS

We say that $Vi$ is *negatively correlated* with f if and only if increments on the value of $Vi$ produce degradations on the purpose of f and decrements on $Vi$ produce improvements on the purpose of f. For example, increments on BER degrade the link quality, which degrades the service continuity of seamlessness, and conversely, decrements on BER improve the link quality, which improves the service continuity of seamlessness. Therefore, BER is negatively correlated with seamlessness.

$\uparrow$BER $\rightarrow$ $\downarrow$LINKQUALITY $\rightarrow$ $\downarrow$SEAMLESSNESS
$\downarrow$BER $\rightarrow$ $\uparrow$LINKQUALITY $\rightarrow$ $\uparrow$SEAMLESSNESS

The set $Vf$ is partitioned in two subsets $Vf^+$ and $Vf^-$ where $Vf^+$ is the set of variables positively correlated with f and $Vf^-$ is the set of variables negatively correlated with f.

Every $Vi$ has a weight $Wi$ associated to its priority where $Wi \in \Re[0,1]$ and $\sum Wi = 1$. Let **v** represent the vector of variables **v** = (V1, V2, …, Vm), then the *objective function* for the desired handoff feature f is defined by

$$f(\mathbf{v}) = \sum(K+Wi)\log(Vi^+) - \sum(K+Wi)\log(Vi^-) \quad (1)$$

where K is a scaling factor so that small changes on the context variables reflect big changes on f(**v**).

In this general objective function, $Vi^+$ and $Vi^-$ are positively and negatively correlated variables of f. The objective function $f(\mathbf{v}) : \Re^m \rightarrow \Re$ is a utility function that we want to maximize because in desirable features the higher the value the best.

Considering k different objective functions $f_i$ that we want to maximize simultaneously where some of them may

be in conflict, then the multi-objective optimization problem (MOP) can be stated as the problem of

$$\text{Maximize } \{f_1(\mathbf{v}), f_2(\mathbf{v}), \ldots, f_k(\mathbf{v})\} \tag{2}$$
$$\text{Subject to } \mathbf{v}_l \leq \mathbf{v} \leq \mathbf{v}_u ,$$

where $\mathbf{v}_l$ and $\mathbf{v}_u$ represent the vectors of lower and upper values of the tolerance range for each variable.

### F. Tradeoffs between Conflicting Objectives

A cognitive handoff is designed to achieve multiple purposes, objectives, and goals simultaneously. In the space of handoff objectives, we can distinguish between those with complementary nature and those with competitive nature. Complementary objectives can be simultaneously optimized without any conflict between them, but competing objectives cannot be simultaneously optimized, unless we find compromised solutions, largely known as the tradeoff surface, Pareto-optimal solutions, or non-dominated solutions [29]. We describe several tradeoffs to consider in a multi-objective handoff scheme:

a) *(Max. DTIB and Min. HOR)*: There is a tradeoff between maximizing the time to stay always best connected (DTIB) and minimizing the number of handoffs (HOR). The conflict arises because in a dynamic environment the best network is changing frequently and stochastically; thus, to maximize DTIB is necessary to make frequent handoffs as soon as a new best is available. This increase in the number of handoffs creates a conflict with minimizing HOR.

b) *(Min. DLat and Max. SHOR)*: This tradeoff is between minimizing the handoff decisions latency (DLat) and maximizing the number of successful handoffs (SHOR). The conflict emerges because the less time elapsed to make decisions will necessary lead to reduce the number of successful handoffs. For example, in case of imperative handoffs, DLat is reduced but this may lead to select an incorrect target because the selection time is also reduced.

c) *(Max. Sizeof-ContextInfo and Min SO)*: This is a tradeoff between minimizing the handoff signaling overload (SO) and maximizing the amount of handoff context information to be managed by the handoff control entities. The conflict arises because broad handoff information is required to attain multiple desirable features, but this will increase the amount of signaling traffic in the network.

d) *(User and Provider Preferences)*: Several conflicts may appear due to differences between provider and user preferences. For instance, providers may prefer networks within its own administrative domain while users may prefer networks with lower charges even if they are owned by other service providers; users may prefer a Mobile Controlled Handoff (MCHO) while providers may prefer Network Controlled Handoffs (NCHO). Conflicts like these require a balance between different interests. Handoff protocols like Mobile Assisted Handoff (MAHO) and Network Assisted Handoff (NAHO) try to balance the handoff control [7].

## V.  CONCLUSION

Handoffs are an integral component of any mobile-wireless network from past, present, and future. Handoffs are transitions that change the data flows from one entity to another, where these entities may be radio channels, base stations, IP networks, service providers, and user terminals. The handoff process should exhibit several desirable features beyond seamlessness and should consider more context information beyond the signal strength. This is a common requirement to face the handoff scenarios of the future Internet.

The existing handoff schemes are not able to achieve a variety of attractive features and managing arbitrary amounts of context information. Therefore, we proposed a conceptual model to create handoffs of this kind. We characterized a cognitive handoff to be multipurpose, multi-criteria, context-aware, self-aware, and policy-based.

We claimed that our cognitive handoff model is holistic because it considers all the transition entities that may be involved in handoffs, all the external and internal sources of context, and considers many significant desirable features.

Using a functional decomposition approach, we divided the functional behavior of a cognitive handoff into six general modules: control algorithm, network discovery, handoff decisions, handoff execution, handoff evaluation, and context management. Each module has assigned a purpose to every feature and decomposed each purpose into objectives and goals. We applied the cognitive handoff model to define its performance parameters and significant tradeoffs between conflicting objectives.

As a future work, we are preparing another manuscript for presenting a new taxonomy of handoff scenarios and the model-driven methodology that we are using to develop cognitive handoffs. There is still much work to do before we can see cognitive handoffs practically implemented. The cognitive handoff project follows theoretical and practical avenues. A theoretical challenge is to further develop the cognitive handoff MOP to study the structure of the variables in the handoff context (e.g., continuous/discrete, deterministic/stochastic, etc.) and the types of constraints required to create a convex optimization problem. In the practical and Applicability Avenue, we have deployed temporal and geometric simulation models to observe and predict the behavior of cognitive handoffs with two conflictive objective functions; however, further development is required to demonstrate the feasibility and applicability of cognitive handoffs in complex scenarios.

## References

[1]  M. Satyanarayanan, M. A. Kozuch, C. J. Helfrich, and D. R. O'Hallaron, "Towards seamless mobility on pervasive hardware" *Pervasive and Mobile Computing 1*, Elsevier, pp. 157-189, 2005.

[2]  A. Singhrova and N. Prakash, "Adaptive Vertical Handoff Decision Algorithm for Wireless Heterogeneous Networks" *11th IEEE Intl. Conf. on High Performance Computing and Communications*, pp. 476-481, 2009.

[3] J. M. Kang, H. T. Ju, and J. W. K. Hong, "Towards Autonomic Handover Decisions Management in 4G Networks" A. Helmy et al. (Eds.): MMNS 2006, LNCS 4267, IFIP, pp. 145-157, 2006.

[4] K. Sethom, H. Afifi, and G. Pujolle, "Secure and Seamless Mobility Support in Heterogeneous Wireless Networks" *Proc. IEEE Globecom*, pp. 3403-3407, 2005.

[5] K.D. Wong and D.C. Cox, "A Pattern Recognition System for Handoff Algorithms" *IEEE Journal on Selected Areas in Communications* 18 (7), pp. 1301-1312, July 2000.

[6] J. W. So, "Vertical Handoff in Integrated CDMA and WLAN Systems" *International Journal of Electronics and Communications* (AEÜ) 62, Elsevier, pp. 478-482, 2008.

[7] N. D. Tripathi, "Generic Adaptive Handoff Algorithms Using Fuzzy Logic and Neural Networks" Ph.D. dissertation, Virginia Polytechnic Institute and State University, August 21, 1997.

[8] N. Nasser, A. Hasswa, and H. Hassanein, "Handoffs in Fourth Generation Heterogeneous Networks" *IEEE Communications Magazine*, pp. 96-103, October 2006.

[9] M. Ylianttila, "Vertical Handoff and Mobility – System Architecture and Transition Analysis" Ph.D. dissertation, Faculty of Technology, Dept. of Electrical and Information Engineering, University of Oulu, Finland, May 6, 2005.

[10] K. Pahlavan, et al., "Handoff in Hybrid Mobile Data Networks" *IEEE Personal Communications*, pp. 34-47, April 2000.

[11] Z. Zhang and A. Boukerche, "A Novel Mobility Management Scheme for IEEE 802.11-based Wireless Mesh Networks" *Intl. Conf. on Parallel Processing*, pp. 73-78, 2008.

[12] F. Siddiqui and S. Zeadally, "Mobility Management across Hybrid Wireless Networks: Trends and Challenges" *Computer Communications* 29, Elsevier, pp. 1363-1385, 2006.

[13] A. Hasswa, N. Nasser, and H. Hassanein, "Generic Vertical Handoff Decision Function for Heterogeneous Wireless Networks" *IEEE*, 2005.

[14] S. R. Tuladhar, C. E. Caicedo, and J. B. D. Joshi, "Inter-Domain Authentication for Seamless Roaming in Heterogeneous Wireless Networks" *IEEE Computer Society*, pp. 249-255, 2008.

[15] W. T. Chen, J. C. Liu, and H. K. Huang, "An Adaptive Scheme for Vertical Handoff in Wireless Overlay Networks" *Proc. 10th Intl. Conf. on Parallel and Distribution Systems*, 2004.

[16] A. Altaf, F. Iqbal, and M. Y. Javed, "S3H: A Secure, Seamless and Soft Handover between WiMAX and 3G Networks" *Intl.*

*Conf. on Convergence and Hybrid Information Technology*, pp. 530-534, 2008.

[17] M. Emmelmann, et al., "Moving toward seamless mobility: state of the art and emerging aspects in standardization bodies" *Wireless Pers Commun* 43, Springer, pp. 803-816, 2007.

[18] B. A. Black, et al., "Introduction to Wireless Systems" Ch. 4, Prentice Hall, 1st Edition, pp. 125-140, May 2008.

[19] J. L. Salina and P. Salina, "Next Generation Networks: Perspectives and Potentials" John Wiley & Sons, England, 2007.

[20] S. Mohanty and J. Xie, "Performance Analysis of a Novel Architecture to Integrate Heterogeneous Wireless Systems" *Computer Networks* 51, Elsevier, pp. 1095-1105, 2007.

[21] P. Koch and R. Prasad, "The Universal Handset" *IEEE Spectrum*, vol. 6 num. 4 International, pp. 32-37, April 2009.

[22] Q. H. Mahmoud (Edt), "Cognitive Networks: Towards Self-Aware Networks" Foreword 2: S. Dixit, J. Wiley & Sons, 2007.

[23] E. Gustaffson and A. Jonsson, "Always Best Connected: 3G Mobile Network Technologies and Experiences" *IEEE Wireless Communications*, pp. 49-55, February 2003.

[24] M. Pidd (Edt), "Systems Modelling: Theory and Practice," John Wiley & Sons, England, pp. 1-42, 2004.

[25] F. A. González-Horta, R.A. Enríquez-Caldera, J. M. Ramírez-Cortés, J. Martínez-Carballido, E. Buenfil-Alpuche, "Towards a cognitive handoff for the future Internet: Model-driven methodology and taxonomy of scenarios" *2nd International Conference on Advanced Cognitive Technologies and Applications*, COGNITIVE 2010, Lisbon, Portugal, Nov. 2010.

[26] O. Maimon and L. Rokach, "Decomposition Methodology for Knowledge Discovery and Data Mining: Theory and Applications" Series in Machine Perception Artificial Intelligence, Vol. 61, World Scientific Pub., London, 2005.

[27] Q. Song and A. Jamalipour, "A Time-Adaptive Vertical Handoff Decision Scheme in Wireless Overlay Networks" *17th Annual Intl. Symposium on Personal, Indoor, and Mobile Radio Commun*, PIMRC'06, 2006.

[28] J. Heer, M. Bostock, and V. Ogievetsky, "A Tour Through the Visualization Zoo" *Communications of the ACM*, vol. 53, no. 6, pp. 59-67, June 2010.

[29] J. Branke, et al. (Eds.), "Multiobjective Optimization: Interactive and Evolutionary Approaches" Springer, Germany, 2008.

# Complexity and Chaos Analysis of a Predator-Prey Ecosystem Simulation

Yasaman Majdabadi Farahani
Department of Computer Science
University of Windsor
Windsor, Canada
majdaba@uwindsor.ca

Abbas Golestani,                    Robin Gras
Department of Computer Science
University of Windsor
Windsor, Canada
golesta@uwindsor.ca,              rgras@uwindsor.ca

*Abstract— We investigated the complexity level of an agent-based predator/prey ecosystem simulation. The variations of the times series associated to this ecosystem simulation are the result of complex simulation mechanisms. For the purpose of understanding how close our system is to random or chaotic processes, we compare these data with data generated by a Markov chain as a simple process. The parameters of the corresponding Markov matrix are learned from the data generated by our simulation. Then we used the Markov chain to generate data similar to those of the simulation. We show that the Markov chain for all three orders, which we tested, generated prey and predator time series that are more random than their counterpart in the original simulation. Also, we used the largest Lyapunov exponent to determine the chaotic behavior of the simulation. We discuss the largest Lyapunov exponent values for population time series of both prey and predator agents, which indicates chaotic behavior in our agent-based ecosystem simulation.*

*Keywords- agent-based ecosystem; chaos analysis; complexity analysis; Markov chain*

## I.    INTRODUCTION

Few attempts have been made to model a complete ecosystem and analyze its complex behavior using an agent-based approach. A predator-prey model proposed by Ward et al [4] in which the agent model is dedicated to represent schooling behaviors, and the evolution is an offline mechanism using a genetic algorithm. More recently, Ronkko [5] has proposed a high-scale simulation based on a particle system approach. There is, however, no evolution mechanism in this artificial ecosystem. More works has been done on the Avida platform [2], which proposes self-replicating and evolving digital organisms. Each digital organism consists of a virtual CPU that processes a sequential program. The biological complexity in these organisms has been defined by Huang et al [3] as the generic information, which an organism has about its environment. However, none of these papers discuss about the complexity of the overall behavior of the simulation.

We are interested in analysing the complexity of such complex dynamic system. We have created a generic platform capable of simulating complex ecosystem with intelligent agents interacting and evolving in a large and dynamic environment [1]. This is the only simulation modeling the fact that agent behaviors affect evolution and speciation. The agents display very complex behavior using

Fuzzy Cognitive Map (FCM) model [19] to make decisions. We would like to understand how predictable the complex system we have conceived is. The two opposite and extreme situations that lead to a rather unpredictable system are random processes and chaotic processes. Therefore, we would like to investigate how close our complex system is to these two extremes.

The rest of the article is organized as follows. We first review our ecosystem simulation, which uses FCM as a behavior model in Section 2. Predicting population using the Markov chain is explained in Section 3. The Markov chain, transition matrix, and predicting prey and predator population is described in this section as well.  The Lyapanov exponent is described in Section 4, and experiments and results are shown in Section 5. Finally, in Section 6, we draw conclusions about this work and propose an extension to it.

## II.    INDIVIDUAL-BASED EVOLVING PREDATOR-PREY ECOSYSTEM SIMULATION USING FUZZY COGNITIVE MAPS AS A BEHAVIOR MODEL

In this section, the main parts of the already existing predator/prey ecosystem is briefly introduced.

### A.  Fuzzy Cognitive Maps

In general, fuzzy cognitive maps (FCMs) aim to represent the causal relationship between concepts, and to analyze inference patterns (the final states of the system after convergence).  Formally, an FCM is a graph which contains a set of nodes C, each node $C_i$ being a concept, and a set of edges I, each edge $I_{ij}$ representing the influence of the concept $C_i$ on the concept $C_j$ . A positive weight associated with the edge $I_{ij}$ corresponds to an excitation of the concept $C_j$ from the concept $C_i$, whereas a negative weight is related to an inhibition (a zero value indicates that there is no influence of $C_i$ on $C_j$). An activation level $a_i$ is associated to each concept. An FCM allows computing the new activation levels of the concepts of an agent, based on its perception and on the current activation levels of its concepts.

### B.  Agents and Behavior Model

The agents of this simulation are either prey or predators, which act in a dynamic environment with 1000×1000 cells. Each cell may contain several individuals and some amount of food. Each agent has several properties that determine its physical capabilities and its behaviors. The behaviors are

determined by the interaction between the FCM and the environment.

Each agent possesses its own FCM that represents its genome. We use an FCM to model an agent's behavior (structure of the graph) and to compute the next action of the agent (dynamic of the map). In each FCM, three kinds of concepts are defined: sensitive (such as distance to foe or food, amount of energy, etc), internal (fear, hunger, curiosity, satisfaction, etc) and motor (evasion, socialization, exploration, breeding, etc.). The activation level of a sensitive concept is computed by performing a fuzzification of the information the agent perceives in the environment. For an internal concept, the activation level corresponds to the intensity of an internal state of the agent. Note that it enables to distinguish between perception and sensation: the sensation is the real value coming from the environment, and the perception is the sensation modified by the internal states. Activation levels of the motor concepts are used to determine what the next action of the individual will be. The amplitude of the chosen action is then calculated by performing a defuzzification of the value of the corresponding motor concept.

The FCM of an agent is transmitted to its offspring after being combined with the one of the other parents, and after the possible addition of some mutations. The behavior model of each agent is therefore unique.

### C. Update

The time step represents a relatively long period of time, during which agents perform several small actions, which are summarized by a unique high level action. The possible high level actions for the agents are:

1. Evasion (for prey only), which is in the opposite direction of the closest foe within the vision range of the prey. The new position of the prey is computed using the speed of the prey.
2. Search for food, which is near the closest food (grass or meat) within the vision range.
3. Socialization, which is the direction toward the closest possible mate within the vision range.
4. Exploration in which the agent moves at its speed in the random direction.
5. Resting in which nothing happens.
6. Eating, which includes the update of the grass/meat unit in the cell and agents energy and hunger level.
7. Breeding: If the energy levels of both agents are more than a certain threshold, and they both choose the breeding action, then breeding is done.

For each action which requires the agent movement, its speed is computed proportional to the current activation level of the motor concept associated with its action.

At each time step, the values of the states of all the parameters in the model are updated. The three successive phases of the update process are as follows for all agent: Perception of the environment, computation of all its concepts, application of their action and update the energy level. Then some general updates for the whole world are performed such as updating the species and updating the amount of food available in each cell of the world.

An agent has a quite short lifespan (in terms of number of time steps), and performs only a few dozens of actions during its life. This enables us to obtain a high level of population renewal, which is an important criterion for studying an evolutionary process.

Fig. 1 shows the population of prey and predator agents after each time step. As expected with a predator-prey system, it is clear that there is a dependency between the number of prey and the number of predators. The evolution of the number of predators follows that of prey, and vice versa. As a clear period, consider time steps between 2000-4000. When the number of prey grows, the number of predators also grows a few time steps later. But when the number of predators grows too much, the number of preys decreases a few time steps later, leading to a still-later decrease in the number of predators.

### D. Evolution

In our simulation, evolution stems from several mechanisms: mating, mutation and speciation. Since species membership of the agents is evaluated at each time step, births and deaths of individuals influence the general species composition. Thus, a species can emerge or disappear at any time step. This enables us to model the evolution of populations of individuals sharing important genetic properties. Due to our species model, species evolution is derived directly from individual evolution. If the mating is successful, the two parents give birth to a unique offspring. This offspring inherits a combination of the genomic information of its parents, with possible mutations and crossover. The genome of an agent is defined as the set of edges, associated with their weights, of its FCM. More precisely, for each concept, the child inherits all the incident edges of this concept from one of its two parents. During that process, the weights of the edges can be modified, based on a probability of mutation (which is a parameter of the simulation). Moreover, some new edges can be created; and some old edges can be removed (if their weight becomes smaller than a given threshold).The apparition of new edges is a very important mechanism, in the sense that new influences between concepts can emerge during the evolutionary process. This allows the apparition of more complex and potentially more adaptive behaviors.
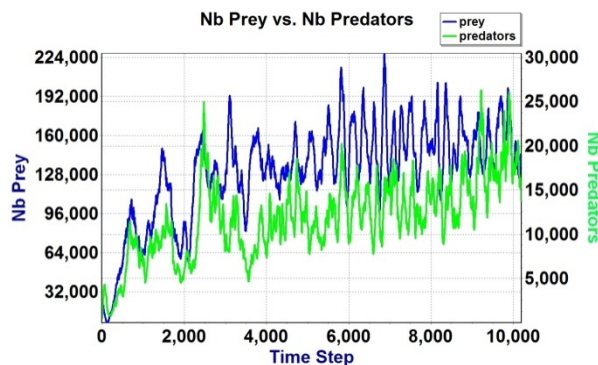


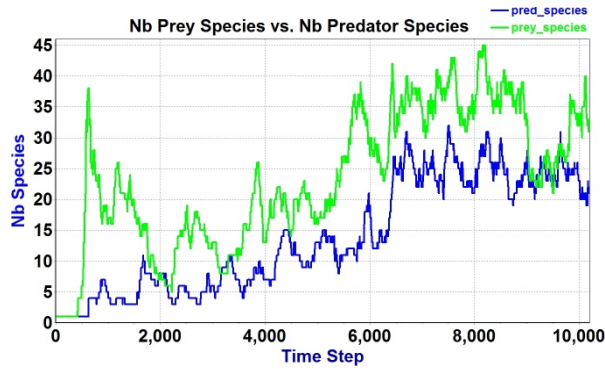Figure 1.    Population of prey and predator agents

Figure 2.    Evolution of the prey and predator species

In general, a single mutation is neutral and cannot produce a very different behavior model. It is the accumulation of neutral mutations during several time steps that allows the apparition of new individual behaviors, and then of new species.

Fig. 2 depicts the evolution of the prey and predator species in different time steps. Comparing Fig. 1 and Fig. 2, it appears that the number of species is closely correlated with the number of agents.

These are the kind of data that we would like to analyze. We consider several time series depicting the variation of different quantities in our system and define a protocol to evaluate the complexity level of our system by an analysis of these time series.

## III.    PREDICTING POPULATION USING THE MARKOV CHAIN

As it appears in Section 2, the agents in our simulation display a very complex behavior using FCM to make a decision. Moreover, the whole behavior of the ecosystem is a very complex system involving interaction between hundreds of thousands of complex agents. However, as it has been shown in [1], the overall system presents interesting correlation patterns.  For instance, the population of prey or predators has strong correlations with the number of predators, prey, level of grass and meat, species distribution, and so on. Despite such regularities, the simulation is far from being easily predictable. The amplitudes and times of inflation and deflation vary considerably, but their correlation is conserved. This means we can suppose that there is no easy way, excluding the simulation itself, to predict the state of the system at time step $t$, knowing the state of the system at time step $t-1$. However, the variations of the time series associated to these numbers are not random as they are the results of the application of the simulation mechanisms. We are interested in evaluating how complex the time series generated by the simulation are. We would like to understand how predictable the complex system we have conceived is. The two opposite and extreme situations that lead to a system, which is hardly predictable are random processes and chaotic processes. Therefore, we would like to investigate how close to these two extremes our complex system is. In that purpose, we have compared the data generated by our system with data generated by a simpler process that generates data similar to the ones of the

simulation. We used a Markov chain as a simple model of our simulation. We have learned the parameters of the corresponding Markov matrix from the data generated by our simulation. Then we used the Markov chain to generate data similar to those of the simulation.  We have then compared the time series data generated by both processes and measured their respective level of randomness/chaos.

### A.  Markov Chain

Define Suppose we generate a sequence of random variables, $\{X_0, X_1, X_2, ...\}$, such that at each time t, the next state $X_{t+1}$ depends only on the current state of the chain, $X_t$. This sequence is called a Markov chain, which is formally denoted as follows [7].

$$Pr(X_{t+1} = x_{t+1} \mid X_1 = x_1, ..., X_t = x_t) = Pr(X_{t+1} = x_{t+1} \mid X_t = x_t)$$

Now suppose we generate a sequence of random variables, $\{X_0, X_1, X_2, ...\}$, such that at each time t, the next state $X_{t+1}$ depends on the current state of the chain, $X_t$ and the previous state of the chain, $X_{t-1}$. This sequence is called a $2^{nd}$ order Markov chain [7].

$$Pr(X_{t+1} = x_{t+1} \mid X_1 = x_1, ..., X_t = x_t) = Pr(X_{t+1} = x_{t+1} \mid X_t = x_t,$$
$$X_{t-1} = x_{t-1})$$

Similarly the $3^{rd}$ order Markov chain is a sequence satisfying

$$Pr(X_{t+1} = x_{t+1} \mid X_1 = x_1, ..., X_t = x_t) = Pr(X_{t+1} = x_{t+1} \mid X_t = x_t,$$
$$X_{t-1} = x_{t-1}, X_{t-2} = x_{t-2})$$

### B.  Transition Matrix

Consider a Markov chain with finite state space $\{1, 2, ..., k\}$. A transition matrix describes the probability of moving form state i to j at each time step using $Pr(j|i) = P_{i,j}$; in other words, the ith  row and jth  column element of the transition matrix P is given by $Pr(j|i)$ [6].

$$P = \begin{bmatrix} P_{1,1} & \cdots & P_{1,k} \\ \vdots & \ddots & \vdots \\ P_{k,1} & \cdots & P_{k,k} \end{bmatrix}$$

where the probabilities of each row sum up to 1.

$$\sum_j P_{i,j} = 1$$

because the overall probability of transitioning from state i to one of all possible states must be 1.

### C.  Markov Chain for Prey/Predator

To show the non-random behavior of the simulation, we have generated the transition matrix for the population of

prey and predator, and used the matrix probabilities to obtain the artificial population. The overall algorithm includes the three following steps:

1) Pre-processing
2) Transition matrix computation
3) Running the Markov chain

The first step includes the smoothing and quantization of the data set. The first part, smoothing, is performed on the dataset by the following linear transformation, which smoothes the original data while preserving its general characteristics and trend.

$$X_t = \frac{10x_t + 5x_{t+1} + 5x_{t-1} + 2x_{t+2} + 2x_{t-2} + x_{t+3} + x_{t-3}}{26}$$

where $x_t$ is prey/predator population size at time t.

To build the transition matrix for the Markov chain, we need to have a finite state space, which means the number of the quantities should be constant and discrete. The interval for the quantization is computed by:

$$Quant\_Interval = floor\left(\frac{(max(X) - min(X))}{k}\right)$$

where $max(X)$ and $min(X)$ are the maximum and minimum values seen in the prey or predator datasets, correspondingly, and $k$ is the number of quantities, which is equal to the number of transition matrix rows in the 1st order Markov chain. The value of $k$ is selected small enough to ensure the nonzero value for quantization interval. The smoothed values are transformed to the quantized data set by the following formula to have a minimum value of 1.

$$X'_t = floor\left(\frac{(X_t - min(X))}{Quant\_Interval}\right) + 1$$

where $X'_t$ is the quantized value of the smoothed value at time *t*.

Transition matrix computation depends on the order of the Markov chain. The transition matrix dimension is equal to the order of the Markov chain plus 1. In the 1st order Markov chain, the transition matrix is 2D, which includes the population at the current time and the interested population at the next time. In the 2nd order Markov chain, the matrix has another dimension, which is the previous population of prey or predator.

Each probability for this matrix is computed by counting how many times the same consecutive population sizes happening in the quantized data set.

By running the Markov chain from an initial value (or several previous values for higher order Markov chain), the artificial data set is created; the next population for each time instance is computed based on the current and previous populations (for the higher order of the Markov chain), and the corresponding probabilities in the transition Matrix.

## IV. LYAPUNOV EXPONENT

Nonlinear signal processing is an important research area with many applications. Specifications and identifications of nonlinear signals can help us to detect nonlinear behavior of dynamical systems [8]. One specification, the discrimination of stochastic and chaotic behaviors of nonlinear time series, is a basic topic in nonlinear dynamic fields [9]. This specification has attracted researchers for a long time [10,12]. As many scientists believe that the natural phenomena have to be considered as deterministic and chaotic systems, it is important that a simulation used to model such a phenomenon generate a complex chaotic pattern [11]. For this reason fractal dimensions and Lyapunov exponents are the most prominent candidates to characterize the chaotic behavior [13], because they express complexity and predictability of a process and are a measure for chaos [14, 15, 16]. In this paper the Lyapunov exponent has been used.

Most experts would agree that chaos is the aperiodic, long-term behavior of a bounded, deterministic system that demonstrates sensitive dependence on initial conditions. For that purpose, we must quantify the sensitivity [17].

Lyapunov exponents quantify the exponential divergence of initially close state-space trajectories and estimate the amount of chaos in a system [18]. A bounded dynamic system with a positive Lyapunov exponent is chaotic [17].

Imagine two nearby initial points $X_0$ and $X_0 + \Delta X_0$, respectively. After one iteration of the map, the points are separated by

$$\Delta X_1 = f(X_0 + \Delta X_0) - f(X_0) \cong \Delta X_0 \, \acute{f}(X_0)$$

where $\acute{f} = df/dX$. Now, we define the local Lyapunov exponent $\lambda$ at $X_0$ such that $e^\lambda = |\Delta X_1/\Delta X_0|$, or

$$\lambda = \ln|\Delta X_1/\Delta X_0| = \ln|\acute{f}(X_0)|$$

To obtain the largest Lyapunov exponent, we average the above equation over large enough iterations.

$$\lambda = \lim_{N \to \infty} \frac{1}{N} \sum_{n=0}^{N-1} \ln|\acute{f}(X_n)|$$

The largest Lyapunov exponent determines the average exponential rate of separation of two nearby initial conditions, or the average expansion of the space. A positive value shows chaos [17].

The different methods that have been proposed for computing Lyapunov exponents from time series can be divided into two classes: Jacobian-based methods and direct methods.

Direct methods directly estimate the divergent motion of the reconstructed states without fitting a model to the data [20]. The method, which has been used in this paper was proposed by Sato et al. [21], and Kurths and Herzel [22]. The average exponential growth of the distance of neighboring orbits is studied on a logarithmic scale, this time via the prediction error below

$$p(k) = \frac{1}{Nt_s} \sum_{n=1}^{N} \log_2 \left( \frac{\|y^{n+k} - y^{nn+k}\|}{\|y^n - y^{nn}\|} \right)$$

where $y^{nn}$ is the nearest neighbor of $y^n$. The dependence of the prediction error $p(k)$ on the number of time steps $k$ may be divided into three phases [19]. Phase 1 is the transient where the neighboring orbit converges to the direction corresponding to the largest Lyapunov exponent. During phase 2 the distance growths exponentially until it exceeds the range of validity of the linear approximation of the flow. Then phase 3 begins where the distance increases more slowly than exponentially until it decreases again because of folding in the state space. In phase 2, a linear segment with slope $\lambda_1$ appears in the p(k) vs. k diagram. This allows an estimation of the largest Lyapunov exponent $\lambda_1$ [22]. Fig. 3 gives an example to determine the largest Lyapunov exponent $\lambda_1$ of data by this method [19].

## V. EXPERIMENTS AND RESULTS

The simulation is implemented in C++ and has been run using the Narwhal cluster on the Sharcnet system and produced 32500 time steps. The resulting prey/predator population has been used as the input data for the Markov analysis explained in Section 3. The Markov chain analysis is implemented in Matlab 7.1 and is run on the AMD dual core processor 3.00 GHz with 3.00 GB RAM.

The analysis is performed on predator/prey's population starting from time step 10,000. This time was provided to ascertain that the simulation reaches its stabilization. The smoothing and quantization is performed on each of the prey and predator datasets, and 40 smoothed and quantized values are obtained. In Fig. 4 the time series analysis for prey dynamics is shown. This figure demonstrates the changes in the size of population in 10 time steps, in which x-axis represents the population value at time t-10, and y-axis represents the corresponding change in the size of the population at time t. Values between 1 and 40 in prey population corresponds to the smoothed and quantized values explained in Section 3.C. Specifically values 1 and 40 are the minimum and maximum values in this dataset, which are 44521 and 219007 respectively. Also each unit in y axis corresponds to the quantization interval given in the formula of Section 3.3, which is 4474 having k as 39.

As a simple example, consider the minimum value, 1, in prey population at time t-10, the changes in the prey population at the time t (10 time steps after), according to the Fig. 4 is 0 and 1 unites. In other words, if the value of smoothed and quantized prey population was 1, then the value of smoothed and quantized prey population 10 time steps after that would either remain unchanged or increase to 1 unit. As it can be seen form Fig. 4, prey dynamics has many variations around 0, meaning that in many time steps prey population size has remained the same in the next 10 time steps.
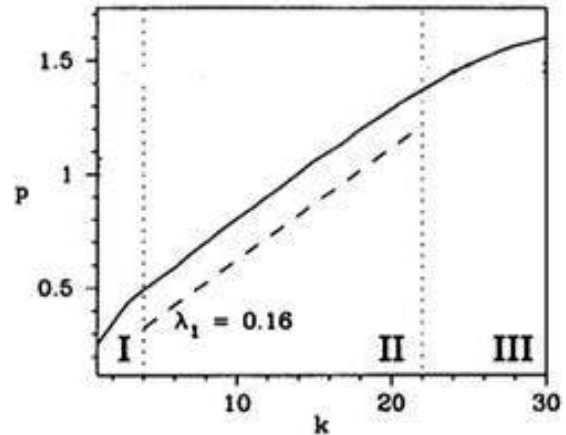


Figure 3. Prediction error p for experimental data vs. the number of time steps k. the slope of the solid line in the intermediate range of k gives the largest Lyapunov exponent $\lambda_1 = 0.16$.



Figure 4. Time series analysis for prey dynamics

The lines in the graph represent transition between successive population states encountered during the simulation process.

Based on the values after preprocessing stage, frequency matrices and transition matrices are obtained. Markov runs are performed using these transition matrices to estimate the prey and predator populations starting with the first 10 values of the population of prey or predator in the simulation. The rest of the population values are estimated by applying the Markov run with the transition matrix. The Markov run for prey dynamics is shown in Fig. 5.

By comparing Fig. 4, which is the prey population dynamics in the simulation, and Fig. 5, which is the estimated prey population given the first 10 values of the simulation prey population, two important conclusions can be made. As expected the pairs of coordinates in Fig. 5 are the subset of pairs of coordinates in Fig. 4 meaning that some of the population pairs (population at time t-10 and its change at time t) are not generated during this Markov run.

Figure 5.    Markov run for prey dynamics

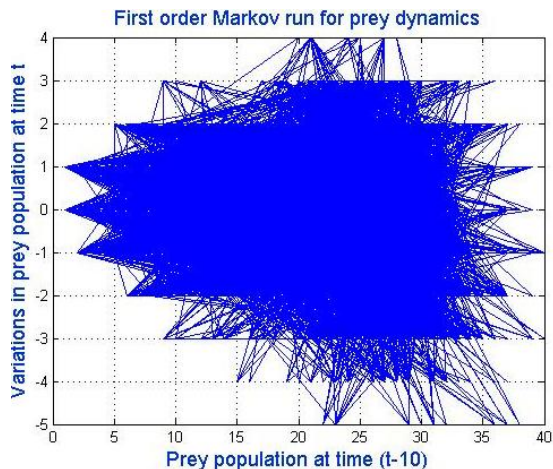For instance, in Fig. 5, unlike Fig. 4, the pair (5, -2) did not appear meaning that if the population size at time t-10 was 5, the population size at time t would never be 3.

On the other hand, although the estimated values depicted in Fig. 5 are obtained using the transition matrix, which is based on the probabilities of values depicted in Fig. 4, the sequence of estimated values in Fig. 5 are clearly more random than the corresponding sequence represented in Fig. 4.

Higher order Markov chains are also implemented and the results are shown in Fig. 6 and Fig. 7.

Comparing Fig. 5 to 7 shows that the Markov chain for all of the order we have tested, generated prey dynamics that are more random than the original simulation prey dynamics. Note that although the average probabilities of having a certain number in the population size at the next interval increase by the order of the Markov chain, the average frequencies and zero rows percentage of the transition matrices are decreased (see Table 1). This is why the difference between different orders of Markov chain runs in population dynamics are not significant. A row in transition matrix with all zero values, indicate an undefined value for the next population interval. This row is created only in the 2nd or higher order Markov chains because it corresponds to the population sequence (sequence of 2 and 3 successive population values for 2nd and 3rd order Markov chain respectively), which was never appeared in the original population time series.

