

# Towards a Method For Evaluating Naturalness in Conversational Dialog Systems

Victor Hung, Miguel Elvir, Avelino Gonzalez, Ronald DeMara

Intelligent Systems Laboratory  
University of Central Florida  
Orlando, Florida

victor@isl.ucf.edu, miguel.elvir@isl.ucf.edu, gonzalez@ucf.edu, demara@mail.ucf.edu

**Abstract**— The evaluation of conversational dialog systems has remained a controversial topic, as it is challenging to quantitatively assess how well a conversation agent performs, or how much better one is compared to another. Furthermore, one of the hurdles which remains elusive in this quandary is the definition of naturalness, as demonstrated by how well a dialog system can maintain a natural conversation flow devoid of perceived awkwardness. As a step towards defining the dimensions of effectiveness and naturalness in a dialog system, this paper identifies existing evaluation practices which are then expanded to develop a more suitable assessment vehicle. This method is then applied to the LifeLike virtual avatar project.

**Keywords**— dialog systems, artificial intelligence, human-computer interaction, software evaluation

## I. INTRODUCTION

Chatbots, or interactive conversation agents, present a special challenge with respect to validation and verification. Specifically, evaluating these programs cannot rely on a process that solely consists of quantitative methods, since there remains a great deal of subjectivity involved in assessing their performance. Hence, the evaluation of chatbots remains a controversial topic, as there is no general method for judging how well a conversation agent performs, in both the relative and the absolute sense. In exploring this subject, a pivotal focus of this paper will be defining *naturalness*, as in how well a chatbot can maintain a natural conversation flow. This paper presents a survey of existing chatbot evaluation methods, as well as a definition for naturalness in relation to Human-Computer Interaction (HCI) applications.

This addresses the needs of a National Science Foundation (NSF) supported endeavor – the LifeLike virtual avatar project [26]. This research involves re-creating an existing NSF figurehead into a digital life form. The resulting avatar will ideally serve as an information-dispensing replacement of the original human it is modeled after, characterizing LifeLike as an assistive conversational agent.

A main aspect of the system is its chatbot-like user interface. Specifically, LifeLike incorporates a conversational dialog system in its user interface whose prime directive is to provide expert decision support to its users. This must be done while maintaining a sense of naturalness in its conversation-based human-to-computer exchanges. Preliminary efforts in

evaluating its dialog system have included both qualitative and quantitative measures. The objective of this paper is to investigate a proper method of chatbot evaluation for the purpose of validating the performance of LifeLike.

With the LifeLike virtual avatar project as a backdrop, we present the following findings in our journey toward providing a suitable validation and verification method for our chatbot research. The remainder of this paper discusses the background technologies involved in chatbot evaluation, followed by a basic framework of the prototypical assessment system to be utilized by LifeLike.

## II. BACKGROUND

To empirically evaluate the naturalness of a dialog system's interaction with human users, we must first revisit the conceptual basis underlying such applications. The proceeding section considers the background issues concerning chatbot technology. This is accomplished by exploring the conclusions drawn by researchers whose applications reveal typical phenomena of naturalness and interaction.

### A. Early Intelligent Systems

Early intelligent applications acted on declarative knowledge to process data. In these production systems, development of the learning framework relied heavily on explicitly defined rules, with the purpose of assimilating new knowledge and exerting conflict resolution schemes. These models operate on and maintain highly domain constrained knowledge bases whereby the user or client becomes the major recipient of the system's conclusion or hypothesis. Thus, early production systems inherently provided immutable information retrieval processes, or fixed contexts, with limited capacity to assess the validity of system output and to modify its actions accordingly.

These intelligent agents simulated human performance of simple tasks by creating "goal-oriented and data-determined behavior." They relied on information processing and problem solving paradigms [11]. Declarative knowledge, in the form of production rules, governed the information retrieval and context selection phases [13]. Within this infrastructure, context can be selected, matched against known scenarios, or traversed in predetermined directions by the agent's use of a set of fixed rules.

