



GENESIS REDUX

Essays in the History
and Philosophy of Artificial Life

EDITED BY JESSICA RISKIN

Essays in the History
and Philosophy of Artificial Life

EDITED BY JESSICA RISKIN

GENESIS REDUX

THE UNIVERSITY OF

CHICAGO PRESS Chicago and London

Cassell, J. 2007. Body language: Lessons from the near-human. In *Genesis redux: Essays in this history and philosophy of artificial life*, ed. Jessica Riskin, 346-374. Chicago: University of Chicago Press.

17 Body Language: Lessons from the Near-Human JUSTINE CASSELL

The story of the automaton had struck deep root into their souls and, in fact, a pernicious mistrust of human figures in general had begun to creep in. Many lovers, to be quite convinced that they were not enamoured of wooden dolls, would request their mistresses to sing and dance a little out of time, to embroider and knit, and play with their lapdogs, while listening to reading, etc., and, above all, not merely to listen, but also sometimes to talk, in such a manner as presupposed actual thought and feeling. | E. T. A. Hoffmann, "The Sandman," 1844

INTRODUCTION

It's the summer of 2005, and I'm teaching a group of linguists in a small Edinburgh classroom. The lesson consists of watching intently the conversational skills of a life-size virtual human projected on the screen at the front of the room. Most of the participants come from formal linguistics; they are used to describing human language in terms of logical formulae and usually see language as an expression of a person's intentions to communicate that issues directly from that person's mouth. I, on the other hand, come from a tradition that sees language as a genre of social practice, or interpersonal action, situated in the space between two or several people, emergent and multiply determined by social, personal, historical, and moment-to-moment linguistic contexts, and I am as likely to see language expressed by a person's hands and eyes as by mouth and pen. As a graduate student pursuing a dual Ph.D. in linguistics and psychology in the 1980s, I had felt profoundly inadequate in the presence of these scholars: their formalized theories belong to a particular kind of technical discourse that is constructed in opposition to everyday language and that had seemed more scientific than my messy relational and embodied understanding of how language looks.¹ Those feelings of inadequacy—along with the experience of having articles rejected from mainstream journals and conferences—led me to try to formalize or "scientific" my work. I undertook a collaboration

with computer scientists in 1993 to build a computational simulation of my hypotheses that took the form of virtual humans who act on the basis of a "grammar" of rules about human communications. In turn, that simulation has, in the manner of all iconic representations, turned out to both reveal and obscure my original goals, depending on what the technical features of the model can and cannot handle. And the simulation has, like many scientific instruments, taken on a life of its own—almost literally in this instance—as the virtual human has come to be a playmate for children, a teaching device for soldiers, and a companion on cell phones—a mode of interacting with computers as well as a simulation that runs on computers.

But back to the classroom in Edinburgh. In the intervening fifteen years since graduate school, I have armed myself with a "sexy" demo to show other scientists, and times have changed so that the notion that language is embodied is somewhat more accepted in linguistics today. And so these formal linguists have chosen to attend the summer school class on "face-to-face pragmatics" that I am co-teaching. In the conversation today I'm trying to convince them of two points: that linguists should study videotapes and not just audiotapes, and that we can learn something important about human language by studying embodied conversational agents—fake humans who are capable of carrying on a (very limited) conversation with real humans—such as the one we call NUMACK, shown in figure 17.1.

I show the students a new video of NUMACK (the Northwestern University Multimodal Autonomous Conversational Kiosk) interacting with a real human, a simulation of our latest work on the relationship between gesture and language during direction-giving. On the basis of an examination of ten people giving directions to a particular place across campus, my students and I have tried to extract generalities at a fine enough level of detail to be able to understand what the humans are doing, and to use that understanding to program our virtual humans to give directions in the same way as humans do. The work exemplified in this particular video has concentrated on the shape of the people's hands as they give directions and on what kind of information they choose to give in speech and what kind in gesture. I'm excited to share this work, which has taken over a year to complete—moment-by-moment investigations into the minutiae of human gesture and language extracted from endless examinations of videotapes that show four views of a conversation (see figure 17.2), followed by complicated and novel implementations of a computer system that can behave in the same way. In fact, this is my own first view of the newly updated system; I've been traveling, and my graduate students finished up the programming and filmed the demo.

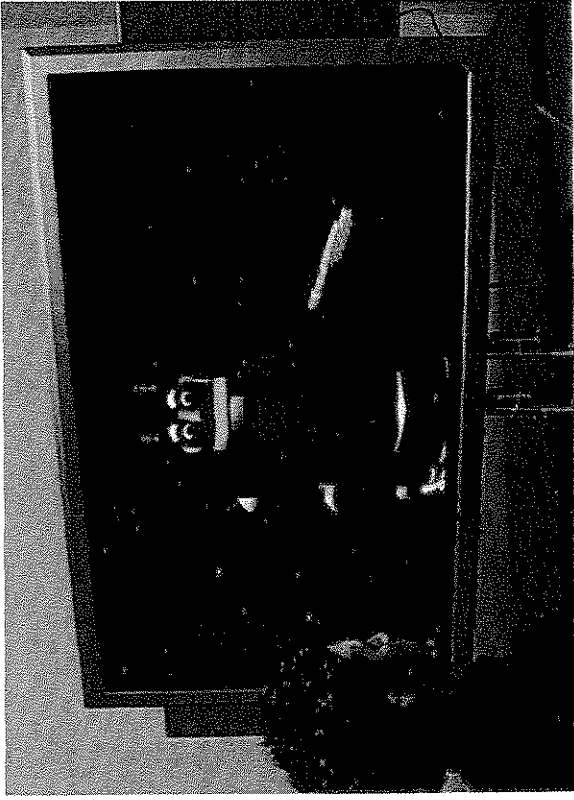


FIGURE 17.1 | NUMACK, the Northwestern University Multimodal Autonomous Conversational Kiosk, giving directions to a real human.

The Edinburgh linguists and I watch the video of NUMACK giving directions to a person, and it looks terrible! The small group of students tries to look down so as not to reveal that they don't think this is a fitting culmination of one year of work. I break the silence and say, "It looks ridiculous! Something is really off here. What is wrong with it? Can anybody help me figure out why it looks so nonhuman?" The students look surprised—after all, NUMACK looks nonhuman along hundreds of dimensions (starting with the fact that it is purple). Accustomed as they are to seeing impeccably animated characters in movies and on Web pages, they have expected to hurt my feelings by criticizing the virtual human's poor rendering of reality. But, as we watch the video over and over again, what stands out is how NUMACK's interaction violates our intuitions about how direction-giving should look. After three or four viewings, one student notes that NUMACK'S two hands operate independently in giving directions. The virtual human says, "Take a right," and gestures with his right hand. He then says, "Take a left," and gestures with his left hand. I've never thought about this before, but in looking at the robot I am struck by the fact that we humans don't do that—we must



FIGURE 17.2 | Analysis of videotapes allows us to draw generalizations about human-to-human direction-giving.

have some kind of cohesion in our gestures that makes us use the same hand for the same set of directions. Another student points out that the virtual human describes the entire route (in roughly fourteen "turn left," "turn right," "go straight ahead" kinds of segments) at once, with only an "uh-huh" on the part of the real human. No real human would do that—the directions are too long and couldn't possibly be remembered in their entirety.

I am thrilled and once again amazed at how much I learn about human behavior when I try to re-create it, especially when, and because, my imitations are partial and imperfect. Only when I try to reproduce the processes in the individual that go into making embodied language, do I get such a clear picture of what I *do* not yet know. For example, here I have realized that we will need to go back to our ten human direction-givers and look at their choice of hands—can I draw any generalizations about the contexts in which they use their right hand or their left? When is the same hand used repeatedly, and when do they switch to a different hand? Likewise, we will need to look further at the emergent properties of their directions. What behaviors signal to the direction-giver when to pause and when to continue,

when to elaborate and when to repeat? What embodied and verbal actions serve to alert the two participants to that the message has been taken up and understood, and that the next part of the message can be conveyed? I am also struck once again at the extent to which people are willing to engage with the virtual human, both as participants in a conversation about how to get to the campus chapel, and as participants in a conversation about the holes in our theory of the relationship between verbal and nonverbal elements in conversation.

