

## The Effects of Robot Voices and Appearances on Users' Emotion Recognition and Subjective Perception

Sangjin Ko\*, Jaclyn Barnes<sup>†</sup>, Jiayuan Dong\*, Chung Hyuk Park<sup>‡</sup>,  
Ayanna Howard<sup>§</sup> and Myounghoon Jeon<sup>\*,†,¶</sup>

*\*Grado Department of Industrial and Systems Engineering  
Virginia Polytechnic Institute and State University  
Blacksburg, VA, USA*

*†Department of Computer Science  
Michigan Technological University  
Houghton, MI, USA*

*‡Department of Biomedical Engineering  
Department of Computer Science  
The George Washington University  
Washington DC, USA*

*§Department of Electrical and Computer Engineering  
The Ohio State University  
Columbus, OH, USA*

*¶myounghoonjeon@vt.edu*

Received 19 March 2022

Revised 1 December 2022

Accepted 19 December 2022

Published 22 February 2023

As the influence of social robots in people's daily lives grows, research on understanding people's perception of robots including sociability, trust, acceptance, and preference becomes more pervasive. Research has considered visual, vocal, or tactile cues to express robots' emotions, whereas little research has provided a holistic view in examining the interactions among different factors influencing emotion perception. We investigated multiple facets of user perception on robots during a conversational task by varying the robots' voice types, appearances, and emotions. In our experiment, 20 participants interacted with two robots having four different voice types. While participants were reading fairy tales to the robot, the robot gave vocal feedback with seven emotions and the participants evaluated the robot's profiles through post surveys. The results indicate that (1) the accuracy of emotion perception differed depending on presented emotions, (2) a regular human voice showed higher user preferences and naturalness, (3) but a characterized voice was more appropriate for expressing emotions with significantly higher accuracy in emotion perception, and (4) participants showed significantly higher emotion

<sup>¶</sup>Corresponding author.

This is an Open Access article published by World Scientific Publishing Company. It is distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 (CC BY-NC-ND) License which permits use, distribution and reproduction, provided that the original work is properly cited, the use is non-commercial and no modifications or adaptations are made.

recognition accuracy with the animal robot than the humanoid robot. A follow-up study ( $N = 10$ ) with voice-only conditions confirmed that the importance of embodiment. The results from this study could provide the guidelines needed to design social robots that consider emotional aspects in conversations between robots and users.

**Keywords:** Social robots; conversational agent; emotive voices; user perception; user preference.

## 1. Introduction

As robots have become prevalent in people's daily lives, expectations for social robots have increased, which has brought numerous studies regarding human–robot interaction (HRI). Robots are expected to play social roles such as a caregiver or companion that might serve as a friend or family member. In this regard, many studies have been conducted to facilitate richer and more natural interaction following human social norms. One of the ways of making the interaction more natural is attributing human characteristics to robots, called anthropomorphism.<sup>58</sup> It can be humanlike appearance (i.e., superficial human characteristics) or humanlike mind (i.e., essential human characteristics).<sup>66</sup> Some researchers have focused more on external design aspects (e.g., Ref. 18), whereas others have investigated more on human mind.<sup>21,66</sup>

Focusing on the appearance and behavior, research has been conducted on interactions between robots and users via multiple modalities incorporating variations in appearances, facial expressions, gestures, verbal communications, non-verbal sounds, and movements.<sup>28,45,47</sup> These modalities convey a wealth of information, influence user perception, and engage in establishing unique relationships between robots and users.

Focusing on the mental state, specifically on emotions, research has been conducted to see which factors influence user perception of robots' emotions. Although these studies have considered robots' facial expressions, voice (speech), body language, and posture as critical factors, the majority of emotion recognition research in HRI has focused on facial expressions.<sup>11,59</sup> Consequently, there has been little research on integrating both superficial and essential characteristics in one study to see interactions among the factors. A few exploratory studies have shown mixed results.<sup>22,23,43,46</sup> As such, to fill this research gap, we investigated the effects of various factors—robots' appearances (robot types), voice types, and emotions on users' perception—clarity, characteristics, naturalness, and preference, as well as emotion recognition accuracy.

## 2. Related Work

### 2.1. *Emotion taxonomy, expression, and perception*

There have been different theories proposed and studies conducted about (1) emotion classification, (2) emotion expression and (3) emotion perception in multiple domains, including psychology, psychiatry, neuroscience, and HRI research.

Largely, there are two types of emotion classification, including a dimensional approach and a categorical approach. In the dimensional approach, the circumplex

model has been widely used with arousal and valence dimensions.<sup>54,55</sup> An individual emotional state can be positioned on the Cartesian coordinate depending on the levels of arousal and valence. In the categorical approach, researchers often assume that people have basic emotions. Ekman's six basic emotions<sup>20</sup> (happiness, sadness, fear, anger, surprise, and disgust) have been one of the most widely mentioned emotion sets in emotion-related research in Psychology, Human Factors, Affective Computing, and HRI.<sup>10,11,52</sup> Basic emotions<sup>20</sup> are known to have unique features such as signal, physiology, and antecedent events, and common characteristics with other emotions such as rapid onset, short duration, unbidden occurrence, automatic appraisal, and coherence among responses. Ekman<sup>19</sup> argued that these basic emotions are expressed and recognized cross-culturally. However, there has been still much criticism about the basic emotion theories.<sup>48</sup> See Ref. 33 for more discussions on generic taxonomy and theories about emotions in the context of human factors and human-computer interaction. In our everyday lives, we typically describe our emotional states using categorical terms, rather than dimensional terms; for example, during a conversation, people usually express happy feelings as "happiness" (categorical) but not "an emotion that is high arousal with positive valence" (dimensional). Therefore, we provided the emotional states using the categorical approach in this study. Research also shows that these basic emotions are pervasive over the world.<sup>20</sup> In addition to Ekman's six emotions, we added 'anticipation', one of the Plutchik's basic emotions<sup>51</sup> because the passage of our stories included anticipation. With the addition of anticipation, we were able to have the second positive emotional state in our study in addition to happiness.

In terms of emotion expression, Darwin and Prodger<sup>15</sup> proposed three causal origins of expressions; immediate benefits (e.g., increasing one's body size to intimidate an opponent), effective communications (e.g., lowering one's body to signal submission), and vestigial byproducts that may not serve a useful role (e.g., trembling in fear). Previous studies also showed that emotion expressions exhibited useful functions (e.g., widening eyes to maximize the visual field during fear) and emotional vocal expressions effectively manipulated the behavior of perceivers.<sup>1,64</sup> Among these, this study focuses more on the effective communications and vocal expressions of emotion.

Emotion perception is the identification of emotionally salient information in the environment, including verbal (lexico-semantic) and nonverbal (intonational, facial, visual, and body movement) cues to the emotions of other people.<sup>50</sup> Emotion is one of the perceptual representations of social cues along with intentionality and eye direction.<sup>16,44</sup> In line with this, human social and emotional behaviors are highly intertwined.<sup>6</sup> Emotion perception is an important source of information about the theory of mind and emotions can be perceived from facial expressions, voices, and whole-body movements.<sup>30</sup>

As provided from previous theories, emotion expression and emotion perception play a critical role in human-robot interactions and are widely studied in a range of disciplines. Researchers commonly argue that these emotion-related expressions and

perceptions can be achieved through both visual and auditory stimuli. However, previous studies have been dominated by facial emotions and other modalities such as vocal and tactile processing have been less frequently considered.<sup>11,59</sup> In this regard, in our work, we focused more on auditory stimuli by including various emotive voices, representing seven different emotions and investigated the differences in users' emotion perception.

## 2.2. User perception on robots from embodiment, appearance, and sounds

There have been studies focused on examining the impact of robots' embodiment, appearance, and auditory displays on HRI.

The physical embodiment of robots could impact user perception positively and promote HRI in many social situations. With the embodiment, social robots brought many benefits to user experience. For example, participants reported higher satisfaction in the shopping mall<sup>57</sup> and higher enjoyment while playing a chess game<sup>49</sup> with the physical embodied robots than the disembodied ones. Many research studies also suggested that the embodiment of social robot engaged longer interaction duration,<sup>53</sup> increased human empathy towards the robots,<sup>38,60</sup> and enhanced compliance with robots' instruction and made the interaction more natural than the virtual or simulated ones.<sup>40</sup> Because the presence of the social robot played an important role in HRI, we used physical robots to emit sounds instead of using just a speaker in this study.

The appearance of robots was considered as an important factor of user perception to support interaction since anthropomorphism allows people to give robots lifelike qualities (e.g., intentions, emotions, etc.).<sup>60,61</sup> Barnes *et al.*<sup>3</sup> and FakhrHosseini *et al.*<sup>24</sup> showed that participants preferred robots which resemble animals or humans over imaginary creatures or robots highly deviating from existing creatures. Barnes *et al.*<sup>3</sup> compared five different robots (Robosapien, Pleo, Zoomer, Romo, and Mindstorm) which are humanoid, zoomorphic, fantastical, and mechanistic. Participants showed different user perception across robots but similar patterns before and after interacting with robots. Another study<sup>56</sup> suggested that a companion robot requires a certain level of emotional expression for a good interaction to occur with children. Also, people accept and trust robots more when the robots show some emotional activities.<sup>42</sup>

The effects of robots' voices have also been investigated in relation to user perception. These studies have employed different types of sounds, such as human voices, TTS voices, and beeping sounds in conjunction with various robots having different form factors. Research showed that participants assumed that a human voice was more capable than a TTS voice, and they anthropomorphized robots with human voices.<sup>62,65</sup> Similar to the pattern in user perception on robots' appearances, people showed a tendency to prefer interacting with robots similar to themselves in voice characteristics, including human-like speech style and accent, and gender.<sup>22,23</sup>

A recent exploratory study<sup>43</sup> showed that gender and naturalness of vocal manipulations strongly affected user perception.

Although various aspects of user perception from visual and auditory cues have been examined through exploratory studies, many of them focused more on users' preferences based on subjective self-report measures.<sup>3,24,25</sup> To tackle these issues, in our work, we applied both qualitative and quantitative measures by examining user perception from broader perspectives.

### 2.3. Emotions in HRI and emotive voices

An effective HRI could be achieved or improved by involving an appropriate emotional communication from social robots.<sup>41</sup> Regarding previous empirical studies on emotive communications in HRI, diverse aspects of communication such as gesture, appearance, style of speech, prosody, and context have been investigated. Implementing emotional features to social robots might enhance children's learning skills and engaged the learning process. Conti *et al.*<sup>13</sup> in their storytelling environment showed that children can memorize more details of a tale if the robot narrates with an expressive social behavior, even compared to the static inexpressive human storyteller. Also, the emotional appearance of robots was proposed for creating a more suitably moral agent<sup>12</sup> or providing interactive interventions for children with autism spectrum disorder (ASD).<sup>4,7</sup> With the results from previous studies, we considered emotion as an indispensable factor in HRI.

