# Establishing Human-Robot Trust through Music-Driven Robotic Emotion Prosody and Gesture

Richard Savery, Ryan Rose and Gil Weinberg[1]

*Abstract*— As human-robot collaboration opportunities continue to expand, trust becomes ever more important for full engagement and utilization of robots. Affective trust, built on emotional relationship and interpersonal bonds is particularly critical as it is more resilient to mistakes and increases the willingness to collaborate. In this paper we present a novel model built on music-driven emotional prosody and gestures that encourages the perception of a robotic identity, designed to avoid uncanny valley. Symbolic musical phrases were generated and tagged with emotional information by human musicians. These phrases controlled a synthesis engine playing back pre-rendered audio samples generated through interpolation of phonemes and electronic instruments. Gestures were also driven by the symbolic phrases, encoding the emotion from the musical phrase to low degree-of-freedom movements. Through a user study we showed that our system was able to accurately portray a range of emotions to the user. We also showed with a significant result that our non-linguistic audio generation achieved an 8% higher mean of average trust than using a state-of-the-art text-to-speech system.

## I. INTRODUCTION

As co-robots become prevalent at home, work, and in public environments, a need arises for the development of trust between humans and robots. A meta-study of human-robot trust [1] has shown that robot-related attributes are the main contributors to building trust in Human-Robot-Interaction, affecting trust more than environmental and human related factors. Related research on artificial agents and personality traits [2], [3] indicates conveying emotions using subtle non-verbal communication channels such as prosody and gesture is an effective approach for building trust with artificial agents. These channels can help convey intentions as well as expressions such as humor, sarcasm, irony, and state-of-mind, which help build social relationship and trust.

In this work we developed new modules for Shimi, a personal robotic platform [4], to study whether emotion-driven non-verbal prosody and body gesture can help establish affective-based trust in HRI. Our approach is to use music, one of the most emotive human experiences, to drive a novel system for emotional prosody [5] and body gesture [6] generation. We propose that music-driven prosody and gesture generation can provide effective low degrees of freedom (DoF) interaction that can convey robotic emotional content, avoid the uncanny valley [7], and help build human-robot trust. Furthermore, trust is highly dictated by the first impression for both human-human and human-robot relations

[1]Georgia Tech Center for Music Technology, Atlanta, GA, USA
rsavery3@gatech.edu

[8], implying methods for gaining trust at the start of a relationship, such as prosody and gesture are crucial.
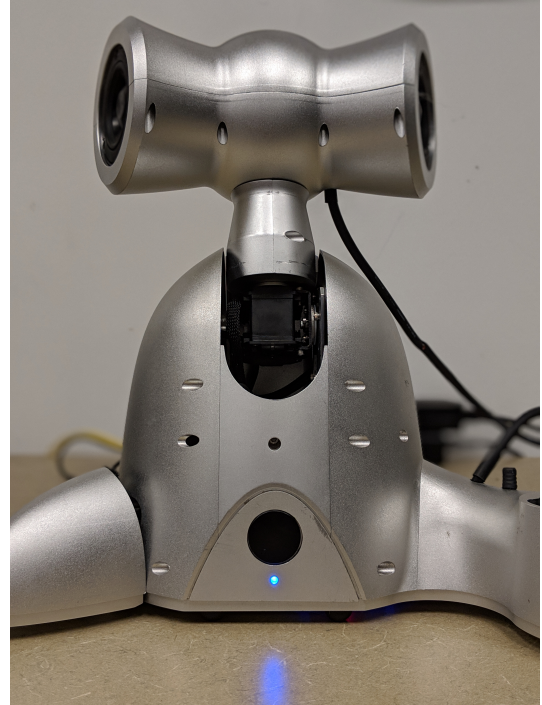


Fig. 1.   The musical robot companion Shimi

We address two research questions, firstly, can we use non-verbal prosody through musical phrases combined with gestures to accurately convey emotions? We present a novel deep learning based musical voice generation system that uses the created phrases to generate gesture through musical prosody. In this way music is central to all interactions presented by Shimi. We evaluate the effectiveness of the musical audio and generative gesture system for Shimi to convey emotion, specified by valence-arousal quadrants. Our second research question is whether emotional conveyance through prosody and gestures driven by music analysis can increase the level of trust in human-robot-interaction. For this we conduct a user study to evaluate prosodic audio and gestures created by our new model in comparison to a baseline text-to-speech system.

## II. BACKGROUND

### A. Trust in HRI

Trust is a key requirement for working with collaborative robots, as low levels of trust can lead to under-utilization

in work and home environments [9]. A key component of the dynamic nature of trust is created in the first phase of a relationship [10], [11], while lack of early trust building can remove the opportunity for trust to develop later on [12]. Lack of trust in robotic systems can also lead to expert operators bypassing the robot to complete tasks [13]. Trust is generally categorized into either *cognitive trust or affective trust* [14]. Affective trust involves emotional bonds and personal relationships, while cognitive trust focuses on considerations around dependability and competence. Perceiving emotion is crucial for the development of affective trust in human-to-human interaction [15], as it increases the willingness to collaborate and expand resources bought to the interactions [16]. Importantly, relationships based on affective trust are more resilient to mistakes by either party [15], and perceiving an emotional identity has been shown to be an important contributor for creating believable and trustworthy interaction [2], [3]. In group interactions, emotional contagion - where emotion is spread between a group - has been shown to improve cooperation and trust in team exercises [17].