The same processes have been performed on predator time series and similar results have been obtained (Table 1).

To show the chaotic behavior observed in the population of prey and predator time series, the largest Lyapunov exponent values after different modifications of dataset are computed. In Table 2, the values of the largest Lyapunov exponent over simulation's prey and predator population data, smoothed time series and first, second and third order of the Markov run for the prey/predator population time series are presented.



Figure 6.    Second order Markov run for prey dynamics



Figure 7.    Third order Markov run for prey dynamics

The same processes have been performed on predator time series and similar results have been obtained (Table 1).

To show the chaotic behavior observed in the population of prey and predator time series, the largest Lyapunov exponent values after different modifications of dataset are computed. In Table 2, the values of the largest Lyapunov exponent over simulation's prey and predator population data, smoothed time series and first, second and third order of the Markov run for the prey/predator population time series are presented.

According to Table 2, we can conclude that the population time series, which has been produced by the simulation, indicate chaotic behavior because the largest Lyapunov exponent values for population time series are greater than zero. Obviously after the smoothing process, the largest Lyapunov exponent value is higher because smoothing removes random behavior. As we expected, the largest Lyapunov exponent of first, second and third order

TABLE I.    DIFFERENCE BETWEEN MARKOV CHAIN ORDER VALUES IN TRANSITION MATRIX

| | | Zero Rows of Matrix | Average Probability of sequence | Average Frequency of sequence |
|---|---|---|---|---|
| Prey | 1st-order Markov chain | 0 | 0.17 | 101.67 |
| | 2nd-order Markov chain | 0.14 | 0.27 | 27.83 |
| | 3rd-order Markov chain | 0.013 | 0.37 | 10.18 |
| Predator | 1st-order Markov chain | 0 | 0.14 | 79.54 |
| | 2nd-order Markov chain | 0.18 | 0.22 | 17.28 |
| | 3rd-order Markov chain | 0.021 | 0.34 | 5.9 |

TABLE II.    LARGEST LYAPUNOV EXPONENT VALUES FOR DIFFERENT TIME SERIES

| Data (time series) | Largest Lyapunov Exponent for prey | Largest Lyapunov Exponent for predator |
|---|---|---|
| Real data | 0.5208 | 0.4862 |
| Smoothed data | 0.5603 | 0.5132 |
| 1st-order Markov chain | 0.0117 | 0 |
| 2nd-order Markov chain | 0 | 0 |
| 3rd-order Markov chain | 0 | 0 |

Markov run for population time series is almost zero due to the random behavior of the Markov chain.

## VI.    CONCLUSION AND FUTURE WORK

We have conceived a protocol to evaluate the complexity level of our agent-based ecosystem simulation. The variations of the time series associated with the ecosystem simulation are the results of complex simulation mechanisms. To understand how close our system is to the random or chaotic processes, we have compared the data generated by our system with data generated by a Markov chain as a simple process.

As explained in Section 3, the parameters of the corresponding Markov matrix have been learned from the data generated by our simulation. Then we used the Markov chain to generate data similar to those of the simulation. As shown in Section 5 the Markov chain for all of the order we have tested, generates prey and predator time series that are more random than their counterpart in the original simulation.

We also used the largest Lyapunov exponent to determine the chaotic behavior of the simulation. Our experiments show that the largest Lyapunov exponent values for population time series of both prey and predator are positive, indicating a chaotic behavior in our ecosystem simulation, which is a good indication of a high complexity level.

As it appears in Section 4, the most prominent candidates to characterize chaotic behavior in a system are Fractal dimensions and Lyapunov exponents [13]. As the next step,

Fractal dimensions can also be applied on the time series produced by our ecosystem simulation to determine the complexity and predictability of the system and measure the chaotic behavior. We would like also to analyze other time series, depicting variations of other quantities like quantities of food or number of species, to have a better understanding of the overall complexity of our system.

## REFERENCES

[1]  R. Gras, D. Devaurs, A. Wozniak, and A. Aspinall, "An Individual-based Evolving Predator-Prey Ecosystem Simulation using Fuzzy Cognitive Map as Behavior Model", Journal of Artificial Life, Number 4, Volume 15(4), pp. 423-463, 2009.

[2]  C. Ofria and C. Wilke. Avida, "A Software Platform for Research in Computational Evolutionary Biology", Journal of Artificial Life, Volume 10, pp. 191-229, 2004.

[3]  W. Huang, C. Ofria, and E. Torng, "Measuring Biological Complexity in Digital Organisms", Ninth International Conference on Artificial Life , Boston MA, Sept 12-15, 315-321, 2004.

[4]  C. R. Ward, F. Gobet, and G. Kendall, "Evolving Collective Behavior in an Artificial Ecology", Journal of Artificial Life, Volume 7(2), pp. 191–209, 2001.

[5]  M. Ronkko, "An artificial Ecosystem: Emergent Dynamics and Lifelike Properties", Journal of Artificial Life, Volume 13(2), pp. 159–187, 2007.

[6]  G. Latouche and V. Ramaswami, "Introduction to Matrix Analytic Methods in Stochastic Modelling", 1st edition, "Chapter 2: PH Distributions", Philadelphia, PA: ASA SIAM", 1999.

[7]  W. R. Gilks, Walter R. Gilks, Sylvia Richardson, and D. J. Spiegelhalter, "Markov Chain Monte Carlo in Practice", 1st edition, Chapman and Hall/CRC, 1996.

[8]  J. Hubbard and B. West, "Differential Equations: A Dynamical Systems Approach: Ordinary Differential Equations", Volume 5, Texts in Applied Mathematics. Springer, 1991.

[9]  Peter Grassberger, Thomas Schreiber, and Carsten Schaffrath, "Nonlinear Time Sequence Analysis", International Journal of Bifurcation and Chaos, Volume 1(3), pp 521-547, 1991.

[10]  Amir H. Omidvarnia and Ali Nasrabadi, "A New Irregularity Criterion for Discrimination of Stochastic and Deterministic Time Series", Fractals, Volume 16(2), pp. 1–12, 2008.

[11]  A. Golestani, M. R. Jahed Motlagh, K. Ahmadian, Amir H. Omidvarnia, and Nasser Mozayani, "A New Criterion for Distinguish Stochastic and Deterministic Time Series with the Poincaré Section and Fractal Dimension", Journal of Chaos: An Interdisciplinary journal of Nonlinear Science, CHAOS 19, Volume 19(1), pp. 1-13, March 2009.

[12]  L. Romanelli, M. A. Figliola, and F. A. Hirsch, Deterministic Chaos and Natural Phenomena. Journal of Statistical Physics, *Vol. 53, Nos. 3/4,* 1988.

[13]  J. Holzfuss and U. Parlitz, "Lyapunov Exponents from Time Series", *Lecture Notes in Mathematics*, Springer, Berlin, 1990.

[14]  J.P. Eckmann and D. Ruelle, "Ergodic Theory of Chaos and Strange Attractors", Reviews of Modern Physics, Volume 57(3), pp. 617-656, 1985.

[15]  W. Lauterborn and J. Holzfuss, "Evidence for a Low-Dimensional Strange Attractor in Acoustic Turbulence", Physics Letters A, Volume 115(8), pp. 369-372, 1986.

[16] W. Lauterborn and J. Holzfuss, "Acoustic Chaos", International Journal of Bifurcation and Chaos, Volume 1(1), pp. 13-26, 1991.

[17] Julien Clinton Sprott, "Chaos and Time Series Analysis", Oxford University Press, 2003.

[18] Michael T. Rosenstein, James J. Collins, and Carlo J. De Luca, "A Practical Method for Calculating Largest Lyapunov Exponents from Small Data Sets", Journal of Physica D: Nonlinear Phenomena, Volume(65), pp. 117-134, 1993.

[19] B. Kosko, "Fuzzy cognitive maps", International Journal of Man-Machine Studies, 24, 65–75, 1986.

[20] U. Parlitz, "Nonlinear Time-Series Analysis", Nonlinear Modeling - Advanced Black-Box Techniques Eds. J.A.K. Suykens and J. Vandewalle Kluwer Academic Publishers, pp. 209-239, 1998.

[21] Sato, S., M.Sano ,and Y. Sawada, "Practical methods of measuring the generalized dimension and largest Lyapunov exponent in high dimensional chaotic systems", Progress of Theoretical Physics, 77, pp. 1-5, 1987.

[22] Kurths, J. and H. Herzel, "An attractor in solar time series", Physica D, 25, pp. 165-172, 1987.

# A Knowledge Development Conception and its Implementation:
# Knowledge Ontology, Rule System and Application Scenarios

Eckhard Ammann

School of Informatics
Reutlingen University, Reutlingen, Germany
Eckhard.Ammann@Reutlingen-University.de

Manuela Ruiz-Montiel, Ismael Navas-Delgado
and José F. Aldana-Montes
E.T.S.I. Informática
University of Málaga, 29071 Málaga, Spain
{mruizmonti, ismael, jfam}@lcc.uma.es

*Abstract--* **Knowledge development in an enterprise is about approaches, methods, techniques and tools, that will support the advancement of individual and organizational knowledge for the purpose of an improvement of businesses. An approach for knowledge development in a company is described in this paper. This approach is based on a new conception of knowledge, with the introduction of three knowledge dimensions and conversions between knowledge assets. This conception is implemented in the form of a knowledge ontology. Thus, we can take advantage of reasoning and rules processing, provided by a reasoner in combination with a rule engine. Important scenarios for knowledge development in a company are identified and it is shown how these scenarios can be supported by processing the developed rules. For example, knowledge requirements for a new or existing employee can be gained once the appropriate requirements for a planned project are known as well as the learning options in the company.**

*Keywords-- Conception of knowledge, knowledge development, knowledge ontology, rule system, application scenarios.*

## I. INTRODUCTION

Knowledge development in an enterprise is about approaches, methods, techniques and tools, that will support the advancement of knowledge for the purpose of an improvement of businesses. This notion includes as well individual knowledge as group and organizational knowledge. It can be seen as integral part of knowledge management; see [1], [2] and [3] for a description of several existing approaches for knowledge management. One specific approach for enterprise knowledge development is EKD (Enterprise Knowledge Development), which aims at articulating, modeling and reasoning about knowledge, which supports the process of analyzing, planning, designing, and changing your business; see [4] and [5] for a description of EKD. EKD does not provide a conceptual description of knowledge and knowledge development.

In this paper, we present a new conception of knowledge and knowledge development and describe an implementation of this conception based on a knowledge ontology, reasoning support and a rule system.

For the conception part, there exists one well-known approach by Nonaka/Takeuchi [6], which is built on the distinction between tacit and explicit knowledge and on four knowledge conversions between the knowledge types

(SECI-model). However, is explicit knowledge still bound to the human being, or already detached from him? Also the linear spiral model of knowledge development is limiting.

An approach for knowledge access and development in firms is given by Boisot [7]. Here, development scenarios of knowledge in the Information Space are provided.

Our conception of knowledge is represented by a three-dimensional model of knowledge with types, kinds and qualities. General knowledge conversions between the various knowledge variants are introduced as a model for knowledge dynamics in the enterprise. First a basic set of such conversions is defined. Building on this set, general knowledge conversions can be defined, which reflect knowledge transfers and development and do not suffer from the restrictions of the SECI-model.

A knowledge ontology is described in this paper, which implements this conception of knowledge and knowledge conversions. It has been developed in the web ontology language OWL [8]. The reasoning support in combination with a rule system allows for a formal treatment of important knowledge development scenarios.

Application scenarios for knowledge development are classified and described in this paper. They can be represented by general knowledge conversions, which are subject to rule processing. A set of corresponding rules for addressing these scenarios and their representations has been developed and is described in this paper. Therefore, possible solutions for those scenarios can be gained.

## II. A CONCEPTION OF KNOWLEDGE AND KNOWLEDGE DYNAMICS

In this section, a conception of knowledge and knowledge dynamics in a company is described. More details of this conception are given in [2].

### A. Knowledge Conception

We provide a conception of knowledge with types, kinds and qualities. As our base notion, knowledge is understood as justified true belief (in the propositional kind), which is (normally) bound to the human being, with a dimension of purpose and intent, identifying patterns in its validity scope,

brought to bear in action and with a generative capability of new information, see [3] and [9]. It is a perspective of "knowledge-in-use" [10] because of the importance for its utilization in companies and for knowledge management. In contrast, information is understood as data in relation with a semantic dimension, but without the pragmatic and pattern-oriented dimension, which characterizes knowledge.

*1) Type Dimension of Knowledge*

The type dimension is the most important for knowledge management in a company. It categorizes knowledge according to its presence and availability. Is it only available for the owning human being, or can it be communicated, applied or transferred to the outside, or is it externally available in the company's organizational memory? It is crucial for the purposes of the company, and hence a main goal of knowledge management activities, to make as much as possible knowledge available, i.e. let it be converted from internal to more external types.

Our conception for the type dimension of knowledge follows a distinction between the internal and external knowledge types, seen from the perspective of the human being. As third and intermediary type, explicit knowledge is seen as an interface for human interaction and for the purpose of knowledge externalization, the latter one ending up in external knowledge. Internal (or implicit) knowledge is bound to the human being. It is all that, what a person has "in its brain" due to experience, history, activities and learning. Explicit knowledge is "made explicit" to the outside world, e.g., through spoken language, but is still bound to the human being. External knowledge finally is detached from the human being and may be kept in appropriate storage media as part of the organizational memory. Figure 1 depicts the different knowledge types.

Internal knowledge can be further divided into tacit, latent and conscious knowledge, where those subtypes do partly overlap with each other; see [9]. Conscious knowledge is conscious and intentional, is cognitively available and may be made explicit easily. Latent knowledge has been typically learning as a by-product and is not available consciously. It may be made explicit, for example in situations, which are similar to the original learning situation, however. Tacit knowledge is built up through experiences and (cultural) socialization situations, is specific in its context and based on intuition and perception.

*2) Kind Dimension of Knowledge*

In the second dimension of knowledge, four kinds of knowledge are distinguished: propositional, procedural and strategic knowledge, and familiarity. It resembles to a certain degree the type dimension as described in [10]. Propositional knowledge is knowledge about content, facts in a domain, semantic interrelationship and theories. Experience, practical knowledge and the knowledge on "how-to-do" constitute procedural knowledge. Strategic
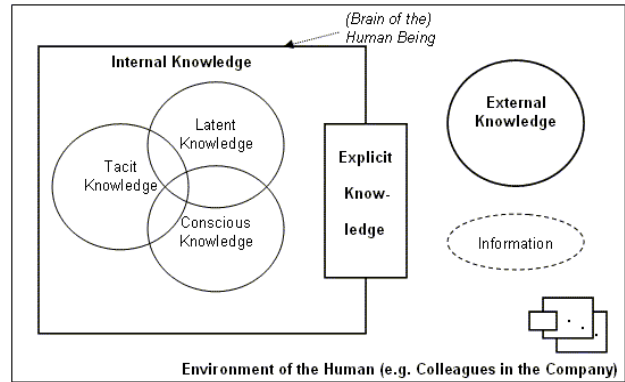


Figure 1.  Conception of knowledge types

knowledge is meta-cognitive knowledge on optimal strategies for structuring a problem-solving approach. Finally, familiarity is acquaintance with certain situations and environments; it also resembles aspects of situational knowledge, i.e. knowledge about situations, which typically appear in particular domains.

*3) Quality Dimension of Knowledge*

The quality dimension introduces five characteristics of knowledge with an appropriate qualifying and is independent of the kind dimension; see [10]. The level characteristics aims at overview vs. deep knowledge, structure distinguishes isolated from structured knowledge. The automation characteristic of knowledge can be step-by-step-doing by a beginner in a domain of work or automated fast acting by an expert.

Modality as the fourth quality of knowledge asks for the representation of it, be it words versus pictures in situational knowledge kinds, or propositions versus pictures in procedural knowledge kinds. Finally, generality differentiates general versus domain-specific knowledge. Knowledge qualities apply to each knowledge asset.

*4) The Knowledge Cube*

Bringing all three dimensions of knowledge together, we gain an overall picture of our knowledge conception. It can be represented by the knowledge cube as shown in Figure 2.
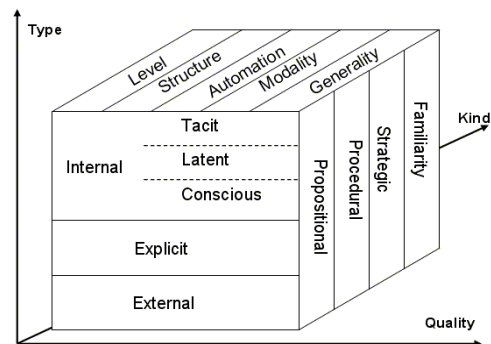


Figure 2.  The knowledge cube

Note, that the dimensions in the knowledge cube behave different. In the type and kind dimensions, the categories are mostly distinctive (with the mentioned exception in the sub-types), while in the quality dimension each of the given five characteristics are always present for each knowledge asset.

### B. Knowledge Dynamics

Here we give a conception of knowledge conversions. The transitions between the different knowledge types, kind and qualities are responsible to a high degree for knowledge development in an organization. These general knowledge conversions are the building blocks to model knowledge dynamics, i.e., all of acquisition, conversion, transfer, development and usage of knowledge, in an enterprise.

Most important for knowledge management purposes are conversions between the knowledge types, especially those making individual and internal knowledge of employees usable for a company. The explicitation and externalization conversions described in this section achieve this. Implicitly, socializations between tacit knowledge of different people also contribute to this goal.

### 1) Basic Knowledge Conversions

Five basic knowledge conversions in the type dimension are distinguished here: socialization, explicitation, externalization, internalization and combination. Basic conversion means, that exactly one source knowledge asset is converted into exactly one destination knowledge asset and exactly one knowledge dimension (i.e. the type dimension in this case) is changed. More complex conversions may be easily gained by building on this set as described in the next sub-section. They will consist of n-to-m-conversions and include information assets in addition.

Socialization converts tacit knowledge of a person into tacit knowledge of another person. This may succeed by exchange of experience or in a learning-by-doing situation. Explicitation is the internal process of a person, to make internal knowledge of the latent or conscious type explicit, e.g. by articulation and formulation (in the conscious case) or by using metaphors, analogies and models (in the latent case). Externalization converts from explicit knowledge to external knowledge or information and leads to detached knowledge as seen from the perspective of the human being, which can be kept in organizational memory systems. Internalization converts either external or explicit knowledge into internal knowledge of the conscious or latent types. It leads to an integration of experiences and competences in your own mental model. Finally, combination combines existing explicit or external knowledge in new forms.

Basic knowledge conversions in the kind dimension of knowledge do not occur. Those in the quality dimension are mostly knowledge developments aiming at quality improvement. Examples are basic conversions changing the overview, structure and automation quality, respectively.

### 2) General Knowledge Conversions

Our conception allows the generalization of the basic five knowledge conversions described above. General knowledge conversions are modeled converting several source assets (possibly of different types, kinds and quality) to several destination assets (also possibly different in their knowledge dimensions). In addition, information assets are considered as possible contributing or generated parts of general knowledge conversions.

For example, in a supervised learning-by-doing situation seen as a complex knowledge conversion, a new employee may extend his tacit and conscious knowledge by working on and extending external knowledge in a general conversion, being assisted by the tacit and conscious knowledge of an experienced colleague. As a result of the conversion we have extended internal knowledge of the new employee and extended external knowledge.

## III. THE KNOWLEDGE ONTOLOGY

In this section we present the Knowledge Ontology, which implements the conception of knowledge and knowledge dynamics as described in Section 2. As one main goal the ontology will enable the discovery of the crucial knowledge conversions for a company. The ontology (as visually shown in Figure 3) is divided in four core concepts: *Knowledge*, *Information*, *Knowledge_Conversion* and *Knowledge_Dimension*. The three different knowledge dimensions are represented as: *Type_Dimension, Kind_Dimension* and *Quality-Dimension*. *Knowledge* is defined according to these dimensions. Properties are used to model the relationships between *Knowledge* and *Dimensions*: *hasType*, *hasKind* and *hasQuality*. For example, *Explicit_Knowledge* is defined as every piece of knowledge, which is related to the instance *Explicit_Type* via the *hasType* property. In the same way, *Knowledge* in general must be related to every quality sub-dimension through the *hasQuality* property.

In the case of the type dimension of knowledge, we have defined disjoint axioms in order to make explicit the fact that a piece of knowledge cannot be simultaneously external and internal - except in the case of *Latent*, *Conscious* and *Tacit Knowledge*, which can actually overlap (compare with Figure 1). There are also disjoint axioms for the kind dimension, since a propositional piece of knowledge cannot be *Procedural*, neither *Strategic* nor *Familiarity*.

Two properties have been defined to model the knowledge conversions: *hasSource* and *hasDestination*, with knowledge conversions as ranges, and pieces of knowledge and information as domains.

A General Conversion is modeled through the *Knowledge Conversion* concept, and its only restriction the fact that it must have at least one source asset and one destination asset. *Basic Conversions* are more specific, in the sense that they have only one source and only one destination. Eight basic conversions (five in the type
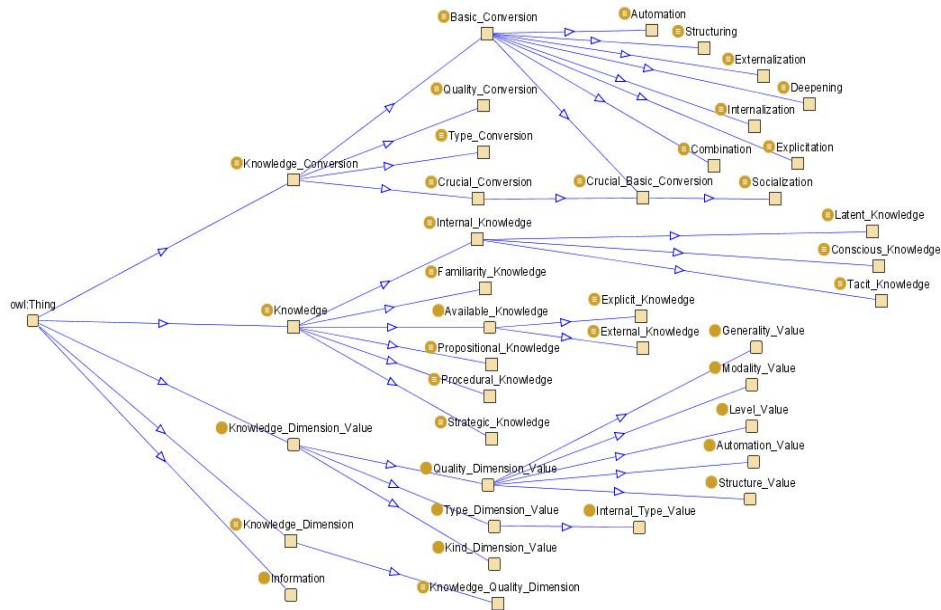
Figure 3. Knowledge ontology hierarchy

dimension, three in the quality dimension) are defined in the ontology. The concept *Crucial_Conversion* gathers those conversions that contribute to the goal of making the knowledge available for the company.

### A. Restrictions and Reasoning

Basic reasoning is based on subsumption mechanisms that deal with the ontology hierarchy. However, ontologies can contain more complex elements to enable advanced reasoning. In this way, the Knowledge Ontology has been extended with OWL restrictions to enable new ways of generating interesting new knowledge.

Here we will only describe some of the most interesting restrictions. Let us imagine that we have two pieces of knowledge in our company: *knowledge1* and *knowledge2*. Both pieces of knowledge have as type *Explicit* (is related to the instance of *Type_Dimension_Value* called *Explicit* through the property *hasTypeValue*). Additionally we have defined *Explicit_Knowledge* as follows:

Available_Knowledge AND
∃ hasTypeValue has Explicit

Thus, a reasoner will identify both pieces of knowledge as *Explicit_Knowledge* (and using subsumption also as *Available_Knowledge*).

We can consider two different conversions *conversion1* and *conversion2*: one that converts *knowledge1* in *knowledge2* and vice versa. Then, we have defined a *Crucial_Conversion* as:

Knowledge_Conversion AND
∃ hasDestination some Available_Knowledge

Thus, we can infer that *conversion1* is a *Crucial_Conversion* for the company.

### B. Rules

Ontology restrictions allow us to infer new characteristics of a given concept or instance. However, in some cases we could require to generate new instances in the ontology depending on certain situations. In this case we have used rules, so the knowledge ontology will be able to infer all the possible conversions given some pieces of knowledge. First, the rule engine will create basic conversions with all the possible source-destination pairs, and then, the same engine will characterize these conversions, inferring the changing dimension for each case.

SWRL [11] rules have been defined and the Jess rule engine [12] has been used for testing purposes. The main rule for our model is the one that creates new conversions for the knowledge assets that we have stored in our ontology:

Knowledge(?k1) ^ Knowledge(?k2) ^
hasDimensionValue(?k1, ?v1) ^
hasDimensionValue(?k2, ?v2) ^
differentFrom(?k1, ?k2) ^ differentFrom(?v1, ?v2) ^
swrlx:makeOWLThing(?c, ?k1, ?k2)
→
Knowledge_Conversion(?c) ^ hasSource(?c, ?k1) ^
hasDestination(?c, ?k2)

Thus, this rule is activated when we have two different pieces of knowledge with different dimensions values. In this case, a new instance is created for providing a new knowledge conversion between both pieces of knowledge.

Then, we have six rules to infer the changing dimensions of each of the new discovered conversions: one for the type

dimension and five for the quality ones. For example, the rule for the type dimension is as follows:

Knowledge(?k1) ^ Knowledge(?k2) ^
hasTypeValue(?k1, ?v1)  ^ hasTypeValue(?k2, ?v2) ^
differentFrom(?v1, ?v2) ^ Knowledge_Conversion(?c1) ^
hasSource(?c, ?k1) ^ hasDestination(?c, ?k2)
→
hasChangingDimension(?c,
              Knowledge_Type_Dimension)

Suppose that we have two pieces of knowledge in our company (*knowledge1* and *knowledge2)*, which are related through the *hasTypeValue* property to *Explicit* and *External*, respectively. Both are related to the values *Familiar* and *Step by step*. Using the defined rules, new instances are produced. Thus, the rule engine has inferred two conversions, one for "knowledge1 → knowledge2", and another for "knowledge2 → knowledge1". Then, the reasoner can infer additional facts:

- About the pieces of knowledge:
    - They are both *Familiar_Knowledge*.
    - One of them is *External_Knowledge,* the other is *Explicit_Knowledge*.
    - Both are *Available_Knowledge*.
- About the conversions:
    - They are both *Basic_Conversion*.
    - Both are *Crucial_Knowledge* (since they have *Available_Knowledge* as destination).
    - Both are *Type Conversions* (since they change the type dimension).

## IV.  APPLICATION SCENARIOS

Application scenarios for knowledge development in a company can be related with our model of knowledge dynamics. Two categories of scenarios exist. The first one is constructive and builds knowledge development chains (see [2] for a modeling approach). Here we focus on the second scenario category, which consists of analytic scenarios. They can be represented by general knowledge conversions and are subject to rule processing as described in Section 3. In these scenarios we face gaps in knowledge dynamics chains as provided by knowledge conversions. These gaps will be closed by applying appropriate rules to the relevant instances of knowledge assets and conversions, which have been instantiated in our knowledge ontology.
Figure 4 explains our approach. The bold arrow in the first line indicates the knowledge development activity, which is needed in order to resolve an application scenario with unknown part. Our approach first represents the application scenario as a general knowledge conversion, applies an appropriate rule of our rule system to it, and finally interprets the completed knowledge conversion as solved application scenario.

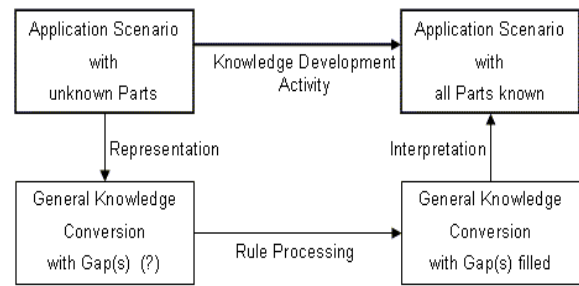For example, the knowledge requirements for a project



Figure 4.  Support of knowledge development scenarios

are known as well as the learning options in the company. From that, one would try to identify minimal knowledge requirements for a new employee, who should work in the project and should be able to fulfill the requirements at least after some learning efforts. Our representation of this scenario is, that we know the result of a knowledge conversion as well as the conversion itself, but we do not know the source knowledge asset. A rule application should deliver the missing knowledge asset.

### A.  Analytic Application Scenarios and their Representation

Analytic application scenarios for knowledge development are characterized by gaps in the corresponding knowledge dynamics chains. Without restriction of generality, we focus on simple scenarios, which can be represented by a single general knowledge conversion. More complex scenarios should be composed of simple ones.
A representation as a general knowledge conversion leads to a set of eight possible scenarios. In the conversion definition with sources, conversion and destinations we can apply zero or more question marks, i.e. gaps of unknown parts, to the conversion. Out of the eight possible scenarios, we do not further consider two of them. The case with no gap is a constructive scenario really, while the case without any known part is not a realistic one. The other six scenarios are outlined in the following and shown in Figure 5.
Scenarios with known destination parts of the conversion and with gaps on the sources side represent situations, where the target of knowledge development activities is known. A known conversion part in the knowledge conversion in this scenario would indicate existing knowledge development options in the company, while a gap indicates missing development support (Scenarios 1 and 2). Scenario 5 describes known sources and destination parts, but missing development options and support in the company. Scenarios 3 and 4 have a complete sources part of the knowledge conversion and gaps in the destinations part. If existing knowledge development options are available, then the scenario would ask for the potential of evolving knowledge applying these options (Scenario 3). If no such options exist, the question of the scenario would be, which knowledge development activities
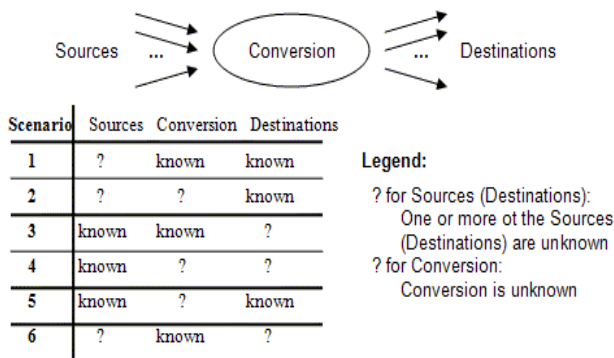
Figure 5. Application scenarios and representations

should be initiated and to which possible result in extended and new knowledge this could lead (Scenario 4). Finally, Scenario 6 assumes existing knowledge development options in the company, but incomplete sources and destinations parts. If only very few out of the sets of sources or destinations are unknown, this scenario can be partly handled with our approach also. Otherwise, especially in the case of completely exclusively unknown sources and destinations, no further treatment is possible.

### B. Rules Application to Representations of Scenarios

As described in Section 3, a rule system has been developed, which is applied to instances of knowledge and conversions introduced in the knowledge ontology.

Only rules for basic knowledge conversions in the type dimension with only one gap exist until now. We therefore are restricted currently to the corresponding 1-to-1 cases of scenarios 1, 3 and 5 as described before in Figure 5. A rule for Scenario 5 case has been given in Sub-section 3.2.. For each of Scenarios 1 and 3, there exist five such 1-to-1 cases, because the known conversion part must be one of the five basic knowledge conversions in the type dimension. Here we analyze two cases and provide appropriate rules.

1) Scenario: Source → Socialization → ?
   The following rule produces a new destination *Tacit_Knowledge:*

   Knowledge(?k1) ^ Socialization(?s) ^
   hasSource(?s, ?k1) ^
   swrlx:makeOWLThing(?k2, ?k1)
   →
   Tacit_Knowledge(?k2) ^ hasDestination(?s, ?k2)

2) Scenario: ? → Combination → Destination
   The following rule produces a new source *Available_Knowledge*, it cannot decide on a specific type of *Explicit_Knowledge* or *External_Knowledge*:

   Knowledge(?k2) ^ Combination(?c) ^
   hasDestination(?c, ?k2) ^
   swrlx:makeOWLThing(?k1, ?k2)
   →
   Available_Knowledge(?k1) ^ hasSource(?c, ?k1)

## V. SUMMARY AND OUTLOOK

A conception of knowledge development in an enterprise has been given. It is based on a concept of knowledge and knowledge dynamics. In order to implement this conception, a knowledge ontology has been built and described in this paper, together with reasoning support and in combination with a rule engine. This has opened the path, to solve open questions in application scenarios for knowledge development. With the help of representations, these scenarios can be mapped to general knowledge conversions, which are subject to rule processing in relation to the knowledge ontology. A final interpretation steps leads back to the solved scenario.

Until now only simple application scenarios and their representations are covered by the set of developed rules. In more complex scenarios, possible solutions are no longer unique. With the help of heuristics, which have to be developed, good or acceptable solutions may be identified.

## ACKNOWLEDGEMENT

## REFERENCES

[1] Ammann, E.: A Meta-Model for Knowledge Management, In: 5th Int. Conf. on Intellectual Capital and Knowledge Management (ICICKM), New York 2008, pp 37-44.

[2] Ammann, E.: The Knowledge Cube and Knowledge Conversions, In: World Congress of Engineering, International Conference on Data Mining and Knowledge Engineering (ICDMKE), London, UK 2009, pp.319-324.

[3] Lehner, F.) Wissensmanagement (in German), 3rd ed., Hanser, München 2010.

[4] Bubenko, J.A., Jr., Brash, D., and Stirna, J.: EKD User Guide, Dept. of Computer and System Science, KTH and Stockholm University, Elektrum 212, S-16440, Sweden.

[5] EKD – Enterprise Knowledge Development, ekd.dsv.su.se/home.html. Last access: August 16, 2010.

[6] Nonaka, I. and Takeuchi, H. The Knowledge-Creating Company – How Japanese Companies Foster Creativity and Innovation for Competitive Advantage, Oxford University Press, London 1995.

[7] Boisot, M.H. "Knowledge Assets", Oxford University Press, 1999.

[8] OWL Web Ontology Language Reference, http://www.w3.org/standards/history/owl-ref. Last access: August 16, 2010.

[9] Hasler Rumois, U. Studienbuch Wissensmanagement (in German), UTB orell fuessli, Zürich 2007.

[10] De Jong, T., Fergusson-Hessler, M.G.M. "Types and Qualities of Knowledge", Educational Psychologist, 31(2), 1996, pp.105-113.

[11] SWRL: A Semantic Web Rule Language Combining OWL and RuleML, http://www.w3.org/Submission/SWRL/ . Last access: August 16, 2010.

[12] Jess Rule Engine, http://www.jessrules.com. Last access: August 16, 2010.

# URBANO: A Tour-Guide Robot Learning to Make Better Speeches

José Javier Rainer, Ramon Galán

Intelligent Control Group. Universidad Politécnica de Madrid
C/ José Gutiérrez de Abascal 2. 28006 Madrid SPAIN
e-mail: javier.rainer@upm.es; ramon.galan@upm.es

*Abstract*—**Thanks to the numerous attempts that are being made to develop autonomous robots, increasingly intelligent and cognitive skills are allowed. This paper proposes an automatic presentation generator for a robot guide, which is considered one more cognitive skill. The presentations are made up of groups of paragraphs. The selection of the best paragraphs is based on a semantic understanding of the characteristics of the paragraphs, on the restrictions defined for the presentation and by the quality criteria appropriate for a public presentation. This work is part of the ROBONAUTA project of the Intelligent Control Research Group at the Universidad Politécnica de Madrid to create "awareness" in a robot guide. The software developed in the project has been verified on the tour-guide robot Urbano. The most important aspect of this proposal is that the design uses learning as the means to optimize the quality of the presentations. To achieve this goal, the system has to perform the optimized decision making, in different phases. The modeling of the quality index of the presentation is made using fuzzy logic and it represents the *beliefs* of the robot about what is good, bad, or indifferent about a presentation. This fuzzy system is used to select the most appropriate group of paragraphs for a presentation. The *beliefs* of the robot continue to evolving in order to coincide with the opinions of the public. It uses a genetic algorithm for the evolution of the rules. With this tool, the tour guide-robot shows the presentation, which satisfies the objectives and restrictions, and automatically it identifies the best paragraphs in order to find the most suitable set of contents for every public profile.**

*Keywords- Cognitive systems; learning; autonomous robot; fuzzy systems; decision making.*

## I. INTRODUCTION

Autonomous robots are intelligent machines capable of performing tasks in the world by themselves, without explicit human control over their actions [1].