These production systems found their way into some early beneficiaries, namely, chatbot technology, HCI experiences whose machines strived to mimic text-based human responses. One of the most successful and recent of these, ALICEbot [12], could maintain short realistic conversations before revealing its computerized identity. Future implementations of ALICE-based bots offered positive results from domain limited dialogs, suggesting improved performance by incorporating user-initiated system corrections as well as providing several thousand response rules [12]. Nevertheless, while these applications appeared to maintain a feeling of realism and a coherent shallow common sense, they lacked the ability to exploit symbolism in human understanding.

### B. *Naturalness in Dialog Systems*

As the state of Automatic Speech Recognition (ASR) technology improves, there remains an increasing need to provide sophisticated response systems that both convey more natural dialogs and actively acquire new information from said dialogs [14]. Gurevych demonstrates, through an evaluation of the semantic coherence of ontology-based ASR systems [14], that a gap exists in recognition that is effectively and semantically coherent, partially due to the arbitrary nature of human speech and understanding of context [15].

Currently, spoken interaction may not be as efficient at accomplishing tasks as text-based interaction. Regardless, Le Bigot et al [15] suggest that spoken interaction promotes collaboration rather than placing emphasis on the task itself and its performance, without regards to the dialog quality. They argue that this may be a result of both the lower informational density of speech and the elimination of essential terms for grounding 'shared knowledge' that occurs in speech-based HCI. For example, consider the simple case of a computer prompting the user for the date on which he or she is arriving at a conference. The user's response may be as succinct as 'The twelfth,' indicating a vague temporal sense. Here, it is necessary for intelligent speech applications to assert confirmation of any declarative knowledge it acquires throughout the interactive session in a manner consistent with these constraints.

Context retrieval experimentation, in the form of meta-cognitive application development, revealed that the acquisition of novel skills by an application can be facilitated by monitoring the state of knowledge, rationalizing goals and implementing an adept instructional structure [16]. Intelligent tutoring systems [16] based on the Adaptive Control of Thought (ACT) theories of knowledge suggest that some level of meta-cognition could improve a system's performance level. Extending these findings to ASR knowledge frameworks implies that the internal structure of the knowledge corpus and the system's awareness of its state and quality will directly impact the effectiveness of HCI. In particular, it is important that, in extracting the relevant segments of a conversation, the agent discovers whether the new knowledge enriches the context of the interaction or whether it is detrimental to it.

A Knowledge Acquisition agent depends on the quality of the information received to identify the conversational domain [17]. An obvious impediment to obtaining contextually relevant data arises from imperfect transcripts from speech recognition.

Semantic checks on the retrieved audio will hinder the system's interpretation of facts and its ability to validate context [14] [15]. A simple experiment would show the effectiveness of a chatbot's coherence, given an input of an erroneously transcribed script. Within the chatbot, the loss of information from transcription and structural organization surfaces as the conversation progresses. While chatbots communicate directly via textual means, the conversation structure permits the information to maintain higher data density than transcripts of spoken communication [12]. Hence, our research may need to address the acquisition of sufficient spoken data to construct domain models. Retrieving subsumed themes from previous and current conversations imposes on the intelligent system the additional task of verifying the accuracy of its inferences and responses.

Schumaker asserts that conversation length is an important metric in maintaining dialog quality [12]. ALICE-based bots support the need for evaluating information quality on the part of the user and computer to quantitatively assess the relevance of new data to the current or emerging contexts. Such a principle falls within the metadata frameworks advocated by [18]. From the results of Gurevych [15], the gold standard for this would be the consensus of human judges with the system's interpretation of the domain.

### C. *Recent Advances and Integration of Realism*

Since the inception of ELIZA [8], an influx of chatbot research has resulted. Evidence of such can be seen in the projection of human cognitive behavior and realism models onto applications with chat-based roots. In this section, we chronologically demonstrate the direction of chatbot research, describing various dialog systems and their associated pursuits in advancing naturalness.

Mateas [19] comparatively provides an overview of advances in chatbot related technologies for the late Nineties. Specifically, he demonstrates the initial departure, in this timeframe, from ELIZA-based [8] chatbots that employ sentence-based template matching. What is highlighted instead is the increased importance of developing simple conversational memory. Accordingly, early conversational memory-equipped systems include multi-user dungeons (MUDs), such as Carnegie Mellon University's Julia project, and Extempo's Erin. Mateas, however, notes several key differences between the conversational characteristics between these chatbots and more believable agents. Namely, interaction occurs in a reactive manner, with no regard for pursuit of a goal by the chatbot. Additionally, these systems were intended to perform under a constrained version of the Turing test, only briefly fooling its human users.

