# Learning Cultural Conversational Protocols with Immersive Interactive Virtual Humans

Sabarish V. Babu1 [1], Evan Suma [2], Larry F. Hodges [1] and Tiffany Barnes[3]

[1] School of Computing, Clemson University
[2] Institute of Creative Technologies, University of Southern California
[3] College of Computing and Informatics, University of North Carolina at Charlotte

*Abstract*—**This paper reports on a study conducted to investigate the effects of using immersive virtual humans in natural multi-modal interaction to teach users cultural conversational verbal and non-verbal protocols in south Indian culture. The study was conducted using a between-subjects experimental design. We compared instruction and interactive feedback from immersive virtual humans against instruction based on a written study guide with illustrations of the cultural protocols. Participants were then tested on how well they learned the cultural conversational protocols by exercising the cultural conventions in front of videos of real people. Subjective evaluations of participants' performance was conducted by three south Indian reviewers who were blind to the condition the participants were assigned. Objective evaluations of participants' performance were conducted on the motion tracking log data recorded during the testing session. We also measured the participants' pre and post positive and negative affect of training in both conditions, as well as the effect of co-presence with the life-size virtual south Indians. The results of our subjective evaluation suggest that participants who trained with the virtual humans performed significantly better than the participants who studied from literature. The results also revealed that there were no significant differences in positive or negative affect between conditions. However, overall for all participants in both conditions, positive affect increased and negative affect decreased from before to after instruction.**

*Index Terms* —**Embodied Agents, Human-Computer Interaction, Social Interaction, Virtual Environments, Virtual Humans;**

## I.    INTRODUCTION

Research evidence has demonstrated that advances in virtual human interfaces and immersive virtual reality technology are effective in supporting multi-modal interaction with users for a variety of tasks. Anthropomorphic immersive virtual characters can use several modalities for communicating information, such as gestures and facial expressions, which are "transparent" to the user [5]. Virtual humans have been used successfully in diverse inter-personal collaborative tasks such as patient interview training [11, 25], virtual tutoring [12], mission rehearsal exercises in high stress scenarios [8], and virtual task-oriented companions in social settings [5, 1]. The benefits of using virtual human interface agents comes from the strength of the virtual human metaphor and leverages people's experience

with real social interaction to enrich the human-virtual human interaction.

Virtual humans have the potential to engage and train human users in tasks that involve interpersonal verbal and non-verbal conversational behaviors and face-to-face social contact. In order to assess the potential of using immersive virtual humans in natural multi-modal interaction to train users in social verbal and non-verbal behaviors, we investigated two basic questions in this study:

1. Is it possible to train users in the verbal and non-verbal behaviors associated with real life novel cultural conversational protocols through interaction with immersive virtual humans?

2. How effective is the use of immersive virtual humans as a tool for training users in cultural conversational protocols as compared to a text-based approach using a written study guide with illustrations?

This study explores the use of immersive virtual humans to train users in cultural conversational behaviors pertaining to conversation initiation and disengagement in south Indian culture. Conversational behaviors in south Indian culture are highly structured and specific to the gender, age, and status of the interlocutor. The temporality, intensity, and synchronicity of the verbal greetings, non-verbal gestures, and eye gaze are well-defined by the rules of etiquette for conversation initiation and disengagement [20, 7, 10]. Learning of these cultural social behaviors for most people comes through social grounding, interactive feedback, and reinforcement from an immersive experience in the culture.

**Our hypothesis: Natural multi-modal interaction with immersive virtual humans can successfully train naïve users in south Indian social conversational protocols.**

In order to evaluate our hypothesis we performed a novel study where we compared natural multi-modal interaction with immersive virtual humans to reading a written study guide with illustrations of cultural protocols. Participants in both conditions were given equal amounts of training time, and were then asked to demonstrate their ability to greet and say goodbye in response to video presentations of real south Indians. This paper extends results of a previous paper published at the IEEE VR 2007 [2].

In section 2, we address related works in systems that explore the benefits of simulating human-virtual human dialogue towards pedagogy and affective interfaces. Section 3 describes the design of the immersive virtual human system that includes multi-modal interaction and feedback. Section 4 details our study design and procedures. Section 5 presents a discussion of

our results.

Virtual humans in virtual environments that are used to teach and train cultural conversational behaviors in foreign culture are important in applications such as military training simulators, inter-cultural awareness for business professionals, and for general education in inter-cultural communication. Our novel contributions in this work are in the design, development and initial evaluation of immersive virtual humans in interactive virtual environments to teach and train users in inter-personal social verbal and non-verbal behaviors in a foreign culture.

## II. RELATED WORK

### 2.1 Virtual Humans in Pedagogy

Researchers have shown that a virtual human interface can provide feedback to human users using multiple verbal and non-verbal channels such as speech, gestures, and facial expressions. Anecdotal evidence suggests that human communication consists of a high bandwidth of modalities such as gestures, facial expressions, speech, and body language [5]. In addition, researchers have found that users can learn a task from demonstrations, far more effectively than learning to perform a task from text-based instructions alone, especially when that task involves spatial motor skills [13].

Many virtual human interfaces have been developed for training, pedagogy, and education. Thorisson et al, presented Gandalf, a communicative humanoid to guide planetary exploration [21]. Gandalf's behavior rules for face-to-face conduct are derived from psychology literature on human-human interaction. Rea, built by Cassell et al. is a virtual real estate agent capable of understanding speech and gaze, and is capable of planning multimodal utterances from propositional abstract representations [5]. Rea also keeps a model of interpersonal distance with the user and uses small talk to reduce interpersonal distance if she notices a lack of closeness with the users. Slater et al. examined the extent to which virtual humans could be used by actors and a director to rehearse for a live performance. The authors suggest that a performance level was reached which led to a successful live performance [17]. The Mission Rehearsal Exercise (MRE) system is an immersive virtual reality system with life-size virtual humans that was created to teach users leadership skills in task oriented social situations using fictional scenarios [8]. The MRE uses fictional scenarios in virtual environments in communicative training. Immersive virtual humans are also being used to allow medical students to experience the interaction between a patient and a medical doctor using natural methods of interaction with a high level of immersion [11]. Babu et al. have shown that an immersive virtual human physiotherapist can be used in training users in rehabilitation exercises by engaging users in natural multi-modal communication [3].

Each of these virtual human interfaces includes characteristics that are important for both perception of user behavior and conveying information. Nass has suggested that any interface that ignores a user's emotional state or fails to manifest the appropriate emotion and social behavior can dramatically impede performance and risk being perceived as cold, socially inept, untrustworthy, and incompetent [16]. Using both speech

and gestures also contributes to making virtual human interfaces more lifelike and believable [5]. Thorisson and Cassell also point out that non-verbal behaviors are important in supporting conversation, e.g. gaze gives cues for turn taking, a nod conveys understanding, and propositional hand gestures and facial expressions can direct the user's attention [22].

