

# Audio is all in one: speech-driven gesture synthetics using WavLM pre-trained model

Fan Zhang<sup>ID</sup>, Naye Ji<sup>(✉) ID</sup>, Fuxing Gao<sup>ID</sup>, Siyuan Zhao<sup>ID</sup>, Zhaohan Wang<sup>ID</sup>, Shunman Li

**Abstract**—The generation of co-speech gestures for digital humans is an emerging area in the field of virtual human creation. Prior research has made progress by using acoustic and semantic information as input and adopting classify method to identify the person’s ID and emotion for driving co-speech gesture generation. However, this endeavour still faces significant challenges. These challenges go beyond the intricate interplay between co-speech gestures, speech acoustic, and semantics; they also encompass the complexities associated with personality, emotion, and other obscure but important factors. This paper introduces “diffmotion-v2,” a speech-conditional diffusion-based and non-autoregressive transformer-based generative model with WavLM pre-trained model. It can produce individual and stylized full-body co-speech gestures only using raw speech audio, eliminating the need for complex multimodal processing and manually annotated. Firstly, considering that speech audio not only contains acoustic and semantic features but also conveys personality traits, emotions, and more subtle information related to accompanying gestures, we pioneer the adaptation of WavLM, a large-scale pre-trained model, to extract low-level and high-level audio information. Secondly, we introduce an adaptive layer norm architecture in the transformer-based layer to learn the relationship between speech information and accompanying gestures. Extensive subjective evaluation experiments are conducted on the Trinity, ZEGGS, and BEAT datasets to confirm the WavLM and the model’s ability to synthesize natural co-speech gestures with various styles.

**Index Terms**—Gesture generation, Gesture synthesis, Cross-Modal, Speech-driven, Diffusion model, Transformer.

## I. INTRODUCTION

RECENTLY, The utilization of 3D virtual human technology has witnessed a notable surge in popularity, paralleling the emergence of the metaverse. This technology finds extensive applications in various domains of real-world society, encompassing animation, gaming, human-computer interaction, VTuber platforms, virtual guidance systems, digital receptionists, presenters, and various other areas.

To create realistic and engaging virtual humans, a crucial objective is the integration of non-verbal (co-speech) gestures

Fan Zhang is with the Faculty of Humanities and Arts, Macau University of Science and Technology, Macau, China; The College of Media Engineering, Communication University of Zhejiang, China; Research Center for Artificial Intelligence and Fine Arts, Zhejiang Lab, Zhejiang, China (e-mail: fanzhang@cuz.edu.cn)

Naye Ji, Fuxing Gao are with the College of Media Engineering, Communication University of Zhejiang, China (e-mail: jinaye@cuz.edu.cn; fuxing@cuz.edu.cn)

Siyuan Zhao is with the Faculty of Humanities and Arts, Macau University of Science and Technology, Macau, China; (e-mail: 2109853jai30001@student.must.edu.mo)

Zhaohan Wang is with the School of Animation and Digital Arts Communication University of China, Beijing, China (e-mail: 2022201305j6018@cuc.edu.cn)

Shunman Li is with Zhejiang Institute of Economics and Trade, Zhejiang, China



Fig. 1. Our proposed system autonomously perceives and extracts information, such as acoustic, semantic, emotional, and personality, exclusively from the raw speech audio signals. This information is then harnessed to stimulate the generation of natural and diverse co-speech gestures. Our approach performs without necessitating any annotations or intricate multimodal processing.

that appear natural and align with human communication patterns. Although motion capture systems have been developed to fulfil this requirement, their implementation necessitates specialized hardware, dedicated space, and trained actors, resulting in significant expenses. As an alternative, automatic gesture generation presents a cost-effective approach that eliminates the need for human intervention during the production phase. Among the potential solutions, speech-driven gesture generation emerges as a viable option. Nevertheless, a major challenge in this endeavour lies in effectively matching and synchronizing relevant gestures with the input speech, given the inherent complexities of cross-modal mapping, many-to-many relationships, and the diverse and ambiguous nature of gesture patterns. Furthermore, the same utterance often elicits distinct gestures at different temporal instances, even when uttered by the same or different individuals. [1].

The close relationship between gestures and the acoustic signals of speech is widely acknowledged in scholarly discourse. Consequently, considerable research has been dedicated to extracting pertinent features from speech audio signals, such as Mel-Frequency Cepstrum Coefficients(MFCCs). These extracted features serve as input for neural networks, thereby facilitating the generation of corresponding co-speech gestures. However, it is important to recognize that gestures are not exclusively tied to speech acoustic features. Rather, they exhibit intricate associations with various other factors, including well-established aspects such as personalities, emotions, and speech context, among others, as well as a multitude of unknown variables. This intricate interplay presents significant challenges in the pursuit of generating co-speech gestures that

---

possess enhanced naturally and realism.

Due to the multifaceted nature of speech signals, encompassing aspects such as speaker personality, acoustic, etc., our objective is to exclusively extract the signal originating from raw speech audio while abstaining from processing various modalities concurrently.

The recent emergence of expansive Pre-training models offers a promising opportunity to improve pre-training outcomes significantly and can potentially be transferred effectively to a wide range of subsequent tasks. In our exploration, we have identified WavLM, a noteworthy system that learned on massive amounts of unlabeled speech data to acquire universal speech representations. WavLM demonstrates remarkable adaptability across diverse speech-processing tasks, from Automatic Speech Recognition(ASR) to non-ASR, further validating its efficacy and potential for practical applications.

This paper introduces DiffMotion-v2, an innovative transformer and diffusion-based probabilistic architecture designed to generate speech-driven gestures. The model operates on extensive collections of unstructured gesture data, alleviating the need for manual annotation. The model synthesises co-speech gestures through this approach, incorporating corresponding styles such as affective expressions, gender-specific characteristics, age-related nuances, and semantic information. These stylistic features are automatically extracted from the speech audio, facilitating a seamless integration of multimodal cues within the generated gestures.