To investigate the impact of emotion expressions in HRI, there have been various research projects regarding emotional conversations that are driven by either internal states, behaviors, or situations.<sup>26,34,63</sup> These studies were based on communication theories about emotion expressions: (1) a robot's internal state drives expressions, (2) specific robot behaviors are related to specific user reactions, and (3) the situation is an important driver of emotion expressions.<sup>27</sup>

Regarding emotive voices on social robots, crucial features such as the style of speech, gender, and prosody have been widely investigated through exploratory studies in HRI. FakhrHosseini *et al.*<sup>24</sup> emphasized the importance of the congruency between anthropomorphism in the appearances and the style of speech. Their study showed that only when the human-like robot speaks with emotional expressions, participants perceive the robot as their social companion. Kishi *et al.*<sup>35</sup> showed that the integration of dynamic emotional expressions and movements made the humanoid robot more attractive, more favorable, more useful, and less mechanical-like. Gender stereotypes were also examined with the explicit gender (from name and voice) and implicit gender (from personality) in a previous study.<sup>9,37</sup> For example, in Kruas *et al.*'s study, no gender stereotypes were found for the explicit gender, but implicit gender showed a strong effect on trust and likability in the stereotypical male task. Participants perceived that the male personality robot (dominant, confident and assertive utterances) is more trustable, reliable, and competent than the female personality robot (agreeable and warm utterances), while the female personality

robot is more likable. A social robot's voice type could also play a critical role in emotive conversation. Eyssel *et al.*<sup>22</sup> examined the effects of vocal cues that reflected both the gender of robot voices (male, female) and voice types (robot-like, human-like). It showed a human voice was rated more likable than the synthetic voice. Jeon and Rayan<sup>33</sup> examined the effects of expressing affective prosody from a zoomorphic robot (Pleo) and showed a higher accuracy of emotion perception in a physical one than a virtual one. Half of the participants mentioned that the human voice generated from the zoomorphic robot was awkward and a characterized or a cartoon-like voice might be more appropriate. Recently, Ko *et al.*<sup>36</sup> investigated the effects of different voice types with two types of robots (same as in this study) on robot emotion perception. Text-to-speech (TTS) condition showed significantly lower emotion recognition accuracy than other human voices, but the robot type (humanoid vs. animal) did not influence emotion recognition accuracy or other robot perceptions. However, in their study the voice was recorded by *female* students, not voice experts, which might have led to different results from this study.

Overall, emotive voice associated with social robots is still veiled in various aspects such as acoustic characteristics, voice types, gender, and prosody. Since previous studies found contrasting results toward voice types in social robots, we narrowed down the scope and focused on the differences in emotion recognition accuracy and user perception on four different voice types in this study.

#### 2.4. Research questions and hypotheses

From this background, we tried to attain a deeper understanding of the effects of robot types, voice types, and emotion types on users' perception towards robots and their emotions. Especially, we aimed to answer the research questions as follows:

- RQ1: How do robot types, voices, emotions, and their interactions have impacts on participants' recognition of different robots' emotional states?
  - H1a: There will be no effects of robot types on emotion recognition accuracy.<sup>36</sup>
  - H1b: Participants will show higher emotion recognition accuracy in the human voice over TTS voice.<sup>36</sup>
  - H1c: There will be no emotion recognition accuracy difference between regular human and characterized human voices.<sup>36</sup>
  - H1d: Different emotions will show different emotion recognition accuracy.<sup>33,36</sup>
- RQ2: How do robot types, voices, and their interactions have impacts on participants' perception of robots' warmth, honesty, and trustworthiness?
  - H2a: Participants will show higher ratings on the humanoid robot than the animal robot in warmth, honesty, and trustworthiness ratings.<sup>3,24,25</sup>
  - H2b: There will be differences in warmth, honesty, and trustworthiness ratings among the different voice conditions.<sup>36</sup>

- RQ3: How do robot types, voices, and their interactions have impacts on participants' preference of robots?
  - H3a: Participants will prefer the humanoid robot over the animal robot.<sup>36</sup>
  - H3b: Participants will prefer the human voice over TTS.<sup>22,23</sup>
  - H3c: There will be no preference difference between regular human and characterized human voices.<sup>36</sup>

To address these research questions, we conducted an experimental study with young adults. Our participants read the two fairy tales to two types of robots each (human-like and animal-like). The robots made emotional comments using four different voices (regular human, characterized human-like, characterized animal-like, and TTS) with seven emotions (six basic emotions + anticipation).

### 3. Method

#### 3.1. Experimental design

Twenty university students participated in the study (age:  $M = 22.1$ ,  $SD = 2.97$ ). Twelve participants identified as male and the other eight participants identified as female. Participants were ethnically diverse (6 Asians, 1 Hispanic, 11 Caucasian, and 2 Multiracial). Participants participated in the experiment for at most two hours and participants were compensated with \$20 (\$10 per hour). All participants agreed to participate after reviewing the consent form approved by the university Institutional Review Board (IRB).

A 2 (robots)  $\times$  4 (voice types)  $\times$  7 (emotions) within-subjects design was applied. Therefore, eight different combinations of robots and voice types were provided to each participant with all seven emotions. Two social robots, NAO and Pleo, were used in the experiment. Four voice types were referred to two Characterized voices (NAO and Pleo), a Regular voice, and a TTS voice. There were two human voices and two TTS engines (Group A and Group B in Table 2) used. They were alternatively mapped to both robots and both stories across participants. More details were explained in the Procedure section.

#### 3.2. Robotic systems and stimuli

Two robots, NAO and Pleo, having different appearances and features were employed in the experiment (Fig. 1). We used these two robots, which represent a humanoid robot and zoomorphic robot each, to contrast the effects that robotic appearance has on people's emotion perception. NAO is a small-size humanoid robot (height: 57.4 cm, length: 27.4 cm, width 31 cm) having similarity to human and Pleo is a zoomorphic robot (height: 20.3 cm, length: 38.1 cm, width 10.2 cm) which looks like a little dinosaur. Both robots played recorded auditory feedback, which were emotive utterances, to participants following the storylines. The task selected to provide structure to the interaction and a more realistic context for conversational



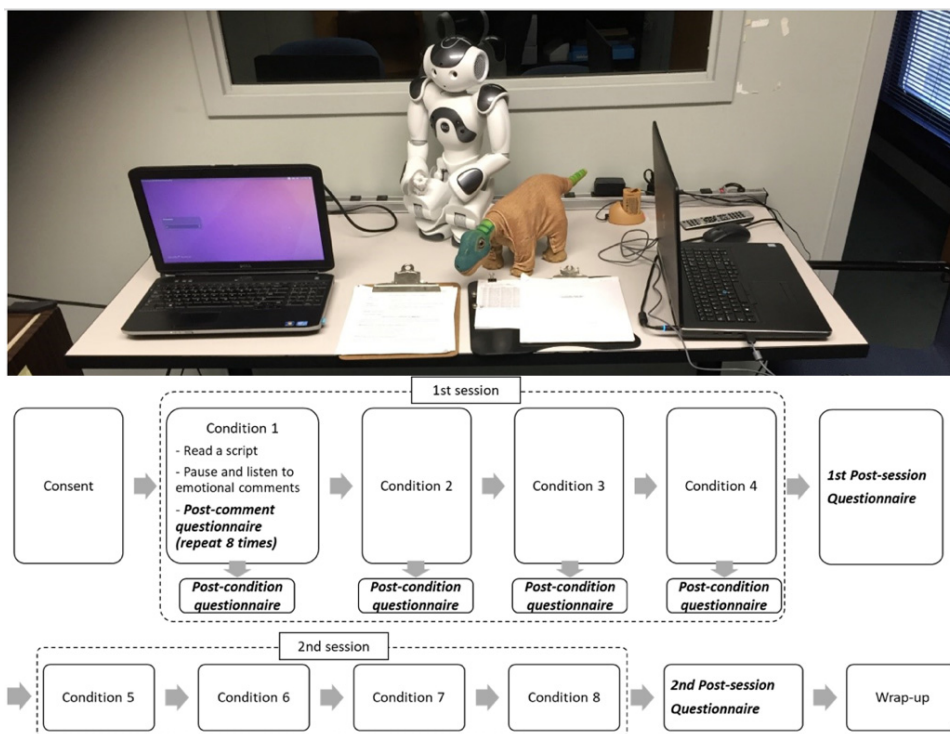


Fig. 1. Experiment settings with NAO (left) and Pleo (right) (upper) and experimental procedure including each step (lower).

emotions was to read fairy tales to the robots. Two different stories (“The three little pigs” and “The boy who cried wolf”) were used in this experiment. These two stories are simple narratives with easy vocabulary and globally well-known so that participants can easily read to the robots even if they are not native speakers. Crucially, we could include all of the emotions we wished to study within the framework of each story. Fairy tales seemed fitting given the childlike appearances of both robots and are suitable for use with a broad range of other populations for replication of this study.

Four voice types were created for seven emotional expressions. We first categorized different voice types as a TTS voice and a recorded human voice. The human voices were provided by two male voice actors and all the voices were speaking American English with American accents. Next, the recorded human voice was subdivided into three categories that included a regular voice and a characterized voice for each robot (i.e., characterized NAO voice and characterized Pleo voice). The TTS voices were generated using text-to-speech<sup>67</sup> engines. Microsoft’s David voice and the iOS Alex voice were used, which were provided by default with the respective operating systems. These TTS voices included no emotional information beyond the words themselves. Characterized voices for each NAO and Pleo were



designed to exaggerate emotional expressions with the robots' characters. These characterized voices were provided by voice actors who majored in performing arts while envisioning the characteristics of robots from their appearances. Direction for the characterization process, vocal performances, and recording was provided by a professional voice actor and professor of theatre who teaches voice and acting in the Department of Visual and Performing Arts. To control for gender effects, only characteristically male voices were used. While the same control effect could have been achieved using female voices, male voices were chosen based on the availability of the actors while designing the study. The example recordings of each voice type are provided on the web for other researchers and educators to get an idea of what participants heard during evaluation: <https://osf.io/m8h64/>.

Seven different emotions were presented throughout each story including Ekman's six basic emotions. The six basic emotions (anger, disgust, fear, happiness, sadness, and surprise) were chosen for their prevalence in psychology. Ekman's basic emotions have four negative emotions (anger, disgust, fear, and sadness), but have only one positive emotion (happiness); surprise can be either. A previous study showed that valence might influence people's emotion recognition accuracy.<sup>36</sup> In Bänziger, Grandjean, and Scherer's study,<sup>2</sup> participants were examined to recognized emotions, and the emotion recognition results showed a higher emotion recognition accuracy score on positive emotions, such as happiness, than the negative emotions, such as anxiety, sadness, and disgust. To make a balance between positive and negative emotions, the seventh emotion, anticipation, was chosen from Plutchik's eight basic emotions.<sup>51</sup> Its inclusion allowed us to add one more positive emotion in addition to happiness. The seven emotions fit into both stories ("The three little pigs" and "The boy who cried wolf") as depicted in Table 1. The content of these emotional phrases was not considered as an experimental factor in this study because all participants received the same treatments (eight combinations of robots and voice types) during the study.

3.3. Procedure

A single participant participated in each session. Note that this study was completed before the COVID pandemic. Thus, there was no COVID-relevant procedure.

Table 1. Dialogues in stories for presenting different emotions.