### 2.2 Social Effects of Virtual Humans

Researchers have investigated how people respond to computers and virtual humans. Nass and Moon have shown that people react to and attribute very human characteristics to computers, such as the computer's helpfulness, expertise, and friendliness [16]. Using a virtual human interface minimizes the need for training users, since they already know how to interact with other people [22]. Zanbaka et al., found that people respond to virtual humans similarly to the way they respond to real humans. The authors were able to elicit social inhibition from female participants in response to a virtual human observer [24]. Mel Slater's group at UCL has conducted studies of the social ramifications of having avatars in virtual environments. They were able to elicit emotions such as embarrassment, irritation, and self-awareness in virtual meetings. They found that the presence of avatars was important for social interaction, task performance, and presence [19]. Raij et. al. examined perceived similarities and differences in experiencing an interpersonal scenario with a real and virtual patient [18]. They found lower ratings on participants' rapport and conversational flow with the virtual patient was attributed to the limited expressiveness of the virtual patient. Level of immersion and natural interaction also facilitated the participants' ability to perform a training task with a virtual patient as effectively as with a real patient.

### 2.3 Virtual Humans in Inter-Cultural Communication Education

We have found little work that directly focuses on using virtual humans in training users in performing social verbal and non-verbal behaviors in a foreign culture. The research that is closest to ours is the Virtual Environment for Operational Readiness (VECTOR) [6] and ELECT BiLAT [9]. VECTOR trains solders in the critical communication skills for survival and mission success. The effects of the trainee's positive and negative conversational input with indigenous virtual Iraqi civilians result in behavioral outcomes (such as hostility, helpfulness, and aggression) in virtual Iraqi civilians based on the cultural expectations or norms. ELECT BiLAT is a PC based virtual environment for US army students to practice their skills in conducting meetings and negotiations in a specific cultural context. The trainee has to engage in bi-lateral meetings with local leaders to achieve their mission objectives. The system also features an intelligent coach to teach and train students in the appropriate cultural communicative behaviors when engaging in negotiations.

Our novel contribution is on training users in performing the verbal and non-verbal (gestures and gaze) social conversational behaviors in a foreign culture, using immersive virtual reality technology with life size virtual humans and multi-modal interaction. The existing training systems described in the literature are focused on training users in social communication skills that are primarily verbal, or in decision making regarding

the use of verbal and non-verbal behaviors during negotiations to achieve higher level goals. In contrast, the novelty of our work primarily focuses on learning how to perform non-verbal behaviors such as gestures, gaze, and facial expressions in addition to verbal phrases during discourse in a foreign culture. We specifically explore a subset of social conversational protocols pertaining to conversational initiation and disengagement pertaining to south Indian culture, and study the pedagogy and training of the co-occurrence of timing, synchronicity, and intensity of the verbal and non-verbal cultural elements of conversation. Our innovation allows users to interactively learn with an intelligent tutor and practice acting out the verbal and non-verbal behaviors with a conversation partner in a foreign culture in a natural manner.



Fig. 1. Shows a participant (right) greeting Anita, a virtual south Indian of the same gender (left), while Radha, the virtual south Indian instructor (middle) observes and shows the trainee how the greeting is performed (Color Plate 5).

## III. VIRTUAL REALITY SETUP

### 3.1 Immersive Virtual Reality Conversational Protocol System

The immersive virtual reality conversational protocol training system was housed in an office where participants could experience training with the virtual characters with no one else present (Fig. 1).

Our system used two networked PCs. One computer perfomed speech and gesture recognition, while the second handled the visual rendering component of the system. A data projector was used to display the virtual humans at life-size. The rendering PC was an Alienware Aurora with dual nVidia 7900 GT SLI graphics cards. The virtual humans were rendered at 40-45 FPS. Fig. 2 shows a schematic of the hardware and software infrastructure of our system.

The virtual humans were projected on a large screen in front of the participant. The dimensions of the projected image measured 1.8 meters in width by 1.35 meters in height. The participant stood 2 meters from the projection screen. The trainee's head, hands, and waist were tracked in 6 Degrees of Freedom (Position and Orientation) using a Polhemus Fastrack electro-magnetic tracker. This tracking data served as input for non-verbal gesture and gaze recognition (Fig. 2). Additionally, the head tracker data was also used in rendering the image according to the correct perspective warping effects. Speech

input was taken through a microphone attached to a head band. Audio output of the system was provided by speakers positioned on either side of the screen. Participants were requested to stand on a marked spot on the floor facing the screen when training with the virtual humans. A student-apprentice pedagogical model was chosen as the basis of the instruction, similar to Johnson et al. [11]. The virtual south Indian rendered on the right (Fig. 1), who is also the same gender as the participant, acts as the virtual instructor. The virtual south Indian rendered on the left (Fig. 1), acts as a conversational partner for the user, with whom the user practices the conversational protocols while observed by the instructor

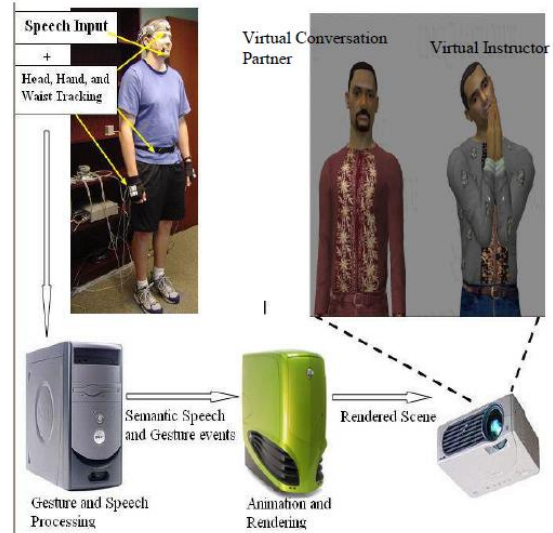### 3.2 Virtual Human Software Framework



Fig. 2. System Hardware Infrastructure

Our system was built using our Virtual Human Interface Framework (VHIF), as described in [1]. VHIF integrates best-existing, widely available components and agent technologies to ensure high quality graphics, speech recognition and generation, animation, and virtual human representation. The speech recognition module was built using Microsoft SAPI 5.1 and Dragon Speech Recognition engine 9.0. Interactive 3D characters from Haptek Inc. were used to create and animate the virtual humans, and Simple Virtual Environments framework (SVE) [14] was used to process motion tracking input and renders the graphics. To ensure high visual fidelity the virtual humans were carefully modeled using pictures of real south Indians using tools provided by Haptek Corp. for virtual human authoring. Speech utterances of the virtual humans were implemented using AT&T Natural Voices text-to-speech. Speech utterances were tailored to provide appropriate intonation and pitch. Verbal and non-verbal behaviors of the virtual humans including timing, and synchronization of gestures, body movements, postures, facial expressions based on emotion, and speech utterances were implemented as a finite state machine of behavior states in VHIF. VHIF also featured a cognitive model, which consisted of pre-scripted tree of behavior states towards instruction, and a discourse planner which selected a set of behavioral actions to be dynamically executed during interactive feedback. Implementation details of the underlying Virtual

Human Interface Framework are given in [1].