Our contributions can be summarized as follows:

- We involved pioneering utilising the WavLM Generative Pre-Trained Transformer Large model for extracting the speech audio features.
- We extended our previous autoregressive Diffmotion model to a non-autoregressive variant known as Diffmotion-v2. This extension encompasses the proposal of a novel diffusion model, which adopts a transformer-based architecture.
- We introduced adaptive layer norm(adaLN) in transformer architecture, following the widespread usage in GANs and diffusion models with U-Net backbones.
- Finally, we conducted comprehensive subjective evaluation experiments using the Trinity, ZEGGS, and BEAT datasets.

This research builds upon our prior architecture, Diffmotion [2]. However, this work extends the previous paper significantly by introducing novel features and improvements. Firstly, we introduced a non-autoregressive generative model that employs a transformer-based architecture. Unlike its predecessor, Diffmotion-v2 generates the entire sequence of full-body gestures instead of generating them frame by frame, resulting in more coherent and holistic gesture synthesis. Secondly, we enhance the feature extraction process by replacing the traditional Mel-frequency cepstral coefficients (MFCC) with the WavLM Generative Pre-Trained Transformer Large model. Lastly, unlike Diffmotion, which generated redundancy gestures and necessitated post-processing techniques to mitigate jitter-induced inconsistencies in the gesture sequence, Diffmotion-v2 overcomes this limitation and produces more stable and refined gesture sequences.

The remainder of this paper is organized as follows. The related works about co-speech gesture generation are described in Section II. Then elaborate on the DiffMotion-v2 schedule in Section III. The experimental results of three baseline detection algorithms on our dataset are presented in Section IV. We conclude this paper in Section V.

## II. RELATED WORK

The alignment of non-verbal communication, specifically co-speech gestures, with the communicative intent of virtual agents requires the establishment of a meaningful correspondence between the two modalities. The investigation of automated co-speech gesture generation, relying on speech information, can be broadly categorized into two primary domains: rule-based methods and data-driven approaches. In light of the notable success of deep learning techniques in various computer tasks, synthesising co-speech gestures has shifted from rule-based approaches(extensively reviewed by Wanger et al., [3]) towards data-driven approaches, particularly with the introduction of deep learning methodologies. Moreover, within the realm of deep learning, a distinction exists between deterministic and generative models. This discussion will briefly focus on generative models for speech-driven gesture generation.

### A. Data-driven Generative approaches

In real-life scenarios, the same utterance can be accompanied by varying gestures, even when repeated by the same speaker at different time points, highlighting the lack of coherence in gesture production. This presents a significant challenge for deterministic models, which struggle to capture the extensive variation between speech and gestures. Consequently, there has been a shift in research focus from deterministic models to probabilistic generative models. Generative adversarial networks (GANs) have shown promise in generating persuasive random samples. Accordingly, Ylva et al. [4] attempted to explore GANs [5] with multiple discriminators to convert speech into 3D gesture motion. However, this approach requires manual dataset annotation, and the results still lack realism. In contrast, Wu et al. [6] verified the effectiveness of conditional and unrolled GANs, showing that they outperformed existing deterministic models.

*Normalizing flows* [7], built on unsupervised learning algorithms such as NICE [8] and RealNVP [9], are capable of constructing complex distributions and approximating the true posterior distribution. Impressively, Alexanderson et al. [10] demonstrated the effectiveness of a network called MoGlow, based on normalizing flows, in generating a diverse set of plausible gestures given the same input speech signal, without the need for manual annotation. Li et al. [11] employed a conditional variational autoencoder (VAE) model to capture the strong correlation between audio and motion, enabling the random generation of diverse motions. Taylor et al. [12] extended normalizing flows by combining them with a variational autoencoder called Flow-VAE. Their evaluation demonstrated that this approach produces expressive body motion close to the ground truth while utilizing fewer trainable

parameters. However, it should be noted that normalizing flows require the imposition of topological constraints on the transformation [8], [9]. Furthermore, the MoGlow method employs an LSTM architecture, necessitating the generation of the entire sequence of gestures frame by frame in an autoregressive manner. This approach inevitably results in an obvious increase in the overall generation time.

Diffusion models [13]–[15] represent an alternative class of generative models that leverage a Markov chain to transform a simple distribution into a complex data distribution gradually. These models can be efficiently trained by optimizing the variational lower bound (ELBO). They have been successfully applied in image synthesis [14] and multivariate probabilistic time series forecasting [16], with connections to denoising score matching [17]. Our previous work proposed DiffMotion, a diffusion model-based framework with an LSTM architecture that generates co-speech gestures frame by frame [2]. Furthermore, Alexanderson et al. [18] adapted the DiffWave architecture, replacing dilated convolutions with Conformers [19] to enhance the modelling power and incorporating classifier-free guidance to adjust the strength of stylistic expression. Another diffusion model-based framework called GestureDiffuCLIP [20] learns a latent diffusion model to generate high-quality gestures and incorporates large-scale Contrastive-Language-Image-Pre-training (CLIP) representations for style control. However, this system requires learning a joint embedding space between corresponding gestures and transcripts using contrastive learning, which provides semantic cues for the generator and an effective semantic loss during training. In the context of style control, Ghrobani et al. [21] presented ZeroEGGS, an innovative framework that enables precise control over style gestures using a single, concise example motion clip. This approach has made notable advancements in the field by leveraging the encoder to capture the latent representation of style from a gesture sequence instead of proposing to define styles. For a comprehensive review of data-driven co-speech gesture generation, refer to Nyatsanga et al. [22].

### B. Condition Encoding Strategy

Co-speech gesture generation systems have recently incorporated conditional information as input, including audio, transcripts, style labels, and other relevant factors. This approach allows the system to consider additional contextual information during the gesture generation process. By incorporating these conditional inputs, the system can generate gestures more aligned with the speech content, style, and other specified conditions, leading to more contextually appropriate and expressive gestures.