Wlodzislaw et al [20] further expand on the naturalness restrictions evident in the template matching approach of the ELIZA-styled programs. Earlier systems lacked domain expandability and could not fully exploit memory and reasoning components. Furthermore, they suggest that reliance on template matching can be associated with three key aspects of chatbots: 1) focus on the Loebner prize, 2) template-based AIML techniques, and 3) slow development of reasoning from natural language in dialog systems. From [20], we learn that the development of cognitive modules and human interface realism

for chatbot systems distinguishes avatars from ELIZA- or ALICE-based agents. As an example, Wlodzislaw et al [20] cite the use of ontologies, concept description vectors, semantic memory models, and CYC [24] as tools that can serve to replace AIML templates and to increase the impression of understanding by the agent.

Further research into chatbots saw a shift towards enhanced immersive reality for dialog systems, emphasizing face-to-face avatar presentations and dialog evaluation improvements [20] [21]. Traum and Rickel [21] identify two considerations that present challenges to dialog management: 1) multi-modal interaction, and 2) multi-party conversations. Becker and Wachsmuth [22] explore the representation and actuation of coherent emotional states in a virtual conversational agent. Lars et al [23] extends this research by presenting a model for sustainable conversation in a real-world application. They discuss several cognitive modules that increase the system's awareness of the human users and the conversation topics. The downside here is that the system relies on textual input similar to that of ELIZA.

Some interest has been generated on the use of natural language processing (NLP) for reasoning about human speech. However, several NLP applications may not be mature enough for implementation in conversational agents. Furthermore, the tasks involved differ from those of natural language generation, such as those tasks concerning agent knowledge acquisition. Moreover, a sense of dialog-based reasoning using NLP techniques can be gleaned by analyzing the works referencing such systems as CYC [24] and WordNet [25].

From the aforementioned approaches, we perceive that an emphasis exists on developing goal-oriented dialog systems that respond naturally. It is also important to note the breadth of research in which chatbot technology has embraced. We see that the principal efforts of this movement focus on creating more sophisticated interpretative conversational modules. Given the differences in techniques used to develop these bots, a need exists for generalizable metrics that evaluate the quality of a conversation in addition to the bot's performance. Hence, the underlying theme of the survey in this section dictates that the conversational agent topic has been widely experimented with, but it has been lacking a basic framework for universal performance comparison.

The following section takes this final sentiment to heart and frames it toward providing a solution for assessing an existing chatbot-based project – the LifeLike virtual avatar. Hence, we present an overview of the issues considered to build an appropriate evaluation method our conversational agent, with universality in mind.

### III. APPROACH

The development of LifeLike, as with any software creation, calls for a proper method of evaluating its performance. The challenging aspect of LifeLike, however, results from its identity as a vehicle of human behavior emulation. This means that the approach we will use for its evaluation process must incorporate elements of subjectivity from its human operators. This section discusses the duality of

qualitative and quantitative aspects needed for chatbot evaluation.

#### A. Previous Approaches

Previous attempts at evaluating conversation agents all reflect a mix of quantitative and qualitative measures. Typically, subjective matters have involved a human user questionnaire. Semeraro et al [3] employ this technique for their bookstore chatbot. In the questionnaire, seven characteristics were appraised: impression, command, effectiveness, navigability, ability to learn, ability to aid, and comprehension. Users would assess their associated satisfaction for each of these metrics, ranging from 'Very Unsatisfied' to 'Very Satisfied.' Semeraro et al recognize the fact that this subjective evaluation does not provide statistically verified conclusiveness, but rather it serves as a general indicator of performance.

Shawar and Atwell [4] propose a universal chatbot evaluation system. They suggest three metrics, which were applied upon an ALICE-based Afrikaans conversation agent. The first metric concerns dialog efficiency, which deals with: atomic matching types, first word matching types, most significant matching types and no matching types. These matching methods establish how effectively a chatting agent can respond to user input. In their testing, Shawar and Atwell saw that first word matching and most significant matching were the most competent techniques. The second metric is the dialog quality metric, which qualitatively categorizes, by human judgment, a chatting agent's responses into three bins: reasonable, weird but understandable, and nonsensical. The final metric is users' satisfaction, which is also qualitatively measured. Feedback from the chatting software end-users is collected and used to directly evaluate the agent's performance.