### 3.3 Virtual Human Instruction and Interactive Feedback

#### 3.3.1 Pedagogical Instruction

The immersive virtual human conversational protocol training system presented two life-size virtual humans to the user. Female participants were presented with a female virtual instructor (Radha), depicted on the right of the screen. Male participants were presented with a male virtual instructor (Sameer). The content of the instructions came from literary sources describing south Indian social customs, norms, and etiquette [20, 10].
Four inter-personal cultural conversational tasks were chosen:

1. Greeting someone of the same gender
2. Saying goodbye to someone of the same gender
3. Greeting someone of the opposite gender
4. Saying goodbye to someone of the opposite gender

The sequence and presentation of each protocol instruction were carefully designed to ensure consistent descriptions of each step. Each protocol task description consisted of a sequence of five steps. The instructions provided by the virtual south Indian instructor in virtual condition (VR) were also consistent with the instructions in the handout provided in the literature condition (L). In condition VR, the virtual instructor also demonstrated each step of the protocol to the participant using animated gestures and facial expressions. Images of the virtual instructor's non-verbal behaviors were also included along with the instructions in the literature provided to participants in condition L (Table 1). Each participant learned a total of 20 steps (4 tasks x 5 steps) consisting of verbal and non-verbal behaviors.

A sample protocol excerpt of verbal and non-verbal instructions for a female participant greeting a person of the opposite gender, taken from the instruction manual provided to female participants in condition L, is shown in Table 1.

#### 3.3.2 Virtual Human Response and Feedback

For every cultural conversational task, the virtual instructor first describes the steps involved in the task using a combination of verbal speech instructions as well as non-verbal gestures and facial expressions in a step-by-step manner. A male or female virtual partner either enters the scene or leaves depending on the nature of the social conversational protocol being currently simulated. Every participant in the study learned to greet and say goodbye to someone of the same gender and opposite gender. The order of instruction was randomized based on gender of the virtual interlocutor. Hence, each participated either learned how to greet and say goodbye to a person of the same gender first or learned how to greet and say goodbye to a person of the opposite gender first. For each cultural conversational protocol, the virtual instructor provides instruction, demonstration, and interactive feedback in the following steps to the participant/trainee.

Step 1: The virtual instructor provides verbal and non-verbal instructions similar to the list of instructions detailed in the instructional handout provided for participants in condition L.

Step 2: After providing the instructions, the virtual instructor

TABLE 1: SAMPLE CONVERSATIONAL PROTOCOL FOR THE TASK OF GREETING A PERSON OF THE OPPOSITE GENDER, TAKEN FROM THE INSTRUCTION MANUAL PROVIDED FOR FEMALE PARTICIPANTS. EACH STEP IS DESCRIBED ON THE LEFT, AND AN IMAGE OF THE VIRTUAL INSTRUCTOR PERFORMING THE STEP IS SHOWN ON THE RIGHT.

| Description | Illustration |
|---|---|
| 1. Stand facing the woman at two arms length. During the greeting the posture should be erect; both feet should be placed together. | |
| 2. Gaze at the other person's mouth. | |
| 3. Place both hands together in front of your neck with the tip of your fingers touching your chin. The elbows should be kept close to your body, and your arms should be tucked close to your chest. | |
| 4. Smile and say "namaste" (nam-as-TAY), while bowing your head slightly forward, and shifting your gaze to the other person's feet. And then bring your head back to the normal position. | |
| 5. Wait with your hands together, until the other person has completed greeting you similarly. Then bring your hands down, and relax your posture. | |

demonstrates to the participant by either greeting or saying goodbye to the virtual conversation partner standing next to him/her (Fig. 3).

Step 3: After providing instructions and showing how it is performed with the aid of the virtual conversation partner, the virtual instructor then asks the participant to initiate the conversational protocol to the virtual partner. At this time the virtual instructor observes the participant carrying out the conversational protocol with the virtual conversational partner.

Step 4: After a certain delay the virtual conversation partner responds similarly, and then the virtual instructor provides feedback on what the participant had done wrong, and reiterates the step of the conversational protocol the participant had performed incorrectly. If the participant had carried out the protocol correctly, then the virtual instructor commends the trainee on his/her performance.

Step 5: Next, steps 2 through 4 are repeated two more iterations. The only exception is that on the second iteration, in step 3 instead of initiating the conversational protocol, the trainee is asked to respond to the social conversational protocol from the virtual conversation partner.

Step 6: After learning how to greet and say goodbye to either

a person of the same gender or a person of the opposite gender. The virtual conversation partner corresponding to the gender learned leaves the scene and a male or a female virtual conversation partner corresponding to the gender not learned then enters the scene. Next steps 1 through 5 are carried out in learning how to greet and say goodbye interactively with the new virtual conversation partner.

Hence, for each of the four cultural conversational tasks (greeting a person of the same gender, saying goodbye to a person of the same gender, greeting a person of the opposite gender, and saying goodbye to a person of the opposite gender), the participant receives instruction, a demonstration, and interactive practice and feedback iteratively three times.

### 3.3.3 Gesture and Gaze Processing

One of the novel aspects of our system was the ability of users to practice the verbal and non-verbal conversational conventions in a foreign culture in a natural manner in simulated inter-cultural situations with a virtual instructor and virtual conversational partner. The virtual instructor upon observing the user's performance of the cultural conversational protocols with the virtual conversational partner would then provide feedback to the trainee. In order to provide this personalized feedback of the user's practice performance of the inter-personal cultural protocols, our system performed a method of automatic gesture and gaze processing in a manner specified below.

The participant's head, hands, and waist were tracked in 6 Degrees of Freedom (Position and Orientation) using a Polhemus Fastrack electro-magnetic tracker. The head-tracked data were used to control the gaze direction of the virtual humans. Gesture events were represented as a finite state machine



Fig. 3. Screenshot of Sameer on the right, the virtual instructor for male participants, demonstrating how to greet someone of the opposite gender in south India. Radha, the virtual conversation partner, is shown on the left.

of behavior states, as described in [1]. Each behavior state encapsulates the timing and sequence of non-verbal behaviors such as gestures, pose, and gaze pertaining to each step of the conversation task. Using the state machine of behavior states as well as the tracking data of the user's head, hands, and waist, the following types of non-verbal behaviors were evaluated:
- Did the participant bring his/her hands together?
- What was the orientation of the hands when brought together? Were they positioned forwards, backwards, left, or right?

- Were the hands kept too high or low, relative to the participant's body?
- Did the participant maintain appropriate gaze with the virtual conversational partner?
- Did the participant pitch his/her head forward or tilt his/her head to one side at appropriate times?
- Did the verbal greetings and goodbye occur in sync with the other non-verbal behaviors?
- Was the timing of the verbal and non-verbal behaviors correctly executed?
- Did the participant perform the verbal and non-verbal behaviors in the right sequence?