1) *Audio Representation*: The most suitable audio representation is an open research question [23]. One of the most common audio speech representations chosen in previous work is Mel Frequency Cepstral Coefficients(MFCCs) [24] [25] [2], which is for frequency better approximates how humans perceive sounds. Another approach, such as ZeroEGGS [26], this work combines the log amplitude of the spectrogram, mel-frequency scale and the energy of the audio as speech

audio features. While in GestureDiffuCLIP [20], the speech audio  $A = [a_i]_{i=1}^L$  is parameterized as a sequence of acoustic features, each  $a_i$  encodes the onsets and amplitude envelopes that reflect the beat and volume of speech, respectively. Although these approaches have provided impressive results, these approaches only represent acoustic information; there is scope for more descriptive features.

2) *Semantic Representation*: For semantically-aware gesture generation, in several recent works mapping from text transcripts to co-speech gestures, Yoon et al. [27] learned a mapping from the utterance text to gestures using a recurrent neural network. Taras et al. [28] presented a model designed to produce the arbitrary beat and semantic gestures together, which takes both acoustic and semantic representations of speech as input. The semantic features are encoded from the text by BERT [29]. Uttaran et al. [30] obtain the word embeddings using the GloVe model pre-trained on the Common Crawl corpus [31]. This method marginally outperformed other similar-dimensional embedding models, such as Word2Vec [32] and FastText [33], and had similar performance as much higher dimensional embedding models, e.g., BEAT.

Integrating acoustic and semantic features in gesture generation has the potential to enhance the appropriateness and contextual relevance of generated gestures. However, this approach requires manual alignment and fusion of both modalities for accurate prediction. Additionally, expanding the range of signal representations to include modalities like personality and emotion introduces additional complexities in effectively processing and modelling the different modalities.

To address this challenge, we employed WavLM [34], a highly advanced pre-trained model that utilizes a vast amount of unlabeled raw speech audio data to learn universal speech representations. WavLM has demonstrated its effectiveness in various speech-processing tasks, including ASR and non-ASR, such as speaker diarization, speech separation, speech recognition, emotion recognition and etc. This study aims to investigate the potential of leveraging WavLM for speech-driven gesture generation.

## III. PROPOSED APPROACH

The task in this paper is to generate a sequence of human poses  $x_{1:T}$  given a raw speech audio waveform  $a_{1:T}$  for the same time instances.

### A. Problem Formulation

We first provide the co-speech gesture generation problem definition. We denote the gesture features and the acoustic signal as  $x^0 = x_{1:T}^0 \in [x_1^0, \dots, x_t^0, \dots, x_T^0] \in \mathbb{R}^{T \times D}$  and  $a = a_{1:T} \in [a_1, \dots, a_t, \dots, a_T] \in \mathbb{R}^T$ , where  $x_t^0 = \mathbb{R}^D$  is 3D skeleton joints angle at frame  $t$ , and  $D$  indicates the number of channels of the skeleton joints, the superscript present the diffusion time step.  $a_t$  is the current subsequence audio waveform signal at frame  $t$ , and  $T$  is the sequence length. Let  $p_\theta(\cdot)$  denote the Probability Density Function(PDF), which aims to approximate the actual gesture data distribution  $p(\cdot)$  and allows for easy sampling. We are tasked with generating the whole sequence of pose  $x \sim p_\theta(\cdot)$  non-autoregressive according to