I have learned something about the particularities of human communication here despite the fact that what I am viewing is a freak of artificial nature—a virtual human that is both generic and very particular, general and very detailed. In fact, for the experiment to work, we depend in part on the not-so-laudable schemas and expectations of our viewers and ourselves—that there can be such a thing as a generic human, which probably entails, for a direction-giving robot, that it is male and humanoid (albeit purple) and that its voice is Caucasian and American. As Clifford Nass and Scott Brave point out, violating cultural assumptions about expertise and gender or race produces distrust in users.² In the art world, Lynn Herschman Leeson, among others, has violated exactly these assumptions by synthesizing an infinitely smart female robot (or “bot”) whose body is present only in certain contexts, and who reproduces herself. But, in the current case, these largely unconscious assumptions on the part of the linguists examining the simulation are what allow them to identify as failings not a lack of personality or cultural identity in the virtual human, but simply that the hands are not synchronized. Thus I have learned something about human communication despite all of the ways in which this virtual human is not very human at all. I question these assumptions below, but for the moment let us return to the fundamental questions that guide this work.

Artificial Intelligence (AI) investigators and their acolytes, like the creators of automata before them, ask, “Can we make a mechanical human (or, in the weaker version, a human-like machine)?” I would rather ask “What can we learn about humans when we make a machine that evokes humanness in us—a machine that acts human enough that we respond to it as we respond to another human (where I mean both respond to us in our status of interlocutor or of scientist)?” Some researchers are interested in stretching the limits and capabilities of the machine or in stretching the limits of what we consider human by building increasingly human machines. Such is the case for the work described by Evelyn Keller in this volume. In my own work, at the end of the day I am less interested in the

properties of machines than in the properties of humans. For me there are two kinds of “aha!” moments: those in which I learn from my successes by watching a person turn to one of my virtual humans and unconsciously nod and carry on a conversation replete with gestures and intent eye gaze; and those in which I learn from my failures by watching the ways in which the real human is uncomfortable in the interaction, or the interaction looks wrong, as I illustrated in the Edinburgh classroom. These simulations serve as sufficiency proofs for partial theories of human behavior—what Keller has described as the second historical stage in the use of simulation and computer modeling³—and thus my goal is to build a virtual human to which people cannot help reacting as if it *were* actually human, to which people cannot prevent themselves from applying native-speaker intuitions. And key to the enterprise is the fact that those theories of human behavior and those native-speaker intuitions refer to the whole body, as it enacts conversations with other bodies in the physical world.

In the remainder of this chapter, I discuss my work on one particular kind of virtual human called an embodied conversational agent (ECA) in terms of its dual function as a simulation and as an interface. That is, I describe how these virtual humans have allowed me to test hypotheses about human conversation and what they have taught me by their flaws. I also describe the life that ECAs have acquired when they leave the lab—the uses to which companies and research labs have put them. In this way, I hope to illuminate the kinds of conversations that these virtual humans engage in when scientists use them as tools to study conversational phenomena and when ordinary people use them as tools to access information, dial phone numbers, learn languages, and so forth.

EMBODIED CONVERSATIONAL AGENTS AS CONVERSATIONAL SIMULATIONS

Just to be clear about our terms, *embodied conversational agents* are cartoonlike, often life-size, depictions of virtual humans that are projected on a screen. They have bodies that look more or less human, they are capable of initiating and responding in (very limited) conversations (in preset domains) with real humans, and they have agency in the sense that they behave autonomously, in the moment of their deployment, without anybody pulling the strings. Of course, this agency relies on a prior preset network of interactions among their inventors, their users, and the sociotechnical context of their deployment. As a point of contrast, consider chat bots or chatterbots.

Chat bots (such as the popular Alice, at <http://www.alicebot.org/>, which readers can try out for themselves) rely on a mixture of matching input sentences to templates, stock responses, and conversational tricks (such as “What makes you say X [where X is what the user typed in]?” or “I would need a more complicated algorithm to answer that question” when they don’t understand). Chat bots are increasingly employed by artists such as Lynn Hershman Leeson, STELAR, or Kirsten Geisler because they are relatively easy to program and thus allow the artist to concentrate on the aesthetic experience she or he wishes to provoke in the viewer. Chat bots often communicate with viewers only through text, but when embodied, they usually have only a head, which displays only the most rudimentary of behaviors (blinking, looking left and right). Embodied conversational agents, on the other hand, are by definition models of human behavior, which means that at least along some dimension they must function in the same way humans do. Thus, the pedagogical agent of Wang and his colleagues and the virtual actor created by Walker, Cahn, and Whittaker both rely on Brown and Levinson’s theory of politeness and language use.⁴ Poggi and Pelachaud base the facial expressions of their ECA on Austin’s theory of performatives.⁵ Likewise, ECAs are fully functioning artificial intelligence systems in the sense that they understand language by composing meanings for sentences out of the meanings of words; they deliberate over an appropriate response, deliver the response, and then remember what they said so as to make the subsequent conversation coherent. They mostly have both heads and bodies, and their behavior is based on observation of human behavior.

Figure 17.3 shows an ECA named REA (for Real Estate Agent) who was programmed on the basis of a detailed examination into the behavior of realtors and clients. Over a period of roughly five years, various graduate students, postdocs, and colleagues in my research group studied different aspects of house-buying talk and incorporated their findings into the ECA. Hao Yan looked at what features of a house description were likely to be expressed in hand gestures and what features in speech. Yukiko Nakano discovered that posture shifts were correlated with shifts in conversational topic and shifts in whose turn it was to talk. Tim Bickmore examined the ways in which small talk was employed to establish trust and rapport between realtor and client. Earlier work by Scott Prevost on intonation and by Obed Torres on patterns of eye gaze also went into the implementation. As our research into human conversation progressed, we also came to better understand some of the overall properties of human conversation, and those insights were also incorporated.⁶



FIGURE 17.3 | REA, the virtual real estate agent.

The result was a virtual woman who tried to sell a home to whomever approached her. A small camera on top of the screen allowed her to detect the presence of real humans and initiate a conversation with them. Her knowledge of appropriate realtor conversation led her to ask questions about a person’s housing needs and then nod, seem to reflect, and pull up data on appropriate properties, describing them using a combination of descriptive hand gestures, head movements, and spoken language.

A SAMPLE INTERACTION WITH AN ECA

Figure 17.4 illustrates an actual dialogue between a human user and REA.

Conversational Properties

As well as demonstrating the use of several conversational modalities, such as speech, hand gestures, and head movements, in this example REA is engaging in some very subtle, humanlike behavior that demonstrates four of the key properties of human face-to-face conversation. Those four properties are (1) the distinction between interactional and propositional functions of language and conversation; (2) the distinction between conversational behaviors (such as eyebrow raises) and conversational functions (such as

Mike approaches the projection screen. Rea is currently turned side on and is gazing idly about. As Mike moves within range of the two cameras mounted above the screen, Rea turns to face him and says:

REA> Hello. How can I help you?

MIKE> I'm looking to buy a place near MIT.

Rea nods, indicating that she is following.

REA> I have a house to show you.

A picture of a house appears on-screen behind Rea who

blinks and looks at the house and then at Mike.

REA> It's in Somerville.

MIKE> Tell me about it.

Rea briefly looks up and away while she plans what to say.

REA> It's big.

Rea makes an expansive gesture with her hands.

Mike brings his hands up as if he is about to speak and so

Rea does not continue; instead waiting for him to speak.

MIKE> Tell me more about it

REA> Sure thing. It has a nice garden.

