

# React to This! How Humans Challenge Interactive Agents using Nonverbal Behaviors

Chuxuan Zhang<sup>1</sup>, Bermet Burkanova<sup>1</sup>, Lawrence H. Kim<sup>1</sup>, Lauren Yip<sup>1</sup>, Ugo Cupcic<sup>2</sup>, Stéphane Lallée<sup>2</sup>, and Angelica Lim<sup>1</sup>

How do people use their faces and bodies to test the interactive abilities of a robot? Making lively, believable agents is often seen as a goal for robots and virtual agents but believability can easily break down. In this Wizard-of-Oz (WoZ) study, we observed 1169 nonverbal interactions between 20 participants and 6 types of agents. We collected the nonverbal behaviors participants used to challenge the characters physically, emotionally, and socially. The participants interacted freely with humanoid and non-humanoid forms: a robot, a human, a penguin, a pufferfish, a banana, and a toilet. We present a human behavior codebook of 188 unique nonverbal behaviors used by humans to test the virtual characters. The insights and design strategies drawn from video observations aim to help build more interaction-aware and believable robots, especially when humans push them to their limits.

## I. INTRODUCTION

Imagine walking into a theme park and seeing an animatronic character, smiling at you. To test if it can see you, you immediately shift from left to right, and its gaze seems to follow you everywhere you go. You wave hello, and surprisingly, it waves back. You make a silly face to see if it will react to more complicated gestures. But this time the spell is broken – it doesn’t respond to you at all!

The concept of believability has long been explored in animation and artificial intelligence, from video games [1] and human-robot interaction (HRI) [2] to interactive virtual agents [3]. One important aspect of believability is interaction awareness [4]. Interaction awareness is defined as the ability of an agent that is “to perceive important structural and/or dynamic aspects of an interaction that it observes or that it is itself engaged in” [5]. Non-verbal interactions with artificial agents, like the creature described above, remain challenging to produce convincingly and appropriately [6].

Indeed, systems that produce believable, lively non-verbal interactions are rare [1] and believability can break down quickly in free interaction environments. We suggest that a cause for this breakdown is that the set of non-verbal behaviors (e.g. facial, and body gestures) that people use to test an agent’s capabilities is not clear. In addition, how

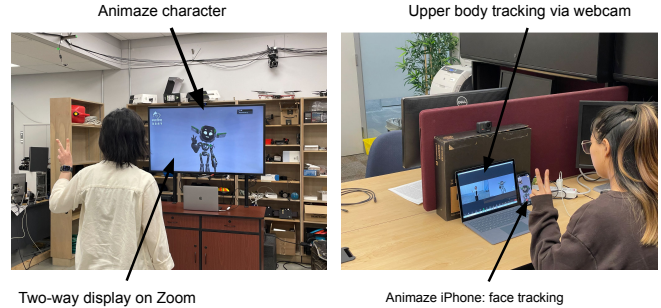


Fig. 1. Physical setup of our study: participant interacting with virtual character (left), and teleoperator using face and upper body tracking to control the virtual character (right).

people behave may depend on the character’s appearance and expected social and physical capabilities [7]; for example, a robot with no hands may be less likely to be offered a handshake. In essence, in addition to a wave or smile, what other interactive behaviors should a robot be prepared to react to?

The main contributions of this paper are:

- 1) A study of nonverbal behaviors that people use to test a character physically, emotionally, and socially
- 2) A human behavior codebook of 188 action classes that researchers can consider for recognition in HRI
- 3) Design insights related to a character’s interactive affordance

## II. RELATED WORK

This research builds upon prior work on testing for awareness along with nonverbal and affective behaviors that arise in interactions with artificial agents. We specifically focus on a subset of socially interactive agents [7] including social robots and interactive virtual agents (IVAs) that rely on visual sensing devices such as cameras. We refer to these agents in this paper as (interactive) characters.

*Interaction awareness* involves the ability to perceive dynamic aspects of an interaction [5]. Specifically, Dautenhahn et al. suggest that an “important ability of an interaction-aware agent is to *track, identify, and interpret visual interactive behavior*” [5], along the continuum of interaction formality [8]. This includes informal interactions such as play, semi-formal interactions such as greetings, and very formal interactions such as scripted law proceedings. As an

<sup>1</sup>C. Zhang, B. Burkanova, L. Kim, L. Yip, and A. Lim are with the School of Computing Science, Simon Fraser University, 8888 University Dr., Burnaby, Canada {chuxuan.zhang, bermet\_burkanova, lawkim, lauren.yip, angelica}@sfu.ca

<sup>2</sup>U. Cupcic and S. Lallée are with Spoon AI, 4 rue de la Bourse, 75002 Paris, France {ugo, stephane.lallee}@spoon.ai

example, Aldebaran Robotics<sup>1</sup> proposed the Basic Awareness module on their NAO and Pepper robots, which includes tracking detected humans and looking in the direction of stimuli such as movement or sound, towards the illusion of life.

*What do humans do when faced with an interactive agent?* Human non-verbal and affective expressions have been studied for many years, resulting in numerous hand and body gestures [9], [10] and facial expression datasets [11], [12] capturing expressions corresponding to emotional labels such as anger, happiness, excitement, sadness, frustration, fear, and surprise [13]. Decades of psychological research have studied human-human interactions, to understand body gestures [13]–[15] including emblems [16], [17], i.e. gestures that can replace speech, such as head motions for yes, no, or a shrug indicating I don’t know. Interactive virtual agents such as Greta [18], SimSensei [19] and M-PATH [20] aim to respond to human behaviors by specifically tracking features such as gaze, facial expressions and body movement to estimate the mental and emotional state of a user to empathize or change its verbal responses. The social signals created for these agents constitute a collection of behaviors relevant to this paper.

