



progress and give personalized feedback. The role of the simulated characters is to engage the learner face-to-face as part of an interactive story, dynamically tailored to the learner's level of skill.

The overarching research goal is to study and employ computer aided learning techniques that effectively prepare learners for basic spoken social interaction in a foreign language environment. Key pedagogical research themes, that we believe are crucial in attaining effectiveness, include adaptive instruction and engaging practice. These themes further break down into perception and evaluation of learner input, learner modeling, feedback generation, modeling social interaction, interactive story and believable characters. This work combines and extends previous CARTE work on socially intelligent pedagogical agents [1,2,3,4] and interactive pedagogical dramas [5], and draws from previous work on the coordination of verbal and nonverbal behavior in ECAs [6,7].

This paper reports on work in progress (the project was launched about a year ago). Thus, while advanced prototypes have demonstrated strengths of the approach and provided solutions to many of the technical issues involved, the system is expected to evolve over iterative cycles of design, development and evaluation.

## **2 Requirements**

The requirements for the MSB tutoring agent and the MPE character agents differ and will be addressed in turn.

### **2.1 The MSB Tutoring Agent**

The relationship between the learner and the agent in the MSB is modeled after one-on-one tutoring, which naturally involves social interaction including conversation. The agent is given a single persistent identity, that of a native speaker of the language, to emphasize the one-on-one nature of the tutoring. Because language tutoring involves clearly showing with your mouth how the sounds of the language are made, a face for the agent was essential. By keeping the tutor's face on the screen at all times, not just when helping with pronunciation, it is possible to provide guidance and feedback, both verbal and nonverbal, at all times. Moreover, the mere visual presence may act as a reminder that the tutor is indeed keeping track of what the learner is doing. An ECA approach therefore lends itself well to this kind of pedagogical agency.

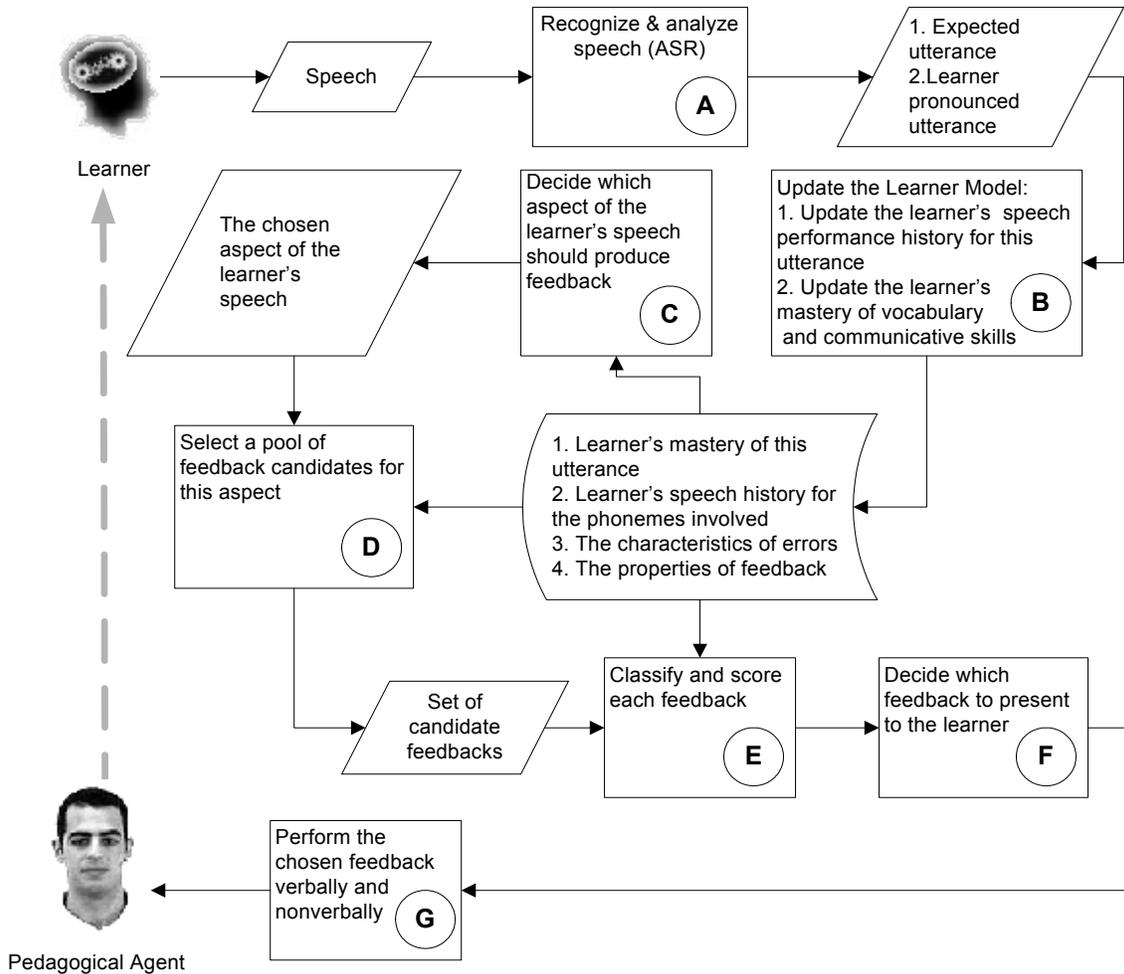
While the tutoring agent is responsible for introducing the learner to the various language and culture lessons, the most important task the agent currently performs is