Presented emotions	Robots' utterance in a story	
	The boy who cried wolf	The three little pigs
Anger	That's not nice!	They shouldn't tease him like that
Anticipation	This should be good.	I wonder what's going to happen!
Disgust	Gross!	He can't want to EAT them!
Fear	He's going to eat the sheep!	Oh no!
Happiness	That sounds nice!	Good!
Sadness	All his sheep are gone	He destroyed their homes
Surprise	Why didn't they help?	Woah, that's fast!

After the consent form procedure, each participant interacted with all eight conditions of robots and voice types and all seven presented emotions. The eight conditions were separated into two sessions to help participants recall and compare four different conditions each. The presented order of each condition was counter-balanced. In each condition, the participant was instructed to read the script aloud in front of a robot and wait for and listen to the robots' emotional comments at various points in the story, which were marked down in the given script. Before reading the script and listen to the robots, participants were explained about all possible voice types they would interact with during the experiment. All voice clips were embedded in each robot and the voice was triggered by a remote controller which was controlled by an experimenter. Participants were aware that the robots were not acting autonomously. Other than vocal communication, the participants did not do any physical interaction with the robot.

The experimental environment (upper) and the whole procedure including each step (lower) are depicted in Fig. 1.

The participants were asked to fill out several questionnaires after listening to each comment generated from the robots, after finishing reading each full story, and after experiencing four conditions. Specifically, after each response to seven emotions, each condition, and each session, the surveys were conducted for measuring the accuracy of emotion perception and characteristics (warmth, honesty, trustworthiness), **naturalness and preferences (likability, attractiveness) of presented emotions**. The questionnaire consisted of open questions, seven-point Likert scales, and single-choice questions. Related questions were asked and each category was rated using a 1–7 Likert-scale (1: lowest, 7: highest) (Appendix A).

Presented orders for emotions in the two stories were different but the order in each story was fixed to maintain the storylines. Two different stories having the same seven emotions presented and two different voice groups having the same characteristics but recorded by different voice actors and two different TTS engines were employed to generalize the results. Each participant experienced both human voice

Table 2. Examples of the presented order.

PID	Start	Trial 1	Trial 2	Trial 3	Trial 4	Trial 5	Trial 6	Trial 7	Trial 8
1	Robot	NAO	Pleo	NAO	Pleo	NAO	Pleo	NAO	Pleo
	Voice type	Regular	Characterized NAO	TTS	Characterized Pleo	Characterized Pleo	TTS	Characterized NAO	Regular
	Story*	Pigs	Wolf	Pigs	Wolf	Pigs	Wolf	Pigs	Wolf
	Voice group**	Group A	Group A	Group A	Group A	Group B	Group B	Group B	Group B
2	Robot	Pleo	NAO	Pleo	NAO	Pleo	NAO	Pleo	NAO
	Voice type	Characterized NAO	Characterized Pleo	Regular	TTS	Characterized Pleo	TTS	Characterized NAO	Regular
	Story*	Pigs	Wolf	Pigs	Wolf	Pigs	Wolf	Pigs	Wolf
	Voice group	Group B	Group B	Group B	Group B	Group A	Group A	Group A	Group A

Notes: \*Pigs: The three little pigs, Wolf: The boy who cried wolf.

\*\*Group A and Group B had the same characteristics but were recorded by different voice actors and TTS engines.

Table 3. Accuracy, clarity, suitability, and preference over stories and voice groups.

		Accuracy	Clarity	Suitability	Preference
Story	The boy who cried wolf	57.0%	5.13	4.64	4.10
	The three little pigs	56.1%	5.25	4.78	4.38
Voice group	Group A	58.6%	5.05	4.53	4.16
	Group B	53.0%	5.11	4.68	4.33

actors and both TTS sounds. The examples of the presented order are depicted in Table 2. To validate the equivalence in emotion recognition accuracy, clarity, suitability, and preference, after the experiment, the results were analyzed (Table 3) showing similar results in all categories. The experiment took two hours at most as approved by IRB. Most participants completed it within 1.5–2 h.

## 4. Results

### 4.1. Data collection

The answer to open questions regarding emotions was interpreted by two examiners. Each examiner categorized all the answers into seven pre-defined emotions or marked as ‘indistinguishable’ if the answers do not fall into any categories. Two examiners worked independently, and the inter-rater reliability test showed the high coefficient value of Cronbach Alpha using variance ( $=0.86$ ). If interpretations from examiners were different, a third examiner reviewed the answers and decided which emotion the answer fell into.

### 4.2. Emotion perception: Accuracy, clarity, suitability, and features

First, the emotion recognition accuracy, defined as the proportion of correct emotion answers, was analyzed. Figure 2 and Table 4 show the descriptive statistics of emotion recognition accuracy across presented emotions, voice types, and robots. Regarding presented emotions, anger, disgust, and fear showed lower accuracies than positive emotions, such as anticipation and happiness. The accuracies for anger, disgust, and fear were 37.5%, 41.9%, and 25.6%, which were all lower than 50%. These three extreme conditions were excluded in statistical analysis to minimize the effects of biased data sets. Results were analyzed with the aligned rank transform (ART)<sup>68</sup> for factorial analyses since there are three factors (Robots, Voice Types, and Emotions) and dependent variable (1: correct, 0: wrong) is not normally distributed. To apply ART, we first computed residuals and estimated effects for all main and interaction effects. After computing aligned response, we assigned averaged ranks. With this data, we could perform a full-factorial repeated measures analysis of variance (ANOVA) following the guidelines of Wobbrock *et al.*<sup>68</sup> The ART allowed analyzing the aligned-ranked data with a 2 (robots)  $\times$  4 (voice types)  $\times$  4 (emotions) repeated measures ANOVA and testing all main effects and interaction effects. The result revealed a statistically significant difference across

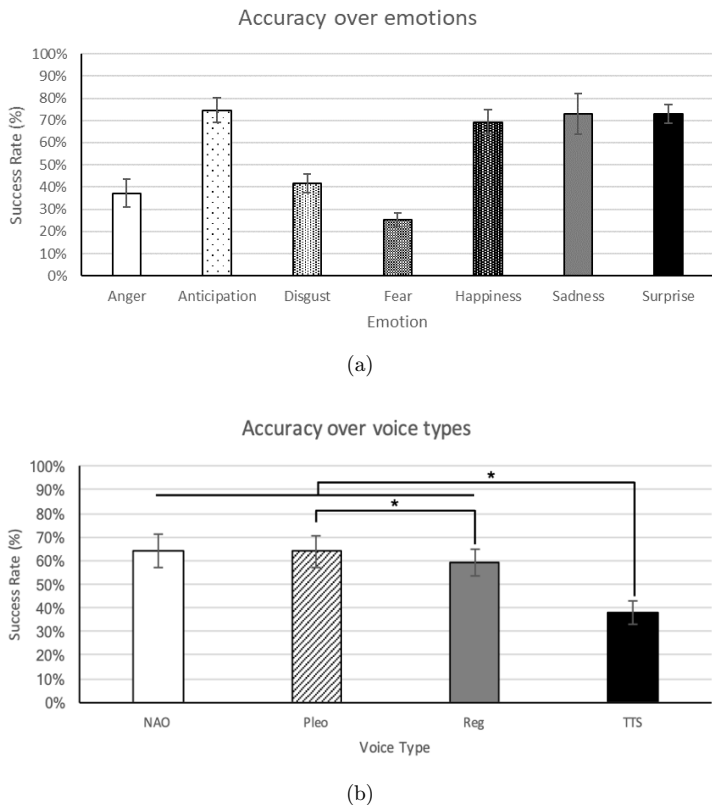


Fig. 2. Accuracy of perceiving emotions over emotions (a) and voice types (b) (\* $p < 0.0083$ ).

robots and voice types. However, there was no significant interaction effect between robots and voice types. For the multiple comparisons among voice types, paired-samples  $t$ -tests were conducted. All pairwise comparisons applied a Bonferroni adjustment to control for Type-I error in this study, which meant that we used more

Table 4. Statistics for emotion recognition (accuracy).

Measures	Conditions	Statistics
Accuracy (%)	Main effect for robots	$F(1, 607) = 4.27, p = 0.0393$
	NAO robot	
	$M = 0.68, SD = 0.47$	
	Pleo robot	
	$M = 0.76, SD = 0.43$	
	Main effect for voice types	$F(3, 607) = 16.07, p < 0.0001$
	Characterized NAO	$t^{19} = 5.78, p < 0.0001$
	$M = 0.64, SD = 0.48$	
	Characterized Pleo	$t^{19} = 6.15, p < 0.0001$
	$M = 0.64, SD = 0.48$	
	Regular	$t^{19} = 3.34, p = 0.0009$
	$M = 0.59, SD = 0.49$	
	Characterized Pleo	$t^{19} = 2.80, p = 0.0053$
	$M = 0.64, SD = 0.48$	
	Regular	$M = 0.59, SD = 0.49$

conservative alpha levels (critical alpha level = 0.0083 (0.05/6)). Participants recognized emotions more accurately with Pleo than NAO. Participants showed significantly lower emotion recognition accuracy in the TTS voice than all other three voice types. Moreover, the characterized Pleo voice showed significantly higher emotion recognition accuracy than the regular voice.

Table 5 shows the confusion matrix between presented and perceived emotions. Anger was mostly misclassified as sadness (32.50%), disgust was mostly misclassified as surprise (18.75%) or undistinguished (14.38%), and fear was mostly misclassified as anticipation (28.75%). Interestingly, 21.25% of happiness was also undistinguished even though it showed higher emotion recognition accuracy than anger, disgust, and fear did.

Second, clarity and suitability of perceived emotions over robots, voice types, and presented emotions were computed with the results as shown in Fig. 3 and Table 6. Clarity and suitability were rated using a 1–7 Likert-scale (1: Lowest, 7: Highest). We considered only responses with correctly recognized emotions. The clarity and suitability scores were measured for the present emotions; therefore, participants had to first recognize the emotions correctly to have their rating scores to be considered for the clarity and suitability measurements without bias. Overall, there were differences found in clarity over emotions and voice types and suitability over voice types. For robots, there were no significant differences found in both categories. Results were analyzed with a 2 (robot) × 4 (voice type) × 7 (emotions) repeated measures analysis of variance (ANOVA). The result revealed a statistically significant difference in clarity ratings among voice types and presented emotions. For the multiple comparisons among voice types, paired-samples *t*-tests were conducted. The TTS voice had a significantly lower clarity rating than the characterized and regular voices. In addition, the characterized Pleo voice had a significantly lower clarity

Table 5. The confusion matrix between presented and perceived emotions (grey: most misclassified).