### B. Emotion, Music and Prosody

Emotion conveyance is one of the key elements for creating believable agents [2], and prosody has been proven to be an effective communication channel to convey such emotions for humans [18] and robots [19]. On a related front, music which shares many of the underlying building blocks of prosody such as pitch, timing, loudness, intonation, and timbre [20], has also been shown to be a powerful medium to convey emotions [21]. In both music and prosody, emotions can be classified in a discrete categorical manner (happiness, sadness, fear, etc.) [22], and through continuous dimensions such as valence, arousal, and less commonly, dominance, and stance [23], [24]. While some recent efforts to generate and manipulate robotic emotions through prosody focused on linguistic robotic communication [19], [25] no known efforts have been made to use models from music analysis to inform real-time robotic non linguistic prosody for collaboration, as we propose here.

### C. Emotion, Music and Gesture

Human bodily movements are embedded with emotional expression [26], [27], which can be processed by human "affective channels" [28]. Researchers in the field of affective computing, have been working on designing machines that can process and communicate emotions through such channels [29]. Non-conscious affective channels have been demonstrated to communicate compassion, awareness, accuracy, and competency, making them vital components of social interaction with robots [30]. Studies have shown clear correlations between musical features and movement features, suggesting that a single model can be used to express emotion through both music and movement [31], [32], [33]. Certain characteristics of robotic motion have been shown to influence human emotional response [34], [35]. Additionally, emotional intelligence in robots leads to

facilitated human-robot interactions [36], [25], [37]. Gesture has been used to accompany speech in robots to help convey affective information [38], however, to our knowledge, there is no prior work that attempts to integrate physical gestures and music-driven prosody to convey robotic emotional states.

### III. Shimi and Emotion

### A. Prosody

The goal of this project was to create a new voice for Shimi, using musical audio phrases tagged with an emotion. We aimed to develop a new voice that could generate phrases in real-time. This was achieved through a multi-layer system, combining symbolic phrase generation using MIDI controlling a synthesis playback system.
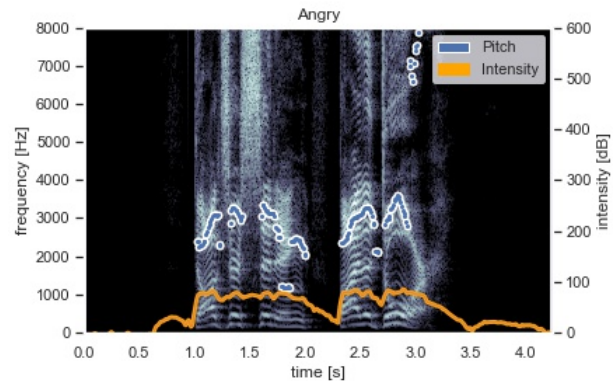


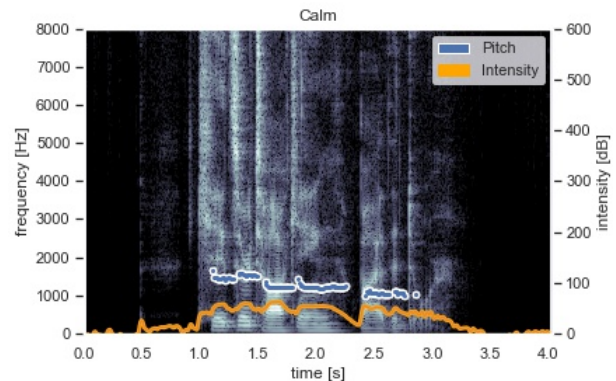Fig. 2.   Angry Speech Pitch and Intensity



Fig. 3.   Calm Speech Pitch and Intensity

*1) Dataset and Phrase Generation:* To control Shimi's vocalizations we generate MIDI phrases that drive the synthesis and audio generation described below and lead the gesture generation. MIDI is a standard music protocol, where notes are stored by pitch value with a note on velocity, followed by a note off to mark the end of the note. With the absence of appropriate datasets we chose to create our own set of MIDI files tagged with valence and arousal by quadrant. MIDI files were collected from eleven different improvisers around the United States, each of whom tagged

their recorded files with an emotion corresponding to a quadrant of the valence/arousal model. Phrases were required to be between 100ms and 6 seconds and each improviser recorded between 50 to 200 samples for each quadrant. To validate this data we created a separate process whereby the pitch range, velocities and contour were compared to the RAVDESS [39] data set, with files removed when the variation was over a manually set threshold. RAVDESS contains speech files tagged with emotion, Figure 2 and Figure 3 clearly demonstrate the variety of prosody details apparent in the RAVDESS dataset (created using [40], [41]) and the variation between a calm and angry utterance of the same phrase.

