# DETECTING HUMAN BEHAVIOR EMOTIONAL CUES IN NATURAL INTERACTION

*George Caridakis, Stylianos Asteriadis, Kostas Karpouzis and Stefanos Kollias*

Intelligent Systems, Content and Interaction Lab
National Technical University of Athens
Iroon Polytexneiou 9, 15780 Zografou, Greece
gcari, stiast, kkarpou, stefanos@image.ntua.gr

## ABSTRACT

Current work focuses on the detection of human behavior emotional cues and their incorporation into affect aware Natural Interaction. Techniques for extracting emotional cues based on visual non verbal human behavior are presented. Namely, gesture qualitative expressivity features and head pose and eye gaze estimation are derived from hand and facial movement respectively. Extracted emotional cues are employed in expressive synthesis on virtual agents, based on the analysis of actions performed by human users, in a Human-Virtual Agent Interaction setting and in Assistive Technologies aiming to infer in real time the degree of attention or frustration of children with reading difficulties.

*Index Terms—* Affective computing, Natural Interaction, Gesture expressivity, Eye Gaze

## 1. INTRODUCTION

Affective computing has been a topic of great interest during the last few years. Research has been performed in various disciplines associated with interaction, including perception, interpretation, cognition and expression. International conferences have been organised on this topic, including ACII series LREC workshops, recently IEEE Transactions on Affective Computing (TAC) has published its first issues while two new books have been published [1] and [2]. IST projects and networks have been funded at European level for investigating different issues of affective interaction, such as theories and models of emotional processes, computational modeling, emotional database, input signal analysis, emotion recognition, generation of embodied conversational agents; projects like Interface, Ermis, Safira, Humaine, Semaine.

Various results have been obtained, in Europe and worldwide (US, Asia) by different projects, researchers and industry, regarding affective interaction. These mostly refer to the derivation and analysis of affective and emotional theories and related computational models, the extraction of affective cues from single or multi-sensorial inputs  mainly aural and visual  the modeling of affective states, the analysis and recognition of user states based on extracted cues, the generation of synthetic characters that communicate different expressive states and attitudes, the generation of databases with affective interactions for training and testing the analysis and synthesis techniques, the inclusion of the above in interactive environments. The aforementioned activities have produced a variety of systems that model and analyze single or multimodal affective cues; they have extracted and used statistical information and rules for this purpose; they have created data sets and environments which have been used next to perform user state detection, Embodied Conversational Agent (ECA) synthesis and interaction.

## 2. RELATED WORK

For exhaustive surveys of existing work in machine analysis of affective expressions, readers are referred to [3] and [4]. Recent advancements and research directions in Affective Computing are also discussed in [1] and [2].

As has been proven by an abundance of experimental studies, incorporating multiple modalities into affective analysis systems enhances their performance and robustness; Audiovisual fusion can make use of complementary information incorporated into these channels. Such reliability improvement comes with the cost of introducing additional challenges related to the multimodal aspect of affective analysis and synthesis. Multimodal fusion techniques, synchronization issues and absence or unreliability of information channels are challenges that are encountered most frequently. Information loss during processing and feature extraction is fairly common in naturalistic recordings, either due to technical implications or due to uncontrolled user behavior. Fusing multiple modalities alleviates such problems by combining multiple flows of information at a feature level (early fusion) or at decision level (late fusion). Combining early and late fusion hybrid or ensemble techniques have been proposed recently [5]. The architecture of affective analysis systems should cater for input from multiple modalities, which vary in several aspects.

Much research work has been carried out on automatic detection of basic, acted and extreme emotions recorded in

posed, controlled interaction environments, reporting highly accurate results [3]. Affective Computing has progressed significantly since the primary attempts of automatic emotion recognition. The area has now matured, so as to validate the hypotheses - initially tested in controlled interaction - onto natural interaction. To present affective analysis of interaction taking place in less constrained settings is still a very challenging problem due to the fact that deliberate behavior differs in terms of affective cues from spontaneously occurring behavior. The shift of focus of the research in the field to the automatic emotional analysis of spontaneously displayed, natural affectively enhanced behavior is essential.