Within the development of multiple applications for a mobile robot, probably one of the first real world applications of indoor service robots has been mobile robots serving as tour guides in museums or exhibitions. We have developed our own interactive mobile robot called Urbano specially designed to be a tour guide in exhibitions [2]. The basic characteristics of Urbano are described in Section 3.

Our goal has been to create an automatic presentation generator that allows the flexible and dynamic display of information depending on the distinct kinds of audiences and other parameters that characterize the presentation. The automatic selection of contents for composing sophisticated presentations is a non-trivial task. If the aims of the presentation, preference and interest for a particular subject have to be taken into account, it becomes even more complex.

The knowledge is based on an ontology of domain-specific concept words. Ontologies have been known in computer science as consensual models of domains of discourse, usually implemented as formal definitions of the relevant conceptual entities [3].

Some of the most recent works about automatic generation are: [4][5][6][7][8]. These works propose different architectures and methodologies than those presented here. What is original is the introduction of fuzzy logic in the pruning of the resulting tree and in the quality index of the presentations.

This paper is focused on the automatic generation of presentations by a robotic system. The MINERVA robot from Carnegie Mellon University was one of the first service robots for guiding tours. Others were Rhino, from Bonn University, or Xavier, built to participate in AAAI Robotic Competition in 1993.

Aspects concerning the grounding of the symbols used by the robot, e.g., paragraphs in the presentation or the learning of new concepts – have been circumvented for problems of space. In [9], the authors explore the problem of learning and the symbol-grounding problem, and propose a systemic and integrative approach both problems.

## II. AUTOMATIC TEXT GENERATION

The history of natural language processing is characterized by the influence of artificial intelligent [10][11]. A natural language generator generally has access to a wide knowledge field from which it must select information to present to the users in various formats. Generating text is, therefore, a decision-making problem with multiple restrictions: knowledge, available linguistic tools, the objectives of the user to whom the text will be directed, the situation, and past discourse. It deals with identifying the factors involved in this process and determining the factors that intervene in this process and its consequences [11][12][13].

An additional contribution of our work is to provide a dynamic framework that allows us to take into account each possible scenario. On the other hand, there can be no single general-purpose presentation format for all users, because

each user differs in all aspects of interests and expertise levels, and in the devices used to visualize the requested information. These aspects are used as restrictions.

As a result, to prepare a presentation, the candidate items can be multiple, depending on knowledge server. The items selected in the generation process contain information about the theme of the presentation, but must as well include the items that are semantically related to the presentation. Thus, for example, a discussion of the painter Velazquez will also speak about the Italian painter Caravaggio due to his influence on Velazquez.

In our case a presentation generator is considered as a Cognitive Skill. It is assumed that skill development, e.g., giving presentations progressively more suitable to the public being addressed – is a fundamental architectural epiphenomenon. Rather than viewing a presentation as a mere form of communication, the focus here is on leveraging it as a means of expanding socio-communicative skills.

Several research projects have been undertaken to develop software tools for generating narratives, histories or presentations, and they have described many characteristics for generating presentations, but the quality criteria vary [14][15].

The decision rules to establish the fact that a particular swap in the presentation strategy is useful and required are not clear yet. An interesting situation arises when the data changes and the environment are dynamic.

### A. Quality criteria of a presentation

The quality of a presentation is defined by the different aspects that it characterizes, referred to in [16], as: nature, purpose, duration, and number of participants; for other authors, such aspects as connecting with the audience, the interests of the participants, the first and last minute, changes of rhythm, and "being natural and having fun" are very important. All the aspects have in common the difficulty of being effectively quantified for a computer program.

The prototype as developed has proposed the following *quality criteria*:

- Differences in time used and the anticipated time
- Time dedicated to the theme
- Time dedicated to entertainment
- Difficulty in understanding the theme
- Interest in the theme
- Time dedicated to related information (anecdotal)
- Time for interaction with the public
- Original information
- Non-conventional focus

These criteria are also used to obtain the assessment of the public.

### III. URBANO

This Section describes the Urbano robot system, its hardware software and the experience we have obtained through its development and use until its actual mature stage.

### A. URBANO, an interactive Mobile Tour-Guide Robot

This Section doesn´t pretend to be an exhaustive technical description of algorithms, mathematical or implementation detail, but just an overview of the system.

Urbano robot is a B21r platform from iRobot, equipped with a four wheeled synchrodrive locomotion system, a SICK LMS200 laser scanner mounted horizontally in the top used for navigation and SLAM, and a mechatronic face and a robotic arm used to express emotions as happiness, sadness, surprise or anger.

The robot is also equipped with two sonar rings and one infrared ring, which allows detecting obstacle at different heights that can be used for obstacle avoidance and safety. The platform has also two onboard PCs and one touch screen.

The software is structured in several executable modules to allow a decoupled development by several teams of programmers, and they are connected via TCP/IP. Most of these executables are conceived as servers or services providers, as the face control, the arm control, the navigation systems voice synthesis and recognition, and the web server. The client-server paradigm is used, being the only client a central module that we call the Urbano Kernel. This kernel is the responsible of managing the whole system [2].

URBANO robot has a technology based on distributed application software. The recent version is an agent based on architecture that uses a specific CORBA approach as an integration tool. The robot has many functions: speaks, listens, navigates through the environment, moves his arm, responses to stimuli that affect its feelings.



Figure 1. Urbano Tour-Guide Robot.

### B. URBANOntology

The knowledge server consists of a Java application developed using the libraries of Protégé-OWL API. The tool is capable of reading and editing files in ".owl" format where the knowledge is stored in the form of ontologies and the management of the information from the kernel is made by means of messages that codify the request of specific information, and the reply is obtained from the server or the introduction of new data.

The functions of the knowledge server are: loading and saving ontologies; creating, renaming, and deleting classes or instances; displaying properties of a class; showing subclasses or superclasses; showing or entering the value of a property; integrating one ontology into another; handling queries.

## IV. APG AGENT

APG agent software has been developed to be integrated in the architecture based on the agents that constitute the software of the Urbano robot. The developed computer application is activated on receiving a request from a user that selects a file that contains the pattern to be developed.

APG will request all the information from the knowledge server agent, using restrictions that it needs to generate a file with the best presentation to be used by the system in the next performance by the robot.

Figure 2 shows this flowchart. It must be mentioned that presentations have been stored with quality indices (QI) for each of them, and in some cases, when faced with a request very similar to a presentation already performed with high quality indices; this presentation will be used without needing to generate a new one.
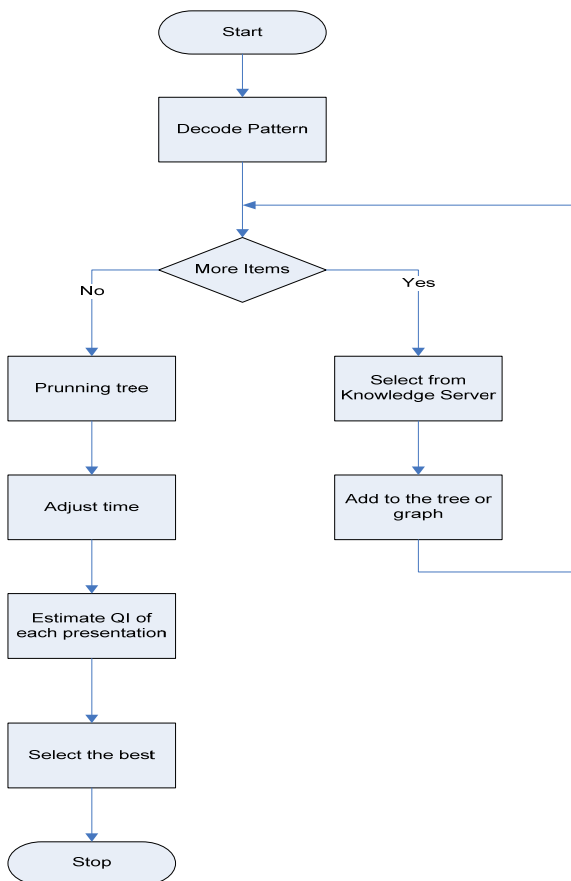


Figure 2. Flowchart APG

## V. PARAGRAPHS

The paragraph is used by the robot as the minimum element of expression. It is assumed that the paragraphs have a size such that coincidences or references between them lack meaning, and that they express a meaningful content. The following paragraph introduces the painting "Las Meninas" by Diego Velazquez:

*"A portrait of the infanta Margarita, daughter of Felipe IV (1605-1665), surrounded by her servants or 'family' in a hall of Madrid's Alcázar Palace.*

*This, the most famous of Velasquez's works, offers a complex composition built with admirable skill in the use of perspective, the depiction of light, and the representation of atmosphere.*

*There have been innumerable interpretations of this subject and later references to it. The most numerous emphasize a defense of the nobility of painting versus craft. Velasquez portrays himself, painting the painting itself, on the left of the canvas, thus affirming the supremacy of the art of painting. The infanta Margarita (1651-1673), wears white and appears in the center of the composition, surrounded by her ladies in waiting, the "meninas" María Agustina de Sarmiento and Isabel de Velasco, along with two court buffoons, María Bárbola and Nicolasito Pertusato, and a mastiff. Behind her, the duenna Marcela de Ulloa converses with the quartermaster, José Nieto, who is in the doorway.*

*The King and Queen, Felipe IV and Maria de Austria (1634-1696) are reflected in the mirror at the back of the room, leading to series of extraordinarily complex spatial relations."*

For each of the paragraphs there also exists a script that specifies the facial expressions and arm movements of the robot, the tone of voice and type of voice to use, the position in the floor plan, and various details for interaction with the public.

But in the ontology, each paragraph is "related" to the categories, which it belongs (e.g., in the previous paragraph: Is_a Painting, Is_a Description, Is_a Adult Level, Painted_by Diego Velazquez, etc. These relationships are used for the selection of the paragraph).

## VI. PATTERNS AND TREES

The pattern, defined by the user, signals the elements that should form part of the presentation. For each item, these elements are established: its identification, its priority, its numerical order, its reference to the theme within available knowledge. There could be a limitation of time and a very large patter, in that case the priority index indicates the most important content to be included, but there could be the opposite case, that there is time left, then the system incorporates content that it is not initially considered.

It uses XML as the language to represent the patterns, which guarantees an easy use with different tools and programming languages. XML has emerged as a de facto standard for encoding and sharing data between various applications. XML is also useful for structured information

management, including information contained in knowledge server [18]. Our proposal is to design a software tool that helps the user to create and maintain the patterns.

APG requests from the knowledge server the paragraphs available for each theme, identified by its reference. The knowledge server will submit one or more paragraphs, because of a same concept can have several versions. Thus the pattern indicated in Figure 3 shows that first item of the presentation will be "museum_presentation" identified on the knowledge server as "Greetings" and it could involve, for example, three possible paragraphs, which will be included in the tree of possibilities, as is shown in Figure 4.

Some global restrictions including the level of audience will be used to prune the tree, eliminating the paragraphs that do not correspond to the requested level.

Three typical alternative heuristic searches have been tested to trim the tree. The first uses "brute force" to generate all the possible combinations and to group all the numeric values of the "quality criteria" of the paragraphs that form the presentation, and then, using a set of fuzzy rules, estimates the quality index. It selects the presentation with the highest index. See Figure 5.

The second alternative uses "best-first search" so that as it goes along it takes the option that partially presents the best index. This alternative is without a doubt the fastest, but it cannot guarantee the selection of the best option.

The third alternative modifies the previous one so that it generates "backtracking" if the quality index falls below a minimum.



Figure 4.   Tree Data Structure



Figure 5.  Different options

The presentation generated in this way analyzes the estimated time for its execution, and if this is greater than anticipated, it eliminates the paragraphs with the least necessary priority. On the other hand, if there is enough time, it includes some socially-oriented paragraphs such as jokes or comments about sports, politics, or local events.

## VII.   DECISION-MAKING

Decision-making is a part of the paradigm proposed by Zadeh [19] that has been currently examined in [20]. In a dynamic scenario as ours, and because of the nature of the information that the system will handle, proper tools are needed to provide the intelligence for decision-making and supervision.

Decision-making is the cognitive process of selecting a course of action from multiple alternatives. Fuzzy set approaches to decision-making are usually most appropriate when human evaluations and the modeling of human knowledge are needed [21].

The proposed solution uses fuzzy rules to prune items to submit from the tree. It will use a variable that indicates the likelihood of inclusion in the submission of a particular content. The fuzzy system will generate these values.

The fuzzy rules enable more flexibility. These rules will be adjusted and expanded.

All information available at the moment about the quality criteria and its influence on the quality index is stored in the ontology of the knowledge server.

The semantic network will indicate that the influence of time dedicated to the theme, expressed in a percentage, is

```
<pattern>
        <id>children_visit_museum</id>
        <date_creation> … </date_creation>
        <date_lastused>…</date_ lastused>
        <restrictions>
                <age><10</age>
                <technical_level>low</technical_level>
                <time>60</time>
        </restrictions>
        <item>
                <item_id>Museum_ Presentation</item_id>
                <item_order>1</item_order>
                <item_priority>10</item_priority>
                <item_data>Greetings</item_data>
        </item>
        <item>
                <item_id>Tour_Guide</item_id>
                <item_order>2</item_order>
                <item_priority>20</item_priority>
                <item_data>Tour_Guide</item_data>
        </item>
        <item>
        <item_id>Painting_presentation</item_id>
                <item_order>3</item_order>
                <item_priority>40</item_priority>
                <item_content>
                        <item_data_id>key_picture</item_data_id>
                        <subitem> title</subitem>
                        <subitem>date</subitem>
                                <subitem_multiple>description</subitem_multiple>
                        <subitem_multiple>period</subitem_multiple>
                </item_content>
        </item>
</pattern>
```
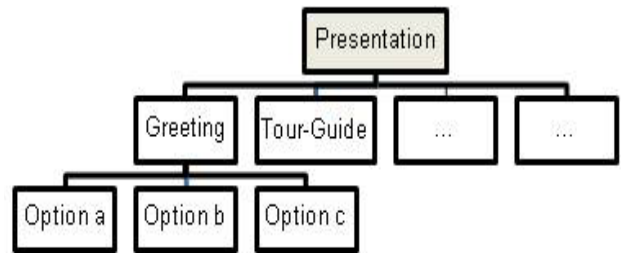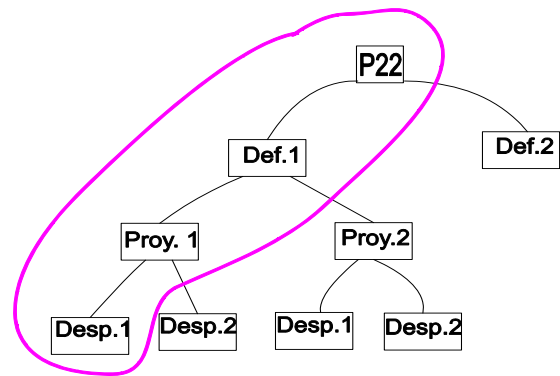
Figure 3.   XML Pattern Example

VERY favorable if the time is VERY_HIGH but LOW favorable if the time is HIGH or NORMAL, and VERY unfavorable in any other case. This relationship is defined by belonging to specified categories. In the developed prototype, the fuzzy rules used to calculate the quality index are obtained by consulting the categories that belong to a criteria in the knowledge server.

Five linguistic terms are defined: VERY_HIGH (VH), HIGH (H), NORMAL (N), LOW (L), VERY_LOW (VL). The fuzzyfication phase uses the function of membership to initially equidistant triangles, but in the learning phase their centers can vary. The exit variable quality_index is also modeled with five terms and triangular functions. The technique of centroid method is used in the defuzzyfication phase.
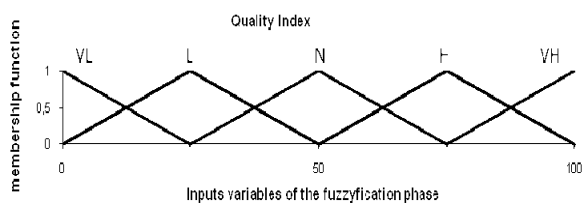


Figure 6. Inputs variables of the fuzzyfication phase

An interesting aspect is the partial quality estimation for the presentation. When considering three paragraphs (#p1,#p2,#p3), for example, it may happen that some criteria have unrepresentative values; if, for example, #p1,#p2,#p3 correspond to paragraphs of technical description, the entertainment criteria will be null and the presentation will have a low estimation. In the tests made using pruning the best-first search, presentations are obtained with very poor estimations, and to compensate for this effect, an option of the best-first was introduced, but rejecting presentations with a low index.

If the categories of information expressed in the pattern are very generic, the number of possible presentations increases enormously, but it permits the robot to generate higher quality presentations.

## VIII. LEARNING PHASE

The most important characteristic of the proposal is the ability of the robot to learn. Initially it is thought that the robot will have a small number of quality criteria available to evaluate some presentations as good and others as bad, corresponding to the minimum level of education for a professional guide, in order to guarantee a minimum level of quality in its presentations.

A simple questionnaire has been designed that the public can fill out after attending a presentation by the robot. It asks for an evaluation of each quality criteria known at the time, indicating whether the robot should spend more or less time on each item, and a percentage evaluation of what the public considers valuable in the presentation. The Table I shows an example.

A proper statistical treatment of the questionnaires, eliminating the extremes and requiring a minimal quantity of data is performed.

TABLE I. QUALITY CRITERIA

| Quality Criteria | Should be |
|---|---|
| Time spent on the theme | - |
| Time spent on entertainment | + |
| … | |
| **New criteria to bear in mind** | **%** |
| Answers to questions | 25 |
| | |
| **Global evaluation** | **%** |
| | 60 |

A genetic algorithm is used, an adjustment of linguistic terms and the membership functions, will permit the quality index to be the closest to the average expressed by the public. To get a greater accuracy, the genetic algorithm is simultaneously used over several presentations as an attempt to offset the local minimals for a presentation.

Learning is realized when new criteria of quality are incorporated in the paragraphs.

The genetic algorithm realizes a readjustment of the rules when it produces a disparity between public opinion and the IC. It utilizes a genetic algorithm for the evolution of the rules

The + and – indications are used to eliminate individual cases generated by the genetic algorithm in which, while still generating a correct quality index, the evaluation of the criteria runs contrary to public opinion.

In the tests carried out, it was shown that the "beliefs of the robot" about the quality indices converge toward public opinion, and as a result generate presentations that are evaluated more favorably.

## IX. CONCLUSION

As a conclusion, special attention has been given to a mechanism for automatic generation of presentations, analysis of search algorithms, learning phase and optimization of fuzzy logic rules, taking into account the intrinsic difficulty of natural language processing and automatic generation.

The generator is a dynamic system where knowledge will increase and, therefore, it will do the quality of the presentations. The robot will become increasingly better when making tour-guided visits. It also adds the ability to gesticulate while is conducting the presentation: it specifies the facial expressions and arm movements of the robot, the tone and type of voice to use, the position in the floor plan, and several other details for interaction with the public.

This development will allow tour-guide robots to offer more affective learning and a dynamic tour-guide visit, because the public generally has become more sophisticated

and, also, its expectations and demands. It provides a better use of knowledge management. At the same time, advances in a new visitor-oriented approach, progress towards the creation of modular and scalable scenarios.

It presents a novel approach to the use of a computer ontology to represent the corpus that the robot works with, and the quality criteria for a presentation.

Methodology designed for the automatic generation of presentations can be used for other applications related with decision making for autonomous robots. In the Urbano project design; there is a scheduler agent that selects the best task between different tasks, which the system has to do in a day. The criteria used for selecting the adequate task in short-term plan can be revised to optimize the long-term objective. A happiness model, as long-term objective, and a decision-taking mechanism, as short-term planner, was modeled for Urbano. Both use fuzzy logic and are adjusted by genetic algorithms that use the public opinion to learn to be a good tour-guide robot.

These systems have special importance to develop learning support for environments that require greater motivation and commitment, such as classrooms and workshops for students with special educational needs.

## ACKNOWLEDGMENT

## REFERENCES

[1] G. Bekey, Autonomous Robots: from biological inspiration to implementation and control, MIT Press books, 2005.

[2] D. Rodríguez-Losada, F. Matia, R. Galán, M. Hernando, J. M. Montero, and J. M. Lucas, Urbano, an Interactive Mobile Tour-Guide Robot. Advances in Service Robotics. Ed. H. Seok. In-Teh, 2008, pp 229-252.

[3] M. Uschold and M. Grüninger, "Ontologies: Principles, Methods, and Applications," Knowledge Eng. Rev., vol. 11, 1996, pp. 93–155.

[4] P. Barnaghi and S. Kareem, "Ontology- Based Multimedia Presentation Generation," IEEE TenCon, 2005.

[5] S. Bocconi, "Automatic Presentation Generation for Scholarly Hypermedia," 1st International Workshop on Scholarly Hypertext at the fourteenth conference on Hypertext and Hypermedia (HyperText 2003), August 30, 2003.

[6] A. Celentano and O. Gaggi, "Schema Modelling for Automatic Generation of Multimedia Presentations," SEKE, 02 Ischia Italy, 2002.

[7] F. Frasincar and G. Hoube, "XML-Based Automatic Web Presentation Generation," Proc. WebNet World Conference on Www and Internet, 2001.

[8] S. Little, J. Geurts and J. Hunter, "Dynamic Generation of Intelligent Multimedia Presentations though Semantic Inferencing," The Sixth European Conference on Research and Advance Technology for Digital Libraries, Rome, Italy EDCL, 2002.

[9] J. Gómez, R. Sanz and R. Galán, "Learning in Technical Systems: a Sign based Approach," In proceedings 1st International Workshop on Cognition for Technical Systems, Munich, 6- 7 October 2008.

[10] R. Cole , J. Mariani , H. Uszkoreit, A. Zaenen, and V. Zue, Survey of the State of the Art in Human Language Technology. Cambridge: Cambridge University Press, 1996.

[11] K. Sparck Jones, Natural languages processing: a historical review. University of Cambridge, 2001.

[12] J. Allen, Natural Language Understanding, Redwood City, Ca.: Benjamin /Cummings, 2ª ed., 1995.

[13] J. Allen, "AI Growing Up. The Changes and Opportunities", AI Magazine, 19 (4), 1998.

[14] D. Jurafsky, and J. Martin, Speech and Language Processing. Prentice-Hall, 2000.

[15] K. Brooks, Metalinear Cinematic Narrative: Theory, Process and Tool. Phd Thesis, MIT, 1999.

[16] M. Cavazza, F. Charles, and S. J. Mead, "Character-based interactive storytelling," IEEE Intelligent Systems: Special Issue on AI in Interactive Entertainment, 17(4):17–24, 2002.

[17] F. Laure, Técnicas de Presentación. Cecsa, 2002.

[18] M. Y. Maarouf and S. M. Chung, "XML Integrated Environment for Service-Oriented Data Management," 20th IEEE International Conference on Tools with Artificial Intelligence, 2008.

[19] L. A. Zadeh, "Fuzzy sets," Information and Control, 8, 338–353, 1965.

[20] L. A. Zadeh, "Is there a need for fuzzy logic?" Information Sciences, 178(13), 2751–2779, 2008.

[21] C. Kahraman, "Fuzzy set applications in industrial engineering", Information Science, Volume 177, Issue 7, Pages 1531-15321, 2007.

# Graphical Modelling in Mental Health Risk Assessments

Olufunmilayo Obembe, Christopher D. Buckingham
School of Engineering and Applied Science
Aston University
Aston Triangle, Birmingham B4 7ET, UK
obembeo@aston.ac.uk, c.d.buckingham@aston.ac.uk

*Abstract* - **Probabilistic models can be a combination of graph and probability theory that provide numerous advantages when it comes to the representation of domains involving uncertainty. In this paper, we present the development of a chain graph for assessing the risks associated with mental health problems, which is a domain that has high amounts of inherent uncertainty. The Galatean mental health Risk and Social care Tool, GRiST, has been developed to support mental-health risk assessments by using a psychological model to represent the expertise of mental-health practitioners. It is a hierarchical knowledge structure based on fuzzy sets for reasoning with uncertainty. This paper describes how a chain graph can be developed from the psychological model to provide a probabilistic evaluation of risk that complements the one generated by GRiST's clinical expertise.**

*Keywords- mental health risk assessment; probability graphs; chain graphs.*

## I. INTRODUCTION

Risk assessment is a fundamental part of life, whether it be a mundane decision about the chance of rain or a much more vital one about the risk of a nuclear power station malfunctioning. In the mental-health domain, predicting whether someone is going to commit suicide or engage in an act of violence is extremely difficult, partly because the likelihoods are so low but also because of the lack of statistical data. The Galatean mental-health Risk and Social care assessment Tool (GRiST, [1]) was developed to address these problems by modelling how expert mental-health practitioners make risk assessments. However, its accumulating database of risk data has become a resource for more probabilistic approaches such as probability graphs, which are well-suited for capturing and reasoning with uncertainty where there is prior knowledge structuring [2].

In the past, mental health risk assessment was predominantly carried out using unstructured clinical approaches but it has since been realised that the best results can be obtained by using a combination of both structured clinical judgements and actuarial tools, such as one based on a probability graphical model [3]. This paper explores how a probability chain graph can be developed from the GRiST model of expertise to perform mental health risk assessments. In Section II, a brief overview of related work is given. In Section III, the background to mental-health risk assessment and GRiST is briefly introduced, followed by a discussion of the types of building blocks used in probability graphs in Section IV. In Section V, the development of the composite GRiST probability chain graph from these building blocks is presented. The paper then concludes with an outline of future work for the research.

## II. RELATED WORK

The challenge of providing effective risk assessment in the mental health domain is not a new one. The two broad approaches used are the clinical and actuarial approaches. The clinical approach can be split into structured and unstructured methods. Both are based on a clinician's experience but the unstructured method has no other input and is thus highly subjective. The structured clinical approach is more formal because it links the clinician's judgement with data-gathering tools that help guide the collection of information. In contrast, the actuarial approach uses statistical methods to provide risk assessments and is the least subjective of the approaches. Proponents of each method have argued about the various advantages and disadvantages over many years [4-6] but current policy is to integrate them where possible.

One of the actuarial methods (the category to which this paper belongs) uses an important technique from computer science: the multiple iterative classification tree (ICT) model, which is part of the MacArthur Violence Risk Assessment method [7]. This tool was developed to predict the risk of violence behaviour among recently discharged patients [7]. The results obtained from the ICT model had a high level of accuracy for the specific population group but the tool is resource and time intensive [8]. Another model uses tree mining [9] but a problem with all these approaches is the difficulty of obtaining clinical data that covers multiple risks and probabilities that can be generalised across populations.

Although GRiST is a model based on structured clinical judgement, it collects comprehensive and precisely defined data for all risks that are automatically stored in a database and thus available for probabilistic analyses. The idea is to link its clinical judgements to actuarial analysis and create a risk tool that explicitly connects the two risk-assessment approaches. Furthermore, the hierarchical structuring of GRiST's knowledge base provides the potential for

informing the structure of corresponding probabilistic graphs, which have limited current use in the mental health domain. GRiST allows both causal and non-causal relationships to be modelled between the various risk factors (cues), which helps provide a more accurate domain representation. This paper explores how to exploit GRiST for creating a probabilistic graph model of risk assessment and thus link probabilistic analyses to its fuzzy-set modelling of clinical reasoning. See [10] and [11] for a review of some identified advantages for using graphical models in risk assessments.

### III. THE GALATEAN RISK SCREENING TOOL, GRIST

GRiST is a decision support system for mental health risk assessments [12] and is based on a psychological model of classification known as the Galatean model [13]. The variables modelled in the GRiST structure and the relationships between them are represented by a hierarchical tree structure (see Figure 1) and uncertainty is processed using fuzzy-set membership grades (MGs). In the GRiST knowledge structure (GKS), there are two main types of nodes (variables), namely concept nodes and datum nodes. Datum nodes are leaf nodes in the tree and thus do not have children. They represent measureable input data such as a person's *insight into behaviour* (Figure 1). Concept nodes are made up of two or more component nodes, which could be datum nodes or other concept nodes (e.g. *substance misuse* in Figure 1)
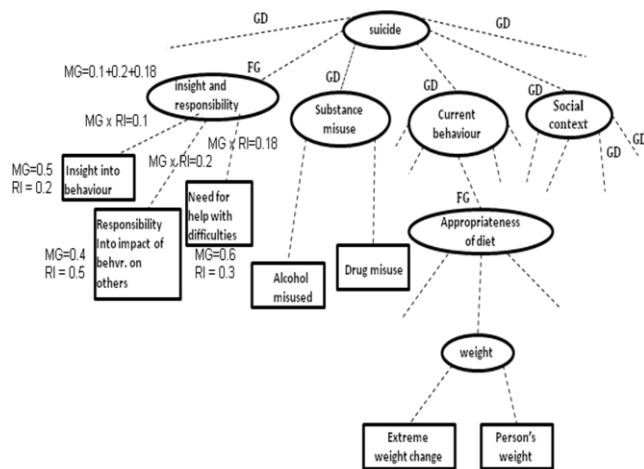
.



Figure 1. Small subsection of the GRiST Knowledge structure with datum and concepts nodes represented as rectangles and ovals respectively and generic nodes by g and generic distinct nodes by gd.

Both datum and concept nodes can be subcategorised into types that describe their locations and contextual behaviour in the GKS. Nodes that occur only once are called non-generic whereas those that have multiple occurrences are named generic. Some generic nodes are additionally distinguished as "generic distinct" because these have different uncertainty parameters for different locations in the tree whereas the plain generic nodes have exactly the same internal uncertainty representation wherever they occur.

The full GKS was originally elicited from 46 domain experts and comprised of 3,026 nodes with 338 unique concept nodes and 692 unique datum nodes [12]. Subsequent knowledge validation reduced it considerably so that only about 220 datum nodes are identified for data collection.

#### A. GRiST Data Structures and Uncertainty Processing

In GRiST, uncertainty is encapsulated by MGs and relative influences (RIs). An MG represents the degree of membership of an object in a node of the tree, with each nodes's MG ultimately contributing to the top-level risk membership (e.g., *suicide* and *self harm*). The actual contribution depends on the node's RI, which weighs the node's influence compared to its siblings. For example, from Figure 1, *insight into behaviour* has three children with RIs of 0.2, 0.5, and 0.3, which filter their MG contributions to the parent node as shown.

The MG distributions of the datum nodes were specified by clinical experts and enable patient cues to map to MGs that feed through the entire GRiST tree, from leaf nodes to the top-level risks. Equation 1 shows the formula used but see [13] for further details on the MG propagation process:

$$MG(X) = \sum_{p=1}^{P} (MG(datum_p) \times \prod_{l=1}^{L} RI_{lp}) \qquad (1)$$

Using Equation (1), the MG of each concept node is calculated by multiplying the MG of the datum node along each path $p$ with all the RIs along that path and up each level $l$ leading to the concept, and then summing the value obtained with all the corresponding values obtained along each path to the concept. An example of this is depicted on the left branch of the tree in Figure 1, to show how the MG is generated for the *insight and responsibility* concept.

#### B. GRiST Knowledge Structure (GKS) Constraints and Independence Properties

Having described the GKS and its data types, we now briefly outline the constraints incorporated in this knowledge structure and their correlation to various independence properties. A brief description is given below but see [14] for a more in-depth discussion. In this paper, we extend the work presented in the earlier paper [14] and expand on the constraint mapping and structure combination rules. Semantically there are three different types of relationships that can exist between any two nodes in the GKS, as follows:

- *IS-A* relationships refer to a 'kind-of' relationship, where the children nodes are a type of the parent node and are thus associated through their common parent. An example of this type of relationship is the parent node *substance misuse* and its children nodes *alcohol misuse* and *drug misuse*.
- *Contribute-to* relationships refer to those where the children nodes 'contribute to' and influence the parent node. For example the relationship type between the parent node *constraints on suicidal behaviour* and its internal nodes *insight and responsibility* and *religious values/beliefs affecting suicide risk* (see Figure 1) are of type *contribute-to*

because the children nodes directly contribute to the value of the parent node. Of the three types of relationship, this is the only causal one.

- *Wrapper* relationships occur when the parent node serves as a form of container for the children nodes (the parent 'wraps' the variables together) rather than being a cohesive variable in its own right. In this type of relationship there is a correlation between the children and parent nodes but not one that can be assumed to be causal. For example the relationship types between the parent concept *general current behaviour* and its children nodes *appropriateness of diet*, *challenging behaviour*, *daily activity*, *reckless risk taking*, *sleep disturbance*, *unintentional risk making* and *uncharacteristic recent change in behaviour* is of the type *wrapper* (it is partially depicted in Figure 1).

For the purpose of mapping to the probability graphs, it is important to note that the *is-a* and *wrapper* relationship types map to non-causal links, whereas the *contribute-to* relationship type is causal. As a result of these semantic relations between the various node types in the GKS, it is possible to give a concise list of the various constraints that exist within the model. These constraints in turn aid in the formal definition of the set of component structures that the GKS can be decomposed into. These component structures are then mapped into probability graphs that will eventually form the building blocks for the resultant probability chain graph that will provide the model for inferring the final risk assessments.

Detailed discussion of the GRiST component structures can be found in [14] but for clarity some of the constraints are outlined next. In the definitions, a root concept refers to the highest ancestor node in the structure under discussion. So, for example, if the entire diagram depicted in Figure 1 is taken to be one structure, then its root concept node is the *suicide* node. The constraints are for certain aspects of the structure that must hold wherever it is located in the GKS.

**Constraints Related to Generic Root node structures:**
- The RI value of the root node varies.
- The MG value of the root node is fixed.
- The RI values of the internal nodes are fixed.
- Given that the MGs and RIs of all internal nodes are fixed then the point of reference (i.e. the context) for the internal nodes is their root node.
- Generic root structures must be kept as cohesive wholes. It is possible to have a cohesive whole within another cohesive node i.e., root concept of type generic with internal nodes of type generic.
- As every node within a root concept of type generic has a fixed RI and MG everywhere the root concept occurs, if one of the internal structures is of type generic distinct (defined later), it will also need to have fixed RIs and MGs values within the context of the root concept everywhere it occurs. This is not

the default or usual behaviour of these nodes, as explained next, but is in fact a special case.

**Constraints Related to Generic Distinct Roots**
- The RI value of the root node varies.
- The MG value of the root node varies.
- The RI values of the internal nodes vary.
- If a generic distinct node has at least one node of generic type as an ancestor then the context (or point of reference) for the generic distinct node is the nearest ancestor of type generic. Otherwise, the context for the generic distinct node is its neighbouring nodes.
- In the case where all the internal nodes of a generic distinct root concept are of type generic, if all the MGs and RIs of these internal nodes are always fixed it is obvious that the root concept MG value cannot vary and will itself always be fixed too, which is incorrect behaviour for a generic distinct root concept. This therefore leads to the constraint that a root concept of type generic distinct cannot have all its internal nodes to be of type generic. To make it possible for the variation in the root concept's MG value in various locations, there must be at least one internal node of type generic distinct. This is seen to be true in the GKS and is a good test of the validity for generic distinct node definitions.

From these constraints we obtain two structures that the GKS can be broken down into and for which probabilistic equivalents can be acquired.

1. For the generic root node that always has the same uncertainty values regardless of its location in the GKS, the context (i.e., point of reference) for all contained nodes is its own root concept. Within (and only within) the root concept the uncertainty values of the internal nodes are fixed and always remain the same regardless of location.

The name given to this type of structure is the **fixed generic component structure** (FG). An example from Figure 1 is the generic root node *insight and responsibility* and its internal nodes.

2. The generic distinct structures have varying internal RIs and varying root concept MG. The context for these nodes are the neighbouring nodes where the neighbouring nodes refer to the root concept, all its internal nodes (descendents), the root concept siblings and the root concept parents (as the root concept MG varies). It is also dependent on the top risk in which it occurs (e.g., suicide, self harm and so on). This type of structure has been named the **generic distinct component structure** (GD). This structure has no generic ancestor or, more to the point, if it did, it behaves as an FG node and can be ignored as a GD concept.