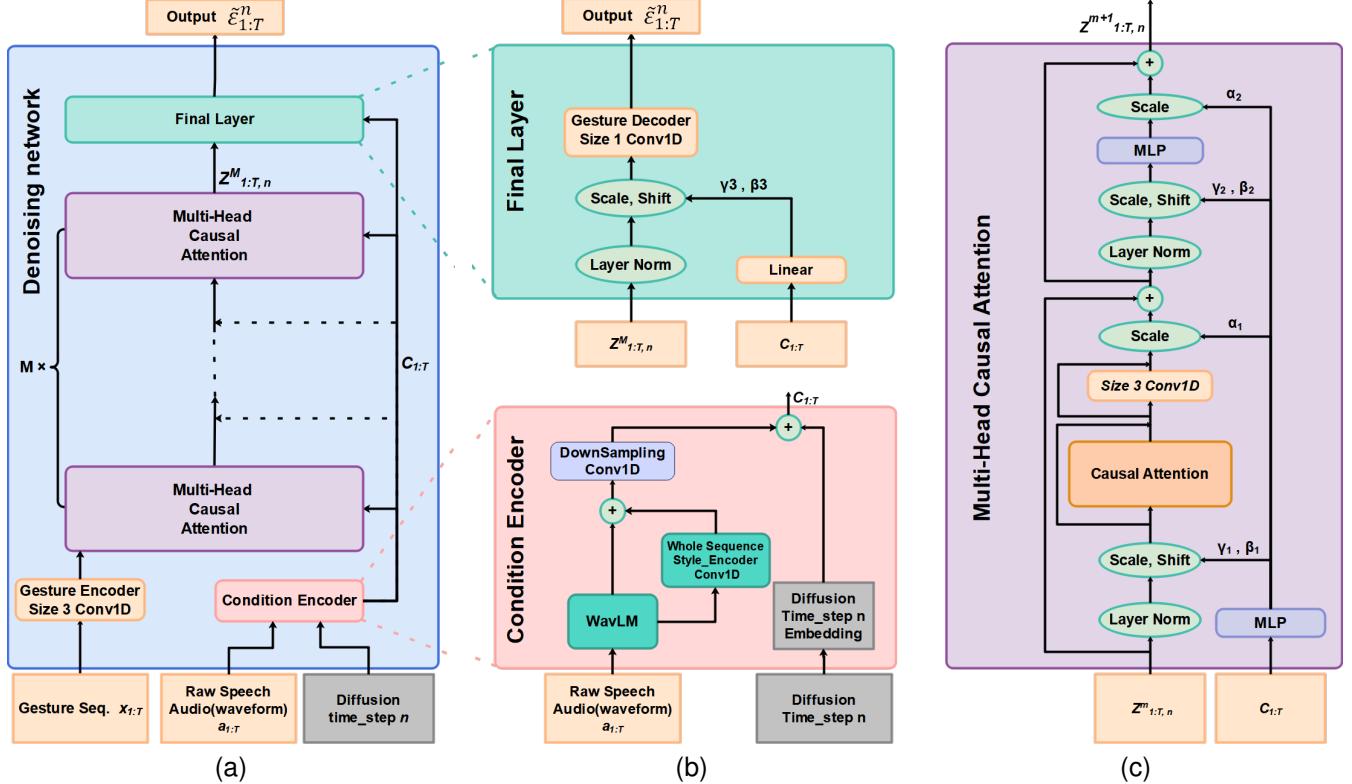


Fig. 2. Architecture of the Diffmotion-v2. The model is a multi-block causal attention structure with an adaptive layer norm for incorporating conditioning. The condition encoder takes the raw speech raw audio features extracted by WavLM as input and feeds to the multi-head causal attention blocks to learn the relation between the co-speech gestures and the audio features, estimating the diffusion noise. (a) The whole architecture. (b) The condition encoder and Final layer. (c) multi-head causal attention with adaptive layer norm.

its conditional probability distribution given acoustic signal  $a$  as covariate:

$$x^0 \sim p_\theta(x^0|a) \approx p(\cdot) := p(x^0|a) \quad (1)$$

where the  $p_\theta(\cdot)$  aims to approximate  $p(\cdot)$  trained by the denoising diffusion model. We discuss these two modules in detail in Sec.III-B

### B. Denoising Diffusion Probabilistic Model

We introduced the Denoising Diffusion Probabilistic Model (DDPM), which is a variant of diffusion models [35] of the form  $p_\theta := \int p_\theta(x^{0:N}) dx^{1:N}$ , where  $x^1, \dots, x^N$  are latent of the same dimensionality as the data  $x^n$  at the  $n$ -th diffusion time stage. The model contains two processes, namely *diffusion process* and *generation process*. At training time, the diffusion process gradually converts the original data( $x^0$ ) to white noise( $x^N$ ) by optimizing a variational bound on the data likelihood. At inference time, the generation process recovers the data by reversing this noising process through the Markov chain using Langevin sampling [36]. The whole sequence of the gestures can be generated no-autoregressive by sampling from the conditional data distribution. The Markov chains in the diffusion process and the generation process are:

$$\begin{aligned} p(x^n|x^0) &= \mathcal{N}\left(x^n; \sqrt{\bar{\alpha}^n}x^0, (1-\bar{\alpha}^n)I\right) \quad \text{and} \\ p_\theta(x^{n-1}|x^n, x^0) &= \mathcal{N}\left(x^{n-1}; \tilde{\mu}^n(x^n, x^0), \tilde{\beta}^n I\right), \end{aligned} \quad (2)$$

where  $\alpha^n := 1 - \beta^n$  and  $\bar{\alpha}^n := \prod_{i=1}^n \alpha^i$ . As shown by [14],  $\beta^n$  is a increasing variance schedule  $\beta^1, \dots, \beta^N$  with  $\beta^n \in (0, 1)$ , and  $\tilde{\beta}^n := \frac{1-\bar{\alpha}^{n-1}}{1-\bar{\alpha}^n} \beta^n$ .

The training objective is to optimize the parameters  $\theta$  that minimize the NLL via Mean Squared Error(MSE) loss between the true noise  $\epsilon \sim \mathcal{N}(0, I)$  and the predicted noise  $\epsilon_\theta$ :

$$\mathbb{E}_{x_{1:T}^0, \epsilon, n} [\|\epsilon - \epsilon_\theta (\sqrt{\bar{\alpha}^n}x^0 + \sqrt{1-\bar{\alpha}^n}\epsilon, a_{1:T}, n)\|^2], \quad (3)$$

Here  $\epsilon_\theta$  is a neural network, which uses input  $x_t^0, a_{t-1}$  and  $n$  that to predict the  $\epsilon$ , and contains the similar architecture employed in [37]. The complete training procedure is outlined in Algorithm 1.

---

#### Algorithm 1: Training for the whole sequence gestures

---

**Input:** data  $x_{1:T}^0 \sim p(x^0|a_{1:T})$  and  $a_{1:T}$   
**repeat**

    Initialize  $n \sim \text{Uniform}(1, \dots, N)$  and  $\epsilon \sim \mathcal{N}(0, I)$   
    Take gradient step on

$$\nabla_\theta \|\epsilon - \epsilon_\theta (\sqrt{\bar{\alpha}_n}x_{1:T}^0 + \sqrt{1-\bar{\alpha}_n}\epsilon, a_{1:T}, n)\|^2$$

**until** converged;

---

---

**Algorithm 2:** Sampling  $x_{1:T}^0$  via annealed Langevin dynamics

---

**Input:** noise  $x_{1:T}^N \sim \mathcal{N}(0, I)$  and raw audio waveform  $a_{1:T}$

```

for  $n = N$  to 1 do
    if  $n > 1$  then
         $z \sim \mathcal{N}(0, I)$ 
    else
         $z = 0$ 
    end if
     $x_{1:T}^{n-1} =$ 
         $\frac{1}{\sqrt{\alpha^n}} \left( x_{1:T}^n - \frac{\beta^n}{\sqrt{1-\alpha^n}} \epsilon_\theta(x_{1:T}^n, a_{1:T}, n) \right) + \sqrt{\sigma_\theta} z$ 
end for
Return:  $x_{1:T}^0$ 

```

---

After training, we expect to use variational inference to generate new gestures matching the original data distribution ( $x_t^0 \sim p_\theta(x_t^0 | x_{t-\tau:t-1}^0, a_t)$ ). We follow the sampling procedure in Algorithm 2 to obtain a sample  $x_t^0$  of the current frame. The  $\sigma_\theta$  is the standard deviation of the  $p_\theta(x^{n-1} | x^n)$ . We choose  $\sigma_\theta := \tilde{\beta}^n$ .

