

# Speech-Gesture GAN: Gesture Generation for Robots and Embodied Agents

Carson Yu Liu<sup>1</sup>, Gelareh Mohammadi<sup>1</sup>, Yang Song<sup>1</sup> and Wafa Johal<sup>2</sup>

**Abstract**—Embodied agents, in the form of virtual agents or social robots, are rapidly becoming more widespread. In human-human interactions, humans use nonverbal behaviours to convey their attitudes, feelings, and intentions. Therefore, this capability is also required for embodied agents in order to enhance the quality and effectiveness of their interactions with humans. In this paper, we propose a novel framework that can generate sequences of joint angles from the speech text and speech audio utterances. Based on a conditional Generative Adversarial Network (GAN), our proposed neural network model learns the relationships between the co-speech gestures and both semantic and acoustic features from the speech input. In order to train our neural network model, we employ a public dataset containing co-speech gestures with corresponding speech audio utterances, which were captured from a single male native English speaker. The results from both objective and subjective evaluations demonstrate the efficacy of our gesture-generation framework for Robots and Embodied Agents.

## I. INTRODUCTION

As a result of the ongoing improvement of humanoid robots and computer graphics, conversational embodied agents, including social robots and virtual agents, have emerged as effective interaction instrumentality. The ESI (Evaluation of Social Interaction) [1], a human evaluation instrument, identifies important social skills such as approaching, speaking, turn-taking, gazing and gesturing. Therefore, in human-agent interactions, social agents also need these social capabilities similar to humans. In particular, human gestures are a form of nonverbal cues utilised with utterances in interpersonal interaction. Secondly, researchers revealed that in certain cultures, speech and gestures are tightly linked in time [2]. Therefore, it is crucial to create gestures and briefly integrate them with speech while designing embodied agents. In fact, the danger for embodied agents is a mismatch between verbal and nonverbal information, which may cause extraordinary unpleasantness to the communicators [3]. Thirdly, gestures may be used to emphasize words, demonstrate purpose, depict things more vividly, and aid understanding of a conversation [4]. In human-robot interaction, it has been discovered that common language gestures strengthen the robot’s attraction and prospective contact motivation [5]. However, considering the diversity

of embodied agents and the physical limits of robots, it does not seem feasible to manually create gestures for each possible speech. Linguists [6] suggested a categorisation system with four classes: 1) Iconic (expressing an object’s features or behaviours); 2) Deictic, or pointing (indicating an object’s position); 3) Metaphoric (representing abstract concepts with a concrete form); 4) Beat (keeping with the rhythm of speech). Only the Beat gestures are audio signal-dependent (speech acoustic), while other types of gestures rely on the speech context (speech semantic). Therefore, gesture generation frameworks with single modal input can lack some types of gestures.

Motivated by the accomplishments of GAN (Generative Adversarial Network) [7] in generative models, we propose a GAN-structured neural network model to generate gestures from speech. We trained our model on a gesture dataset with English speech. The subjective evaluation demonstrates the proposed model is effective, showing a good performance when compared with the ground truth. Also, the objective evaluation results confirm our model is highly effective when compared with other state-of-the-art gesture generation models.

The contribution of our work is two-fold: 1) We propose a novel GAN-based generative framework that can use multimodal inputs to extract semantic and acoustic features as conditional information for adversarial training and generate multiple gestures from the same speech input using different input noises. 2) A comprehensive evaluation of the full model from objective to subjective with ablation studies of the outcomes of various designs and crucial modelling options;

The rest of this paper is organised as follows: We first introduce the background and related work in Section II. Then, Section III describes our proposed speech-based gesture generation framework, including features extraction, model architecture and its implementation. Next, sections IV and V explain the quantitative metrics used in our proposed model and quantitative result with validation on an extra user study. Finally, we conclude our work with a brief discussion.

## II. RELATED WORK

Several gesture generation approaches, ranging from rule-based to innovative data-driven, have been created in recent years. Initially, most approaches were rule-based; however, rule-based approaches result in a repetitious and monotonous experience in the lifetime human-agent connection. Recent innovative approaches are data-driven, enabling more variety in gesture production but making it more difficult to adapt to the physical limits of the embodied agents.

<sup>1</sup>Carson Yu Liu, Gelareh Mohammadi and Yang Song are with the Faculty of Engineering, School of Computer Science and Engineering, University of New South Wales, Sydney, NSW 2052, Australia carson.liu@unsw.edu.au; g.mohammadi@unsw.edu.au; yang.song1@unsw.edu.au

<sup>2</sup>Wafa Johal is with the Faculty of Engineering and Information Technology, School of Computing & Information Systems, University of Melbourne, Melbourne, VIC 3010, Australia wafa.johal@unimelb.edu.au

### A. Rule-based gesture generation

The primary concept behind rule-based generation approaches is to correlate speech syllables, and words with gestures as a straightforward way to produce gestures from speech content [8]. The rules for generating gestures in these studies were hand-defined by specialists [9], [10], [11]. One study [9] derived punctuation marks from a sentence using a dialogue sentence analysis methodology. Using image processing and clustering approaches, unique research [12] produced its own gestures dictionary for speech gestures from internet images, although the processing of gesture production is still governed by rules. For rule-based gesture generation, the greatest drawback is that manually defining a gesture pattern for each word requires an enormous amount of time and effort. By utilising the machine learning technique, the issue of repeated and labor-intensive generation of a speech gestures dictionary might be addressed.

### B. Data-driven gesture generation

Recent Data-driven studies focus on learning mapping functions from speech text or speech audio or both of them to speech gestures.