We chose to use a Recurrent Neural Network, Long Short Term Memory (RNN-LSTM) as described in [42] to generate musical phrases using this data. RNN-LSTM's have been used effectively to generate short melodies, as they are sequential and consider the input as the output is generated. Other network structures were considered however we found RNN-LSTM's very effective for this task when compared to alternate approaches developed by the authors [43], [44], [45], [46], more detail is provided in [47]. While using samples directly from the data was possible, using a neural network allowed infinite variation but also allowed for generated phrases utilizing all the musical features created by the improvisers, not just one improviser per sample.

*2) Audio Creation and Synthesis:* The generated MIDI phrases contain symbolic information only without audio. To create audio we developed a synthesis system to playback the MIDI phrase. As we desired to create a system devoid of semantic meaning a new vocabulary was constructed. This was built upon phonemes from the Australian Aboriginal language Yuwaalaraay a dialect of the Gamilaraay language. Sounds were created by interpolating four different synthesizer sounds with 28 phonemes. Interpolation was done using a modified version of WaveNet. These samples are then time stretched and pitch shifted to match the incoming MIDI file. For a more detailed technical overview of audio processing read [47].

*B. Gestures*

In human communication gestures are tightly coupled with speech [48]. Thus, Shimi's body language is implemented in the same way, derived from its musical prosody and leveraging the musical encoding of emotion to express that emotion physically. Music and movement are correlated, with research finding commonalities in features between both modes [31]. Additionally, humans demonstrate patterns in movement that is induced from music [49]. Particular music-induced movement features are also correlated to perceived emotion in music [50]. After a musical phrase is generated for Shimi's voice to sing, the MIDI representation of that phrase is provided as input to a gesture generation system. Musical features such as tempo, range, note contour, key, and rhythmic density are obtained from the MIDI through Python libraries `pretty_midi` [51] and `music21`[1]. These

[1] https://github.com/cuthbertLab/music21

features are used to create mappings between Shimi's voice and movement: for example, pitch contour is used to govern Shimi's torso forward and backward movement. Other mappings include beat synchronization across multiple subdivisions of the beat in Shimi's foot, and note onset-based movements in Shimi's up-and-down neck movement.

After mapping musical features to low-level movements, Shimi's emotional state is used to condition the actuation of the movements. Continuous values for valence and arousal are used to influence the range, speed, and amount of motion Shimi exhibits. Some conditioning examples include limiting or expanding the range of motion according to the arousal value, and governing how smooth motor direction changes are through Shimi's current valence level. In some cases, the gestures generated for one degree of freedom are dependent on another degree of freedom. For example, when Shimi's torso leans forward, Shimi's attached head will be affected as well. As such, to control where Shimi is looking, any neck gestures need to know the position of the torso. To accommodate these inter-dependencies, when the gesture system is given input, each degree of freedom's movements are generated sequentially and in full, before being actuated together in time with Shimi's voice. Video examples of gesture and audio are available at www.richardsavery.com/shimitrust.

## IV. METHODOLOGY

We designed an experiment to identify how well participants could recognize the emotions shown by our music-driven prosodic and gestural emotion generator. This part of the experiment aimed to answer our first research question, can non-verbal prosody combined with gestures accurately portray emotion. After watching a collection of stimuli, participants completed a survey measuring the trust rating from each participant. This part of the experiment was designed to answer the second question, can emotion driven, non-semantic audio generate trust in a robot.

We hypothesized that through non-semantic prosodic vocalizations accompanied with low-DoF robotic gesture humans will be able to correctly classify Shimi's portrayed emotion as either happy, calm, sad, or angry, with an accuracy consistent with that of text-to-speech. Our second hypothesis was that we will see higher levels of trust from the Shimi using non-speech.

