

The Anonymous entry to the GENE Challenge 2023-A Diffusion Model for Co-speech Gesture generation

Anonymous Author(s)

ABSTRACT

This paper describes the Anonymous entry to the GENE Challenge 2023. We formulate the gesture generation problem as a co-speech gesture generation problem and a semantic gesture generation problem. We solve the co-speech gesture generation problem by denoising the diffusion probabilistic model with text, audio, and pre-pose conditions. We use the U-Net with cross-attention architecture as a denoising model, and we propose a gesture autoencoder as a mapping function. The collective evaluation released by GENE Challenge 2023 shows that our model successfully generates co-speech gestures. Remarkably, our system receives the highest interlocutor speech appropriateness (53.5% matched) among all conditions except natural motion. We also conduct an ablation study to measure the effects of the pre-pose. By the results, our system contributes to the co-speech gesture generation for natural interaction.

CCS CONCEPTS

• Computing methodologies → Animation; • Human-centered computing → Human computer interaction (HCI).

KEYWORDS

co-speech gesture generation, diffusion, neural networks, generative models

ACM Reference Format:

Anonymous Author(s). 2023. The Anonymous entry to the GENE Challenge 2023-A Diffusion Model for Co-speech Gesture generation. In *Proceedings of 25th ACM International Conference on Multimodal Interaction (ICMI'23)*. ACM, New York, NY, USA, 8 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

Synthesizing synchronized and human-like gestures performs crucial roles to improve immersion, engagement, and naturalness for embodied virtual agents and humanoid robots. During the human-computer interaction(HCI) process, human uses both verbal and non-verbal expressions to provide their intent to the interlocutor. Gesture generation, which is one of the main challenges for non-verbal interaction, aims to synthesize natural-looking and meaningful human gestures. The task can be separated whether verbal expression exists or not. When verbal expressions, such as audio

or text, are given, the gesture generation model focuses on making gestures that emphasize the meaning of verbal expressions. In the other case, the model should generate gestures that deliver the intent whether verbal expressions are given or not. We define the task with verbal information as co-speech gesture generation and the task that focuses on synthesizing meaningful body motions that deliver intent as semantic gesture generation. In this research, we focus on generating high-fidelity co-speech gestures.

There are many challenges for the co-speech gesture generation. The first is timing synchronization. Since the speech and gestures are shown to the interlocutor sequentially, he or she will be confused if gestures depart from speech. For example, if the start and end timing of the gestures slightly differs from speech, the users will think that it is an implemental error. A more detrimental situation is traffic jams during continuous generation. Once the timing is out of sync, the timing between speech and gestures is continually departed and the discomfort will be gradually increased. With similar thinking, semantic synchronization, which is the second challenge, is also important to deliver proper intent. For example, when people say "I disagree." by nodding, the interlocutor will be confused that it is positive or negative.

The third obstacle is noise robustness. 3D pose estimation or motion capture is utilized to acquire gesture data. However, the quality of raw data obtained by 3D pose estimation is not enough because the algorithm is basically image-to-3D reconstruction, which is a one-to-many problem. The motion capture is better, but it is too expensive and time-consuming. To secure quality, the cost is increased exponentially. Therefore, the raw data may contain noise. Since training with noisy data hurts both quantitative and qualitative performance, a workaround such as pre-processing or noise-robust training is needed.

To tackle these problems, deep learning-based approaches have been applied to generating co-speech gestures, recently. There are three types of training strategies: reconstruction-based method[13, 16, 32], generative adversarial network(GAN)[6] based method[23, 31], and diffusion[5, 10] based method[36]. The reconstruction-based co-speech gesture generation methods directly estimate gestures from text or audio. Although the methods induce reasonable results in terms of joint error, disadvantages are seen in terms of diversity. To generate various results without quantitative performance degradation, GAN-based co-speech gesture generation models are trained by controlling the weight between reconstruction loss and adversarial loss. Recently, denoising diffusion probabilistic models(DDPMs) are achieving huge success in the generative model and computer vision fields and expanding to other research fields[12, 22]. Especially, the diffusion model could synthesize various images that reflect input conditions, even if its semantic space is large. Since the semantic space of the speech for co-speech gesture generation is large, the diffusion model may help to synthesize various and synchronized results. Therefore, the goal of the paper

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ICMI'23, October 09–13, 2023, Paris, France

© 2023 Association for Computing Machinery.
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00
<https://doi.org/XXXXXXX.XXXXXXX>

initial poses \mathbf{x}_{-1} of the M frames. The learning objective of the problem can be formulated as $\arg \min_{\theta} \|\mathbf{x} - G_{\theta}(\mathbf{a}, \mathbf{s}, \mathbf{x}_{-1})\|$.