1) *Gesture Generation with Speech Text*: Yoon et al. [13] presented a seq2seq-based autoencoder model which employed speech text as input to generate 2D co-gestures; they also implemented their model on the NAO robot. Another work [14] also extracted speech text features as input for their probabilistic model. However, both of them observed an unusual mapping issue in which the synthesised audio and produced gestures could not be closely synchronised.

2) *Gesture Generation with Speech Audio*: Hasegawa et al. [15] extracted the MFCCs (Mel-Frequency Cepstral Coefficients) from the inputted audio as the speech representation; they used a bi-directional LSTM (Long Short-Term Memory) based recurrent neural network to generate co-speech gestures and then went through a noise filter as smoothing step. With the same speech gesture database, Kucherenko et al. [16] presented an autoencoder, which is used for representation learning to align the audio with gestures. Ferstl et al. [17] also used bi-directional LSTM regression with adversarial training to generate gestures from acoustic features (MFCCs with audio pitch); they also utilized multiple discriminators in adversarial training to improve the results from the generator. Our proposed approach varies from prior systems in that it generates co-speech gestures using text transcription and audio utterances.

3) *Gesture Generation with Multimodal Input*: Single modality systems have clear limitations; as mentioned before, the lack of either acoustic or semantic features resulting from the single modal input is currently a significant hurdle to achieving outstanding results. However, multimodality systems could address this problem. Kucherenko et al. [18] proposed the first multimodal input autoregressive neural network model on co-speech gesture generation. Yoon et al. [19] added speaker identity as third modal input to achieve style control.

## III. PROPOSED SPEECH-BASED GESTURE GENERATION FRAMEWORK

### A. Speech and Gesture Dataset

Unlike previous studies that used non-English gesture datasets [20], small gesture datasets [21], datasets with low-quality gestures [13] or multi-language datasets [22] our proposed speech-based gesture generation framework is specifically trained with the Trinity Dataset [23], that captured from a single male actor who is an English native speaker with 20 Vicon cameras (a sort of motion capture cameras). This gesture dataset contains 244 minutes of speech and gesture data from among a variety of topics, e.g., daily activities, hobbies and movies. First, we removed lower body data, because our work is aiming at co-speech gestures. Then, in order to save training time, for the upper body data, we used 4 joints from the spine, 2 joints from the neck, 3 joints from the left and right side arm, 2 joints from both side shoulder, 1 joint from the head. In addition, the fingers data is removed due to two reasons: 1) poor quality of data and 2) many common humanoid robots like NAO and Pepper from Softbank Robotics do not have enough fingers like human beings. Finally, we have speech audio utterances in the form of 44 kHz Waveform Audio file format, and speech text transcripts in the form of JavaScript Object Notation file format with timestamps and corresponding gestures in the form of Biovision Hierarchy file format.

### B. Data Pre-processing and Feature Extraction

Based on the experiments of previous work[18], we employ frame synchronization at 20 FPS (Frames Per Second) during feature extraction. The gesture data in Biovision Hierarchy consist of Euler angles and offsets of each joint in a hierarchical structure. Unlike previous studies that adopted conversion of Euler angles and absolute position in 3D coordinates [24], we converted Euler angles to exponential maps [25] because it is easy to convert exponential maps back to Euler angles and will not introduce potential discontinuities issues. After frame conversion from 60 FPS to 20 FPS, we get 45 features for each frame of gestures.

As for the acoustic features extraction, similar to other state-of-the-art in speech-based gesture generation[15], [16], [26], [27], in order to align with the gesture features, we get feature vectors in 26 dimensions (for 26 Mel-spaced filterbanks) by calculating the MFCCs of the audio utterances waveform with the same frame rate, which is a representation of an utterance's short term power spectrum.

However, a sequence of the speech audio utterances and its corresponding speech text transcripts generally have different lengths. In order to address this problem, we first encoded the words with semantic information as 768-dimensional vectors by using the BERT[28] pre-trained model, a state-of-the-art neural network model that uses surrounding text to assist computers in grasping the meaning of ambiguous words in a text. As for the words that do not have semantic information, we encoded them as fixed vectors that have the same dimensions as the BERT features. Then, we used the exact

utterance time information of each word to upsample the text features. Therefore, the text and audio feature sequences get aligned and uniform.

### C. Problem Formulation

The problem of the co-gesture generation from speech can be defined as a mapping function  $\mathbf{F}_{Generation}$ , which is shown in Equation 1 for a segment of the input speech length  $T$ , where  $\mathbf{s}_a = [s_a]_{t=1:T}$  are the features extracted from speech audio utterances. Likewise, the features extracted from speech text are  $\mathbf{s}_t = [s_t]_{t=1:T}$ , with multiple noise  $\mathbf{n}$ . The corresponding result  $\mathbf{g} = [g_{t=1:T}]$  can be a sequence of Euler angles of selected joints in the form of  $g_t = [pitch_t^i, raw_t^i, yaw_t^i]_{i=1:J}$ , where  $J$  is the number of selected joints. Furthermore, we define  $\mathbf{g} = [g_{t=1:T}]$  as a sequence of 3D (three-dimensional) coordinates of selected joints, with  $g_t = [x_t^i, y_t^i, z_t^i]_{i=1:J}$ . The object of our problem is to achieve the maximization of the conditional probability  $p(\mathbf{g}|\mathbf{s})$  to match well with the given speech input, where  $\mathbf{s}$  is the concatenation of  $\mathbf{s}_a$  and  $\mathbf{s}_t$ .

$$\mathbf{g} = \mathbf{F}_{Generation}(\mathbf{s}_a, \mathbf{s}_t, \mathbf{n}) \quad (1)$$

### D. Model Architecture

Speech features extracted from audio utterances and text transcripts are used as the condition in our proposed model, which is a conditional GAN-based architecture. Figure 1 shows the overview of the architecture.