Presented \ Perceived		Anger	Anticipation	Disgust	Fear	Happiness	Sadness	Surprise
Anger	Count	60	1	7	6	0	7	5
	Col %	37.50	0.63	4.38	3.75	0.00	4.38	3.13
Anticipation	Count	15	120	14	46	13	2	11
	Col %	9.38	75.00	8.75	28.75	8.13	1.25	6.88
Disgust	Count	8	1	67	0	2	0	0
	Col %	5.00	0.63	41.88	0.00	1.25	0.00	0.00
Fear	Count	0	0	14	41	0	0	1
	Col %	0.00	0.00	8.75	25.63	0.00	0.00	0.63
Happiness	Count	1	9	1	0	111	1	3
	Col %	0.63	5.63	0.63	0.00	69.38	0.63	1.88
Sadness	Count	52	1	4	27	0	118	9
	Col %	32.50	0.63	2.50	16.88	0.00	73.75	5.63
Surprise	Count	5	2	30	10	0	7	117
	Col %	3.13	1.25	18.75	6.25	0.00	4.38	73.13
Indistinguishable	Count	19	26	23	30	34	25	14
	Col %	11.88	16.25	14.38	18.75	21.25	15.63	8.75

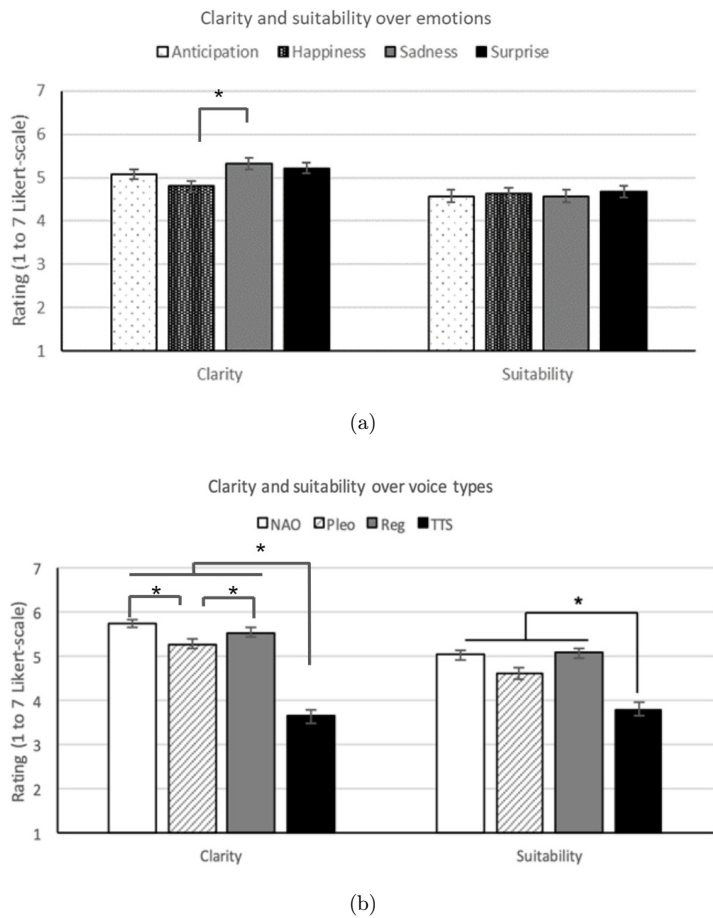


Fig. 3. The rating scores of clarity and suitability over emotions (a) and voice types (b) (\* $p < 0.05$ ).

Table 6. Statistics for clarity and suitability.

Measures	Conditions	Statistics
Clarity	Main effect for voice types	$F(3, 52.86) = 18.32, p < 0.0001$
	Characterized NAO	$t^{19} = 9.89, p < 0.0001$
	TTS	
	$M = 5.61, SD = 1.05$	$M = 3.63, SD = 1.67$
	Characterized Pleo	$t^{19} = 6.52, p < 0.0001$
	$M = 5.10, SD = 1.38$	
	Regular	$t^{19} = 11.36, p < 0.0001$
	$M = 5.76, SD = 1.22$	
	Characterized NAO	$t^{19} = 3.39, p = 0.0010$
	Characterized Pleo	
	$M = 5.61, SD = 1.05$	$M = 5.10, SD = 1.38$
	Regular	$t^{19} = 3.82, p = 0.0002$
	$M = 5.76, SD = 1.22$	
Main effect for emotions		$F(6, 115.1) = 3.25, p = 0.0055$

Table 6. (Continued)

Measures	Conditions		Statistics
Suitability	Sadness	Happiness	$t^{19} = 2.02, p = 0.0456$
	$M = 5.41, SD = 1.47$	$M = 5.00, SD = 1.45$	
	Interaction between voice types and emotions		$F(18, 312.3) = 2.77, p = 0.0002$
	Main effect for voice types		$F(3, 57.58) = 6.59, p = 0.0007$
	Characterized NAO	TTS	$t^{19} = 3.96, p = 0.0002$
	$M = 5.02, SD = 1.59$	$M = 3.79, SD = 1.63$	
	Characterized Pleo		$t^{19} = 3.07, p = 0.0032$
	$M = 4.61, SD = 1.77$		
	Regular		$t^{19} = 3.86, p = 0.0003$
	$M = 5.07, SD = 1.47$		

rating than the characterized NAO and regular voices. Participants reported Sadness as having a significantly higher clarity rating than Happiness. There was also a significant interaction effect between voice types and presented emotions. It is assumed that the relatively too low rating score of TTS voice compared to the other three voices caused the interaction effects. In suitability ratings, the result revealed a statistically significant difference among voice types. There were no significant interaction effects between emotions and voice types. For the multiple comparisons among voice types, paired-samples  $t$ -tests were conducted. Participants showed significantly lower rating scores in the TTS voice than all other three voice types.

Finally, the features by which to perceive emotions were analyzed with the results as shown in Table 7. The answers were collected from an open question (“What characteristics of the voice brought to mind that emotion?”) and the number of occurrences of words was counted. Each participant was allowed to provide multiple answers for each comment. After reading through each participant’s answer, we categorized their comments into different feature groups. Terms used in the participant’s answers that fell into specific features were counted. Most of the emotions were perceived from tone by 29.53%, words by 19.29%, and pitch by 17.72%. For each emotion, speech tone highly influenced perceiving anger (29.58%), anticipation (32.12%), happiness (32.56%), sadness (32.89%), and surprise (27.97%). Different from these emotions, disgust was mostly perceived by words (26.19%). Fear was perceived by different features such as pitch (24.49%), words (22.45%), and tone (20.41%).

### 4.3. Characteristics: Warmth, honesty, and trustworthiness

Figure 4 and Table 8 show the rating scores in warmth, honesty, and trustworthiness over voice types and robots. For robots, there were no significant differences found in three categories. Because by definition, emotions are short-lasting “states”, not long-lasting “traits”, the factor emotion was not analyzed in the following perception sections. Results were analyzed with a 2 (robot)  $\times$  4 (voice type) repeated measures analysis of variance (ANOVA). First, the result revealed a statistically significant



Table 7. The result of surveys on features that used to perceive emotions (grey: most used).

Feature		Anger	Anticipation	Disgust	Fear	Happiness	Sadness	Surprise	Total
Context	Count*	2	9	1	3	6	8	7	36
	Col %**	2.82%	6.57%	1.19%	6.12%	4.65%	5.37%	4.90%	4.72%
Familiarity	Count	3	7	5	7	5	9	6	42
	Col %	4.23%	5.11%	5.95%	14.29%	3.88%	6.04%	4.20%	5.51%
Length	Count			7		2	4	4	17
	Col %	0.00%	0.00%	8.33%	0.00%	1.55%	2.68%	2.80%	2.23%
Loudness	Count	8	5	4	2	3	3	5	30
	Col %	11.27%	3.65%	4.76%	4.08%	2.33%	2.01%	3.50%	3.94%
Mood	Count	3	5	5	1	8	6	6	34
	Col %	4.23%	3.65%	5.95%	2.04%	6.20%	4.03%	4.20%	4.46%
Pitch	Count	12	26	10	12	26	31	18	135
	Col %	16.90%	18.98%	11.90%	24.49%	20.16%	20.81%	12.59%	17.72%
Pronunciation	Count	4	1	3	2	1	4	8	23
	Col %	5.63%	0.73%	3.57%	4.08%	0.78%	2.68%	5.59%	3.02%
Speed	Count	2	5	4	1	2	15	9	38
	Col %	2.82%	3.65%	4.76%	2.04%	1.55%	10.07%	6.29%	4.99%
Tone	Count	21	44	19	10	42	49	40	225
	Col %	29.58%	32.12%	22.62%	20.41%	32.56%	32.89%	27.97%	29.53%
Words	Count	9	28	22	11	27	14	36	147
	Col %	12.68%	20.44%	26.19%	22.45%	20.93%	9.40%	25.17%	19.29%
Vague	Count	7	7	4		7	6	4	35
	Col %	9.86%	5.11%	4.76%	0.00%	5.43%	4.03%	2.80%	4.59%
Total	Count	71	137	84	49	129	149	143	762
	Col %	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%

Notes: \*The total number of answers.  
\*\* The proportion of the count in each column.

difference in warmth among voice types. There was no interaction effect between robots and voice types. For the multiple comparisons among voice types, paired-samples *t*-tests were conducted. In all three categories, the results commonly showed the lowest score in a TTS voice. Also, there were no significant differences among the characterized NAO, Pleo, and regular voices.

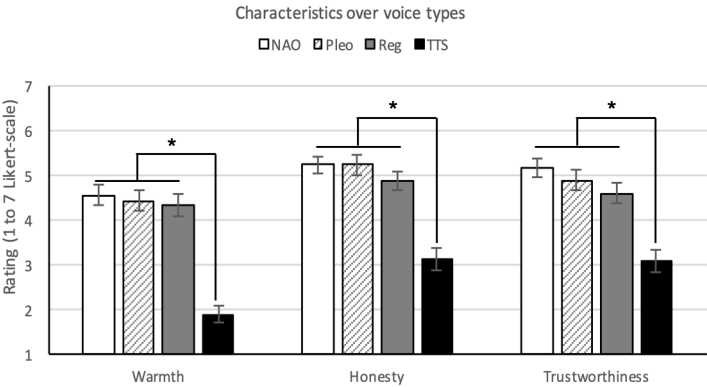


Fig. 4. The rating scores of characteristics (\**p* < 0.05).

Table 8. Statistics for characteristics (warmth, honesty, trustworthiness).

Measures	Conditions		Statistics
Warmth	Main Effect for Voice Types		$F(3, 57) = 33.84$ , $p < 0.0001$ , $\eta_p^2 = 0.640$
	Characterized NAO $M = 4.55$ , $SD = 1.52$	TTS $M = 1.88$ , $SD = 1.18$	$t^{19} = 7.48$ , $p < 0.0001$
	Characterized Pleo $M = 4.32$ , $SD = 1.55$		$t^{19} = 7.14$ , $p < 0.0001$
	Regular $M = 4.33$ , $SD = 1.49$		$t^{19} = 7.14$ , $p < 0.0001$
Honesty	Main Effect for Voice Types		$F(3, 57) = 32.24$ , $p < 0.0001$ , $\eta_p^2 = 0.630$
	Characterized NAO $M = 5.23$ , $SD = 1.19$	TTS $M = 3.10$ , $SD = 1.60$	$t^{19} = 6.67$ , $p < 0.0001$
	Characterized Pleo $M = 5.23$ , $SD = 1.40$		$t^{19} = 6.87$ , $p < 0.0001$
	Regular $M = 4.88$ , $SD = 1.34$		$t^{19} = 5.70$ , $p < 0.0001$
Trustworthiness	Main Effect for Voice Types		$F(3, 57) = 20.19$ , $p < 0.0001$ , $\eta_p^2 = 0.515$
	Characterized NAO $M = 5.15$ , $SD = 1.33$	TTS $M = 3.08$ , $SD = 1.54$	$t^{19} = 5.61$ , $p < 0.0001$
	Characterized Pleo $M = 4.88$ , $SD = 1.44$		$t^{19} = 5.11$ , $p < 0.0001$
	Regular $M = 4.58$ , $SD = 1.45$		$t^{19} = 4.17$ , $p < 0.0001$

4.4. Naturalness

Figure 5 and Table 9 show the rating scores in naturalness over voice types and robots. For voice types, the regular voice showed the highest scores in naturalness. For robots, there were no significant differences found in both categories.