Human Computer Interaction continuously introduces new means of communication and interaction with systems [6] [7]. Gesture based HCI and alternative Natural Interaction is increasingly attracting the attention of researchers in related research areas. Popular methods in gesture recognition encounter problems such as arbitrary or experimental parameters initialization, high computational cost consisting approaches unsuitable for real time applications, user dependency and scaling issues for applicability to large scale gesture lexicons. Additionally, human motion qualitative aspects associated with expressivity have been extensively explored by modeling of human body motion and animation of virtual characters able to convey emotional content. But computational modeling of expressivity parameters has not been investigated adequately from the perspective of automatic analysis. On the other hand, from a synthesis point of view, Embodied Conversational Agents and Virtual Character procedural and parametric animation are being incorporated into systems within several domains (E-Learning, Virtual Museums and Serious Games).

Understanding intentions of others is an important ability that involves representing the mental states of others in ones own mind; in the case of an affect-aware system, a machine learning infrastructure. The principles and techniques that humans deploy in order to understand, predict, and manipulate the behavior of other humans is collectively referred to as a theory of mind (ToM - [8]). In a Human-Computer Interaction framework, this would translate to tracking the behavior of users, based on how they behave when interacting with the system, identifying certain characteristics (e.g. staring at a fixed point for some seconds or looking away from the screen) and adapting the means of interaction to cater for the detected user state. In literature, various works exist for approaching the problem of attention estimation, varying in terms of applications, cues and hardware set-up they use. A lot of work has been done for estimating the degree of concentration in driving conditions [9] and [10]. For example, in [10], the authors use stereoscopic techniques to estimate head and eye directionality, in order to simulate attentiveness in driving conditions, while, in [9], the authors utilize the eye blink o model driver's behavior. Similar, the problem of Human-Virtual agent interaction is under intense research [11] with a

lot of works focused on the issue of creating believable ECA's behaviors, able to enrich the conversation, using as input both human and agent interaction non-verbal features.

## 3. MULTIMODAL AFFECTIVE ANALYSIS

Current article focuses on techniques used for extracting emotional cues from non verbal human behavior aiming for application in Natural Interaction. Video input is adopted in all three techniques for processing hand movement in order to extract qualitative expressivity features (Section 3.1) and classify gestures (Section 3.2) while facial regions are being processed for an head pose and eye gaze estimation (Section 3.3).

### 3.1. Gesture expressivity analysis

Features and cues of non verbal behavior are an integral part of the communication process since they provide information on the current emotional state and the personality of the interlocutor [12]. Common classification schemes include binary categories such as slow/fast, restricted/wide, weak/strong, etc. Our gesture expressivity modeling is close to these schemes in the sense that they provide formulation and quantitative measurement of the respective aspects of the gesture. Adopting a subset of the gesture synthesis expressivity modeling parameters [13] we define five expressivity features: Overall activation, Spatial extent, Temporal, Fluidity and Power.

Starting from a computational formulation of these parameters described in [14] and in an attempt to provide a more strict definition let us consider a gesture $G$ as a sequence, of $T$ frames, consisting of coordinates of the left and right hand, $(x_{li}^G, y_{li}^G)$ and $(x_{ri}^G, y_{ri}^G)$, respectively and $i \in [1, T]$. The coordinates of hands are relative to the position of the head which is defined as the center of the bounding box of the region of head as provided by the head detection module and normalized with reference to the diagonal of this box which is considered indicative of the size of the head. These transformations are required in order to ensure that the coordinates are invariant to the position and the distance of the user in front of the camera, parameters that are not known a priori. Thus a gesture is formally defined as:

$$
G = [((x_{l1}^G, y_{l1}^G), (x_{r1}^G, y_{r1}^G)), ((x_{l2}^G, y_{l2}^G), \\ (x_{r2}^G, y_{r2}^G)), \cdots, ((x_{lT}^G, y_{lT}^G), (x_{rT}^G, y_{rT}^G))] \tag{1}
$$