*A. Stimuli*

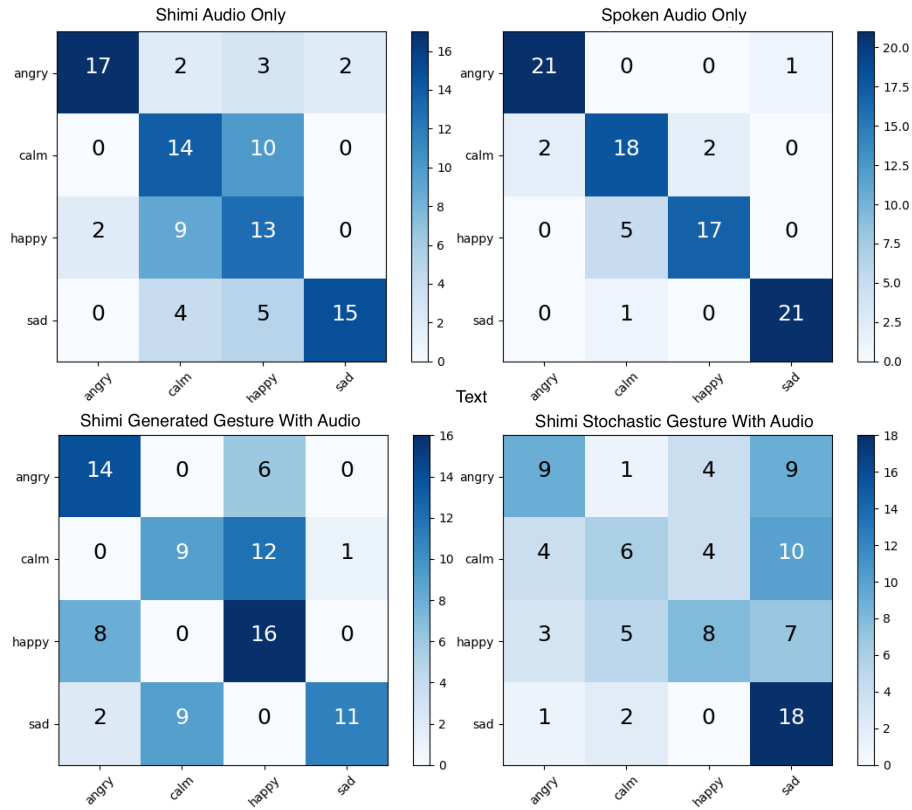| Name | Audio | Stochastic | Experimental |
|---|---|---|---|
| Audio Only | X | | |
| Stochastic Gesture, audio | X | X | |
| Stochastic Gesture, no audio | | X | |
| Experimental Gesture, audio | X | | X |
| Experimental Gesture, no audio | | X | X |

TABLE I
EXPERIMENT STIMULI

Fig. 4. Confusion Matrix

The experiment was designed as a between-subjects study, where one group would hear the audio with the Shimi voice, while the other would hear pre-rendered text-to-speech. Both groups saw the same gesture and answered the same prompts. The text-to-speech examples were synchronized in length and emotion to Shimi's voice. The stimuli for the Speech Audio experiment used CereProc's Meghan voice[2]. CereProc is a state of the art text to speech engine. The text spoken by Meghan was chosen from the EmoInt Dataset [52], which is a collection of manually tagged tweets.

### B. Emotion

The generated gestures were either deterministic gestures created using the previously described system, or deterministic stochastic gestures. Stochastic gestures were implemented by considering each DoF separately, restricting their ranges to those implemented in the generative system, and specifying random individual movement durations up to half of the length of the full gesture. The random number generator used in these gestures were seeded with an identifier unique to the stimuli such that they were deterministic between participants. Gesture stimuli were presented both with and without audio.

### C. Procedure

Participants were gathered from the undergraduate student population at the Georgia Institute of Technology (N=24).

Subjects participated independently, with the group alternating for each participant, culminating with 12 in each group. The session began with an introduction to the task of identifying the emotion displayed by Shimi. Participants responded through a web interface that controlled Shimi through the experiment and then allowed the user to select the emotion they thought Shimi was expressing. Stimuli were randomly ordered for each participant. Table I shows the order of stimuli used, each category contained 8 stimuli, 2 for each valence arousal quadrant. After identifying all stimuli participants were directed to a Qualtrics survey to gather their trust rating.

To measure trust, we used the Trust Perception Scale-HRI [12]. This scale uses 40 questions, each one using a rating scale between 0-100%, to give an average trust rating per participant. The questions take between 5-10 minutes to complete and include questions such as how often the robot will be reliable or pleasant. After completing the trust rating, participants had several open text boxes to discuss any observations in regards to emotion recognition, trust or the general experiment. This was the first time trust was mentioned in the experiment.

## V. RESULTS

### A. Gestures and Emotion

After data was collected, two participant's emotion prediction data was found to be corrupted due to a problem with the testing interface, reducing the number of participants in
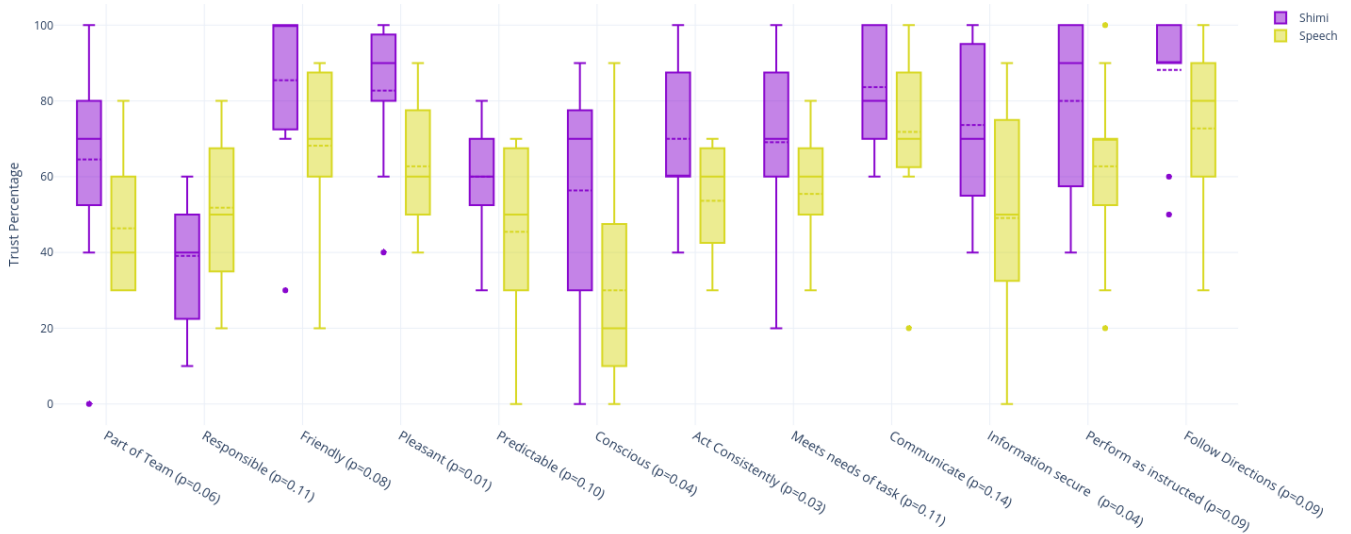
Fig. 5. Questions with P less than 0.1