During inferencing, instead of inputting the concatenated data with the past poses  $[x_{t-\tau-1}^0, \dots, x_{t-1}^0]$  and acoustic features  $[a_{t-\tau}, \dots, a_{t+r}]$  in our previous work [2], we only sent the raw audio to condition encoder(WavLM), then the WavLM output feed to Diffusion Model to generate the whole sequence of the accompanies gesture( $x^0$ ).

### C. Model Architecture

We proposed an extension to the diffusion models utilizing a transformer architecture as the backbone. The architecture referred to as Diffmotion-v2, is depicted in Figure 2. The architecture comprises three main components: (1) a Condition Encoder, (2) a Gesture Encoder and a Gesture Decoder, (3) stacks of Multi-head causal attention Blocks with adaptive layer norm, and (3) a Final Layer.

The Condition Encoder is a vital component for extracting and embedding speech audio features using the WavLM model while incorporating the embedded diffusion time step  $n$ . Simultaneously, the Gesture Encoder processes the input gesture sequence and transforms it into a latent representation. To capture the complex relationship between speech and gestures, the model employs stacks of Multi-head causal attention Blocks with adaptive layer norm Blocks, which are stacked  $M$  times. This architecture enables the model to effectively capture the dependencies and interactions between speech and gestures, facilitating the generation of realistic and coherent co-speech gestures. To ensure precise generation of the gesture sequence, the Final Layer incorporates an adaptive layer-norm mechanism [38] for output noise prediction. This architectural design enhances the generation process, enabling the production of realistic and diverse gestures within the context of co-speech communication.

1) *Condition Encoder:* The condition encoder, illustrated in Figure 2b, converts raw audio input into a sequence of

speech embedding space by WavLM large-scale pre-trained model [14]. In our study, we integrated the WavLM model due to its ability to handle the intricate nature of speech audio signals effectively. The WavLM model has been extensively trained on a large-scale dataset consisting of unlabeled speech audio data, covering a wide range of tasks, Automatic Speech Recognition (ASR) and non-ASR tasks, such as Speaker Verification, Speech Recognition, Paralinguistics, Spoken content, and Emotion Recognition. The model was trained on a substantial amount of English audio, totalling 94k hours, featuring diverse speakers, topics, speaking styles, and scenarios. We believe that the WavLM model exhibits enhanced robustness and possesses the capability to extract various features, including acoustic characteristics, speaker personalities, affective information, and more, from the speech audio data. The pretraining process equips the model with the ability to capture universal latent representations, denoted as  $Z_a$ , which encapsulate the essential information contained within the speech signals.

By leveraging the capabilities of the WavLM model, we aim to enhance the performance of co-speech gesture generation tasks, as shown in Figure 2b. This approach stands in contrast to the conventional methodology that relies solely on Mel-frequency cepstral coefficients (MFCC) for audio feature extraction, as observed in our previous Diffmotion model and other related studies. The WavLM model offers promising prospects due to its ability to go beyond acoustic information and incorporate knowledge from various tasks, leading to more comprehensive and contextually relevant gesture generation.

A downsampling module is seamlessly integrated into the architecture to ensure alignment between each latent representation and the corresponding sequence of poses. This module takes the form of a Conv1D layer with a kernel size of 201, which means that each 201 lengths of target labels output of WavLM is mapped to one frame of gesture sequence. Its primary objective is to facilitate the synchronization of latent representations with the gesture sequence, enabling the generation of coherent and contextually relevant gestures.

In our investigation, we observed that utilizing only the short length of target labels to map to each gesture frame was effective locally. However, it did not capture the entirety of the sequence's style, such as the overall emotion or context. We introduced an additional Conv1D layer for the entire sequence to capture the global style information. This layer is responsible for extracting the overall style of the speech audio and producing a single token denoted as  $Z_s \in \mathbb{R}^{1 \times D}$ . We then broadcast this token and add it to the universal speech representation  $Z_a$  before the down-sampling processing. By incorporating this approach, we aim to enhance the representation of the entire sequence by considering both local and global information in the generation of co-speech gestures. Finally, the condition encoder outputs  $C_{1:T}$ , representing the combined encoded audio feature and the diffusion time step  $n$ .

2) *Gesture Encoder and Decoder:* We employ Convolution 1D with a kernel size of 3 to extract the sequence of gestures from sequential data. Convolution 1D operates by sliding the kernel across the input sequence and performing element-wise

multiplication and summation to generate feature maps [6], [19].

The selection of a kernel size of 3 is driven by its efficacy in capturing local patterns and dependencies within the sequence. It enables the model to consider neighbouring elements and effectively capture short-term temporal dependencies [39]. This is particularly advantageous in gesture sequences, where adjacent frames often exhibit specific patterns or transitions contributing to overall motion dynamics.

By utilizing a kernel size of 3, we strike a balance between capturing fine-grained details and avoiding excessive parameterization. Smaller kernel sizes may overlook important contextual information, while larger ones can introduce a higher number of parameters and increase computational complexity [40]. In our experimental analysis, we found that using a kernel size of 1 resulted in an animation jitter, underscoring the importance of an appropriate kernel size for gesture sequence extraction.