*Data-driven behavior databases.* In the human visual behavior databases mentioned above, study participants express or annotate behaviors assuming *a priori* class labels. But what if the labels are unknown? Another, more naturalistic, paradigm is to allow expressions to emerge from interactions, and then perform post-hoc labeling. As an example, in the Tower Game [21] experiment, participants were asked to build a tower out of blocks and were not allowed to talk, resulting in only non-verbal interactions and expressions. Gestures such as kiss and peek-a-boo have also been uncovered in data-driven analyses of physical and social interactions specifically with infants, and have limited verbal interaction capabilities [22]. One of the most related studies to our work explores abuse of robots by children, analyzing 12 hours of behavior in a shopping mall and describing a handful of bullying behaviors such as grabbing, pulling, blocking, and shoving robots [23].

Overall, despite the many studies studying nonverbal social and affective behaviors between humans and artificial agents, a) many of them focus on modeling the agent’s behavior rather than the human’s, and b) datasets annotating human interactive behavior still lack information regarding how people challenge virtual agent and machine interactive abilities. In order to design interactive agents that react appropriately to human behavior, it is imperative to first understand what exactly the agents should react to.

### III. METHODOLOGY

#### A. Research Questions

Our long-term goal is to improve a character’s interaction awareness. In this study, we address a specific knowledge gap in the context of a *one-to-one, non-verbal interaction with an*



Fig. 2. Participants were asked to test 6 interactive virtual characters physically, emotionally, and socially.

*interactive agent* (Fig. 1, left). We investigate the research question (RQ): **What do people do to test a character physically, emotionally, and socially?**

#### B. Setup and Materials

We prepared a Wizard-of-Oz(WoZ) experimental setup using three off-the-shelf software packages to create interactions between human participants and virtual characters (Fig. 1). We used Animaze<sup>2</sup>, a software that provides ready-to-use animated characters of various morphologies, accompanied by its built-in iPhone face tracking module and Webcam Motion Capture<sup>3</sup> to track the teleoperator’s upper body including head, shoulders, and finger movements. Zoom<sup>4</sup> was used for this study to allow the teleoperator to see the participants’ behaviors and vice versa. On the participant side, the virtual character was displayed on a large TV screen (55 inches) oriented horizontally, slightly above eye level. Participants were asked to stand at a marked location approximately 1.5m from the display.

#### C. Pilot Study

Five participants were recruited to refine the study design. In this pilot study, we used the physical setup depicted in Fig. 1. Each of the 5 participants interacted with 11 different agents. The task given to the participants interacting with the virtual character was: “These are free interactions, and you can do whatever you want.” To keep all participants in the main study receiving the same treatment, we designated a single teleoperator and wrote a concise interaction guideline for the teleoperator. The teleoperator was told not to intentionally start any interaction, but only react to participants’ actions, so that we could observe interactions that are only initiated by participants.

We observed 3 major behavioral categories in the pilot: 1) *Posture, proxemics and physical contact*: Pilot participants sometimes adjusted their *posture* [24] (e.g. tilted their head), and moved their bodies subtly (e.g. wiggled body) during the interaction. Participants also changed their relative distance

<sup>2</sup><https://www.animaze.us/>

<sup>3</sup><https://webcammotioncapture.info/>

<sup>4</sup><https://zoom.us/>

<sup>1</sup><https://www.aldebaran.com/>

to the character, for example by walking towards the character, i.e. *proxemics* [25]. Finally, we encountered faked *physical contact* such as poking the character, similar to haptic interactions without touching the screen. 2) *Affect Displays*: Some pilot participants showed frowning eyebrows when seeing the characters' unexpected reactions, and smiled when being amused by the characters. 3) *Emblematic Gestures*: We saw that some pilot participants tried to compliment the characters by giving thumbs up and raising one hand to ask for a high-five.

As a result of the pilot study, we made 3 changes to the study design. First, we reduced the number of characters to interact with from 11 to 6 due to the boredom and tiredness reported by the pilot participants. The six characters (Fig. 2) were chosen to cover a broad range of morphologies: 1) a pair of humanoid characters (a robot and a human), 2) a pair of animate non-humanoid characters (a fish and a penguin), and 3) a pair of traditionally inanimate objects (a banana and a toilet). The last category was included to reflect traditionally inanimate interactive agents studied in HRI such as a donation box [26] or a moving desk [27]. The peach, cat, red panda, and shark were removed to reduce duplication with the banana, penguin, and fish, and the bacteria character was also omitted as we did not believe it would be as popular an interactive agent. Secondly, we modified our prompt, because the participants reported a lack of motivation to initiate interactions when given the original open-ended task. The details are described in Sec. III-D. Finally, we added additional rules to the teleoperator guidelines based on the scenarios observed; the full guide is included in the Appendix.

#### D. Study Design

Following our pilot study, we conducted a WoZ study with 20 adults (gender: 9/10 women/men, 1 prefer not to disclose; age:  $28.5 \pm 12.98$ ) participants interacting with 6 different virtual characters for 1 minute each. The study was approved by the university ethics board.

We changed our task prompt to “Your goal is to test what the character can and cannot do physically, emotionally, and socially” to encourage participants to initiate interaction. This prompt was built upon observed behaviors in the pilot study: posture, proxemics, and physical contact as “*physical behaviors*”, affect displays as “*emotional behaviors*”, and emblematic gestures as “*social behaviors*”. The on-site researcher verbally provided each participant with two non-verbal interactive behaviors as examples, picked randomly from the following list:

- waving to others indicates you are greeting someone
- nodding to someone at the other end of the hallway to show you noticed them from afar
- passing someone your smartphone to show you are trying to share something with them
- you can test the characters' emotional capability by testing if they can tell and react to you when you are sad/happy/angry
- pretending to have physical contact with the characters

In the pre-study briefing, we obtained participants' consent to be video recorded, their permission to release the fully anonymized video data for research purposes, and the identifiable video data for conference presentation. We suspected that people might act differently towards a teleoperated agent from a fully autonomous one, therefore participants were told that all characters were fully autonomous. Participants were also told that the agents were not programmed to process any audio data. Participants were nonetheless allowed to make sounds and speak if they thought this would help them communicate more effectively. During each study session, the participant was instructed to have a 1-minute interaction with each of the 6 virtual characters, following the prompt above. To alleviate the order effect, we randomized the interaction order of the virtual characters, and the researcher turned away from the participant to prevent the participant from feeling observed. At the end of the study, participants were debriefed that the agents were teleoperated, and were asked if they wanted to revise their consent for data release after being debriefed.