this portion of the study to 22. First, we considered classification statistics for the isolated predictions of Shimi's voice and text-to-speech (TTS) voice. While TTS outperformed Shimi's voice (F1 score $TTS = 0.87$ vs. $Shimi = 0.63$), the confusion matrices show errant predictions in similar scenarios (see figure 4). For example, both audio classes struggle to disambiguate happy and calm emotions.

Our hope was that adding gestures to accompany the audio would help to disambiguate emotions. To test that our gestures properly encoded emotion, we compared predictions for Shimi's voice accompanied by generated gestures with predictions accompanied by stochastic gestures, the results of which can also be seen in figure 4.

While the confusion matrices show a clear prediction improvement in using generated gestures over stochastic, the results are not statistically significant. A two-sided T-test provides a p-value of 0.089, which does not reject the null hypothesis at $\alpha = 0.05$. Disambiguities from the audio-only cases were not mitigated, but the confused emotions changed slightly, following other gesture and emotion studies [53].

Some experimental error may have accrued through the mixing of stimuli when presented to participants. Each stimuli was expected to be independent but some verbal user feedback expressed otherwise, such as: "the gestures with no audio seemed to be frequently followed by the same gesture with audio, and it was much easier to determine emotion with the presence of audio." The presentation of stimuli may have led participants to choose an emotion based on how we ordered stimuli, rather than their perceived emotion of Shimi.

### B. Trust

As per the trust scale, a mean percentage for trust was calculated on combined answers to 40 questions from each participant. A t-test was then run on each group mean. The average score variation between speech and Shimi audio
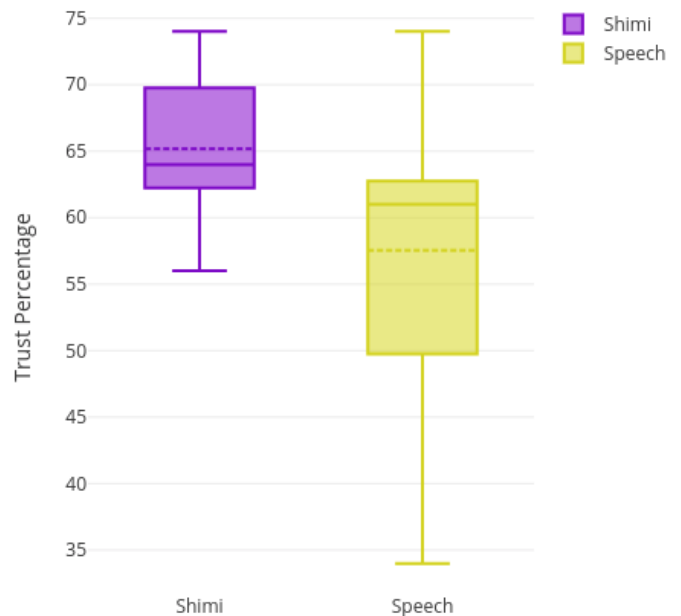


Fig. 6. Participants Trust Mean

showed a significant result (p=0.047), proving the hypothesis. Figure 6 shows the variation in average scores from all participants. The difference of mean between groups was 8%. Results from the text entries were positive for the prosodic voice, and generally neutral or often blank for speech. A common comment from the participants for the Shimi voice was "Seemed like a trustworthy friend that I would be fine confiding in."
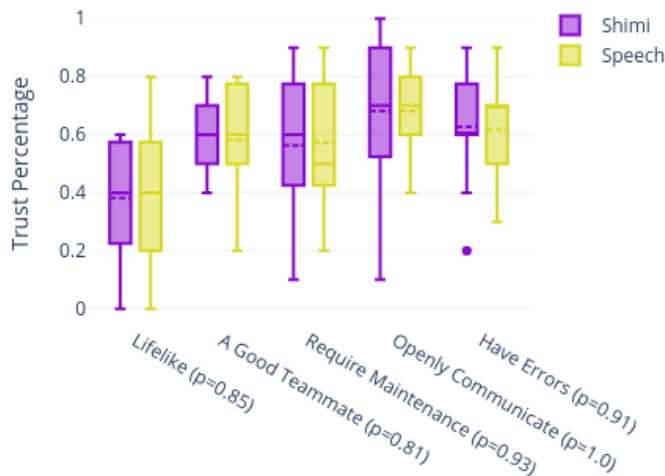
Fig. 7.   Not Significant Trust Results

