# FREETALKER: CONTROLLABLE SPEECH AND TEXT-DRIVEN GESTURE GENERATION BASED ON DIFFUSION MODELS FOR ENHANCED SPEAKER NATURALNESS

*Sicheng Yang[1,*], Zunnan Xu[1], Haiwei Xue[1], Yongkang Cheng[2], Shaoli Huang[3], Mingming Gong[4,5], Zhiyong Wu[1,†]*

[1] Tsinghua Shenzhen International Graduate School, Tsinghua University [2] Northwest A&F University
[3] Tencent AI Lab [4] University of Melbourne [5] Mohamed bin Zayed University of Artificial Intelligence

yangsc21@mails.tsinghua.edu.cn, shaol.huang@gmail.com, mingming.gong@unimelb.edu.au, zywu@sz.tsinghua.edu.cn,

## ABSTRACT

Current talking avatars mostly generate co-speech gestures based on audio and text of the utterance, without considering the non-speaking motion of the speaker. Furthermore, previous works on co-speech gesture generation have designed network structures based on individual gesture datasets, which results in limited data volume, compromised generalizability, and restricted speaker movements. To tackle these issues, we introduce FreeTalker, which, to the best of our knowledge, is the first framework for the generation of both spontaneous (e.g., co-speech gesture) and non-spontaneous (e.g., moving around the podium) speaker motions. Specifically, we train a diffusion-based model for speaker motion generation that employs unified representations of both speech-driven gestures and text-driven motions, utilizing heterogeneous data sourced from various motion datasets. During inference, we utilize classifier-free guidance to highly control the style in the clips. Additionally, to create smooth transitions between clips, we utilize DoubleTake, a method that leverages a generative prior and ensures seamless motion blending. Extensive experiments show that our method generates natural and controllable speaker movements. Our code, model, and demo are are available at https://youngseng.github.io/FreeTalker/.

***Index Terms***— Motion processing, gesture generation, multimodal learning, human-computer interaction

## 1. INTRODUCTION

In various applications like virtual agents, animation, and human-computer interaction, the motions of a speaker are of paramount importance [1, 2, 3, 4]. These motions can be primarily divided into two segments: co-speech gestures that are inherently tied to the spoken content and non-spontaneous motions exhibited during talks [1, 5].

In recent years, substantial focus has been dedicated to the generation of co-speech gestures. ZeroEGGS [6] emphasizes naturalness and zero-shot style control. [7] adapts DiffWave for audio-driven motion synthesis, highlighting distinctive styles and control. DiffuseStyleGesture [8] and GestureDiffuCLIP [9] generate stylized gestures with exceptional human likeness and appropriateness. However, existing works primarily focus on global style control of co-speech gestures and do not facilitate free movement of the speaker, such as walking around the stage, pointing or looking in specific directions, or interacting with the audience. These aspects are crucial in presentations and speeches. In the domain of non-spontaneous motions [10], some works such as MDM [11], M2DM [10], and MotionDiffuse [12], have focused on text-controlled motion generation, achieving improvements in realism and controllability. PriorMDM [13] introduces composition methods for denoising diffusion models.

Despite these notable advancements, a significant gap remains. To our knowledge, there hasn't been an effort that coherently integrates both of these motion categories. Challenges arise from varied motion representations, and multi-modal learning intricacies. MoFusion [14] addresses dataset harmonization through pretraining for multi-task learning. Similarly, [15] offers a framework for motion retargeting. UDE [16] introduces an engine for human motion sequences from diverse inputs. UnifiedGesture [17] employs further improvements in speech-driven gestures across multiple datasets. It's important to recognize the inherent challenges in utilizing multiple datasets.