teaching how to pronounce Arabic phrases. This is where an effective perception/feedback loop is essential. The tutor first needs to demonstrate the correct pronunciation and then ask the learner to try. After listening to the learner's pronunciation, the agent needs to give feedback that is both sensitive to any pronunciation errors being made and to the learner's overall progress. These two concerns call for specialized speech recognition and learner modeling, which will both be discussed under the Perception Loop section of this paper.

### **2.2 The MPE Character Agents**

The agents in the MPE are essentially virtual actors playing roles in an interactive story under the influence of a high-level director agent that makes sure overall pedagogical and dramatic goals are met. The MPE is implemented as a multi-agent system where each agent can perceive, maintain beliefs and act. The fact that the agents represent people in a complete virtual environment requires of course that they have bodies and that they exhibit believable verbal and nonverbal behavior. Considering that the learner also needs to be able to have conversations with these agents, the ECA framework seems to perfectly apply.

In most current games the interaction with non-player characters is not well balanced. The user typically "says" things to them by choosing text from a menu and has no way of providing accompanying nonverbal behavior while the characters respond with full speech and gesture. The MPE brings the interaction to a whole new level of symmetry by requiring the learner to actually speak with the characters and select the appropriate gesture. It is important that the learners get the impression that they are in face-to-face dialog with the characters because making them feel they really are capable of engaging someone face-to-face in a foreign language contributes to their confidence and feeling of accomplishment. The agents, of course, need to be able to understand what the learner is saying and respond appropriately with the rapidity typical of natural conversation. Appropriate speech and gesture responses need to be determined from the speech and gesture input in the context of the ongoing story and with regard to what language skills the learner has mastered so far. A special consideration is that an important part of the training is how to exhibit proper nonverbal behavior given the cultural setting. Therefore, the agents need to provide believable and accurate social behavior, both in and out of conversation, and be able to clearly react to the learner's behavior choices. These aspects, adaptive characters and gesture interaction in the MPE, will be further discussed under the Perception Loop section.



**Figure 3: A flowchart demonstrating the flow from perception to action in the pedagogical tutoring agent during a pronunciation exercise.**

### 3 Perception Loop

#### 3.1 The MSB Tutoring Agent

Figure 3 is a flowchart that demonstrates how the perception of the learner's speech during a pronunciation exercise, together with information from a learner model, produce appropriate feedback from the pedagogical tutoring agent. The rest of this section describes the processes involved in more detail.