Overall activation is considered as the quantity of movement during a dialogic discourse and is formally defined as the sum instantaneous quantities of motion: $OA_G = \sum_{i=1}^{T-1} D_{li}^G + D_{ri}^G$, $D_i = \left| \overrightarrow{(x_i, y_i)(x_{i+1}, y_{i+1})} \right|$. Spatial extent is expressed with the expansion or the condensation of the used space in front of the user (gesturing space). In order to provide a strict definition of this expressivity feature spatial extent is considered as the maximum

value of the instantaneous spatial extent during a gesture: $SE_G = \max e_i, i \in [1,T], e_i = \left| \overrightarrow{(x_{ri}, y_{ri})(x_{li}, y_{li})} \right|$. The temporal expressivity parameter is defined as the as the arithmetic mean of this quantity and since $OA_G$, as defined earlier corresponds to the discrete integral: $TE_G = \frac{OA_G}{T}$. On the other hand, the energy expressivity parameter refers to the movement of the hands at during the stroke phase of the gesture. The formalization of the energy expressivity feature according to this definition however is far from trivial since the automatic detection of the gesture phases is a quite challenging task. Alternatively we opted to associate this parameter qualitatively with the first derivative of the norm of $D$ which refers to the acceleration of hands during a gesture: $PO_G = |D|'$. Fluidity differentiates smooth / elegant from the sudden / abrupt gestures. This concept attempts to denote the continuity between hand movements and is suitable for modeling modifications in the acceleration of the upper limbs. Under this prism, we formally define as the gesture's fluidity the variance of the energy expressivity parameter as described in the previous paragraph: $FL_G = var(PO_G)$. The reader is prompted to note that quantity $FL_G$ corresponds to an expressive quality that is reversely proportional to the notion of fluidity.

## 3.2. Gesture recognition

Affective cues of non verbal behavior are often context dependent. This viewpoint for gesture expressivity establishes the need for a gesture recognition module essential in order to provide adaptable affective analysis. Self-organizing maps and Markov chains are incorporated into the adopted classification scheme. Extracted features describing hand trajectory, region and shape are used as input to separate classifiers, forming a robust and adaptive architecture whose main contribution is the optimal utilization of the neighboring characteristic of the SOM during the decoding stage of the Markov chain, representing the gesture class. Although an abundance of architectures have been proposed [15] for automatic gesture recognition achieving true user independence and user performance invariant recognition still remains a challenging issue. Current work adopts a novel approach based on a combination of Self Organizing Maps (SOM) and Markov models for hand trajectory classification [16]. The extracted features used in the trajectory module include the trajectory of the hand and the direction of motion in the various stages of the gesture. This classification scheme is based on the transformation of a gesture representation from a series of coordinates and movements to a symbolic form and on building probabilistic models using these transformed representations. The abstraction of the symbolic form enriches the classification scheme with adaptability, due to the incorporation of the neighboring function of the SOM during the decoding phase. Our study indicates that, although each of the two sets of features (trajectory and hand shape information) can provide distinctive information in most cases, only an appropriate combination can result in robust and confident user independent gesture recognition. First transformation is based on position and results in a set of Markov models for each gesture class. Equivalently another transformation is based on hand movement direction during gesture duration resulting in a set of Markov models, corresponding to the motion, rather than position, aspect of the gesture. An additional model that is created per sign class is the Generalized Median of the gesture set. In general, a generalized median of a set of sequences is defined as the sequence, that consists of a combination of all or some of the symbols used in the set that minimizes the sum of distances to every string of the set.

During decoding a class is assigned to a gesture instance based on evaluation of the models derived from the three transformations described earlier. In an abstraction level, the overall decoding stage operates as an energy maximization algorithm, constantly seeking at every step to maximize the next transition at the possible cost of choosing a different, from the one originally mapped, but more probable node for every model. The cost is located at the neighboring relation between the two nodes, since by choosing a different node the evaluation of a specific instance/model pair is penalized by this term in order to compensate for the node replacement. The algorithm intuitively strives to converge to the model's most probable path, penalizing each deviation from this path with the respective nodes' neighboring relation.

## 3.3. Gaze estimation

Estimating gaze directionality is of high importance for recognizing behavioral states related to attention estimation. Although not directly related to affect recognition, eye gaze and head posture can be of equal importance with facial expressions for transmitting non-verbal signals related to emotion recognition, alertness, attention estimation, context-dependent human-machine interaction [10], [17] , [18]. Very important is the role of gaze estimation in conditions of multi-party conversations where, people's attention is an important indicate of their role, their disposition or the degree of their engagement to the event (shared attention). A system's applicability is, most of the times, dependent on its simplicity, portability, user independence, as well as its applicability in unknown environments. The scheme has been built taking into account the aforementioned constraints. It is completely user and scale independent, and has been built in order to be used in uncontrolled environments with only demand in terms of specialized hardware a off-the-shelf web-camera. Furthermore, it combines eye gaze directionality and head rotation estimation, using non-intrusive mechanisms. The combination of the above cues has not been thoroughly studied in bibliography, with a few exceptions [19], [20].