Rea sketches a curved gesture with her hands indicating that the garden extends along two sides of the house.

MIKE> How far is it?

REA> It is five minutes to the Porter Square T station.

Rea makes it clear that it is five minutes on foot from the

T station by making a walking gesture with her fingers.

MIKE> How big is the house?

REA> It has four bedrooms, three bathrooms. . .

And the house tour continues.

FIGURE 17.4 | Transcript of an interaction between REA and a person.

turn-taking); (3) the importance of timing among conversational behaviors; (4) the deployment of each modality to do what it does best. Our insights into each of these properties has come through the cycle of watching real humans, attempting to model what we see in virtual humans, and observing the result or observing people interacting with the result.

DIVISION BETWEEN PROPOSITIONAL AND INTERACTIONAL FUNCTIONS | Some of the things that people say to one another move the conversation forward, while others regulate the conversational process. Propositional information corresponds to the content (sometimes referred to as transmission of information) and includes meaningful speech as well as hand gestures that represent something, such as punching a fist forward while saying “she gave him one” (indicating that the speaker’s meaning is that she punched

him, not that she gave him a present). Interactional information regulates the conversational process and includes a range of nonverbal behaviors (e.g., quick head nods to indicate that one is following, bringing one’s hands to one’s lap and turning to the listener to indicate that one is giving up one’s turn), as well as sociocentric speech (“Huh?” or “Do go on”). It should be clear from these examples that both functions may be filled by either verbal or nonverbal means. Thus, in the dialogue excerpted above, REA’s nonverbal behaviors sometimes contribute propositions to the discourse, such as the gesture that indicates that the house in question is five minutes on foot from the T stop, and sometimes they regulate the interaction, such as the head nod that indicates that REA has understood Mike’s utterance.

DISTINCTION BETWEEN FUNCTION AND BEHAVIOR | When humans converse, few of their behaviors are hard-coded. That is, there is no mechanism or database “look-up table” that gives the appropriate response for every possible conversational move on the part of one’s partner. Every day we hear thousands of phrases that we have never heard before, assembled through the infinite creativity of language use, and we reply to each of these phrases in just a couple of milliseconds, with an equally creative response. Gestures and head movements are no more likely to be routinized—head nods will look different if we are looking up at a taller interlocutor or down at somebody short, if we are wearing a hat or bareheaded. And other than the small number of culturally meaningful gestures (such as “V for victory” or “Up yours”), gestures display a great variety across people and even for one person across time. In observing human-human conversation, our group discovered that speakers do not always nod when they understand. Instead they sometimes signal that they are following along by making agreement noises such as “uh-huh.” In our simulation of this behavior, then, instead of hard-coding, the emphasis is on identifying the high-level structural elements that make up a conversation. We describe these elements in terms of their role or function in the exchange. Typical discourse functions include conversation invitation, turn-taking, providing feedback, contrast and emphasis, and breaking away. Each function can be filled through a number of different behaviors, in one or several modalities. The form given to a particular discourse function depends on, among other things, current availability of modalities such as the face and the hands, type of conversation, cultural patterns, and personal style.

REA generates speech, gestures, and facial expressions based on the current conversational state, the conversational function she is trying to convey, and the availability of her hands, head, and face to engage in the desired

Although it has long been known that the most effortful part of a gesture occurs with the part of an utterance that receives prosodic stress,⁷ it wasn't until researchers needed to generate gestures along with speech in an ECA—and therefore needed to know the details of the context in which meaningful gestures were most likely to occur—that they discovered that a gesture is most likely to occur with the rhematic, or new contribution, part of an utterance.⁸ This means that if a speaker is pointing to her new vehicle and saying, “This car is amazingly comfortable. In fact, this car actually has reclining seats,” the phrase “amazingly comfortable” would be the rheme in the first sentence, because “car” is redundant (since the speaker is pointing to it) and “reclining seats” would be the rheme in the second sentence, because “car” has already been mentioned. Therefore, the speaker would be most likely to produce hand gestures with “amazingly comfortable” and “reclining seats.”

USING THE MODALITIES TO DO WHAT THEY DO BEST | E-mail obliges us to compress all of our communication goals into textual form (plus the occasional emoticon). In face-to-face conversation, on the other hand, humans have many more modalities of expression at their disposal, and they depend on each of them, as well as various combinations of them, to communicate what they want to say. They use gestures to indicate things that may be hard to represent in speech, such as spatial relationships among objects,⁹ and they depend on the simultaneous use of speech and gesture to communicate quickly. In this sense, face-to-face conversation may allow us to be maximally efficient or, in other instances, to use conversation to do other kinds of work than information transmission (for example, we may use the body to indicate rapport with others, while language is getting task work done). In the dialogue reproduced above, REA uses the hands' ability to represent spatial relations among objects and places by indicating the shape of the garden (sketching a curved gesture around an imaginary house) while her speech gives a positive assessment of it (“it has a nice garden”). However, in order to produce this description, the ECA needs to know something about the relative representational properties of speech and gesture, something about how to merge simultaneous descriptions in two modalities, and something about what her listener does and does not already know about the house in question.

The need to understand how speech, gestures, and movements of the head and face can be produced together by ECAs has forced me to design experimental and naturalistic methodologies to examine the nature of the

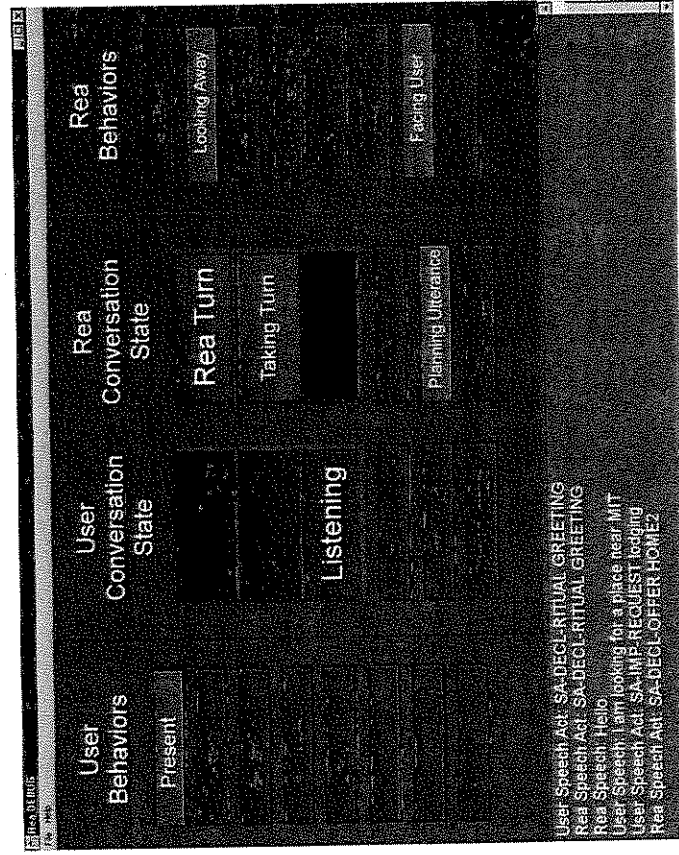


FIGURE 17.5 | Visualization of ECA and human conversational state.

behavior. For example, when the user first approaches REA (“user present” state), she signals her openness to engage in conversation by looking at the user, smiling, and/or tossing her head. Figure 17.5 shows a visualization of REA’s internal state with respect to conversational behaviors and conversational states.