However, samples in the training data often have a long duration. To reduce the computational cost and memory usage, every modality of the sample is cropped into segments $\mathbf{x} = \{\mathbf{x}_1, \dots, \mathbf{x}_i\}$, $\mathbf{a} = \{\mathbf{a}_1, \dots, \mathbf{a}_i\}$, and $\mathbf{s} = \{\mathbf{s}_1, \dots, \mathbf{s}_i\}$, where \mathbf{x}_i has N frames and $\mathbf{a}_i, \mathbf{s}_i$ have the same time length as \mathbf{x}_i . Now the generative model G estimates \mathbf{x}_i from the audio \mathbf{a}_i , text \mathbf{s}_i , and the M pose frames from previous segment $\mathbf{x}_{i-1}^{(N-M):N}$, instead of synthesizing \mathbf{x} at once. Finally, the generative model G synthesizes the gestures $\{\mathbf{x}_1, \dots, \mathbf{x}_i\}$ continuously.

The model is autoregressive because the poses generated by the previous segment are used to synthesize the current segment, and stochastic because the initial diffusion feature map is random noise.

3.2 Gesture Autoencoder

In the Stable Diffusion[26], the latent diffusion model provides flexible, computationally tractable, and sometimes achieving quality improvement. The gesture autoencoder focus on finding good latent embedding space projected from gesture space. The gesture autoencoder consists of two autoencoder models: pose autoencoder and motion autoencoder. Since the gesture is the sequential pose data, we design the pose autoencoder for projecting the raw pose space to latent space, and the motion autoencoder to find correlations along the time axis.

The pose encoder and decoder consist of 3 fully-connected layers with dropout[27] and GELU activation function[9] each. The input poses sequence $\mathbf{x}^{N \times 3J}$ is projected to $\mathbf{z}'^{N \times D}$ by the pose encoder, where \mathbf{z}' denotes mid-level hidden representation, J is the number of joints, and D is the dimension of \mathbf{z}' , and the pose decoder performs reverse projection. The pose autoencoder is first trained with L1 reconstruction loss. Once the pose autoencoder is optimized, the parameters are frozen in the rest training stages.

The motion autoencoder aims to capture sequential information of the data. Thus, the motion encoder and decoder consist of 3 gated recurrent units(GRU) layers[4] and 3 multi-head self-attention layers[29], which have strong capacity in sequential data modeling. The motion encoder is formulated

$$\mathbf{z} = \text{MHSA}(\text{GRU}(\mathbf{z}')) \quad (1)$$

where $\text{MHSA}(X) = \text{Attention}(X, X, X)$. The attention mechanism is

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right) \cdot V \quad (2)$$

where Q, K , and V are the query, key, and value from the feature matrix, d is the channel dimension, and T is the matrix transpose operation.

The mid-level hidden representation $\mathbf{z}'^{N \times D}$ is projected to $\mathbf{z}^{N \times D}$ by the motion encoder, where \mathbf{z} denotes hidden representation in feature space, and the motion decoder performs reverse projection. The motion autoencoder is trained with L1 reconstruction loss. The parameters of the motion autoencoder are also frozen after this training stage.

3.3 Conditioning

The diffusion models are theoretically capable of modeling the conditional distribution $p(\mathbf{z}|y)$. This can be implemented with a conditional denoising autoencoder $\epsilon_{\theta}(\mathbf{z}_t, t, y)$, where $y \in \{\mathbf{a}, \mathbf{s}, \mathbf{z}_{i-1}\}$, to address the generation process through inputs y . To combine conditional information and latent vector in the U-Net backbone, we use a cross-attention mechanism, which is used in Stable Diffusion[26].

The three modalities, which are audio, text, and pre-pose, are used as conditions in the diffusion process. The pre-processed audio features, text features, and pre-pose features are projected to the embedding vectors by fully-connected layers. These three embedding vectors are added to the time embedding vector and propagate the information of each modality to the denoising U-Net model.

3.4 Diffusion

DDPMs define the latent variable models of the form $p_{\theta}(x_0) = \int p_{\theta}(x_{0:T}) dx_{1:T}$, where $x_{1:T}$ are latent variables in the same sample space as x_0 with the same dimensionality.

The forward process, which is also called the diffusion process, approximates the posterior distribution $q(x_{1:T}|x_0)$ by the Markov chain that gradually adds Gaussian noise to the data according to the variance schedule β_1, \dots, β_T :

$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1}), \quad (3)$$

where

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t\mathbf{I}). \quad (4)$$

The forward process variances β_t can be learned by reparameterization or held constant as hyperparameters. Since our model uses gesture autoencoder for mapping from pose to latent embeddings, the latent embeddings are gradually corrupted by noise, which finally leads to a pure white noise when T goes to infinity. Therefore, the prior latent distribution of $p(x_T)$ is $\mathcal{N}(x_T; \mathbf{0}, \mathbf{I})$ with only information of Gaussian noise.

The reverse process estimates the joint distribution of $p_{\theta}(x_{0:T})$. It is defined as a Markov chain with learned Gaussian transitions starting at $\mathcal{N}(x_T; \mathbf{0}, \mathbf{I})$:

$$p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t), \quad (5)$$

where

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_{\theta}(x_t, t)). \quad (6)$$

The corrupted noisy latent embedding x_t is sampled by $q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I})$, where $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$.