In this paper, we propose a novel framework for generating both spontaneous and non-spontaneous speaker motions. Specifically, we first develop a diffusion-based model [18] for speaker motion generation, utilizing heterogeneous data from various motion datasets. Then, we employ classifier-free guidance [19] during inference for highly controllable style in the generated clips. Additionally, we adopt DoubleTake [13] to create smooth transitions between clips and ensure seamless motion blending. The main contributions of our work are: (1) Proposing FreeTalker, the first framework to the best of our knowledge for generating both spontaneous and non-spontaneous speaker motions trained on multiple datasets. (2) Incorporating classifier-free guidance and DoubleTake in our

---

diffusion-based model for enhanced flexibility and control in gesture generation. (3) Demonstrating improved naturalness in generated speaker motions through extensive experiments, surpassing existing approaches in terms of motion quality.
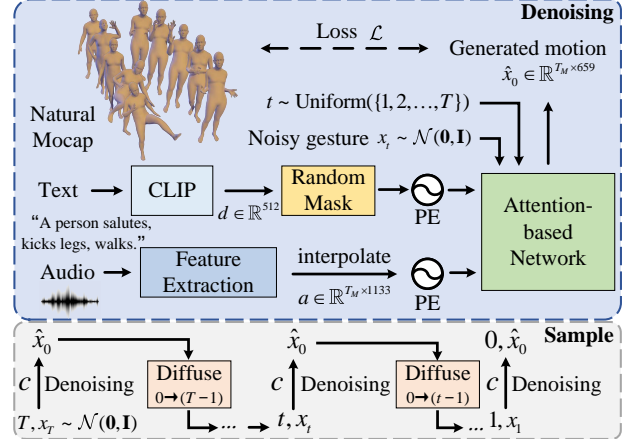
## 2. PROPOSED APPROACH

We aim to generate free-motion speakers using heterogeneous data from diverse motion datasets. In this section, we first describe the preprocessing steps required to integrate various motion data. Building on this, we introduce the diffusion model for motion generation. We then illustrate the controlled text-guided gesture generation method and explore long motion generation. Together, these components form a comprehensive system for effective generation of natural gestures.

### 2.1. Motion Processing

We expect that the features of the different motion datasets are correctly preserved. In contrast to [16] and [17], where [16] represents human motions with discrete codes and [17] retargets human motion to a homograph consisting of five terminal joints (head, hands, and feet), potentially losing important detailed information such as shoulders and fingers, our approach addresses this issue and preserves motion details. We first convert the rotation matrix of the motion capture (BVH format) data to an axis angle representation of SMPL-X [20]. For the 3D position dataset, we fit it to the SMPL-X representation using VPoser [20]. We then scale the 3D translations of the root joint appropriately and adjust the initial orientation to be uniform across the different datasets as [17]. With the SMPL-X model forward computation, we can obtain the 3D position of the SMPL-X representation. Then as in [21], we use root height, root linear and rotational velocity, joint rotation, joint position, joint velocity, and foot contact as kinematic feature representations. Each frame of the processed motion sequence has 659 dimensional features, which we denote as $\hat{x}_0 \in \mathbb{R}^{T_M \times 659}$, where $T_M$ denotes the number of motion sequence frames.

### 2.2. Diffusion Model for Motion Generation

As illustrated in Figure 1, we develop a diffusion model [18] inspired by [11] and [8]. For a noising step $t \in T$, we assume that $x_T \sim \mathcal{N}(0, I)$. The model assumes a stochastic process with $T$ noising steps: $q\left(x_t \mid x_{t-1}\right) = \mathcal{N}(\sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)I)$. The denoising process aims to predict the clean motion $\hat{x}_0$ from a noised motion $x_t$, a noise step $t$, a text condition encoded to CLIP [22] space (represented by $d$, $d \in \mathbb{R}^{512}$), and an audio condition. The audio representation, consistent with [23], includes MFCC, Mel Spectrum, Pitch, Energy, WavLM [24], and Onsets. We perform linear interpolation of audio features in the time dimension to match the number



**Fig. 1**. (Top) Denoising module. A noising step $t$ and a noisy motion sequence $x_t$ at this noising step conditioning on $c$ (including text description and audio) are fed into the model. PE indicates the addition of a positional encoding. (Bottom) Sample module. We predict the $\hat{x}_0$ with the denoising process, then add the noise to the noising step $x_{t-1}$ with the diffuse process. This process is repeated from $t = T$ until $t = 0$.