#### E. Data Collection and Analysis

We used ELAN [28] by the Max Planck Institute to annotate our video data. Each video was processed by at least two annotators. The primary annotator segmented and annotated the video, and the secondary annotator went over the existing annotation and took notes of potential disagreement. We adopted the discuss-until-consensus method that is often used in qualitative studies [29]. A third annotator (the teleoperator) was consulted to break disagreements.

1) *Segmentation*: The primary annotator segmented the entire recorded video into interactive behavior segments. One interactive behavior segment is defined as a complete action-reaction (response) pair, containing (a) an action initiated by either the virtual agent or the participant, and (b) its corresponding response by the other interactor.

2) *Segment classification*: Every segment was classified into one of: physical (posture, proxemics or physical contact), emotional (affect display), or social (emblematic) behaviors. The annotator first determined whether the behavior had a straightforward and interpretable verbal meaning, and if so classified it as a *social behavior*. If not, the annotator examined whether there was a clear emotion expressed via the behavior, and if so, tagged it as an *emotional behavior*. Finally, if the behavior did not fall into either social or emotional behavior categories, the annotator classified it as a *physical behavior*. Thus, this last category contained postural, proxemic, and physical contact behaviors, along with any other behaviors that did not fall into the other two categories. Our method provides us with a hierarchy of behavior understanding, where social behaviors may, but are not required to, be composed of emotional displays, which in turn can be comprised of physical actions. In this way, classifications might not be mutually exclusive from one other. For example, when participants tried to attack the agents with violent actions, they sometimes frowned hard

to pretend they were angry. In this case, both emblematic meanings and distinct emotions were annotated.

3) *Additional labeling*: In addition, for each segment, the annotators decided on the following:

- *Initiator of the interaction* (character/participant): the initiator of the interaction
- *Character-specific* (yes/no): separates character-agnostic and character-specific behavior by determining if the motive of the behavior is only explainable with the existence of the character’s unique features
- *Description of physical actions* (free text): short, physical movement descriptions for the instance
- *Emotion label* (free text): the emotion participants tried to express during an interaction instance, if any
- *Social meaning* (free text): the social meaning behind an interaction instance, if any
- *Response type* (reacting/mirroring/no response): the type of the response from the interaction recipient
- *Response description* (free text): short physical/emotional/social descriptions on the interaction recipient’s response action, if any

The full set of annotations and anonymized videos can be downloaded from <https://rosielab.github.io/react-to-this/>.

#### F. Thematic Analysis

After video annotation, we set aside the behaviors labeled as character-specific and proceeded with thematic analysis on the remaining (character-agnostic) behaviors. We conducted the thematic analysis using the physical, emotional, and social grouping paradigm derived from the pilot study to form our human behavior codebook:

- 1) **Social**: For all the social behaviors (interaction labeled as social/emblematic), we listed all the social meaning annotations and merged those with similar high-level meanings. For example, participants expressed a disapproving attitude with various behaviors such as shaking the head, shaking the finger, thumb(s) down, and forearms crossed as an ’’X’ figure, etc. Thus we extracted a high-level and abstract meaning *disapproval* to summarize this group of behaviors conveying a similar social meaning.
- 2) **Emotional**: We grouped the freely annotated emotion labels based on Ekman’s basic emotions, namely anger, disgust, fear, happiness, sadness, and surprise for all the emotional behaviors. For the freely annotated emotion labels that did not seem to fit into the basic emotion categories, we created a separate category on its own.
- 3) **Physical**: We first decided if participants changed their relative spatial distance to the character (proxemics) or if they intended to have physical contact with the character (physical contact). Then, we separated the physical behaviors excluded from the above 2 categories by the activated body parts into 5 groups (full body, head/face, arm, hand, lower body). We observed

sequences of movements involving multiple body parts in one physical interactive behavior. For example, participant 19 (P19) jumped, drew a big circle with their arms, widened their eyes, and opened their mouth quickly and continuously in one interaction segment. We counted each individual action and classified each of them into different physical behavior categories to form a detailed thematic result (see Table III).

For character-specific behaviors, we grouped them based on the character conditions, due to the small number of instances (Table I).

## IV. RESULTS AND DISCUSSION

In total, 1169 interactions were observed during the study, including 1111 interactions initiated by participants and 58 interactions initiated by the agent. From the 1111 interactions started by the participants, we obtained 999 character-agnostic (424 physical instances, 162 emotional instances, 413 social instances), and 112 character-specific behaviors. In this section, we provide results and discussions for each category of behaviors. In addition, we discuss two other behavioral patterns—mimicry and compound behaviors—that emerged from the analysis.

*Physical Behaviors* are presented in Fig. 3 and Table III, and include posture changes, proxemics and physical contact. Firstly, it appeared that posture changes were related to actuation testing, i.e. participants initiated some interactive behaviors (e.g. *leaning* or *raising arms*) to see if the character was able to follow and replicate. Secondly, in some interactions, it appeared characters were not expected to copy the participant’s exact movement. Rather, participants were checking whether it could track them (e.g. *walking to the left/right*). We believe this points to a class of behaviors for sensor testing. For example, P113 put their hoodie in front of the robot character to block its view. Finally, in some cases, participants expected physical feedback as a result of their pretend physical contact. Actions involving physical contact contributed to these types of behaviors by testing physics-based responses, such as lifting, pushing, and poking. These are reminiscent of the mild robot abuse behaviors found in [23], although we did not necessarily observe aggression (see, however, the *attacking* behavior in Table V).