IMPORTANCE OF TIMING | The relative timing of conversational behaviors plays a large role in determining their meaning. That is, for example, the meaning of a nod is determined by where it occurs in an utterance, all the way down to the 200-millisecond scale; consider the difference between “you did a [great job]” (square brackets indicate the temporal extent of the nod) and “you did a [...] great job”). Thus, in the dialogue above, REA says “it is five minutes from the Porter Square T station” at exactly the same time as she performs a walking gesture. If the gesture occurred in another context, it could mean something quite different; if it occurred during silence, it could simply indicate REA’s desire to take the turn.

interaction between modalities, and has resulted in significant advances in my theorizing about the relationship between speech and gesture in humans. Thus, for example, in my current work with the purple virtual robot NUMACK as a simulation, Paul Tepper, Stefan Kopp, and I have become interested in the seeming paradox of how gesture communicates, given that there are no standards of form in spontaneous gesture—no consistent form-meaning mappings. Some gestures clearly depict visually what the speaker is saying verbally, and these gestures are known as *iconics*. What is depicted on the fingers, however, and its relationship to what is said can be more or less obvious. And two speakers' depiction of the same thing can be quite different.

For example, we compared two people describing the same landmark on Northwestern University's campus: an arch at the entrance to the campus that lies at the intersection of Sheridan Road and Chicago Avenue. In order to collect these data, we hid prizes in various spots on campus and asked one student, who knew where the prize was hidden, to give directions to the prize to a second student. If the second student succeeded in finding it, the two shared the prize (and both were entered into a drawing for an iPod, probably the most motivating feature of the experiment!). We used four cameras to videotape the direction-giving, training them on different parts of the bodies of the two speakers, as described above (and shown in figure 17.2), and then we transcribed each gesture and the speech that accompanied it for further study. One speaker in the experiment, describing directions to a church near the arch, said "go to the arch" and, with his fingertips touching one another and pointing upward, made a kind of tepee shape. In this instance, the gesture seemed to indicate a generic arch, but this person, although his fingertips were touching one another, pointed his fingers toward the listener with his thumbs up, making the shape of a right angle. In this instance, the gesture seems to indicate . . . what? An arch lying on its side?! It makes, in fact, no sense to us as observers—unless we know that the arch is located at the right angle formed by Sheridan Road and Chicago Avenue. And the speaker's next utterance, "It's located at the corner of Sheridan," supports this interpretation of the gesture. So, in the absence of the relatively stable form-meaning pairing that language enjoys (the same image may not be evoked for both of us, but when I say "right angle" I can be relatively sure that we will both imagine something similar), how do gestures communicate? The answer to this question (which is outside the scope of this chapter but involves the interpretive flexibility of

gestures, which have meanings only in situated contexts) resulted both in a new computational architecture, through which gesture and speech are computationally generated together, and a new way of understanding of how gestures communicate among humans.

Translating Conversational Properties into Computational Architectures

The four conversational properties discussed in the previous section gave rise in 2000 to the computational architecture represented in figure 17.6. As this diagram makes clear, and like many systems in Artificial Intelligence, ECAs are largely linear and devoid of contingent functionality—a person asks a question that is collected by the input modules of the system (cameras to view the speaker's gestures and posture, microphones to hear the speech) and then interpreted into a unified understanding of what the speaker meant. In turn, that understanding is translated into some kind of obligation to respond. That response is planned out first in "thought" or communicative intention, and then in speech and movements of the animated body, face, and hands through the use of a speech synthesizer, computer graphics engine, and various other output modes. Meanwhile, so as not to wait for all of that processing to be completed before a response is generated, a certain number of hardwired responses are sent to be realized: head nods, phatic noises (mmm, uh-huh), and shifts of the body.

The linear nature of this architecture is one of the constraints imposed by the scientific instrument—like trying to cut out circles with straight blades.

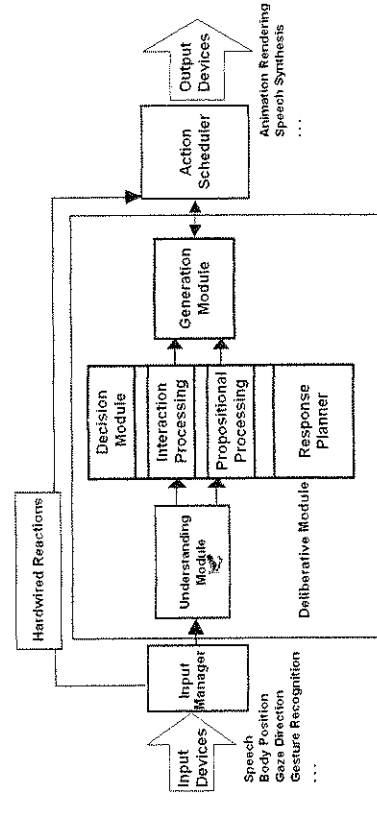


FIGURE 17.6 | Computational architecture of an ECA.

When I first began to collaborate with computer scientists in 1993–94 to build a virtual human, I asked them to build one that was responsive to itself and to its interlocutor in a number of ways. I told them that I wanted the virtual human to be able to “see” its own hands gesturing, and from what it saw decide what it wanted to say in the moment—the way humans often do, for example, when they cannot recall a word until they make the gesture for it. And I told them I wanted some kind of entrainment or accommodation between the different participants in the conversation, such that their language and gestures grew increasingly alike as they came to mirror one another. The response was incredulity and a request for me to be better informed before I went asking for features. The goal, I was told, was autonomy and not codependence. Of course, as Suchman has pointed out about other work in Artificial Intelligence, this means that we have not produced a truly conversational agent, since “interaction is a name for the ongoing, contingent co-production of a shared social/material world.”¹⁰ But the kinds of interdependence that we wish to simulate are hard to achieve with our current models.

In general terms, however, building ECAs has forced researchers in human behavior to attend to the integration of modalities and behaviors in a way that merges approaches from fields that usually do not speak to one another: ethnomethodological interpretive and holistic studies of human communication merge with psycholinguistic, experimental, isolative studies of particular communicative phenomena. To build a human entails understanding the context in which one finds each behavior—and that context is the other behaviors.

During that first collaboration with computer scientists, when we were building the very first of these animated embodied conversational agents, a different researcher was implementing each part of the body. Catherine Pelachaud was writing the algorithms to drive the character’s facial movements (head nods, eye gaze, etc.) based on conversational parameters such as who had the turn. Scott Prevost was writing rules to generate appropriate intonation—the prosody of human language—on the basis of the relationship between the current utterance and previous utterances. I was working on where to insert gestures into the dialogue. After several months of work, we finally had a working system. In those days, ECAs needed to be “rendered”—they were not real-time—and so with bated breath we ran the simulation, copied it to videodisc, and then watched the result. The result was an embodied conversational agent who looked like he was speaking to very small children or to foreigners. That is, the resultant virtual human used so

many nonverbal behaviors that signaled the same thing, that he seemed to be trying to explain something to a listener who didn’t speak his own language or was just very stupid. This system, called Animated Conversation, was first shown at SIGGRAPH, the largest computer graphics conference, for an audience of four thousand researchers and professional animators (the folks who build cartoons and interactive characters), and they found it hilarious. To my mind, on the other hand, we had made a huge advance. We had realized that the phenomena of hand gestures, intonation, and facial expressions were not separate systems, nor was one a “translation” of the others; instead, they had to be derived from one common set of communicative goals. That was the only explanation for the perception of overemphasizing each concept through a multiplicity of communicative means. The result fundamentally changed the way we build embodied conversational agents, but it was an advance in understanding human communication as well. It also led to a design methodology that I have relied on ever since (see figure 17.7). Iteratively, my students and I collect data on human-human conversation, interpret those data in such a way as to build a formal model, implement a virtual human on the basis of the model, confront the virtual human with a real human, evaluate the results, and collect more data on human-human communication