Participant head gaze was detected by finding the intersection of a ray from the tracked head to a virtual sphere representing the head of the virtual conversation partner. Hand clasping gestures were detected by the surface intersections of tracked virtual objects that corresponded to the position and orientation of the participant's real hands.

## IV.    EVALUATION

An initial study was conducted between two conditions to determine the effectiveness of using immersive virtual humans in teaching users verbal and non-verbal cultural conversational protocols, as compared to an existing method such as learning the protocols from a study guide with illustrations (a non-technological, traditional learning approach). Participants were randomly assigned to one of two conditions:

**Condition L:** Participants were provided an 8-page study guide with illustrations (Table 1).

**Condition VR:** Participants received instructions from a virtual south Indian instructor and interactively practiced with a virtual conversation partner.

### 4.1 Measures

The following measures were used to compare the performance of participants in condition L to condition VR.

### 4.1.1 Performance of Cultural Protocols

Participants were tested immediately after completing either the virtual reality (VR) or text-based (L) instruction. During testing, each participant was told that they would be presented with videos of real south Indians on the screen, and were instructed to carry out the appropriate protocol (either greeting or saying goodbye). In every testing scenario the participant always initiated the greeting or the goodbye. They were also told that after a certain delay, the person presented on the screen would respond appropriately. The participant's greetings or goodbyes were recorded by video camera and the video recordings were used to score how well the participants performed the south Indian cultural conversational conventions.

***Subjective evaluation*** was performed by three south Indian raters who were blind to each participant's condition. They viewed the digital videos and scored each participant's performance. The raters were asked to use the training instructions to evaluate each participant's performance. The raters assigned a score between 1 and 7 for each step in the protocol (1 = not at all, 7 = perfectly), to evaluate each step's correctness and

proper order.

***Objective evaluation*** of participants' performance was made using the tracker log data recorded for each participant during the testing session. Each participant was scored on how well he/she performed each step of the south Indian conversational protocol using the same motion analysis technique mentioned in section 3.3.3. The system scored each participant on a scale of 0 to 10 for each task performed by the participant during the testing session (0 = not at all, 10 = perfectly).

Subjective raters were given a smaller rating scale (7 point scale) to reduce the effects of rater bias. A 10 point scale was used for the automated objective evaluation since the score for each participant could be computed accurately based on our gesture processing technique. The scores from the objective and subjective evaluations were later normalized for comparative analysis.

### 4.1.2 Positive and Negative Affect

Learning takes place in a social atmosphere [Watson et al. 1988]. We believe that learning with immersive virtual characters should provide a fun, interesting, and social atmosphere for users to learn the cultural conversational protocols as compared to learning the cultural conventions0 from text. Conversely, if user's experiences of the learning environment were stressful, anxious, or boring then we expect negative affect to be increased. Our hypotheses:

（1） Positive affect measures should be greater in participants who learn the cultural protocols from immersive virtual characters as compared to participants who learn the proto cols from the written study guide.

（2） Negative affect measures should be less in participants who learn the cultural protocols from immersive virtual characters as compared to participants who learn the protocols from a written study guide.

Participants' positive and negative affect were measured prior to the training session and also immediately after the training session using the Watson, Clark, and Tellegan Positive and Negative Affect Test [23]. The test consisted of 10 positive affect questions and 10 negative affect questions were measured on a Likert scale (1 = very slightly or not at all, 5 = extremely).

### 4.1.3 Co-Presence Questionnaire

Co-presence of the immersive virtual humans was measured using the 14 question Slater Co-Presence Questionnaire in-condition VR only [15]. Responses were on a scale of 1 to 7 (1 = not at all, 7 = a great deal).

### 4.1.4 System Usability Questionnaire

The System Usability Scale questionnaire, developed by Digital Equipment Corporation [4], consists of 10 questions on a Likert scale (1 = strongly disagree, 5 = strongly agree) and was administered to participants in condition VR only. This section of the questionnaire included questions such as "I think I would like to use the Virtual Human Training System frequently." The questions came from four categories: Satisfaction, Simplicity, System Design, and Learn-ability. Satisfaction measured the extent to which the participant enjoyed working
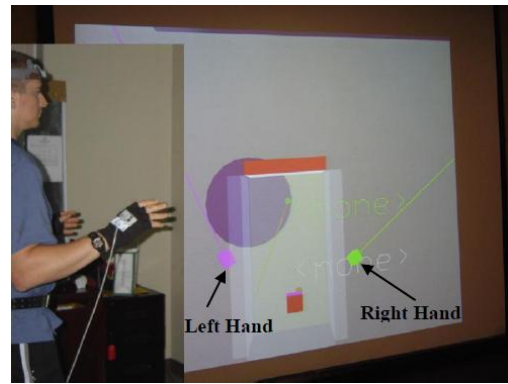


Fig. 4. Shows a participant performing a short training exercise to help the participant become accustomed to moving while wearing the motion tracking equipment. The green cube on the screen represents his right hand; the purple cube represents his left hand. (Color Plate 6 )



Fig. 5. Shows a participant in the testing session, greeting a person of the opposite gender presented on the screen.

with the virtual human social conversational protocol training system.

### 4.1.5 Virtual Humans Questionnaire

This survey was created to evaluate which aspects of the natural interaction and which characteristics of the immersive virtual humans helped or hindered the learning process. It was administered to participants in condition VR only. The set of questions addressed in this questionnaire were scored on a 7-point Likert scale (1 = not at all, 7 = a great deal). This section included questions such as "How clear were the virtual human's responses?" and "How much did you enjoy working with the virtual conversational partner?"

### 4.1.6 Post-Experiment Qualitative Questionnaire

The post-experiment questionnaire consisted of 10 questions about each participant's experience. The questionnaire was administered to participants in both conditions. This section included questions such as, "Did you feel you would have performed better, worse, or the same with another method of practice? If so, what?" The questions allowed us to evaluate the qualitative aspects of the participant's mode of training (L vs. VR) and the participant's subjective experiences.

## 4.2 Participant Information

Participants were randomly assigned to either condition VR or condition L. 40 participants completed the study, with 20 participants in each condition. Participants were recruited by classroom announcements to summer school students, responses to fliers, and by word-of-mouth. The average age of the participants was 27 [18-56]. Participants were required to be of non-Asian culture and able to communicate comfortably in English. Out of a total of 40 participants, 18 participants were female, and 22 participants were male.

## 4.3 Experiment Procedure

The pre-experiment session, training session, testing session, and post-experiment session took each participant approximately one hour to complete.

### 4.3.1 Pre-Experiment Session

The participant in each condition first read the Participant Informed Consent sheet and was asked if he/she had any questions. The participant then read and signed the Informed Conse

### 4.3.2 Training Session

Next, participants in both conditions were asked to take the positive and negative affect (PANAS) pre-training questionnaire. The participants were then trained in cultural conversational protocols based on their assigned condition.