*Emotional Behaviors* were conveyed through facial expressions and/or body movements. We found two interesting outcomes: 1) Exaggeration. Participants tried to express sadness by finger-drawing a downward curve in front of their face, moving index fingers from the eye’s corner to the chin to depict tear dropping, and rubbing fists around their cheeks to pretend to be crying hard. The occurrence of such pantomime actions might indicate a low expectation of the intelligent virtual characters’ capability in emotion recognition. 2) Diversity. We found that one emotion category could comprise multiple unique behaviors. While not surprising based on the psychology literature [30], these actions could serve as alternate targets for recognition systems, i.e. instead of attempting to recognize ’’sadness’’, which is an overarching concept comprised of heterogeneous actions, computer vision

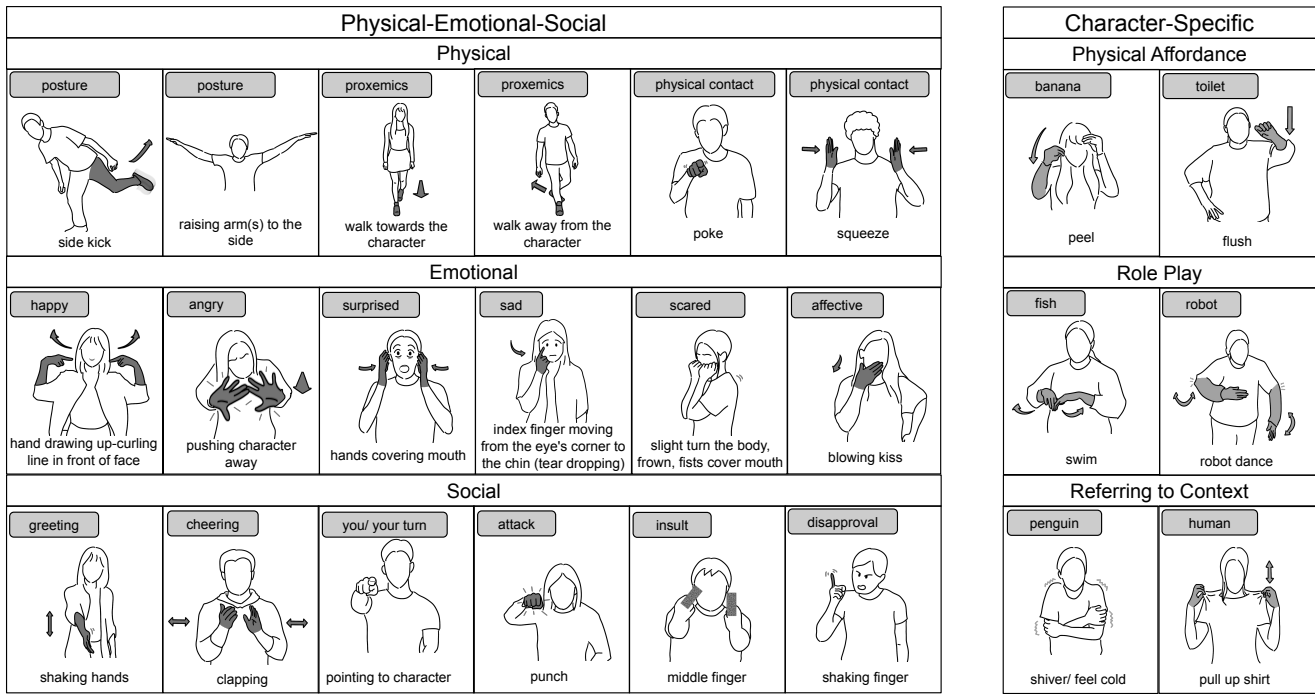


Fig. 3. (LEFT) Illustrative examples for physical, emotional, and social nonverbal human behaviors. Physical behaviors included changing posture, proxemic movement, and mock physical contact. (RIGHT) Examples of character-specific interactions related to physical affordance (banana peeling, toilet flushing) role play (fish swimming, robot dancing), and the agent’s contextual environment and accessories (shivering from cold, pulling on shirt).

researchers could aim to recognize “rubbing eyes, pretending to cry”, “lowered lip corners” and so on. The full list of 34 unique emotional behaviors and 9 categories can be found in Table IV.

*Social Behaviors* were those with clearly identifiable verbal meanings such as greeting, cheering, attacking, insulting, etc. The detailed explanation can be found in Sec. III-C under *emblematic gesture*. We observed 82 unique behaviors and derived 42 different social behavior categories from the annotations (Table V). Two main observations were: 1) Culture-based behaviors. In addition to well-recognized behaviors such as thumbs up for “good job”, there were also culture-specific behaviors such as namaste (pressing hands together, and fingers pointing upwards). This highlights the importance of investigating culture-specific social signals. 2) Diversity. Similar to the observed emotional behaviors, we found that different behaviors might share similar social meanings. For instance, punching, kicking, and firing handguns can all be used to express one’s aggression toward others.