Both the FG and GD structures are composite wholes that can contain other composite variables within them. An example of this is seen in Figure 1 where the GD structure

with root concept node *current behaviour* contains within it the FG structure with root node *appropriateness of diet*.

The next section explores the relationship between these component structures and the independence properties they represent and maps them to different probability graphs, which will serve as building blocks for the final probability chain graph.

## IV. THE BUILDING BLOCK PROBABILITY GRAPHS

In this section we develop the building blocks by examining the independence properties of each of the component structures discussed in Section III.

### A. Mapping the Fixed Generic Structure

As mentioned earlier, because in a FG structure the relevant context for the determination of the uncertainty values is the root concept node of the FG structure itself, the independence property of the FG structure is as follows;

*The nodes in a FG structure (i.e., its root node and all its internal structures) are conditionally independent of all other nodes in the GKS.*

The above can be modelled by a Markov Random Field (MRF) [15]. The local Markov property for a MRF structure (the property outlining its independence status) states that a variable is conditionally independent of all other variables given its neighbours [15]. This directly correlates to the independence property of a FG structure, if we replace the variable in the MRF with the FG structure and in a similar manner also define the neighbours to be the FG root concept and all its descendents. An example FG structure from Figure 1 is the *insight and responsibility* root concept and its internal nodes *insight into behaviour*, *need for help with difficulties* and *responsibility into impact of behaviour on others*.

We, however, also need to consider the different relationship types (i.e., *is-a*, *wrapper* and *contribute-to*) between the nodes. Normally the *contribute-to* type should be represented by a directed edge, because the relationship is causal, and both the *is-a* and *wrapper* types are represented with non-directed edges as there is no implied causality.
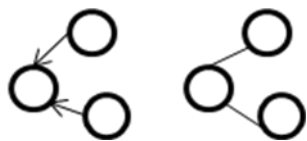


Figure 2. Directed and Undirected Graphs

However, the two diagrams in Figure 2 represent two different independence relations (see [16] for more on independence relations between the different graph types) but for FG structures these are not relevant within the MRF equivalent. Hence the relationships are all represented as undirected edges even when the relationship type between nodes is of type *contribute-to* (i.e., causal).

### B. Mapping the Generic Distinct Structure

The GD structure is highly context sensitive and as such its independence properties are dependent on the identified relationship types between the various nodes in the structure. There are four possible relationship structures that can be obtained from the GD, namely: the non-causal to non-causal; the non-causal to causal; the casual to non-causal; and finally the causal to causal hierarchical relationship links. Recall that *is-a* and *wrapper* relationships are non-causal and the *contribute-to* relationship type is causal.

The independence properties of both the causal to non-causal and the non-causal to causal map to causally linked MRFs, which are in effect chain graphs (these will be discussed in more detail in the next section). The non-causal to non-causal on the other hand map to an undirected graph (MRF). For cases where this mapping is hierarchical, tree structured MRFs (TS-MRFs) [17] allow us to model such cases. And finally, the causal to causal links map to directed graphs (Bayesian Belief Networks) [18].

### C. Summary of Building Block Probability Approach

To summarise the building block probability approach, a two step method is used. Initially each identified component structure (i.e., FG or GD) is represented in the overall graphical model as a composite variable. This results in an embedded model where each node itself represents a graphical model. The second layer, is then reached when we model and consider an individual sub-tree. For the FG structures we map the nodes to MRFs or their variants (as discussed earlier). The uncertainty contribution from each node (i.e., embedded graphical model) is then plugged into the overall graph. See [19] and [20] for other work involving different aspects of embedded graphs and mixture trees.

## V. COMPOSITE PROBABILITY CHAIN GRAPH

Chain graphs are graphical models, which allow both directed and undirected graphs with the constraint that they do not have semi-directed cycles [21]. Linking two variables in a chain graph with a directed edge implies that the relationship between them is causal, and the direction of the edge is from cause to effect. On the other hand variables that are linked with an undirected edge do not have a causal relationship but however have an associative relationship (in a similar manner to MRFs). As a result of the inherent causal and associative relationships contained within the GKS, which are also clearly seen in the mapping to the building block probability graphs (discussed in Section IV), it makes logical sense to model this knowledge structure using a probability chain graph. More in-depth discussions on the chain graph can be seen in [22]

### A. Development of the GRiST Chain Graph

The GRiST chain graph was developed from the building block probability graphs using a step by step combination of the graphs. The combining of multiple probabilistic graphical models requires care and one of the important considerations

is to ensure that the independence properties represented in the different models being combined is maintained in the composite graphs (see [23] for a discussion on a framework for probabilistic graphical model combination).

### B. Maintaining the Integrity of the Independence Properties

To ensure that the composite structure we have developed correctly models both a chain graph and a chain graph with the correct independence properties, some conditions have to be fulfilled.

- During the combination process, cycles and semi cycles must be avoided (to circumvent violating the chain graph constraint). This condition can be fulfilled because as a direct result of the GKS constraints and the GRiST hierarchical form, the graphical building blocks that are to be combined do not contain any cycles or semi-cycles. As such the only way that cycles can be introduced into the combined structure is if the components are combined in an opposite direction to the uncertainty propagation, which would not make logical sense.
- Secondly a chain graph needs to map into subsections (or blocks), with variables within a subset being linked via undirected edges and variables between blocks being linked via directed edges. Again this is possible because of the hierarchical form of the GRiST structure (and hence the building blocks' probability structures). The different levels in these knowledge structures map directly to the notion of subsets in chain graphs.

In addition to the above, to ensure that the integrity of the independence properties are maintained throughout the development process, the possible combination options between the various graphical building block graphs and their combination rules must be clearly defined, as described next.

### Combining two FG Structures

Any two structures that need to be combined must be directly linked in the GKS and one of the two FG root nodes will be an ancestor to the other. Semantically this means that given the ancestor root concept and all the other nodes comprising the two FG structures (including the other FG root concept, which in effect is seen as an internal node to the ancestor root concept) the combined MRF structure (i.e., from both FG structures) is independent of all other nodes in the model. Therefore regardless of the relationship types between the nodes in these two structures, they should be linked using an undirected edge to obtain an MRF or TS-MRF (depending on the type of FGs combined). For the FG and FG structure combinations, if the combination rule does not override the relationship types between nodes in some instances, a directed graph might be used to combine two FGs and this will completely change the independencies represented in the GKS. Figure 3 depicts an example of the change in conditional independencies caused by the use of the wrong link type.
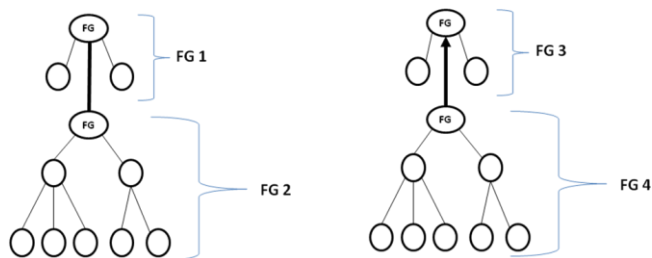


Figure 3. Change in conditional independence as a result of wrong link

In Figure 3, the combination of structures FG 1 and FG 2 results in a composite FG structure that is independent of all other nodes in the model. However, in the second structure (right hand diagram) as a result of the arrow used to link FG structures FG 3 and FG 4, a set of nodes are obtained within the composite combined FG that are conditionally independent of each other given their prior and concurrent nodes.

### Combining a FG Structure and GD Structure

Semantically, when a FG structure is linked to a GD structure, from the definition of the FG structure, we see that it remains independent of the GD structure (recall that the composite FG structure is independent of all other nodes in the model). The GD structure on the other hand is dependent on its neighbours and thus will not be independent of the FG structure. The challenge here is defining the link in such a way that the GD dependence relations with the FG structure remain consistent. In this case the order of the combination is important. Where the FG structure is an ancestor to the GD structure, if an undirected link type is used to combine these two structures, the composite structure will be a FG structure (see Figure 4).



Figure 4. Combination of FG and GD structures resulting in a composite FG structure.

However where the GD structure is the ancestor structure, the relationship type between the two structures is needed to determine the link type. It is a directed link for the *contribute-to* relationship and undirected links for all other relationship types.

### C. Summary of GRiST Chain Graph Development Steps

The following summarises the translation process from the GKS to the final composite GRiST chain graph.

**Step 1:** Partition the GKS into graphical building blocks (i.e., FG or GD).
**Step 2:** For each type of structure, identify the relationship types that exist between the various nodes (i.e., *is-a*, *wrapper* or *contribute-to*).

**Step 3:** Next identify the graphical building blocks that each component structure maps to, based on its independence properties.

**Step 4:** Repeat the above for the entire structure, until you are left with just graphical building blocks needing combination.

**Step 5:** Identify the type of combination to be carried out.

**Step 6:** Apply the relevant combination rules.

**Step 7:** The resultant chain graph should contain a mixture of explanatory, intermediate and response variables, and also maintain the independence properties of the original GKS.

## VI. CONCLUSION AND FUTURE WORK

The development of a probability chain graph for mental health risk assessment has been presented in this paper. We have shown how the knowledge encapsulated in the psychological fuzzy-set based GRiST model can be mapped into initial component structures based on the inherent constraints in the model and how these in turn can be mapped to building block probability structures that can be eventually moulded together to construct a chain graph for mental health risk assessments.

The present solution is ongoing and future work will focus on using data collected by the tool in clinical use to validate the chain graph structure and learn its parameter settings. Examination of the dependencies between variables and their ontological definitions within GRiST will help determine whether the structural relationships are correct. Parameter estimation will then take place in two stages: estimating the potential functions for related cues in the MRF structures followed by estimating the conditional probability distributions for the nodes in the directed segments of the chain graph. Finally the identification of the most effective and efficient inference algorithms for the developed structure will be carried out. The methods discussed in this paper could be applicable to other systems based on hierarchical expertise, especially ones that contain both causal and non-causal relations.

## REFERENCES

[1] See www.galassify.org/grist

[2] P. J. Lucas, "Bayesian Networks in biomedicine and health-care", Artificial Intelligence in Medicine, vol 30, 2004, pp. 201-214.

[3] Department of Health, Best Practice in Managing Risk, Department of Health, London, 2007.

[4] B. Littlechild and C. Hawley, "Risk Assessments for Mental Health Service Users : Ethical, Valid and Reliable?", Journal of Social Work, vol 10(2), pp 211-229, April 2010.

[5] T. R. Litwack, "Actuarial versus Clinical Assessments of Dangerousness", Psychol, Pub Pol, Law, vol 7, pp 409-443, 2001.

[6] M. E. Rice, G. T. Harris, and V. L. Quinsey,"The Appraisal of Violence Risk", Current Opinion in Psychiatry, vol 15(6), pp 589-593, 2002.

[7] J. Monahan, H. J. Steadman, P. C. Robbins, P. Appelbaum, S. Banks, T. Grisso, K. Heilbrun, E. P. Mulvey, L. Roth, and E. Silver, "An Actuarial Model of Violence Risk Assessment for Persons With Mental Disorders", Psychiatric Services, vol 56(7), pp 810 – 815, 2005.

[8] J. Monahan, H. J. Steadman, P. S. Appelbaum, P. C. Robbins, E. P. Mulvey , E. Silver , L. Roth , and T. Grisso, "Developing a Clinically Useful Actuarial Tool for Assessing Violence Risk", British Journal of Psychiatry, vol 176, pp 312-319., 2000.

[9] M. Hadzic, F. Hadzic, and T. Dillon, "Tree Mining in Mental Health Domain", Proceedings of the Proceedings of the 41st Annual Hawaii International Conference on System Sciences, pp 230, 2008.

[10] R. Anderson, F. Camacho, F. Camacho, and R. Balkrishnan, "A Comparison of Risk Classification Methods in Medicare HMO Enrollee Health Risk Assessment", Academy Health Update Health Services Research, vol 19(5), 2002.

[11] S. Ferson, "Bayesian Methods in Risk Assessment", Applied Biomathematics Technical Report, 2005, Available: http://www.ramas.com/bayes.pdf. Last accessed 21.08.2010.

[12] C. D. Buckingham, A. Ahmed, and A. E. Adams, "Using XML and XSLT for Flexible Elicitation of Mental-health Risk Knowledge", Medical Informatics and the Internet in Medicine, vol 32(1), pp 65-81, 2007.

[13] C. D. Buckingham, "Psychological Cue Use and Implications for a Clinical Decision Support System", Medical Informatics and the Internet in Medicine, vol 27(4), pp. 237-251, 2002.

[14] O. Obembe and C. D. Buckingham, "Developing a Probabilistic Graphical Structure from a Model of Mental-Health Clinical Risk Expertise", in Knowledge-Based and Intelligent Information and Engineering Systems, 14th International Conference, KES2010, R. Setchi, I. Jordanov, R. J. Howlett and L. C. Jain, Eds. in press.

[15] P. Perez, "Markov Random Fields and Images", CWI Quarterly, vol 11(4), pp. 413-437, 1998.

[16] C. M. Bishop, Pattern Recognition and Machine Learning, Springer, 2006.

[17] C. D'Elia, G. Poggi, and G. Scarpa, "A Tree-Structured Markov Random Field Model for Bayesian Image Segmentation", IEEE Transactions on Image Processing, vol 12(10), October 2003, pp. 1259-1273.

[18] F. V. Jensen and T. D. Nielsen, Bayesian Networks and Decision Graphs, 2nd ed., Springer Publishing Company, Incorporated , 2007.

[19] K. Murphy and A. Nefian, "Embedded Graphical Models", Intel Research Technical Report, June 2001.

[20] M. Meila and M. I. Jordan, "Learning with Mixture of Trees", The Journal of Machine Learning, vol 1, pp. 1 – 48, 2001.

[21] M. Drton, "Discrete Chain Graph Models", Bernoulli, vol 15(3), pp. 736-753, 2009.

[22] S. L. Lauritzen and N. Wermuth, "Graphical Models for Associations between Variables, Some of which are Qualitative and Some Quantitative", Ann. Statist., vol 17, pp. 31-57, 1989.

[23] C. A. Jiang, T. Y. Leong, and K. L. Poh, "PGMC: A Framework for Probabilistic Graphic Model Combination", in C. P. Friedman, J. Ash and P. Tarczy-Hornoch, Eds. Proceedings of the American Medical Informatics Association Annual Symposium (AMIA) 2005, pp 370-374.

# A Tool for Experimenting with a Theorem Prover

Foteini Grivokostopoulou

School of Engineering Department of Computer Engineering & Informatics

University of Patras

6500 Patras, Hellas (Greece)

grivokwst@ceid.upatras.gr

Ioannis Hatzilygeroudis

School of Engineering Department of Computer Engineering & Informatics

University of Patras

6500 Patras, Hellas (Greece)

ihatz@ceid.upatras.gr

*Abstract*—**In this paper, we presented EX-ACT-P, a tool for experimenting with a theorem prover. The heart of EX-ACT-P is ACT-P (A Configurable Theorem-Prover). ACT-P is a resolution-based theorem prover, which has a unique characteristic: allows the user to configure its resolution control regime. EX-ACT-P is an extension of ACT-P that allows user to experiment with theorem proving aspects, such as: configure a suitable resolution control regime, translate and automatically solve problems from the TPTP library, create and prove his/her own problems, display the proof steps in both text-based and graphical way and give information relevant to the proof process. So, EX-ACT-P can be useful to students for learning and to tutors for teaching aspects related to use of logic as a knowledge representation and reasoning language, by creating the right cognitive models, as well as to researchers for experimenting with theorem proving in ACT-P. A small scale evaluation for students has shown promising results.**

*Keywords- Automated reasoning; Theorem prover; Resolution control strategies; TPTP library; Teaching logic assistant*

## I. INTRODUCTION

Logic is one of the fundamental topics taught in computer science and/or engineering departments. In most such departments, logic is taught as a means for constructing formal proofs in a natural deduction style. However, teaching logic, especially first-order logic (FOL), as a knowledge representation and reasoning (KR&R) vehicle is also basic in all introductory artificial intelligence (AI) courses. We have constructed some tools for helping students in learning and tutors in teaching logic as a KR&R language (e.g., tools for translating natural language sentences into FOL ones or for converting FOL sentences into clause form, etc) [6]. However, we haven't constructed a tool for assisting in learning or teaching automated logic-based reasoning. There have been tools assisting students in using logic as a natural deduction tool, usually called proof assistants or proof editors (e.g., like [1, 2, 3, 4, 7, 8]). Most of them give an emphasis on the user interface. Some of them are built on top of an interactive theorem prover, like Isabelle [11]. However, there are no tools for helping students in learning and teachers in teaching automated

logic-based reasoning aspects. Even further, there are no tools for experimenting with controlling automated logic-based proof processes.

Theorem proving is a subfield of automated reasoning where logic is used as a KR&R vehicle. Resolution-based reasoning is a fundamental mechanism for logic-based reasoning, namely theorem proving. Many automated theorem provers (ATPs) are based on that mechanism. One of them based on first-order logic (FOL) is Prover9 [9]. So, such an ATP system could be the basis for a system assisting in teaching or generally in experimenting with automated logic-based reasoning.

We introduce here a tool for experimenting with logic-based reasoning via an ATP, namely ACT-P (A Configurable Theorem-Prover). ACT-P was initially introduced in [5] and is configurable in the sense that one can (re)define its resolution control regime. We call our tool EX-ACT-P (EXtended ACT-P), since it extends the interactive and presentation facilities of ACT-P.

The structure of the paper is as follows. Section II presents an overview of ACT-P. In Section III, EX-ACT-P is described. Section IV deals with using EX-ACT-P, whereas Section V briefly refers to system evaluation. Finally, Section VI concludes the paper.

## II. ACT-P: AN OVERVIEW

ACT-P (A Configurable Theorem-Prover) is the heart of EX-ACT-P. It is a skeleton resolution-based theorem prover, where specific steps are programmable by the user. ACT-P is based on a meta-level architecture. This means that it has an object-level language and a meta-level language. Its object-level language is a classical FOL. Its syntax is based on Cambridge Polish notation, using as connectives {~, &, V, =>} instead of {¬, ∧, ∨, ⇒} respectively and as quantifiers {forall, exists} instead of {∀, ∃} respectively. For example, the following are two well-formed ACT-P expressions:

((forall ?x) ((exists ?y) (=> (& (hunter ?x) (animal ?y))
　　　　　　　　　　(kills ?x ?y))))
(~ ((exists ?x) (V (cat ?x) (~ (kills-mice ?x)))))

which correspond to the following conventionally notated formulas:

$$(\forall x)\,(\exists y)\,(\text{hunter}(x) \wedge \text{animal}(y)) \Rightarrow \text{kills}(x, y)$$
$$\neg\,(\exists x)\,(\text{cat}(x) \vee \neg\text{kills-mice}(x))$$

In the object-level, ACT-P is a classical theorem prover. It can accept a set of FOL sentences (the axioms), which converts into clause form (CF) resulting in a set of clauses. Then it can accept a FOL sentence (the theorem) to be proved from the axioms. To that end, ACT-P negates it, converts it into its CF and tries to prove it using binary resolution refutation. It finally returns the solution (T or variable bindings) or 'no-solution'.

ACT-P internally represents its search space as an OR tree. Each branch of the tree represents a resolution step. Each node of the tree represents the updated set of clauses with the produced resolvent after execution of the incoming branch's step. We store clauses using the discrimination-tree indexing [10], which is an improvement to the version in [5]. This has significantly improved the efficiency of ACT-P.

The meta-level language of ACT-P is Common Lisp extended with a number of meta-primitives. Meta-primitives are predefined Lisp functions that implement a default resolution control regime and help users in implementing other resolution control strategies. That is, the user can redefine the bodies of the meta-primitives, which results in resolution control regime changes. Thus, a large number of different control strategies can be implemented.

## III. EX-ACT-P

### A. An Overview

We have extended ACT-P into EX-ACT-P (EXtended ACT-P) to be able to (a) experiment with and (b) practice with in a friendly manner. The first direction concerns tutors (even researchers), whereas the second concerns students. So, tutors/researchers can experiment with EX-ACT-P in order to test new possible exercises/problems related to KR&R or try new (combinations of) resolution strategies or test efficiency of them or even ACT-P itself. Students can practice by comparing their hand-made proofs with the ones provided by the system, see the differences when trying different strategies, study the steps of a proof etc. Students can practice with proofs that can't be made by hand, because they take many steps.

More specifically, one can
- Edit a problem in ACT-P FOL language
- Convert a TPTP library problem into ACT-P FOL language
- Determine different combinations of strategies
- Produce an automated proof
- See a linear text-based proof
- See a graphical representation of a proof

- See the number of produced clauses and the time required for a proof

### B. System Architecture

The architecture of EX-ACT-P is illustrated in Fig. 1. It consists of six main units: the *user interface* (UI), the ACTRANS, the *problem editor* (PE), the ACT-P, the *problem collection* (PC) and the *control strategies pool* (CSP).
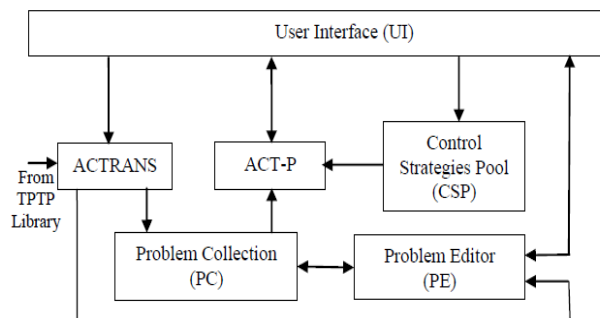


Figure 1.    EX-ACT-P Architecture.

PE is used to code a problem in ACT-P FOL language. Then can be added to PC, where it is available for experimentation. ACTRANS (Automated Code TRANSformer) is used for transforming TPTP library problems [12] into ACT-P language. Then the problems can be directly put into PC or be edited first through PE and then stored into PC. From CSP the user can select a proper combination of resolution control strategies to determine the overall control regime of ACT-P.

ACT-P is used to produce the proof for a selected problem from PC. Finally, UI is the means of interacting with the rest units of the system and presenting proof results.

### C. ACTRANS

ACTRANS (Automated Code TRANSformer) is a very useful tool of EX-ACT-P. It takes as input a TPTP problem file (e.g., puz006-1.p), through a direct link to TPTP, and creates two Lisp files. The first (PUZ006-1.lsp) contains whatever the TPTP file contains except that the problem formulas are converted into ACT-P language. The second (PUZ006-1PR.lsp) contains the query of the problem, i.e. the theorem to be proved.

Let see an example. Following are the axioms and the theorem representing problem "puz002-1" from Puzzle category in the TPTP library:

```
cnf(only_cats_in_house,axiom,
    ( ~ in_house(Cat)
    | cat(Cat) )).
cnf(gazers_are_suitable_pets,axiom,
    ( ~ gazer(Gazer)
    | suitable_pet(Gazer) )).
```

```
cnf(avoid_detested,axiom,
    ( ~ detested(Detested)
    | avoided(Detested) )).
cnf(carnivores_are_prowlers,axiom,
    ( ~ carnivore(Carnivore)
    | prowler(Carnivore) )).
cnf(cats_are_mice_killers,axiom,
    ( ~ cat(Cat)
    | mouse_killer(Cat) )).
cnf(in_house_if_takes_to_me,axiom,
    ( ~ takes_to_me(Taken_animal)
    | in_house(Taken_animal) )).
cnf(kangaroos_are_not_pets,axiom,
    ( ~ kangaroo(Kangaroo)
    | ~ suitable_pet(Kangaroo) )).
cnf(mouse_killers_are_carnivores,axiom
,
    ( ~ mouse_killer(Killer)
    | carnivore(Killer) )).
cnf(takes_to_me_or_detested,axiom,
    ( takes_to_me(Animal)
    | detested(Animal) )).
cnf(prowlers_are_gazers,axiom,
    ( ~ prowler(Prowler)
    | gazer(Prowler) )).
cnf(kangaroo_is_a_kangaroo,axiom,
    ( kangaroo(the_kangaroo) )).
cnf(avoid_kangaroo,negated_conjecture,
( ~ avoided(the_kangaroo) ))
```



Figure 2.   ACTRANS user interface.

After passing through ACTRANS, they are converted to the following ACT-P sentences:

```
(premise '(V (~(in_house ?Cat)) (cat
?Cat)))
(premise '(V (~(gazer ?Gazer))
(suitable_pet ?Gazer)))
(premise '(V (~(detested ?Detested))
(avoided ?Detested)))
(premise '(V (~(carnivore ?Carnivore))
(prowler ?Carnivore)))
```

```
(premise '(V (~(cat ?Cat)) (mouse_killer
?Cat)))
(premise'(V (~(takes_to_me ?Taken_animal))
(in_house ?Taken_animal)))
(premise '(V (~(kangaroo ?Kangaroo))
(~(suitable_pet ?Kangaroo))))
(premise '(V (~(mouse_killer ?Killer))
(carnivore ?Killer)))
(premise '(V (takes_to_me ?Animal)
(detested ?Animal)))
(premise '(V (~(prowler ?Prowler)) (gazer
?Prowler)))
(premise '(kangaroo the_kangaroo))
(setf query '(~(~(avoided the_kangaroo))))
```

All above sentences, except the last one, are included in the PUZ002-1.lsp file, while the last one is included in the PUZ002-1PR.lsp file. Figure 2 depicts the user interface of ACTRANS.

### D.   Control Strategies

We have categorized resolution control strategies according to the heuristics they use to reduce the search space and control resolution refutation process. That categorization helped us to design the reasoning cycle of ACT-P. We distinguish three major classes of heuristics. The first major class, *resolution restricting* strategies, concern generation of the search space and are used to increase efficiency by restricting the size of the search space. Resolution restricting strategies comprise two subtypes, namely *parent selection* strategies and *clause elimination* strategies. Parent selection strategies are based on the observation that not all possible resolvents have to be constructed to be able to derive the empty clause. They therefore impose restrictions on the clauses to be selected for resolution. Parent selection strategies have also been called 'refinement strategies' or 'restriction strategies'. The second type of resolution restricting strategies, clause elimination strategies, aim to eliminate clauses that will not be useful in further search. They are also called 'simplification strategies' or 'deletion strategies'.

The second main class of resolution control heuristics, *resolution search* strategies, concern the way the search space is searched. We distinguish between general search strategies and resolution ordering strategies. *General search* strategies traverse the search space in a blind way, without taking into account any resolution or domain or problem specific knowledge. Unlike general search strategies, *resolution ordering* strategies, aim to increase the efficiency of the theorem proving process by judiciously ordering potential resolutions. Clearly, such ordering strategies presuppose some ordering criterion. Best-first type strategies belong to this class.

Whereas resolution search strategies concern the way in which the search space is searched, our third main class of heuristics, *process oriented* strategies, concern the way in
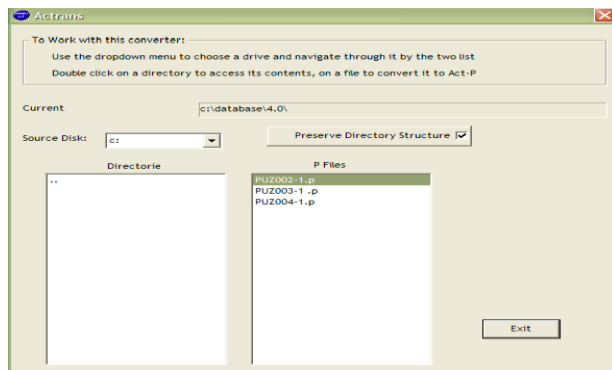
which resolutions are performed. They are typically used in conjunction with a specific parent selection strategy. We draw a distinction between *resolving preparation* and *resolving operation* strategies. The former concern operations on clauses before the actual application of the resolution rule. Resolving operation heuristics are used during actual resolution.

We have implemented a number of strategies from each category and stored in CSP. A user can select an appropriate combination of strategies from some or all the categories, through UI, to determine the overall control regime of ACT-P. A parent selection strategy should be first selected, since it will be the core strategy.
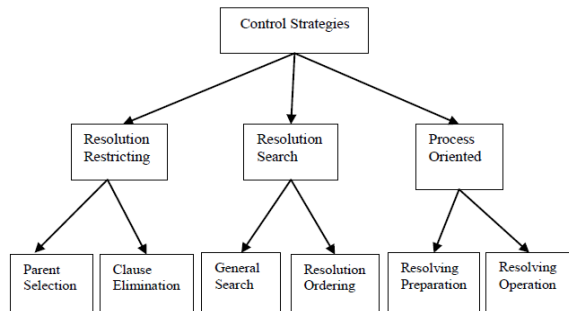


Figure 3.   Classification of resolution control heuristics.

### E.   Implementation Issues

EX-ACT-P has been developed in Lispworks 4.4. The choice of Lispworks as the programming environment was made on the basis of its features, such as its support for projects that are complex or need rapid prototyping and delivery. The user interface of ACTRANS has been developed in MS Visual Studio, an environment for constructing graphical user interface

## IV.   USING EX-ACT-P

The main scenario for using EX-ACT-P is as follows:
1.   Connect to TPTP and choose a problem file
2.   Use ACTRANS to transform it into corresponding ACT-P files
3.   Alternatively, choose a problem form PC or create a new problem via PE
4.   Choose a combination of resolution strategies
5.   Call ACT-P to activate proof process and produce the proof file
6.   See the displayed proof information
7.   If satisfied, stop; otherwise, go to any of steps 1, 3 or 4.

Of course, one can start from step 3, by choosing a problem from PC or creating his/her own problem files.

In step 3, the user can determine a combination of strategies from some of the drop-down menus, illustrated in Table I.

A combination of strategies can include, for example, a parent selection, a clause elimination and a search strategies, or any two of them. Another type of combination can include a process oriented strategy and a search strategy etc.

The screenshot in Fig. 4 presents the situation after completion of step 5. In the main window, the problem description file in ACT-P language is presented. Also, the proved theorem and proof information (number of axioms, number of produced clauses, number of used clauses and the CPU time required) are also displayed. At the right-hand side the selected problem and the chosen strategies are also displayed.

TABLE I.       RESOLUTION CONTROL STRATEGIES IN EX-ACT-P

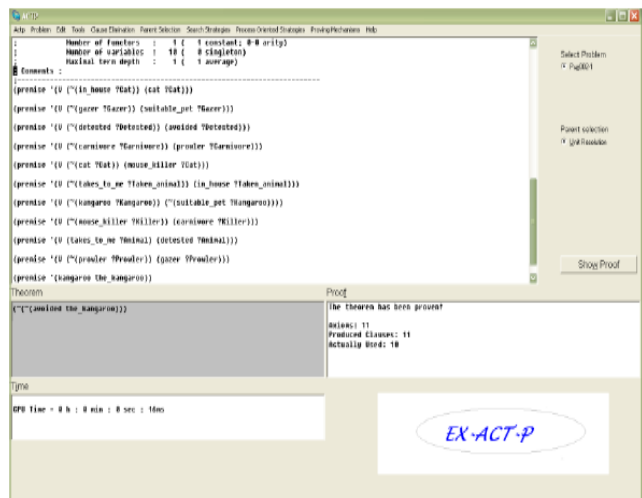| Resolution Control Strategies | | | |
|---|---|---|---|
| *Parent Selection* | *Clause Elimination* | *Search* | *Process Oriented* |
| Input Resolution<br><br>Linear Resolution<br><br>Linear Input Resolution<br><br>Set of Support Resolution<br><br>Unit Resolution<br><br>P1 Resolution<br><br>N1 Resolution<br><br>Hyperresolution | Tautology Elimination<br><br>Pure Literal Elimination<br><br>Backward Subsumption<br><br>Forward Subsumption | Fewest Literals Preference<br><br>Unit Preference | Ordered Resolution<br><br>OI-Resolution<br><br>Lock Resolution<br><br>Model Elimination |



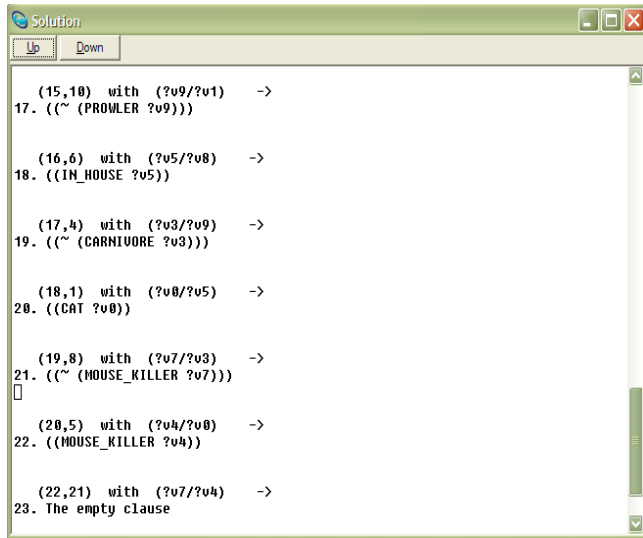Figure 4.    A screenshot of EX-ACT-P in use.

Figure 5.   A screenshot of the Proof window.

By clicking on the "Show Proof" button, two windows, namely the Proof and the Proof Graph, appear. In Figure 5, a snapshot of the Proof window is presented. The Proof window contains step by step the proof process in a text-base style.
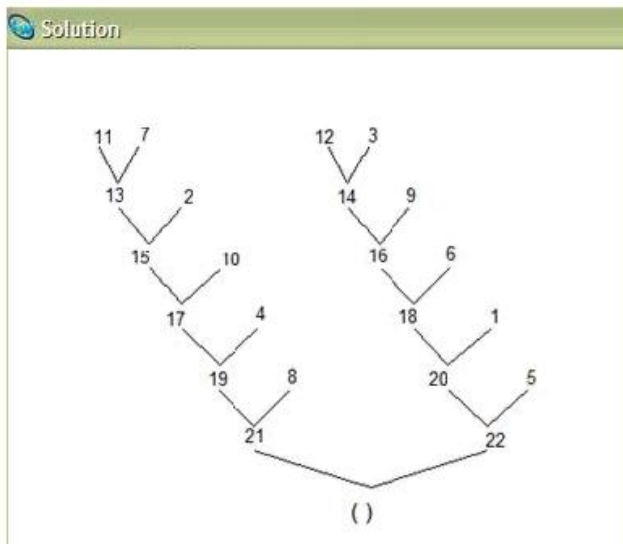


Figure 6.   A Screenshot of the Proof Graph window using "Unit Resolution"

Figure 6 and Figure 7 present two snapshots of the Proof Graph window, after having used two different control strategies. Figure 6 presents a proof graph of the problem "puz002-1" after having used the "Unit Resolution" strategy, whereas Figure 7 the one after having used "P1-resolution'. The user can now have a better view for understanding and comparing proofs. In this case, ACT-P solves the problem with fewer steps by using "P1

Resolution" than by using "Unit Resolution". The numbers in Figures 6 and 7 represent corresponding clauses (which are resolved in each step).



Figure 7.   A  Screenshot of the Proof Graph window using  "P1 Resolution"

V.    SYSTEM  EVALUATION

The system was used by the class of the Artificial Intelligence course, in our Department, which consisted of twenty-five senior computer engineering students. The students had been taught about logic as a knowledge representation and automated reasoning during the course lectures. They were instructed to use the system. Then, they were asked to fill in a questionnaire including questions for evaluating usability and learning.

TABLE II.        QUESTIONNAIRE RESULTS

| Q | Questionnaire | | | | | |
|---|---|---|---|---|---|---|
| | Questions | Answers (%) Total Students 25 | | | | |
| | | 1 | 2 | 3 | 4 | 5 |
| 1 | How much did the system help you to learn  automated reasoning ? | 0 | 0 | 20 | 36 | 44 |
| 2 | How much did the system help you to learn about resolution control strategies? | 0 | 4 | 24 | 32 | 40 |
| 3 | Will you suggest the system to next year students? | 0 | 0 | 16 | 32 | 52 |
| 4 | Will you use the system in the future to prove  problems? | 0 | 4 | 8 | 42 | 44 |
| 5 | Did you find the interface easy to use? | 0 | 4 | 32 | 48 | 16 |

The questionnaire included eight questions. The questions 1-5 were based on Likert scale (1: not at all, 5:

very much). The result is shown as Table 1. Finally, questions 5-8 were of open type and concerned strong and weak points or problems faced in using the system. Twenty five students filled in the questionnaire. Their answers showed that the students in general were helped in learning automated reasoning with the system (Table II). Also, they found that the user interface is easy to use. On the other hand, 72% the students agreed that the system helped them in learning the resolution control strategies.