Head Pose is estimated based on the position of the eyes midpoint with regards to its position when the user is fac-

ing the camera frontally, also discriminating between frontal and rotated views of the head. The system restarts based on expected geometrical relations among facial points and natural human motion criteria. The eyes midpoint distortions are normalized with the inter-ocular distance, as calculated every time the user is facing the camera frontally. In this way, the system is scale independent and can give reliable results for various distances of the user with regards to the camera, whenever re-initialization occurs. For tracking facial features coordinates, a three-pyramid Lucas-Kanade tracker is used, in order to handle large and sudden point movements. The above system can perform real time and the rules employed for re-initialization render it robust to sudden changes in lighting and spontaneous movements. For eye gaze detection, the area defined by the four points around the eye is used. Prototype eye areas depicting right, left, upper and lower gaze directionality are used to calculate mean grayscale images corresponding to each gaze direction. The areas defined by the four detected points around the eyes, are then correlated to these images. The differences of the correlation values of the eye area with the left and right, as well as upper and lower mean gaze images are calculated and used for estimating eye gaze directionality. Further details of the system, from face and facial feature detection to head pose and eye gaze estimation, are given in [17]

## 4. AFFECT AWARE NATURAL INTERACTION

The above mentioned techniques (Section 3) are employed in Natural Interaction scenarios such as Human-Virtual Agent Interaction and Assistive technologies.

### 4.1. Gesture analysis and Human-Virtual Agent Interaction

Gestural analysis as described in section 3 has been applied to Embodied Conversational Agents by multimodal and expressive synthesis on virtual agents, based on the analysis of actions performed by human users. As input we consider the image sequence of the recorded human behavior and an appropriate combination of image processing techniques lead to robust head and hand detection and tracking.

Based on the gesture expressivity features extracted as described in 3.1 the same gestures are performed by the agent in a qualitatively different way depending on this set of parameters. The ECA animation is generated either by defining some key frames or by computing the interpolation curves passing through these frames. This figure illustrates clearly the effect of the Spatial Extent parameter, since the gesturing space is significantly broadened in the second instance and condensed in the third. The effect of other expressivity parameters is not readily illustratable in still images and the reader is reffered to video samples located at `http://www.image.ece.ntua.gr/~gcari/videos/`, `seq1_all.avi`

and `seq2_all.avi`. The animation is specified by a sequence of keyframes as defined by MPEG–4 facial and body animation parameters (FAP and BAP) and the interpolation between these keyframes. Interpolation is performed using TCB (Tencion, Continuity, Bias) splines. The expressivity parameters are implemented by changing the TCB parameters of the interpolating splines and by scaling the values and changing the timing or scaling of the keyframes points. For example, the $SE$ parameter will influence the value of the keyframes by scaling them. The higher $SE$ will be, the wider the interpolating curves will be and so facial expressions will be more visible on the face and gestures wider. The $FL$ parameter will modulate the Continuity parameters of the splines, making them becoming smoother (high $FL$) or jerkier (low $FL$).

Remaining in the framework of virtual agent natural interaction user focus of attention estimation has been studied in [21], gaze estimation of humans is correlated with non-verbal behavior of an Embodied Conversational Agent (ECA) for inferring those conditions where shared attention is maximised. A virtual agent tries to sell products to the human observer, posing different behavior styles. Following a pyramidal scheme for inferring user attention (starting from simple screen coordinates up to context-aware engagement recognition), human participant's focus of attention is modelled and a series of conclusions are drawn regarding the most appropriate behavior the virtual seller is supposed to employ. A typical example of the application can be seen in figure 1, where, in the center, visual feedback of user's attention is given, while, the rest of the screen, shows the virtual seller and shop.