Since the problem is co-speech gesture generation, which is a conditional generation problem, we have to provide additional inputs \mathbf{a}, \mathbf{s} , and \mathbf{z}_{i-1} to the model. Therefore, these conditions are injected into the generation process. The reverse process of each timestep can be updated for our problem as:

$$p_{\theta}(z_{t-1}|z_t, y) = \mathcal{N}(z_{t-1}; \mu_{\theta}(z_t, t, y), \beta_t\mathbf{I}). \quad (7)$$

The reverse process is started by sampling a Gaussian noise $z_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and following the Markov chain to iteratively denoise the latent variable x_t via Eq. 7 to get the original latent vector z_0 .

Table 1: Detailed hyperparameters setting

Hyperparameter	Value
# of joints (J)	26
# of pre-pose frames (M)	8
# of frames of the segment (N)	128
Denoising diffusion steps	1000
Feature dimension (D)	128
Condition vector dimension	512
# of residual blocks per up/downsampling layer	2
# of up/downsampling layers	4
# of attention heads	4
N-FFT	4096
Hop length [ms]	33
Text embedding dimension	1024

To optimize the diffusion model, the variational lower bound on negative log-likelihood. We follow [10] to simplify the training objective to the ensemble of MSE losses as:

$$L(\theta) = \mathbb{E}_{t, x_0, \epsilon} [|\epsilon - \epsilon_\theta(\sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t}\epsilon, y, t)|^2], \quad (8)$$

where t is uniformly sampled between 1 and T , and ϵ is initialized as $\mathcal{N}(\mathbf{0}, \mathbf{I})$. The diffusion model is trained by the gradient descent steps on Eq. 8 until converged.

4 EXPERIMENT

4.1 Data Processing

We train our model using the GENE Challenge 2023 dataset[33], which is based on the Talking with Hands 16.2M dataset[19]. The dataset includes a training set of 371 clips, a validation set of 40 clips, and a test set of 70 clips. Each clip contains audio, transcript, and gesture motion for the main agent, gesture motion for the interlocutor, and metadata. The sampling rate of the audio is 44100kHz. The gesture motion is in BVH format and its FPS is 30.

Our system only used main-agent audio and text. The metadata and interlocutor information were ignored. We extracted the mel-spectrogram, mel-frequency cepstrum coefficients, and prosody with $n\text{-fft}=4096$ and $\text{hop length}=33\text{ms}$. We used Librosa[25] package and Parselmouth[11] library to extract audio features. The output from the network was joint angles relative to a T-pose. These joint angles were parameterized using the exponential map[7], with each dimension normalized to have zero mean and unit standard deviation over the official challenge training set. We selected all joints in the full-body expression and the number of selected joints is 26. We smooth the generated gestures using the Savitzky-Golay filter[24] with a window length of 9 and polynomial order of 3. The text segment is embedded by pre-trained text embedding model[30], which has 1024 dimensions for each sentence. We used sentence embedding because we assume that sentence embedding can address semantic information compared to word embedding. Since the audio, text, and gestures are cropped from the same timestamps, The timing of the audio features, text embeddings, and pose sequences are synchronized.

Table 2: Summary of the collective perception study with a 0.05 confidence interval about human-likeness. Our entry is SA.

Condition	Human-likeness	
	Median	Mean
NA	71 ∈ [70, 71]	68.4±1.0
SG	69 ∈ [67, 70]	65.6±1.4
SF	65 ∈ [64, 67]	63.6±1.3
SJ	51 ∈ [50, 53]	51.8±1.3
SL	51 ∈ [50, 51]	50.6±1.3
SE	50 ∈ [49, 51]	50.9±1.3
SH	46 ∈ [44, 49]	45.1±1.5
BD	46 ∈ [43, 47]	45.3±1.4
SD	45 ∈ [43, 47]	44.7±1.3
BM	43 ∈ [42, 45]	42.9±1.3
SI	40 ∈ [39, 43]	41.4±1.4
SK	37 ∈ [35, 40]	40.2±1.5
SA	30 ∈ [29, 31]	32.0±1.3
SB	24 ∈ [23, 27]	27.4±1.3
SC	9 ∈ [9, 9]	11.6±0.9

5 DISCUSSION

In this section, we discuss evaluation results. The submitted co-speech gestures are measured by three aspects: human likeness, appropriateness for agent speech, and appropriateness for the interlocutor. The natural motion, monadic baseline, and dyadic baseline are labeled NA, BM, and BD, respectively. Our submitted entry name is named SA. Our gesture generation system is tested on a desktop with a 3.20GHz i9-12900K CPU, 128GB RAM, and a RTX 3090 GPU.