Process A maps microphone sound input to transliterated Arabic words. Most commercial automated speech recognition (ASR) systems are not designed for learner language [8], and commercial computer aided language learning (CALL) systems that employ speech tend to overestimate the reliability of the speech recognition technology [9]. To support learner speech

recognition in the TLT system, our initial efforts focused on acoustic modeling for robust speech recognition especially in light of limited domain data availability [10]. Recognition, here, is accomplished by a Hidden Markov Model ASR system that has been bootstrapped from Modern Standard Arabic speech and enhanced with data from native and learner Lebanese Arabic speech. Our speech recognition engine departs from the standard ASR task in two notable ways: First, we not only recognize true Arabic words, but also mispronounced and misuttered Arabic words. Secondly, because we are dealing with learner speech, we need deal only with a smaller subset of language that is applicable to that which the learner has been taught, so we can safely reduce our recognition vocabulary size, simplifying the problem. This simplification is necessary, because supplementing our base Arabic recognition grammar with disfluencies upon each item in the original grammar increases the

HMM state size to the degree that robust detection would be untenable otherwise.

The speech recognizer assigns confidence values to each utterance it recognizes, classifying the sound sequence of the learner's utterance using two grammars—one of the set of all Arabic words that the learner could possibly use, and one grammar of the set of mistakes that the learner could be expected to make when attempting to say the specific utterance in question [11]. By considering the confidence value scores of the utterance as it fits into these two classification grammars, a wrapper around the speech recognizer is able to choose the most serious error that lead to the learner's speech failure. This approach compensates for inaccurate recognition by the ASR, avoiding misdetection (and therefore mis-correction) of nonexistent mistakes made by learners.

We approximate the different errors we can expect at any point in the user's speech using a probabilistic model derived from real learner speech data [11]. This modeling takes into account many different aspect of language-learner speech. Our current system supports:

- *Phonological errors:*  
Mistakes arising from interference with the learner's native language phoneme set.
- *Orthographic/Phonological errors*  
Mistakes arising from interference between transcription systems and native and new phoneme sets.
- *Lexical/Semantic errors:*  
Mixing up words in a lexicon because of their ontological or phonemic similarity.
- *Cognitive errors:*  
Dropping different words in a sentence because the cognitive load is too high to remember them all at once.

Additionally, development is underway to address the following:

- *Morphological errors:*  
Mistakes due to confusion of tense, number/gender marking.
- *Syntactic errors:*  
Mistakes due to differences between syntactic parameters in the native and new languages.

This model is applied to our base recognition grammar and augments each fluent word or phrase to be detected with a supplemental top  $n$  most likely disfluent words or phrases.

Information about the agent's knowledge on Arabic teaching materials and perception of the learner's learning progress is organized hierarchically into two relational databases in the TLT system. The curriculum database stores the content and description for lessons, exercises, vocabulary, cultural and grammar items (from general objectives to fine-grained units), serving as the ontology

of the Pedagogical Agents. Our tutor agent's pre-ordered feedback and possible learner errors are also characterized and kept in this database.

The learner model database is similarly hierarchically structured from lesson to phoneme level, which enables the agent to capture the learner's progress in depth. Process **B** in Figure 3 depicts the agent's internal process of updating the learner model database. Once a recognized utterance is received from the ASR, the agent inserts a new record of the learner pronunciation into the database, including the learner's phonetic error(s) and a timestamp. Meanwhile, the agent updates its knowledge about the learner's mastery of vocabulary and skills, and checks if the last utterance by the learner completes the current exercise or the whole lesson.

In process **C** we are given a set of errors associated with the recognized utterance, and must decide which (if any) to give feedback on. This problem is often nontrivial because of the indefinite and overlapping nature of these errors. First, we must deal with both the possibility that the error was falsely detected due to noise in the speech recognition. Additionally, we must deal with the fact that multiple kinds of errors can appear the same. Given a disfluent utterance, for example, we do not know if the learner has misremembered the vocabulary, or if she is only having problems producing certain sounds.