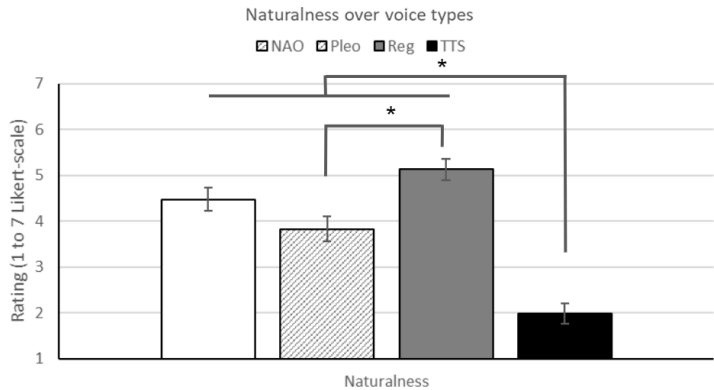


Fig. 5. The rating scores of naturalness (\* $p < 0.05$ ).

Table 9. Statistics for naturalness.

Measures	Conditions		Statistics
Naturalness	Main Effect for Voice Types		$F(3, 57) = 37.67, p < 0.0001,$ $\eta_p^2 = 0.665$
	Characterized NAO $M = 4.48, SD = 1.58$	TTS $M = 1.98, SD = 1.40$	$t^{19} = 6.75, p < 0.0001$
	Characterized Pleo $M = 3.83, SD = 1.71$		$t^{19} = 5.09, p < 0.0001$
	Regular $M = 5.13, SD = 1.42$		$t^{19} = 8.49, p < 0.0001$
	Characterized Pleo $M = 3.83, SD = 1.71$	Regular $M = 5.13, SD = 1.42$	$t^{19} = 3.45, p = 0.0011$

Results were analyzed with a 2 (robot)  $\times$  4 (voice type) repeated measures analysis of variance (ANOVA). Since there was no interaction effect between robots and voice types, paired-samples  $t$ -tests were conducted for the multiple comparisons among voice types. First, the result revealed a statistically significant difference in the rating scores in naturalness among voice types. Participants showed significantly lower rating scores in the TTS voice than all other three voice types. The regular voice showed significantly higher rating scores than the characterized Pleo voice.

4.5. Preferences: Likability and attractiveness

Figure 6 and Table 10 show the rating scores in likability and attractiveness over voice types and robots. Among voice types, the TTS voice commonly showed the lowest rating scores in both categories. For robots, there were no significant differences found in both categories.

Results were analyzed with a 2 (robot)  $\times$  4 (voice type) repeated measures analysis of variance (ANOVA). First, the result revealed a statistically significant

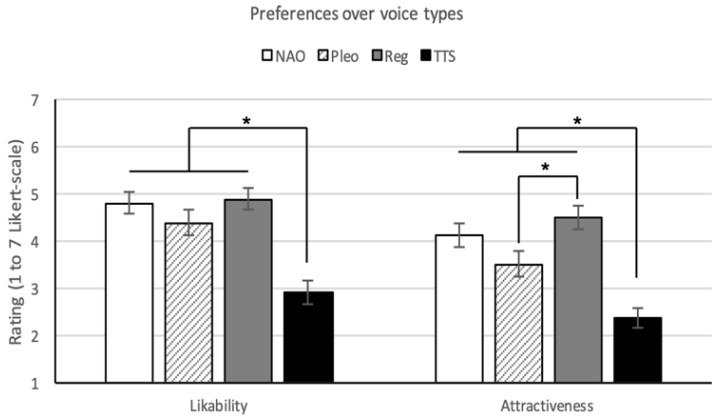


Fig. 6. The rating scores of preferences ( $*p < 0.05$ ).

Table 10. Statistics for preferences (likability, attractiveness).

Measures	Conditions		Statistics
Likability	Main Effect for Voice Types		$F(3, 57) = 18.91$ , $p < 0.0001$ , $\eta_p^2 = 0.499$
	Characterized NAO	TTS	$t^{19} = 4.84$ , $p < 0.0001$
	$M = 4.80$ , $SD = 1.44$	$M = 2.90$ , $SD = 1.57$	
	Characterized Pleo		$t^{19} = 3.90$ , $p = 0.0003$
	$M = 4.38$ , $SD = 1.64$		
Attractiveness	Regular		$t^{19} = 5.19$ , $p < 0.0001$
	$M = 4.88$ , $SD = 1.42$		
	Main Effect for Voice Types		$F(3, 57) = 18.65$ , $p < 0.0001$ , $\eta_p^2 = 0.495$
	Characterized NAO	TTS	$t^{19} = 4.85$ , $p < 0.0001$
	$M = 4.10$ , $SD = 1.53$	$M = 2.38$ , $SD = 1.33$	
	Characterized Pleo		$t^{19} = 3.18$ , $p = 0.0025$
	$M = 3.50$ , $SD = 1.63$		
	Regular		$t^{19} = 6.14$ , $p < 0.0001$
	$M = 4.50$ , $SD = 1.53$		
	Characterized Pleo	Regular	$t^{19} = 2.97$ , $p = 0.0045$
	$M = 3.50$ , $SD = 1.63$	$M = 4.50$ , $SD = 1.53$	

difference in likability among voice types. There was no interaction effect between robots and voice types. For the multiple comparisons among voice types, paired-samples *t*-tests were conducted. Participants showed significantly lower rating scores in the TTS voice than all other three voice types. Next, the result revealed a statistically significant difference in attractiveness among voice types. There was no interaction effect between robots and voice types. For the multiple comparisons among voice types, paired-samples *t*-tests were conducted. Same as shown in a likability category, participants showed significantly lower rating scores in the TTS voice than all other three voice types. The regular voice showed significantly higher rating scores than the characterized Pleo voice.

5. Discussions

In the experiment, 20 participants experienced verbal interactions with robots while reading scripts of fairy tales to robots. Humanoid and zoomorphic robots used four different voice types and seven emotions were presented to participants through robots’ verbal comments. Each participant interacted with all eight conditions of robots and voice types and all seven presented emotions. The participant was instructed to read the script in front of a robot and listen to the emotional comment from the robot at various points in the story. The participant filled out the questionnaire after listening to each emotional comment, completing each condition and completing four conditions. The emotion recognition accuracy and subjective ratings such as characteristics, naturalness, and user preferences were measured.

Referring to the research questions and hypotheses in Sec. 2.4, the results are listed as follows:

- RQ1:
  - H1a (rejected): A significantly higher emotion recognition accuracy was reported from Pleo robot than NAO robot.
  - H1b (supported): The TTS voice showed significantly lower emotion recognition accuracy than the characterized NAO, characterized Pleo, and regular voices.
  - H1c (rejected): The characterized Pleo voice showed significantly higher emotion recognition accuracy than the regular voice.
  - H1d (supported): Anger, disgust, and fear had significantly lower emotion recognition accuracy with lower rating scores in clarity and suitability than other emotions.
- RQ2:
  - H2a (rejected): No significant difference was found among robot types for different characteristics ratings.
  - H2b (supported): The TTS voice showed significantly lower rating scores in warmth, honesty and trustworthiness than the characterized NAO, characterized Pleo, and regular voices; **and the regular voice showed significantly higher rating scores in naturalness than the characterized Pleo and TTS voices.**
- RQ3:
  - H3a (rejected): There were no significant differences found in both likeability and attractiveness ratings for robot types.
  - H3b (supported): The regular voice showed significantly higher rating scores in attractiveness than the TTS voice.
  - H3c (rejected): The regular voice also showed significantly higher rating scores in attractiveness than the characterized Pleo voice.

The critical points and explanations in each category are described below by dependent variables.

### 5.1. Accuracy, clarity, and suitability

The result showed that the emotion recognition accuracy significantly differed depending on presented emotions (H1d). As shown in Table 5, overall, unpleasant emotions with high arousal levels such as anger, disgust and fear showed significantly lower emotion recognition accuracy than other emotions such as anticipation, happiness, surprise and sadness did. There might be possible explanations about why some emotions were not accurately perceived. First, the emotion recognition accuracy results aligned with our previous study<sup>36</sup> that negative emotions received lower emotion recognition accuracy than positive emotions. Those two fairy tales used in the experiments were well-known for children and thus, participants might expect

pleasant emotions more than unpleasant emotions. The most misclassified three emotions were all unpleasant emotions with high arousal levels.<sup>54</sup> Next, the intensity of emotions might be different, which causes inequivalence among emotions. For example, among auditory stimuli used in the experiment, the intensity of unpleasant emotions might be lower than the one of positive emotions. Lastly, the mixed result was possible because there were many emotions presented through auditory cues. As shown in Ref. 8, although emotion recognition can be fairly accurate when listeners choose from a limited set of emotion categories, agreement drops significantly as more categories of emotion become available. Note that in our experiment, the participants freely guessed each emotion without preset options. Also, fewer emotions can be perceived from voice<sup>14</sup> compared to facial expressions.

For voice types (H1b & H1c), as expected, the TTS voice showed significantly lower emotion recognition accuracy than all other human voice types — \*\*\*characterized NAO, characterized Pleo, and regular voices — did. Furthermore, the TTS voice also showed significantly lower rating scores in clarity and suitability. It suggests that these TTS voices are inappropriate for emotive expressions since the intended emotions might not be delivered correctly to listeners even though they have the same semantic content. Instead, recorded human voices such as characterized NAO, characterized Pleo, and regular voices are more suitable for robots to express emotive voices and deliver emotions correctly. Most interestingly, the characterized Pleo voice showed significantly higher emotion recognition accuracy than the regular voices did. There was a possibility that these results suggest that a characterized voice might be more appropriate for emotive expressions delivering intended emotions more accurately and facilitating the interactions than just a regular voice. However, because only characterized Pleo voice showed a higher emotion recognition accuracy in this study, more research should be conducted to determine if characterized voice types are more effective than the regular voice in expressing the emotions more accurately. It also suggests that there may be value in creating TTS engines that exaggerate emotional characterization for use in contexts where highly recognizable emotional signals are desired. Mimicking a natural speaking style may not be the optimal approach for delivering emotional information via synthetic speech from a robot. The results provide additional guidance on designing robot speech to deliver different emotions more effectively. As shown, other emotions can be sufficiently conveyed by affective tones, but disgust and fear require more semantic contents.