## VI. CONCLUSION AND FUTURE WORK

In this paper, we present EX-ACT-P system, a tool for experimenting with an ATP system, namely ACT-P. Experimentation can be at student, tutor or even researcher level. A student can try to solve as many problems as required, compare hand-made solutions with automatic ones, see graphical representations of solutions and even try his/her own problems. A tutor can test any candidate exercises related to FOL based reasoning, check the complexity of an exercise, gain experience from using different strategies combinations etc. Finally, a researcher can try any problem from the TPTP library, even difficult ones, experiment with different combinations of strategies and see which combinations are more appropriate for which category of problems etc. We have used ACT-P in a number of difficult problems from the TPTP library. The results are more than promising, but we want to re-check it.

EX-ACT-P can be improved in a number of ways. A first direction is the enhancement of the CSP, both in number of implemented strategies and its structure. At the moment, the structure of the CSP and corresponding interface menus do not absolutely reflect the categorization of the strategies of Fig. 3. Furthermore, the categorization of Fig. 3 is not final. It can be more sophisticated (like the one in [5]).

A more interesting direction is to construct an intelligent system advising the user about which combinations of strategies have or have no meaning. Finally, a large scale experimentation with problems of various categories from TPTP problem library targeting at comparing various combinations of strategies, could create a guide for determining control regimes in resolution-based ATP systems.

## REFERENCES

[1] A. Asperti, C. Sacerdoti Coen, E. Tassi, and S. Zacchiroli. User interaction with the Matita proof assistant. Journal of Automated Reasoning, 2007. Special Issue on User Interfaces for Theorem Proving.

[2] D. Aspinall. Proof General: A generic tool for proof development. In Susanne Graf and Michael Schwartzbach, editors, Tools and Algorithms for the Construction and Analysis of Systems, Lecture Notes in Computer Science 1785, pages 38–42. Springer, 2000

[3] W. Billingsley and P. Robinson. Student Proof Exercises Using MathsTiles and Isabelle/HOL in an Intelligent Book. Journal of Automated Reasoning, 39:181–218, 2007.

[4] R. Bornat and B. Sufrin.. Jape: A calculator for animating proof-on-paper. In Mc-Cune, W. (ed.) Automated Deduction - CADE-14. LNCS, vol. 1249, pp. 412–415, 1997. Springer, Heidelberg

[5] I. Hatzilygeroudis and H. Reichgelt. ACT-P: A Configurable Theorem Prover, Data and Knowledge Engineering 12 (1994) 277-296.

[6] I. Hatzilygeroudis, C. Giannoulis and C. Koutsojannis. A Web Based Education System for Predicate Logic. Proceedings of the IEEE International Conference on Advanced Learning Technologies (ICALT'04), 106-110.

[7] C. Kaliszyk, F. Wiedijk, M. Hendriks and F. van Raamsdonk. Teaching logic using a state-of-the-art proof assistant, in: H. Geuvers and P. Courtieu, editors, International Workshop on Proof Assistants and Types in Education (PATE'07), 2007.

[8] C. Kaliszyk, F. van Raamsdonk, F. Wiedijk, H. Wupper, M. Hendriks & R. de Vrijer, "Deduction using the ProofWeb system", report ICIS-R08016, Radboud University Nijmegen, 2008

[9] W.McCune. Prover9: Automated theorem prover for first-order and equational logic, 2008. http://www.cs.unm.edu/~mccune/mace4/manual-examples.html.

[10] W. McCune. Experiments with discrimination-Tree and Path Indexing for Term Retrieval. Journal of Automated Reasoning 9:147-167, 1992

[11] T. Nipkow, L. C. Paulson and M. Wenzel. Isabelle/HOL - A Proof Assistant for Higher-Order Logic, Lecture Notes in Computer Science 2283, Springer, 2002.

[12] G. Sutcliffe and C. B. Suttner. The TPTP Problem Library: CNF Release v1.2.1. Journal of Automated Reasoning, 21(2): 177-203, 1998.

# Algorithm to Solve Web Service Complex Request Using Automatic Composition of Semantic Web Service

Brahim Batouche, Yannick Naudet,
Public Research Center Henri Tudor,
Luxembourg
{brahim.batouche, yannick.naudet}@tudor.lu

Frédéric Guinand
University of Le Havre,
France
frederic.guinand@univ-lehavre.fr

*Abstract* -- **Automatic composition of web services supports the solving of complex user request. The set of possible solutions can be represented by a graph, modeling the composition. Usually, this kind of approach is highly simplified by considering only sequences of services. This paper proposes an algorithm for automatic semantic web services composition, which generates a graph taking into account any composition structure. The request resolution process identifies possible composition structures and selects relevant services based on their semantic description. The resulted composition graph answers all requested functionality with coherent composition structures.**

*Keywords - semantic web service; composition graph; automatic compositio; web service composition structure.*

## I.    INTRODUCTION

Web services composition is a classical approach to answer complex queries that cannot be solved with one single service. Answering to such requests requires several steps: (1) finding suitable services; (2) finding how they can be composed together to answer the request; (3) create the corresponding composite service; (4) invoke it; (5) maintain it so that it can be reused later. The structure of the composite service depends obviously on the request, but also of the available services.

Composing services can be useful in many different domains, such as, e.g., tourism, transport, multimedia, etc. Some of them involve a dynamic environment where events at any time can affect previously computed compositions answering a request. A fundamental issue is then how to repair failures in a composite service execution, which can occur in dynamic environments. A typical example is when one of the services involved in the composition is faulty or can no more be executed. This fail is translated to a complex request and then use our algorithm to find another composition alternative.

In this paper we propose an algorithm for automatically finding all candidate compositions answering a complex request, without a priori knowledge of the composition structure. When the request does not formally specify any chaining between the requested elements, the algorithm must found suitable composition structures based on the available services. This problem is not trivial because there are many possible services combinations and composition structures. To determine the composition structure we base in the existing functionalities, which are automatically determined because the available services are supposed described semantically by OWL-S [1].

In section 2, we present related works. In Section 3, we first formalize the problem and detail it. Section 4 presents the composition structures and their semantics. In Section 5, we present our algorithm for automatic construction of a composition graph. In Section 6, we provide the experimentation results, and finally, we conclude in Section 7.

## II.    RELATED WORKS

Solving a complex request by services composition in dynamic environments, can be tackled by different approach.

The algorithm presented in [2] builds a composition graph answering a request. The algorithm identifies first the input and output of the request and search for a matching service. When none can be found a service having a matching output is selected and recursively, subsequent services having output matching with the input of the latter service and input matching with the request input are sought. The algorithm ends when a sequence of services starting with the request input and ending with its output is found, or when the set of available services has been searched. The provided composition graph does not allow the direct invocation of services. Also, it is still limited to sequences structure of composition.

In [3], the flooding algorithm is used. Once again, the proposed approach is limited to services sequences and does not allow composition execution.

In [3], an architecture for automatic web service composition is proposed. This architecture allows fast composition of OWL-S service. However, while authors

provide interesting ideas for the design of the composite service and automating service invocation, they only consider sequences of services. Kazhamiakin and Pistore [5] proposed a model to answer a request using the composition of web services, the model supposes the available services are described by BPEL-WS. A request requires much functionality, which are identified and then used in a finite state-machine, implementing the composition structure. The state-machine provided does not allow composition execution.

In [5], a multi-agent system is used to automate the composition of services. The agents collaborate to provide the composition needed, an agent is presented by the OWL-S service and its functional parameters describe the agent role. With this system, we can consider the compositional structures: sequence, parallel and conditional. The conditional structure concerns only the functional parameter of the service.

According to the state of the art, many methods of automatic composition focus to find the needed functionalities to answer a request and there order but do not give the link to execute them. So, they usually consider only the sequence structure. To exceed these limits, we automate the detection of the composition structure needed and the selection of the services requested. To select automatically the services, we use the I/O dependence basing in the matching function. The matching function uses only the IO parameter [6] or uses the IOPE [8], which provides more accurate results.

### III. PROBLEM FORMALIZATION

A typical example of complex user request, which we will use as a basis to present our approach, is the following: "I want travel from City A to City B, reserve several hotel rooms in destination city where each book is billed separately, rent a car for six people, have the weather and plan for the destination city". Such request needs first to be formalized in a machine processable way.

#### A. Request Formalization

A complex request is a combination of more focused or atomic sub-requests, which concerns each a service or functionality. We write: $R = F_R = \{F_{r_i}\}$, where $F_{r_i}$ is a functionality requested. Our example requires four functionalities: transport, booking hotel, rent a car, city information. Each functionality has input/ output $(I_{F_{r_i}}/O_{F_{r_i}})$. Formally, we write a request as a triple: $R = <I_R, O_R, C>$, where $I_R = \bigcup_{i=1}^{card(F_R)} I_{F_{r_i}} = (i_1, \dots, i_k)^T$ is

the set of inputs, $O_R = \bigcup_{i=1}^{card(F_R)} O_{F_{r_i}} = (o_1, \dots, o_l)^T$ is the set of outputs, and $C = (c_1, \dots, c_m)^T$ is the set of conditions or constraints (related to data, service or composition). Conditions differ from constraints in that they must be verified to instance the input parameter of service, but the constraints to filter the set of available services, data provided by services or composition paths.

All the sets elements are URIs of concepts defined in ontologies. While $I$ and $O$ correspond to functional parameters which describe a domain ontology, $C$ concerns both functional and non functional parameters. Quality of Service is an example of such parameter, as well as price, cardinality of some services output, etc. The non functional parameters values are found in the service description. The functional parameters values are identified after execution the informative service, which provide information without modify its source database.

- From the fail execution of service to request

The execution of composition can fail if one of its services fails. The fail can then be translated to a new request, which depends on the functionalities realized at the moment of fail. These functionalities correspond to a set of Terminated Input / Output (*TI/TO*). The new request formalized as $R = < \overline{I_R \cap TI}, \overline{O_R \cap TO}, C >$. To configure the original composition graph (or find another alternative) we use reclusively our algorithm with the new request.

#### B. Composition Graph Formalization

The composition of services presents a set of functionalities and there structures, but usually does not give the execution possibility, e.g. [5] [9]. This brings us to define two types of composition graph.

**Definition:** ***The executable composition*** *graph allows the service execution, thus making the composition executable.* ***The abstract composition*** *graph represents only the structure of the composition and cannot be executed.*

This definition based in the existing (or not) the link to invoke the service. But [10] defines the abstract composition according to the existing (or not) the sub-service I/O of the composition. According to our definitions, an executable composition graph is an abstract composition graph whose nodes integrate services identifiers, (URIs), instead of input / output parameters only. Abstract graphs represent only composition of functionalities fitting a request, while

executable ones describe actual services chaining. We focus here on finding automatically suitable services composition structures that we model with an executable graph.

The executable composition graph corresponding to a complex request; represents a set of services paths which constitute possible answers. It is formulated as: $G = <N, V>$, where, $V$ is a set of directed arcs and $N$ is a set of nodes. We distinguish four types of nodes, *IS, AS, DA and SW,* where *IS* is an informative service, *AS* is an active service, which provides an action and modify its source database; *DA* is the data (information) provided by an *IS;* and *SW* is a switch node that represents a conditional structure, specifying a condition formula.

We define a node as $n = <NT, id, URI_S, I_s, O_s, URI_{DATA}>$, where $NT$ is the node type, $id$ is the identifier of starting parallel structure node ($id = \emptyset$ if $n$ does not belong to a parallel structure), $URI_S$ is the URI of OWL-S service ($URI_S = \emptyset$ if $NT = DA$ *or SW*), $I_s/O_s$ are respectively the Input and output of the service, they are defined from $URI_S$, and $URI_{DATA}$ is the URI of data ($URI_{DATA} = \emptyset$ if $NT = IS, AS$ *or SW* ). The special node switch $n_{SW} = <NT = SW, c_i, LF>$, where $c_i$ is the condition provided by the request and it is verified by the node, and $LF$ is the linked node because the multiple paths in the graph can meet in one SW node, then a SW node embeds a hash function recording authorized successors of nodes. For this reason, the SW node is a kind of meta-node containing several nodes.

We add the node *SN* and *EN* which respectively starting and ending the composite graph, $SN = <O = I_R>$, which its output corresponds to the request input, and $EN = <\emptyset>$.

The functional parameters of the answer composite of services match with the functional parameters of the request. So, the non functional parameters values are calculated according to the parameter type and the composition structure used, for example, see [11] .

## IV. STRUCTURES OF WEB SERVICE COMPOSITIONS

Existing web services languages supporting composition model different structures in different ways. Taking the most commonly known, we observed the following. The structures modeled by OWL-S are: "sequence", "any-order", "if-then-else", "choice", "while", "until", "split" and "split-joint". Differently, BPEL4WS [12] uses: "sequence", "switch" "while" "Pick" and

"flow". A mapping between the two representations involves three operators: equivalence (e.g., if-then-else is equivalent to switch; choice is equivalent to pick); composition (e.g., the flow structure in BPEL4WS can be decomposed into two structures of OWL-S: split and split-joint); and identity (for constructs that cannot be realized with structures of the other representation, e.g., any-order is not identity (see Section 4.2)). In order to insure interoperability with the different representations and keeping a generic approach, we focus on elementary structures (sequence, if-then-else, split, split-joint), from which many others can be modeled.

### A. Composition Structures Illustration

A composition may comprise several different structures, which can themselves contain combinations of structures. A tree representation helps understanding and visualizing the composition: the leaves are services; the nodes and the root are the compositions structures. The path corresponds to read of composition tree which follow a prefixed depth approach. The Figure 1 shows for our example the composition tree and the corresponding composition flow. The used services are: available train (AT), available flight (AF), book train (BT), book flight (BF), available hotel (AH), book hotel (BH), available rentals car (ARC), rent car (RC), plan touristic map (PT), city weather (CW).
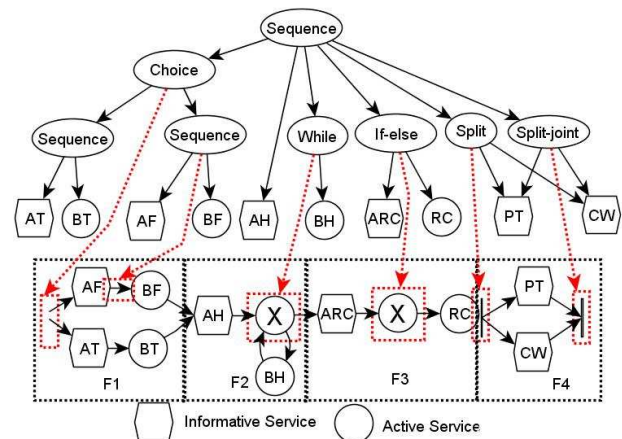


Figure 1: Service composition illustrated by tree and flow

### B. Characteristic of Composition Structures

In the following, we detail the characteristics of structures we retained and explain how they identified from the request.

*Sequence* "→": This structure defines an order between services. The order can be detected directly or indirectly. There are two ways to detect the order directly: (1) - Checking the match between services IOPE; (2) – Checking the priority between the services answering the question: *which service cancels the other when it is cancelled?*

Since the order operation is transitive: $(A \rightarrow B) \wedge (B \rightarrow C) \Longrightarrow (A \rightarrow C)$. To detect indirectly the order between A and C. We base on the order of services (e.g., B) which have the order with A, C.

*Choice* "+": (or or-split): This structure represents a choice between several services that have a same functionality. $Choice(A, B_1, B_2, ..., B_k) \equiv (A \rightarrow B_1) \vee (A \rightarrow B_2) \vee ... \vee (A \rightarrow B_k)$, knowing that service "A" precedes services $B_i$ and the services $B_i$ have not the different functionality.

*Any-Order* "⊙": This structure is not elementary and represents a random invocation of services. This structure can be expressed using choice and sequence structures: $A \odot B \equiv (A \longrightarrow B) + (B \longrightarrow A)$. Therefore this structure is replicable.

*If-then-else* "$\otimes_c$": This structure checks a condition of request to instance the functional parameter of the service. The structure follows a service if ones of its parameter (functional/ non functional) correspond to a condition.

*Split* "⊢": This structure indicates a simultaneous start of multiple services (or services chains). Services that can be parallelized have the same predecessor and provide different types of outputs. Each service starts a new sub-path in the composition. All services chains starting at a split will be executed in parallel and ended with a split-joint. $Split(A, B_1, B_2, ..., B_k) \equiv (A \rightarrow B_1) \wedge (A \rightarrow B_2) \wedge ... \wedge (A \rightarrow B_k)$.

*Split-joint* "⊣": This structure ends a parallel structure, where the sub-composition paths belong to a same "split". The last services $B_i$ in parallel chains have the same successor $A$. $split - Joint (B_1, B_2, ..., B_k, A) \equiv (B_1 \rightarrow A) \wedge (B_2 \rightarrow A) \wedge ... \wedge (B_k \rightarrow A)$, where services "$B_i$" end the parallel sub-composition paths. It is possible that all services chains in a same "split" do not end in the same "split-joint".

*While* "$\circledast_c$" and *until* "$\odot_c$": These structures are not elementary and used for iterative service invocation. They can be constructed with if-then-else and sequence structures: $\circledast_c (A) = \otimes_c \rightarrow A \wedge A \rightarrow \otimes_c$ and $\odot_c (A) = A \rightarrow \otimes_c \wedge \otimes_c \rightarrow A$.

The compositional structures $(St)$ are illustrated in the graph by arc or node, $St = < N, V >$ (see Figure 1). The structures while and until $= \{SW, V\}$, the structure if-then-else$= \{SW, \emptyset\}$ and the structure sequence, choice, split and split-joint$= \{\emptyset, V\}$, where $V$ is respectively a sequence-arc, set of sequence-arc, split-arc and split-joint-arc. This gives to distinct three types of arc $\{\rightarrow, \vdash, \dashv\}$. An arc is defined by its type, departure node and destination node.

## V. ALGORITHM GENERATING THE COMPOSITION GRAPH

Our algorithm processes progressively the request to build the executable composition graph. In the following, we define our terminology.

We name in the graph *current layer* $l_k$, the set of nodes in the graph having a same depth level, currently being processed: $l_k = \{n_i\}$. Initially $l_k = \{SN\}$. One step of the algorithm corresponds to full covers of $l_k$. The node of $l_k$ being processed named *current node.*

The temporary buffer is used to store the set of nodes following the current node, and not preceding *EN*. When precede *EN* are placed directly in the *end layer* of the graph.

The algorithm is illustrated in Figure 2. From a request, it fills the current layer and processes it. For each node in the current layer, selects the next services according to their matching with functionalities $F_R$. The set of nodes created from next services is first put in the temporary buffer, which is later put in the next layer. When the current node output matches with one of the $F_{r_i}$ outputs, the algorithm carries on with next node in current layer. That has the input of a $F_{r_i}$ not yet covered. Otherwise, services having inputs matching the current node output are selected. Corresponding nodes are created after checking does not already exist in the set of nodes $N$ in $G$.

An arc-sequence is created between the current node and the next nodes. When a selected service is an IS, it is invoked to obtain the information it provides before creating the arc. When the data are obtained, the algorithm creates an arc sequence between the node of the service and data nodes created, then it replaces the service node by the set of data nodes.

When a next node has been newly created, the algorithm checks the existence of a condition. The next node contains a condition if the output of the service it represents corresponds to one of request conditions. In

this case we create a SW-node and linked to the node by an arc-sequence. The node following the SW-node is then selected according to the first node output.

In case all $F_R$ have been covered, the next node is affected to the *end layer*. Otherwise, it is put into the temporary buffer, and later to the next layer.

The checking of split-structures is performed when the temporary buffer is full, containing all the nodes matching the current node. The checking of split-joint-structure is performed when the next node is selected. Therefore, the algorithm checks split-joint structure before the split structure.

The process checks the existence of a split-joint structure starting from next node. If it is selected from *N*, then it is possible to find a node which can precede the next node. In this case a complete check is performed, otherwise only a partial check is necessary. The complete check considers all nodes of the current layer. The partial check considers a current node and current layer nodes which have not been yet processed.

The algorithm creates a split-joint-arc when the split-joint is verified, i.e., the follow conditions are verified:

-The starting nodes of the split-joint-arc have a same split-structure, i.e., they contain the same identifier of split structure $id(n')$, where $n'$ is the node starting the parallel structure.

-The nodes have the same succeeding node $n^+$, where $n^+$ ends the parallel structure.

$$\forall n_i, n_j \; prcced \; n^+, \textbf{\textit{if}} \; (n_i, n_j) \text{contains} \left(id(n')\right) \textbf{\textit{then}}$$

$$CreatArcSplitJoint \left(n_i, n_j, n^+\right), n_i.delet(id(n')), n_j.delet(id(n^*))$$

Concatenation of parallel structures is possible. When nodes of same split-structure don't regroup in a same split-joint structure, the node $n^+$is included in the structure split, so it can be grouped with the remaining nodes. $\exists \; n_i.contains\left(id(n')\right), n_i \in l_k : n^+.add(id(n'))$.

The checking for the existence of a split structure is performed between the current node and the nodes in the temporary buffer. If these nodes have different functionalities i.e., different output, then we create a split-arc and add the identifier of the split structure $id(n')$ to these nodes. $\forall \; n_i, n_j \; follow \; n', \textbf{\textit{if}} \; M\left(out_{n_i}, out_{n_j}\right) > \varepsilon \; \textbf{\textit{then}}$

$$n_j.add\left(id(n')\right), n_i.add\left(id(n')\right), CreatArcSplit(n', n_i, n_j)$$

When a node $n_k$ follows the node $n_j$ which appears in parallel structure, $n_j contains\left(id(n')\right)$, then we affect $n_k$ to this structure, $n_k.add\left(id(n')\right), n_j.delet(id(n'))$.

When all nodes of the current layer are processed, the next layer becomes the current layer and so on until the

next layer is empty. The algorithm terminates when this state is reached.
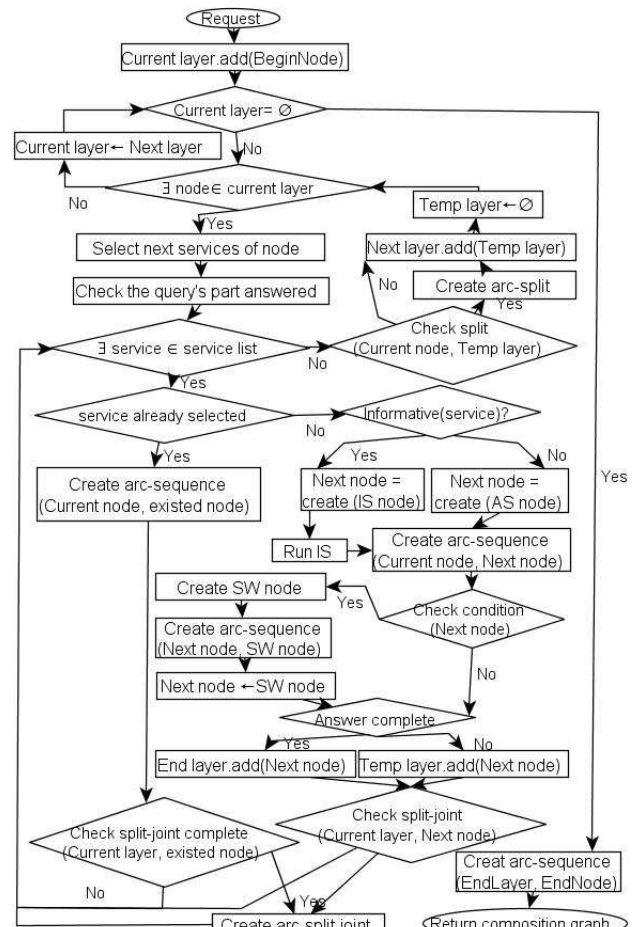


Figure 2: Algorithm of solving complex request

Finally, the complexity of each step of the algorithm graph construction composition is about $O\left(|l_k|.|S|\right)$, where $S$ is the set of selected service. Since, the algorithm is based in the flooding algorithm. To assure the process logic, we check the composition structures according to the flooding algorithm step.

## VI. EXPERIMENTATION AND RESULTS

We have tested our example request on a base, containing the services OWL-S descriptions and varying all the functionality needed in the request.

After running the algorithm, we verify: - the service composite answering a request has all requested functionalities, - its internal composition structure is coherent, i.e., there is no false detection of structures.

Different APIs were used Jena [13], SPARQL [14]; to check the data constraint and the conditions, OWL-S API [15], to check the service constraints, and Pellet [16], to

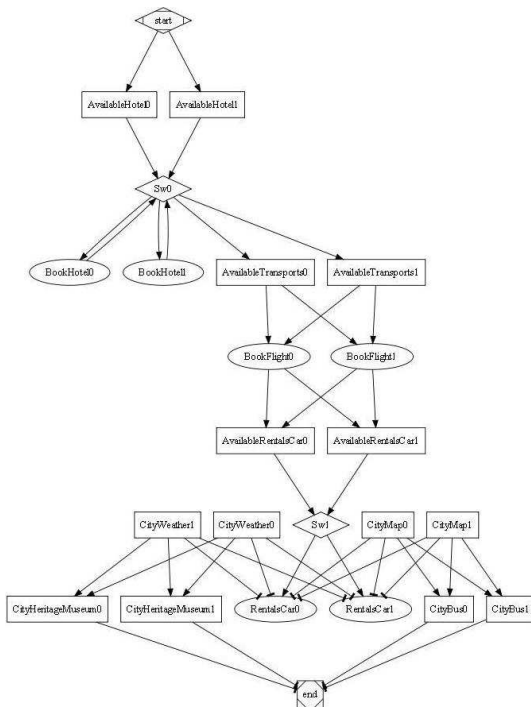check the matching level between services I/O and request I/O.



Figure 3: Resulted executable composition graph

The Figure 3 illustrates the composition graph given by the algorithm. The composition path is semantically correct because it contains all requested functionalities: transport, booking hotel, etc. And the composition structures used are coherent with the used services. E.g., choice: between the service "Available Hotel 0" and "Available Hotel 1". Sequence: between "Book Flight" and "Available Rentals Car". While: the node *sw0* checks the number of booking hotel. If the condition is true then another booking is made; else the loop is left. If-then else: the node *sw1* checks the number of car rented according to the type of car provided, if the available rental car does not take six people then rent two cars. Split/Split-joint: the service "City Heritage Museum" and "City Bus" will be executed in parallel.

## VII. CONCLUSION AND PERSPECTIVES

In this paper, we have proposed an algorithm for multi-structure web services composition. It allows answering a user request by composing available matching services using all possible composition structures.

The composition graph provided by the algorithm will mainly be used as input for giving a search space

authorized to optimize the composition of services. Additionally, we have also shown how to deal with composition execution failures (in this case, the composition graph can be adapted).

Finally, the solutions to a request contained in the composition graph can be formalized using classical languages like, e.g., BPEL-WS, OWL-S, etc., and stored in the services base for re-use.

In future works, we consider all ways to detect a sequence between services and we integrate the precondition/effects to calculate a level of matching between the services and request.

REFERENCES

[1] D Martin, et al.: OWL-S: Semantic Markup for Web Services. *W3C Member Submission*, 22, 2004.

[2] G. Silva, F. Pires, and V. Sinderen. An Algorithm for Automatic Service Composition, *1 st International Workshop on Architectures, Concepts and Technologies for service Oriented Computing.* pp. 65-74, Barcelona Spain. July 2007.

[3] S. Oh, B. On, E.J. Larson, and D. Lee. BF*: Web Services Discovery and Composition as Graph Search Problem, 6-8, *e-Technology, e-Commerce, and e-Services, IEEE International Conference on*, 784-786, 2005.

[4] K. Matthias and G. Andreas, Semantic web service composition planning with OWLS-XPlan, *In Proceedings of the 1st Int. AAAI Fall Symposium on Agents and the Semantic Web, pp. 55-62, 2005.*

[5] R. Kazhamiakin *a*nd M. Pistore, A Parametric Communication Model for the Verification of B*PEL4WS* Compositions*, Formal Techniques for Computer Systems and Business Processes,* 318-332, Trento, Italy. 2005.

[6] D. Pellier and H. Fiorino. Un modèle de composition automatique et distribuée de services web par planification, *Revue d'Intelligence Artificielle ,v23,13-46, 2009.*

[7] *M. Klusch, B. Fries, M. Khalid, and K. Sycara, OWLS-MX: Hybrid OWL-S Service Matchmaking, In Proceedings of 1st Intl. AAAI Fall Symposium on Agents and the Semantic Web*. 2005.

[8] A.B. Bener, V. Ozadali, and E.S.Ilhan. Semantic matchmaker with precondition and effect matching using SWRL. *Expert Systems with Applications*, 36, 9371-9377, 2009.

[9] S.V. Hashemian, and F. Mavaddat. A Graph-Based Approach to Web Services Composition. *Proceedings of Symposium on Applications and the Internet*, 183-189, 2005.

[10] M. Mihhail, M. Riina, and T. Enn. Compositional Logical Semantics for Business Process Languages. Pro of *ICIW*. 2007.

[11] C. Wan, C. Ullrich, L. Chen, R. Huang, J. Luo, and Z. Shi. On Solving QoS-Aware Service Selection Problem with Service Composition, *Grid and Cooperative Computing,* 2008.

[12] T. Andrews, et al: 'Business Process Execution Language for Web Services Version 1.1', IBM, May, 2003.

[13] http://jena.sourceforge.net/ (11/03/2010)

[14] http://www.w3.org/TR/rdf-sparql-query/ (11/03/2010)

[15] http://www.mindswap.org/2004/owl-s/api/ (11/03/2010)

[16] http://www.mindswap.org/2003/pellet/ (11/03/2010)

# Recognition and Understanding Situations and Activities with Description Logics for Safe Human-Robot Cooperation

Jürgen Graf, Stephan Puls and Heinz Wörn

Institute of Process Control and Robotics (IPR)

Karlsruhe Institute of Technology (KIT)

Karlsruhe, Germany

{graf, puls, woern}@ira.uka.de

*Abstract*—**Recognition of human activities and situation awareness is an important basis for safe human-robot-cooperation. In this paper, a recognition module is presented and discussed. The usage of Description Logics allows for knowledge based representation of activities and situations. Furthermore, reasoning about context dependent actions enables conclusions about expectations for robot behavior. This approach represents a significant step towards a full-fledged cognitive industrial robotic framework.**

*Keywords – cognitive robotics, Description Logics, ambient intelligence, situation and action recognition, human-robot cooperation.*

## I. INTRODUCTION

Industrial robotics is a challenging domain for cognitive systems, especially, when human intelligence meets solid machinery like most of today's industrial robots.

Hence, guaranteeing safety for human workers, safety fences are installed to separate humans and robots. As a consequence no real interaction or cooperation sharing time and space can be found in industrial robotics.

Some progress has gained in the past so that some modern working cells are equipped with laser scanners performing foreground detection. But with these systems one is not able to know what is going on in the in the scene and therefore could not contribute something meaningful for challenging tasks like safe human-robot cooperation.

We are conducting research on recognition of and reasoning about actions and situations in a human centered production environment, in order to enable interactive and cooperative scenarios.

This paper focuses on using Description Logics (DLs) [8] as means for representation of knowledge and as reasoning facilities for inference about activities and situations. Furthermore, conclusions about user expectations about robotic behavior can be drawn.

In Section II, some related research work on reasoning about scenes and situations will be presented. In Section III, the framework will be introduced, which enables the sensor data processing and subsequent knowledge based reasoning. In Section IV, DLs will be briefly introduced and the module realizing the communication with a Description Logics reasoner, knowledge base management and reasoner result management will be presented in detail. Also the modeled situations and activities are explained. Section V discusses experimental results which have been carried out for both, predetermined test cases and under real-life conditions. In Section VI, a summary is given. Finally, some hints for future work are also mentioned.

## II. RELATED WORKS

There are a lot of approaches for action recognition systems based on probabilistic methods, e.g., hidden Markov Models (HMMs) [16, 17, 18], as their theoretic foundation is well understood and applications in speech recognition have shown their capabilities.

Based on arguments, that HMMs are not suitable for recognition of parallel activities, instead propagation networks [19] have been introduced. The propagation network approach associates each node of the network with an action primitive, which incorporates a probabilistic duration model. Also conditional joint probabilities are used to enforce temporal and logic constraints. In analogy to HMMs, many propagation networks are evaluated, in order to approximate the observation probability.

In [20], arguments are put forward, that recognition of prolonged activities is not feasible based on purely probabilistic methods. Thus, an approach is presented which uses parameterized stochastic grammars.

The application of knowledge based methods for action recognition tasks is scarce, but work on scene interpretation using DLs has been conducted.

In [9], DLs are used for reasoning about traffic situations and understanding of intersections. Deductive inference services are used to reduce the intersection hypotheses space and to retrieve useful information for the driver.

In [10], scene interpretation was established using DLs. Table cover scenes are analyzed and interpreted based on temporal and spatial relations of visual aggregate concepts. The interpretation uses visual evidence and contextual information in order to guide the stepwise process. Additionally probabilistic information is integrated within the knowledge based framework in order to generate preferred interpretations. This work is widened to cope with general multimedia data in [11], in which a general interpretation framework based on DLs is presented.

In [12], a comprehensive approach for situation-awareness is introduced, which incorporates context capturing, context abstraction and decision making into a generic framework. This framework manages sensing devices and reasoning components which allows for using

different reasoning facilities. Thus, DLs can be used for high level decision making.

These last examples show that the usage of DLs bears great potential. Hence its adoption in the situation and action recognition task incorporated into the MAROCO framework.

To the best of our knowledge this is the first paper which incorporates description logics in the domain of cognitive robotics. For reasons of this, it was not possible to compare the runtime analysis results to concurrent research groups.

There are some investigations concerning runtime analysis of descriptions logic reasoners (see [21], e.g.) but they are far away from the robotics community and finally they show that the pellet system which was used in this publication is one of the best with respect to the given constraints of the software architecture of MAROCO.

The main motivation writing this paper is introducing the description logics approach into the domain of cognitive robotics. There are just a few other research groups which are dealing with description logics in a similar research domain and the most related ones were referenced in this paper. Most attention was spent on extending the cognitive robotic system MAROCO with description logics and building a knowledge base for action and gesture recognition.

The markerless tracking of a human body in real time is not at the core this paper. But this paper is the first which brings together markerless real time tracking of a human body together with a safe robot path-planning module and the description logic approach. Thus, this paper intends to present interesting results that are gathered from experimental investigations using description logics.

## III. THE MAROCO FRAMEWORK

The MAROCO (human robot cooperation) framework [2,3] is an implemented architecture that enables human centered computing realizing a safe human-robot interaction and cooperation due to advanced sensor technologies and fancy algorithms [6,7].
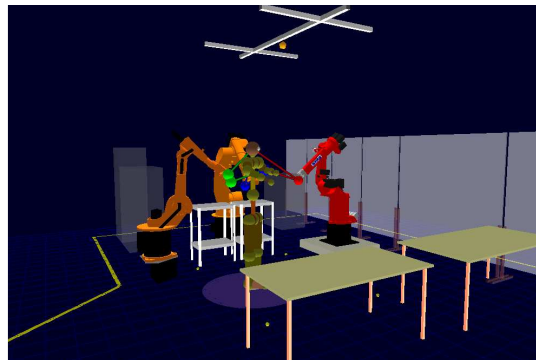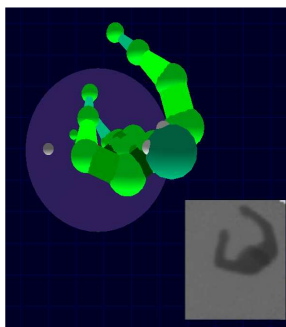


Figure 1. (Up) Reconstructed human model from depth images. (Down) Environmental scene model consisting of several kinematical chains. Three different industrial robots and a human model. All agents and robots have been reconstructed by MAROCO and are integrated into the virtual model in real-time including safety features extraction, risk estimation and path planning.