**Condition VR:**

Participants were fitted with the electro-magnetic tracking equipment for head, hands, and waist. The trackers were affixed to a pair of gloves, a head band, and a belt. Each participant then trained the speech recognition program by reading from a short passage. A short training exercise was then performed to help the participant become accustomed to moving while wearing the motion tracking equipment (Fig. 4).

The participant was then told that he/she would now be trained in south Indian cultural protocols by virtual humans and that the training session would last for approximately 20 minutes. During the training session, the participant was left alone in the experiment room with the virtual reality training system. At the conclusion of the training session, the participant was administered the post-training PANAS questionnaire, and the Co-Presence questionnaire.

**Condition L:**

Participants were given a written study guide with illustrations of the cultural conversational protocols, and were told that they had a maximum of 20 minutes to study the material provided. During this time, the participant was left alone in the experiment room. If the participant felt that he/she was ready with less than 20 minutes of study, he/she could knock on the door of the experiment room to alert the experimenter, who then initiated the next part of the experiment. At the conclusion of the training session, the participant was administered the post-training PANAS questionnaire.

### 4.3.3 Testing Session

Participants were tested on how well they learned the cultural conversational protocols using videos of real south Indians as described in section 3.1.1 (Fig. 5).

### 4.3.4 Post-Experiment Session

**Condition VR:**

Participants filled out a System Usability Questionnaire, followed by the Virtual Humans Evaluation Questionnaire and the Qualitative Questionnaire.

**Condition L:**

Participants filled out the Qualitative Questionnaire.

Finally, participants in both conditions were orally debriefed.

## V. RESULTS AND ANALYSIS

### 5.1 How well did participants in condition VR learn the cultural conversational protocols as compared to participants in condition L?

#### 5.1.1 Subjective Evaluation

The video recordings were rated by three independent evaluators who were blind to the subjects' condition. The ratings from the three subjective evaluators for each participant nt Form. were summed and translated to provide a score on a scale from 0 to 96. Participants in condition VR ($M = 91.97$, $SD = 2.41$) scored higher than participants in condition L ($M = 84.90$, $SD = 4.79$). A t-test was used to compare the differences in video evaluation scores between conditions L and VR (Table 2). Levene's test for equality of variances was significant, $F = 5.04$, $p = .031$, and so results were generated without assuming homogenous variances. There was a significant difference between the two groups, $t(28.02) = 5.90$, $p < .001$, indicating that the participants who were instructed in condition VR were able to learn the relevant cultural conventions better than those in condition L. Additionally, there was less variation in scores for those in condition L. This experimental design provided an estimated power of .46 to detect medium-size effects.

TABLE 2: SHOWS THE DESCRIPTIVE STATISTICS FOR PARTICIPANTS IN CONDITIONS L AND VR ON PERFORMANCE IN THE TESTING SESSION.

| Conditions | N | Mean | Std. Dev. | StD. Error |
|---|---|---|---|---|
| L | 20 | 84.90 | 4.79 | 1.070 |
| VR | 20 | 91.97 | 2.41 | 0.538 |

$$t(28.02) = 5.90, p < 0.001$$

An assessment of *inter-reviewer difference* was performed using a one-way Analysis of Variance (ANOVA). The scores of participants across three reviewers were evaluated. An ANOVA revealed a significant effect due to rater on the mean score on the ratings of each reviewer $F(2, 37) = 5.82$, $p = 0.004$. Multiple group comparison tests, using the Tukey HSD test ($p = 0.05$), indicated that the mean scores from reviewer 1 ($M = 105.25$, $SD = 4.93$) and reviewer 3 ($M = 105.92$, $SD = 5.18$) were not significantly different. However, reviewer 2 ($M = 102.12$, $SD = 5.78$) scores were significantly lower than reviewer 1 ($p = 0.26$) and also significantly lower than reviewer 3 ($p = 0.005$). This indicates, that reviewer 2 was significantly stricter than reviewer 1 and 3. The relationships among the three reviewer's scores were assessed using Pearson's Correlation coefficients. There was a significant positive relationship
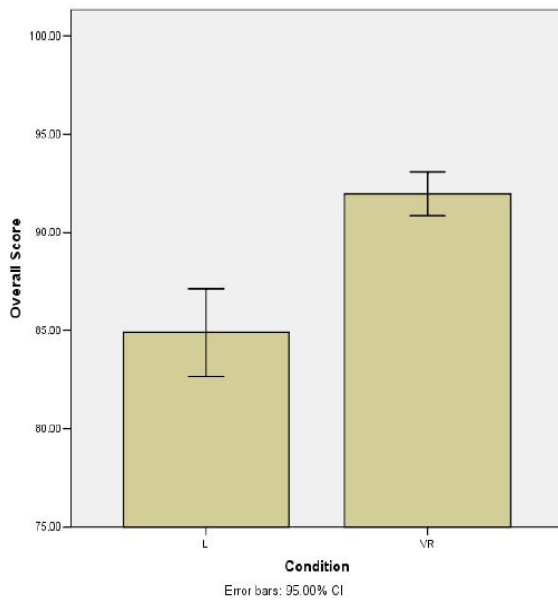
Fig. 6. Bar graph showing the mean score and standard deviations of the evaluations based on video for conditions L and VR.

between the scores of reviewer 1 and reviewer 2, r(4) = +0.92, *p* < 0.001, a significant positive relationship between the scores of reviewer 1 and reviewer 3, r(4) = +0.93, *p* <0.001, and a significant positive relationship between the scores of reviewer 2 and reviewer 3 scores, r(4) = +0.93, *p* < 0.001. Hence, Pearson's correlation analysis revealed that the scores of all three reviewers were significantly positively correlated and consistent.

The average scores for both conditions L and VR were high. Participants learned the cultural conversational protocols with either type of training (Fig. 6). However, participants in condition VR performed significantly better than participants in condition L. An important difference between the two conditions was that the variance in test performance scores for the VR condition was four times lower than the variance for the L condition. This suggests that training with virtual humans provides a more consistent and reliable result.

### 5.1.2 Objective Evaluation

The objective evaluation was performed on the tracker log data collected during the testing session. Due to experimenter error, data from only 34 participants were recovered; data from 6 participants were lost (3 participants in condition VR and 3 participants in condition L). The log data was scored by the system using the same criteria used in interactive feedback for participants in condition VR. Scores for each task were summed to provide a total score on a scale from 0 to 50. Participants in condition VR (*M* = 41.97, *SD* = 3.34) scored higher than participants in condition L (*M* = 39.12, *SD* = 2.46). A t-test was used to compare the differences in objective evaluation scores between conditions L and VR. The differences in scores between the two groups was not significant, t(18.16) = 2.07, *p* = 0.066, however there seemed to be a strong trend with participants in condition VR performing better than participants in condition L.