This decision-making is accomplished by means of a rule-based system that looks at, among other things, the following:

- Our learner's history of *not making* the mistake in question—if the learner has had few problems with this sort of errors in the past, we give more credence to the possibility that an error in speech recognition has occurred.
- Our learner's history of making the mistake in question—while common errors should be corrected, if our learner has made a mistake too many times, we will temporarily “give up” and move on to more correctable errors.
- The seriousness of an error, taken in context of offending cultural norms—a learner making a morphological error and addressing a man as a woman, for example, is a much more serious cultural error than mispronouncing a sound, and should be dealt with accordingly.
- The seriousness of an error, taken in context of understanding—some disfluencies more seriously influence listener understanding than others (for example, a learner misordering a noun and its adjective does not affect a listener's understanding as much as the learner

switching two nouns in a sentence). Special care needs to be taken to eliminate these errors.

These factors, and others, combine together to help us decide both the agent's confidence that a learner has committed a given error, and also how much this error merits correction. From them, we select the error that maximizes both, and pass it on to the next step.

Process **D** shows the process of selecting a pool of feedback candidates given the error chosen for correction. This process takes the properties of the feedback and the predefined relation between the errors and feedback as an input. Its outcome, a set of feedback candidates, are classified and scored in the next process (Box E). The objective is to select a most suitable feedback to present to the learner. Currently, there are four classifiers involved in the decision-making process: *Polarity*, *Strength*, *Type* and *Freshness*. The classifier "Polarity" has three possible values: positive, negative and neutral, which indicates whether to select the feedback to motivate learners when they speak right, or to correct learners when they make mistakes in speech, or to remind the learner about other problems, e.g. a misplaced microphone that blocks the speech input. The classifier "Freshness" is used to make decision based on when the feedback is called for. In our feedback system, some feedback is prepared for fresh learners, who have never practiced a specific utterance, some for learners who have practiced it before and some for both kinds of learners. The classifier "Type" identifies the purpose of the correction. For instance, feedback like "Good job", "Please try it again" are labeled "general" while feedback for correcting lexicon errors is of the "word" type and feedback for correcting phonetic errors is of the "phoneme" type. Classifiers like "Freshness" and "Strength" mainly consider learner's speech history; classifiers like "Polarity" and "Type" mainly focus on the property of the feedback itself.

A preliminary feedback selection is performed based on a relational table in the curriculum database, where the relationship between errors and feedback is grossly specified. Given an error in the user's speech, the agent then winnows the first batch of chosen feedback by filling in a form of requirements that restricts the favored feedback properties and submitting it to a feedback model. "Bad" feedback that does not meet requirements is scored low, while the scores of "good" feedback is boosted by the classifiers. Currently, the agent is using a rule-based mechanism to fill the requirement form. A probabilistic model has been planned to replace it in the next version. Box **F** describes a simple process that selects the feedback with the highest score, which heuristically is the most desired feedback that the agent should provide to the user. Process **G** refers to the presentation of the feedback to the user. The feedback

includes a prerecorded audio/video clip of the tutor's instruction and screen text (we will be looking into using an animated character, but for now the recording provides better fidelity). Before playing the tutor's instruction clip, the agent allows the learner to hear how their speech got recorded, directly followed by the agent's feedback. The textual feedback prints out the text form of the learner's pronunciation detected by ASR, and the text form of the utterance the agent expected the learner to say.