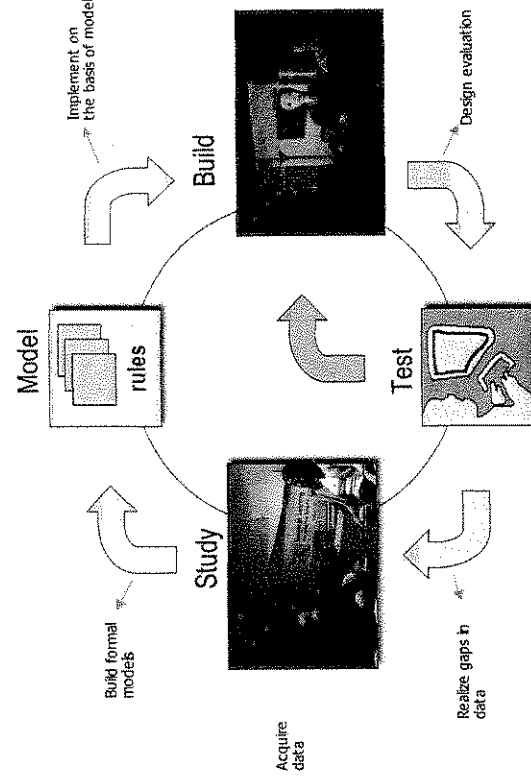


FIGURE 17.7 | Methodology for modeling human conversation and building ECAs.

our own performance in the world.¹⁴ They allow us to evaluate our hypotheses about the relationship between verbal and nonverbal behavior, and to see what gaps exist in our knowledge about human communication, by seeing ourselves and our conversational partners in the machine. How do we go about evaluating our hypotheses? As described above, we watch the ECAs and observe our own reactions. But we also put others in front of the ECAs and examine the differences between their behavior with ECAs and their behavior with other humans. This second kind of experiment relies on the supposition that correctly implemented virtual humans evoke natural human responses. Mechanisms that seem human make us attribute humanness and aliveness to them and make us react in natural human ways. Successful virtual humans evoke distinctly human characteristics in our interaction with them. The psychological approach to artificial life leads to functional bodies that are easy to interact with, “natural” in a particular sense: they evoke a natural response.

In an early experiment, for example, Kris Thorisson and I compared reactions to three versions of an ECA called Gandalf (this was 1996, and the ECA consisted of a head with one disembodied hand, as shown in figure 17.8). Our goal was to demonstrate, in those early days, that interactional behaviors—ones that did not move the conversation forward—could be simulated computationally, and that those behaviors in virtual humans would elicit similar behaviors on the part of human interlocutors. An additional goal was to demonstrate that if we had to choose only certain nonverbal functions to implement computationally, they should be interactional (what we called “envelope” functions) rather than emotional functions. We felt

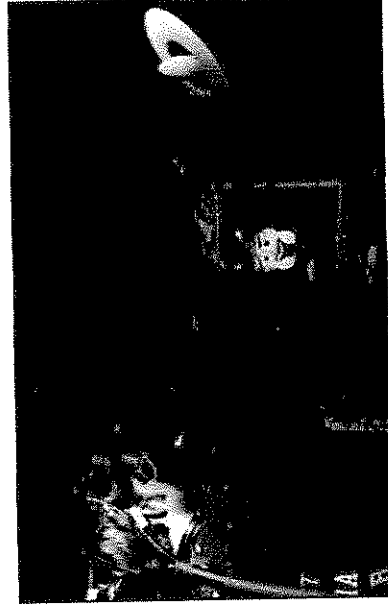


FIGURE 17.8 | Person interacting with Gandalf.

if needed. (A side effect of this methodology is the need to confront the response of lay viewers to the necessary flaws and lacunae in the implementation, but I try to think of that as character building.)

I should repeat that building a computational system has traditionally demanded a formal or predictive model. That is, in addition to being able to interpret why a particular experience occurs in a particular context, one must also be able to predict in the future what set of conditions will give rise to a particular experience so that one can generate those behaviors in the ECA in response to the appropriate conditions.¹¹ Unfortunately, predictive models come with their own baggage, as they tend to underscore the way in which fixed sets of conditions give rise to fixed outputs, as opposed to highlighting the very contingent, coproduced nature of human conversation where, on the fly, hearers and speakers influence one another's language and indeed their very thinking patterns, as Suchman has forcefully argued.¹² In this sense, I sometimes worry that building computational simulations of this sort may set back the computational study of language; that phenomena that cannot yet be modeled in virtual people will be ignored. On the other hand, for the most part, before the advent of ECAs, computational linguistics and work on dialogue systems (which arose from the cognitive sciences—psychology, linguistics, philosophy, computer science) concentrated on the propositional functions of language, which many linguists saw as the primary if not the only functions of language. Before ECAs, computational models of language were, for the most part, capable only of simulating task talk bereft of social context and bereft of nonverbal behavior. And given the power of these computational models, perhaps the arrival of ECAs, with their attention to the non-informational and socially contextualized functions of language, has played some positive role in the cognitive sciences.

Now that there has been a decade of research on ECAs, several researchers, including myself, are beginning to explore other kinds of computational architectures and techniques that do not require deterministic, formal, input-output style models of conversation. Probabilistic techniques—such as spreading activation, Bayesian reinforcement learning, and partially observable Markov decision processes—are being applied to the newest phenomena to be modeled with ECAs. These phenomena, which tend to have more to do with social context than local linguistic context, include the effect of emotion on verbal and nonverbal behavior in conversation, the role of personality and cultural differences, social influence, etiquette, and relationship-building.¹³

In all of these implementation experiments, ECAs are tools to think with, much like other computer software and hardware that allow us to evaluate

that emotional reactions should be studied only after these very ubiquitous interactional behaviors had been simulated.

In the first version, called “content-only,” the ECA spoke but used no nonverbal expressions of any kind. An example of an interaction with such an ECA follows:

Gandalf: “What can I do for you?” (Face looks at user. Eyes do not move.)

User: “Will you show me what Mars looks like?” (User looks at Gandalf.)

Gandalf: “Why not—here is Mars.” (Face maintains orientation. No change of expression.

Mars appears on monitor.)

User: “What do you know about Mars?” (User looks at map of solar system.)

Gandalf: “Mars has two moons.” (Face maintains orientation. No change of expression.)

In the second version, called “content + envelope,” the ECA spoke and also used eye gaze, eyebrow and head movements, and an occasional hand wave to moderate turn-taking, feedback, and other conversational-envelope processes. An example of an interaction with such an ECA follows:

User: “Is that planet Mars?”

Gandalf: “Yes, that’s Mars.” (Gandalf raises eyebrows and performs head gesture while saying “yes,” turns to planet, and points at it while saying “that’s Mars,” and then turns back to face user.)

User: “I want to go back to Earth now. Take me to Earth.” (User looks at map of solar system, so Gandalf looks at solar system.)

Gandalf: “OK. Earth is third from the sun.” (Gandalf turns to planet as he brings it up on the screen, then turns to user and speaks.)

User: “Tell me more.” (Gandalf takes about 2 seconds to parse the speech, but he knows within 250 ms when the user gives the turn, so he looks to the side to show that he’s taking the turn, and his eyebrows go up and down as he hesitates while parsing the user’s utterance.)

Gandalf: “The Earth is 12,000 kilometers in diameter.” (Gandalf looks back at the user as he speaks.)

In the final version, called “content + emotional,” the ECA spoke and also smiled, frowned, and looked puzzled as the occasion warranted. An example of an interaction with such an agent follows:

Gandalf: “What can I do for you?” (Gandalf smiles when user’s gaze falls on his face, then stops smiling and speaks.)

User: “Take me to Jupiter.” (User looks at screen and then back at Gandalf, and so Gandalf smiles.)

Gandalf: “Sure thing. That’s Jupiter.” (Gandalf smiles as he brings Jupiter into focus on the screen.)

User: (Looks back at Gandalf. Short pause while deciding what to say to Gandalf.)

Gandalf: (Looks puzzled because the user pauses longer than expected. Waits for user to speak.)

User: “Can you tell me about Jupiter?”