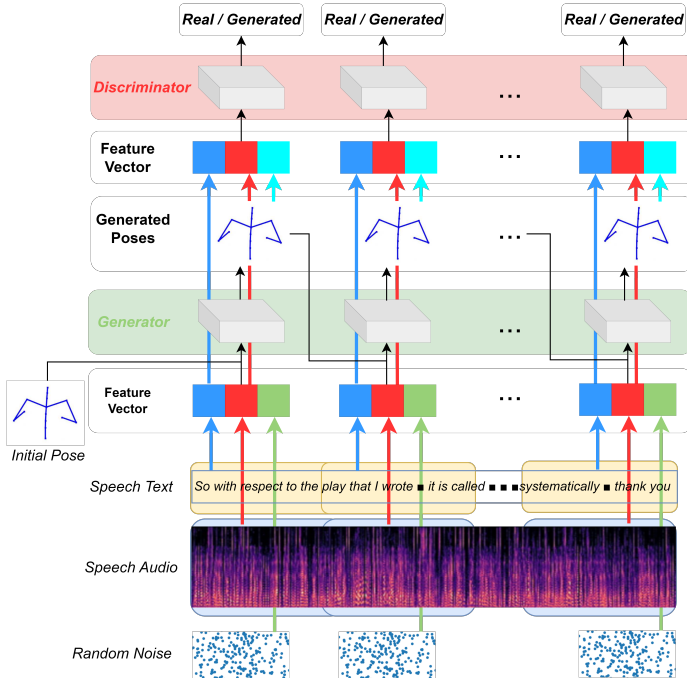


Fig. 1. The architecture of the proposed gesture generation model

In the generation step, a random noise  $\mathbf{n}$  from a normal distribution is reproduced in the same length as the speech features. Then, the noise, the text embeddings  $\mathbf{s}_t$  and the MFCCs values  $\mathbf{s}_a$  are concatenated as the feature vector,

take it into the generator to get the corresponding sequence of gestures. Specifically, we employed the initial pose for the previous frames to improve the continuity during gesture generation. In order to improve the generator, we concurrently trained the discriminator to calculate the difference between the real distribution and fake distribution on the speech features condition. Next, after getting the sequence of generated gestures, we concatenated the generated gestures or real gestures with accompanying audio and semantic features and then sent them into the discriminator. The output value shows if the input gestures were real or fake for the corresponding speech features condition.

### E. Gesture Generator

Our gesture generator  $G$  generates gestures using encoded semantic and acoustic features as input. The structure of the generator  $G$  is shown in Figure 2.

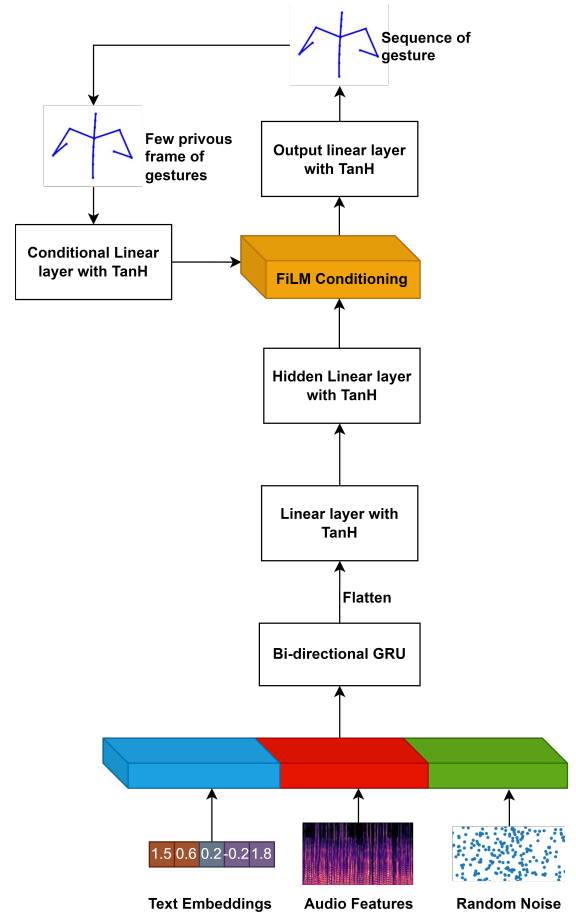


Fig. 2. Gesture Generator

First, we concatenate the text embedding, MFCCs and random noise as a long vector, then send them through to the two-layer bi-direction GRU (Gated recurrent unit) with 0.2 dropouts. Next, the vector passes through the following linear layer with the TanH activation function to reduce the dimensionality of the feature. In order to ensure the continuity of generated gestures, we used the few frames of previously generated gestures as condition information

TABLE I  
GENERATOR ARCHITECTURE

Layer	Layer Type	Hyperparameter
1	GRU	$C_{in} = 814, C_{out} = 128, L_{num} = 2$
2	Linear	$C_{in} = 3840, C_{out} = 512$
3	Linear	$C_{in} = 135, C_{out} = 512$
4	Linear	$C_{in} = 512, C_{out} = 256$
5	Linear	$C_{in} = 256, C_{out} = 45$

to feed back to the FiLM (Feature-wise Linear Modulation) layer [29], as another state-of-the-art work [18] did. Finally, the output layer is a linear layer with the TanH activation function to get a possible range of results. The layers detail of the gesture generator is shown in Table I, where  $C_{in}$ ,  $C_{out}$  are dimensions of in and out channels, and  $L_{num}$  is the number of GRU layers.