of gesture frames, denoted as $a$, where $a \in \mathbb{R}^{T_M \times 1133}$. Subsequently, the textual description is spliced together as the first frame along with the speech embedding, the noise step, and the noisy action, feeding it into the self-attention [25] layer, to yield the generated motion sequence. The denoising process is expressed as $\hat{x}_0 = Denoising\left(x_t, t, c\right)$, where $c = [d, a]$. In practice, due to the lack of datasets with both non-spontaneous speaker motion and co-speech gestures, we blend datasets with speech-driven gestures and text-driven motions, and the missing modalities are set to zero during training. The model is trained using Huber loss [26] function:

$$\mathcal{L} = E_{x_0 \sim q(x_0|c), t \sim [1,T]} \left[\|x_0 - \hat{x}_0\|_2^2\right] \quad (1)$$

During inference, at each noising step $t$, the original sample $\hat{x}_0$ is predicted and noised back to $x_{T-1}$. This process is iteratively repeated, starting from $t = T$ and continuing until $t = 0$ is reached, resulting in more natural motion generation.

### 2.3. Controllable text-guided gesture generation

Generating gestures that are both expressive and consistent with textual descriptions is a challenge. Our diffusion model addresses this problem by extending the core idea of the classifier-free approach [19, 7, 8] to adjust the strength of the non-spontaneous motion. As illustrated in Figure 1, a random mask is added to the textual embedding for classifier-free learning. The classifier-free guidance of gesture generation is achieved by combining the predictions of the text-conditioned model $Denoise\left(x_t, t, c_1\right)$, where $c_1 = [d, a]$, and the audio-conditioned model $Denoise\left(x_t, t, c_2\right)$, where $c_2 = [\varnothing, a]$, as

follows:

$$\hat{x}_{0,\gamma,c_1,c_2} = \gamma \, \text{Denoise}\,(x_t, t, c_1) + (1 - \gamma)\, \text{Denoise}\,(x_t, t, c_2) \tag{2}$$

where $\hat{x}_{0,\gamma,c_1,c_2}$ represents the combined output, and $\gamma$ is a parameter controlling the balance between the text-conditioned and audio-conditioned models. In this work, the Denoising module learns both text-conditioned and audio-conditioned distributions by randomly masking 10% of the samples using Bernoulli masks.
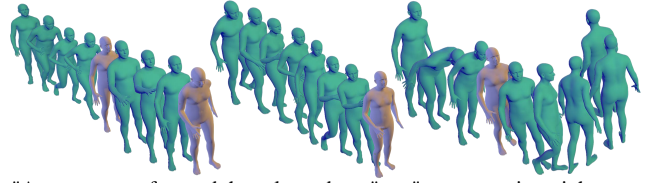
## 2.4. Long Motion Generation

In tasks involving time series, a major challenge is generating long and coherent motion sequences. Traditional approaches leveraging seed poses [8, 27] in generative tasks with non-time-aware sequences (e.g., text-to-motion) do not work well, so we use DoubleTake [13] to generate long-distance motion. Specifically, we first generate samples conditioned on its own text, audio and a handshake $\tau$ with its neighboring intervals through the denoising process, formulated as $\tau_i = (1-\vec{\alpha})\odot M_{i-1}[-h:]+\vec{\alpha}\odot M_i[:h]$, where $h$ is the length of $\tau$, $M_i$ indicates the $i^{th}$ sequence $\alpha_j = j/h, \forall j : j \in [0:h]$ and $\odot$ indicates a element-wise multiplication. Then we refine the transitions by reshaping the batch and focusing on the 'transition sandwich' $(M_i, \tau_i, M_{i+1})$. We apply a soft-masking feature, using soft mask $M_{soft}$ and hard mask $M_{hard}$ for the sequence $M$ and handshake $\tau$. The masks ensure a gradual transition between the mask values, allowing $b$-frame-long linear masking between $M_{hard}$ and $M_{soft}$. This process refines the originally generated motion (suffix or prefix) to fit the transition during the second take at every denoising step. We partially noise and denoise the sandwich $T'$ noising steps: $M'' = M' + M_{hard}\odot M_{soft}\odot (M'_{noisy} - M')$. Here, $M''$ is the refined transition of the sequence, $M'$ is the original. Finally, we construct the long motion by unfolding the refined sequences and transitions, resulting in a smooth motion.