Every system implementing machine intelligence needs sensors. The MAROCO system analyzes image sequences that are gathered from a 3D vision system [1] based on time-of-flight principle which is mounted to the top of the ceiling of the working cell (see Fig. 1). Modules dedicated to image sequence analysis make it possible to estimate more than a dozen of kinematical parameters, e.g., head orientation, upper body orientation, arm configuration, etc., of a human model without using any markers (Fig. 1). The technical details of the methods realizing the real-time reconstruction of the kinematical model are not in the focus of this paper. Details can be found in [3,6,7].

As safety is one of the most demanding features when industrial robots get in contact with human workers, MAROCO is focused on estimating the risk for the human worker depending on the scene configuration. A variety of methods are integrated into the framework like pure functional evaluation, machine learning tools, e.g., support vector machines, and a two-threaded adaptive fuzzy logic approach, which at the moment makes the race [7].

Having estimated the risk, one is interested in finding a procedure minimizing the risk for both, the worker and machinery. Re-planning is an efficient tool minimizing the risk. A method for re-planning the path of the robot with respect to safety and real-time capability is presented in [4].

The kinematical model also allows for recognition of human activities and situations inside the robot working area. Using Description Logic (DL) reasoning facilities, conclusions about occurring situations, actions, their temporal relations and expectations about robot behavior can be drawn. This is what will be shown in the next sections.

## IV. THE RECOGNITION MODULE

This section is dedicated to discuss the recognition module including its components and modeled knowledge base after a very brief introduction to DLs.

### A. Description Logics

In this paper, DLs [8] are used to formalize knowledge about situations, actions and expectations. DL is a 2-variable
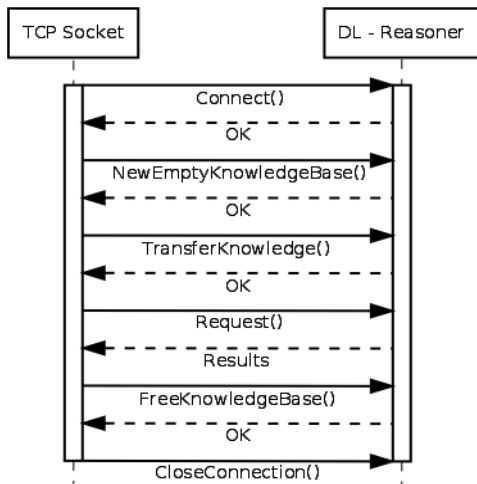
Figure 3. Communication between interface component and DL reasoner.

fragment of First Order Logic and most DLs are decidable. Thus, sound, complete and terminating reasoning algorithms exist.

A DL knowledge base is divided distinctly into general knowledge and knowledge about the individuals in the domain. The former defines the terminology of the domain and its axioms are declared in the terminology box, hence TBox. The latter defines assertions about individuals and, therefore, is declared in the assertion box, hence ABox. This allows for modular and reusable knowledge base and thus for more efficient coding of knowledge [9].

Due to DL's open world assumption, it can deal naturally with incomplete information, which is essential in reasoning about sensor data.

### B. The Module Design

The recognition module needs to fulfill at least the tasks of establishing a communication interface with the Description Logics reasoner, managing the knowledge base and managing the reasoner results.

In Figure 2, components of the module are presented. The communication via TCP and the XML parsing are done by the components marked as DIG-interface. The DIG-interface is a W3C standard developed by the Description Logic Implementation Group for communication with Description Logics reasoners in the realm of the semantic web and is introduced in [5]. Many reasoners [13,14,15] support this interface definition, which allows the separation of application and reasoner by the means of programming language and execution place.
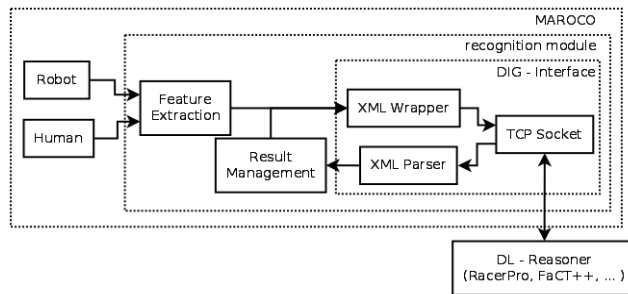


Figure 2. Components of the recognition module.

The DIG-interface follows a functional approach called *Tell&Ask* [8]. After defining a knowledge base – the *tell* operation – reasoner results and information can be retrieved – the *ask* operation. The modification of an existing knowledge base after using an *ask* operation is not defined by the DIG-interface. Therefore in each run time cycle the recognition module creates a complete knowledge base, which will be released in the end (see Fig. 3).

As a consequence the recognition module needs to manage an up-to-date model of the knowledge base, which consists of domain specific knowledge and assertions dependent on the current kinematical human model and robot specific parameters. This distinction corresponds in Description Logics with TBoxes and ABoxes even though the DIG-interface does not distinguish between them. The domain specific knowledge is modeled a priori, the assertional knowledge is created in each runtime cycle afresh. The modeled knowledge base will be explained in more detail in Section IV C.

As the assertional knowledge depends on kinematical parameters a feature extraction component is applied in order to fill the attribute values of the assertions. The following features are important w.r.t. the component *Human*: Angles of both elbows, Angles of both arms to shoulder respective to the up-axis, Angle difference between head orientation and robot, Walking velocity and used tool.

The feature *used tool* is not supported by existing sensors at the moment and is therefore simulated. It can have one of the following values: *none*, *measurement tool* or *working tool*. The simulation of this parameter can be influenced directly by user input using standard human machine interfaces. As a result complex working scenarios can be modeled and analyzed.

The component *Robot* provides the parameters for: gripper status and movement status.

During feature vector creation, extracted values are mapped onto sharp sets. The knowledge base is then populated with corresponding set strings which can be used for comparative operations during reasoning.

One major aspect of understanding human activity is modeling temporal relations between different actions. In this work, these relations are introduced by defining an *after*-role. Hence a certain action can only be recognized if certain other actions occurred prior. This *after*-role can be regarded as defining preconditions onto actions. Previously recognized actions need to be included in the knowledge base in order to allow for correct recognition of current actions. All recognized actions are stored by the reasoner result management component and are retrieved during recreation of the knowledge base.

### C. The Knowledge Base

In Figure 4, the ontology about situations and activities which are modeled by the knowledge base are presented. The concept *Situation* has the attribute *Number Humans* to distinguish between the concepts *Robot alone* and *Human present*.

In the situations of *Human present*, or its sub-concepts, *Activities* can *take place*, which are *done by* a *Human*. This defines the corresponding concepts and relating roles.
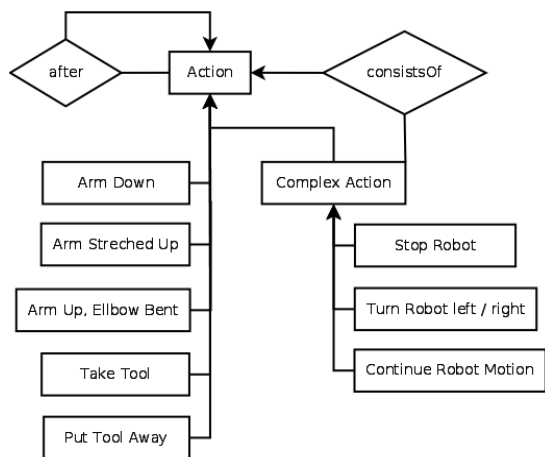


Figure 5. ER model of the action ontology.

In Figure 5, the ontology concerning *Actions* and *complex Actions* is shown. As pointed out above, actions can have a temporal relation expressed as *after*-role. The action *Put Tool Away* can only occur after the action *Take Tool*. This role is also exploited in complex actions, e.g., *Continue Robot Motion* can only be signaled after *Stop Robot*.

Actions can be regarded as atomic concepts, whereas complex actions consist of other actions, regardless of atomicity. The concepts *Take Tool* and *Put Tool Away* are considered atomic, because they are defined by and based on the single attribute *Used Tool*. This attribute is directly altered by user input, thus does not result from sensor data analysis. The role *doneBy* which is defined for activities is

also modeled for actions. For reasons of readability this relation is not depicted.

The occurrence of the situation *Cooperation* implies that there are *expectations* towards the robot behavior. Moreover, an expectation can be *triggered by* an action (see Fig. 6). This allows for reasoning about expectations without necessarily recognizing a triggering action. This implicit relation is also exploited between the activities *Monitor*, *Hold Tool* and *Actions*.

## V. EXPERIMENTAL RESULTS

For reasons of experimental analysis of the implemented activity and situation recognition different courses of action were executed and the recognition results were recorded.

In order to analyze different scenarios efficiently means of automated feature value presetting were implemented. The overall analysis is based on these presets and on actual sensor data processing. Hence natural movements and
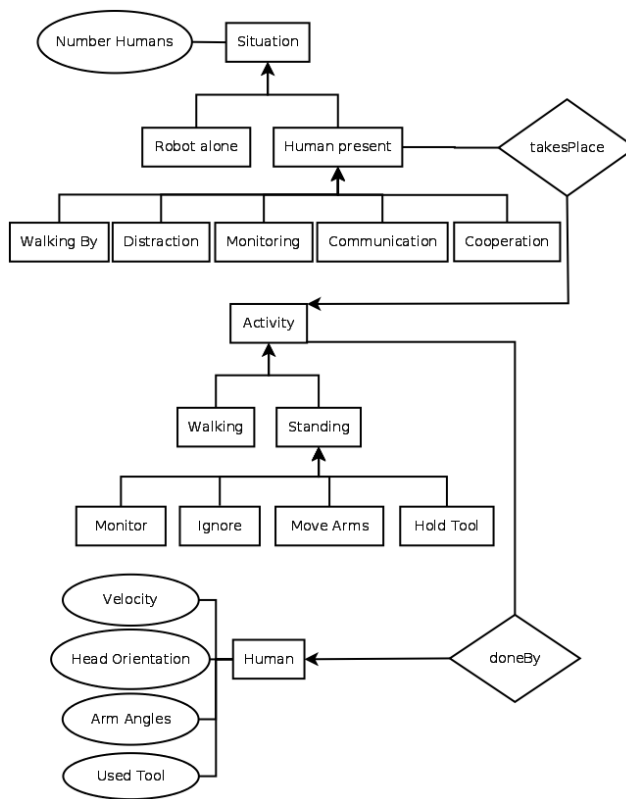


Figure 4. ER model of activity and situation ontology.

transitions between actions can be tested and special use cases can be investigated.

In this section, recorded recognition results will be illustrated and discussed.
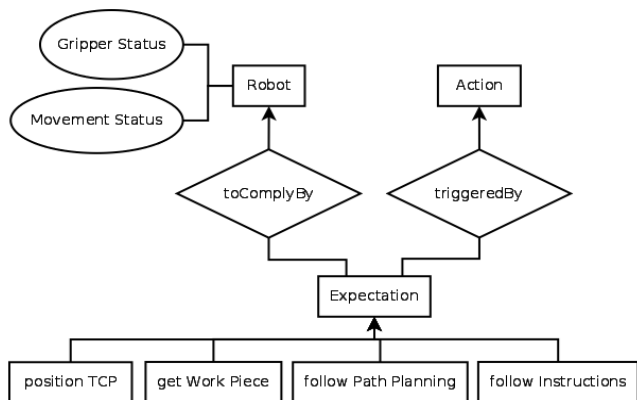
Figure 6. ER model of the expectation ontology.

## A. Examplary Result Records

The recorded experimental results contain a timestamp which indicates the starting time of the recognition cycle in milliseconds since program start. This timestamp is then followed by the extracted feature values if there is a human worker in the supervised area. The components of the feature vector are listed in following order: Angle arm left, angle arm right, angle elbow left, angle elbow right, walking velocity, angle difference between head orientation and robot, holding tool, gripper status and robot movement status.

The next number is the timestamp of the final result message from the DL reasoner (Tab. I). Results will be recorded whenever there are new insights. Thus, the last two lines of Table I have no special entries past the last return timestamp.

TABLE I.          EXAMPLE RECORD BASED ON SENSOR DATA

```
29009 29395 RobotAlone
29396 0   0 0 0 1 84 0 0 1 29797 Distraction Ignore
29799 0   0 0 0 1 86 0 0 1 30212
30213 0 15 0 8 1 56 0 0 1 30642
```

Table II demonstrates the recognition of different situations and activities. Furthermore an additional action and expectation are reasoned and recognized.

During a recognition cycle all recognized concepts are returned from the DL reasoner in a single flush, therefore, the number of lines in the records represents the number of returned responses.

TABLE II.          EXAMPLE RECORD BASED ON PRESETS

```
16160 90 0 0 0 20 0 0 0 1 16965 WalkingBy Walking
. . .
22061 90 0 0 0 20 0 0 0 1 22447
22448  0 0 0 0  0 0 1 0 1 22834 Cooperation
                         HoldTool TakeTool getWorkPiece
```

## B. Results

Tables I and II already indicate that the processing time of a recognition cycle varies around 500 ms. This indication can be shown to hold true by analysis of a large amount of cycles.

TABLE III.          RESULTS FROM EVALUATION

| # Recognition cycles | 2140 | Max [ms] | 9705 |
|---|---|---|---|
| Ø Response time [ms] | 551.78 | # > 1000 ms | 17 (0.79%) |
| Min [ms] | 216 | # > 5000 ms | 4 (0.18%) |

In Table III, the results of 2140 recognition cycles are summarized. It shows that the average processing time is approximately 550 ms. The lower bound is 216 ms. The casual outliers take up to 10 seconds in worst case scenarios. The number of cycles taking more than 1 second reaches 0.79% of all cycles. The amount of processing cycles consuming more than 5 seconds is 0.18%.

TABLE IV.          RECORD FOR ANALYSIS OF LONG RUNTIMES

```
60260 90 0 0 0 0   0 0 0 1 60740 Comm. MoveArms 480
. . .
64501 90 0 0 0 0   0 0 0 1 64940                439
64940 90 0 0 0 0   0 0 0 1 66475               1535
66475 90 0 0 0 0   0 0 0 1 67017                542
67017 90 0 0 0 0   0 0 0 1 72300               5283
72300 90 0 0 0 0   0 0 0 1 72750                450
72750 90 0 0 0 0 60 0 0 1 73221 Distr. Ignore  471
```

In Table IV, cycle run times are noted at line's end. These numbers show that long cycle times cannot be related directly to changes in the feature vector. Thus, the recognition process itself might not cause the outliers. This will need further investigation.

By using the kinematical human model, recognition of gestures and human motion can be analyzed (see Fig. 7). Table V shows an example in which a human first watches the robot. This concludes the expectation, that the robot shell follow a planned path. After some time the human moves his arms which results in a communicative situation. Because the arms are moved differently by the human, a *Stop Robot* instruction is recognized in the next recognition cycle. The reasoning results in the expectation that the robot shell comply with the instructions.

Consequently natural movements and actions can be recognized despite the average cycle processing time of approx. 550 ms.

TABLE V.          EXAMPLE RECORD FOR NATURAL MOVEMENT

```
103607 0 0 1 0 3 5 0 1 2 104135 Monitoring Monitor
                                  followPathPlanning
. . .
112169 0 0 1 0 1 9 0 0 1 112706
112707 62  9 26 21 3 0 0 0 1 113193 Comm. MoveArms
113194 74 70 21 23 2 6 0 0 1 113823 StopRobot
                                  followInstructions
113824 76 88 20 35 5 9 0 0 2 114473
```

Tables II and V demonstrate that depending on situation and actions expectations are generated. The generation of expectation is also dependent on the robot movement status. Table VI shows that at first a cooperative situation is recognized and a generated expectation *get Work Piece*. At this moment the robot was following a planned path, which is signaled as 1 in the feature vector. In the simulation incorporated in MAROCO, this generated expectation leads to a change of the robot movement status which sets the

corresponding feature value to 2, meaning the robot is obeying instructions. This change allows the reasoning to conclude the new expectation to position the robot's tool center point in order to ease the work that the user is about to do with the work piece.

TABLE VI.        EXAMPLE FOR DYNAMIC EXPECTATION REASONING

```
96795 75 0 21 0 0 3 1 0 1 97287 Coop. HoldTool
                               TakeTool getWorkPiece
97289 75 0 22 0 0 0 1 1 2 97799 positionTCP
```

This process of interaction between reasoner results and robotic behavior demonstrates the dynamic abilities of the presented approach to recognize and understand situations and actions.
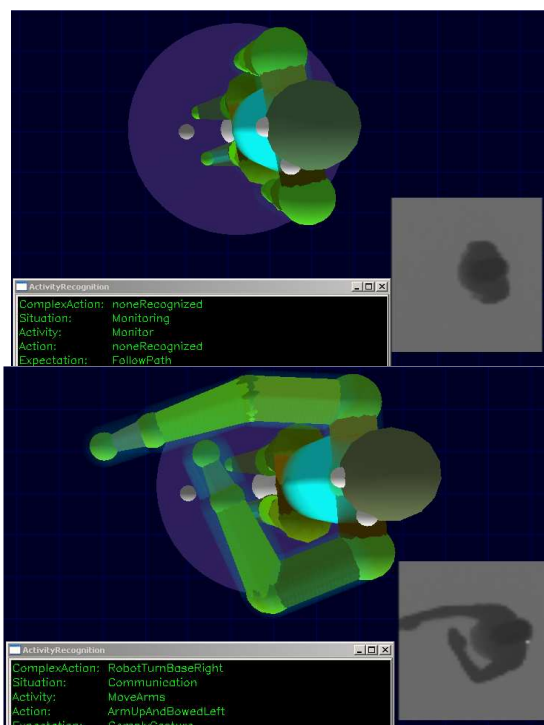


Figure 7. (Top) Human watching the robot. Recognized situation: Monitoring. Recognized activity: Monitor. No specified action recognized. The robot is expected to carry on with its task of following its planned path. (Bottom) Human is communicating with the robot. The complex action to signal a right turning movement is recognized. Recognized situation: Communication. Activity: Move arms. The robot is expected to comply with the users instructions.

## C.  Evaluation of Results

The results demonstrate that the capabilities of the presented approach reach beyond sole activity and situation recognition. By generating expectations towards robot behavior, an understanding of the situation can be achieved. This induction of relations between concepts can hardly be realized by purely probabilistic methods.

The achieved processing cycle time of approx. 550 ms does not allow for safe cooperation based only on the recognition module. Thus, the MAROCO framework uses its implemented techniques and algorithms to enforce safety and real-time capabilities during robot motion. Nevertheless, the measured results will be used to quantify improvements of later developments. To the best of our knowledge, there are no such time related results made available in the field of industrial human-robot cooperation or another related field close to it so far.

## VI.    SUMMARY AND FUTURE WORK

In this paper, a situation and action recognition module was implemented, which is capable of generating expectations towards robotic behavior.

A knowledge base containing domain and assertional knowledge was modeled. It defines concepts about situations, activities, actions and expectations. These concepts are linked and related by role definitions. Temporal associations of actions are modeled by an *after*-role, which allows preconditioning the recognition of certain actions.

Description Logics are used to define the knowledge base. By implementing the DIG-interface, Description Logics reasoning facilities can be used independently of programming language and execution space.

In order to express value constraints on concept attributes, the feature extraction process maps feature values onto sets, which can be represented as strings in the knowledge base. This allows additionally for support of a wide range of Description Logic reasoners.

During evaluation the effectiveness was shown. Situations, activities and naturally conducted actions are recognized. Expectations are generated and can influence dynamically subsequent processing cycles.

The here presented experimental results are promising for further research in the field of cognitive industrial robotics.

The next steps will be modeling a broader knowledge base in order to incorporate multi-robot setups. Also, the implementation of action plan recognition will deepen the understanding of situations and enable the analysis of complex cooperation scenarios.

It was taken a stand against the probabilistic way of estimating actions from image sequences in the beginning of the related work section. But it is suggested to evaluate different approaches in the near future which also take probabilistic methods into account or maybe fuse different methods bringing together the best of both worlds.

### REFERENCES

[1]   http://www.pmdtec.com/products-services/pmdvisionr-cameras/pmdvisionr-camcube-30/ [Last visited on 2010-08-13]

[2]   J. Graf and H. Wörn, "An Image Sequence Analysis System with Focus on Human-Robot-Cooperation using PMD-Camera", in VDI Proc. of Robotik 2008, June 2008, pp. 223-226.

[3]   J. Graf and H. Wörn, "Safe Human-Robot Interaction using 3D Sensor", in  Proc. of VDI Automation 2009, June 2009, pp. 445-456.

[4] J. Graf, S. Puls, and H. Wörn, "Incorporating Novel Path Planning Method into Cognitive Vision System for Safe Human-Robot Interaction", in Proc. of Computation World, pp. 443-447, 2009.

[5] S. Bechhofer, "The DIG Description Logic Interface: DIG/1.1.", in Proc. of the 2003 Description Logic Workshop, 2003.

[6] J. Graf, F. Dittrich, and H. Wörn, „High Performance Optical Flow Serves Bayesian Filtering for SafeHuman-Robot Cooperation", in Proc. of the Joint 41th Int. Symp. on Robotics and 6th German Conf. on Robotics, pp. 325-332, Munich, 2010.

[7] J. Graf, P. Czapiewski, and H. Wörn, "Evaluating Risk Estimation Methods and Path Planning for Safe Human-Robot Cooperation", in Proc. of the Joint 41th Int. Symp. on Robotics and 6th German Conf. on Robotics, pp. 579-585, Munich, 2010

[8] F. Baader, D. Calvanese, D. McGuinness, D. Nardi, and P. Patel-Schneider, "The Description Logic Handbook", 2nd Edition, Cambridge University Press, 2007.

[9] B. Hummel, W. Thiemann, and I. Lulcheva, "Description Logic for Vision-Based Intersection Understanding", in Proc. of Cognitive Systems with Interactive Sensors (COGIS), Stanford University, CA, 2007.

[10] B. Neumann and R. Möller, "On Scene Interpretation with Description Logics", in Image and Vision Computing, vol. 26, pp. 81-101, 2008.

[11] R. Möller and B. Neumann, "Ontology-Based Reasoning Techniques for Multimedia Interpretation and Retrieval", in Semantic Multimedia and Ontologies, part 2, pp. 55-98, Springer London, 2008.

[12] T. Springer, P. Wustmann, I. Braun, W. Dargie, and M. Berger, „A Comprehensive Approach for Situation-Awareness Based on Sensing and Reasoning about Context", in Lecture Notes in Computer Science, vol. 5061, pp. 143-157, Springer, Berlin, 2010.

[13] V. Haarslev, R. Möller, and M. Wessel, "RacerPro User's Guide and Reference Manual", Version 1.9.1, May 2007.

[14] D. Tsarkov and I. Horrocks, "FaCT++ Description Logic Reasoner: System Description", in Lecture Notes in Computer Science (LNCS), vol. 4273, pp. 654-667, 2006.

[15] E. Sirin, B. Parsia, B. Grau, A. Kalyanpur, and Y. Katz, "Pellet: A practical OWL-DL reasoner", in Web Semantics: Science, Services and Agents on the World Wide Web, vol.5 (2), pp. 51-53, 2007.

[16] V. Krüger, D. Kragic, A. Ude, and C. Geib, „The Meaning of Action: A Review on action recognition and mapping", in Proc. of Advanced Robotics, Vol. 21, pp. 1473-1501, 2007.

[17] P. Raamana, D. Grest, and V. Krueger „Human Action Recognition in Table-Top Scenarios : An HMM-Based Analysis to Optimize the Performance", in Lecture Notes in Computer Science (LNCS), Vol. 4673, pp. 101-108, 2007.

[18] Y. Wu, H. Chen, W. Tsai, S. Lee, and J. Yu, „Human action recognition based on layered-HMM", in IEEE Inter. Conf. on Multimedia and Expo (ICME), pp.1453-1456, 2008.

[19] Y. Shi, Y. Huang, D. Minnen, A. Bobick, and I. Essa, „Propagation Networks for Recognition of Partially Ordered Sequential Action", in Proc. of Computer Vision and Pattern Recognition (CVPR), Vol. 2, pp. 862-869, 2004.

[20] D. Minnen, I. Essa, and T. Starner, „Expectation Grammars: Leveraging High-Level Expectations for Activity Recognition", in Proc. of Computer Vision and Pattern Recognition (CVPR), Vol. 2, pp. 626-632, 2003.

[21] T. Gardiner, I. Horrocks, and D. Tsarkov, "Automated Benchmarking of Description Logic Reasoners", in Proc. of the 2006 Intern. Workshop on Description Logics (DL2006), Windermere, Lake Districrt, UK, 8 pages, June 2006.

# An Estimation of Distribution Algorithm using the LZW Compression Algorithm

Orawan Watchanupaporn and Worasait Suwannik

Department of Computer Science

Kasetsart University

Bangkok, Thailand

orawan.liu@gmail.com, worasait.suwannik@gmail.com

*Abstract*-**This paper proposes a new evolutionary algorithm called LZWCGA. LZWCGA is an algorithm that combines the LZW compressed chromosome encoding and compact genetic algorithm (cGA). The advantage of LZW encoding is to reduce the search space thus speed up the evolutionary search. cGA is one of Estimation of Distribution Algorithms. Its advantage is compact representation of the whole binary-string genetic algorithm population.**

*Keywords-Estimation of Distribution Algorithms; Lempel-Ziv-Welch Algorithm; Compression Algorithm; Compact Genetic Algorithm*

## I. INTRODUCTION

Genetic Algorithm (GA) is an algorithm that solves problems by simulating natural evolution [1]. To solve a problem using GA, a candidate solution must be encoded into a binary string. The length of this string represents the size of the problem. As the length of the binary string increases, the size of the search space also increases at an exponential rate. For example, the size of search space for 10-bit chromosome is $2^{10}$. While the size of search space for 100-bit chromosome is $2^{100}$.

To reduce the search space, one approach is to utilize a compressed encoding chromosome. Kunasol et. al. proposed LZWGA, which is a GA that uses LZW compressed chromosomes [2]. An LZWGA chromosome has to be decompressed by an LZW decompression algorithm before its fitness can be evaluated. LZWGA can solve very large problem such as one-million-bit OneMax, RoyalRoad and Trap functions.

Estimation of Distribution Algorithm (EDA) is a new approach in evolutionary computation [3][4]. EDA models highly-fit individuals in each generation by assuming a particular distribution. After the model is created, EDA generates new individuals from the model and inserts them to the population. Modeling and generating can avoid the disruption of partial solution resulted from genetic operations such as crossover and mutation. EDAs include Compact Genetic Algorithm (cGA) [5], Mutual Information Maximization for Input Clustering (MIMIC) [6], Bayesian Optimization Algorithm (BOA) [7], etc.

In this paper, we combine LZW compressed encoding with cGA. cGA has an advantage of a compact representation. A chromosome in cGA is a probability vector which represents the whole GA's binary string population. cGA considers all variables independently. Each item in the probability vector represents the probability that the gene

will be 0 or 1. However, because the LZW encoded chromosome is an integer array, we have to modified cGA to handle the integer value.

The remainder of this paper is organized as follows. Section II presents technical background. Section III gives details about LZWCGA. Section IV describes the experiments. Section V shows experimental results and discussion. Finally, we conclude our work and suggest future work in Section VI.

## II. TECHNICAL BACKGROUND

### A. Lempel-Ziv-Welch (LZW) Algorithm

The LZW is a lossless data compression algorithm [8]. The compression algorithm starts with a dictionary in which each entry contains one character. During the compression, the algorithm dynamically expands the dictionary and outputs codes that refer to strings in the dictionary. Normally, the number of bits of the code is less than that of the variable length string in the dictionary. Data is compressed because the algorithm replaces the whole string with its code.

A nice property of LZW is that the dictionary does not have to be packed with a compressed data. LZW decompression does not require a dictionary because the algorithm can reconstruct the dictionary while decompressing data. When using LZW to decompress an English text, the dictionary is initialized with all English characters and symbols. However, when this algorithm is used with GA, the dictionary is initialized with the number 0 and 1 because the output of the decompression algorithm must be a binary string.

A pseudo code for LZW decompression used in LZWGA is shown next page.

### B. LZWGA

The main difference between LZWGA and GA is that an LZWGA chromosome is in a compressed format. Therefore, the LZWGA chromosome has to be decompressed before its fitness can be evaluated. In [2], LZWGA is compared with the simple GA using the same parameters except the length of individuals (compressed vs no compression). For OneMax problem, by using the same amount of time, the best chromosome that simple GA can find is a little more than half of solution fitness (LZWGA can find a solution). LZWGA requires less memory and time to transfer data from one generation to the next generation. For example, to solve

one-million-bit problem, each chromosome in LZWGA have 40,000 genes or 640,000 bits (40,000 × 16) but GA used $10^6$ bits per each chromosome. LZWGA spends less time than GA during genetic operations (e.g., crossover, mutation, and reproduction). The pseudo code of LZWGA is shown below.

```
Algorithm LZWGA
        Z ← create_first_generation()
        repeat
                P ← decompress(Z)
                evaluate(P)
                Z ← create_next_generation(Z)
        until is_terminate()
```

The variable $Z$ is the population of compressed chromosome.

The variable $P$ is the population of uncompressed binary chromosomes.

The algorithm begins by creating the first generation of compressed chromosomes. Before evaluating the fitness of a chromosome, the compressed chromosome is decompressed using LZW Decompression algorithm. The fitness evaluation is performed on the uncompressed chromosome.

After that, the new population is created to replace the old population. The algorithm repeats the process of decompression, fitness evaluation, and creating a new population until the termination criterion is met. The algorithm terminates when a solution is found or a maximum generation is reached.

```
Algorithm LZW Decompress
        add entries 0 and 1 to the dictionary
        read one code from input to c
        output str(c)
        p = c
        while input are still left
                read one code from input to c
                if the code c is not in the dictionary
                        add str(p)+fc(str(p)) to the dictionary
                        output str(p)+fc(str(p))
                else
                        add str(p)+fc(str(c)) to the dictionary
                        output str(p)
                end if
                p = c
        end while
```

The variable $c$ is used to store a code read from input.

The variable $p$ is the previous value of $c$.

The function str($code$) returns a string associated with $code$.

The function fc($string$) returns the first character in $string$.

### 1) Creating the First Generation

Unlike a canonical GA, a chromosome in LZWGA is encoded as integers. The chromosome in LZWGA is in a compressed format. LZWGA chromosome is an array of integer. Each integer is a code for an index of an entry in the dictionary. Chromosomes in the first generation are created as a random integer strings with the constraint that the $i^{th}$ integer of a chromosome must not have value greater than $i+2$.

For example, an LZWGA chromosome that can be successfully decompressed is (1,2,3). The decompression algorithm will output a binary string 111111. After decompression, a dictionary has the entries (0,0), (1,1), (2,11), and (3,111). Another valid chromosome is (0,1,2). The decompression algorithm will output a binary string 0101.

If the $i^{th}$ integer in an LZWGA chromosome is invalid, the dictionary look up in will be failed after the $(i+1)^{th}$ integer is read. An example of an invalid chromosome is (1,3,3). Before entering the loop, the input "1" (the $0^{th}$ integer in the chromosome) is read and the algorithm output 1. In the first iteration, the algorithm reads "3" (the $1^{st}$ integer), adds to dictionary the string 11 at the entry 2, and outputs 11. In the second iteration, the algorithm reads "3" (the $2^{nd}$ integer), and fail when trying to execute str("3").

In order to generate the value of the $i^{th}$ integer, a random non-negative integer is modulo with $i+2$.

### 2) Decompression

Because the chromosome in LZWGA is compressed, it has to be decompressed before its fitness evaluation. A compressed chromosome is decompressed using LZW decompression algorithm. The result is a binary chromosome.

The length of the decompressed chromosome is varied. If the length is more than the size of the problem size, the excess bits are discarded. If the length is less than the problem size, LZWCGA will evaluate the fitness of available bits. After decompression, the decompressed binary string is evaluated. A fitness of a compressed chromosome is equals to the fitness of the decompressed chromosome.

### 3) Creating the Next Generation

LZWGA creates the population of the next generation by selecting, recombining, and mutating compressed chromosomes. A highly fit chromosome is likely to be selected using any selection method such as tournament or roulette-wheel selection. Compressed chromosomes can be recombined using single-point, two-point, or uniform crossover. Because each of these crossover methods does not change the position of each integer, it automatically creates valid chromosomes that each integer satisfies the constraint. Therefore, the offspring can be decompressed. Mutation changes an integer in uncompressed chromosome to a random value that satisfies the constraint.

## C. Compact Genetic Algorithm (cGA)

Harik et al. [5] introduced a compact genetic algorithm (cGA). The performance of cGA is comparable to GA with uniform crossover. cGA is a graphical representation of the probability model of EDAs without independencies. This algorithm uses a single probability vector to represent the whole GA population. Therefore cGA consumes less memory than traditional GA.

## III. LZWCGA

LZWCGA combines LZWGA with cGA. cGA uses a probability vector to represent the whole GA population. In contrast, LZWCGA uses a probability matrix instead of a single probability vector because LZWGA's chromosome is an array of integer. Each column of the probability matrix is a probability that a particular value will occurs for each gene. An example of a probability matrix is shown in Fig. 1.

The main difference between LZWCGA and cGA are initializing and updating the probability matrix process. The sequence of LZWCGA process is shown below.

Step 1. Initialize the probability matrix
Step 2. Generate two individuals
Step 3. Decompress both individuals
Step 4. Evaluate both individuals
Step 5. Update the probability matrix
Step 6. Check if the probability matrix has converged or the solution is found, if not return to Step 2

The first step in LZWCGA is to initialize the probability matrix. The pseudo code is shown below. The sum of the probability in one column of the matrix is 1.

```
Algorithm Initialize Probability Matrix
    for i = 1 to l do
        for j = 1 to i + 1 do
            p[i][j] = 1 / (i + 1)
        end for
    end for
```

The variable $l$ is length of an individual.

Then, we randomly generate two individuals $a$ and $b$ from the probability matrix using the pseudo code shown below.

```
Algorithm Generate Individuals
    for i = 1 to l do
        r = random()
        interval = 0
        for j = 1 to  i+1 do
            interval += p[i][j]
            if (r ≤  interval)
                lzwChromosome[i] = j
                break
            end if
        end for
    end for
```

Next, we decompress both individuals using LZW. Then, we evaluate their fitness. The individual with higher fitness score is called the *winner*, whereas the other is called the *loser*. The probability matrix is updated according to values from *winner* and *loser*. The main idea is to increase the probability value at the winner's position by $1/n$ (the variable $n$ is the population size) and decrease value in *loser* positions by $1/n$. The pseudo code for updating the probability matrix is shown in the right column of this page.

By way of illustration, the initial probability matrix is shown in Fig. 1. The probability matrix after updating using values from *winner* and *loser* is shown in Fig. 2.

| 0.50 | 0.33 | 0.25 | 0.20 | 0.17 |
|------|------|------|------|------|
| 0.50 | 0.33 | 0.25 | 0.20 | 0.17 |
|      | 0.33 | 0.25 | 0.20 | 0.17 |
|      |      | 0.25 | 0.20 | 0.17 |
|      |      |      | 0.20 | 0.17 |
|      |      |      |      | 0.17 |

Figure 1. The initial probability matrix of LZWCGA population when the length of each individual is 5

| winner | 0 | 1 | 2 | 3 | 5 |
|--------|---|---|---|---|---|

| loser | 1 | 0 | 1 | 0 | 2 |
|-------|---|---|---|---|---|

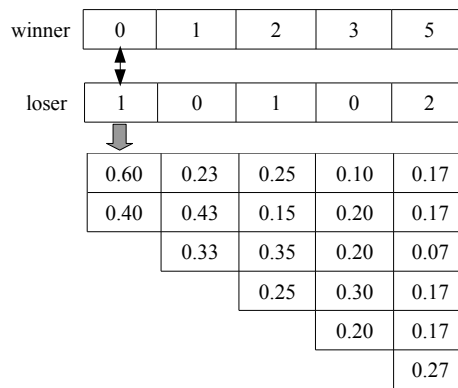| 0.60 | 0.23 | 0.25 | 0.10 | 0.17 |
|------|------|------|------|------|
| 0.40 | 0.43 | 0.15 | 0.20 | 0.17 |
|      | 0.33 | 0.35 | 0.20 | 0.07 |
|      |      | 0.25 | 0.30 | 0.17 |
|      |      |      | 0.20 | 0.17 |
|      |      |      |      | 0.27 |

Figure 2. The probability matrix after updating (population size $n$ is 10)

```
Algorithm Update Probability Matrix
    for i = 1 to l do
        indexW = winner[i]
        indexL = loser[i]
        if (indexW ≠ indexL)
            if (p[i][indexW] + (1/n) ≥ 1.0)
                p[i][indexW] = 1.0
                for j = 1 to i+1 do
                    if (j ≠ indexW)
                        p[i][j] = 0.0
                    end if
                end for
            else
                if (p[i][indexL] - (1/n) ≤ 0.0)
                    p[i][indexW] += p[i][indexL]
                    p[i][indexL] = 0.0;
                else
                    p[i][indexW] += (1/n)
                    p[i][indexL] -= (1/n)
                end if
            end if
        end if
    end for
```

The last step of LZWCGA is to check whether the probability matrix has been converged or the solution is found. If not, the evolution process is repeated starting from step 2.