### F. Adversarial Scheme

In order to optimize our gesture generator, a discriminator  $D$  is used in our adversarial scheme. Figure 3 illustrates the structure of our discriminator. First, the sequence of generated gestures from the generator, text embeddings and MFCCs both individually go through two linear layers: one with the Leaky ReLU activation function and the next one without the activation function. Inspired by the work [30], we take the vector of the concatenated gestures, audio and text features and then feed them into five layers of the 1D convolutional block, which consists of one 1D convolutional layer with Leaky ReLU and layer normalization, finally followed by an extra 1D convolutional layer. Next, the vector passes through two linear layers with Leaky ReLU for vector dimensional reduction. At the end of the discriminator, using a sigmoid activation function, the result is compressed between 0 and 1. These values could determine if the input gestures are real and well-matched with the condition features. The layers detail of the discriminator is shown in Table II, where  $k$ ,  $s$  and  $p$  are kernel size, stride and padding, respectively.

### G. Training

The losses listed below are used to train the proposed framework. The gesture generator is trained by using the loss  $\mathbf{L}_G$  in Equation 2, while the loss  $\mathbf{L}_D$  in Equation 6 is used for training the discriminator.

$$\mathbf{L}_G = \alpha \cdot \mathbf{L}_G^{mse} + \beta \cdot \mathbf{L}_G^{continuity} + \lambda \cdot \mathbf{L}_G^{WGAN} \quad (2)$$

$$\mathbf{L}_G^{mse} = \frac{1}{n} \sum_{i=1}^n (\mathbf{g}_i - \hat{\mathbf{g}}_i)^2 \quad (3)$$

$$\mathbf{L}_G^{continuity} = \frac{1}{n} \sum_{i=1}^n (\mathbf{S}_i - \hat{\mathbf{S}}_i)^2 \quad (4)$$

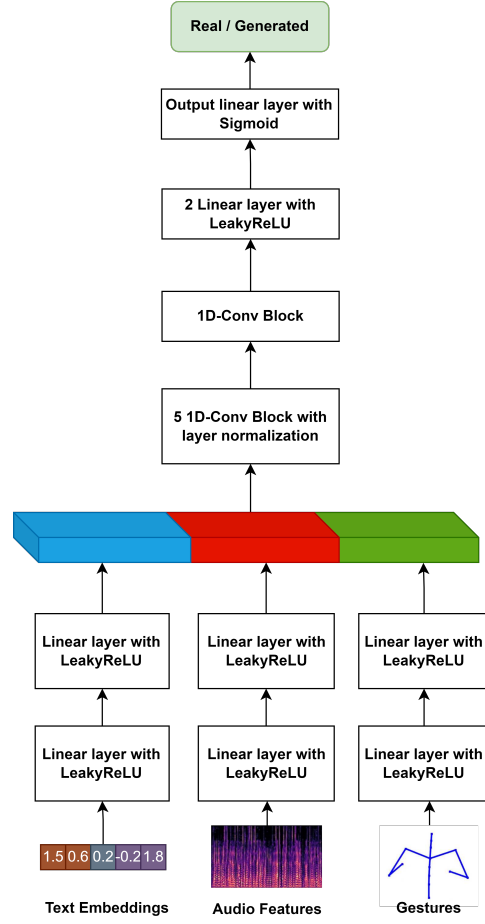


Fig. 3. Discriminator

$$\mathbf{L}_G^{WGAN} = -\frac{1}{N} \sum_{i=1}^n D(\mathbf{s}_a, \mathbf{s}_t, \hat{\mathbf{g}}_i) \quad (5)$$

$$\mathbf{L}_D = \frac{1}{N} \sum_{i=1}^n D(\mathbf{s}_a, \mathbf{s}_t, \hat{\mathbf{g}}_i) - \frac{1}{N} \sum_{i=1}^n D(\mathbf{s}_a, \mathbf{s}_t, \mathbf{g}_i) \quad (6)$$

Where  $\mathbf{s}_a$ ,  $\mathbf{s}_t$  represent the speech audio and text features, respectively. Specifically,  $n$  is the total duration of the gesture sequence,  $\mathbf{g}_i$  and  $\hat{\mathbf{g}}_i$  are the  $i$ th original gesture and  $i$ th generated gesture, respectively. Using MSE (mean squared error) in Equation 3 and continuity loss in Equation 4, we reduced the gap between original gestures in training samples and the matching generated gestures while training our gesture generator. This loss  $\mathbf{L}_G^{continuity}$  can be construed as the mean squared error for the current speed difference of  $i$ th original gesture speed  $\mathbf{S}_i$  and  $i$ th generated gesture speed  $\hat{\mathbf{S}}_i$ . The adversarial losses  $\mathbf{L}_G^{WGAN}$  in the Equation 5, where  $G$  is the generator and  $\mathbf{L}_D$  where  $D$  is discriminator come from the WGAN (Wasserstein Generative Adversarial Networks) [31], an improved generative model that makes the training more stable when compared with the traditional GAN model. As in the GAN training,  $\mathbf{L}_G$  and  $\mathbf{L}_D$  are alternately used to update the gesture generator and discriminator. The trained