# 3. EXPERIMENTS

## 3.1. Experiment setup

In our experiments, we use the text-driven (non-spontaneous) motion generation dataset HumanML3D [21] and speech-driven (spontaneous) gesture generation dataset BEAT [28]. All motion data are first resampled to 20 FPS. For HumanML3D dataset, we only use data with motion frame counts between 40 and 180 frames, and the maximum text length for CLIP encoding is set to 20. During training, the motion sequence length is set to $T_M = 180$ frames, with zero-padding for shorter sequences. And for BEAT, we use four English speaker gestures as described in [28] and randomly select a 180-frame segment of speech and corresponding gestures from a continuous gesture sequence. To balance



"A person runs forward then slows down" → "a person raises right hand" → Make a speech → "a person bows" → "a person turns around"

**Fig. 2**. Visualization of FreeTalker generation. We can control the speaker's non-spontaneous motion through text, while the speaker generates spontaneous co-speech gestures from speech. The light yellow color indicates the model's ability to smoothly transition between motion segments.

the number of motion data samples from both datasets, we employ weighted sampling to construct the dataloader. All motion data are normalized by subtracting the mean and dividing by the standard deviation. The data is split into training, validation, and testing sets in an 8:1:1 ratio. For the diffusion model, we use $T = 1000$ noising steps and a cosine noise schedule. The self-attention layer has a hidden space dimension of 256. The batch size is 256, the learning rate is 2e-4, and the total number of training steps is set to 1 million. The model is trained on a V100 GPU for five days. The DoubleTake method with a handshake size $h = 20$, a blend length $b = 10$, a maximum $M_{hard}$ value of 85%, a minimum $M_{soft}$ value of 15%, and $T' = 900$ denoising steps for $M'_{noisy}$.

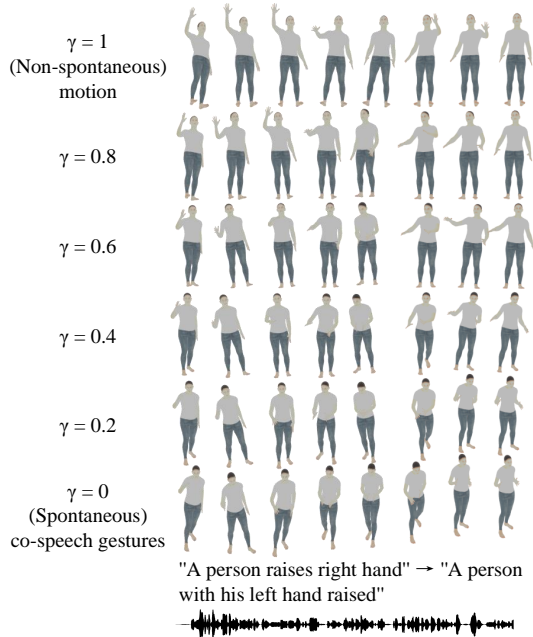## 3.2. Experimental results and analysis

### 3.2.1. Visualization

As illustrated in Figure 2, FreeTalker generates a sequence of motions, including the speaker walking on stage, waving to the audience, delivering a speech, and finally bowing before leaving the stage. The generated motions exhibit smooth transitions between segments, allowing the speaker to move and speak in a natural manner. As shown in Figure 3, when $\gamma$ in Equation (2) is set to 0, the gesture generation is conditioned only on speech input, enabling the model to produce co-speech gestures. As $\gamma$ gradually increases from 0 to 1, the model generates non-spontaneous gestures while maintaining alignment with speech. This allows us to freely edit the generated gestures and motions according to the text description.