The study consisted in asking people to interact with Gandalf and then examining the real human’s conversational-envelope and emotional behaviors during the interaction, as well as asking subjects to fill out a questionnaire assessing “lifelikeness.” What we discovered was that participants tended to mimic the virtual human: if he stood rigid, so did they; if he was animated, so were they. In fact, the people standing in front of the content-only version of Gandalf were most animated in their expressions of frustration—sighs and the occasional request for signs of life (“Gandalf, are you there?”). People interacting with the content + envelope version, on the other hand, started off wary as Gandalf’s head and single hand began to describe the solar system, and then, after an utterance or two, became more animated, gesturing and nodding to Gandalf in much the same way as they had to the experimenter before the experiment started.¹⁵ Finally, we discovered no difference in the people’s interaction with, nor in their assessment of, the ECA between the content-only version and the content + emotion version.

More recently, Yukiko Nakano and I carried out a study of the role of nonverbal behaviors in grounding and how these behaviors could be implemented in a virtual human.¹⁶ Common ground is the sum of mutual knowledge, mutual beliefs, and mutual suppositions necessary for a particular stage of a conversation.¹⁷ Grounding refers to the ways in which speakers and listeners ensure that the common ground is updated, such that the participants understand one another. Grounding may occur by nodding to indicate that one is following, by asking for clarification when one does not understand, or by asking for feedback, as in “You know what I mean?” Here, too, an extensive study of human-human behavior in the domain of direction-giving paved the way for the implementation of an ECA that could ground while giving directions using a map and hand gestures. And here, too, we evaluated our work by comparing people’s reactions to two versions of the ECA: one that demonstrated grounding behaviors and one that had the grounding “turned off.” When the ECA’s grounding behaviors were turned off, the person simply acted as if she were in front of a kiosk and not another human—not gazing at the ECA or looking back and forth between him and the map. When the ECA did engage in grounding behaviors, the human acted strikingly . . . human, looking back and forth between the map and the ECA, as shown in figure 17.9.

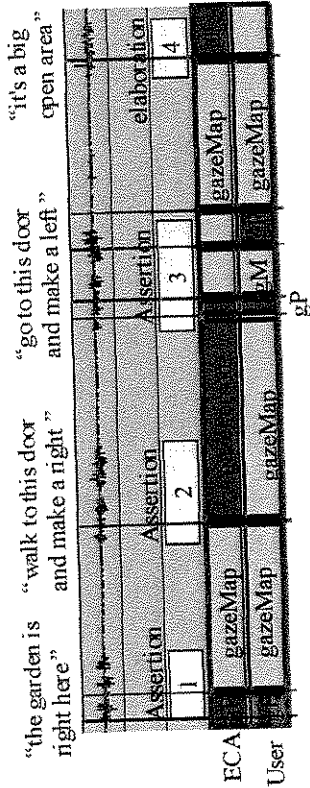


FIGURE 17.9 | Analysis of grounding behaviors in human-ECA conversation.

A final example, and perhaps the one that best illustrates how ECAs—properly constructed, on the basis of theories elaborated from human observation—can elicit natural behavior from humans who interact with them (and how this behavior can be illuminating along both positive and negative dimensions), is an experiment in which we endowed REA with social chitchat skills.¹⁸ As mentioned above, Tim Bickmore carried out an extensive study of small talk among realtors and traveling salesmen. The results indicated that small talk was not randomly inserted into conversation, but served specific purposes (e.g., to minimize the potential face threat of personal questions such as “How much do you earn?”). These functions of small talk could be simulated in a way that would allow us to implement a small-talking realtor who used chitchat to smooth the progress of a house-selling transaction with a human. In order to test our model of human conversation, we asked people to interact with one of two versions of the ECA. One used task-talk only, while the other added social chitchat at key places in the interaction. Figures 17.10 and 17.11 show the actual conversational contributions by REA under the two conditions.

The people who interacted with each ECA were asked to evaluate their experience: how natural they felt the interaction to be, how much they liked the ECA, how warm they felt she was, and how trustworthy. We also tested the subjects on their own social skills, dividing them into extroverts and introverts by using a common psychological scale. The results showed that extroverts preferred the small-talk version of the ECA, while introverts preferred the ECA to keep to the task (we also discovered that it was difficult to find extroverts among MIT students, but that’s another story).

An introvert working with the small-talk version remarked:

- 1 How about this weather?
- 2 I think winters in Boston are awful.
- 3 How do you like Boston?
- 4 I have lived in Boston all my life. Come to think of it, I have lived inside this room all of my life. It is so depressing.
- 5 Boston is certainly more expensive than it used to be.
- 6 So, where would you like to live?
- 7 How many bedrooms do you need?
- 8 Do you need access to the subway?
- 9 Is one bath enough?
- 10 You know, I keep showing the researchers here the same houses, over and over again. Maybe one day I will get lucky.
- 11 Have you been in the Media Lab before?
- 12 Do you know that the Media Lab is going to expand into another building. Things are really going well for the researchers here.
- 13 It is pretty cool do you think?
- 14 They are doing some crazy things in here.
- 15 I have shown houses to lots of students and faculty from MIT. But I always enjoy talking to them.
- 16 Anyway, What can you afford?
- 17 What kind of down payment can you make?
- 18 Let me see what I have available.

FIGURE 17.10 | Small talk in ECAs.

REA exemplifies some things that some people, for example, my wife, would have sat down and chatted with her a lot more than I would have. Her conversational style seemed to me to be more applicable to women, frankly, than to me. I come in and I shop and I get the hell out. She seemed to want to start a basis for understanding each other, and I would glean that in terms of our business interaction as compared to chitchat. I will form a sense of her character as we go over our business as compared to our personal life. Whereas my wife would want to know about her life and her dog, whereas I really couldn’t give a damn.

An extrovert working with the same version had a very different response: “I thought she was pretty good. You know, I can small-talk with somebody for a long time. It’s how I get comfortable with someone and how I get to trust them and understand how trustworthy they are, so I use that as a tool for myself.”

Clearly, the people in this experiment are evaluating the ECA’s behaviors in much the same way as they would evaluate a flesh-and-blood realtor. And clearly, our unexamined implementation of the realtor as a woman instead of a man has played into those evaluations, as much as have any

- 1 So, where would you like to live?
- 2 What can you afford?
- 3 What kind of down payment can you make?
- 4 How many bedrooms do you need?
- 5 Do you need access to the subway?
- 6 Is one bath enough?
- 7 Let me see what I have available.

FIGURE 17.11 | Task-only talk in ECAs.

of our carefully examined decisions about small talk, hand gestures, and body posture. Although our goal was to obtain input for a theory of the role of small talk in task talk, this response from one of REA's interlocutors effectively demolishes the claim that human identity can be denuded of its material aspects. Much of the previous work on responses to ECAs as inter-faces has concentrated on exactly this sort of effect, with some researchers advising industry executives to use a female ECA to sell phone service, but a male ECA to sell cars.¹⁹ In response to this unintended research finding in our small-talk study, my students and I have begun to use the virtual human paradigm to investigate explicitly which linguistic, nonverbal, and visual cues signal aspects of identity. Some have suggested that the race of ECAs be paired to the putative race of the user; my students and I have begun to look at the complex topic of racial identity and how a person's construction of his/her own race, and recognition of the racial identity of others, may be conveyed not just by skin color, but (also) by aspects of linguistic practice, patterns of nonverbal behavior, and narrative style.²⁰

Embodied Conversational Agents as Interfaces

I've alluded to other ways in which ECAs are used, where they serve not as scientific instruments or tools to think with, but as interfaces to computers. In this function, ECAs might take the place of a keyboard, screen, and mouse—the user speaks to them instead of typing. Or they might represent the user in an online chat room. ECAs can also serve as teachers or tutors in combination with educational software—so-called “pedagogical agents.” Research in this applied science examines whether ECAs are preferable to other modalities of human-computer interaction such as text or speech; what kinds of behaviors make the ECAs most believable and most effective (as tutors, information retrievers, avatars); and what personas the ECA should adopt in order to be accepted by their users. My students and I have also

conducted some of this research, looking at whether virtual children are effective learning companions for literacy skills, whether people are willing to be represented by ECAs in online conversations, and whether tiny ECAs—small enough to fit on a cell phone—still evoke natural verbal and nonverbal responses in the people speaking with them.²¹ Even here, however, our research on virtual peers has led us back to an exploration of human-human communication, as we attempt to identify the features that signal to children that somebody else is a peer, is good friendship material, is worth listening to and telling stories with. Our related exploration of the pragmatics of the body has led us to some key features of social interaction—how rapport and friendship are negotiated—which, in turn, have led us to a better understanding of peer learning. Figure 17.12 shows one of our virtual peers.