TABLE II  
DISCRIMINATOR ARCHITECTURE

Layer	Layer Type	Hyperparameter
1	Linear	$C_{in} = 768, C_{out} = 32$
1	Linear	$C_{in} = 26, C_{out} = 32$
1	Linear	$C_{in} = 45, C_{out} = 32$
2	Linear	$C_{in} = 32, C_{out} = 64$
3	1D-Conv	$k = 3, s = 1, p = 0, C_{in} = 192, C_{out} = 192$
4	1D-Conv	$k = 4, s = 2, p = 0, C_{in} = 192, C_{out} = 256$
5	1D-Conv	$k = 3, s = 1, p = 0, C_{in} = 256, C_{out} = 256$
6	1D-Conv	$k = 4, s = 2, p = 0, C_{in} = 256, C_{out} = 512$
7	1D-Conv	$k = 3, s = 1, p = 0, C_{in} = 512, C_{out} = 512$
8	1D-Conv	$k = 4, s = 2, p = 0, C_{in} = 512, C_{out} = 1024$
9	Linear	$C_{in} = 1024, C_{out} = 512$
10	Linear	$C_{in} = 512, C_{out} = 256$
11	Linear	$C_{in} = 256, C_{out} = 1$

result of  $D()$  is 1 for original gestures and 0 for generated (fake) gestures.

We split the Trinity dataset into three parts: 84% for the training set (205 minutes), 7.4% for the validation set (18 minutes), and 8.6% for the test set (21 minutes), and every set has its own audio, text transcript, and co-speech motion files. We trained the proposed model for 100 epochs. The batch size was 64, while the learning rate was 0.0001. The optimizer for both gesture generator and the discriminator is Adam with  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . The weights for loss functions ( $\alpha = 1, \beta = 0.6, \lambda = 0.3$ ) were set experimentally. The model was trained for approximately 7 hours on a GPU (NVIDIA RTX 3070) with CPU (Intel 12900k). For a 30-second speech input, the overall compilation time from loading the speech input to feature extraction to final motion file generation takes about 12.3 seconds in total, either by loading the pre-trained model on CPU or loading it on GPU.

#### IV. EVALUATION METRICS

##### A. Quantitative Evaluation

Objectively evaluating generated gestures is difficult due to the lack of suitable measures for measuring the perceived quality of co-speech gestures. The recent review works [32], [33] indicated that there is no consensus in previous works on which quantitative evaluation could be applied to assess the quality of the generated gestures. Objective assessment measures are still necessary for fair and trustworthy comparisons of various models. We mainly utilised measurements suggested by earlier studies as a trend towards developing standard assessment measures in the area of gesture generation. Thus, we followed the step from the state-of-the-art model Gesticulator [18], used Acceleration of gestures, Jerk of gestures which is Acceleration changing rate to evaluate the average value of a sequence of gestures. Additionally, RMSE (Root-Mean-Square Error) is also included, which is a common metric of the discrepancies between results produced by a model, shown in the equation 7:

$$\text{RMSE}(\mathbf{g}_i, \hat{\mathbf{g}}_i) = \sqrt{\frac{1}{n} \sum_{i=1}^n (\mathbf{g}_i - \hat{\mathbf{g}}_i)^2} \quad (7)$$

Where  $\mathbf{g}_i$  and  $\hat{\mathbf{g}}_i$  are the coordinates of  $i$ th original gestures and  $i$ th generated gesture, respectively. The  $n$  is the total length of the sequence of a gesture.

##### B. User Study

The most important purpose in the field of gesture generation is to produce gestures that make humans feel comfortable and natural in communication. However, the quantitative metrics are difficult to evaluate these characteristics, which need human perception. For example, some gestures that score very low in objective evaluations may look natural and comfortable. Hence, we conducted a user study to measure the generated gestures against the ground truth.

Three criteria are used in our user study, including naturalness, time consistency and semantics of the gestures. As shown in Table III, we used the three questions for each criterion, which is frequently used in other works [15], [16], [18].

TABLE III  
THE CRITERION USED IN USER STUDY

Criterion	Description
Naturalness	Gestures were natural
	Gestures were smooth
	Gestures were comfortable
Time Consistency	Gestures timing was matched to speech
	Gesture speed was matched to speech
	Gesture pace was matched to speech
Semantics	Gestures were matched to speech content
	Gesture well described speech content
	Gesture helped people to understand the content

#### V. RESULTS

##### A. Objective Evaluation

In order to benchmark against the state-of-the-art model, we compared our proposed model with the Gesticulator [18], the first multimodality speech gesture generation model. As shown in Table IV, the results are averaged values of Acceleration, Jerk and RMSE over 50 samples from the original test dataset. Fig 4 shows an example. In accordance with the circumstances of the speech, various gestures are used. It depicts a metaphorical motion of spreading the arms to represent the idea of "all that kind of" and "kind of". The iconic gestures that depict "ultimate powers" is generated by the framework. For the speech "don't know", the skeleton makes the shrug to depict an iconic gesture. The framework correctly recognised characteristic words and produced a deictic gesture for "I" and "end".

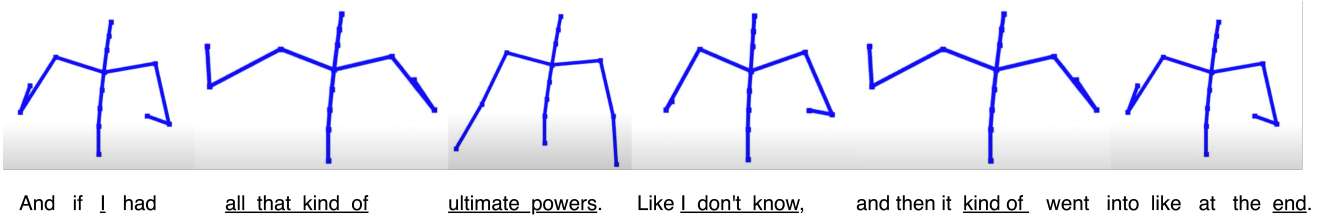


Fig. 4. Qualitative results.