Despite their efforts to establish a set of generic metrics, Shawar and Atwell [4] discourage the use of such a universal conversation agent evaluation mechanism. Instead, they conclude that the proper assessment of chatbots is the end result in how successfully it accomplishes its intended goals.

Evaluation of maintaining naturalness in a conversation similarly suffers from the same inherent problems of the general chatbot assessment system. Again, subjectivity plays a large role in judging the naturalness of a conversation. Rzepka et al [2] used a 1-to-10 scale for two metrics: a "naturalness degree," and a "will of continuing a conversation degree." In this study, human judges used these measures to evaluate a conversation agent's utterances. While their assessment system did not identify a concrete baseline for universal naturalness, they were able to make relative measurements of naturalness between different dialog management approaches, such as comparing an ELIZA-based [8] manager with a world wide web-based commonsense retrieval system.

Chatbot evaluation remains an open problem, especially because of its dependence on subjective assessment. Researchers use questionnaire-based methods to provide general insight on the effectiveness of their conversation agents. Similarly, measuring conversational naturalness also relies on user subjectivity. The major pitfall of these evaluation methods is their lack of quantitative universality, as no set of

chatbot performance metrics has been defined. Nevertheless, current research has found success in using these techniques to make relative comparisons between conversation agents. Conversation agent evaluation, with emphasis on naturalness, plays a substantial role in appraising the performance of the work in this paper.

The remainder of this section gives a more in-depth treatment of the chatbot evaluation process, pointing out the primary factors that delineate the effectiveness of such dialog-based system software.

### B. Chatbot Objectives

A dialog system, especially those of the *assistive* nature (as in LifeLike) proves its effectiveness under the light of two primary objectives: 1) dialog performance, and 2) task success. Each of these aims reflects different aspects of a human-computer conversation. Dialog performance relates to the experience of the interaction, while task success is concerned with the utility of the dialog exchange. Basically, these two objectives separately assess the effectiveness of the means (dialog performance) and the ends (task success).

The main goal of a dialog system is to achieve task success and dialog performance levels that are: 1) better than other dialog system solutions, and 2) similar to a human-to-human interaction. The latter stipulation defines the measure of naturalness, where a dialog system that can conduct a conversation that is similar to one between two people is considered natural. The next sub-section provides the metrics necessary to measure task success and dialog performance.

TABLE I. CHATBOT METRICS

| Metric                                   | Type        | Data Collection Method |
|--|-------------|------------------------|
| Total elapsed time                       | Efficiency  | Quantitative Analysis  |
| Total number of user/system turns        | Efficiency  | Quantitative Analysis  |
| Total number of system turns             | Efficiency  | Quantitative Analysis  |
| Total number of turns per task           | Efficiency  | Quantitative Analysis  |
| Total elapsed time per turn              | Efficiency  | Quantitative Analysis  |
| Number of re-prompts                     | Qualitative | Quantitative Analysis  |
| Number of user barge-ins                 | Qualitative | Quantitative Analysis  |
| Number of inappropriate system responses | Qualitative | Quantitative Analysis  |
| Concept Accuracy                         | Qualitative | Quantitative Analysis  |
| Turn correction ratio                    | Qualitative | Quantitative Analysis  |
| Ease of usage                            | Qualitative | Questionnaire          |
| Clarity                                  | Qualitative | Questionnaire          |
| Naturalness                              | Qualitative | Questionnaire          |
| Friendliness                             | Qualitative | Questionnaire          |
| Robustness regarding misunderstandings   | Qualitative | Questionnaire          |
| Willingness to use system again          | Qualitative | Questionnaire          |

TABLE II. ATTRIBUTE-VALUE CONFUSION MATRIX [1]

| DATA | Departure City |           |          |          |
|------|----------------|-----------|----------|----------|
|      | ATL            | BOS       | CLT      | DEN      |
| ATL  | <u>16</u>      |           | 1        |          |
| BOS  | 1              | <u>20</u> | 1        |          |
| CLT  | 5              | 1         | <u>9</u> | 4        |
| DEN  | 1              | 2         | 6        | <u>6</u> |
| SUM  | 23             | 23        | 17       | 10       |

### C. Evaluation Metrics

The evaluation system featured in this paper is derived from the PARAdigm for Dialogue System Evaluation (PARADISE) [1]. Table I depicts the structure of the objectives and their corresponding metrics within PARADISE. Under this model, the master objective is user satisfaction, which is comprised of task success and dialog costs. Walker et al [1] further break down the dialog costs as a combination of efficiency measures and qualitative measures. These PARADISE-based objectives directly reflect the task success and dialog performance objectives mentioned in the previous section. The next sections discuss the metrics involved in task success and dialog costs.