Most recently, Andrea Tartaro and I have begun to look at how children with autism can play the role of scientist—learning about the gaps in their knowledge of communication and social interaction by authoring virtual people and watching them interact with others.²² Mostly, however, our work is focused on the minutiae of human interaction and is therefore sometimes less useful to interface designers. In fact, computer scientists sometimes respond to my talks about NUMACK, the direction-giving robot, by asking, “But wouldn't it just be more effective to display a map on the computer screen and skip the virtual human?” When I respond that such an interface

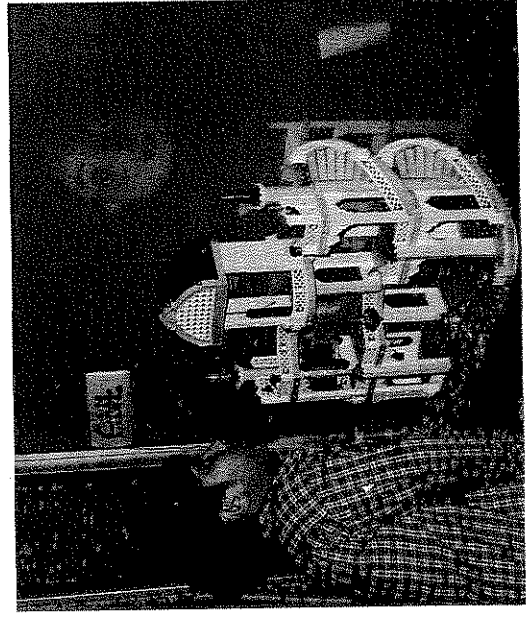


FIGURE 17.12 | A child playing with Sam, the virtual peer.

wouldn't teach us anything about human communication, those same questioners often nod sagely, as if they knew all along that my interest was only in humans. Others have taken the ECA much further as an interface—the most extreme work (and most studied by historians of science and technology) is being done at the Institute for Creative Technologies (ICT) at the University of Southern California. Funded in equal parts by the U.S. Army and Hollywood, the ICT has created a vast, immersive, video-game-like room intended to teach soldiers before they enter the field—what Tim Lenoir has called a “military entertainment complex.”²³

The development of the ECA from a scientific instrument that simulates human behavior to an attractive interface bears interesting parallels to the history of mechanical automata. Automata makers of the sixteenth century, such as the one who built the perpetually praying monk described so elegantly by Elizabeth King in this volume, depended on the gaze of the perceiver to confer lifelikeness on the machine. Automata makers of the eighteenth century intended to find out in what way the human activities of drawing and writing and playing an instrument differed, if at all, when machines performed them.²⁴ Droz's writing boy, whose pen moves across the page just as real writers' pens move, is an example. The ECAs that I build today are likewise a way to compare conversation among humans with conversation between a human and a humanlike machine in order to discover what we know and do not know about human communication, and that simulation only works because of the “life” conferred on the virtual human by the interlocutor. Mechanical automata of the later nineteenth century, however, were meant to entertain, and not illuminate. An example of such a pretty virtual body as entertainment is the Pierrot automaton doll that writes—but simply by moving an inkless pen smoothly across a page—while sighing deeply and progressively falling asleep by the lamplight. These latter examples of mechanical humans did sustain relationships with real humans in that people wanted to own the pretty mechanical toys and were fascinated by them. But in these instances, the gaze of the viewer was one of concupiscence, not that of an interlocutor. Likewise, the tiny virtual human on a cell phone is meant to evoke the greedy desire of the collector more than the unconscious gaze of a partner in conversation.

CONCLUSIONS

These exercises in building virtual people have led to advances in what we know about the interaction between verbal and nonverbal behavior in

humans, about the role of small talk in task talk, about the kinds of functions performed in conversation by the different modalities of the body, and about how learning is linked to rapport in children. Learning what must be implemented in order to make embodied conversational agents evoke a natural response and learning what the technology can and cannot do at the present time have also given me a sense of the meaning of humanness through human behavior. It is the ensemble of behaviors in all of their minuteness and unconscious performance that create the sense of humanness. Flaws and lacunae in that ensemble of behaviors give the scientist interlocutor a sense of what we do not know about human communication. Strengths and continuities in the theory that underlies the implementation lead people to respond to the virtual human as if it were another human being. The sufficiency criterion in cognitive science consists of explaining human cognitive activity by showing how a computer program may bring about the same result when the computer is provided with the same input.²⁵ In virtual human simulations, cognitive activity is not sufficient. I know that my model successfully explains human behavior when it evokes human behavior, because human communicative behavior is intrinsically relational and cannot be understood without two humans.

To come back to the anecdote with which this essay began, it is important to note the essential role of the physical body in the study of both language and social experience (insofar as those might be distinguishable). Language has traditionally been relegated to taking place purely in the head. But I hope the examples of communicative functions given above have made it clear that language is spread throughout the whole body—the hands, the torso, the eyes—and across two bodies in interaction. My original goal in building virtual humans was to focus attention on the whole-body aspects of language and on its intrinsically relational nature. As Descartes points out, the difference between real men and artifacts that only have the shape of men exists both in word and movement: imitation and gesture are as constitutive of humanness and social interaction as spoken language:

And suppose there existed machines built in the image of our bodies, and capable of imitating our actions, as far as morally possible, there would still remain two certain tests by which to know that they were not really men. The first is that these automata could never use words or other signs in conversation, as we are able to do in order to convey our thoughts to others; for even if we can easily conceive of a machine that can emit the sounds of speech, or that can respond to external action such that, for example, if touched in one particular place it

may ask what we wish to say to it; if touched in another it may cry out that it is hurt, and so forth; we nevertheless cannot imagine a machine that can answer to what is said in its presence, as even fools can do. The second test is that even though such machines may carry out actions as well or even more perfectly than we humans can, they still will fail in executing other actions, by which we can discover that they did not act from knowledge, but from a particular arrangement of their organs.¹⁶

NOTES

The research reported in this chapter would not have been possible without the hard work, persistence, and insight of graduate students who were so generous with their time and knowledge that they quickly became my colleagues and teachers. Sincere thanks go to the members of Animated Conversation, Gesture and Narrative Language, and the ArticLab. Thanks to Jessica Riskin for including me in this volume, and abiding gratitude to Ken Alder, Pablo Boczkowski, Sid Horton, Jessica Riskin, Dan Schwartz, Matthew Stone, and two reviewers for careful and perceptive comments that greatly improved the quality of the manuscript.