### 5.1.3 Correlation between Subjective and Objective scores

The relationship between the two scores measured on the performance of each participant (Objective Evaluation Score and Subjective Evaluation Score) was assessed using Pearson Correlation Coefficients. The Correlation analysis was performed on Objective and Subjective Evaluation Scores of 34 participants. There was a significant positive relationship between Objective Evaluation Scores (*M* = 40.21, *SD* = 3.48) and Subjective Evaluation Scores (*M* = 87.64, *SD* = 4.86), r(2) = +0.694, p=0.021, indicating higher Objective Evaluation Scoreswere associated with higher Subjective Evaluation Scores for each participant. In summary the Pearson Correlation analysis revealed that the Objective measures, obtained from automated analysis of logs, were indeed a reliable metric for measuring the participant's performance of the social protocols in the testing session.

TABLE 3: SHOWS DESCRIPTIVE STATISTICS FOR POSITIVE AND NEGATIVE AFFECT SCORES FOR PRE- AND POST- TRAINING IN CONDITIONS L AND VR.

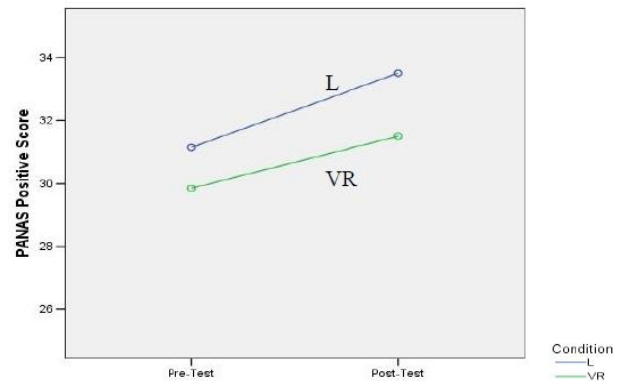| | Condition | Mean | Std. Deviation | N |
|---|---|---|---|---|
| PANAS Positive Pretest | L | 31.15 | 8.683 | 20 |
| | VR | 29.85 | 6.385 | 20 |
| | Total | 30.50 | 7.552 | 40 |
| PANAS Positive Posttest | L | 33.50 | 9.000 | 20 |
| | VR | 31.50 | 7.810 | 20 |
| | Total | 32.50 | 8.379 | 40 |
| PANAS Negative Pretest | L | 13.05 | 3.486 | 20 |
| | VR | 12.00 | 2.772 | 20 |
| | Total | 12.53 | 3.154 | 40 |
| PANAS Negative Posttest | L | 13.80 | 3.708 | 20 |
| | VR | 11.35 | 1.814 | 20 |
| | Total | 12.58 | 3.137 | 40 |



Fig. 7. Graph showing the trends for PANAS Positive score for Conditions L (Blue Line) and VR (Green Line).

### 5.2 Was there a difference in affect of participants learning cultural conversational protocols with immersive virtual humans as compared to learning from a study guide?

To interpret influence on positive and negative affect scores (PANAS), two 2x2 analyses of variance were performed, testing the within-subjects effect of the respective PANAS score before and after instruction (ranging from 10 to 50, with higher numbers corresponding to greater affect) and the between-subjects effect of type of instruction (study guide vs.

immersive virtual humans). This experimental design provided an estimated power of 0.67 to detect medium-size effects. Table 3 summarizes the results of PANAS positive and negative scores. Fig.7 shows the interaction plot for PANAS positive score, and Fig. 8 shows the interaction plot for PANAS negative score.

Analysis of PANAS positive score revealed a non-significant interaction between the score and type of instruction, $F(1,38) = .28$, $p = .603$, partial $eta_2 = .007$, and a significant main effect for time (pre- vs. post- instruction) of PANAS score, $F(1,38) = 9.01$, $p = .005$, partial $eta_2 = 0.192$. The main effect for instruction type was not significant, $F(1,38) = .453$, $p = .505$, partial $eta_2 = .012$. These results indicate that participants who were instructed by virtual humans reported similar positive affect to those who were instructed using the study guide. Overall, positive affect increased from before instruction ($M = 30.50$, $SD = 7.55$) to after instruction ($M = 32.50$, $SD = 8.38$), though the effect size was small.

Analysis of PANAS negative score revealed a non-significant interaction between score and type of instruction, $F(1,38) = 3.91$, $p = .055$, partial $eta_2 = .093$, and a non-significant main effect for time of PANAS score, $F(1,38) = 0.02$, $p = .888$, partial $eta_2 = .001$. The main effect for instruction type was not significant, $F(1,38) = 3.84$, $p = .057$, partial $eta_2 = .092$, indicating that participants who were instructed by virtual humans reported similar negative affect to those learning from a study guide. Additionally, overall negative affect did not change from before instruction ($M = 12.53$, $SD = 3.14$) to after instruction ($M = 12.58$, $SD = 3.14$).
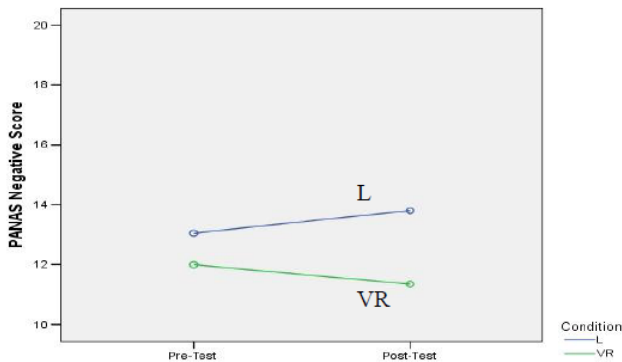


Fig. 8. Graph showing the trends for PANAS Negative score for Conditions L (Blue Line) and VR (Green Line).

In discussion, positive affect corresponds to positive emotions such as excitement and interest, and negative affect corresponds to negative emotions such as stress and anxiety. Overall, results suggests that participants in condition VR and condition L experienced increase in positive affect, and a decrease in negative affect, between pre- and post- instruction. There was a close to significant ($p = 0.057$) decrease of negative affect in participants in condition VR as compared to participants in condition L.

### 5.3 Co-Presence

Co-presence for participants in condition VR was assessed using the questionnaire proposed by Mortensen et al in 2002 [19]. Co-presence measured the extent to which participants felt that they had been interacting with another human being. The questionnaire contained 6 questions pertinent to co-presence, each rated from 1 (not at all) to 7 (a great deal) on a Likert-type scale. The co-presence score for each participant is a count of the number of high scores (responses of 5, 6, or 7) on these questions after adjusting direction so that a higher score is correlated with higher co-presence. The scores varied across the scale ($M = 3.16$, $SD = 1.19$) as there were participants who received a score of 0 and others who received the maximum score of 6. The raw mean was slightly higher ($M = 4.36$, $SD = 1.18$) suggesting that at some moderate level, users treated the virtual conversation partner and the virtual instructor as other humans as opposed to software agents.

### 5.4 System Usability Scale

System usability for participants in condition VR was evaluated using the System Usability Scale developed by the Digital Equipment Corporation [4]. Overall, the system exhibited a high degree of reported usability ($M = 81.75$, $SD = 11.15$) with a highest usability score of 97.5 and a lowest score of 60.0 on a scale of 0 to 100 (100 being the highest usability score).