We employ a kernel size of 1 convolution instead of a fully connected layer for several reasons in the context of gesture decoding. A kernel size of 1 convolution 1D enables us to capture local dependencies and interactions within the sequence while preserving the spatial dimensionality of the data. By convolving a 1D kernel with each position in the input sequence, the model can extract meaningful features and relationships between adjacent elements [41], [42]. In contrast, a fully connected layer would necessitate connecting each input element to every output neuron, resulting in a significantly larger number of parameters and a loss of the spatial structure of the data. Furthermore, using a kernel size of 1 convolution provides flexibility in capturing local patterns and fine-grained details within the sequence. This enables the model to learn non-linear relationships between neighbouring elements, which is particularly crucial in tasks such as gesture decoding where short-term dependencies significantly contribute to understanding motion dynamics [43].

*3) The multi-head causal attention with adaptive layer norm:* In line with the prevailing trend of employing adaptive normalization layers in GANs [44] and diffusion models featuring UNet backbones [45], we opted to substitute the conventional layer norm layers within the transformer blocks with adaptive layer normalization (adaLN). The adaLN further corroborates its superiority in terms of computational efficiency, as it introduces the least number of Gflops compared to cross-attention and in-context conditioning methods [38]. Furthermore, it stands out as the sole conditioning mechanism that is constrained to apply the same function to all tokens universally.

#### D. Final Layer

The final layer consists of an adaLN and a gesture decoder with a conv  $1 \times 1$ , and finally outputs the predicted noise.

## IV. EXPERIMENTS

To demonstrate the capabilities of our proposed approach, we compare it to the best available alternatives on three

gesture-generation datasets via careful user studies. All experiments exclusively focus on full-body 3D motion generation. This particular choice poses a more arduous challenge compared to generating only upper body motion. The increased dimensionality of the output space and the visual prominence of artifacts, such as foot-skating and ground penetration.

#### A. Dataset and Data Processing

*1) Datasets:* We train and test our system on three high-quality speech-gesture datasets: Trinity [46], ZeroEGGS(ZEGGS) [26], and BEAT [47]. The Trinity dataset includes 244 minutes of motion capture and audio of a male native English speaker producing spontaneous speech on different topics. The ZEGGS dataset contains 135 minutes of full-body motion capture and speech audio from a monologue performed by a female actor speaking in English and covers 19 different motion styles. The Body-Expression-Audio-Text dataset(BEAT) has 76 hours of high-quality, multi-modal data capture from 30 speakers talking with eight motions and in four languages. We only use the mainly English speech dataset, which amounts to about 35 hours in total.

*2) Speech Audio Data Process:* The audio in Trinity was recorded at a sampling rate of 44 kHz, and in ZEGGS and BEAT was recorded at 48 kHz. Due to the WavLM large model being pre-trained on 16 kHz sampled speech audio, we resampled all the audio to 16 kHz. For MFCC, we used 16-dimensional features.

*3) Motion Data Process:* In all conducted experiments, the focus was exclusively on full-body motion. However, considering the hierarchy and data quality variations across different motion datasets, we made specific adaptations by selecting different joints for each dataset. For the Trinity Gesture Dataset, we included the whole spine, head, hands, legs, and the root joint in our analysis. In the case of the ZEGGS and BEAT datasets, we expanded our selection to include hand joints in addition to the same set of joints used in the Trinity dataset. All datasets were augmented with a mirrored version of the joint angles together with the unaltered speech. The frame rate of all datasets was downsampled to 20 frames per second (fps). To represent each joint angle, we utilized the exponential map technique [48], which helps avoid discontinuities in the joint angles. Furthermore, to mitigate constant movements that may affect the learning process, we applied a constant remover to eliminate joints with a small range of motion. This preprocessing step helps ensure the model focuses on relevant and informative joint movements.

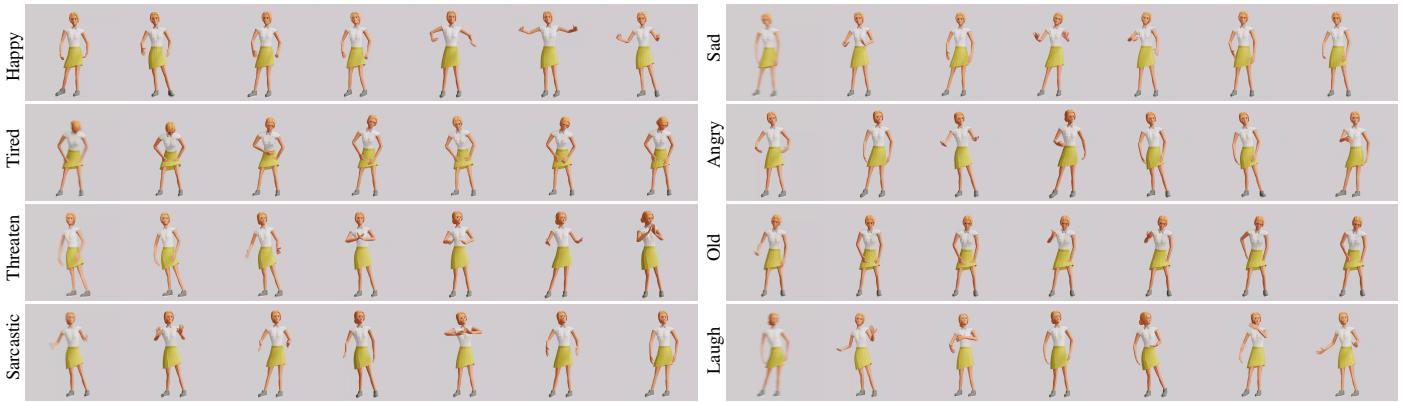
For the purpose of training, evaluation, and testing, we divided the dataset into clips, each spanning a duration of 20 seconds. This allows for a manageable and consistent input size for the model, facilitating the learning process and evaluation of performance.