5.1 Human-likeness

The evaluation results are shown in Table 2 and Figure 2. Our submitted system receives a median human-likeness score of 30 and a mean human-likeness score of 32.0. There is a gap in human likeness between our entry and natural motion. We think one major reason is that our system does not use any structural information about the joints. Since the model does not catch a relationship between the joints, the model generates gestures by focusing on the movements of the arms, which have a large movement compared to the head or body joints. Moreover, our system does not include finger motions. One other potential issue is smoothing methods. The motions generated by our system seem not to be smooth, despite the smoothing filter being applied. There can be some reasons, such as the smoothing filter is not fully optimized and the number of pre-pose is not enough.

5.2 Appropriateness for main agent speech

Regarding speech appropriateness for the main agent, Table 3 and Figure 4 describes that our entry obtains 54.8% matched. The speech appropriateness score for the main agent is not outstanding compared to most of the other conditions. We think one potential reason is semantic conditioning. Our system uses a pre-trained sentence

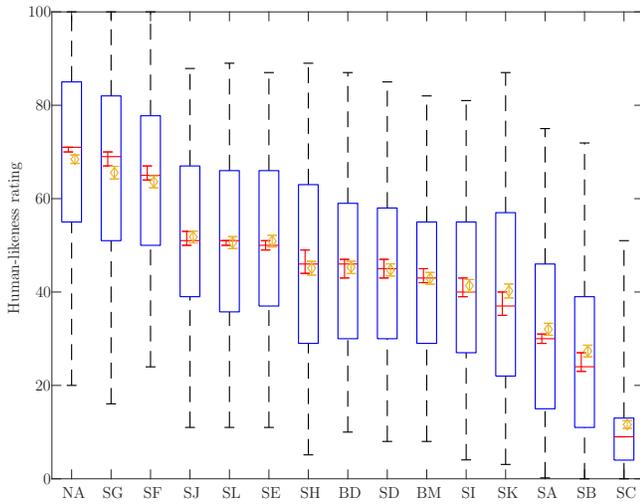


Figure 2: Box plot visualizing the rating distribution in the human-likeness study. Red bars are the median ratings (each with a 0.05 confidence interval); yellow diamonds are the mean ratings (also with a 0.05 confidence interval). Box edges are at 25 and 75 percentiles, while whiskers cover 95% of all ratings for each condition.

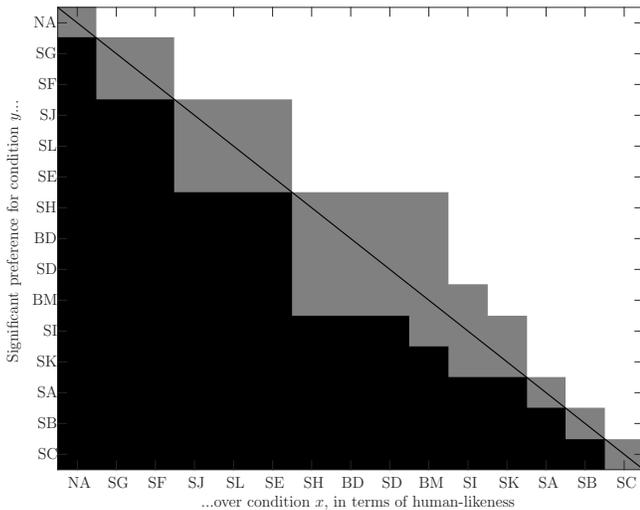


Figure 3: Significance of pairwise differences between conditions. White means that the condition listed on the y-axis rated significantly above the condition on the x-axis, black means the opposite (y rated below x), and grey means no statistically significant difference at the level $\alpha = 0.05$ after Holm-Bonferroni correction.

embedding model without fine-tuning. However, lots of text segments in the data are not satisfy sentence structure. Therefore, embeddings may not express the semantics of the text segment correctly. To solve the problem, we suggest that the model uses longer segments or a word embedding model. Another issue is the naturalness of the generated gestures. The gesture results generated by our

Table 3: Summary statistics of user-study responses from appropriateness for main agent speech, with confidence intervals for the mean appropriateness score(MAS) at the level $\alpha = 0.05$. "Pref. matched" identified how often test-takers preferred matched motion in terms of appropriateness, ignoring ties.

Condition	MAS	Pref. matched	Raw response count					Sum
			2	1	0	-1	-2	
NA	0.81±0.06	73.6%	755	452	185	217	157	1766
SG	0.39±0.07	61.8%	531	486	201	330	259	1807
SJ	0.27±0.06	58.4%	338	521	391	401	155	1806
BM	0.20±0.05	56.6%	269	559	390	451	139	1808
SF	0.20±0.06	55.8%	397	483	261	421	249	1811
SK	0.18±0.06	55.6%	370	491	283	406	252	1802
SI	0.16±0.06	55.5%	283	547	342	428	202	1802
SE	0.16±0.05	54.9%	221	525	489	453	117	1805
BD	0.14±0.06	54.8%	310	505	357	422	220	1814
SD	0.14±0.06	55.0%	252	561	350	459	175	1797
SB	0.13±0.06	55.0%	320	508	339	386	262	1815
SA	0.11±0.06	53.6%	238	495	438	444	162	1777
SH	0.09±0.07	52.9%	384	438	258	393	325	1798
SL	0.05±0.05	51.7%	200	522	432	491	170	1815
SC	-0.02±0.04	49.1%	72	284	1057	314	76	1803

system are sometimes not smooth and seem to odd. We believe that this situation occurred by the structure of the pose autoencoder, the number of pre-pose, and the degree of the smoothing filter.