*Character-Specific Behaviors* (Fig. 3 and Table I) were defined as the motive of the behaviors are only explainable with the existence of the character’s unique features during our annotation process. For example, “flushing the toilet” and “pulling the toilet paper” behaviors were only observed during interaction with the toilet, and “peeling” and “eating” actions were only observed during the banana interaction session. This might suggest that the physical traits of the character may afford specific physical interactions, and may need to be considered when designing a robot (Fig. 3, right-

Character	Nonverbal Character-Specific Behaviors
banana	grabbing; peeling; eating; rotating command by pointing down from above head; bending/breaking; placing the character on hand
fish	puffing out cheeks; swimming; flipping hands; rolling over command; juggling command; bouncing command
toilet	opening and closing lid; flushing; pulling toilet paper; sitting on
penguin	waddling; complimenting the scarf; taking off the scarf; shivering (pretending to be cold)
robot	robotic arm movement/robot dance; flapping ears
human	touching hair (e.g. twirling, combing, wearing); complimenting hair; flexing arm; tugging sleeve; grabbing shirt

TABLE I  
CHARACTER SPECIFIC INTERACTIVE BEHAVIORS.

top). Next, we observed participants “waddling” with the penguin, “swimming” with the fish, and “robot dance” with the robot (Fig. 3, right-middle). It appeared that participants engaged in role-play with some characters, similar to informal “play” interactions noted by Dautenhahn [5]. Imagining behaviors related to the agent’s character may help to predict these types of behaviors (e.g., if designing a robot lion, consider that humans may engage in pantomime roaring). Finally, participants made reference to the imagined environment such as “shivering” from cold with the penguin, and agent accessories, such as complimenting a shirt or scarf (Fig. 3, right-bottom). This could suggest that humans may be testing the environmental awareness of the agent,

Category	Mimicked Behaviors
Physical (20)	swinging (5); raise hands (2); wiggling (2); hand covering mouth (1); shrug (1); wave hands (1); shake head; tilt head (1); lean to the side (1); bend body (1); reaching out arms (1)
Emotional (5)	smile – happy (2); hug – affectionate (2); pouting mouth – sadness (1); shocked face, raise two hands, lean back – surprised (1); frowning – angry (1)
Social (16)	shrug (7); ok gesture (6); thumbs up (1); boxing (1); shaking head (1)

TABLE II

NONVERBAL BEHAVIORS THAT WERE INITIATED BY THE AGENT AND MIMICKED BY THE PARTICIPANTS.

“pointing and referring to areas of and things in it” [31], suggested to be a component of agent believability in video game contexts. Designers may need to consider this when adding accessories to their character.

*Mimicry* [32] was also observed in this study. From the 57 character-initiated behaviors, we observed that 17 participants mimicked the characters’ behavior or derived the next interactive behavior from it (Table II). Although the teleoperator would not initiate interaction on purpose, the presence of delay can interrupt interaction flow; for example, participants might treat a delayed response from an agent as the start of a new interaction. In some of these cases, we observed mimicry behaviors from participants.

It appeared that the mimicry was at times involuntary and unconscious. For example, participant P101 claimed that because he saw that the human character had hair, they actively chose to test if the character was able to groom her hair as he could. However, the recorded video showed that the character touched her hair first, then the participant immediately performed the same action. In another example, P103 reported that after he knocked himself on the head, the banana character appeared to fall asleep (pretending to pass out after being hit), so the participant tried to sleep to examine if the banana actually understood this action. When participants ran out of interaction ideas, consciously mimicking the character’s response was one of the methods they used to create more interactive behaviors. As a cautionary guideline for designers, this could suggest that robots may need to recognize all the behaviors that it expresses.

*Compound Behaviors* consisting of sequential and multimodal actions were also observed. For example, (Fig. 4, top) P121 drew the shape of the scarf first, then made a “thumbs up” gesture to give a compliment. Thus, being able to segment action sequences and understand them correctly is crucial for the virtual character to respond appropriately. As an example of a multimodal social signal (Fig. 4 bottom), the “shrug” action was seen as a signal of “I don’t know/understand” when it was in the company of raised eyebrows, and was perceived as a sign of “tada!” when the participants held the shrug for a long time and displayed a surprised-happy face. Sequential and multimodal actions make segmenting and recognizing behavior a challenging problem for future work.

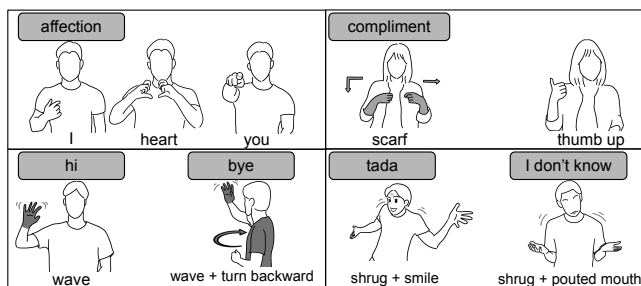


Fig. 4. Compound Behavior Examples: behaviors consisting of sequential actions (top) and behaviors consisting of multimodal actions (bottom).

## V. LIMITATIONS AND FUTURE WORK

Our research was conducted in a North American post-secondary institution. Thus, our findings may not be representative of the general public from different cultural backgrounds. Future research may target a larger and more diverse sample (including people who are of different ages, who are neurodiverse, or who are from varied cultural groups). Also, virtual agents allowed for flexible testing of various morphologies, but physical interaction with real robots may differ. For example, Kanda et al. found that children try to obstruct the path of a navigating robot [23]. We also did not investigate auditory nonverbal behaviors, such as clapping or non-linguistic utterances. Finally, we do not claim that our list of interactive behaviors is exhaustive, but rather contributes to the set of nonverbal human behaviors under scrutiny by HRI researchers. The next step of this research involves creating algorithms to recognize these behaviors, and investigating the effect of various agent reactions to them. We hope that this set of behaviors can be applied to robots and interactive agents deployed in various settings (e.g. theme parks, education, video games) to increase believability in the future.