For robot types (H1a), NAO showed significantly lower emotion recognition accuracy than Pleo for happiness (NAO:  $M = 0.61$ ,  $SD = 0.49$ ; Pleo  $M = 0.76$ ,  $SD = 0.43$ ,  $p < 0.05$ ). However, there was no difference between voice types of the two robots. We can cautiously infer that the participants might expect happy expressions from Pleo more than Nao and it caused higher emotion recognition accuracy in happiness. According to the previous findings,<sup>17,29,32</sup> people perceive that Pleo manifested positive emotions (e.g., love, grateful) more than NAO (e.g., uneasy, fear). However, to the best of our knowledge, the relationships between perceived

emotions (e.g., happiness) and robots' appearances have not been comprehensively studied. The overall underlying cognitive process of recognizing emotions from form factors should be investigated in the future.

### 5.2. *Characteristics, naturalness, and preferences*

Surprisingly, no significant difference was found on participants' perception of robot's characteristics and preferences (H2a & H3a). This result might suggest that participants perceived both robots as similar, or they evaluated the auditory portion of the social robots more than the embodiment and appearance regarding the ratings for each category. Because participants reported a significantly higher emotion recognition accuracy in Pleo than NAO robot, this might imply that performance and perception might not always be congruent. In the results, the TTS voice showed the significantly lowest rating scores across all characteristics and preferences including likability, attractiveness, warmth, honesty, and trustworthiness (H2b & H3b). The TTS voice showed a significantly lower rating score in the naturalness feature and the result might be because it had basically a flat voice without variations in pitch and speed. Other recorded voices such as characterized and regular voices having intonations and variations in speech showed significantly higher scores in the naturalness rating than the TTS voice.

A regular voice showed significantly higher rating scores in naturalness and attractiveness than a characterized Pleo voice (H3c). The results indicate that a regular voice might be more suitable for general use with higher user preferences and naturalness than characterized or TTS voices.

Overall, these results indicate that the characterized voice might lead to the highest emotion recognition accuracy, but the regular voice is the most preferred. It is assumed that characterized voices might be appropriate for emotional expressions. On the other hand, regular voices which show the highest attractiveness and naturalness might be suitable for general use. For example, for the first stage of human-robot interaction, regular voices might be appropriate to facilitate the interaction. However, for the next step for in-depth and emotion-related interactions, a characterized voice might be helpful to express emotional states and establish a unique relationship between users and robots since this stage involves personal familiarity with the other person and strong emotional commitment to the relationship.<sup>39</sup> To further generalize our results, more experiments are required to consider possible other variables.

### 5.3. *Anecdotal findings*

Interestingly, there were no significant effects of the appearance of robots on all dependent variables except for emotion recognition accuracy. This might be because the given tasks were mostly focused on conversation which requires reading aloud and listening to verbal feedback but were not relevant to visual cues as much as auditory cues. According to Ref. 30, emotions are perceived by facial expressions and



whole body movements instead of fixed features such as appearances, but these dynamic visual cues were not applied in this experiment.

There were interesting comments on auditory feedback from participants. A participant said, “(P2) *The final robot seemed to be happy at the start of the wolf story. My brain was saying it shouldn’t be that but that’s all my emotions were getting*”, which indicates the individual differences in expectation. Other comments such as “(P15) *The robots sounded more surprised/happier than showing signs of any other emotion*” and “(P18) *When Pleo would say “What!” in a shocked tone, it was easy to recognize his surprise in both the natural sounding voice and robotic sounding voice,*” which showed that the intensity of emotions could vary for different participants.

#### 5.4. Limitations

There are limitations and improvements that need to be considered in the next experiment to broaden this study and draw more reliable results. First, twenty participants may not be enough to generalize the results of this study. We plan to replicate the study with more participants and expand it to other populations (e.g., children and older adults). Because this study includes multiple factors (robot types, emotions, and voice types), a different approach of statistical tests could be used (e.g., a linear mixed effect model), to investigate the effects of multiple factors on one measurement. In the future study, we will explore more appropriate statistical tests for further analysis.

The equivalence among the intensity of emotions should be secured. We used one of the most widely used emotion sets, Ekman’s basic emotions, but the result showed that some of them were not clearly distinguished by participants. This study excluded the selected negative emotions with poor emotion recognition accuracy due to potential biases, but again using a different statistical model or analysis will help us understand the deviation. Using the only two phrases for each emotion might have provided biases to the participants’ emotion recognition. Also, it may not be sufficient to ensure the generalizability of the finding. Depending on the content of the phrase, emotional semantics or strength might have been changed. However, as our results indicated, even with those same phrases, the participants showed significantly different emotion recognition accuracy depending on the robot type and voice type. In future research, we will diversify the phrases more with the similar length. The order of presentation might also have influenced the participants’ responses. However, it is an intrinsic limitation because we were not able to change the storyline every time. If we randomly change the order of emotions without the context of the story, the experiment might lack external validity. We believe that people perceive emotions in the context.

Next, the characteristics of voice types should be more specifically studied to figure out which factors cause differences. In this study, characterized NAO and Pleo voices were generated by voice actors to exploit their expertise. It was a first attempt

to produce the voice that well expresses the characteristics of NAO and Pleo. Regarding the emotion recognition accuracy results, participants reported a significantly higher emotion recognition accuracy in the characterized Pleo voice (but not in the characterized Nao voice) than the regular voice. The reason for this result might be that different appearances of the robots (animal versus humanoid) impacted participants' emotion recognition, because participants recognized emotions significantly more accurately in the Pleo robot than the NAO robot. In the follow-up study (Appendix B), participants reported a higher emotion recognition accuracy in both characterized voices (NAO and Pleo) than the regular voice. In the next experiment, the acoustic characteristics with specific physical properties (e.g., frequency range, speed, intensity) will also be considered when the representative voice types were designed so that the influential factors for different voice types will be investigated in depth. This approach will enable us to quantify the relationship between voice parameters and perception effects and model the robot voices. The gender effects will also be investigated. In this experiment, only male voices were used to control the gender effect and female voices were not included. We will design female voices for all four voice types and compare the gender differences in the following experiment.

There might have been some novelty effects. The participants did not have any previous opportunity to interact with or see the robots used in this study. To minimize any novelty effects, the orders of the robots and voice types were counter-balanced across participants. Therefore, while interacting with the robots, the plausible novelty effects might have been reduced. We also had a standardized introductory section and minimized features used in the experiment (i.e., we used only the "speech" function and did not use other features, such as moving robot arms or its head). We are conducting separate experiments to see the effects of robot gestures and facial expressions. Taking all together of these experiments, we will be able to see the separate and overall effects.

## 6. Future Work

Throughout this study, various aspects of social robots such as appearances, emotive expressions, and voice types were investigated. Based on the results and experimental settings, follow-up studies will be conducted with two complementary approaches. First, the research scope will be narrowed down to focus more on the acoustic characteristics of voice types having distinct features. This approach will help in-depth understanding in emotive and interactive robotic systems and developing computational models for emotional and conversational human-robot interactions. Gender-specific factors such as the user's gender and the gender of robot voice will also be considered based on the previous result.<sup>22</sup> Meanwhile, other factors such as ages and modalities will be included to widen the research scope to investigate the multiple influential factors. As provided from previous studies,<sup>28,45,47</sup> considering that the interactions take place via various modalities, facial expressions,

gaze and gestures,<sup>31</sup> and even non-verbal sounds can be included as independent variables. The results will provide a design guideline for emotional and trustworthy robots, especially employing emotive expressions and facilitate the relationship between people and social robots such as assistive robots, voice assistants, and any other conversational agents.

## **Acknowledgements**

This work was partly supported by National Institutes of Health (US) (No.1 R01 HD082914-01).

## **Appendix A. Questionnaires**

- Post-comment questionnaire
  - What emotion do you feel the robot expressed? (Open question)
  - What characteristics of the voice brought to mind that emotion? (Open question)
  - How clearly did the robot express this emotion? (1–7 Likert scale)
  - How suitable was this emotion coming from the robot? (1–7 Likert scale)
- Post-condition questionnaire
  - How likable is the voice? (1–7 Likert scale)
  - How attractive is the voice? (1–7 Likert scale)
  - How warm is the voice? (1–7 Likert scale)
  - How honest is the voice? (1–7 Likert scale)
  - How trustworthy is the voice? (1–7 Likert scale)
  - How natural does the voice sound? (1–7 Likert scale)
- Post-session questionnaire
  - Thoughts about first, second, third, and fourth voices (Open question)
  - Which story was your favorite? (Open question)
  - What is your sex? (Open question)
  - What is your age? (Open question)
  - What is your race and/or ethnicity? (Multiple-choice, Open question)

## **Appendix B. Voice Types Validation Study**

To further investigate the impact of robot embodiment on participants' perception towards different voice types, we conducted a follow-up validation study for voice types only. Based on the results of the main study, TTS voice showed significantly lower score on the most subjective ratings. Therefore, this validation study used only human voices, which made a three (voice types) by seven (emotions) within-subjects design. Ten new participants (age:  $M = 22.5$ ,  $SD = 4.12$ ) were recruited for the

follow-up study. Six participants identified as male and four participants identified as female with five Asians, four Caucasian, and one Hispanic. They listened to all recordings and evaluated three voice types: Characterized NAO voice, Characterized Pleo voice, and Regular human voice. Because the suitability rating subjectively determined how suitable the voice types were on a certain robot, we excluded the scale in the validation study because there was no robot or physical embodiment involved with this follow-up study.

### **B1. Accuracy**

Following the main study, the emotion recognition accuracy data were transformed with the aligned rank transform (ART).<sup>68</sup> Then, the aligned-ranked data were analyzed with a 3 (voice types)  $\times$  7 (emotions) repeated measures ANOVA, followed by paired samples *t*-tests with a Bonferroni correction for pairwise comparisons. A significant difference was found in the main effects of voice types,  $F(2, 18) = 11.68$ ,  $p < 0.001$ ,  $\eta_p^2 = 0.567$  emotions,  $F(6, 54) = 4.61$ ,  $p < 0.001$ ,  $\eta_p^2 = 0.339$  and the interaction effect between voice types and emotions,  $F(12, 108) = 4.48$ ,  $p < 0.001$ ,  $\eta_p^2 = 0.342$ . The average accuracy of emotion recognition in both characterized voices (NAO and Pleo) were significantly higher than the regular voice. The average accuracy was significantly higher in happiness (65.7%), sadness (77.6%), and surprise (67.6%) than anger (41.6%), disgust (37.6%), and fear (37.1%), which is similar to the main study. However, the average accuracy of anticipation (58.9%) was much lower compared to the percentage of the main study (75%). It might not be appropriate to compare the absolute percentage between the main study and the follow-up study because of different population and different number of participants. However, the average emotion recognition accuracy of the main study (56.61%) is numerically higher than that of the follow-up study (55.16%). The emotion recognition accuracy of the four emotions (happiness, anticipation, surprise, and disgust) was numerically higher in the main study than in the follow-up study. This might imply that when the voice is presented with embodied robots, emotion recognition accuracy might increase depending on different emotions. Further analysis of the interaction effects showed that the accuracy of emotion recognition was higher when characterized voices were paired with emotions that are positive and high arousal, such as happiness and surprise, or negative and low arousal, such as sadness than the regular voices paired with the emotions with opposite valence and arousal, such as anger, disgust, and fear. These results might suggest that the characterized voices improve participants' emotion recognition capabilities for certain emotions compared to regular human voices when there was no physical embodiment.