TABLE IV  
OBJECTIVE EVALUATION OF PROPOSED MODEL WITH THE STATE-OF-THE-ART. FOR METRICS: CLOSER TO THE GROUND TRUTH IS BETTER. ACCELERATION(ACC.).

Model	Acc.( $cm/s^2$ )	Jerk( $cm/s^3$ )	RMSE( $cm$ )
Gesticulator	$63.8 \pm 8.3$	$1332 \pm 192$	$13.0 \pm 14.7$
Proposed Model	$94.48 \pm 19.64$	$2187.76 \pm 611.97$	$4.21 \pm 4.54$
Ground Truth	$144.7 \pm 36.6$	$2322 \pm 538$	0

### B. Subjective Evaluation

Our user study was delivered via an anonymous online questionnaire with video clips<sup>1</sup>. The questionnaire asked participants to rate the statements from strongly disagree value (1) to strongly agree value (7) after watching gesture videos. We made 10 sets of videos by using different speeches. Each set contains two 10s video clips: ground truth and generated gestures from our proposed model. The order in which the videos appear is random, and the entire questionnaire takes about 15 minutes to complete. Our user study is supported by UNSW Research Ethics Compliance Support<sup>2</sup>. From the social media, 20 native English speakers (13 male, 7 female, mean = 24.1, standard deviation = 1.8 years old) participated in our user study. Fig 5 below presented the results.

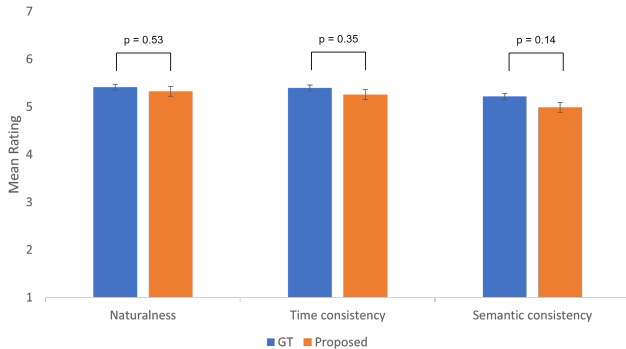


Fig. 5. Results of the user study

A two-tailed T-test was used to determine if there was a statistically significant difference in the scores of the GT

<sup>1</sup>Sample from proposed group and sample from GT group

<sup>2</sup>HC No: HC220411

and proposed groups. Although the mean rating scores of the proposed model are lower than the ground truth, especially in semantic consistency, there was no statistically significant difference among these three criteria. For the naturalness, between the ground truth group (M = 5.41, SD = 1.52) and the proposed group (M = 5.33, SD = 1.56),  $t = 0.6210$ ,  $p = 0.5349$ , and the result is not significant at  $p < 0.05$ . For the time consistency, between the ground truth group (M = 5.40, SD = 1.64) and the proposed group (M = 5.26, SD = 1.59),  $t = 0.9317$ ,  $p = 0.3520$ , and the result is not significant at  $p < 0.05$ . For the semantic consistency, between the ground truth group (M = 5.22, SD = 1.73) and the proposed group (M = 4.99, SD = 1.70),  $t = 1.48494$ ,  $p = 0.1382$ , and the result is not significant at  $p < 0.05$ .

Overall, by conventional criteria, we select a significance threshold of  $p$  value: 0.05. We observed all  $p$  values of different criteria are more than 0.05. Then, it indicated the difference between the means of the proposed model and the ground truth is not probably the result of chance. We have no basis in the data to infer that the population means of the proposed model and GT group are different because of the lack of proof of difference. Hence, the difference is considered to be not statistically significant. The results mean the performance of the proposed model is similar to the ground truth.

## VI. ABLATION STUDY

In this section, we conducted two ablation studies. One is to evaluate the difference between various input speech features, and another is to focus on the various framework structures. Both of them are evaluated by objectively.

### A. Audio Features Experiments

According to the previous data-driven method for gesture generation, they tend to use MFCCs, prosodic and Mel spectrogram as speech audio features. In order to get a better understanding of the impact of the audio feature's type, we proposed five models that used different features input. Detail settings and results are shown in Table V.

Same as the quantitative measurement, we trained 100 epochs for each type of model. From the results, we found the MFCCs-based model got the best result in RMSE and Jerk metrics, and a suboptimal result in the Acceleration metric. Although the MFCCs + Prosodic-based model achieved the best performance in the Acceleration metric when compared with the ground truth, it only showed slightly higher than the MFCCs-based model. Hence, MFCCs based model

TABLE V

OBJECTIVE EVALUATION OF AUDIO FEATURES. FOR METRICS: CLOSER TO THE GROUND TRUTH IS BETTER. ACCELERATION(ACC.).

Model	Acc.( $cm/s^2$ )	Jerk( $cm/s^3$ )	RMSE( $cm$ )
MFCCs	$94.48 \pm 19.64$	$2187.76 \pm 611.97$	$4.21 \pm 4.54$
Mel Spectrogram	$69.18 \pm 12.90$	$1234.70 \pm 234.05$	$4.50 \pm 5.12$
Prosodic	$60.54 \pm 8.90$	$948.00 \pm 138.98$	$4.42 \pm 4.93$
MFCCs + Prosodic	$94.99 \pm 22.88$	$2157.39 \pm 607.94$	$4.28 \pm 4.69$
Mel Spectrogram + Prosodic	$70.39 \pm 13.45$	$1273.65 \pm 242.76$	$4.29 \pm 4.77$
Ground Truth	$144.7 \pm 36.6$	$2322 \pm 538$	0

is the best one, which is much closer to the ground truth when comparing other models we trained.