**Fig. 1**: HVAI: Human-Virtual Agent Interaction

### 4.2. Gaze estimation and Assistive technologies

The statistical significance of eye gaze, head pose, head rotation speed and facial feature coordinates has been examined [17] and a series of fuzzy rules have been learnt for inferring the degree of attention of a child, towards an electronic material. To this aim, Takagi-Sugeno-Kang [22] fuzzy inference engines have been trained using videos of Greek and Danish children with learning difficulties (between 8 and 10 years old), annotated by experts, in terms of their behavioral

state. The engines used head pose, eye gaze and inter-ocular distance changes as inputs to estimate the degree of attentiveness of a person towards electronic material they had in front of them. Eye gaze, head rotation and distance variations of a person from the camera, during a learning procedure have been mapped to clusters corresponding to attention or distraction. Furthermore, head pose, eye gaze and head speed were shown to be statistically important features for training fuzzy inference engines for estimating frustration. A clustering technique [23] was employed for defining the number of fuzzy rules for each network (levels of attention and levels of frustration), as a pre-processing step. The above, instead of using a grid partition of the data, guarantees that fuzzy rules will only be created where a large concatenation of data exists. Subsequently, the network's parameters (membership function centers and widths, and inference weights) are acquired following the hybrid model described in [22]. Acco rding to the behavioral state attributed to each child, at a particular instance, word highlighting changes, adapting the presentation to the child's capabilities and needs. A total of 10,000 and 12,250 frames were used for attention and frustration experiments, respectively and, for training the engines, a leave-one-out cross-validation protocol was followed for each child, training the networks with instances of the rest of the children. More details regarding the methodology can be found in [17]. Inferring attention and frustration based on non-verbal, facial cues, has been applied on the Agent Dysl FP6 project [24] system, a software created for assisting children with learning difficulties, instances of which are shown in Figure 2.



**Fig. 2**: Inferring attention based on non-verbal, facial cues

## 5. CONCLUSIONS AND ONGOING RESEARCH

We have presented the incorporation of emotional cues in expressive synthesis on virtual agents, based on the analysis of actions performed by human users, in a Human-Virtual Agent Interaction setting and in Assistive Technologies aiming to infer in real time the degree of attention or frustration of children with reading difficulties.The detected human behavior emotional cues include gesture qualitative expressivity features and head pose and eye gaze estimation derived from hand and facial movement respectively, and based on visual non verbal human behavior.

Ongoing research work includes further investigation of Natural Interaction settings such as Human Robot Interaction (figure 3) and Serious Games. Enhancing robotic systems with the ability to sense and convey emotions through non-verbal cues such as body postures, gestures and facial expressions is one of the most promising direction in HRI. Finally, analyzing player behavior in Serious Gaming would allow for adoption of deeper emotional models that include cognitive concepts that are not trivial to detect and process in other contexts.
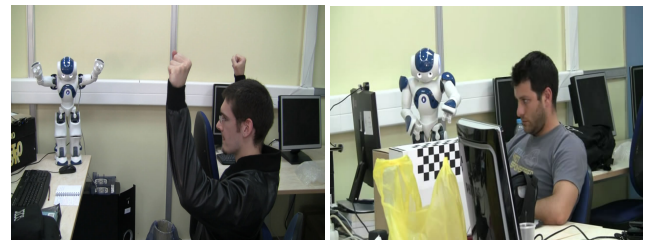


**Fig. 3**: Human Robot Interaction

## 6. REFERENCES

[1] P. Petta, C. Pelachaud, and R. Cowie, Eds., *Emotion-Oriented Systems, The Humaine Handbook*, Springer, Series: Cognitive Technologies, February, 2011.

[2] Klaus R. Scherer, Tanja Banziger, and Etienne Roesch, Eds., *A Blueprint for Affective Computing, A sourcebook and manual*, Oxford University Press, November, 2010.

[3] Z. Zeng, M. Pantic, G.I. Roisman, and T.S. Huang, "A survey of affect recognition methods: Audio, visual, and spontaneous expressions," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 31, no. 1, pp. 39–58, 2008.

[4] Z. Zeng, M. Pantic, G.I. Roisman, and T.S. Huang, "A survey of affect recognition methods: Audio, visual, and spontaneous expressions," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 1, pp. 39–58, 2009.

[5] J. Kim and F. Lingenfelser, "Ensemble approaches to parametric decision fusion for bimodal emotion recognition," in *Int. Conf. on Bio-inspired Systems and Signal Processing (Biosignals 2010)*, J. Filipe A. Fred and H. Gamboa, Eds. BIOSTEC, 2010, pp. 460–463, INSTICC Press: Portugal.