## IV. EXPERIMENTS

We conducted experiments to compare the performance of LZWCGA and LZWGA on OneMax and Trap problems.

### A. OneMax Problem

The OneMax problem [9] (or bit counting) is a widely used problem for testing the performance of various genetic algorithms. Formally, this problem can be described as finding a string $\vec{x} = \{x_1, x_2, ..., x_k\}$, where $x_i \in \{0,1\}$, that maximizes the following equation:

$$F(\vec{x}) = \sum_{i=1}^{k} x_i \qquad (1)$$

### B. Trap Problem

The general $k$-bit trap functions [9] are defined as:

$$F(\vec{x}) = \begin{cases} f_{high} & ; if\ u = k \\ f_{low} - (u \times f_{low})/(k-1) & ; otherwise \end{cases} \qquad (2)$$

where $\vec{x} = \{0,1\}$, $u = \sum_{i=1}^{k} x_i$ and $f_{high} > f_{low}$. Usually,

$f_{high}$ is set at $k$ and $f_{low}$ is set at $k$-1. The Trap problem

denoted by $F_{m \times k}$ are defined as:

$$F_{m \times k}(K_1 ... K_m) = \sum_{i=1}^{m} F_k(K_i), K_i \in \{0,1\}^k \qquad (3)$$

The $m$ and $k$ are varied to produce a number of test functions. The Trap functions fool the gradient-based optimizers to favor zeros, but the optimal solution is composed of all ones. The Trap function is a fundamental unit for designing test functions that resist hill-climbing algorithms.

### C. Parameters

The parameters for both algorithms are shown in Table I and II. The Table I shows parameters for OneMax problem. Table II shows parameters for Trap problem. The size of compressed chromosome is set to 4, 5 and 6 times on OneMax and 4 times smaller than the size of a decompressed chromosome on Trap problem. We call the ratio the chromosome compression ratio. We compare the performance of LZWCGA and LZWGA for various compression ratios. LZWGA uses tournament selection (tournament size = 4). It uses uniform crossover and does not use mutation.

All experimental results are the average performance obtained from 30 runs.

TABLE I. PARAMETERS OF LZWGA AND LZWCGA FOR ONEMAX PROBLEM

| Parameter | Value |
|---|---|
| Population size | 128, 512, 1024 |
| Problem size (bits) | 1000, 10000, 100000 |
| Chromosome compression ratio | 1/4, 1/5, 1/6 of problem size |
| Max generation (for LZWGA) | 500 |
| Max round (for LZWCGA) | 500 x population size |
| Crossover rate (for LZWGA) | 1 |
| Mutation rate (for LZWGA) | 0 |

TABLE II. PARAMETERS OF LZWGA AND LZWCGA FOR TRAP PROBLEM

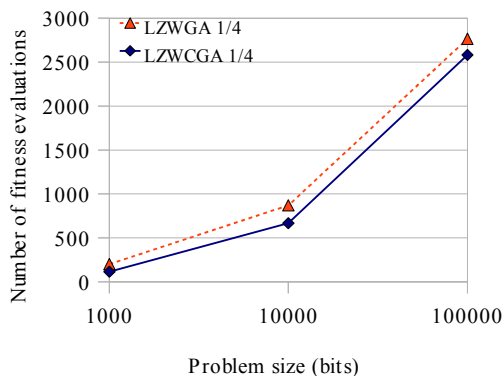| Parameter | Value |
|---|---|
| Population size | 128, 512, 1024 |
| Trap size | 5 |
| Total trap | 100, 1000, 10000 |
| Problem size (bits) | Trap size x Total trap |
| Chromosome compression ratio | 1/4 of problem size |
| Max generation (for LZWGA) | 500 |
| Max round (for LZWCGA) | 500 x population size |
| Crossover rate (for LZWGA) | 1 |
| Mutation rate (for LZWGA) | 0 |

## V. RESULTS AND DISCUSSION

The experimental results show that LZWCGA outperforms LZWGA on both OneMax and Trap problems (see Fig. 3 and Fig. 4). We found that the bigger problem size needs more fitness evaluations. Moreover, higher compression ratio requires more fitness evaluations.
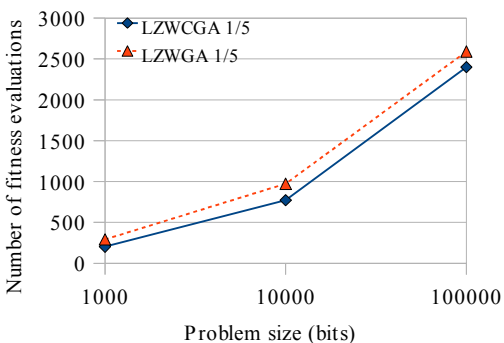
LZWGA's memory requirement depends on chromosome length and population size while LZWCGA depends only on chromosome length. For equal chromosome length, LZWGA will use approximately the same amount of memory as LZWCGA when the population size is equal to the length of the chromosome. For example, when compressed chromosome length is 1000, LZWGA with 1003 individuals uses the same amount of memory as LZWCGA. (Note that each item in an LZWGA individual is 16-bit unsigned integer and each item in an LZWCGA matrix is 32-bit float.)

A visual representation for an LZWCGA probability matrix is shown in Fig. 5. The X-axis represents positions in a chromosome and the Y-axis represents probability that a value can occur in that position. The darker area indicates higher probability. The initial probability matrix is shown in the first sub figure. Each column in the first sub figure has the same shade of gray because the initial probability that each value in each gene will occur is equal. However, during the evolution, the probability is changed. The second, third, fourth sub figure is a probability matrix at 10000, 20000, 30000 fitness evaluations and so on. In the last sub figure, the probability matrix converges. Normally, LZWGA finds a
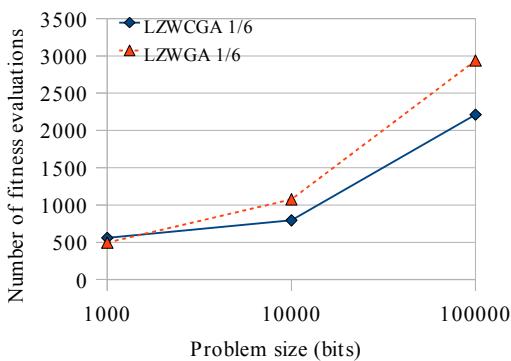
solution before the probability matrix converges. In one experiment, LZWCGA found a solution around 35000 fitness evaluations while the probability matrix converges around 45000 fitness evaluations.



(a) Population size is 128. The compression ratio is 1/4.



(b) Population size is 512. The compression ratio is 1/5.



(c) Population size is 1024. The compression ratio is 1/6.

Figure 3.    The number of fitness evaluations of LZWCGA and LZWGA when solving various sizes of OneMax problem.
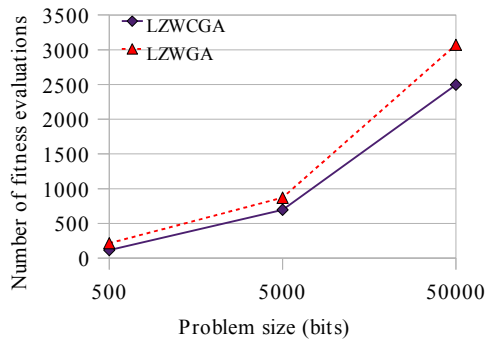


Figure 4.    The number of fitness evaluations when using LZWCGA and LZWGA to solve Trap problem. The compression ratio is 1/4.
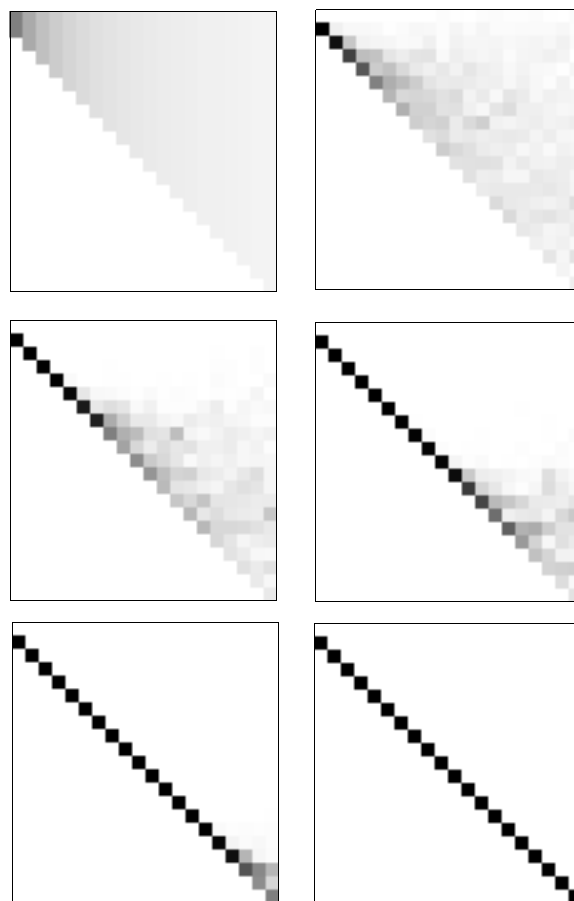


Figure 5.    A visual representation for a probability matrix at 0, 10000, 20000, 30000, 40000, and 50000 fitness evaluations

## VI. Conclusion and Future Work

We proposed the algorithm LZWCGA which combines the compress encoding and probabilistic model building. The main feature of LZWCGA is an ability to reduce the search space which makes the algorithm find the solution more effectively. We found that the LZWCGA's performance is comparable to LZWGA on OneMax and Trap problem. This result is promising because we think that if LZW encoding is integrated with more advanced EDAs, the performance of the new algorithm might be better than the original LZWGA. In the future, we will improve the update process for probability matrix and apply LZW with more advanced EDAs such as MIMIC (Mutual Information Maximization for Input Clustering), which can solve combinatorial optimization problems with bivariate dependencies.

## Acknowledgements

## References

[1] David E. Goldberg, Genetic Algorithms in Search, Optimization, and Machine Learning, Addison-Wesley, Jan. 1989.

[2] Naris Kunasol, Worasait Suwannik, and Prabhas Chongstitvatana, "Solving One-Million-Bit Problems Using LZWGA," Proc. International Symposium on Communications and Information Technologies (ISCIT), Oct. 2006, pp. 32-36.

[3] Pedro Larrañaga and Jose A. Lozano, Estimation of Distribution Algorithms A New Tool for Evolutionary Computation, Ed., Kluwer academic publishers, Boston, 2002.

[4] Topon K. Paul and Hitoshi Iba, "Linear and Combinatorial Optimizations by Estimation of Distribution Algorithms," Proc. 9th MPS Symposium on Evolutionary Computation, IPSJ, 2002.

[5] Georges R. Harik, Fernando G. Lobo, and David E. Goldberg, "The Compact Genetic Algorithm," IEEE Transaction on Evolutionary Computation, vol. 3, no. 4, Nov. 1999, pp. 287-297.

[6] Jeremy S. De Bonet, Charles L. Isbell, Jr., and, Paul Viola, "MIMIC: Finding Optima by Estimating Probability Densities," Advances in Neural Information Processing Systems, vol. 9, MIT Press, Cambridge, 1997, pp. 424-430.

[7] Martin Pelikan, David E. Goldberg, and Erick Cantù-Paz, "BOA: The Bayesian Optimization Algorithm," Proc. The Genetic and Evolutionary Computation Conference (GECCO), 1999, pp. 525-532.

[8] Terry A. Welch, "A Technique for High-Performance Data Compression," IEEE Computer, vol. 17, no. 6, Jun. 1984, pp. 8-19.

[9] Melanie Mitchell, An Introduction to Genetic Algorithms, MIT Press, 1998.

# Learning Approaches to Visual Control of Robotic Manipulators

Paulo J. S. Gonçalves*[†], Pedro M. B. Torres*
*School of Technology,
Polytechnic Institute of Castelo Branco
Av. Empresário, 6000-767 Castelo Branco, Portugal
Email: pedrotorres@ipcb.pt
[†]Technical University of Lisbon, IDMEC / IST,
Av. Rovisco Pais, 1049-001 Lisboa, Portugal
Email: paulo.goncalves@ipcb.pt

*Abstract*—This paper presents learning approaches to model the interaction between a robotic manipulator and its working environment. The approaches used are fuzzy modeling, neural networks and support vector machines. The interaction tackled in this paper is between the robot visual perception of the work environment and its actuators, while performing positioning tasks. This interaction, e.g., model, is obtaining only on measurements. This fact allows to obtain an uncalibrated model of the interaction, minimizing the setup time of the robotic system, not requiring calibrated robot kinematic and camera models. The input-output sample data used to learn the model are visual features from the work environment and the robot joint velocities, respectively. Experimental data, obtained from a IR52C robot and a visual stereo system, was used to validate the obtained models. Due to its accuracy and lower computational complexity, when compared to the other three, the off-line fuzzy model was used to control the robot, which clearly shows the effectiveness of the approach.

*Keywords*-Fuzzy Modeling, Neural Networks, Support Vector Machines, Computer Vision, Robotic Manipulators.

## I. Introduction

Visual servo control of robots is an area in continuous expansion, since vision sensors provide much more information about the working environment of the robot, than any other type of sensors. The area of robotics that addresses this concept is also called Visual Servoing, which essentially defines the control methods for dynamic systems, using information from vision sensors (cameras). The basis of this work came from the idea of modeling the interaction between the motion and vision of an industrial robot manipulator, without a-priori calibration of the system. Industrial manipulators are usually programmed using 3D coordinates of the work environment. This lead us to use stereo vision to obtain the image features, in order to obtain the 3D coordinates to control the robot.

To perform Uncalibrated Visual Servo Control, the robot-camera model must be estimated. Previous work from the authors [5], showed that Uncalibrated Visual Servo Control can be applied to control an industrial manipulator using vision, with the following advantages: no need to calibrate the robot; no need to calibrate the camera(s); the controller

has no singularities. In [5] was developed a system based on off-line fuzzy modeling.

In this paper, the estimation is performed by learning. Four approaches are presented, the first off-line fuzzy modeling [5], the second on-line fuzzy modeling [11], the third neural networks [13] and the fourth support vector machines [14]. Neural networks and support vector machines are two major machine learning approaches, and were not yet applied to the estimation of the robot-camera model. These approaches are used to compare the previous work results [5], obtained using the off-line fuzzy modeling, and also to find a better alternative to fuzzy modeling. On-line fuzzy modeling pretend to be an extension of the off-line fuzzy modeling approach. These four approaches lead to a model capable of controlling the visual servo system. The learning approaches are used to derive the inverse robot-camera model, i.e., the inverse Jacobian, in order to compute the joints and end-effector velocities in a straightforward manner. The models can be directly applied as a controller, which is a simple way to implement a controller in real-time. Note that this feature is very important in robotic systems.

This paper is organized as follows. Section II describes briefly the concept of visual servo control and presents the uncalibrated visual servo control approach. Section III presents very briefly on-line and off-line fuzzy modeling, neural networks and support vector machines. Section IV describes the experimental setup and presents the obtained results, where the identified models are discussed. Finally, Section V presents the conclusions and the future work.

## II. Visual Servo Control

### A. Calibrated Visual Servo Control

In this paper 3D visual servoing with 3D features, [1], is used in an eye-to-hand system, [2], where the camera is fixed and looking the robot and the object. The 3D image features, $s$ are 3D points of the object in the camera frame, $p$. The kinematic modeling of the transformation between the image features velocities, $\dot{s}$, and the joints velocities $\dot{q}$

is defined as follows, [1]:

$$\dot{s} = [\; -I_3 \quad S(p) \;] \cdot {}^cW_e \cdot {}^eJ_R \cdot \dot{q} = J \cdot \dot{q}, \qquad (1)$$

where $I_3$ is the $3 \times 3$ identity matrix, $S(p)$ is the skew-symmetric matrix of the 3D point $p$, ${}^cW_e$ is defined as the transformation between the camera and end-effector frames velocities, and ${}^eJ_R$ is the robot Jacobian. The 3D point is obtained from the captured image using a pose estimation algorithm, [1].

### B. Uncalibrated Visual Servo Control

To derive an accurate Jacobian, $J$, a perfect modeling of the camera, the chosen image features, the position of the camera related to the world, and the depth of the target related to the camera frame must be accurately determined. Even when a perfect model of the Jacobian is available, it can contain singularities, which hampers the application of a control law. Remind that the Jacobian must be inverted to send the camera velocity to the robot inner control loop. When the Jacobian is singular, the control cannot be correctly performed.

There are visual servo control systems that obviate the calibration step and estimate the robot-camera model either online or offline. The robot-camera model may be estimated:

- Analytically, using nonlinear least square optimization [3], and
- By learning or training, using fuzzy membership functions and neural networks [4], [5].

In addition, the control system may estimate an image Jacobian and use the known robot model, or a coupled robot-camera Jacobian may be estimated.

To overcome the difficulties regarding the Jacobian, a differential relationship between the features and camera velocities was proposed in [4]. This approach states that the joint variation depends on the image features variation and the previous position of the robot manipulator:

$$\delta q(k) = F_k^{-1}(\delta s(k+1), q(k)). \qquad (2)$$

In visual servo control, the goal is to obtain a joint velocity, $\delta q(k)$, capable of driving the robot according to a desired image feature position, $s(k+1)$, with an also desired image feature error, $\delta s(k+1)$, from any position in the joint spaces. This goal can be accomplished by modeling the inverse function $F_k^{-1}$, using fuzzy modeling as proposed in this paper and presented in Section III. This new approach to visual servo control allows to overcome the problems stated previously regarding the Jacobian and the calibration of the robot-camera model. It can be applied to all types of visual servo control. It was apllied to 2D in [5], and it will be applied to 3D in this paper.

## III. LEARNING APPROACHES TO MODELING

### A. Fuzzy Models

From the modeling techniques based on soft computing, fuzzy modeling is one of the most appealing. If no a priori knowledge is available, the rules and membership functions can be directly extracted from process measurement. Fuzzy models provide a transparent description of the system, that can reflect a possible nonlinearity of the system. The fuzzy models implemented in the presented toolbox are Takagi-Sugeno fuzzy models [6] where the consequents are crisp functions of the antecedent variables and linguistic or Mandani [7], [8] fuzzy models where both the antecedent and consequent are fuzzy propositions.

*1) Takagi Sugeno:* Takagi-Sugeno (TS) models consist of fuzzy rules describing a local input-output relation, typically in an affine form:

$$R_i : \textbf{If } x_1 \text{ is } A_{i1} \textbf{ and } \ldots \textbf{and } x_n \text{ is } A_{in}$$
$$\textbf{then } y_i = \mathbf{a}_i\mathbf{x} + b_i, \quad i = 1, 2, \ldots, K. \qquad (3)$$

Here $R_i$ is the $i$th rule, $\mathbf{x} = [x_1, \ldots, x_n]^T$ are the antecedent variables, $A_{i1}, \ldots, A_{in}$ are fuzzy sets defined in the antecedent space, and $y_i$ is the rule output variable. $K$ denotes the number of rules in the rule base, and the aggregated output of the model, $\hat{y}$, is calculated by taking the weighted average of the rule consequents:

$$\hat{y} = \frac{\sum_{i=1}^K \beta_i y_i}{\sum_{i=1}^K \beta_i}, \qquad (4)$$

where $\beta_i$ is the degree of activation of the $i$th rule: $\beta_i = \Pi_{j=1}^n \mu_{A_{ij}}(x_j)$, $i = 1, \ldots, K$, and $A_{ij}(x_j) : \mathbb{R} \rightarrow [0, 1]$ is the membership function of the fuzzy set $A_{ij}$ in the antecedent of $R_i$.

To identify the model in (3), the regression matrix X and an output vector $\mathbf{y}$ are constructed from the available data: $\mathrm{X}^T = [\mathbf{x}_1, \ldots, \mathbf{x}_N]$, $\mathbf{y}^T = [y_1, \ldots, y_N]$, where $N \gg n$ is the number of samples used for identification. The number of rules, $K$, the antecedent fuzzy sets, $A_{ij}$, and the consequent parameters, $\mathbf{a}_i, b_i$ are determined by means of fuzzy clustering in the product space of the inputs and the outputs [9]. Hence, the data set Z to be clustered is composed from X and $\mathbf{y}$: $\mathrm{Z}^T = [\mathrm{X}, \mathbf{y}]$. Given Z and an estimated number of clusters $K$, the Gustafson-Kessel fuzzy clustering algorithm [10] is applied to compute the fuzzy partition matrix U.

The fuzzy sets in the antecedent of the rules are obtained from the partition matrix U, whose $ik$th element $\mu_{ik} \in [0, 1]$ is the membership degree of the data object $\mathbf{z}_k$ in cluster $i$. One-dimensional fuzzy sets $A_{ij}$ are obtained from the multidimensional fuzzy sets defined point-wise in the $i$th row of the partition matrix by projections onto the space of the input variables $x_j$. The point-wise defined fuzzy sets $A_{ij}$ are approximated by suitable parametric functions in order to compute $\mu_{A_{ij}}(x_j)$ for any value of $x_j$.

The consequent parameters for each rule are obtained as a weighted ordinary least-square estimate. Let $\theta_i^T = \left[ \mathbf{a}_i^T ; b_i \right]$, let $X_e$ denote the matrix $[X; \mathbf{1}]$ and let $W_i$ denote a diagonal matrix in $\mathbb{R}^{N \times N}$ having the degree of activation, $\beta_i(\mathbf{x}_k)$, as its $k$th diagonal element. Assuming that the columns of $X_e$ are linearly independent and $\beta_i(\mathbf{x}_k) > 0$ for $1 \le k \le N$, the weighted least-squares solution of $\mathbf{y} = X_e \theta + \epsilon$ becomes

$$\theta_i = \left[ X_e^T W_i X_e \right]^{-1} X_e^T W_i \mathbf{y}. \qquad (5)$$

Previous work of the first author have already stated that this approach can obtain a model capable of controlling an image based visual servoing system [5].

*2) Evolving:* The model obtained from the techniques presented in the previous two sections is assumed to be fixed, since it is learned in off-line mode. Recently attention is focused in on-line learning [11], where in a first phase, input-output data is partitioned using unsupervised clustering methods and in a second phase, parameter identification is performed using a supervised learning method.

In On-Line Fuzzy Modeling and according to [11], also rule-based models of the TS type, are considered. Typically in the affine form described in (3), where the input-output data is acquired continuously. The new data, arriving at some time instant, can bring new information from the system, which could indicate a change in its dynamics. This information may change an existing rule, by changing the spread of the membership functions, or even introduce a new one. To achieve this, the algorithm must be able to judge the informative potential and the importance of the new data.

In the following is briefly presented the on-line fuzzy modeling algorithm, proposed in [11], called evolving fuzzy systems. The first step is based on the subtractive clustering algorithm [12], where the input-output data is partitioned. The procedure used must be initialized, i.e., the focal point of the first cluster is equal to the first data point and its potential is equal to one. Starting from the first data point, the potential of the next data point is calculated recursively using a Cauchy type function of first order:

$$P_k(z_k) = \frac{1}{1 + \frac{1}{k-1} \sum_{i=1}^{k-1} \sum_{j=1}^{n+1} (d_{ik}^j)^2}, \quad k = 2, 3, \dots \quad (6)$$

where $P_k(z_k)$ denotes the potential of the data point $z_k$ calculated at time $k$; $d_{ik}^j = z_i^j - z_k^j$, denotes projection of the distance between two data points ($z_i^j$ and $z_k^j$) on the axis $z^j$.

When a new data point arrives it also influences the potential of the already defined center of the $K$ clusters ($z_i^*, i = 1, 2, ..., K$). A recursive formula for the update of the cluster centers potential is defined in [11]:

$$P_k(z_i^*) = \frac{(k-1)P_{(k-1)}(z_i^*)}{k - 2 + P_{(k-1)}(z_i^*) + P_{(k-1)}(z_i^*) + \sum_{j=1}^{n+1} (d_{ik}^j)^2},$$

where $P_k(z_i^*)$ is the potential at time $k$ of the cluster center, related to the rule $i$.

The next step of the algorithm is to compare the potential of the actual data point to the potentials of the existing cluster centers.

If the potential of a new data point is higher than the potential of the existing cluster centers, then the new data point is accepted as a new cluster center and a new rule is formed. If in addition to the previous condition the new data point is close to an old cluster center, the old cluster center is replaced. The decision to create or remove rules was based on the following principles:

1)The sample has a high potential is legible to be a focal point of a fuzzy rule:

$$P_k(z_k) > \max(P_k(z_i^*)) \qquad (7)$$

2)A sample that is over an area of spatial data are is not covered by other rules, is also eligible to form a rule:

$$P_k(z_k) < \min(P_k(z_i^*)) \qquad (8)$$

3)To avoid overlap and redundancy of information in creating new rules, the following condition is also checked:

$$\exists i, i = [1, R]; \mu_{ij}(x(k)) > e^{-1}; \forall j; j = [1, n] \qquad (9)$$

R denotes the number of fuzzy rules up to the moment k. The membership function are gaussian, with the form:

$$\mu_{ij} = e^{-r\|x_j - x_{ij}^*\|^2}, \qquad (10)$$

The consequents of the fuzzy rules are obtained using the global parameter estimation procedure based on the weighted recursive least squares, presented in [11].

*B. Neural Networks*

Neural networks have found profound success in the area of pattern recognition, function approximation, optimization, pattern matching and associative memories. By repeatedly showing a neural network inputs classified into groups, the network can be trained to discern the criteria used to classify, and it can do so in a generalized manner allowing successful classification of new inputs not used during training [13]. In the present paper a feedforward backpropagation network with 5 neurons, with sigmoid activation functions in the hidden layer and a linear one in the output layer, is used to obtain the model of the interaction.

*C. Support Vector Machines*

The support vector machine (SVM) maps an input vector x into a high-dimensional feature space Z through some nonlinear mapping, chosen a priori. In this space, optimal separating hyperplanes are constructed. In the case of regression, SVM performs modeling between several clusters by finding decision hyper surfaces determined by certain points of the training set, termed Support Vectors [14].
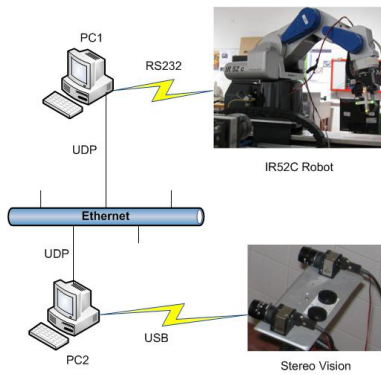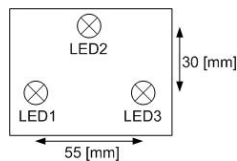
Figure 1.    The eye-to-hand experimental setup.



Figure 2.    Configuration of the markers placed on the end-effector.

## IV.  EXPERIMENTAL RESULTS

### A.  Experimental Setup

To validate the proposed approach, we use the IR52C Robotic Manipulator and a stereo vision system, with *U-Eye* cameras, in eye-to-hand configuration. The experimental setup is presented in Fig. 1. The PC2 acquire and process [15] the images from cameras and sends 3D data to the network using UDP protocol. Computer PC1 receives UDP packets coming from the PC2, and implement the visual servo control algorithm, to control the robot. To extract the 3D features for visual servo control, an object with tree LEDs, was placed in end-effector, with the marker configuration depicted in Figure 2. Computer PC2 acquires images from cameras, perform color segmentation of images, and extract the 3D coordinates, in real-time, corresponding to the 3D position of the three markers, [15]. These nine features are send to the network.

*1) Modeling Results:* For the robotic application in this paper, the models are identified using input-output data from the inputs $\dot{q}(k)$ and the outputs $\delta s(k+1)$, following the procedure described in [5]. In this paper, the approach presented in [16] to obtain the training set was used. Note that we are interested in the identification of an inverse model as in equation 2.

To obtain the data for model identification, the robot must move in the 3D workspace within the field of view of cameras, making a 3D spiral with a center point, (Fig. 3). The variables needed for identification, $\dot{q}(k)$ and $\delta s(k+1)$, are obtained from the spiral. This allows to cover a wide range of values for $\dot{q}(k)$ and $\delta s(k+1)$, by the equations 11 and 12.
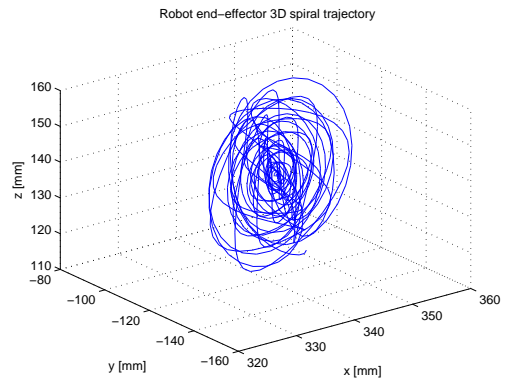


Figure 3.    3D Spiral path for model identification.

Table I
RESULTS OF THE OFF–LINE FUZZY MODEL

|         | Rules | VAF    | MSE  |
|---------|-------|--------|------|
| Joint 1 | 3     | 98,2%  | 0,23 |
| Joint 2 | 3     | 97,3%  | 0,93 |
| Joint 3 | 3     | 94,4%  | 2,81 |
| Joint 4 | 3     | 93,2%  | 1,21 |
| Joint 5 | 3     | 98,2%  | 0,23 |

$$\delta s(k+1) = s^* - s(k+1) \tag{11}$$

$$\dot{q}(k) = \frac{q^* - q(k)}{\Delta t} \tag{12}$$

To estimate the modeling accuracy, we use the VAF (Variance Accounted For), defined in equation 13, where "cov" represent the covariance vector, and "MSE" (Mean Squared Error), defined in equation 14. A perfect match occurs, when VAF is 100% and MSE have value 0.

$$VAF = 1 - \frac{cov(y_i - \widehat{y}_i)}{cov(y_i)} \times 100\% \tag{13}$$

$$MSE = \frac{1}{n} \sum (\widehat{y}_i - y_i)^2 \tag{14}$$

In Table 1, are presented the values of VAF and MSE for the off-line fuzzy modeling. With only three rules, excellent values of VAF and MSE were obtained, meaning that the model is good for estimating the joint velocities.

In Off-Line Fuzzy Modeling the number of clusters (rules) must be defined a priori in order to obtain a model. In On-Line Fuzzy Modeling, evolutionary algorithms are used after initialization, and will estimate the number of rules required in accordance with the potential associated with each data. The results from On-Line Fuzzy Modeling are presented in Table 2. The variable $\Omega$ is the initialization parameter of the algorithm, that varies with the type of data.

Table II
RESULTS OF THE ON-LINE FUZZY MODEL

|         | $\Omega$ | VAF   | MSE  | Rules |
|---------|------|-------|------|-------|
| Joint 1 | 400  | 97,9% | 0,04 | 82    |
| Joint 2 | 215  | 94,1% | 0,37 | 74    |
| Joint 3 | 250  | 95,6% | 0,36 | 75    |
| Joint 4 | 250  | 93,1% | 0,20 | 66    |
| Joint 5 | 342  | 64,3% | 1,18 | 82    |

Table III
RESULTS OF THE NEURAL METHOD MODEL

|         | VAF   | MSE  |
|---------|-------|------|
| Joint 1 | 98,3% | 0,22 |
| Joint 2 | 97,3% | 0,94 |
| Joint 3 | 94,4% | 2,84 |
| Joint 4 | 93,9% | 1,08 |
| Joint 5 | 98,2% | 0,23 |

Table IV
RESULTS OF THE SUPPORT VECTOR MACHINE MODEL

|         | VAF   | MSE  |
|---------|-------|------|
| Joint 1 | 99,7% | 0,04 |
| Joint 2 | 99,6% | 0,14 |
| Joint 3 | 99,2% | 0,39 |
| Joint 4 | 99,3% | 0,13 |
| Joint 5 | 99,6% | 0,04 |

The results presented in Table 2, show very good results with the exception of the last joint, but at the expense of a high number of rules, which will hopefully be minimized in future works. In Table 3 the results obtained with the neural network approach are presented, with similar results to the fuzzy off-line approach with only 5 neurons in the hidden layer. The SVM approach lead to better results, presented in Table 4, but at expenses of the complexity of the model since there are almost 900 support vectors. This fact does not allow a real time control of the robot, because of the computational time.

From the presented results of the four learning approaches, the on-line fuzzy model gives the best results, when taking only into account the error parameters, MSE or VAF. Since the model must be implemented to control the robot, the computational time is very important. Taking this parameter into account, the off-line fuzzy model must be used for control because it only has 3 if-then rules when compared to the 82 rules, only for joint 1, of the on-line model. The computational complexity of the neural networks and the SVM is similar to the on-line model case.

*2) Control Results:* The modeling results obtained lead us to implement the off-line approach because of the simplicity of the model (only 3 if-then rules), when compared to the neural network and the support vector machine. The on-line fuzzy model have not achieved good results for control, specially due to joint 5.

The control results were obtained using the Off-Line fuzzy model based control, defined in Fig. 4.
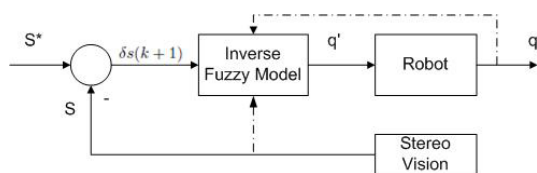


Figure 4.   Uncalibrated Visual Servo Control Loop.

To test the fuzzy models estimated, some trajectories were set within the workspace of the robot. The results obtained are quite satisfactory. Although some initial oscillations of the 3D positions error, the robot could stabilize and stop at a position close to the desired value of presenting a small error, within 3mm. Figure 5, Figure 6 and Figure 7, show the error for each marker, with respect to the 3D coordinates X, Y, Z, respectively, obtained in one of several trajectories performed with the robot. The control approach can stabilize the robot, as depicted in Figures 5 to 7, which shows the evolution of the robot joint velocities during the trajectory.

## V.  CONCLUSIONS AND FUTURE WORK

This paper presented a comparison between four learning approaches to obtain the interaction between the robot manipulator actuators and vision, when the robot performs positioning tasks. Four methods are presented and compared: Off-line Fuzzy Modeling, On-Line Fuzzy Modeling, Neural Networks and Support Vector Machines. The Off-line Fuzzy Modeling approach proven to be the adequate choice to control the robotic manipulator. With the Off-Line Model, was implemented a controller based on the learned fuzzy model, to control the IR52C robot to perform trajectories in its workspace. This controller presented very good results, with errors within 3mm of the desired position. The future goal is to implement a procedure that can update on-line the off-line learned model. For that, the first steps were already accomplished, i.e., a model based on Evolving Takagi-Sugeno Fuzzy Systems was obtained in this paper. With this approach very good results of VAF and MSE were obtained, with very satisfactory accuracy. The main objective of the future work is to reduce the number of rules of the on-line model, to allow that an On-Line Fuzzy model can control the IR52C Robotic Manipulator.

## REFERENCES

[1] E. Cervera and P. Martinet, "Combining pixel and depth information in image-based visual servoing," in *Proceedings of the Ninth International Conference on Advanced Robotics*, Tokyo, Japan, 1999, pp. 445–450.