**Learner (L):** maHaba

**Pedagogical Agent (A):** I was expecting you to say "marHaba".

**A:** It takes some time to learn how to pronounce the r correctly. 'rrrr'

**L:** marHaba

**A:** Correct! - Good Job! Try Again!

**L:** marHaba

**A:** Correct! - Good Job! One more time!

**L:** marHaba

**A:** Correct! - Good Job! [*A green tick mark appears next to this utterance*]

**L:** maHaba

**A:** I was expecting you to say "marHaba". It takes some time to learn how to pronounce the r correctly. 'rrrr'

**L:** marhaba

**A:** I was expecting you to say "marHaba". Remember, capital h is 'H', and is quite different from lowercase h. 'H', 'h'.

**L:** marhaba

**A:** I was expecting you to say "marHaba". The trick is to roll the tongue in the mouth: 'rrrr'.

**L:** marhaba

**A:** I was expecting you to say "marHaba". Don't worry too much about trilling the r, you will pick it up soon.

[*Learner tries more than 10 times and gets stuck on the rolling sound 'rrrr'*]

**A:** Please go on for now.

#### **Figure 4: Sample speaking turns between the learner and agent during an exercise in saying "hello" in Arabic**

To illustrate the system in action, Figure 4 provides a transcript of several turns between a learner and the agent. In this example the learner is practicing to say "marHaba" (hello), but first leaves out the 'rrrr' sound (or pronounces it indistinctly). The agent finds a feedback that is negative (Polarity) with less serious tone (Strength) and that is prepared for phonetic errors (Type) and for the first mispractice (Freshness). Then the learner's two successes with pronouncing the utterance "marHaba" improves the learner's performance history for this phoneme, in terms of total correct pronunciation ratio and the timing of correct pronunciation. Thus, when the learner makes the

same mistake again, the agent still offers feedback with a light tone. The third mistake made by this learner is a mispronounced "H" as "h" (these two are different phonemes in Arabic). The agent is able to detect this error and provides feedback related to this phonetic error. This error increases the user's unsuccessful tries on "marHaba", but doesn't affect the learner's success history on phoneme 'r'. Thus, when the learner fails to speak "marHaba" correctly again, the agent chooses a feedback with higher strength to encourage the learner. After correcting the learner on the same mistake for 10 times, the agent actually gives up correcting the learner and lets him skip practicing this utterance for now.

### 3.2 The MPE Character Agents

In MPE, the learner controls a character in the game environment that represents him. The learner can walk around and interact with other characters through conversation and gestures. In the scene depicted in Figure 2, the learner is having a conversation with two characters (multi-party conversation), a young man (sitting) and an older man (standing). The soldier standing right behind the learner is an aide that follows the learner around and can provide assistance when needed. Finally, there are "extras" sitting or standing around the café that provide a dynamic backdrop much like one would experience in a film.

The characters are all controlled by a multi-agent system where there is an agent for each of the characters. The system receives events from the environment, such as when the learner's character walks into the café, speaks or gestures. The learner speaks by toggling into a recording mode with a mouse click (see the red icon in Figure 2), which then generates a speech event with the output from a speech recognizer. The learner can select a gesture to go along with the speech with the mouse wheel (see the green icon in Figure 2). We plan to experiment with vision-based gaze tracking as another way getting multi-modal input for further achieving balanced face-to-face communication. The system also has access to the same learner model that was described in the previous section. Based on these inputs, the multi-agent system as a whole decides how to respond, based on what each agent wants to do.

The running phase of this simulation can be divided into three stages. During the first stage, each agent updates its beliefs about the state of the world and the other agents based on available inputs. For example, if the learner has just announced that his mission is peaceful, both the old man and the young man will believe that they can trust him more than before. The second stage involves deciding how the agents should respond to the event. Currently, we have a simple mechanism for deciding who takes the dialog turn. After the learner speaks, each

character (in a predefined order) can decide whether it wants to respond (based on the recent dialog). This may mean for example, that the learner says something, the old man may respond and then the young man has the option to say something as well (in response to either the student or the old man). However, this strict order can be altered by arousal level. Each character has an arousal level, which indicates how angry or worried this character is. The arousal level is updated each time a new action is perceived. Occasionally, a character may have a very high arousal level, guaranteeing that the character gets the next turn, no matter whom it normally belongs to. For example in the café scene, the young man starts out trusting the learner less than the old man does. If the learner fails to build sufficient trust (for example, failing to describe his mission or using the inappropriate gesture), the young man will interrupt the dialog between the learner and the old man and accuse the learner of being a spy. In the third and last stage the character that has the turn decides what action to take. This character's action will then be added to the queue of actions that can be perceived by all characters.