### 3.2.2. Objective Evaluation

Due to the lack of methods capable of generating both spontaneous co-speech gestures and non-spontaneous motions, we evaluate each type of motion separately. We select [8] and [11] as our baseline models, as they have recently achieved excellent results. For co-speech gesture generation, we assess jerk, acceleration [29], and FID [30]; on the other hand, for textual description-driven motion generation, we evaluate

**Table 1**. Quantitative results of comparison with the baseline models and ablation studies. '→' denotes the closer to the real motion the better. 'Naturalness' denotes the "Ours vs. Compared model" of the user study. '\*' denotes "Ours vs. Ground Truth", implying a more rigorous evaluation, while entries without an asterisk are in reference to comparisons with other models. '\*' denotes the proposed model. '-' and '+' denote the removal and addition of the component, respectively.

| Name | Co-speech gesture generation | | | | Motion Generation | | | Free-motion | |
|---|---|---|---|---|---|---|---|---|---|
| | jerk → | acceleration → | FID ↓ | Naturalness ↑ | SSIM ↑ | FID ↓ | Naturalness ↑ | FID ↓ | Naturalness ↑ |
| Natural Mocap | $135.36 \pm 58.61$ | $12.39 \pm 11.79$ | - | - | - | - | - | - | - |
| DiffuseStyleGesture [8] | $206.52 \pm 83.65$ | $5.68 \pm 2.19$ | 0.008 | 49% | - | - | - | - | - |
| MDM [11] | - | - | - | - | 0.386 | 0.050 | 53% | - | - |
| Ours* | $245.78 \pm 108.27$ | $6.03 \pm 2.55$ | 0.139 | 40%* | 0.457 | 0.226 | 24%* | 0.139 | - |
| - Huber loss* | $226.30 \pm 73.53$ | $5.98 \pm 2.33$ | 0.027 | 52% | 0.389 | 0.041 | 53% | 0.029 | 54% |
| + local attention* | $203.77 \pm 84.45$ | $5.97 \pm 2.51$ | 0.005 | 49% | 0.431 | 0.051 | 54% | 0.032 | 52% |



γ = 1
(Non-spontaneous)
motion

γ = 0.8

γ = 0.6

γ = 0.4

γ = 0.2

γ = 0
(Spontaneous)
co-speech gestures

"A person raises right hand" → "A person with his left hand raised"

**Fig. 3**. Visualization of style editing (non-spontaneous motion control) based on co-speech gestures. From top to bottom, generated motions gradually transition from text description-based control to spontaneous co-speech gestures based on speech, resulting in highly controllable gestures.

SSIM [31] and FID. The results are shown in Table 1. Our method attains competitive results with the baseline models for both generation tasks, demonstrating the effectiveness of our approach. Moreover, our method slightly outperforms the baselines in terms of jerk, acceleration and SSIM metrics.

### 3.2.3. Subjective Evaluation

To further evaluate the quality of the generated motions, we conducted a user study focusing on the naturalness (quality of the generated motions). The study consisted of ten pairs of naturalness scoring, evaluating the naturalness of motions generated solely by co-speech gestures, solely by text-driven motions, and a combination of both. During the evaluation, participants were presented with motion sequences generated by our model and the compared models. Following [11], users were prompted with the question: "Which motion appears more human-like and reasonable?" 25 people participated in the study. The results are shown in Table 1. A score closer to 100% denotes higher naturalness. It can be observed that our model demonstrates commendable performance, often rivaling the baseline models in terms of perceived naturalness. This suggests that expanding the motion database could further improve the performance.

Our method significantly enhances the Speech2Gesture and Text2Motion subtasks, as shown in Table 1. It improves motion accuracy and naturalness, offering a diverse range of gestures, both spontaneous and non-spontaneous. This approach fills gaps in current methodologies and introduces a more adaptable motion generation framework.