#### B. Model Settings

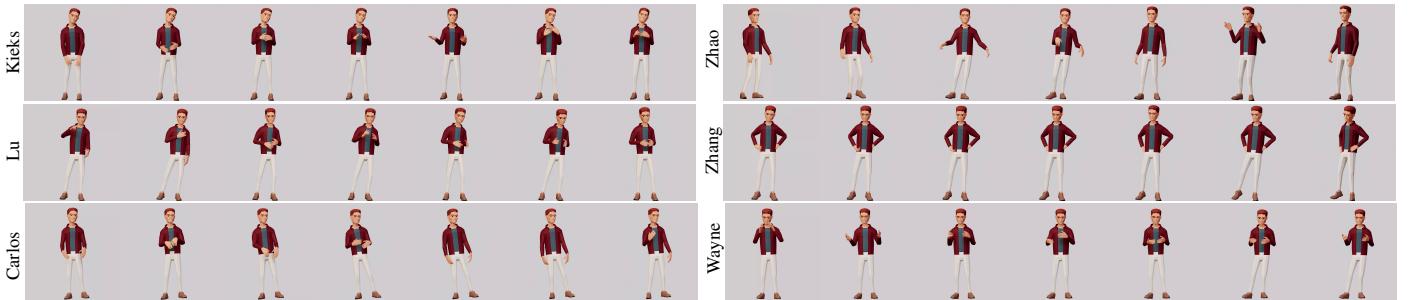
In our experiments, we utilized a stack of 12 Transformer blocks, with each block containing 16 attention heads, as shown in Fig.2c. Each frame of the gesture sequence was encoded into hidden states  $h \in \mathbf{R}^{1280}$ . For the WavLM model,



(a) The full-body gestures generated in response to the audio from Record\_008.wav and Record\_015.wav in the Trinity Dataset. The results suggest that the proposed architecture has the capacity to generate gestures that accompany acoustic, semantic, etc. and the ability to produce more relaxed and non-hyperactive gestures.



(b) The results obtained through the utilization of various styles of audio in the ZEGGS dataset. It illustrates the model's ability to generate diverse motion styles, such as emotions (happiness, sadness, anger, etc.), state (such as tired, laughing and threatening), and age-related nuances. Remarkably, this variability is achieved solely by leveraging the speech audio information extracted by WavLM without requiring any manual annotation.



(c) Extrapolating results on the BEAT dataset confirmed the ability of the proposed method to distinguish between different personalities of different people based solely on speech audio. The method effectively generates co-speech gestures that mirror the distinguishing features of the respective persons.

Fig. 3. The co-speech gesture sequences generated by training and inference on the Trinity, ZEGGS, and BEAT datasets. The results affirm the capability of the proposed approach to generate co-speech gestures that are natural, distinguishable, emotionally expressive, personalized, and stylistically varied, solely based on the raw audio input.

we employed the WavLM Base+ pre-trained model<sup>1</sup>, which consists of 12 transformer encoder layers, 768-dimensional hidden states, and 8 attention heads. This architecture resulted in a total of 94.70M untrainable parameters.

Interestingly, we discovered that there was no need to incorporate an additional positional encoding mechanism. This was because the WavLM model inherently included positional information as part of its design. This eliminates the need for an explicit positional encoding step, simplifying the overall architecture and reducing computational complexity.

The quaternary variance schedule of diffusion model starts from  $\beta_1 = 1 \times 10^{-4}$  till  $\beta_N = 5 \times 10^{-5}$  with linear beat

schedule. The number of diffusion steps is  $N = 500$ , and the training batch size is 32 per GPU.

The model is built on Torch Lightning framework, AdamW optimizer with a learning rate of  $20 \times 10^{-5}$  with LinearScheduler with  $10^3$  warm-up steps. All experiments run on an Intel i9 processor and a single A100. We train the models for about 4h(Trinity), 4h(ZeroEGGS) and 21h(BEAT). During the inference process, we divided the entire sequence into 20-second segments and stacked them together. Our system required approximately 30 seconds to generate all the subsequence. Finally, we reconnected all the sub-sequences to form the final whole sequence. Indeed, the observation that the inference time is not dependent on the length of the audio sequence is consistent with the approach of slicing the whole

<sup>1</sup><https://huggingface.co/microsoft/wavlm-base-plus>

TABLE I  
THE SYSTEM OVERVIEW AND THE SUBJECT MEAN PERCEPTUAL RATING SCORE.

Dataset	AudioEncode	Param.Count	System			Human-likeness	Appropriateness	Subject Evaluataion Metric
			Train.time	Inference Time(20s lengths)				
Trinity	GT	/	/	/		4.57±0.12	4.62±0.20	/
	WavLM	1B	4h30m	43s		4.08±0.10	4.13±0.13	4.11±0.10
	MFCC	442M	1h20m	25s		4.05±0.19	4.01±0.26	3.86±0.21
ZEGGS	GT	/	/	/		4.63±0.07	4.53±0.17	/
	WavLM	1B	5h25m	39s		4.11±0.15	4.18±0.16	4.20±0.09
	MFCC	442M	2h34m	21s		4.07±0.20	4.01±0.18	3.89±0.17
BEATS	GT	/	/	/		4.54±0.08	4.55±0.22	/
	WavLM	1B	1d3h	41s		4.13±0.09	4.11±0.15	4.10±0.22
	MFCC	442M	13h40m	20s		4.11±0.13	4.01±0.23	3.86±0.07

sequence and stacking them together for parallel inferencing.