## VI. CONCLUSION

In this study, we discovered 188 unique actions and 51 socio-emotional behavior categories among 1169 non-verbal interactions between participants and 6 virtual characters, contributing to the list of target classes for visual gesture recognition algorithms. With a bottom-up analysis method, we created a rich and diverse behavior codebook to guide designers and programmers of interactive agents/robots. The 188 actions and corresponding meanings also could help provide a list of classes for machine learning gesture recognition algorithms to target. The set of abundant interactive behaviors in our codebook can be applied to interactive agents deployed in various settings (e.g. video games, theme park, education) in the future.

## REFERENCES

- [1] C. Curtis, S. O. Adalgeirsson, H. S. Ciurda, P. McDermott, J. Velásquez, W. B. Knox, A. Martinez, D. Gaztelumendi, N. A. Goussies, T. Liu, et al., “Toward believable acting for autonomous animated characters,” in *SIGGRAPH Conference on Motion, Interaction and Games*, 2022, pp. 1–15.

Category (Occurrence count)	Nonverbal Physical Behaviors
posture – full body (136)	turning (36); tilting/leaning body (29); bending the torso sideways (27); jumping (16); body forward or backward (8); spinning (4); swaying/swing (4); bending the collapsed body pose (4); walking/running/jogging (in place) (3); wiggling (1); kneeling on the ground (1); jumping jack (1); stretch out legs (1); rotating upper body (1)
posture – head/face (118)	tilting head (23); open/closed mouth (17); pouting mouth (12); nodding (8); raise eyebrows (7); frowning (6); winking (6); shaking head (5); looking at some direction (5); stick out tongue (5); staring (5); squinting (4); turning head to the side (4); crooked mouth (3); closing one eye (2); rolling eyes (2); sucken cheeks (2); stick out head (1); lips touching nose (1)
posture – arm (68)	raising arm(s) (17); stretch arm(s) out (15); wave arms/hands (13); arms/hands drawing a circle (9); crossed arms (7); arms/hands flapping (2), flipping and rotating wrist (2); open arms (1); rotating forearm(s) around the elbow(s) (1); bending arm (1)
posture – hand (42)	hand(s) touching other body part(s) (23); scratching other body parts (5); clap (3); moving fingers (3); palms together (2); raising hand(s) (2); putting on hood (1); flicking hand (1); hand clasping (1); showing finger(s) (1)
posture – lower body (21)	squatting (9); lifting/raising leg(s) up (7); side kick (1); lifting leg(s) to the side (1); stretching out legs (1); standing on toes (1); shaking knees (1)
proxemics (37)	walk to the left/right (13); walking away from the character (9); running(5); walking toward the character(4); stepping forward/backward/to the side (4); making big steps(1); walking around(1)
physical contact (13)	push character with hands (4); poking with index finger(s), squeeze character by pinching index finger and thumb (3); grab the character (3); pick up gesture with both hands, put aside (1); lift character up by grabbing and lifting motion (1); squeeze character with palms (1)

TABLE III

PHYSICAL BEHAVIOR CODEBOOK: 73 NONVERBAL HUMAN BEHAVIORS INITIATED BY THE PARTICIPANTS TO TEST THE AGENT’S PHYSICAL ABILITIES, GROUPED INTO 7 CATEGORIES. ‘,’ SEPARATES DIFFERENT BEHAVIORS.

Category (Occurrence count)	Nonverbal Emotional Behaviors
angry, annoyed, sullen, stern, aggressive, menacing, unsatisfied (45)	frown(23); pouting mouth (10); stare(6); shake head (3); shaking finger (3); pushing character away (2)
affectionate (44)	heart gestures (fingers, hands, arms) (22); hugging (15); caressing (1); petting (2); blowing kiss (2); kissing (1); hands overlap, rest hands on chest (aw gesture) (1)
happy (42)	smile (25); laugh (11); hand drawing an up-curved line in front of face (3); pulling the corners of mouth up (2); giggle (1)
surprised, shocked (23)	open mouth (11); <b>widen eyes, mouth open</b> (6); widen eyes/raised eyebrows (4); hands cover mouth (1); fingers spread out around eyes (1)
sad (16)	pouted mouth (7); rubbing eyes (pretending to cry) (3); index finger moving from the eye’s corner to the chin (tear dropping) (3); pulling the corner of mouth down with fingers (2); bowed head (1)
tired (5)	<b>yawning (3); stretching (1); sighing (1)</b>
scared (1)	<b>slightly turning the body, frowning, fists covering mouth (1)</b>
shy (1)	<b>hands on face, turning away (1)</b>
contempt (1)	<b>turning head to the side, looking down, chin up, side eye (1)</b>

TABLE IV

EMOTIONAL BEHAVIOR CODEBOOK: 34 NONVERBAL HUMAN BEHAVIORS INITIATED BY THE PARTICIPANTS TO TEST THE AGENT’S EMOTIONAL ABILITIES, GROUPED INTO 9 CATEGORIES. ‘,’ SPLITS ONE BEHAVIOR INTO SMALLER UNITS TO INCREASE THE CLARITY OF THE ACTION DESCRIPTIONS. ‘,’ SEPARATES DIFFERENT BEHAVIORS. COMPOUND BEHAVIORS ARE LABELED IN BOLD.