Once an agent has decided on the action to take, a layer termed the social puppet layer coordinates the realization of the action by the graphical puppet that represents the agent in the virtual environment. This layer is responsible for planning the actual verbal and nonverbal behavior that appropriately and expressively realizes the agent's communicative intent given the social context, based on a model of particular culture and language. The plan is that this layer will also generate appropriate reactive behavior in all the puppets involved in the scene according to social rules, such as glances and posture shifts, to reflect the tight coordination of behavior by all members of any social gathering. This layer finally hands precisely timed behavior descriptions down to articulated puppets inside the game engine (we are using the Unreal Tournament 2003 game engine) that animate the behavior using a mixture of procedural and key-framing techniques.

The difficulty of a scene can be adjusted according to learner's language ability as reported by the learner model, so that learner will have an experience that is hard enough to provide good practice, but not so hard it leads to frustration. The difficulty of scenes can be adjusted by altering the personalities and goals of the virtual characters. For example, for a learner that is starting out with low language ability, it will take less to convince the virtual characters that the purpose of learner's mission is to help their village, and thus the above confrontation would be avoided.

The underlying agent technology used for the characters is the PsychSim multi-agent system [12]. PsychSim was chosen for several reasons. It models social relationships and reasoning, a key requirement for

realizing the social and cultural interactions required in TL. For example, PsychSim models factors such as trust and support between agents. Agents also have mental models of other agents and can employ those models to inform their decision-making about whether to believe another agent, what action to take, etc. In addition, the agents are realized as Partially Observable Markov Decision Processes (POMDP). Partial observability provides a mechanism to populate an entire “world” where agents may not have access to complete set of observations.

Note that the way we are regulating the interaction in general and the turn taking specifically is currently too simple. But we are designing a director agent that can access and change all other agent’s beliefs, motivations and models of others as well as manipulate the turn taking. Its goal is to initialize all other agent’s mental models appropriately, so that both dramatic and pedagogical goals will be achieved as the story unfolds. The director will also be a PsychSim agent.

## 4 Evaluation

System and content evaluation is being conducted through a staged, systematic process that involves both critiques from second language learning experts and feedback from learners. So far, learners at the US Military Academy and at USC worked through the first set of lessons and scenes and provided feedback, which in turn was used to refine the interaction model. The interaction model described in this paper was motivated in part by the feedback from those evaluations, e.g., the need for more detailed feedback on learner errors. Meanwhile recordings of learner speech are automatically recorded during use, to provide data to support further improvement of the speech recognition and error models.

A formative evaluation of the version of the system described in this paper was performed with a set of five college-age subjects at USC in May 2004. This evaluation revealed some problems with the current version. Subjects were somewhat reluctant to enter the MPE, for fear that they would not know how to communicate with the non-player characters and not know how to play the game. Although the subjects were informed that there was an aide character who could assist them, this was insufficient to give learners confidence to try the game. We are therefore modifying the game to include an introductory practice mode, in which learners can become familiar with the multimodal conversational interface and practice conversation until they are comfortable entering the game story. The evaluation revealed some problems in the MSB Tutoring Agent’s interaction with the learners. The tutor would critique the learners’ pronunciation until it was completely error-free; this was frustrating for beginning learners, who had some

difficulty even hearing the phonetic distinctions that the tutor was asking them to make, e.g., the distinction between emphatic and non-emphatic Arabic consonants. This was complicated by the fact that many of the practice utterances were quite short, often single phrases, and speech recognition confidence tends to be lower with very short utterances. As a result, the tutor would repeatedly reject learner pronunciation attempts for these short utterances. We are therefore modifying the Tutoring Agent to permit lower levels of pronunciation accuracy in learner speech in beginning lessons, and adjusting the length of the practice utterances so that better recognition confidence can be achieved.