### **B2. Other subjective ratings**

The results from other subjective ratings of this validation study were similar to the results in the main study. The main effect of voice types was found significant in the scale of warmth,  $F(2, 832) = 3.65$ ,  $p = 0.0466$ ;  $\eta_p^2 = 0.297$ ; trustworthiness,  $F(2, 832) = 5.38$ ,  $p = 0.0147$ ,  $\eta_p^2 = 0.375$ ; naturalness,  $F(2, 832) = 17.57$ ,  $p < 0.0001$ ,

$\eta_p^2 = 0.664$ ; likeability,  $F(2, 832) = 10.20$ ,  $p = 0.0011$ ,  $\eta_p^2 = 0.532$ ; and attractiveness,  $F(2, 832) = 12.42$ ,  $p = 0.0004$ ,  $\eta_p^2 = 0.586$ .

Participants rated higher scores of warmth, and trustworthiness in regular voices than just the characterized Pleo voice. However, participants reported higher scores of naturalness, likeability, and attractiveness in regular voices than both characterized NAO and characterized Pleo voices. Note that in the main study, regular voice did not show higher scores of warmth and trustworthiness than the characterized voices. This might suggest that the appearance and embodiment of the robots can improve participants' perception toward the characterized voice positively such as increasing the warmth and trustworthiness of the robot. It is interesting to see that the validation study results of naturalness aligned with the results in the main study because it might imply that naturalness did not necessarily influence warmth and trustworthiness of the robot.

In sum, when there is embodiment of the robots, overall, people may recognize the same voice's emotions better. Also, they may perceive the characterized voice more positively (e.g., warm and trustworthy). The results of the validation study once again revealed the importance of the robot appearance and embodiment in HRI.

## References

1. J. A. Bachorowski and M. J. Owren, Sounds of emotion: Production and perception of affect-related vocal acoustics, *Ann. NY. Acad. Sci.* **1000**(1) (2003) 244–265.
2. T. Bänziger, D. Grandjean and K. R. Scherer, Emotion recognition from expressions in face, voice, and body: The multimodal emotion recognition test (MERT), *Emotion* **9**(5) (2009) 691.
3. J. Barnes, S. M. F. Hosseini, M. Jeon, C. H. Park and A. M. Howard, The influence of robot design on acceptance of social robots, in *Paper presented at the 2017 14th Int. Conf. Ubiquitous Robots and Ambient Intelligence (URAI)* (Jeju, Korea (South), 2017), pp. 51–55.
4. J. A. Barnes, C. H. Park, A. Howard and M. Jeon, Child-robot interaction in a musical dance game: An exploratory comparison study between typically developing children and children with autism, *Int. J. Hum. Comput. Interact.* **37**(3) (2021) 249–266.
5. J. Barnes, E. Richie, Q. Lin, M. Jeon and C. H. Park, Emotive voice acceptance in human–robot interaction, in *Paper presented at the Proc. 24th Int. Conf. Auditory Display* (Houghton, MI, 2018), pp. 271–274.
6. J. S. Beer and K. N. Ochsner, Social cognition: a multi level analysis, *Brain Res.* **1079**(1) (2006) 98–105.
7. R. Bevill, C. H. Park, H. J. Kim, J. W. Lee, A. Rennie, M. Jeon and A. M. Howard, Interactive robotic framework for multi-sensory therapy for children with autism spectrum disorder, in *2016 11th ACM/IEEE Int. Conf. Human-Robot Interaction (HRI)* (IEEE, New York, 2016), pp. 421–422.
8. P. Birkholz, L. Martin, K. Willmes, B. J. Kröger and C. Neuschaefer-Rube, The contribution of phonation type to the perception of vocal emotions in German: An articulatory synthesis study, *J. Acoust. Soc. Am.* **137**(3) (2015) 1503–1512.
9. D. Bryant, J. Bornstein and A. Howard, Why should we gender? The effect of robot gendering and occupational stereotypes on human trust and perceived competency, in

- Paper presented at the ACM/IEEE Int. Conf. Human-Robot Interaction (HRI)* (Cambridge, UK, 2020), pp. 13–21.
10. M. Cakmak, G. Hoffman and A. Thomaz, Computational human-robot interaction, *Found. Trends Robot.* **4**(2–3) (2016) 104–223. doi: 10.1561/23000000049.
  11. R. A. Calvo and S. D’Mello, Affect detection: An interdisciplinary review of models, methods, and their applications, *IEEE Trans. Affect. Comput.* **1**(1) (2010) 18–37. doi: 10.1109/t-affc.2010.1.
  12. M. J. E. Coeckelbergh and I. Technology, Moral appearances: Emotions, robots, and human morality, **12**(3) (2010) 235–241.
  13. D. Conti, C. Cirasa, S. Di Nuovo and A. Di Nuovo, “Robot, tell me a tale!”: A social robot as tool for teachers in kindergarten, *Interact. Stud.* **21**(2) (2020) 220–242.
  14. D. T. Cordaro, D. Keltner, S. Tshering, D. Wangchuk and L. M. Flynn, The voice conveys emotion in ten globalized cultures and one remote village in Bhutan, *Emotion* **16**(1) (2016) 117.
  15. C. Darwin and P. Prodger, *The Expression of the Emotions in Man and Animals* (Oxford University Press, USA, 1998).
  16. J. Decety, The neurodevelopment of empathy in humans, *Develop. Neurosci.* **32**(4) (2010) 257–267.
  17. M. Díaz, N. Nuño, J. Saez-Pons, D. E. Pardo and C. Angulo, Building up child-robot relationship for therapeutic purposes: From initial attraction towards long-term social engagement, in *Face and Gesture 2011* (IEEE, New York, 2011), pp. 927–932.
  18. C. F. DiSalvo, F. Gemperle, J. Forlizzi and S. B. Kiesler, All robots are not created equal: The design and perception of humanoid robot heads, in *Proc. Conf. Designing Interactive Systems: Processes, Practices, Methods, and Techniques, DIS* (London, UK, 2002), pp. 321–326.
  19. P. Ekman, An argument for basic emotions, *Cognit. Emot.* **6**(3–4) (1992) 169–200.
  20. P. Ekman and D. Cordaro, What is meant by calling emotions basic, *Emot. Rev.* **3**(4) (2011) 364–370.
  21. N. Epley, A. Waytz and J. T. Cacioppo, On seeing human: A three-factor theory of anthropomorphism, *Psychol. Rev.* **114**(4) (2007) 864–886.
  22. F. Eyssel, L. De Ruiter, D. Kuchenbrandt, S. Bobinger and F. Hegel, ‘If you sound like me, you must be more human’: On the interplay of robot and user features on human-robot acceptance and anthropomorphism, in *Paper presented at the 2012 7th ACM/IEEE Int. Conf. Human-Robot Interaction (HRI)* (Boston, MA, 2012), pp. 125–126.
  23. F. Eyssel, D. Kuchenbrandt, F. Hegel and L. de Ruiter, Activating elicited agent knowledge: How robot and user features shape the perception of social robots, in *Paper presented at the 2012 IEEE RO-MAN: The 21st IEEE Int. Symp. Robot and Human Interactive Communication* (2012).
  24. S. M. Fakhrosseini, S. Hilliger, J. Barnes, M. Jeon, C. H. Park and A. M. Howard, Love at first sight: Mere exposure to robot appearance leaves impressions similar to interactions with physical robots, in *Paper presented at the 2017 26th IEEE Int. Symp. Robot and Human Interactive Communication (RO-MAN)* (2017).
  25. S. M. Fakhrosseini, D. Lettinga, E. Vasey, Z. Zheng, M. Jeon, C. H. Park and A. M. Howard, Both “look and feel” matter: Essential factors for robotic companionship, in *Paper presented at the 2017 26th IEEE Int. Symp. Robot and Human Interactive Communication (RO-MAN)* (2017).
  26. J. Feldmaier, M. Stimpfl and K. Diepold, Development of an emotion-competent SLAM agent, in *Paper presented at the Proc. Companion of the 2017 ACM/IEEE Int. Conf. Human-Robot Interaction* (2017).

27. K. Fischer, M. Jung and L. C. Jensen, Emotion expression in HRI: When and why, in *Paper presented at the 2019 14th ACM/IEEE Int. Conf. Human-Robot Interaction (HRI)* (Daegu, Korea (South), 2019), pp. 29–38.
28. T. Fong, I. Nourbakhsh and K. Dautenhahn, A survey of socially interactive robots, *Robot. Auton. Syst.* **42**(3–4) (2003) 143–166.
29. M. R. Fraune, S. Sherrin, S. Sabanović and E. R. Smith, Rabble of robots effects: Number and type of robots modulates attitudes, emotions, and stereotypes, in *Proc. Tenth Annual ACM/IEEE Int. Conf. Human-Robot Interaction* (Portland, OR, 2015), pp. 109–116.
30. C. D. Frith and U. Frith, The neural basis of mentalizing, *Neuron* **50**(4) (2006) 531–534.
31. J. Ham, R. H. Cuijpers and J. J. Cabibihan, Combining robotic persuasive strategies: The persuasive power of a storytelling robot that uses gazing and gestures, *Int. J. Soc. Robot.* **7**(4) (2015) 479–487.
32. K. S. Haring, K. Watanabe and C. Mougenot, The influence of robot appearance on assessment, in *2013 8th ACM/IEEE Int. Conf. Human-Robot Interaction (HRI)* (IEEE, New York, 2013), pp. 131–132.
33. M. Jeon and I. A. Rayan, The effect of physical embodiment of an animal robot on affective prosody recognition, in *Paper presented at the Int. Conf. Human-Computer Interaction (HCII)* (Orlando, FL, 2011), pp. 523–532.
34. M. F. Jung (2017). Affective grounding in human-robot interaction, in *Paper presented at the 2017 12th ACM/IEEE Int. Conf. Human-Robot Interaction (HRI)* (2017).
35. T. Kishi, T. Kojima, N. Endo, M. Destephe, T. Otani, L. Jamoneand and S. Cosentino, Impression survey of the emotion expression humanoid robot with mental model based dynamic emotions, in *Paper presented at the 2013 IEEE Int. Conf. Robotics and Automation (ICRA)* (IEEE, New York, 2013), pp. 1663–1668.
36. S. Ko, X. Liu, J. Mamros, E. Lawson, H. Swaim, C. Yao and M. Jeon, The effects of robot appearances, voice types, and emotions on emotion perception accuracy and subjective perception on robots, in *Int. Conf. Human-Computer Interaction* (Springer, Cham, 2020), pp. 174–193.
37. M. Kraus, J. Kraus, M. Baumann and W. Minker, Effects of gender stereotypes on trust and likability in spoken human-robot interaction, in *Paper presented at the Proc. Eleventh Int. Conf. Language Resources and Evaluation (LREC-2018)* (Miyazaki, Japan, 2018), pp. 112–118.
38. S. S. Kwak, Y. Kim, E. Kim, C. Shin and K. Cho, What makes people empathize with an emotional robot?: The impact of agency and physical embodiment on human empathy for a robot, in *2013 IEEE RO-MAN* (IEEE, New York, 2013, August), pp. 180–185.
39. J. D. Lewis and A. Weigert, Trust as a social reality, *Soc. Forc.* **63**(4) (1985) 967–985.
40. J. Li, The benefit of being physically present: A survey of experimental works comparing copresent robots, telepresent robots and virtual agents, *Int. J. Hum.-Comput. Stud.* **77** (2015) 23–37.
41. Z. T. Liu, F. F. Pan, M. Wu, W. H. Cao, L. F. Chen, J. P. Xu and M. T. Zhou, A multimodal emotional communication based humans-robots interaction system, in *2016 35th Chinese Control Conf. (CCC)* (IEEE, New York, 2016), pp. 6363–6368.
42. R. Lowe, E. Barakova, E. Billing and J. Broekens, *Grounding Emotions in Robots—An Introduction to the Special Issue* (Sage Publications, London, England, 2016).
43. C. McGinn and I. Torre, Can you tell the robot by the voice? An exploratory study on the role of voice in the perception of robots, in *Paper presented at the 2019 14th ACM/IEEE Int. Conf. Human-Robot Interaction (HRI)* (Daegu, Korea (South), 2019), pp. 211–221.
44. R. L. Mitchell and L. H. Phillips, The overlapping relationship between emotion perception and theory of mind, *Neuropsychologia* **70** (2015) 1–10.