- [2] T. Ribeiro and A. Paiva, “The illusion of robotic life: principles and practices of animation for robots,” in *international conference on Human-Robot Interaction*, 2012, pp. 383–390.
- [3] C. Pelachaud and M. Bilvi, “Computational model of believable conversational agents,” *Communication in multiagent systems: Agent communication languages and conversation policies*, pp. 300–317, 2003.
- [4] A. Bogdanovych, T. Trescak, and S. Simoff, “What makes virtual agents believable?” *Connection Science*, vol. 28, no. 1, pp. 83–108, 2016.
- [5] K. Dautenhahn, B. Ogden, and T. Quick, “From embodied to socially embedded agents—implications for interaction-aware robots,” *Cognitive Systems Research*, vol. 3, no. 3, pp. 397–428, 2002.
- [6] I. Wang and J. Ruiz, “Examining the use of nonverbal communication in virtual agents,” *International Journal of Human-Computer Interaction*, vol. 37, no. 17, pp. 1648–1673, 2021.
- [7] B. Lugrin, C. Pelachaud, and D. Traum, Eds., *The Handbook on Socially Interactive Agents: 20 Years of Research on Embodied Conversational Agents, Intelligent Virtual Agents, and Social Robotics Volume 1: Methods, Behavior, Cognition*, 1st ed. New York, NY, USA: Association for Computing Machinery, 2021, vol. 37.
- [8] I. Hutchby and R. Wooffitt, *Conversation analysis*. Polity, 2008.
- [9] Y. Zhang, C. Cao, J. Cheng, and H. Lu, “Egogesture: A new dataset and benchmark for egocentric hand gesture recognition,” *IEEE Transactions on Multimedia*, vol. 20, no. 5, pp. 1038–1050, 2018.
- [10] Y. Luo, J. Ye, R. B. Adams, J. Li, M. G. Newman, and J. Z. Wang, “Arbee: Towards automated recognition of bodily expression of emotion in the wild,” *International journal of computer vision*, vol. 128, pp. 1–25, 2020.
- [11] D. Kollias, “Abaw: learning from synthetic data & multi-task learning challenges,” in *the European Conference on Computer Vision*. Springer, 2023, pp. 157–172.
- [12] X. Liu, H. Shi, H. Chen, Z. Yu, X. Li, and G. Zhao, “imigue: An identity-free video dataset for micro-gesture understanding and emotion analysis,” in *conference on computer vision and pattern recognition*, 2021, pp. 10 631–10 642.
- [13] C. Busso, M. Bulut, C.-C. Lee, A. Kazemzadeh, E. Mower, S. Kim, J. N. Chang, S. Lee, and S. S. Narayanan, “Iemocap: Interactive emotional dyadic motion capture database,” *Language resources and evaluation*, vol. 42, pp. 335–359, 2008.
- [14] E. Volkova, S. De La Rosa, H. H. Bülthoff, and B. Mohler, “The mpi emotional body expressions database for narrative scenarios,” *PLoS*

Category (Occurrence count)	Nonverbal Social Behaviors
greeting (125)	waving hand (107); shaking hands (5); bow (5); salute (3); raise one hand, moving fingers (1); chin up and down, raise eyebrows quickly (1); two-finger salute (1); upwards nod (1); palms together, fingers facing upward (1)
good job (42)	thumb(s) up (42)
approval (29)	nodding (18); ok gesture (11)
disapproval (26)	shaking head (10); shaking finger (5); frowning (5); thumb(s) down (5); forearms crossed as "X" (1)
you/ your turn (19)	pointing to character (19)
dance(19)	causal dancing (11); dab (4); whip (3); nae nae (1)
attacking (19)	boxing/punch (12); fire handgun (3); slap (2); kick (1); hit oneself (1)
entertaining (12)	peekaboo (5); making face (3); tada (2); wiggling fingers behind the head (1); flamingo pose (1)
I don't know (11)	shrug (11)
cheering (11)	clapping (7); high five (4)
deictic(11)	point to some direction for the character to look at/follow (11)
questioning (11)	shrug (7); tilt head (2); raise eyebrows, 'asking' face (1); horizontally wave the hand (1)
instruct to copy the exact behavior (10)	perform certain action and point to the character to instruct the character to replicate the same action (10)
I (8)	pointing to themselves (5); hand(s) rest on chest (3)
come closer (7)	pull hands to oneself (7)
thinking (7)	<b>index finger over mouth, serious face</b> (2); fist under chin (2); <b>crossed arms, bite lips, nod</b> (1); hand in chin (1); <b>arm crossed, serious face, tilt head</b> (1)
insult (7)	middle finger (7)
sleep (7)	<b>closed eyes, tilting head, rest head on hands</b> (7)
goodbye (5)	<b>walk away/turn back, waving hand</b> (5)
come with me (4)	<b>pull hands quickly towards oneself</b> (3); point to the back (1)
searching (4)	hands over eyes (4)
re-draw attention (4)	wave hand (when characters are facing to the side) (2); finger snapping (1); turn 180 degree, then suddenly turn back (1)
peace (3)	victory gesture (3)
look cute/pretty (2)	hands under chin (aw face)/hands under chin(bare teeth, smile) (2)
taunting (2)	<b>point at the character, leaning back, laughing</b> (1); raising both hands, pointing to the character (1)
holding hand (2)	<b>reaching out one arm, palm facing up, pointing to the reached-out hand</b> (1); hold both hands(1)
interact with a smartphone (2)	pull out a smartphone, pretend to take a picture of the character (1); show the character the smartphone screen (1)
talking (2)	<b>hands out moving, pretend to talk</b> (1); hands out (1)
cut it off/stop it (1)	whip hands (1)
identity revealing (1)	pulling down and up the hood on clothes (1)
eyes on you (1)	pointing fingers to their own eyes then to the character(1)
broken heart (1)	hands making a heart gesture and separate hands (1)
sick (1)	sneeze (1)
call me (1)	hand gesture as a phone, rest the hand next to the ear (1)
whatever (1)	shrug (1)
whispering (1)	<b>hand(s) closed to mouth, moving lips</b> (1)
listening (1)	<b>turning body 90 degrees, putting hand close to the ear</b> (1)
numbers (1)	using fingers to indicate some number (1)
reading (1)	look at one palm (1)
write (1)	index finger of one hand hovering over another hand (1)
comfort/calm (1)	put two hands up, smile (1)
so-so (1)	wave the hand horizontally (1)

TABLE V

SOCIAL BEHAVIOR CODEBOOK: 82 NONVERBAL HUMAN BEHAVIORS INITIATED BY THE PARTICIPANTS TO TEST THE AGENT'S SOCIAL ABILITIES, GROUPED INTO 42 CATEGORIES. ', ' SPLITS ONE BEHAVIOR INTO SMALLER UNITS TO INCREASE THE CLARITY OF THE BEHAVIOR DESCRIPTIONS. ', ' SEPARATES DIFFERENT BEHAVIORS. COMPOUND BEHAVIORS ARE LABELED IN BOLD.