### 5.5 Virtual Humans Evaluation Questionnaire

Overall, the results for the quantitative evaluation scores were quite high ($M = 5.41$, $SD = 0.663$). Users indicated that they thought the virtual human training system would be useful for training in cultural conversational protocols in a foreign culture ($M = 6.15$, $SD = 1.04$). Users also indicated that the virtual human training system provided instructions with clarity ($M = 6.45$, $SD = 0.61$). Users also suggested that the virtual human training system was also appropriate for the task of teaching and training in inter-personal cultural conversational conventions ($M = 6.21$, $SD = 1.08$).

Scores on other dimensions such as likelihood that the participant will use the system if made available, appearance of the virtual humans, intuitiveness of the interface, enjoyment of the interface, and appearance of the virtual humans all had mean scores greater than 5.0 on a scale of 1 to 7.

### 5.6 Post-Experiment Qualitative Evaluation

Participants' post-experiment comments, feedback, and suggestions for conditions L and VR are summarized as follows:

While some participants in the L condition felt that the training was adequate, many participants indicated that this method of training was not sufficient.

*"I think I would have performed better if I had a person to practice with. I tend to perform better with lots of practice and feedback."*

*"I probably did worse than if someone just showed me by example. But the guide was probably better than anything I'd find on the internet if I wanted to learn it on my own."*

Many participants in the VR condition felt that the training adequately prepared them for interaction using South Indian social conversational protocols.

*"The information seemed basic, but I feel I could greet someone of South Indian culture."*

However, many participants reported that they would require additional practice with a real person before they would feel confident in exercising the conversational protocols.

*"I think I would want to practice with real people before I*

*naturally used what I learned."*

Interestingly, some participants indicating that they preferred interacting with a virtual human, mostly due to the possibility of embarrassment, shyness, or social anxiety.

*"I think I would have performed worse if I was practicing with a real person. I would have been more shy and self-conscious."*

*"I felt like this was a way to learn cultural behavior without feeling embarrassed about mistakes as I might have made with an actual human trainer."*

Overall, users indicated that the virtual humans were very realistic, and that this facilitated their interaction.

*"They look realistic. It is easy and enjoyable to communicate with them."*

*"Throughout the experiment, I did not feel that they are VH. They seemed truly real."*

While the feedback on the visual fidelity of the virtual humans was generally positive, several participants indicated that deficiencies in the appearance or realism of the virtual humans hindered them from completing the task.

*"I found myself focusing on what was not life-like about the interaction."*

*"Movements are not yet fluid enough to allow the user to forget that they are not interacting with a human."*

Several participants also indicated that the delay in the virtual humans' response time was too long. Additionally, imprecision in the head-tracked data occasionally caused the virtual human to incorrectly report that the user failed to maintain eye gaze.

*"At the onset of the testing I was consistently told that I was not maintaining eye-contact when I believed very much that I was."*

Finally, almost all participants in the VR condition indicated that they would use a system for training social conversational etiquette if it were made available to them.

*"I think it could be very helpful for people who travel frequently and are concerned about interacting with people in the proper way."*

## VI. FUTURE WORK

Based on the results of this initial study, we have shown that immersive virtual humans in natural multi-modal interaction can be used as a tool for teaching cultural conversational protocols. However, in future studies we would like to investigate the following:

☐ (1) We would like to determine the effectiveness of instruction and demonstration with virtual humans alone vs. instruction and demonstration with interactive practice and feedback with virtual humans, in order to ascertain what aspect of our VR system contributes most effectively to learning the cultural conventions.

☐ (2) We would also like to investigate the impact of desk-top based VR systems such as ones with instruction and minimal interactive feedback, as well as scenario based training systems with virtual humans as a part of a narrative scenario in teaching cultural conversational customs.

☐ (3) We are also interested in using game platforms to investigate the effectiveness of context specific training in teaching cultural conventions in an interactive game-like scenario.

☐ (4) In future studies, we are also interested in comparing virtual human based teaching and demonstration vs. videos of real south Indians teaching and demonstrating the cultural conversational protocols.

## VII. CONCLUSION

In this initial study, we explored the applicability of immersive virtual humans in training and teaching users cultural conversational protocols in natural interaction, specifically in the verbal and non-verbal behaviors associated with conversation initiation and disengagement in south Indian culture. Our virtual human social conversational protocol training system provides instruction as well as interactive multi-modal practice and feedback from the virtual instructor and virtual conversation partner in real-time.

Results of our study suggest that participants who trained with and gained interactive feedback from the immersive virtual humans performed significantly better than participants who learned from the written study guide. The average scores for both conditions L and VR were high. Participants learned the cultural conversational protocols with either type of training. However, participants in condition VR performed significantly better than participants in condition L. An important difference between the two conditions was that the variance in test performance scores for the VR condition was four times lower than the variance for the L condition. This suggests that training with virtual humans provides a more consistent and reliable result. Our results also revealed that there was no significant difference in positive and negative affect between participants in conditions VR and L, although, overall, positive affect increased and negative affect decreased from before to after instruction in both conditions.

Our evaluations on co-presence of participants in condition VR suggest that to a moderate level, users treated the virtual instructor and the virtual conversation partner as another human as opposed to software components. Our quantitative evaluations suggest that participants in the VR condition generally enjoyed interacting with the virtual humans. Participants commented that the multi-modal interaction with the virtual humans to be intuitive, and the instructions provided by the virtual instructor to be clear. Participants also found the appearance of the virtual instructor and virtual conversation partner to be human-like. Users also suggested that the virtual human training system was appropriate for the task of inter-cultural conversational conventions training. Comments and suggestions from our participants indicate that a great number of participants in condition L found the written study guide not as useful as learning from example, practice, and feedback. Although movements of the virtual humans were sometimes not as fluid or the feedback sometimes not as accurate, almost all participants claimed that they found the system very helpful for users who travel frequently and are concerned about interacting with people according to appropriate cultural etiquette. Based on our participants' comments, we believe that our system can also be useful in training users in

cultural conversational protocols in other cultures as well.

## ACKNOWLEDGEMENT

## REFERENCES

[1]    S. Babu, S. Schmugge, R. K. Inugala, S. Rao, T. Barnes, L. F. Hodges, 2005. Marve: a prototype virtual human interface framework for studying human-virtual human interaction. *Proceedings of the 5th International Working Conference on Intelligent Virtual Agents (IVA 2005)*, Kos, Greece, pp. 120-133.

[2]    S. Babu, E. Suma, T. Barnes and L. F. Hodges, 2007. Can Immersive Virtual Humans Teach Social Conversational Protocols? *Proceedings of the IEEE International Conference on Virtual Reality 2007*, Charlotte, North Carolina, pp. 215 – 218.

[3]    S. Babu, C. Zanbaka, J. Jackson, T.-O. Chung, B. Lok, M. C. Shin, L. F. Hodges, 2005. Virtual Human Physiotherapist Framework for Persona-lized Training and Rehabilitation. *Proceedings of the International Conference on Graphics Interface 2005 (GI 2005)*, Victoria, British Co-lumbia, Canada, May 9 - 11.