### D. Task Success

The tasks involved with a dialog system are of a multiple-goal nature. Thus, for any conversation, all of these goals must be recognized and satisfactorily serviced for the entire task to be considered successful. Conversations are modeled as a set of attribute-value pairs. Every user goal (and sub-goal) corresponds to an attribute, and the dialog agent's response to those goals represents a value.

As in PARADISE [1], an attribute-value matrix is created for both the expected response and the actual agent response in a conversation. A confusion matrix is produced to identify the discrepancies between the expected and actual attribute-value pairings. Table II gives an excerpted version of Walker et al's example attribute-value confusion matrix [1].

Walker et al present an attribute-value confusion matrix for a travel schedule system, with Departure and Arrival attribute-value pairings [1]. Table II gives a representative depiction of this matrix. Let us assume the question asked to the travel scheduling chatbot is, "Which city has a departure time of  $X$  o'clock?" The rows represent the *actual* responses from the agent, and the columns reflect the *expected* values.

In this matrix, there are four possible values for the departure city question. The value ATL was correctly identified 16 out of 23 times, while DEN was agreed upon 6 out of 10 times. This type of accuracy data may be extrapolated from the attribute-value confusion matrix. From this information, task success,  $\mathcal{K}$  is computed as the percentage of 'right' responses given by the agent.

### E. Dialog Costs

Dialog performance is defined as a function of two types of dialog costs: efficiency and quality. Efficiency costs refer to

the resource consumption needed to accomplish a single task or sub-task. These attributes can be measured in a solely quantitative manner. Qualitative costs measure the actual conversational content. These metrics may be recorded quantitatively or qualitatively. For qualitative assessments, users are given a Likert scale-based questionnaire following their interactions, providing feedback on the dialog system’s naturalness, friendliness, etc. Walker et al [1], Stibler and Denny [5], Charfuelán et al [6], and Hassel and Hagen [7] provide some examples on suitable dialog costs. Table I delineates the relevant cost metrics for this paper.

#### F. Performance Function

To evaluate the total effectiveness of a dialog system in relation to its task success,  $\mathcal{K}$ , and its dialog costs,  $c_i$ , the following PARADISE-based [1] performance function is used

$$Performance = (\alpha * \mathbf{N}(\mathcal{K})) - \sum_{i=1}^n w_i * \mathbf{N}(c_i)$$

In this relationship,  $\mathcal{K}$  is weighted by  $\alpha$ , and each  $w_i$  is a weight on  $c_i$ . The weight assignments are established in an arbitrary, yet meaningful manner. The function,  $\mathbf{N}$ , uses a Z-score normalization process to balance out the effects of  $\mathcal{K}$  and  $c_i$  on the overall system performance.

This performance function allows for a normalized method of comparing two different dialog systems using the same conversational task goals. A dialog system will be considered performing in a natural manner when its performance level matches that of a human-to-human dialog.

The approach considerations presented in this section preludes into the proposed evaluation system prescribed for the LifeLike project. The following section provides a description of the system’s prototype.

#### IV. EVALUATION SYSTEM PROTOTYPE

The purpose of establishing a chatbot evaluation system is to devise an experimentation infrastructure that collects data to support the idea of an improved human-computer interaction experience over other conversation agent systems. Furthermore, this interaction must reflect closely to the characteristics of a human-to-human exchange under the same situational premises. The quality of the interactive experience is judged using the previously mentioned metrics and task success assessment method, ultimately using the collaborative performance function to give a single measure of its effectiveness.