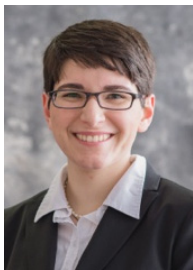


45. S. Nabe, S. J. Cowley, T. Kanda, K. Hiraki, H. Ishiguro and N. Hagita, Robots as social mediators: Coding for engineers, in *Paper presented at the ROMAN 2006- The 15th IEEE Int. Symp. Robot and Human Interactive Communication* (Hatfield, UK, 2006), pp. 384–390.
46. C. Nas, U. Foeh, S. Brave and M. Somoza, The effects of emotion of voice in synthesized and recorded speech, in *Paper presented at the Proc. AAAI Symp. Emotional and Intelligent II: The Tangled Knot of Social Cognition* (2001).
47. S. Nonaka, K. Inoue, T. Arai and Y. Mae, Evaluation of human sense of security for coexisting robots using virtual reality. 1st report: Evaluation of pick and place motion of humanoid robots, in *Paper presented at the IEEE Int. Conf. Robotics and Automation, Proc. ICRA '04*, Vol. 3. (New Orleans, LA, 2004), pp. 2770–2775.
48. A. Ortony, Are all “basic emotions” emotions? a problem for the (basic) emotions construct, *Perspect. Psychol. Sci.* **17**(1) (2021) 41–61. doi: 10.1177/1745691620985415.
49. A. Pereira, C. Martinho, I. Leite and A. Paiva, iCat, the chess player: The influence of embodiment in the enjoyment of a game, in *Proc. 7th Int. Joint Conf. Autonomous Agents and Multiagent Systems*, Vol. 3 (Estoril, Portugal, 2008), pp. 1253–1256.
50. M. L. Phillips, Understanding the neurobiology of emotion perception: Implications for psychiatry, *Br. J. Psychiatr.* **182**(3) (2003) 190–192.
51. R. Plutchik, A general psychoevolutionary theory of emotion, in *Theories of Emotion* (Academic Press, Cambridge, 1980), pp. 3–33.
52. R. Reisenzein, E. Hudlicka, M. Dastani, J. Gratch, K. Hindriks, E. Lorini and J.-J. C. Meyer, Computational modeling of emotion: Toward improving the inter- and intra-disciplinary exchange, *IEEE Trans. Affective Comput.* **4**(3) (2013) 246–266. doi: 10.1109/t-affc.2013.14.
53. E. Rodriguez-Lizundia, S. Marcos, E. Zalama, J. Gómez-García-Bermejo and A. Goraliza, A bellboy robot: Study of the effects of robot behaviour on user engagement and comfort, *Int. J. Human-Comput. Stud.* **82** (2015) 83–95.
54. J. Russell, A circumplex model of emotions, *J. Pers. Soc. Psychol.* **39** (1980) 1161–1178.
55. J. A. Russell (2017). Cross-cultural similarities and differences in affective processing and expression, in *Emotions and Affect in Human Factors and Human-Computer Interaction*, eds. M. Jeon (Academic Press, Cambridge), pp. 123–141.
56. S. Saint-Aimé, B. Le-Pevédic, D. Duhaut and T. Shibata, EmotiRob: Companion robot project, in *Paper presented at the RO-MAN 2007: The 16th IEEE Int. Sym. Robot and Human Interactive Communication* (Jeju, Korean (South), 2007), pp. 919–924.
57. K. Sakai, Y. Nakamura, Y. Yoshikawa and H. Ishiguro, Effect of robot embodiment on satisfaction with recommendations in shopping malls, *IEEE Robot. Autom. Lett.* **7**(1) (2021) 366–372.
58. T. S. S. Schilhab, Anthropomorphism and mental state attribution, in *Animal Behavior* (Academic Press, Cambridge, 2002).
59. A. Schirmer and R. Adolphs, Emotion perception from face, voice, and touch: Comparisons and convergence, *Trends Cogn. Sci.* **21**(3) (2017) 216–228. doi: 10.1016/j.tics.2017.01.001.
60. S. H. Seo, D. Geiskkovitch, M. Nakane, C. King and J. E. Young, Poor thing! Would you feel sorry for a simulated robot? A comparison of empathy toward a physical and a simulated robot, in *Paper presented at the 2015 10th ACM/IEEE Int. Conf. Human-Robot Interaction (HRI)* (Portland, OR, 2015), pp. 125–132.
61. M. Sharma, D. Hildebrandt, G. Newman, J. E. Young and R. Eskicioglu, Communicating affect via flight path: Exploring use of the laban effort system for designing affective locomotion paths, in *Paper presented at the Proc. 8th ACM/IEEE Int. Conf. Human-Robot Interaction* (2013).

62. V. K. Sims, M. G. Chin, H. C. Lum, L. Upham-Ellis, T. Ballion and N. C. Lagattuta, Robots' auditory cues are subject to anthropomorphism, in *Paper presented at the Proc. Human Factors and Ergonomics Society Annual Meeting* (2009).
63. S. Song and S. Yamada, Expressing emotions through color, sound, and vibration with an appearance-constrained social robot, in *Paper presented at the Proc. 2017 ACM/IEEE Int. Conf. Human-Robot Interaction* (Vienna, Austria, 2017), pp. 2–11.
64. J. M. Susskind, D. H. Lee, A. Cusi, R. Feiman, W. Grabski and A. K. Anderson, Expressing fear enhances sensory acquisition, *Nat. Neurosci.* **11**(7) (2008) 843.
65. M. L. Walters, D. S. Syrdal, K. L. Koay, K. Dautenhahn and R. Te Boekhorst, Human approach distances to a mechanical-looking robot with different robot voice styles, in *Paper presented at the RO-MAN 2008-The 17th IEEE Int. Symp. Robot and Human Interactive Communication* (Munich, Germany, 2008), pp. 707–712.
66. A. Waytz, J. Heafner and N. Epley, The mind in the machine: Anthropomorphism increases trust in an autonomous vehicle, *J. Exp. Soc. Psychol.* **52** (2014) 113–117.
67. J. M. G. Williams, F. N. Watts, C. MacLeod and A. Mathews, *Cognitive Psychology and Emotional Disorders* (John Wiley & Sons, New York, 1988).
68. J. O. Wobbrock, L. Findlater, D. Gergle and J. J. Higgins, The aligned rank transform for nonparametric factorial analyses using only ANOVA procedures, in *Paper presented at the Proc. SIGCHI Conf. Human Factors in Computing Systems* (Vancouver, BC, Canada, 2011), pp. 143–146.



**Sangjin Ko** received his MEng in Industrial and Systems Engineering at Virginia Tech in 2020. During his tenure at VT, he was working on the emotive robot voice project, the gesture-based drone controller project, the new auditory display for level 4 automated vehicles project, and the computational modeling of driver behavior in semi-automated vehicle projects. He is currently working in the robot industry as a Technical Project Manager.



**Jaclyn Barnes** received her B.S. in Computer Science from Madonna University in 2014. At Michigan Technological University, she received her M.S. in 2017 and is currently a Ph.D. candidate, both in Computer Science. Her research interests include social robotics, child-robot interaction, assistive technology, and STEAM education.



**Jiayuan Dong** received her B.S. degree in Industrial and Systems Engineering from Virginia Tech in 2021. Since 2021, she is a Ph.D. student in Industrial and Systems Engineering with a concentration in Human Factors from the same university. Her current research project specifically focuses on emotions and trust in human-computer Interaction. She has authored over 7 scholarly publications, including conference proceedings and journal articles. Her research interests include human-robot interaction, human-computer interaction, in-vehicle interface design, and user experience.



**Chung Hyuk Park** received his Ph.D. in Electrical and Computer Engineering (ECE) from the Georgia Institute of Technology in 2012, M.S. in Electrical Engineering and Computer Science (EECS) in 2002 and B.S. in ECE in 2000 from Seoul National University. He is an Associate Professor in the Department of Biomedical Engineering, School of Engineering and Applied Science (SEAS), at the George Washington University (GW) and the director of the Assistive Robotics and Tele-Medicine (ART-Med) Lab. He is also affiliated with the Department of Computer Science and the Autism and Neurodevelopmental Disorders Institute (ANDI) at GW. His current research interests are: 1) Multimodal human-robot interaction and robotic assistance for individuals with disabilities or special needs, 2) Robotic learning and humanized intelligence, and 3) Tele-medical robotic assistance and AI-based reasoning for medical perception and decision-making.



**Ayanna Howard** received her M.S. and Ph.D. degrees in Electrical Engineering from the University of Southern California. She is the Dean of Engineering at The Ohio State University. Previously she was the Chair of the School of Interactive Computing at the Georgia Institute of Technology. Dr. Howard's research encompasses advancements in artificial intelligence (AI), assistive technologies, and robotics, and has resulted in over 275 peer-reviewed publications. She is a Fellow of IEEE, AAI, AAAS, the National Academy of Inventors, and elected member of the American Academy of Arts and Sciences. Prior to Georgia Tech, Dr. Howard was at NASA's Jet Propulsion Laboratory where she held the title of Senior Robotics Researcher and Deputy Manager in the Office of the Chief Scientist.



**Myounghoon Jeon** received his M.S. and Ph.D. degrees in Engineering Psychology and Human-Computer Interaction from the Georgia Institute of Technology, and M.S. degree in Cognitive Science from Yonsei University. He is an Associate Professor in the Department of Industrial and Systems Engineering and the Department of Computer Science (by courtesy) at Virginia Tech, where he is directing the Mind Music Machine Lab. His research interests include human emotions and sound in the application areas of automotive user experiences, assistive robotics, and art and technology integration. His research has resulted in over 250 peer-reviewed publications.