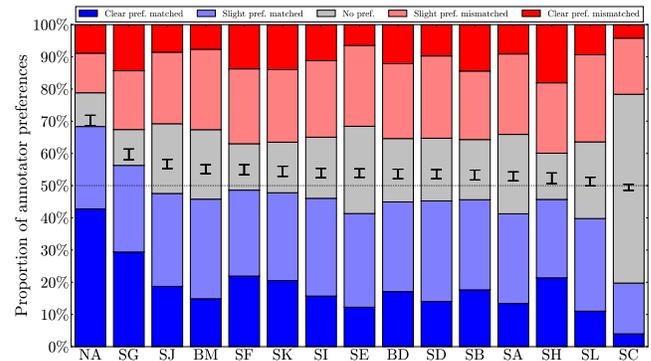


Figure 4: Bar plots visualizing the response distribution in the appropriateness for main agent speech. The blue bar(bottom) represents responses where subjects preferred the matched motion, the light grey bar(middle) represents tied responses, and the red bar(top) represents responses preferring mismatched motion, with the height of each bar being proportional to the fraction of each category. Lighter colors correspond to slight preference, and darker colors to clear preference. On top of each bar is also a confidence interval for the mean appropriateness score, scaled to fit the current axes. The dotted black line indicates chance-level performance.

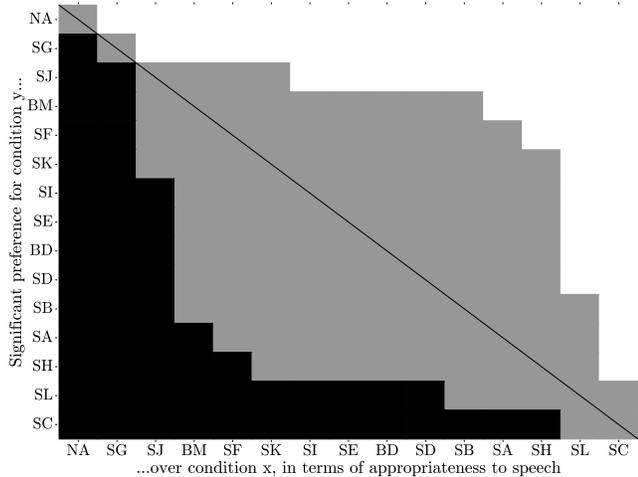


Figure 5: Significant differences between conditions in the appropriateness for main agent speech. White means the condition listed on the y -axis achieved a MAS significantly above the condition on the x -axis, black means the opposite (y scored below x), and grey means no statistically significant difference at level $\alpha = 0.05$ after correction for the false discovery rate.

5.3 Appropriateness for interlocutor speech

The evaluation results are shown in Table 4 and Figure 6 about interlocutor speech appropriateness. Our system receives a 53.5% preference matched score, which is the highest appropriateness score for interlocutor speech among all conditions, except for the natural motion entry. Interestingly, although our system did not have good scores for the main agent, it achieved a good appropriateness score for interlocutor speech. Here we summarize several settings of our condition that we think might be beneficial to improving interlocutor speech appropriateness. The first reason may be timing synchronization. In this evaluation setting, the video with the main agent and interlocutor gives two times of perceptive information, such as gestures and speech, to the referee compared to the video with the main agent only. Therefore, the referee may focus on timing rather than semantics because timing is easy to perceive and intuitive. Since our system uses speech features, including the mel-spectrogram, MFCC, and prosody, to find timing from audio, the model learns to properly align audio with gestures. Thus, our predicted gestures can control the timing of when the speaking starts and pauses. The second reason is the characteristics of generated gestures. In particular, there were many movements of extending the arms to the interlocutor, and these movements seem like a conversation with the interlocutor.

We suggest potential improvement methods for our system in the main agent and interlocutor speech appropriateness. First, the model can use the interlocutor gestures, audio, and text as conditions. Second, putting the longer previous features of the main agent and interlocutor into the conditions may generate better gestures. Third, carefully designing the text embedding model and gesture autoencoder improves semantic conditioning and the naturalness

Table 4: Summary statistics of user-study responses from appropriateness for interlocutor speech, with confidence intervals for the mean appropriateness score(MAS) at the level $\alpha = 0.05$. "Pref. matched" identified how often test-takers preferred matched motion in terms of appropriateness, ignoring ties.