During the experimentation process, a single conversational scenario is employed on five different dialog systems. The first is the fully operational Context-Based Reasoning (CxBR)-based dialog system [10] featured in the LifeLike avatar project. The second dialog system is a crippled version of the first one, where the dynamic context-switching functionality is turned off. In a third experiment, an ELIZA-based [8] dialog system is tested. The fourth dialog system is modeled after an automated-phone operator, and the final system utilizes what is known as a Wizard of Oz (WOZ) experiment [9], where a human interlocutor replaces the normally machine-based agent.

In each experiment, the user is assigned five pre-specified goals to achieve during his or her dialog system interaction. A verbatim log of each conversation is retained for quantitative analysis and the user fills out a system quality questionnaire at the conclusion of the experiment. These data sources are used to compute the performance measure of the dialog system. Five different users will test each system. The user base will be selected under the assumption that cultural bias should not be a major factor when compiling results.

After executing all 25 trials (5 systems, 5 trials per system), the performance of each agent is compiled and evaluated for comparative analysis. Upon careful examination of these results, conclusions regarding the dialog system can be made. Future iterations of this experimentation process can be used for comparing later builds or enhanced versions of the same chatbot. This is analogous to “normal” software engineering practices, where unit testing is employed to verify that baseline functionalities are still intact between iterations.

#### V. CONCLUSION

This paper focuses on the inherent challenge of providing a proper evaluation process for conversation agent software. An overview of historical background technologies were presented to attest to the idea that such a problem is truly a challenge in the software engineering realm. We specifically sought out the validation requirements of the LifeLike virtual avatar to frame the chatbot evaluation problem in a real-world treatment. A proposed assessment method for LifeLike was presented, a prototypical framework derived from the PARADISE [1] infrastructure.

The advancement of chatbots provides hope for a future of HCI exchanges that go beyond the keyboard and mouse. A more immersive and personalized touch can be added to human-machine relations when such technologies are paired with the appropriate application. In relation to LifeLike [26], it is our hope that a proper information-dispensing system can be utilized without jeopardizing the ‘talking with a human’ experience. While the current LifeLike system deals primarily in the NSF realm, it is assumed that future iterations will have implications in medical assistance, military training, and education facilitation. The work described in this paper harkens the call for a bolstering of conversational agent evaluation, as to strengthen the effectiveness and efficiency of chatbot research.

#### VI. FUTURE WORK

The LifeLike virtual avatar is still in its prototypical stage, which precludes the immediate requirement for a verification and validation process. Preliminary evaluations of the avatar software’s prototypes have been made, with much of the aforementioned material taken into consideration. A formal treatment of the evaluation process has yet to be implemented; thus, any formidable results have yet to be prepared in publishable form.

As our work with LifeLike progresses, so will our need to provide a reliable evaluation system. Establishing such a method will allow us to better judge the evolutionary direction of our software, as well as any other chatbot software outside of LifeLike.

#### ACKNOWLEDGMENT

This research is supported by NSF Collaborative Research award 0703927.