### 3.2.4. Ablation study

To investigate the effectiveness of different components of our method, we designed the following ablation experiments. The results are detailed in the bottom two rows of Table 1. When the model is trained without Huber loss and instead uses MSE loss, the overall performance experiences a slight decline. Huber loss is more robust to outliers, generalizes better, and is better suited for smoothing the gradient to obtain a more coherent and natural sequence of actions. Furthermore, it converges to better results with fewer iterations. When we feed $a$ into the local attention network [32] with relative position encoding [33] to extract the local information related to the gesture before the self-attention layer as [8], the performance of co-speech gesture generation decreases slightly. However, the performance of motion generation improves. This illustrates the necessity of balancing different motion generation tasks to maintain optimal performance.

### 4. CONCLUSIONS

In this paper, we presented FreeTalker, a simple yet effective framework for generating both spontaneous and non-spontaneous speaker motions. Leveraging a diffusion-based

model, our method is trained on heterogeneous data sourced from various motion datasets. The incorporation of classifier-free guidance and DoubleTake during inference stage ensures the natural, highly controllable and long-range motion generation. Moreover, our approach lays the foundation for future work on large-scale motion datasets and more sophisticated models, paving the way for further advancements in speaker motion generation and enhancing talking avatars' naturalness in various applications.

We intend to elaborate on extending our work to the generation of fully digital humans, encompassing motions, facial expressions, and lip movements. We also aim to explore a more unified approach to digital human generation.

# Acknowledgments

## 5. REFERENCES

[1] Simbarashe Nyatsanga, Taras Kucherenko, Chaitanya Ahuja, et al., "A comprehensive review of data-driven co-speech gesture generation," in *Computer Graphics Forum*, 2023, vol. 42, pp. 569–596.

[2] Taras Kucherenko, Rajmund Nagy, Youngwoo Yoon, et al., "The genea challenge 2023: A large scale evaluation of gesture generation models in monadic and dyadic settings," *arXiv preprint arXiv:2308.12646*, 2023.

[3] Haolin Zhuang, Shun Lei, Long Xiao, et al., "Gtn-bailando: Genre consistent long-term 3d dance generation based on pretrained genre token network," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2023, pp. 1–5.

[4] Zunnan Xu, Yachao Zhang, Sicheng Yang, Ronghui Li, and Xiu Li, "Chain of generation: Multi-modal gesture synthesis via cascaded conditional control," *arXiv preprint arXiv:2312.15900*, 2023.

[5] Wentao Zhu, Xiaoxuan Ma, Dongwoo Ro, et al., "Human motion generation: A survey," *arXiv preprint arXiv:2307.10894*, 2023.

[6] Saeed Ghorbani, Ylva Ferstl, Daniel Holden, et al., "Zeroeggs: Zero-shot example-based gesture generation from speech," in *Computer Graphics Forum*, 2023.

[7] Simon Alexanderson, Rajmund Nagy, Jonas Beskow, et al., "Listen, denoise, action! audio-driven motion synthesis with diffusion models," *ACM Transactions on Graphics (TOG)*, vol. 42, no. 4, pp. 1–20, 2023.

[8] Sicheng Yang, Zhiyong Wu, Minglei Li, et al., "Diffusestylegesture: Stylized audio-driven co-speech gesture generation with diffusion models," *International Joint Conference on Artificial Intelligence*, 2023.

[9] Tenglong Ao, Zeyi Zhang, and Libin Liu, "Gesturediffuclip: Gesture diffusion model with clip latents," *ACM Trans. Graph.*, 2023.

[10] Hanyang Kong, Kehong Gong, Dongze Lian, et al., "Prioritycentric human motion generation in discrete latent space," *arXiv preprint arXiv:2308.14480*, 2023.

[11] Guy Tevet, Sigal Raab, Brian Gordon, et al., "Human motion diffusion model," in *The Eleventh International Conference on Learning Representations*, 2023.