[4]    J. Brooke, 1996. SUS: A quick and dirty usability scale. *Usability Eval-uation in Industry*, P. Jordan, B. Thomas, B. Weerdmeester, and I. McClelland, Eds. Taylor and Francis, London, pp. 189-194.

[5]    J. Cassell, 2000. Embodied Conversational Interface Agents. *Communi-cations of ACM*, vol. 43, pp. 70-78.

[6]    J. Deaton and C. Mccollum, 2004. Applying a cognitive architecture to control of virtual non-player characters. *Proceedings of the 2004 Winter Simulation Conference*, pp. 883-889.

[7]    P. Ekman and W. V. Friesen, 1969. The Repertoire of Nonverbal Beha-vior: Categories, Origins, Usage and Coding. *Proceedings of Semiotica*, vol. 1, pp. 49-97.

[8]    R. W. Hill Jr, J. Gratch, S. Marsella, R J. ickel, W. Swartout and D. Traum 2003. Virtual Humans in the Mission Rehersal Exercise System, *Kynstliche Intelligenz (KI Journal)*, vol. 17, 2003.

[9]    R. W. Hill, Jr, J. Belanich, C. L. Lane, M. Core, M. Dixon, E. Forbell, J. Kim  and J. Hart, 2006. Pedagogically Structured Game-Based Training: Development of The ELECT BiLAT Simulation, in Proceedings of the 25th Army Science Conference (2006).

[10]  F. E.Jandt, *An Introduction to Intercultural Communication: Identities in a Global Community*. Sage Publications, Thousand Oaks, California.

[11]  K. Johnson, R. Dickerson, A. Raij, B. Lok, J. Jackson, M. Shin, J. Her-nandez, A. Stevens and D. S. Lind, 2005. Experiences in Using Immer-sive Virtual Characters to Educate Medical Communication Skills. *Pro-ceedings of IEEE Virtual Reality 2005 (VR 2005)*, Bonn, Germany.

[12]  W. L. Johnson, J. W. Rickel and J. C. Lester, 2000. Animated Pedagog-ical Agents: Face-to-Face Interaction in Interactive Learning Environ-ments. *The International Journal of Artificial Intelligence in Education*, vol. 11, pp. 47-78, 2000.

[13]  W. L. Johnson and J. Rickel, 1998. Steve: An Animated Pedagogical Agent for Procedural Training in Virtual Environmnents. *SIGART Bulle-tin*, vol. 8, pp. 16-21.

[14]  G. Kessler, D. Bowman and L. F. Hodges, 2000. The Simple Virtual Environment Library: An Extensible Framework for Building VE Ap-plications. *Presence: Teleoperators and Virtual Environments*, Vol 9(2): p. 187-208.

[15]  J. Mortensen, V. Vinayagamourty, M. Slater, A. Steed, B. Lok and M. Whitton, 2002. Collaboration in Tele-Immersive Environments. *Pro-ceedings of the Eighth Eurographics Workshop on Virtual Environments*.

[16]  C. Nass and Y. Moon, 2000. Machines and Mindlessness: Social Res-ponses to Computers. *Journal of Social Issues*, vol. 56, pp. 81-103.

[17]  D. Pertaub, M. Slater and C. Baker 2001. An Experiment on Public Speaking Anxiety in Response to Three Different Types of Virtual Au-dience. *Presence: Teleoperators and Virtual Environments*, vol. 11, pp. 68-78.

[18]  A. Raij, K. Johnson, R. Dickerson, B. Lok, M. Cohen, A.Stevens, T. Bernard, C. Oxendine, P. Wagner and D. S. Lind, 2006. Interpersonal Scenarios: Virtual § Real? *Proceedings of IEEE Virtual Reality 2006*, Alexandria, VA.

[19]  M. Slater, M. Sadagic, M. Usoh and R. Schroeder, 2000. Small-Group Behavior in a Virtual and Real Environment: A Comparative Study. *Presence: Teleoperators and Virtual Environments*, vol. 9, pp. 37-51.

[20]  M. S. Thirumala, 2003. Communication via Gesture. *Language in India*, Strength for Today and Bright Hope for Tomorrow, vol 3.

[21]  K. R. Thorisson, Gandalf: An Embodied Humanoid Capable of Real-Time Multimodal Dialogue with People. *Proceedings of The First ACM International Conference on Autonomous Agents*, Marina del Rey, California, 1997.

[22]  K. R. Thorisson and J. Cassell, 1996. Why Put an Agent in a Body: The Importance of Communicative Feedback in Human-Humanoid Dialogue. *Proceedings of Lifelike Computer Characters '96*, Snowbird, Utah.

[23]   D. Watson, L. Clark and A. Tellegan, 1988. Development and Validation of brief measures of Positive and Negative Affect: The PANAS Scales. *Journal of Personality and Social Psychology*, Vol 54(6), pp. 1063-1070.

[24]  Z C. Anbaka, A. Ulinski, P. Goolkasian, L. F. Hodges, 2007. Social responses to virtual humans: Implications for future interface design. *Proceedings of CHI 2007*, ACM Press, 1561 – 1570.

[25]  C. Ziemkiewicz, A. Ulinski, C. Zanbaka, S. Hardin, and L. F. Hodges, 2005. Digital Patient for Triage Nurse Training. *Proceedings of HCI In-ternational 2005 (HCI 2005)*, Las Vegas, USA, 2005.

**Sabarish V. Babu** is an Assistant Professor in the Divi-sion of Human Centered Computing in the School of Computing at Clemson University.  He received his PhD in the Department of Computer Science at the University of North Carolina at Charlotte in 2007. Prior to joining Clemson University, he was a Post-Doctoral Fellow in the Department of Computer Science at the University of Iowa.  His research interests are in Virtual Environments, Applied Perception and Cognition in Virtual Reality, Virtual Humans, and 3D Human Computer Interaction.  For more information see his webpage at: http://people.clemson.edu/~sbabu

**Evan A. Suma** is a Postdoctoral Research Associate in the Institute for Creative Technologies at the University of California. He received his Ph.D. in 2010 from the University of North Carolina in Charlotte. His research interests include virtual environments, 3D user inter-faces, and human-computer interaction.

**Larry F. Hodges** is C. Tycho Howle Endowed Chair and Director of the School of Computing at Clemson University. His research interests include virtual envi-ronments, 3DUI, and human-virtual human interaction. In 2006 he received the *IEEE Virtual Reality Career Award* for his contributions to clinical applications of virtual reality.

**Tiffany Barnes** is an Associate Professor in the De-partment of Computer Science at The University of North Carolina at Charlotte.  She received her PhD in Computer Science from The North Carolina State Uni-versity in 2003.  Her research interests are in Algorithms, Computer Based Education, Knowledge Modeling, Data Mining, Bioinformatics and Intelligent Systems.