## 5 Acknowledgements

The project team includes, in addition to the authors, CARTE members Catherine M. LaBore, Dimitra Papachristou, Carole Beal, David V. Pynadath, Nicolaus Mote, Shumin Wu, Ulf Hermjakob, Mei Si, Nadim Daher, Gladys Saroyan, Hartmut Neven, Chirag Merchant, Brett Rutland, Alfred Au, Prasan Samtani, Michael Bedernik, Tomer Mor-Barak and Andrew Marshal. From the US Military Academy COL Stephen Larocca, John Morgan and Sherrie Bellinger. From the USC School of Engineering Shrikanth Narayanan, Naveen Srinivasamurthy, Abhinav Sethy, Jorge Silva, Joe Tepperman and Larry Kite. From Micro Analysis and Design Anna Fowles-Winkler, Andrew Hamilton, Beth Plott, and Ursula Lauper. From the USC School of Education Harold O’Neil, and Sunhee Choi, and Eva Baker from UCLA CRESST. This project is part of the DARWARS initiative sponsored by the US Defense Advanced Research Projects Agency (DARPA).

## 6 References

- [1] Johnson, W.L., Rickel, J., & Lester, J. (2000). Animated pedagogical agents: face-to-face interaction in interactive learning environments. *IJAIED* 11, 47-78.
- [2] Johnson, W.L. (2003). Interaction tactics for socially intelligent pedagogical agents. *IUI 2003*, 251-253. New York: ACM Press.
- [3] Johnson, W.L., Rizzo P. (2004a). Generating socially appropriate tutorial dialog. *Affective Dialog Systems 2004*, in press.
- [4] Johnson, W.L., Rizzo, P., Pain, H. (2004b). Politeness in tutoring dialogs: “Run the factory, that’s what I’d do.” *ITS 2004*, in press.
- [5] Marsella, S., Johnson, W.L. and LaBore, C.M (2003). An interactive pedagogical drama for health interventions. In Hoppe, U. and Verdejo, F. eds. *Artificial Intelligence in Education: Shaping the Future of Learning through Intelligent Technologies*, Amsterdam: IOS Press.

- [6] Cassell, J., Vilhjalmsson, H., and Bickmore, (2001). BEAT: the Behavior Expression Animation Toolkit. SIGGRAPH 2001, Los Angeles, August 12-17, p.477-486.
- [7] Cassell, J, Bickmore, T., Campbell, L., Vilhjalmsson, H., and Yan, H. (2001). More Than Just a Pretty Face: Conversational Protocols and the Affordances of Embodiment. *Knowledge Based Systems* 14: 55-64.
- [8] LaRocca, S.A., Morgan, J.J., & Bellinger, S. (1999). On the path to 2X learning: Exploring the possibilities of advanced speech recognition. *CALICO Journal* 16 (3), 295-310.
- [9] Wachowicz, A. and Scott, B. (1999). Software That Listens: It's Not a Question of Whether, It's a Question of How. *CALICO Journal* 16 (3), 253-276.
- [10] Srinivasamurthy, N. and Narayanan, S. (2003). Language-adaptive Persian speech recognition. Eurospeech (Geneva,Switzerland).
- [11] Mote, N., Sethy, A., Silva, J., Narayanan, S., Lewis Johnson, W.L. (2004). Detection and modeling of learner speech errors: The case of Arabic tactical language training for American English speakers. InSTIL 2004, in press.
- [12] Marsella, S.C. & Pynadath, D.V. (2004). Agent-based interaction of social interactions and influence. Proceedings of the Sixth International Conference on Cognitive Modelling, Pittsburgh, PA, in press.