[12] Mingyuan Zhang, Zhongang Cai, et al., "Motiondiffuse: Text-driven human motion generation with diffusion model," *arXiv preprint arXiv:2208.15001*, 2022.

[13] Yonatan Shafir, Guy Tevet, Roy Kapon, and Amit H Bermano, "Human motion diffusion as a generative prior," *arXiv preprint arXiv:2303.01418*, 2023.

[14] Jianxin Ma, Shuai Bai, and Chang Zhou, "Pretrained diffusion models for unified human motion synthesis," *arXiv preprint arXiv:2212.02837*, 2022.

[15] Kfir Aberman, Peizhuo Li, Dani Lischinski, et al., "Skeleton-aware networks for deep motion retargeting," *ACM Transactions on Graphics (TOG)*, vol. 39, no. 4, pp. 62–1, 2020.

[16] Zixiang Zhou and Baoyuan Wang, "Ude: A unified driving engine for human motion generation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 5632–5641.

[17] Sicheng Yang, Zilin Wang, et al., "Unifiedgesture: A unified gesture synthesis model for multiple skeletons," *ACM International Conference on Multimedia*, 2023.

[18] Jonathan Ho, Ajay Jain, and Pieter Abbeel, "Denoising diffusion probabilistic models," *Advances in neural information processing systems*, pp. 6840–6851, 2020.

[19] Jonathan Ho and Tim Salimans, "Classifier-free diffusion guidance," *arXiv preprint arXiv:2207.12598*, 2022.

[20] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, et al., "Expressive body capture: 3d hands, face, and body from a single image," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 10975–10985.

[21] Chuan Guo, Shihao Zou, Xinxin Zuo, et al., "Generating diverse and natural 3d human motions from text," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 5152–5161.

[22] Alec Radford, Jong Wook Kim, Chris Hallacy, et al., "Learning transferable visual models from natural language supervision," in *International conference on machine learning*. PMLR, 2021, pp. 8748–8763.

[23] Sicheng Yang, Haiwei Xue, Zhensong Zhang, et al., "The diffusestylegesture+ entry to the genea challenge 2023," *arXiv preprint arXiv:2308.13879*, 2023.

[24] Sanyuan Chen, Chengyi Wang, Zhengyang Chen, et al., "Wavlm: Large-scale self-supervised pre-training for full stack speech processing," *IEEE Journal of Selected Topics in Signal Processing*, 2022.

[25] Ashish Vaswani, Noam Shazeer, Niki Parmar, et al., "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.

[26] Peter J Huber, "Robust estimation of a location parameter," in *Breakthroughs in statistics: Methodology and distribution*, pp. 492–518. Springer, 1992.

[27] Jonathan Tseng, Rodrigo Castellon, and Karen Liu, "Edge: Editable dance generation from music," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 448–458.

[28] Haiyang Liu, Zihao Zhu, Naoya Iwamoto, et al., "Beat: A large-scale semantic and emotional multi-modal dataset for conversational gestures synthesis," in *European Conference on Computer Vision*, 2022.

[29] Taras Kucherenko, Dai Hasegawa, Gustav Eje Henter, Naoshi Kaneko, and Hedvig 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.

[30] Youngwoo Yoon, Bok Cha, Joo-Haeng Lee, et al., "Speech gesture generation from the trimodal context of text, audio, and speaker identity," *ACM Transactions on Graphics (TOG)*, 2020.

[31] Alain Hore and Djemel Ziou, "Image quality metrics: Psnr vs. ssim," in *2010 20th international conference on pattern recognition*. IEEE, 2010.

[32] Aurko Roy, Mohammad Saffar, Ashish Vaswani, et al., "Efficient content-based sparse attention with routing transformers," *Transactions of the Association for Computational Linguistics*, vol. 9, pp. 53–68, 2021.

[33] Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya, "Reformer: The efficient transformer," in *International Conference on Learning Representations*, 2020.