## VI. Discussion and Future Work

We were able to clearly demonstrate participant recognition of the expected emotion from Shimi, confirming our first hypothesis. Our model however did not perform completely as predicted, as audio without gesture lead to the clearest display of emotion. With a small sample size and a p-value close to being significant, we were encouraged by qualitative feedback that provided insight into the shortcomings of the gestures and gave us ideas for future improvements. For instance, emotions on the same side of the arousal axis were often hard to disambiguate. One participant noted that "it was generally difficult to distinguish happy and angry if there were no sounds (similar situation between sad and calm)", while another noted "I had some trouble discerning calm from sad here and there", and "without speaking, it was difficult to decipher between anger and excitement". The general intensity of the emotion was apparent, however." Certain movement features led to emotional connections for the participants, as demonstrated here: "generally, when Shimi put it's head down, I was inclined to say it looked sad. When it moved more violently, particularly by tapping [sic] it's foot, I was inclined to say it was angry or happy", "more forceful movements tended to suggest anger", and "When there was more severe motion, I associated that with anger. When the motion was slower I associated it with sad or calm. If the head was down more I associated it with sad. And I associated it with happy more when there was sound and more motion."

The trust perception scale is designed to give an overall rating, and independent questions should not necessarily be used to draw conclusions. However, there were several interesting results indicating further areas of research. Fig 5 shows all categories with a p value less than 0.10, for which multiple questions showed significant results with a p value under 0.05). Shimi's voice was crafted to be friendly and inviting and as expected received much higher results for pleasantness and friendliness. Unexpectedly, it also showed

much higher ratings for its perception as being conscious. While further research is required to confirm the meaning, we believe that the question on consciousness of Shimi demonstrating a significant result shows that embodying a robot with a personal prosody (as opposed to human speech) creates a more believable agent. Figure 7 shows the categories with very similar distributions of scores. These include Lifelike, A Good Teammate, Have Errors, Require Maintenance and Openly Communicate. While further research is needed, this may imply that these features are not primarily associated with audio. Further work should be done to explore if the same impact can be found by adjusting audio features of a humanoid robot may also lead to interesting results.

In other future work we plan to develop experiments with a broader custom musical data-set across multiple robots. We intend to study emotional contagion and trust between larger groups of robots across distributed networks [54], aiming to understand collaboration and trust at a higher level between multiple robots.

Overall, our trust results were significant and showed that prosody and gesture can be used to generate higher levels of trust in human-robot interaction. Our belief that creating a believable agent that avoided uncanny valley was shown to be correct and was validated through participant comments, including the open text response: "Shimi seems very personable and expressive, which helps with trust".

## References

[1] P. A. Hancock, D. R. Billings, K. E. Schaefer, J. Y. Chen, E. J. De Visser, and R. Parasuraman, "A meta-analysis of factors affecting trust in human-robot interaction," *Human Factors*, vol. 53, no. 5, pp. 517–527, 2011.

[2] M. Mateas, "Artificial intelligence today," M. J. Wooldridge and M. Veloso, Eds. Berlin, Heidelberg: Springer-Verlag, 1999, ch. An Oz-centric Review of Interactive Drama and Believable Agents, pp. 297–328. [Online]. Available: http://dl.acm.org/citation.cfm?id=1805750.1805762

[3] J. Bates, "The role of emotion in believable agents," *Commun. ACM*, vol. 37, no. 7, pp. 122–125, Jul. 1994. [Online]. Available: http://doi.acm.org/10.1145/176789.176803

[4] M. Bretan, G. Hoffman, and G. Weinberg, "Emotionally expressive dynamic physical behaviors in robots," *International Journal of Human-Computer Studies*, vol. 78, pp. 1–16, 2015.

[5] R. Adolphs, D. Tranel, and H. Damasio, "Emotion recognition from faces and prosody following temporal lobectomy." *Neuropsychology*, vol. 15, no. 3, p. 396, 2001.

[6] C. Shan, S. Gong, and P. W. McOwan, "Beyond facial expressions: Learning human emotion from body gestures." in *BMVC*, 2007, pp. 1–10.

[7] K. F. MacDorman, "Subjective ratings of robot video clips for human likeness, familiarity, and eeriness: An exploration of the uncanny valley," in *ICCS/CogSci-2006 long symposium: Toward social mechanisms of android science*, 2006, pp. 26–29.

[8] J. Xu and A. Howard, "The impact of first impressions on human-robot trust during problem-solving scenarios," in *2018 27th IEEE International Symposium on Robot and Human Interactive Communication (RO-MAN)*, Aug 2018, pp. 435–441.

[9] J. D. Lee and K. A. See, "Trust in automation: Designing for appropriate reliance," *Human Factors*, vol. 46, no. 1, pp. 50–80, 2004.