Condition	MAS	Pref. matched	Raw response count					Sum
			2	1	0	-1	-2	
NA	0.63±0.08	67.9%	367	272	98	189	88	1014
SA	0.09±0.06	53.5%	77	243	444	194	55	1013
BD	0.07±0.06	53.0%	74	274	374	229	59	1010
SB	0.07±0.08	51.8%	156	262	206	263	119	1006
SL	0.07±0.06	53.4%	52	267	439	204	47	1009
SE	0.05±0.07	51.8%	89	305	263	284	73	1014
SF	0.04±0.06	50.9%	94	208	419	208	76	1005
SI	0.04±0.08	50.9%	147	269	193	269	129	1007
SD	0.02±0.07	52.2%	85	307	278	241	106	1017
BM	-0.01±0.06	49.9%	55	212	470	206	63	1006
SJ	-0.03±0.05	49.1%	31	157	617	168	39	1012
SC	-0.03±0.05	49.1%	34	183	541	190	45	993
SK	-0.06±0.09	47.4%	200	227	111	276	205	1019
SG	-0.09±0.08	46.7%	140	252	163	293	167	1015
SH	-0.21±0.07	44.0%	55	237	308	270	144	1014

of the gestures, respectively. We will focus on these aspects in our future work.

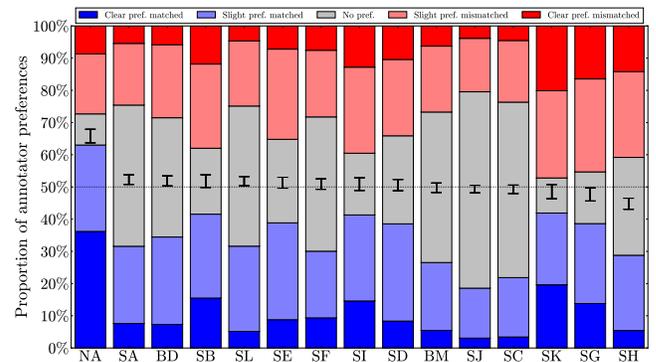


Figure 6: Bar plots visualizing the response distribution in the appropriateness for interlocutor speech. The blue bar(bottom) represents responses where subjects preferred the matched motion, the light grey bar(middle) represents tied responses, and the red bar(top) represents responses preferring mismatched motion, with the height of each bar being proportional to the fraction of each category. Lighter colors correspond to slight preference, and darker colors to clear preference. On top of each bar is also a confidence interval for the mean appropriateness score, scaled to fit the current axes. The dotted black line indicates chance-level performance.

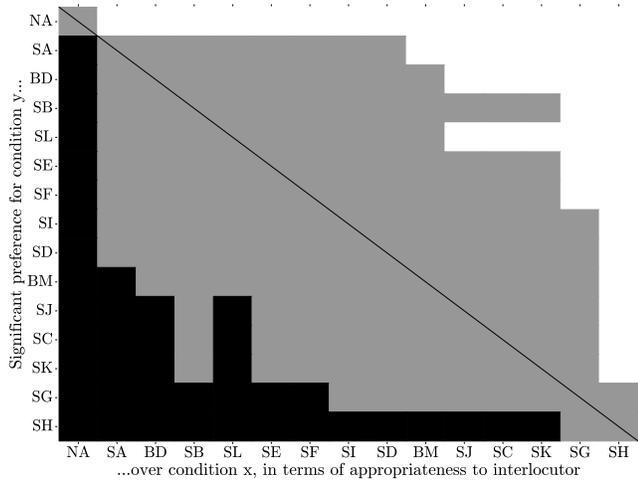


Figure 7: Significant differences between conditions in the appropriateness for interlocutor speech. White means the condition listed on the y -axis achieved a MAS significantly above the condition on the x -axis, black means the opposite (y scored below x), and grey means no statistically significant difference at level $\alpha = 0.05$ after correction for the false discovery rate.

Table 5: Effects of autoregression.

Model	FGD(feature)	FGD (raw)
w/o. pre-pose	154.984	4977.059
w. pre-pose	77.909	2279.612

5.4 ablation study

We conduct an ablation study to ensure that autoregression is helpful to co-speech gesture synthesis. We calculate Fréchet Gesture Distance (FGD), between ground truth and generated motions in the validation set, which are shown in Table 5. As a result, the FGD of discriminator features and raw gestures are improved when using the pre-pose condition.

6 CONCLUSION

In this paper, we presented a novel diffusion-based co-speech gesture generation framework, which is submitted to the GENE Challenge 2023. To generate high-fidelity co-speech gestures, we proposed a gesture autoencoder for domain transfer between gesture space and latent feature space. We also migrated the denoising diffusion probabilistic models to solve the co-speech gesture generation problem. The collective evaluation results indicated that our method is not outstanding in main agent aspects compared to most of the other entries, but our system outperforms all of the other entries in appropriateness for interlocutor speech. We further conducted an ablation study to ensure that autoregression is useful to co-speech gesture synthesis. We conclude that our system has merits in timing synchronization and generating suitable gestures

for the interaction. We also suggest some further challenges, including semantic embedding and gesture embedding model structures, for future work. We hope our method contributes to the research about diffusion-based gesture generation and be applied to various gesture generation applications.