### B. Framework Structures Experiments

In this section, based on the results from the first ablation study, we proposed five framework variants, as described in Table VI. In order to get a better understanding of the proposed framework in detail from the elimination of the key structure of the full gesture generator.

TABLE VI  
THE FIVE FRAMEWORK VARIANTS

Framework	Description
Full model	The proposed model
No Text	Only used Speech Audio as input
No Audio	Only used Speech Text as input
No GRU	Bi-directional GRU layers are not used
No FiLM Conditions	Previous gesture conditions are not used

TABLE VII

OBJECTIVE EVALUATION OF PROPOSED FRAMEWORKS. FOR METRICS: CLOSER TO THE GROUND TRUTH IS BETTER. ACCELERATION(ACC.).

Framework	Acc.( $cm/s^2$ )	Jerk( $cm/s^3$ )	RMSE( $cm$ )
Full model	$94.48 \pm 19.64$	$2187.76 \pm 611.97$	$4.21 \pm 4.54$
No Text	$105.45 \pm 22.98$	$2927.35 \pm 686.52$	$4.23 \pm 4.74$
No Audio	$63.01 \pm 10.40$	$979.33 \pm 162.21$	$4.46 \pm 4.99$
No GRU	$100.36 \pm 22.00$	$2415.85 \pm 588.61$	$4.25 \pm 4.78$
No FiLM Conditions	$173.44 \pm 37.95$	$5327.06 \pm 1239.75$	$4.70 \pm 4.32$
Ground Truth	$144.7 \pm 36.6$	$2322 \pm 538$	0

The results were presented in Table VII. Changing any structure of our proposed framework will cause lower results on the RMSE metric. We note that the results are similar for no GRU compared to the full model, and the reasons could be: 1) The full model may have been too complex for the task. Removing the GRU layer may have resulted in a simpler model that still captures the relevant information from the data. 2) The efficacy of the model may not be

significantly affected if the other layers are very good at catching the necessary patterns of the data. In this instance, the lack of the GRU layer might not affect the other layers' ability to accurately reflect the data. Nevertheless, although no-GRU obtained similar results, the full model produced the best overall performance, especially in RMSE. Removing the speech audio input caused higher Jerk than the ground truth while removing the speech text input resulted in the lowest Acceleration among all frameworks.

## VII. CONCLUSION

We propose a new framework that can generate sequences of joint angles from the speech text and speech audio utterances. Based on a conditional GAN network, the proposed neural network model learns the relationship between the co-speech gestures and both semantic and acoustic features from the speech input. In order to train our neural network model, we employ co-speech gestures with corresponding speech audio utterances dataset, which is captured from a single male native English speaker. Unlike most previous works, our model has the capability to generate continuous gestures associated with the acoustic and semantics of speech. The results from both objective and subjective evaluations demonstrate the efficacy of our gesture generation framework for robots and embodied agents.

## ACKNOWLEDGMENT

Thanks to Commonwealth's contribution for funding this work through an "Australian Government Research Training Program Scholarship".

## REFERENCES

- [1] A. Fisher and L. Griswold, *Evaluation of social interaction (ESI)*. Fort Collins, CO: Three Star Press, 2010.
- [2] J. Streeck, "Gesture as communication i: Its coordination with gaze and speech," *Communications Monographs*, vol. 60, no. 4, pp. 275–299, 1993.
- [3] W. Johal, G. Calvary, and S. Pesty, "Non-verbal signals in hri: Interference in human perception," in *International Conference on Social Robotics*. Springer, 2015, pp. 275–284.
- [4] A. S. Dick, S. Goldin-Meadow, U. Hasson, J. I. Skipper, and S. L. Small, "Co-speech gestures influence neural activity in brain regions associated with processing semantic information," *Human brain mapping*, vol. 30, no. 11, pp. 3509–3526, 2009.
- [5] A. Hamacher, N. Bianchi-Berthouze, A. G. Pipe, and K. Eder, "Believing in bert: Using expressive communication to enhance trust and counteract operational error in physical human-robot interaction," in *2016 25th IEEE international symposium on robot and human interactive communication (RO-MAN)*. IEEE, 2016, pp. 493–500.
- [6] D. McNeill, *Hand and mind: What gestures reveal about thought*. University of Chicago press, 1992.
- [7] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial networks," *arXiv preprint arXiv:1406.2661*, 2014.
- [8] P. Bremner, A. G. Pipe, M. Fraser, S. Subramanian, and C. Melhuish, "Beat gesture generation rules for human-robot interaction," in *RO-MAN 2009-The 18th IEEE International Symposium on Robot and Human Interactive Communication*. IEEE, 2009, pp. 1029–1034.
- [9] J. Kim, W. H. Kim, W. H. Lee, J.-H. Seo, M. J. Chung, and D.-S. Kwon, "Automated robot speech gesture generation system based on dialog sentence punctuation mark extraction," in *2012 IEEE/SICE International Symposium on System Integration (SII)*. IEEE, 2012, pp. 645–647.
- [10] H.-H. Kim, Y.-S. Ha, Z. Bien, and K.-H. Park, "Gesture encoding and reproduction for human-robot interaction in text-to-gesture systems," *Industrial Robot: An International Journal*, 2012.