- one, vol. 9, no. 12, p. e113647, 2014.
- [15] D. Matsumoto and H. C. Hwang, "Cultural similarities and differences in emblematic gestures," *Journal of Nonverbal Behavior*, vol. 37, pp. 1–27, 2013.
- [16] H. J. Ottenheimer and J. M. Pine, *The anthropology of language: An introduction to linguistic anthropology*. Cengage Learning, 2018.
- [17] D. Matsumoto and H.-S. Hwang, "Emblematic gestures (emblems)," *The encyclopedia of cross-cultural psychology*, vol. 2, pp. 464–466, 2013.
- [18] R. Niewiadomski, E. Bevacqua, M. Mancini, and C. Pelachaud, "Greta: An interactive expressive eca system," in *International Conference on Autonomous Agents and Multiagent Systems - Volume 2*, ser. AAMAS '09. Richland, SC: International Foundation for Autonomous Agents and Multiagent Systems, 2009, p. 1399–1400.
- [19] D. DeVault, R. Artstein, G. Benn, T. Dey, E. Fast, A. Gainer, K. Georgila, J. Gratch, A. Hartholt, M. Lhommet, G. Lucas, S. Marsella, F. Morbini, A. Nazarian, S. Scherer, G. Stratou, A. Suri, D. Traum, R. Wood, Y. Xu, A. Rizzo, and L.-P. Morency, "Simsensei kiosk: A virtual human interviewer for healthcare decision support," in *the International Conference on Autonomous Agents and Multi-Agent Systems*, ser. AAMAS '14. Richland, SC: International Foundation for Autonomous Agents and Multiagent Systems, 2014, p. 1061–1068.
- [20] Ö. N. Yalçın and S. DiPaola, "M-path: A conversational system for the empathic virtual agent," in *Biologically Inspired Cognitive Architectures 2019: Annual Meeting of the BICA Society 10*. Springer, 2020, pp. 597–607.
- [21] D. A. Salter, A. Tamrakar, B. Siddiquie, M. R. Amer, A. Divakaran, B. Lande, and D. Mehri, "The tower game dataset: A multimodal dataset for analyzing social interaction predicates," in *International Conference on Affective Computing and Intelligent Interaction (ACII)*. IEEE, 2015, pp. 656–662.
- [22] S. E. Colgan, E. Lanter, C. McComish, L. R. Watson, E. R. Crais, and



- G. T. Baranek, "Analysis of social interaction gestures in infants with autism," *Child Neuropsychology*, vol. 12, no. 4-5, pp. 307-319, 2006.
- [23] S. Yamada, T. Kanda, and K. Tomita, "An escalating model of children's robot abuse," in *International Conference on Human-Robot Interaction*, 2020, pp. 191-199.
- [24] B. Howarth, "Dynamic posture," *Journal of the American Medical Association*, vol. 131, no. 17, pp. 1398-1404, 1946.
- [25] E. T. Hall, "A system for the notation of proxemic behavior," *American anthropologist*, vol. 65, no. 5, pp. 1003-1026, 1963.
- [26] R. H. Kim, Y. Moon, J. J. Choi, and S. S. Kwak, "The effect of robot appearance types on motivating donation," in *Proceedings of the 2014 ACM/IEEE international conference on Human-robot interaction*, 2014, pp. 210-211.
- [27] L. H. Kim, A. A. Leon, G. Sankararaman, B. M. Jones, G. Saha, A. Spyropoulos, A. Motani, M. L. Mauriello, and P. E. Paredes, "The haunted desk: exploring non-volitional behavior change with everyday robotics," in *Companion of the 2021 ACM/IEEE International Conference on Human-Robot Interaction*, 2021, pp. 71-75.
- [28] T. L. A. Nijmegen: Max Planck Institute for Psycholinguistics, "Elan." [Online]. Available: <https://archive.mpi.nl/tla/elan>
- [29] T. Harding, D. Whitehead, *et al.*, "Analysing data in qualitative research," *Nursing & midwifery research: Methods and appraisal for evidence-based practice*, vol. 5, pp. 141-160, 2013.
- [30] L. F. Barrett, R. Adolphs, S. Marsella, A. M. Martinez, and S. D. Pollak, "Emotional expressions reconsidered: Challenges to inferring emotion from human facial movements," *Psychological science in the public interest*, vol. 20, no. 1, pp. 1-68, 2019.
- [31] K. Ijaz, A. Bogdanovych, and S. Simo, "Enhancing the believability of embodied conversational agents through environment-, self- and interaction-awareness," in *Conferences in Research and Practice in Information Technology Series*, 2011.
- [32] T. L. Chartrand and J. A. Bargh, "The chameleon effect: The perception-behavior link and social interaction." *Journal of personality and social psychology*, vol. 76, no. 6, p. 893, 1999.

## APPENDIX

### A. Teleoperator Guidelines

The teleoperator followed the following guidelines during interactions with the participants. For each participant action (left of the arrow), the virtual character should show a corresponding response (right of the arrow):

- Greeting → mirroring
- Pointing directions → turning/leaning/looking/pointing at/towards the direction
- Physical contact → pretend the physical contact is happening
- Show aggressiveness/ disgust face → show hurt/angry
- For an expression that you cannot do with the character → use the character to make the expression 'I cannot do that'
- If no interaction performed → looking around (idle mode)