[10] P. H. Kim, K. T. Dirks, and C. D. Cooper, "The repair of trust: A dynamic bilateral perspective and multilevel conceptualization," *Academy of Management Review*, vol. 34, no. 3, pp. 401–422, 2009.

[11] R. E. Miles and W. D. Creed, "Organizational forms and managerial philosophies-a descriptive and analytical review," *RESEARCH IN ORGANIZATIONAL BEHAVIOR: AN ANNUAL SERIES OF ANALYTICAL ESSAYS AND CRITICAL REVIEWS, VOL 17, 1995*, vol. 17, pp. 333–372, 1995.

[12] K. E. Schaefer, *Measuring Trust in Human Robot Interactions: Development of the "Trust Perception Scale-HRI"*. Boston, MA: Springer US, 2016, pp. 191–218.

[13] P. M. Satchell, *Cockpit monitoring and alerting systems*. Routledge, 2016.

[14] A. Freedy, E. DeVisser, G. Weltman, and N. Coeyman, "Measurement of trust in human-robot collaboration," in *Collaborative Technologies and Systems, 2007. CTS 2007. International Symposium on*. IEEE, 2007, pp. 106–114.

[15] D. M. Rousseau, S. B. Sitkin, R. S. Burt, and C. Camerer, "Not so different after all: A cross-discipline view of trust," *Academy of management review*, vol. 23, no. 3, pp. 393–404, 1998.

[16] T. Gompei and H. Umemuro, "Factors and development of cognitive and affective trust on social robots," in *International Conference on Social Robotics*. Springer, 2018, pp. 45–54.

[17] S. G. Barsade, "The ripple effect: Emotional contagion and its influence on group behavior," *Administrative science quarterly*, vol. 47, no. 4, pp. 644–675, 2002.

[18] Y. T. Wang, J. Han, X. Q. Jiang, J. Zou, and H. Zhao, "Study of speech emotion recognition based on prosodic parameters and facial expression features," in *Applied Mechanics and Materials*, vol. 241. Trans Tech Publ, 2013, pp. 1677–1681.

[19] J. Crumpton and C. L. Bethel, "A survey of using vocal prosody to convey emotion in robot speech," *International Journal of Social Robotics*, vol. 8, no. 2, pp. 271–285, 2016.

[20] A. Wennerstrom, *The music of everyday speech: Prosody and discourse analysis*. Oxford University Press, 2001.

[21] J. Sloboda, "Music: Where cognition and emotion meet," in *Conference Proceedings: Opening the Umbrella; an Encompassing View of Music Education; Australian Society for Music Education, XII National Conference, University of Sydney, NSW, Australia, 09-13 July 1999*. Australian Society for Music Education, 1999, p. 175.

[22] L. Devillers, L. Vidrascu, and L. Lamel, "Challenges in real-life emotion annotation and machine learning based detection," *Neural Networks*, vol. 18, no. 4, pp. 407–422, 2005.

[23] J. A. Russell, "Emotion, core affect, and psychological construction," *Cognition and Emotion*, vol. 23, no. 7, pp. 1259–1283, 2009.

[24] A. Mehrabian, "Pleasure-arousal-dominance: A general framework for describing and measuring individual differences in temperament," *Current Psychology*, vol. 14, no. 4, pp. 261–292, 1996.

[25] C. Breazeal and L. Aryananda, "Recognition of affective communicative intent in robot-directed speech," *Autonomous robots*, vol. 12, no. 1, pp. 83–104, 2002.

[26] M. Inderbitzin, A. Vljame, J. M. B. Calvo, P. F. M. J. Verschure, and U. Bernardet, "Expression of emotional states during locomotion based on canonical parameters," in *Ninth IEEE International Conference on Automatic Face and Gesture Recognition (FG 2011), Santa Barbara, CA, USA, 21-25 March 2011*. IEEE, 2011, pp. 809–814.

[27] H. G. Walbott, "Bodily expression of emotion," *European Journal of Social Psychology*, vol. 28, no. 6, pp. 879 – 896, 1998.

[28] B. de Gelder, "Towards the neurobiology of emotional body language," *Nature Reviews Neuroscience*, vol. 7, pp. 242–249, March 2006.

[29] R. W. Picard, "Affective computing," 1995.

[30] M. Scheutz, P. Schermerhorn, and J. Kramer, "The utility of affect expression in natural language interactions in joint human-robot tasks," in *Proceedings of the 1st ACM SIGCHI/SIGART conference on Human-robot interaction*. ACM, 2006, pp. 226–233.

[31] B. Sievers, L. Polansky, M. Casey, and T. Wheatley, "Music and movement share a dynamic structure that supports universal expressions of emotion," *Proceedings of the National Academy of Sciences*, vol. 110, no. 1, pp. 70–75, jan 2013. [Online]. Available: http://www.pnas.org/content/110/1/70

[32] G. Weinberg, A. Beck, and M. Godfrey, "Zoozbeat: a gesture-based mobile music studio." in *NIME*, 2009, pp. 312–315.