[6] G. Castellano, G. Caridakis, A. Camurri, K. Karpouzis, G. Volpe, and S. Kollias, *A Blueprint for Affective Computing*, chapter Body gesture and facial expression analysis for automatic affect recognition, Oxford University Press, 2010.

[7] A. Braffort, R. Gherbi, S. Gibet, J. Richardson, and D. Teil, *Gesture-based communication in human-computer interaction*, Springer, 1999.

[8] R. Bar-On and S. Cohen, *Mindblindness: As essay on autism and theory of mind*, Wiley Online Library, 1995.

[9] T. D' Orazio, M. Leo, C. Guaragnella, and A. Distante, "A visual approach for driver inattention detection," *Pattern Recognition*, vol. 40, no. 8, pp. 2341–2355, 2007.

[10] A. Doshi and M.M. Trivedi, "Head and Gaze Dynamics in Visual Attention and Context Learning," in *Proceedings of the CVPR Workshop on Visual and Contextual Learning*, 2009.

[11] B. Brandherm, H. Prendinger, and M. Ishizuka, "Interest estimation based on dynamic bayesian networks for visual attentive presentation agents," in *Proceedings of the 9th ACM International Conference on Multimodal Interfaces*, 2007, pp. 346–349.

[12] A. Mehrabian, *Nonverbal communication*, Aldine, 2007.

[13] B. Hartmann, M. Mancini, S. Buisine, and C. Pelachaud, "Design and evaluation of expressive gesture synthesis for embodied conversational agents," in *Proceedings of the fourth international joint conference on Autonomous agents and multiagent systems*. ACM, 2005, pp. 1095–1096.

[14] G. Caridakis, A. Raouzaiou, E. Bevacqua, M. Mancini, K. Karpouzis, L. Malatesta, and C. Pelachaud, "Virtual agent multimodal mimicry of humans," *Language Resources and Evaluation*, vol. 41, no. 3, pp. 367–388, 2007.

[15] S.C.W. Ong and S. Ranganath, "Automatic sign language analysis: A survey and the future beyond lexical meaning," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 27, no. 6, pp. 873–891, 2005.

[16] G. Caridakis, K. Karpouzis, A. Drosopoulos, and S. Kollias, "Somm: Self organizing markov map for gesture recognition," *Pattern Recognition Letters*, vol. 31, pp. 52–59, 2010.

[17] S. Asteriadis, *User attention and interest modeling and recognition in non Intrusive Interaction Enviroments*, Ph.D. thesis, School of Electrical and Computer Engineering, National Technical University of Athens, February, 2011.

[18] G. Castellano, I. Leite, A. Pereira, C. Martinho, A. Paiva, and PW McOwan, "It's all in the game: Towards an affect sensitive and context aware game companion," in *Affective Computing and Intelligent Interaction and Workshops, 2009. ACII 2009. 3rd International Conference on*. IEEE, 2009, pp. 1–8.

[19] Roberto Valenti, Adel Lablack, Nicu Sebe, Chabane Djeraba, and Theo Gevers, "Visual gaze estimation by joint head and eye information," *Pattern Recognition, International Conference on*, vol. 0, pp. 3870–3873, 2010.

[20] U. Weidenbacher, G. Layher, P. Bayerl, and H. Neumann, "Detection of head pose and gaze direction for human-computer interaction," in *Perception and Interactive Technologies*, 2006, pp. 9–19.

[21] C. Peters, S. Asteriadis, and K. Karpouzis, "Investigating shared attention with a virtual agent using a gaze-based interface," *Journal on Multimodal User Interfaces, Springer, 2009, DOI 10.1007/s12193-009-0029-1*, 2009.

[22] J.-S. Roger Jang, "ANFIS: Adaptive-Network-Based Fuzzy Inference System," *IEEE Transactions on Systems, Man, and Cybernetics*, vol. 23, pp. 665–684, 1993.

[23] S. L. Chiu, "Fuzzy Model Identification Based on Cluster Estimation," *Journal of Intelligent and Fuzzy Systems*, vol. 2, no. 3, 1994.

[24] AgentDysl, *The AgentDysl Project*, http://www.agent-dysl.eu/.