REFERENCES

- [1] Simon Alexanderson, Gustav Eje Henter, Taras Kucherenko, and Jonas Beskow. 2020. Style-controllable speech-driven gesture synthesis using normalising flows. In *Computer Graphics Forum*, Vol. 39. Wiley Online Library, 487–496.
- [2] Simon Alexanderson, Rajmund Nagy, Jonas Beskow, and Gustav Eje Henter. 2022. Listen, denoise, action! audio-driven motion synthesis with diffusion models. *arXiv preprint arXiv:2211.09707* (2022).
- [3] Che-Jui Chang, Sen Zhang, and Mubbasir Kapadia. 2022. The IVI Lab entry to the GENE Challenge 2022–A Tacotron2 based method for co-speech gesture generation with locality-constraint attention mechanism. In *Proceedings of the 2022 International Conference on Multimodal Interaction*. 784–789.
- [4] Kyunghyun Cho, Bart van Merriënboer, Çağlar Gulçehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 1724–1734.
- [5] Prafulla Dhariwal and Alexander Nichol. 2021. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems* 34 (2021), 8780–8794.
- [6] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2020. Generative adversarial networks. *Commun. ACM* 63, 11 (2020), 139–144.
- [7] F Sebastian Grassia. 1998. Practical parameterization of rotations using the exponential map. *Journal of graphics tools* 3, 3 (1998), 29–48.
- [8] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al. 2020. Conformer: Convolution-augmented transformer for speech recognition. *arXiv preprint arXiv:2005.08100* (2020).
- [9] Dan Hendrycks and Kevin Gimpel. 2016. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415* (2016).
- [10] Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in neural information processing systems* 33 (2020), 6840–6851.
- [11] Yannick Jadoul, Bill Thompson, and Bart De Boer. 2018. Introducing parselmouth: A python interface to praat. *Journal of Phonetics* 71 (2018), 1–15.
- [12] Yifan Jiang, Han Chen, and Hanseok Ko. 2023. Spatial-temporal Transformer-guided Diffusion based Data Augmentation for Efficient Skeleton-based Action Recognition. *arXiv preprint arXiv:2302.13434* (2023).
- [13] Gwantae Kim, Seonghyeok Noh, Insung Ham, and Hanseok Ko. 2023. MPE4G: Multimodal Pretrained Encoder for Co-Speech Gesture Generation. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 1–5.
- [14] Gwantae Kim, Youngsuk Ryu, Junyeop Lee, David K Han, Jeongmin Bae, and Hanseok Ko. 2022. 3d human motion generation from the text via gesture action classification and the autoregressive model. In *2022 IEEE International Conference on Image Processing (ICIP)*. IEEE, 1036–1040.
- [15] Jihoon Kim, Jiseob Kim, and Sungjoon Choi. 2023. Flame: Free-form language-based motion synthesis & editing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 37. 8255–8263.
- [16] Taras Kucherenko, Dai Hasegawa, Gustav Eje Henter, Naoshi Kaneko, and Hedvig Kjellström. 2019. Analyzing input and output representations for speech-driven gesture generation. In *Proceedings of the 19th ACM International Conference on Intelligent Virtual Agents*. 97–104.
- [17] Taras Kucherenko, Patrik Jonell, Sanne Van Waveren, Gustav Eje Henter, Simon Alexandersson, Iolanda Leite, and Hedvig Kjellström. 2020. Gesticulator: A framework for semantically-aware speech-driven gesture generation. In *Proceedings of the 2020 international conference on multimodal interaction*. 242–250.
- [18] Taras Kucherenko, Rajmund Nagy, Youngwoo Yoon, Jieyeon Woo, Teodor Nikolov, Mihail Tsakov, and Gustav Eje Henter. 2023. The GENE Challenge 2023: A large-scale evaluation of gesture generation models in monadic and dyadic settings. In *Proceedings of the ACM International Conference on Multimodal Interaction (ICMI '23)*. ACM.
- [19] Gilwoo Lee, Zhiwei Deng, Shugao Ma, Takaaki Shiratori, Siddhartha S Srinivasa, and Yaser Sheikh. 2019. Talking with hands 16.2 m: A large-scale dataset of synchronized body-finger motion and audio for conversational motion analysis and synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 763–772.
- [20] Jing Li, Di Kang, Wenjie Pei, Xuefei Zhe, Ying Zhang, Zhenyu He, and Linchao Bao. 2021. Audio2gestures: Generating diverse gestures from speech audio with conditional variational autoencoders. In *Proceedings of the IEEE/CVF International*