#### REFERENCES

- [1] M. A. Walker, D. J. Litman, C. A. Kamm, and A. Abella, "PARADISE: a framework for evaluating spoken dialogue agents," Proc. of the 8th conference on European chapter of the Association for Computational Linguistics, 1997, pp. 271-280.
- [2] R. Rzepka, Y. Ge, and K. Araki, "Naturalness of an utterance based on the automatically retrieved commonsense," Proc. of Nineteenth IJCAI, 2005.
- [3] G. Semeraro, H. H. K. Andersen, V. Andersen, P. Lops, and F. Abbattista, "Evaluation and validation of a conversational agent embodied in a bookstore," Universal Access: Theoretical Perspectives, Practice and Experience. Lecture Notes in Computer Science, 2615, 2003, pp. 360-371.
- [4] B. A. Shawar, and E. Atwell, "Different measurements metrics to evaluate a chatbot system," Proc. of the 2nd Workshop on TextGraphs: Graph-Based Algorithms for Natural Language Processing, 2007.
- [5] Stibler, K., and Denny, J., "A three-tiered evaluation approach for interactive spoken dialogue systems," Proceedings of the Third international conference on Human language technology research, 2001, pp. 1-5.
- [6] Charfuelán, M., Gómez, L. H., López, C. E., and Hensen, H., "A XML-based tool for evaluation of SLDS," Proceedings of the Third International Conference on Language Resources and Evaluation, 2002.
- [7] Hassel, L., and Hagen, E., "Evaluation of a Dialogue System in an Automotive Environment," 6th SIGdial Workshop on Discourse and Dialogue, 2005.
- [8] J. Weizenbaum, "ELIZA-a computer program for the study of natural language communication between man and machine," Communications of the ACM, 9(1), 1966, pp. 36-45.
- [9] Hajdinjak, M. and Mihelic, F., "A Dialogue-management Evaluation Study," Journal of Computing and Information Technology, 15(2), 2007, pp. 111-121.
- [10] B. Stensrud, G. Barrett, V. Trinh, and A. Gonzalez, "Context-based reasoning: a revised specification," FLAIRS Conference, 2004.
- [11] d'Ydewalle, G., and Delhay, P., "Artificial intelligence, knowledge extraction and the study of human intelligence," International Social Science Journal, 1988, pp. 63-72.
- [12] Schumaker, R. P., Liu, Y., Ginsburg, M., and Chen, H., "Evaluating mass knowledge acquisition using the ALICE chatterbot: The AZ-ALICE dialog system," International Journal of Human-Computer Studies, 2006, pp. 1132-1140.
- [13] Anderson, J., "ACT - A Simple Theory of Complex Cognition. American Psychologist," 1996, pp. 355-365.
- [14] Le Bigot, L., Terrier, P., Amiel, V., Poulain, G., Jamet, E., and Rouet, J.-F., "Effect of modality on collaboration with a dialogue system," International Journal of Human-Computer Studies, 2007, pp. 983-991.
- [15] Gurevych, I., Malaka, R., Porzel, R., and Zorn, H., "Semantic coherence scoring using an ontology," Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology - Volume 1. Morristown: Association for Computational Linguistics, 2003, pp. 9-16.
- [16] Pirolli, P., and Recker, M., "Learning Strategies and Transfer in the Domain of Programming," Cognition and Instruction, 1994, pp. 235-275.
- [17] Potter, S., "A Survey of Knowledge Acquisition from Natural Language," TMA of Knowledge Acquisition from Natural Language, 2001.
- [18] Hahn, U., and Schnattinger, K., "An Empirical Evaluation of a System for Text Knowledge Acquisition," EKAW '97: Proc. of the 10th European Workshop on Knowledge Acquisition, Modeling and Management, 1997, pp. 129-144.
- [19] Mateas, M., "An Oz-CentriV Review of Interactive Drama and Believable Agents," In Artificial Intelligence Today: Recent Trends and Developments, 1997, pp. 297-343.
- [20] Wlodzislaw, J., Szymański, J., Sarnatowicz, T., "Towards Avatars with Artificial Minds: Role of Semantic Memory," Journal of Ubiquitous Computing and Intelligence, 2006, 1(in print).
- [21] Traum, D. and Rickel, J., "Embodied agents for multi-party dialogue in immersive virtual worlds," Proceedings of the First International Joint Conference on Autonomous Agents and Multiagent Systems: Part 2, 2002, pp. 766-773.
- [22] Becker, C. and Wachsmuth, I., "Simulating the emotion dynamics of a multimodal conversational agent," Proc. Tutorial and Research Workshop on Affective Dialogue Systems (ADS-04), 2004, pp. 154-165.
- [23] Lars, G., Krämer, N., Wachsmuth, I., "A conversational agent as museum guide - Design and evaluation of a real-world application," The 5<sup>th</sup> International Working Conference on Intelligent Virtual Agents, 2005, pp. 329-343.
- [24] Lenat, D. B., "CYC: A large-scale investment in knowledge infrastructure," Communications of the ACM, 1995, 38(11).
- [25] Miller, G., Beckwith, R., Fellbaum, C., Gross, D., and Miller, K., "WordNet: an on-line lexical database," International Journal of Lexicography, 1990, 3(4).
- [26] DeMara, R. F., Gonzalez, A. J., Jones, S., Johnson A., Leigh, J., Hung, V., Leon-Barth, C., Dookhoo, R., Renambot, L., Lee, S., and Carlson, G., "Towards Interactive Training with an Avatar-based Human-Computer Interface," Interservice/Industry Training Systems and Education Conference, 2008.