[33] G. Weinberg, S. Driscoll, and T. Thatcher, "Jamaa: a percussion ensemble for human and robotic players," in *ACM International Conference on Computer Graphics and Interactive Techniques (SIGGRAPH 2006), ACM Boston, MA*, 2006.

[34] L. D. Riek, T.-C. Rabinowitch, P. Bremner, A. G. Pipe, M. Fraser, and P. Robinson, "Cooperative gestures: Effective signaling for humanoid robots," in *Human-Robot Interaction (HRI), 2010 5th ACM/IEEE International Conference on*. IEEE, 2010, pp. 61–68.

[35] A. Moon, C. A. Parker, E. A. Croft, and H. Van der Loos, "Design and impact of hesitation gestures during human-robot resource conflicts," *Journal of Human-Robot Interaction*, vol. 2, no. 3, pp. 18–40, 2013.

[36] H. Kozima and H. Yano, "In search of otogenetic prerequisites for embodied social intelligence," in *Proceedings of the Workshop on Emergence and Development on Embodied Cognition; International Conference on Cognitive Science*, 2001, pp. 30–34.

[37] G. Castellano, I. Leite, A. Pereira, C. Martinho, A. Paiva, and P. W. McOwan, "Affect recognition for interactive companions: challenges and design in real world scenarios," *Journal on Multimodal User Interfaces*, vol. 3, no. 1, pp. 89–98, 2010.

[38] M. W. Alibali, S. Kita, and A. J. Young, "Gesture and the process of speech production: We think, therefore we gesture," *Language and cognitive processes*, vol. 15, no. 6, pp. 593–613, 2000.

[39] S. R. Livingstone and F. A. Russo, "The Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS): A dynamic, multimodal set of facial and vocal expressions in North American English," *PLOS ONE*, vol. 13, no. 5, pp. 1–35, 2018. [Online]. Available: https://doi.org/10.1371/journal.pone.0196391

[40] Y. Jadoul, B. Thompson, and B. de Boer, "Introducing parselmouth: A python interface to praat," *Journal of Phonetics*, vol. 91, pp. 1–15, 11 2018.

[41] P. Boersma and D. Weenink, "Praat: Doing phonetics by computer (version 5.3.51)," 01 2007.

[42] R. Savery and G. Weinberg, "Shimon the robot film composer and deepscore," *Proceedings of Computer Simulation of Musical Creativity*, p. 5, 2018.

[43] M. Bretan, G. Weinberg, and L. Heck, "A unit selection methodology for music generation using deep neural networks," *arXiv preprint arXiv:1612.03789*, 2016.

[44] R. J. Savery, "An interactive algorithmic music system for edm," *Dancecult: Journal of Electronic Dance Music Culture*, vol. 10, no. 1, 2018.

[45] ——, "Algorithmic improvisers," in *UC Irvine Electronic Theses and Dissertations*, 2015.

[46] R. Savery, M. Ayyagari, K. May, and B. N. Walker, "Soccer sonification: Enhancing viewer experience." Georgia Institute of Technology, 2019.

[47] R. Savery, R. Rose, and G. Weinberg, "Finding Shimi's voice: fostering human-robot communication with music and a NVIDIA Jetson TX2," *Proceedings of the 17th Linux Audio Conference*, p. 5, 2019.

[48] D. McNeill, *How language began: Gesture and speech in human evolution*. Cambridge University Press, 2012.

[49] P. Toiviainen, G. Luck, and M. R. Thompson, "Embodied Meter: Hierarchical Eigenmodes in Music-Induced Movement," *Music Perception: An Interdisciplinary Journal*, vol. 28, no. 1, pp. 59–70, sep 2010. [Online]. Available: http://mp.ucpress.edu/content/28/1/59

[50] B. Burger, S. Saarikallio, G. Luck, M. R. Thompson, and P. Toiviainen, "Relationships Between Perceived Emotions in Music and Music-induced Movement," *Music Perception: An Interdisciplinary Journal*, vol. 30, no. 5, pp. 517–533, Jun. 2013. [Online]. Available: http://mp.ucpress.edu/cgi/doi/10.1525/mp.2013.30.5.517

[51] C. Raffel and D. P. W. Ellis, "Intuitive analysis, creation and manipulation of midi data with pretty_midi," in *Proceedings of the 15th International Conference on Music Information Retrieval Late Breaking and Demo Papers*, 2014.

[52] S. M. Mohammad and F. Bravo-Marquez, "WASSA-2017 shared task on emotion intensity," in *Proceedings of the Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis (WASSA)*, Copenhagen, Denmark, 2017.

[53] A. Lim, T. Ogata, and H. G. Okuno, "Towards expressive musical robots: a cross-modal framework for emotional gesture, voice and music," *EURASIP Journal on Audio, Speech, and Music Processing*, vol. 2012, no. 1, p. 3, Jan. 2012. [Online]. Available: https://doi.org/10.1186/1687-4722-2012-3

[54] G. Weinberg, "Expressive digital musical instruments for children," Ph.D. dissertation, Massachusetts Institute of Technology, 1999.