- 813 *Conference on Computer Vision*. 11293–11302.
- 814 [21] Yuanzhi Liang, Qianyu Feng, Linchao Zhu, Li Hu, Pan Pan, and Yi Yang. 2022. Seeg: Semantic energized co-speech gesture generation. In *Proceedings of the*
815 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10473–10482.
- 816 [22] Songxiang Liu, Yuewen Cao, Dan Su, and Helen Meng. 2021. Diffsvc: A diffusion probabilistic model for singing voice conversion. In *2021 IEEE Automatic Speech*
817 *Recognition and Understanding Workshop (ASRU)*. IEEE, 741–748.
- 818 [23] Xian Liu, Qianyi Wu, Hang Zhou, Yinghao Xu, Rui Qian, Xinyi Lin, Xiaowei Zhou,
819 Wayne Wu, Bo Dai, and Bolei Zhou. 2022. Learning hierarchical cross-modal
820 association for co-speech gesture generation. In *Proceedings of the IEEE/CVF*
821 *Conference on Computer Vision and Pattern Recognition*. 10462–10472.
- 822 [24] Jianwen Luo, Kui Ying, and Jing Bai. 2005. Savitzky–Golay smoothing and
823 differentiation filter for even number data. *Signal processing* 85, 7 (2005), 1429–
824 1434.
- 825 [25] Brian McFee, Colin Raffel, Dawen Liang, Daniel P Ellis, Matt McVicar, Eric
826 Battenberg, and Oriol Nieto. 2015. librosa: Audio and music signal analysis in
827 python. In *Proceedings of the 14th python in science conference*, Vol. 8. 18–25.
- 828 [26] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn
829 Ommer. 2022. High-resolution image synthesis with latent diffusion models. In
830 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*.
831 10684–10695.
- 832 [27] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan
833 Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from
834 overfitting. *The journal of machine learning research* 15, 1 (2014), 1929–1958.
- 835 [28] Guy Tevet, Sigal Raab, Brian Gordon, Yoni Shafir, Daniel Cohen-or, and
836 Amit Haim Bermano. 2023. Human Motion Diffusion Model. In *The Eleventh*
837 *International Conference on Learning Representations*. [https://openreview.net/](https://openreview.net/forum?id=SJ1kSyO2jwu)
838 forum?id=SJ1kSyO2jwu
- 839 [29] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones,
840 Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all
841 you need. *Advances in neural information processing systems* 30 (2017).
842 871
- 843 [30] Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang,
844 Rangan Majumder, and Furu Wei. 2022. Text embeddings by weakly-supervised
845 contrastive pre-training. *arXiv preprint arXiv:2212.03533* (2022).
846 875
- 847 [31] Youngwoo Yoon, Bok Cha, Joo-Haeng Lee, Minsu Jang, Jaeyeon Lee, Jaehong
848 Kim, and Geehyuk Lee. 2020. Speech gesture generation from the trimodal
849 context of text, audio, and speaker identity. *ACM Transactions on Graphics (TOG)*
850 39, 6 (2020), 1–16.
851 877
- 852 [32] Youngwoo Yoon, Woo-Ri Ko, Minsu Jang, Jaeyeon Lee, Jaehong Kim, and Gee-
853 hyuk Lee. 2019. Robots learn social skills: End-to-end learning of co-speech
854 gesture generation for humanoid robots. In *Proceedings of 2019 International*
855 *Conference on Robotics and Automation*. IEEE, 4303–4309.
856 880
- 857 [33] Youngwoo Yoon, Pieter Wolfert, Taras Kucherenko, Carla Viegas, Teodor Nikolov,
858 Mihail Tsakov, and Gustav Eje Henter. 2022. The GENE Challenge 2022: A
859 large evaluation of data-driven co-speech gesture generation. In *Proceedings of*
860 *the 2022 International Conference on Multimodal Interaction*. 736–747.
861 881
- 862 [34] Mingyuan Zhang, Zhongang Cai, Liang Pan, Fangzhou Hong, Xinying Guo,
863 Lei Yang, and Ziwei Liu. 2022. MotionDiffuse: Text-Driven Human Motion
864 Generation with Diffusion Model. *arXiv preprint arXiv:2208.15001* (2022).
865 884
- 866 [35] Chi Zhou, Tengyue Bian, and Kang Chen. 2022. Gesturemaster: Graph-based
867 speech-driven gesture generation. In *Proceedings of the 2022 International*
868 *Conference on Multimodal Interaction*. 764–770.
869 887
- 870 [36] Lingting Zhu, Xian Liu, Xuanyu Liu, Rui Qian, Ziwei Liu, and Lequan Yu. 2023. Taming
871 Diffusion Models for Audio-Driven Co-Speech Gesture Generation. In
872 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
873 *Recognition*. 10544–10553.
874 888
875 889
876 890
877 891
878 892
879 893
880 894
881 895
882 896
883 897
884 898
885 899
886 900
887 901
888 902
889 903
890 904
891 905
892 906
893 907
894 908
895 909
896 910
897 911
898 912
899 913
900 914
901 915
902 916
903 917
904 918
905 919
906 920
907 921
908 922
909 923
910 924
911 925
912 926
913 927
914 928