- [11] I. Mlakar, Z. Kačič, and M. Rojc, "Tts-driven synthetic behaviour-generation model for artificial bodies," *International Journal of Advanced Robotic Systems*, vol. 10, no. 10, p. 344, 2013.
- [12] Y. Kadono, Y. Takase, and Y. I. Nakano, "Generating iconic gestures based on graphic data analysis and clustering," in *2016 11th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. IEEE, 2016, pp. 447–448.
- [13] Y. Yoon, W.-R. Ko, M. Jang, J. Lee, J. Kim, and G. Lee, "Robots learn social skills: End-to-end learning of co-speech gesture generation for humanoid robots," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 4303–4309.
- [14] C. T. Ishi, D. Machiyashiki, R. Mikata, and H. Ishiguro, "A speech-driven hand gesture generation method and evaluation in android robots," *IEEE Robotics and Automation Letters*, vol. 3, no. 4, pp. 3757–3764, 2018.
- [15] D. Hasegawa, N. Kaneko, S. Shirakawa, H. Sakuta, and K. Sumi, "Evaluation of speech-to-gesture generation using bi-directional lstm network," in *Proceedings of the 18th International Conference on Intelligent Virtual Agents*, 2018, pp. 79–86.
- [16] T. Kucherenko, D. Hasegawa, G. E. Henter, N. Kaneko, and H. Kjellström, "Analyzing input and output representations for speech-driven gesture generation," in *Proceedings of the 19th ACM International Conference on Intelligent Virtual Agents*, 2019, pp. 97–104.
- [17] Y. Ferstl, M. Neff, and R. McDonnell, "Multi-objective adversarial gesture generation," in *Motion, Interaction and Games*, 2019, pp. 1–10.
- [18] T. Kucherenko, P. Jonell, S. van Waveren, G. E. Henter, S. Alexandersson, I. Leite, and H. Kjellström, "Gesticulator: A framework for semantically-aware speech-driven gesture generation," in *Proceedings of the 2020 International Conference on Multimodal Interaction*, 2020, pp. 242–250.
- [19] Y. Yoon, B. Cha, J.-H. Lee, M. Jang, J. Lee, J. Kim, and G. Lee, "Speech gesture generation from the trimodal context of text, audio, and speaker identity," *ACM Transactions on Graphics (TOG)*, vol. 39, no. 6, pp. 1–16, 2020.
- [20] K. Takeuchi, D. Hasegawa, S. Shirakawa, N. Kaneko, H. Sakuta, and K. Sumi, "Speech-to-gesture generation: A challenge in deep learning approach with bi-directional lstm," in *Proceedings of the 5th International Conference on Human Agent Interaction*, 2017, pp. 365–369.
- [21] C.-C. Chiu and S. Marsella, "How to train your avatar: A data driven approach to gesture generation," in *International Workshop on Intelligent Virtual Agents*. Springer, 2011, pp. 127–140.
- [22] H. Liu, Z. Zhu, N. Iwamoto, Y. Peng, Z. Li, Y. Zhou, E. Bozkurt, and B. Zheng, "BEAT: A large-scale semantic and emotional multi-modal dataset for conversational gestures synthesis," in *Lecture Notes in Computer Science*. Springer Nature Switzerland, 2022, pp. 612–630. [Online]. Available: [https://doi.org/10.1007/978-3-031-20071-7\\_36](https://doi.org/10.1007/978-3-031-20071-7_36)
- [23] Y. Ferstl and R. McDonnell, "Investigating the use of recurrent motion modelling for speech gesture generation," in *Proceedings of the 18th International Conference on Intelligent Virtual Agents*, 2018, pp. 93–98.
- [24] B. Wu, C. Liu, C. T. Ishi, and H. Ishiguro, "Modeling the conditional distribution of co-speech upper body gesture jointly using conditional-gan and unrolled-gan," *Electronics*, vol. 10, no. 3, p. 228, 2021.
- [25] F. S. Grassia, "Practical parameterization of rotations using the exponential map," *Journal of graphics tools*, vol. 3, no. 3, pp. 29–48, 1998.
- [26] C. Ahuja, D. W. Lee, Y. I. Nakano, and L.-P. Morency, "Style transfer for co-speech gesture animation: A multi-speaker conditional-mixture approach," in *European Conference on Computer Vision*. Springer, 2020, pp. 248–265.
- [27] S. Alexanderson, G. E. Henter, T. Kucherenko, and J. Beskow, "Style-controllable speech-driven gesture synthesis using normalising flows," in *Computer Graphics Forum*, vol. 39, no. 2. Wiley Online Library, 2020, pp. 487–496.
- [28] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.
- [29] E. Perez, F. Strub, H. De Vries, V. Dumoulin, and A. Courville, "Film: Visual reasoning with a general conditioning layer," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.
- [30] K. Vougioukas, S. Petridis, and M. Pantic, "End-to-end speech-driven facial animation with temporal gans," *arXiv preprint arXiv:1805.09313*, 2018.
- [31] M. Arjovsky, S. Chintala, and L. Bottou, "Wasserstein generative adversarial networks," in *International conference on machine learning*. PMLR, 2017, pp. 214–223.
- [32] Y. Liu, G. Mohammadi, Y. Song, and W. Johal, "Speech-based gesture generation for robots and embodied agents: A scoping review," in *Proceedings of the 9th International Conference on Human-Agent Interaction*, 2021, pp. 31–38.
- [33] P. Wolfert, N. Robinson, and T. Belpaeme, "A review of evaluation practices of gesture generation in embodied conversational agents," *IEEE Transactions on Human-Machine Systems*, 2022.