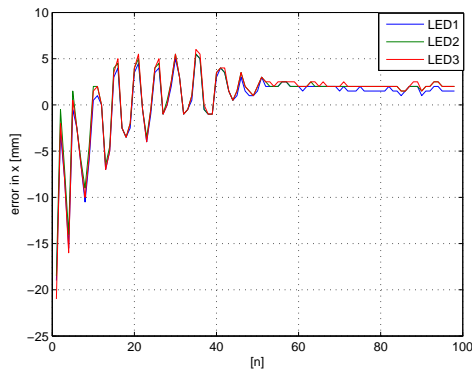ISBN: 978-1-61208-108-3                                    107

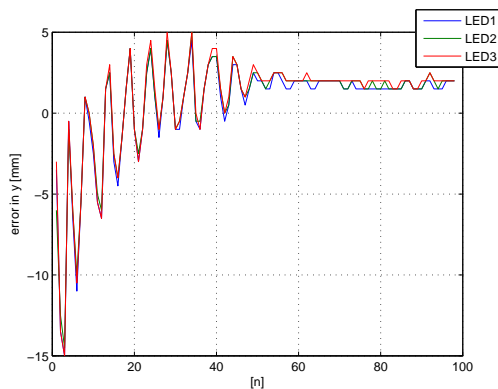Figure 5.  Evolution of the error of the position on the X coordinate.



Figure 6.  Evolution of the error of the position on the Y coordinate.
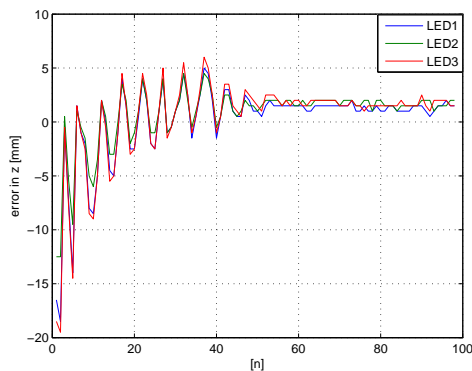


Figure 7.  Evolution of the error of the position on the Z coordinate.

[2] F. Chaumette and S. Hutchinson, "Visual servo control, part i: Basic approaches," *IEEE Robotics and Automation Magazine*, vol. 13, no. 4, pp. 82–90, December 2006.

[3] J. Peipmeier, G. McMurray, and H. Lipkin, "Uncalibrated dynamic visual servoing," *IEEE Trans. on Robotics and Automation*, vol. 20, no. 1, pp. 143–147, February 2004.

[4] I. Suh and T. Kim, "Fuzzy membership function based neural networks with applications to the visual servoing of robot manipulators," *IEEE Transactions on Fuzzy Systems*, vol. 2, no. 3, pp. 203–220, 1994.

[5] P. S. Gonçalves, L. Mendonça, J. Sousa, and J. C. Pinto, "Uncalibrated eye-to-hand visual servoing using inverse fuzzy models," *IEEE Transactions on Fuzzy Systems*, vol. 16, no. 2, pp. 341–353, 2008.

[6] T. Takagi and M. Sugeno, "Fuzzy identification of systems and its applications to modelling and control," *IEEE Transactions on Systems, Man, and Cybernetics*, vol. 15, pp. 116–132, 1985.

[7] L. Zadeh, "Outline of a new approach to the analysis of complex systems and decision processes," *IEEE Transactions on Systems, Man and Cybernetics*, vol. 3, no. 1, pp. 28–44, 1973.

[8] E. Mamdani, "Application of fuzzy logic to approximate reasoning using linguistic systems," *IEEE Transactions on Computers*, vol. 26, no. 12, pp. 1182–1191, 1977.

[9] J. Sousa and U. Kaymak, *Fuzzy Decision Making in Modeling and Control*.  Singapore: World Scientific Pub. Co., 2002.

[10] D. E. Gustafson and W. C. Kessel, "Fuzzy clustering with a fuzzy covariance matrix," in *Proceedings IEEE CDC*, San Diego, USA, 1979, pp. 761–766.

[11] P. Angelov and D. Filev, "An approach to online identification of takagi-sugeno fuzzy models," *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, vol. 34, pp. 484–498, 2004.

[12] S. L. Chiu, "Fuzzy model identification based on cluster estimation," *Journal of Intelligent Fuzzy Systems*, vol. 2, pp. 267–278, 1994.

[13] G. Zhang, "Neural networks for classification: a survey," *IEEE Transactions on Systems, Man and Cybernetics, Part C: Applications and Reviews*, vol. 30, no. 4, pp. 451–462, 2000.

[14] V. N. Vapnik, *Statistical Learning Theory*.  Wiley-Interscience, New York, 1998.

[15] P. Morgado, J. C. Pinto, J. M. M. Martins, and P. Gonçalves, "Cooperative eye-in-hand/stereo eye-to-hand visual servoing," in *Proc. of RecPad 2009 - 15th Portuguese Conference in Pattern Recognition*, Aveiro, Portugal, 2009.

[16] P. S. Gonçalves, A. Paris, C. Christo, J. Sousa, and J. C. Pinto, "Uncalibrated visual servoing in 3d workspace," *Lecture Notes in Computer Science*, vol. 4142, pp. 225–236, 2006.

# A Perspective on Machine Consciousness

Dilip K. Prasad

School of Computer Engineering
Nanyang Technological University
Singapore - 639798
e-mail: dilipprasad@pmail.ntu.edu.sg

Janusz A. Starzyk

School of Electrical Engineering and Computer Science,
Ohio University,
Athens, OH 45701, USA
e-mail: starzykj@ohio.edu

*Abstract*—**Understanding consciousness and implementing it in manmade machines has interested researchers for a long time. Despite the amount of research efforts, the measurable development in this field is very small. Among various technical, philosophical and ethical reasons, the primary reason is the difficulty in understanding consciousness. Most researchers assume that consciousness is a meta-physical phenomenon. In this paper, we consider various perspectives related to cognitive and information science to conclude that consciousness is a physical phenomenon. We also discuss the role of information in making machines conscious and present our own views on consciousness.**

*Keywords — Machine consciousness; models of machine consciousness*

## I. INTRODUCTION

Many researchers (philosophers, psychologists, cognitive neuroscientists, artificial intelligence (AI) community, etc.) have tried to define consciousness [7], [11], [15], [19], [26], [32], [35]-[38], [49], [51]-[53]. A dogmatic theory of consciousness still evades the researchers and consciousness remains an abstract term (similar to intelligence). Many considered consciousness as causal (or non-causal) [32], [34], [35], accessible (or inaccessible) [33], [34], [38], [53], stateless (or having physical state) [1], [3], [38], [47], [52], [53], representational (or non-representational) [3], [11], [13], [18], [20], [23], [24], [26], [39] and so on. Yet, none of the approaches provides a complete and thorough theory of consciousness. Based on the current research work [1]-[50], one can conclude that consciousness is not merely a physical process but a meta-physical phenomenon which builds upon and is manifest in some physical sense.

In this work, we discuss several approaches and present our view on this topic that "consciousness need not be defined in metaphysical terms". We also address the role of information and information processing mechanisms in realizing machine consciousness. We compare the existing practical models of machines consciousness with our perspective. We conclude that consciousness is a physical phenomenon, and can be realized in machines.

## II. MACHINE CONSCIOUSNESS

In order to facilitate further discussion, we introduce the concept of virtual machines and attempts of researchers to link them to consciousness [1], [2], [23], [24], [37], [41].

Virtual machines are biological/non-biological, living/non-living, or physical/non-physical systems (parts of system) that exhibit conscious phenomenon. Such machines are typically expected to:
1. Operate on complex information by acquiring, processing, and using information selectively,
2. Make decisions/selections from available options, and
3. Initiate appropriate actions.

One of the primary perspectives in the definition of machine consciousness is that consciousness is the catalyst of action (consciousness initiates action, or appropriate action as suggested by Sloman [22], [24]). We question whether satisfying the properties of virtual machines listed in (1-3) above are sufficient for consciousness.

In the framework of virtual machines, many day-to-day activities of plants (evident and undeniable examples are pitcher plants and touch-sensitive mimosa pudica) indicate that they are virtual machines too. Yet we do not refer to plants as conscious. Since if plants are conscious, then it can be argued that contemporary AI machines that exhibit some kind of information processing and control should also be referred to as being conscious which, obviously, is an overstatement.

The actions performed by low level animals (like insects) are basically sensory motor skills or results of circadian (sleep/wake) rhythm, and do not require consciousness. These animals are not conscious, though they may partially satisfy (1-3) above.

Another perspective is that consciousness needs to involve a centralized system (unlike plants) as reasoned by Muller [58]. He says that it is someone (I, he, it), a center of being, that is aware or conscious and that experiences. We do partially agree that centralized system is necessary for consciousness. More evolved and conscious animals invariably have a developed central nervous system. It is also notable that the level of consciousness increases with the development of pre-frontal cortex. In our opinion, centralized system is necessary but not sufficient for consciousness and other conditions must be present as well [57].

Different perspective is presented by Blackmore [12] and Velmans [36]. They propose that human-like consciousness is an illusion (it exists but is not what it appears to be). When we try to replicate human/animal behavior, we create this illusion that we are a conscious self having a stream of experiences [11]. We do not agree with this statement. In our

opinion, consciousness is not an illusion; it does exist and is experienced physically.

### A. Levels of machine consciousness

From the perspective of philosophers, 'machine consciousness' is a very complex issue. To make things more tractable, Sloman and Chrisley [13], [15], [22]-[24] suggest that instead of defining and characterizing consciousness, it is better to put down the expected traits of something being conscious (a fly being conscious, a new-born calf being conscious, a file protection system being conscious, etc.). Besides avoiding unnecessary abstraction and unproductive philosophical rigor, it serves another important purpose. It provides a guideline for artificially implementable aspect of consciousness and a manner of defining the requirements for a machine to be conscious.

For instance, for a fly (robot-fly) to be conscious, it needs to be able to identify an impending swat as a threat (something to avoid), the direction of threat, and a possible escape. Its consciousness need not require it to be aware of this action, or to feel happy about escaping the threat. Similarly, a new born calf needs to be aware of mother's presence, and that she can feed him, that he should suckle the mother, and that he needs to reach the mother in order to feed himself. He need not know the exact form of hunger, or what mother exactly is, or that there is a need to traverse the space between him and the mother, and so on. Another example is the requirement of consciousness in a file protection system. Such system needs to be conscious about the types of threats and the actions that it should take in the event of unauthorized access. On the other hand, it does not need to be aware of its capability of doing so. It also need not feel helpful or obliging or enslaved to do a boring job.

The point that is being made is that multiple levels of abstractions may be defined regarding the expected level and form of consciousness. We do not agree with these arguments any more than with the concept of a virtual machine satisfying (1-3) as conscious.

According to us, consciousness is an emerging phenomenon. It involves perception, learning, memory, decision making, and self-awareness. Like intelligence, consciousness may be more or less developed. For instance, social consciousness is less developed in an average 10 years old child than in an adult. The child's frontal area responsible for empathy, responsibility, etc. is still developing. So, we would like to define a scalable concept of consciousness by describing a minimum set of conditions for a conscious mind. We would rather focus on realizing some form of machine consciousness, before delving into the levels and degrees of consciousness.

### B. Problems in the implementation of machine consciousness

Many problems are encountered when implementing machine consciousness; some of them are presented next. First, the knowledge of possibilities is usually limited, and may not suffice for implementing the desired effect [1], [18], [23]. Further, virtual machines themselves are not well understood [24]. Another issue is the availability, adequacy and viability of the technology for implementation of consciousness [22], [42], [57], [59].

A more important problem is the problem of formalism and abstraction. As discussed before, consciousness is built upon, concerns, and affects, the physical phenomena occurring in and around the virtual machines. Thus, various important questions arise. Should, to what extent, and which physical phenomena affect the consciousness? How these physical phenomena should be represented? How should they affect the consciousness? In what form should the consciousness be implemented and how can it manifest itself (actions, alarms, emotions, adaptations, etc.)?

This brings us to the importance of qualia, information, architecture, ontology, and information processing. These are discussed individually in the subsequent sections.

### III. QUALIA, ONTOLOGY AND THEIR IMPORTANCE FOR MACHINE CONSCIOUSNESS

### A. Qualia

Following the difficulties in dealing with consciousness, researchers have come up with 'qualia', which helps us to speak of consciousness in a simplified (though again incomplete) manner. If we think about an experience, we do it (or feel it) in terms of certain phenomenal qualities [15], [24], [52] typically referred to as 'qualia'. Examples of qualia include perceptual experiences, bodily sensations, feelings of reactions/passions/emotions/moods, etc.

Qualia are considered central to the mind-body problem and to the development of a proper understanding of consciousness. They are usually referred to as introspection [60] or awareness, which are different from consciousness [49]. It is notable that Sloman has specifically warned against using just introspection towards the means of knowing (becoming aware of) consciousness [23].

However, if qualia have a functional nature in the form of an intermediate causal occurrence between physical inputs (like body-damage) and outputs (like withdrawal behavior), then qualia should be multiply realizable [61], [62]. In other words, physically very different states may generate the same feeling. From the point of view of functionalism, the internal modeling of physical states becomes an important part of implementing consciousness.

Such functional approach shall involve the capability to derive a connection between the possible sensory inputs, the intentions, and the possible outcomes. Therefore, it must deal with physical information and internal states (which may be in the form of direct or indirect manifestations of previous knowledge, experience, intentions, etc.) to give rise to a physical phenomenon of consciousness.

### B. Ontology

Ontology is defined as the characterization of conceptualization. Philosophers and artificial intelligence scientists differ in the definition of ontology. Since our paper pertains to machine consciousness, we use the definition used in the AI community [23], [25], [38]. According to them, ontology is the study of various categories of abstract

entities, events, processes, matters, etc., and their inter-relationships.

In this sense, ontology is an integral part of machine consciousness, as it formalizes the relationship between physically existing, measurable, decipherable, or deducible information to the abstract mental phenomena that the AI scientists wish to implement. Also, in this sense, it is directly related to the functional nature of qualia.

It is evident that if the chosen ontology is unable to characterize the desired phenomenon, implementation of desired phenomenon may result in failure. Technically, the effect of choosing an unsuitable ontology is called ontological blindness [15], [23], [24]. It refers to the failure of the ontology in visualizing certain abstract characteristics that are pertinent to the understanding and realization of the desired phenomenon.

Information, as used in the machine consciousness, is typically different from Shannon's information and simply means 'content of relevance to something'. The properties of information that are relevant and important for machine consciousness are listed below [15], [23]:

- Information can be false
- Items of information can stand in relations like consequence, contradiction and relevance
- Items of information can be understood or misunderstood.
- Information content is sometimes completely predictable
- Information is non-physical (albeit physically realized), thus, it requires specialized methods for identifying, explaining, and processing.

As agreed by most researchers, the ontology (used to study and link information to consciousness) and the information processing architecture (with consideration of the depth and breadth of information analysis, processing, interpretation, and consequent awareness, learning, planning, action initiation, etc.) are very important for realizing machine consciousness. As mentioned above, if information is picked wrongly, misrepresented, misinterpreted, or wrongly processed, it is expected to at least warp or distort consciousness, and in harsher conditions may lead to failure in the realization of desired phenomenon.

Thus, it is very important to suitably understand, classify, distinguish, and group information. It is also important that information is cast into suitable architectures and processing methodologies so that the pertinent aspects are not neglected.

Sloman [22]-[24] says that for implementation of machine consciousness, it is useful to deal with information in terms of its useful characteristics. Examples are various types of information, the forms it can take, the means of acquiring it, manipulating it, storing it, and communicating it, the purposes for which it can be used, the various ways of using it, etc. It is also important to choose suitable forms of representations, suitable algorithms, and identification of subsystems that process independently and concurrently. Their variations lead to the variation in the resulting consciousness achieved by these choices. For instance, such choices may determine the role of consciousness, in identifying possible events, actions, internal states, internal processes, and causal chain of happenings.

It should be noted that a machine gathers the information from the external signals generated in the sensory units and the internal signals like pains/rewards and machine's motivations. Thus, these physical units of information processing are necessary for consciousness.

## IV. REACTIVE, DELIBERATIVE AND REFLECTIVE MECHANISMS OF INFORMATION PROCESSING

Another impact of the choice of information and information processing is the possibility of reactive [23], [33], [37], [47], [50], deliberative [7], [23], [50], or reflective mechanisms [2], [4], [23], [28], [50]. Of course, the converse also holds, which means, that the type of mechanism desired has an impact on the choice of structural and processing aspects of information.

Reactive mechanism means that the virtual machine always lives in present (never in past or future) and simply reacts or responds to the present stimuli. It is able to make sense of the presently acquired information and choose an action from presently possible options. It never learns from past, nor can it predict future possibilities, and thus it cannot plan, predict, foresee and adapt. It should be noted that if such machines are equipped with internal states they may be capable of influencing the future and may be adaptive.

Deliberative mechanism, on the other hand, is equipped with an advanced ontology that can represent, store, process and relate to the possibilities (of input information, intermediary and internal states, as well as the possible actions). Thus, in some sense (depending upon the richness, depth and breadth, and design of ontology), it has the ability to learn, plan and govern its actions while being aware of the process and having active involvement in it. Though the importance of reactive mechanisms should not be understated, the deliberative mechanism forms the first link between contemporary AI and the desired machine consciousness.

Reflective mechanism enables the virtual machine to reflect on its own actions, to self-monitor, self-examine, and self control. Self-control in reflective mechanism is a result of being aware of one-self and one's own actions. Thus, it is reflective mechanism that shall enable the machine to possess qualia, introspect itself, possess emotions, and relate to itself and its own existence. It also enables a machine to interpolate its self-knowledge in order to understand other machines. Thus, reflective mechanism is very helpful in realizing machine consciousness.

## V. MODELS OF MACHINE CONSCIOUSNESS

We summarize the state of art research in this area by presenting some of the practical achievements in machine consciousness. We mention useful frameworks/architectures and some actually implemented systems for machine consciousness and their models.

We begin with the physical definition of consciousness, its working mechanism and a computational model proposed by us [57]. We define machine consciousness as follows:

"A machine is conscious if besides the required components for perception, action, and associative memory,

it has a central executive that controls all the processes (conscious or subconscious) of the machine; the central executive is driven by the machine's motivation and goal selection, attention switching, learning mechanism, etc. and uses cognitive perception and cognitive understanding of motivations, thoughts, or plans. Thus, central executive, by relating cognitive experience to internal motivations and plans, creates self-awareness and conscious state of mind."

A simplified version of our proposed model of the conscious machine similar to the one presented in [57] is shown in Fig. 1.
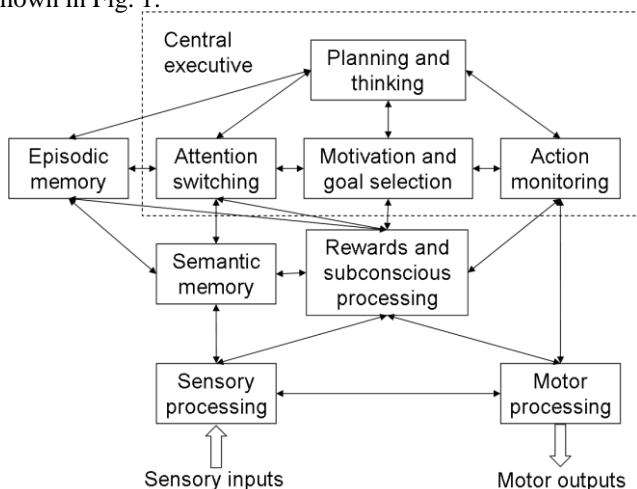


Fig. 1 Proposed computational model of consciousness.

Our current efforts are to build such defined conscious machines, and improve their mental ability and level of intelligence.

Baars and Franklin proposed a model of consciousness called the global workspace theory and developed an agent called Intelligent Distribution Agent (IDA) based on this model [6]-[8]. IDA is an agent that was designed for the U.S. navy for the purpose of collecting information from personnel, assessing the personnel on the basis of their performance and the issues they have had as human beings, and help in new task-allocation and problem identification. The global workspace theory [54] says "Consciousness is accomplished by a distributed society of specialists that is equipped with working memory, called a global workspace, whose contents can be broadcast to the system as a whole" (an argument refuted by Susan Blackmore [11]). Further Baars says [55], "Global workspace theory suggests a fleeting memory capacity in which only one consistent content can be dominant at any given moment". The content of the memory is decided by the consciousness. The exact role of consciousness is to decide the dominant content of the memory. In our model, conscious mechanism that uses planning, thinking, goal creation and goal selection based on machine's motivation, plays a critical role in selecting the dominant content of memory [57].

Rosenthal's Higher Order Thought (HOT) theory [56] says "We don't sense our conscious thoughts and sensations, since there's no distinctive sensory modality or sense organ for doing so. The only alternative is that we are conscious of our conscious thoughts, feelings, and sensations by having thoughts about them. These higher-order thoughts are themselves seldom conscious; so we are typically unaware of them". According to us, there can be only one thought at any moment in the conscious state of a conscious machine. According to us [57], "Conscious machine central executive directs cognitive aspects of machine experiences but its operation is influenced by competing signals representing motivations, desires, and attention switching that are not necessarily cognitive or consciously realized. Central executive does not have any clearly identified decision making center. Instead, its decisions are the result of competition between signals that represent motivations, pains and desires. At any moment, competition between these signals can be interrupted by attention switching signal. Such signals constantly vary in intensity as a result of internal stimuli (e.g., hunger) or externally presented and observed opportunities. Thus, the fundamental mechanism that directs machine in its action is physically distributed as competing signals are generated in various parts of machine's mind. Further, it is not fully cognitive, since, before a winner is selected, machine does not interpret the meaning of competing signals" [57]. Hence, in our opinion, Rosenthal's HOT theory is based on fallacy that multiple thoughts exist in conscious machine simultaneously.

Haikonen has proposed a cognitive architecture based on his theory that if the neural network is large enough and complicated enough, the traits of consciousness will eventually emerge on their own [16], [17], [19], [50]. Thus, he suggests that no algorithm specifically designed for implementing consciousness is necessary. According to us, his claim is vague and too generic. It does not give any concrete direction towards the realization of machine consciousness. According to us, mere presence of large and complicated neural network is not enough to realize machine consciousness. Apart from sufficiently large neural network and suitable cognitive architecture we do require other units such as motivation unit, attention switching unit, action monitoring unit, goal creation system unit, planning and thinking unit [57]. Although these units are fully distributed and are part of large interconnected network, they are functionally different in terms of their role in selecting a cognitive input, governing the planning process, and monitoring the control sequence for motor action.

Sun proposed another architecture called CLARION [27, 29]. CLARION is a two-level architecture in which the lower level concerns things, events, perceptions, reactive information, which cannot be associated with consciousness, and the higher level solely implements consciousness. These two levels are inter-connected with each other in both directions to form complete virtual machine architecture. A similar theory is the supramodular interaction theory proposed by Morsella [48], which models consciousness as integration of high level, specialized, multi-modal systems. Under this model, a distinction between the consciously penetrable and impenetrable modules (systems) is highlighted, and consciousness appears as a cross-talk among these modules. No mechanism is implied how these

cross-talks between different modules or levels are organized and how a machine reflects its own actions.

Taylor proposed corollary discharge of attention movement (CODAM) model, that emphasizes the role of attention and change in attention for implementing consciousness [31]. We do agree that attention and attention switching mechanism is important to realize machine consciousness. However, other factors also need to be considered [57].

Sloman, Chrisley, and their team have proposed the cognitive and affective (CogAff) schema [22]-[24], which is a generalized schema applicable to a wide range of AI and consciousness architectures. It is able to accommodate multiple hierarchical levels as well as independent lateral modules in its architecture. It is capable of representing reactive, deliberative, as well as reflective (meta-management) systems, thus providing a broad framework for comparison, judgment, and possibilities exploration of AI and consciousness architectures.

Some AI robots or systems have also been made to explore, understand and implement machine consciousness, which include CRONOS (a human like musculoskeletal robot) with the aim of phenomenal consciousness [47], Cyber Child [46] (a test bed for machine consciousness), Khepera models/robots [40], etc. A good overview is provided in [50].

## VI. CONCLUSION

Based on our review of the research done on machine consciousness, we can make the following observations.

Like intelligence, consciousness is a property of a physical mind not a meta-physical phenomenon. It can be alive or not but it requires awareness and intelligence as illustrated in Fig. 2.
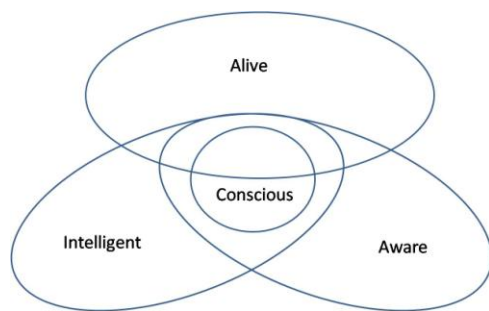


Fig. 2 Consciousness centered view of intelligent systems

Given its complicated nature, a unified explanation of consciousness is difficult. Yet, it is much more difficult to realize it in machines. Thus, in our opinion, though we have to use discretely defined, crisp, and formal theories and architectures, we also have to incorporate the gray area, the possibility for the system to evolve by itself and find its own consciousness. Our work in [57] is a step towards this aim.

Large part of the work in machine consciousness is inspired by biological systems. Our idea is also based on it.

The idea is actually similar to the appearance of consciousness in human beings that we can understand the best. From the conception of a fetus, through its development as a human being, till its death, a human being is taught a lot of things, including moral behavior, importance of social existence, relationship with nature, science, mathematics, sports, art, and so on. But no human being is explicitly taught to be conscious. There might be some intelligence/thinking enhancing exercises. However, these are enhancements only, while the basic mechanism is provided and consciousness develops itself. Furthermore, there is no basic consciousness training, and in its fundamental meaning consciousness is not a thought phenomenon. Higher levels of consciousness (like social awareness) may be trained, but a mechanism for consciousness is a property of the mind and cannot be taught.

In short, a large part of what we think as consciousness emerges through social interactions if a proper architecture, mechanism and functional blocks required for consciousness are available [57]. We conclude that consciousness is a physical phenomenon which can be realized in machines.

## REFERENCES

[1] M. L. Anderson, "Embodied Cognition: A field guide," *Artificial Intelligence*, vol. 149, pp. 91-130, 2003.

[2] M. L. Anderson and T. Oates, "A review of recent research in metareasoning and metalearning," *AI Magazine*, vol. 28, pp. 7-16, 2007.

[3] M. L. Anderson, "Circuit sharing and the implementation of intelligent systems," *Connection Science*, vol. 20, pp. 239-251, 2008.

[4] M. L. Anderson and D. Perlis, "What puts the "meta" in metacognition?," *Behavioral and Brain Sciences*, vol. 32, pp. 138-139, 2009.

[5] J. Newman, B. J. Baars, and S. B. Cho, "A neural global workspace model for conscious attention," *Neural Networks*, vol. 10, pp. 1195-1206, 1997.

[6] S. Franklin, A. Kelemen, and L. McCauley, "IDA: A cognitive agent architecture," in *Proceedings of the IEEE International Conference on Systems, Man and Cybernetics*, San Diego, CA, USA, 1998, pp. 2646-2651.

[7] B. J. Baars and S. Franklin, "How conscious experience and working memory interact," *Trends in Cognitive Sciences*, vol. 7, pp. 166-172, 2003.

[8] B. J. Baars and S. Franklin, "An architectural model of conscious and unconscious brain functions: Global Workspace Theory and IDA," *Neural Networks*, vol. 20, pp. 955-961, 2007.

[9] S. J. Blackmore, "Meme machines and consciousness," *Journal of Intelligent Systems*, vol. 9, pp. 355-376, 1999.

[10] S. Blackmore, "Evolution and memes: The human brain as a selective imitation device," *Cybernetics and Systems*, vol. 32, pp. 225-255, 2001.

[11] S. Blackmore, "There is no stream of consciousness," *Journal of Consciousness Studies*, vol. 9, pp. 17-28, 2002.

[12] S. Blackmore, "Consciousness in meme machines," *Journal of Consciousness Studies*, vol. 10, 2003.

[13] R. Chrisley, "Embodied artificial intelligence," *Artificial Intelligence*, vol. 149, pp. 131-150, 2003.

[14] R. Clowes, S. Torrance, and R. Chrisley, "Machine consciousness - Embodiment and imagination," *Journal of Consciousness Studies*, vol. 14, pp. 7-14, 2007.

[15] R. Chrisley, "Philosophical foundations of artificial consciousness," *Artificial Intelligence in Medicine*, vol. 44, pp. 119-137, 2008.

[16] P. O. A. Haikonen, "Towards associative non-algorithmic neural networks," in *Proceeding of the IEEE International Conference on Neural Networks*, Orlando, FL, USA, 1994, pp. 746-750.

[17] P. O. A. Haikonen, "Modular neural system for machine cognition," in *Proceedings of the International Joint Conference on Neural Networks*, Como, Italy, 2000, pp. 47-50.

[18] P. O. A. Haikonen, "Essential issues of conscious machines," *Journal of Consciousness Studies*, vol. 14, pp. 72-84, 2007.

[19] P. O. A. Haikonen, "Reflections of consciousness: The mirror test," in *AAAI Fall Symposium - Technical Report*, Arlington, VA, 2007, pp. 67-71.

[20] E. T. Rolls, "Consciousness in neural networks?," *Neural Networks*, vol. 10, pp. 1227-1240, 1997.

[21] E. T. Rolls, "A computational neuroscience approach to consciousness," *Neural Networks*, vol. 20, pp. 962-982, 2007.

[22] A. Sloman, "What kind of cognitive architecture does an emotional agent need?," *International Journal of Psychology*, vol. 31, pp. 4001-4001, 1996.

[23] A. Sloman and L. Chrisley, "More things than are dreamt of in your biology: Information-processing in biologically inspired robots," in *International Workshop on Biologically Inspired Robotics*, Bristol, England, 2002, pp. 145-174.

[24] A. Sloman and R. Chrisley, "Virtual machines and consciousness," *Journal of Consciousness Studies*, vol. 10, pp. 133-172, 2003.

[25] J. Chappell and A. Sloman, "Natural and artificial meta-configured altricial information-processing systems," *International Journal of Unconventional Computing*, vol. 3, pp. 211-239, 2007.

[26] R. Sun, "Learning, action and consciousness: A hybrid approach toward modelling consciousness," *Neural Networks*, vol. 10, pp. 1317-1331, Oct 1997.

[27] R. Sun, "Cognitive Architectures and Multi-agent Social Simulation," *in 8th Pacific Rim International Workshop on Mult-Agents (PRIMA 2005)*, Kuala Lumpur, MALAYSIA, 2005, pp. 7-21.

[28] R. Sun, X. Zhang, and R. Mathews, "Modeling meta-cognition in a cognitive architecture," *Cognitive Systems Research*, vol. 7, pp. 327-338, 2006.

[29] R. Sun, "The importance of cognitive architectures: an analysis based on CLARION," *Journal of Experimental & Theoretical Artificial Intelligence*, vol. 19, pp. 159-193, 2007.

[30] J. G. Taylor, "Neural networks for consciousness," *Neural Networks*, vol. 10, pp. 1207-1225, 1997.

[31] J. G. Taylor, "CODAM: A neural network model of consciousness," *Neural Networks*, vol. 20, pp. 983-992, 2007.

[32] M. Velmans, "Is human information-processing conscious," *Behavioral and Brain Sciences*, vol. 14, pp. 651-668, 1991.

[33] M. Velmans, "A reflexive science of consciousness," *Ciba Foundation Symposium*, vol. 174, pp. 81-91, 1993.

[34] M. Velmans, "When perception becomes conscious," *British Journal of Psychology*, vol. 90, pp. 543-566, 1999.

[35] M. Velmans, "Making sense of causal interactions between consciousness and brain," *Journal of Consciousness Studies*, vol. 9, pp. 69-95, 2002.

[36] M. Velmans, "Why conscious free will both is and isn't an illusion," *Behavioral and Brain Sciences*, vol. 27, p. 677, 2004.

[37] M. Velmans, "Where experiences are: Dualist, physicalist, enactive and reflexive accounts of phenomenal consciousness," *Phenomenology and the Cognitive Sciences*, vol. 6, pp. 547-563, 2007.

[38] M. Velmans, "How to define consciousness: And how not to define consciousness," *Journal of Consciousness Studies*, vol. 16, pp. 139-156, 2009.

[39] C. Browne, R. Evans, N. Sales, and I. Aleksander, "Consciousness and neural cognizers: A review of some recent approaches," *Neural Networks*, vol. 10, pp. 1303-1316, 1997.

[40] T. Kitamura and K. Ono, "An architecture of emotion-based behavior selection for mobile robots," in *5th International Conference on Simulation of Adaptive Behavior*, Zurich, Switzerland, 1998, pp. 671-690.

[41] S. Densmore and D. Dennett, "The virtues of virtual machines (Paul Churchland)," *Philosophy and Phenomenological Research*, vol. 59, pp. 747-761, 1999.

[42] R. H. Schlagel, "Why not artificial consciousness or thought?," *Minds and Machines*, vol. 9, pp. 3-28, 1999.

[43] J. McCarthy, "Free will - Even for robots," *Journal of Experimental and Theoretical Artificial Intelligence*, vol. 12, pp. 341-352, 2000.

[44] G. Buttazzo, "Artificial consciousness: Utopia or real possibility?," *Computer*, vol. 34, pp. 24-30, 2001.

[45] F. Crick and C. Koch, "A framework for consciousness," *Nature Neuroscience*, vol. 6, pp. 119-126, 2003.

[46] R. M. J. Cotterill, "CyberChild - A simulation test-bed for consciousness studies," *Journal of Consciousness Studies*, vol. 10, pp. 31-45, 2003.

[47] O. Holland, "The future of embodied artificial intelligence: Machine consciousness?," in *Lecture Notes in Artificial Intelligence* (Subseries of Lecture Notes in Computer Science), Dagstuhl Castle, 2004, pp. 37-53.

[48] E. Morsella, "The function of phenomenal states: Supramodular interaction theory," *Psychological Review*, vol. 112, pp. 1000-1021, 2005.

[49] Z. L. Torey, "The immaculate misconception," *Journal of Consciousness Studies*, vol. 13, pp. 105-110, 2006.

[50] D. Gamez, "Progress in machine consciousness," *Consciousness and Cognition*, vol. 17, pp. 887-910, 2008.

[51] D. Dennett, "The milk of human intentionality," *Behavioral and Brain Sciences*, vol. 3, pp. 428-430, 1980.

[52] D. Dennett, "Recent work in philosophy of interest to AI," *Artificial Intelligence*, vol. 19, pp. 3-5, 1982.

[53] D. Dennett, "Are we explaining consciousness yet?," *Cognition*, vol. 79, pp. 221-237, 2001.

[54] B. J. Baars, *A cognitive theory of consciousness*, Cambridge University Press, 1998.

[55] B. J. Baars, "The conscious access hypothesis: Origins and recent evidence," in *Trends Cogn. Science*, Vol 6, pp. 47–52, 2002.

[56] D. M. Rosenthal, The nature of Mind, Oxford University Press, 1991.

[57] J. A. Starzyk and D. K. Prasad, " Machine Consciousness: A Computational Model," *Brain-inspired Cognitive Systems (BICS 2010)*, Madrid, Spain, 14-16 Jul 2010.

[58] V. C. Muller, "Is there a future for AI without representation?," *Minds and Machines*, vol. 17, pp. 101-115, 2007.

[59] R. Sun and C. X. Ling, "Computational cognitive modeling, the source of power, and other related issues," in *Workshop on Computational Modeling - The Source of Power, at the 13th National Conference on Artificial Intelligence (AAAI-96)*, Portland, Oregon, 1996, pp. 113-120.

[60] P. M. Churchland, "Reduction, Qualia, and the Direct Introspection of Brain States," *The Journal of Philosophy*, vol. 82, pp. 8-28, 1985.

[61] V. S. Ramachandran and W. Hirstein, "Three laws of qualia: What neurology tells us about the biological functions of consciousness," *Journal of Consciousness Studies*, vol. 4, pp. 429-457, 1997.

[62] W. G. Lycan, "Consciousness explained - Dennett, D. C.," *Philosophical Review*, vol. 102, pp. 424-429, 1993.