1. Philip Agre, "Formalization as a Social Project," *Quarterly Newsletter of the Laboratory of Comparative Human Cognition* 14, no. 1 (1992): 25–27.
2. Clifford Ivar Nass and Scott Brave, *Wired for Speech: How Voice Activates and Advances the Human-Computer Relationship* (Cambridge, MA: MIT Press, 2005).
3. Evelyn Fox Keller, "Models, Simulation, and 'Computer Experiments,'" in *The Philosophy of Scientific Experimentation*, ed. H. Radder (Pittsburgh: University of Pittsburgh Press, 2003).
4. See N. Wang, W. L. Johnson, P. Rizzo, E. Shaw, and R. E. Mayer, "Experimental Evaluation of Polite Interaction Tactics for Pedagogical Agents," in *Proceedings of the International Conference on Intelligent User Interfaces* (New York: ACM Press, 2005); Marilyn A. Walker, Janet E. Cahn, and Stephen J. Whittaker, "Improvising Linguistic Style: Social and Affective Bases for Agent Personality," in *Proceedings of the First International Conference on Autonomous Agents* (New York: ACM Press, 1997); and Penelope Brown and Stephen C. Levinson, *Politeness: Universals in Language Usage*, vol. 4 of *Studies in Interactional Sociolinguistics* (New York: Cambridge University Press, 1987).
5. Isabella Poggi and Catherine Pelachaud, "Performative Facial Expressions in Animated Faces," in *Embodied Conversational Agents*, ed. J. Cassell, J. Sullivan, S. Prevost, and E. Churchill (Cambridge, MA: MIT Press, 2000); John Austin, *How to Do Things with Words* (Oxford: Oxford University Press, 1962).
6. See Hao Yan, *Fused Speech and Gesture Generation in Embodied Conversational Agents* (M.S. diss., MIT Media Lab, Massachusetts Institute of Technology, 2000); Justine Cassell, Yukiko Nakano, Timothy Bickmore, Candy Sidner, and Charles Rich, "Non-Verbal Cues for Discourse Structure," in *Proceedings of the Thirty-ninth Annual Meeting of the Association of Computational Linguistics* (Philadelphia: Linguistics Data Consortium, 2002); Timothy Bickmore and Justine Cassell,

"Small Talk and Conversational Storytelling in Embodied Conversational Characters," in *Narrative Intelligence: Papers from the Fall 1999 AAAI Symposium* (Menlo Park, CA: AAAI Press, 1999); Scott Allan Prevost, "Modeling Contrast in the Generation and Synthesis of Spoken Language," in *Proceedings of the Fourth International Conference on Spoken Language Processing [ICSLP '96]* (New York: IEEE, 1996); and Obed E. Torres, Justine Cassell, and Scott Prevost, "Modeling Gaze Behavior as a Function of Discourse Structure," in *Proceedings of the First International Workshop on Human-Computer Conversation* (n.p., 1997).

7. Adam Kendon, "Some Relationships between Body Motion and Speech," in *Studies in Dyadic Communication*, ed. A. W. Siegman and B. Pope (Elmsford, NY: Pergamon Press, 1972).
8. M. A. K. Halliday, *Intonation and Grammar in British English* (The Hague: Mouton, 1967).
9. Justine Cassell, Matthew Stone, and Hao Yan, "Coordination and Context-Dependence in the Generation of Embodied Conversation," in *Proceedings of the INLG 2000* (Münzpe Ramon, Israel: Association of Computational Linguistics, 2000).
10. Lucy Suchman, "Writing and Reading: A Response to Comments on Plans and Situated Actions," *Journal of the Learning Sciences* 12, no. 2 (2003): 299–306.
11. Daniel L. Schwartz and Taylor Martin, "Representations That Depend on the Environment: Interpretative, Predictive, and Praxis Perspectives on Learning," *Journal of the Learning Sciences* 12, no. 2 (2003): 285–97.
12. Lucy Suchman, "Do Categories Have Politics? The Language/Action Perspective Reconsidered," in *Human Values and the Design of Computer Technology*, ed. B. Friedman (Cambridge: Cambridge University Press, 1997).
13. See, respectively, Cristina Conati and Xiaoming Zhou, "A Probabilistic Framework for Recognizing and Affecting Emotions," in *Proceedings of the AAAI Spring Symposium on Architectures for Modeling Emotions* (Stanford, CA: Stanford University Press, 2004); Fiorella de Rosiis, Catherine Pelachaud, Isabella Poggi, Valeria Carofoglio, and Berardina Nadjia De Carolis, "From Greta's Mind to Her Face: Modelling the Dynamics of Affective States in a Conversational Embodied Agent," in "Applications of Affective Computing in HCI," special issue, *International Journal of Human-Computer Studies* 59 (2003): 81–118; Gene Ball and Jack Breese, "Emotion and Personality in a Conversational Agent," in Cassell et al., *Embodied Conversational Agents*; Stacy Marsella, David V. Pynadath, and J. Stephen Read, "PsychSim: Agent-Based Modeling of Social Interactions and Influence," in *Proceedings of the Sixth International Conference on Cognitive Modeling* (Mahwah, NJ: Lawrence Erlbaum Associates, 2004); Timothy Bickmore, "Unspoken Rules of Spoken Interaction," *Communications of the ACM* 47, no. 4 (2004): 38–44; Justine Cassell and Timothy Bickmore, "Negotiated Collusion: Modeling Social Language and Its Relationship Effects in Intelligent Agents," *User Modeling and Adaptive Interfaces* 12 (2002): 1–44; Bas Strömbäck, Anton Nijholt, Paul van der Vet, and Dirk Heylen, "Designing for Friendship: Becoming Friends with Your ECA," in *Proceedings of Embodied Conversational Agents: Let's Specify and Evaluate Them!* (New York: ACM Press, 2002).
14. Sherry Turkle, *Life on the Screen: Identity in the Age of the Internet* (New York: Simon & Schuster, 1995).
15. Justine Cassell and Kristinn R. Thorisson, "The Power of a Nod and a Glance: Envelope vs. Emotional Feedback in Animated Conversational Agents," *Applied Artificial Intelligence* 13 (1999): 519–38.

16. Yukiko I. Nakano, Gabe Reinstein, Tom Stocky, and Justine Cassell, "Towards a Model of Face-to-Face Grounding," in *Proceedings of the Forty-first Annual Meeting of the Association for Computational Linguistics* (Philadelphia: Linguistics Data Consortium, 2004).
17. Herbert H. Clark, *Areas of Language Use* (Chicago: University of Chicago Press, 1992).
18. Cassell and Bickmore, "Negotiated Collusion."
19. Nass and Brave, *Wired for Speech*.
20. Justine Cassell, Andrea Tartaro, Vani Oza, Yolanda Rankin, and Candice Tse, "Virtual Peers for Literacy Learning," in "Pedagogical Agents," special issue, *Educational Technology*, forthcoming.
21. See, respectively, Kimiko Ryokai, Catherine Vaucelle, and Justine Cassell, "Virtual Peers as Partners in Storytelling and Literacy Learning," *Journal of Computer Assisted Learning* 19, no. 2 (2003): 195–208; Justine Cassell and Hannes Vilhjálmsón, "Fully Embodied Conversational Avatars: Making Communicative Behaviors Autonomous," *Autonomous Agents and Multi-Agent Systems* 2 (1999): 45–64; and Timothy Bickmore, "Towards the Design of Multimodal Interfaces for Handheld Conversational Characters," in *Proceedings of the Conference on Human Factors in Computing Systems* (New York: ACM Press, 2002).
22. Andrea Tartaro and Justine Cassell, "Using Virtual Peer Technology as an Intervention for Children with Autism," in *Towards Universal Usability: Designing Computer Interfaces for Diverse User Populations*, ed. J. Lazar (Chichester, UK: John Wiley and Sons, in press).
23. Timothy Lenoir, "All but War: Is Simulation: The Military-Entertainment Complex," *Configurations* 8, no. 3 (2000): 289–335.
24. Jessica Riskin, "Moving Anatomies" (paper presented at the annual meeting of the History of Science Society, 3–7 November 1999, at Pittsburgh, PA).
25. Allen Newell and Herbert A. Simon, *Human Problem Solving* (Oxford, UK: Prentice-Hall, 1972).
26. René Descartes, *Discours de la Méthode*, in *Oeuvres et Lettres* (1637; Paris: Librairie Gallimard, Bibliothèque de la Pléiade, 1953), 164–65; my translation.