### C. Results

Fig.I and Fig.3 show the visualization results of our system generating gestures trained on three different datasets(Trinity Fig.3a, ZeroEGGS Fig.3b, and BEAT Fig.3c). Our system adeptly generates lifelike gestures that harmonize with accompanying speech audio. Furthermore, the characters' gestures reflect an understanding of the acoustic attributes, content, personalities, and emotions encapsulated within the speech audio. For instance, the model trained on the Trinity Dataset demonstrates the capacity to produce gestures that seamlessly align with both acoustic and semantic cues. Interestingly, the ensuing sequences of motion exhibit a more tranquil and non-hyperactive demeanour. Additionally, the model trained on the ZEGGS dataset showcases proficiency in generating diverse motion styles encompassing emotions (such as happiness, sadness, and anger), states (like tiredness, laughter, and threat), and nuanced age-related characteristics. Addressing the dimension of personality, the BEAT dataset-trained model skillfully generates personality-infused gestures in response to the speech audio articulated by different individuals.

### D. Comparison

Following the conventions of recent gesture generation research, we evaluated the generated actions through a series of user studies. we only focus on compare the effects of WavLm encoder and MFCC features in our model architecture on the quality of action generation.

1) *User Study*: The ultimate goal of speech-driven gesture generation is to produce natural and convincing motions. However, objective evaluation of gesture synthesis does not always capture the subjective quality perceived by humans [10], [49], [50]. To address this, we focus on conducting subjective human perception evaluations using a question set consisting of three evaluation aspects: 1) human likeness, 2) appropriateness, and 3) Style-appropriateness. Each aspect was rated on a 5-point Likert scale.

The human-likeness aspect assesses the naturalness and resemblance of the generated gestures to the motion of an actual human, irrespective of the accompanying speech. The appropriateness aspect, on the other hand, evaluates the temporal consistency of the generated gestures, particularly in terms of their alignment with the rhythm of the speech.

Lastly, the style-appropriateness aspect gauges the degree of similarity between the generated and original gestures, thereby evaluating the appropriateness of the generated style.

By incorporating these evaluation aspects, we aimed to comprehensively assess the quality of the generated gestures in terms of their human-likeness, appropriateness, and style-likeness.

First, we trained each model and generated 20 gesture clips with the same speech audio as input. Each clip lasts for 20 seconds. Next, we randomly selected 3 clips generated by each model for valuation. Then, we retargeted the skeleton to a cartoon low-poly character asset provided unity asset store<sup>2</sup> and built up the video by Unity Game Engine<sup>3</sup> for the user study.

30 volunteer participants were recruited, including 18 males and 12 females (aged 19-23. All of them(22 from China, 8 international students from the USA, UK, etc.) are good at English. They were asked to rate the scale for the evaluation aspects. The scores were assigned from 1 to 5, representing the worst to best.

Firstly, we introduced the method to all participants and showed them some example clips which not in the valuation set. After the participants fully understood the process, we started the formal experiment. All participants were instructed to wear headphones and sit in front of a computer screen. The environment was quiet and had no interference. Participants were unaware of which method each video belonged to. The order of videos was random, but each video was guaranteed to appear three times, presented and scored by the participants.

One-way ANOVA was conducted to determine if the models' scores differed on the three evaluation aspects. The results are shown in TableI. The results indicate that there is no statistically significant difference in the human-likeness scores between the WavLM and MFCC methods across the three datasets, suggesting that both features are equally effective in generating natural gestures. However, it is worth noting that there are significant statistical differences in the scores of appropriateness and style-appropriateness, indicating that the WavLM encoding performs better in capturing motion gesture features and stylistic aspects compared to MFCC.

Furthermore, significant and substantial differences are observed between the generated gestures and ground truth (GT)

<sup>2</sup><https://assetstore.unity.com/packages/3d/characters/humanoids/polyart-characters-159791>

<sup>3</sup><https://unity.com>

in terms of human likeness and appropriateness across the three datasets. These differences can be attributed to the fact that the GT contains richer and more characteristic gestures, which might be a long-tail distribution and challenging for the model to replicate during the generation process fully.

## V. CONCLUSION

This paper introduces "diffmotion-v2," an innovative generative model for co-speech gestures that combines a non-autoregressive diffusion process and self-attention mechanisms. Our model possesses the capacity to simultaneously generate a complete sequence of co-speech gestures along with the corresponding raw speech audio waveform as input. We have incorporated the WavLM large-scale pre-trained model to encode the raw audio waveform. Additionally, an adaptive layer norm architecture has been integrated into the transformer-based layer to establish the relationship between speech information and accompanying gestures. Through comprehensive subjective evaluation experiments conducted on the Trinity, ZEGGS, and BEAT datasets, we have effectively demonstrated the natural synthesis of co-speech gestures using our model. Notably, our approach achieves this accomplishment by solely utilizing speech attributes such as acoustic characteristics, semantic information, personality traits, emotions, and more, without necessitating any form of annotation.

## REFERENCES

- [1] B. Matthew, "Voice puppetry," in *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*, 1999, pp. 21–28.
- [2] F. Zhang, N. Ji, F. Gao, and Y. Li, "Diffmotion: Speech-driven gesture synthesis using denoising diffusion model," in *MultiMedia Modeling: 29th International Conference, MMM 2023, Bergen, Norway, January 9–12, 2023, Proceedings, Part I*. Springer, 2023, pp. 231–242.
- [3] P. Wagner, Z. Malisz, and S. Kopp, "Gesture and speech in interaction: An overview," *Speech Communication*, vol. 57, pp. 209–232, 2014.
- [4] F. Ylva, N. Michael, and M. Rachel, "Multi-objective adversarial gesture generation," in *Motion, Interaction and Games*. Newcastle upon Tyne United Kingdom: ACM, 2019, pp. 1–10.
- [5] G. Ian, P.-A. Jean, M. Mehdi, X. Bing, W.-F. David, O. Sherjil, C. Aaron, and B. Yoshua, "Generative adversarial nets," *Advances in neural information processing systems*, vol. 27, 2014.
- [6] Z. Kong, W. Ping, J. Huang, K. Zhao, and B. Catanzaro, "Difwave: A versatile diffusion model for audio synthesis," *arXiv preprint arXiv:2009.09761*, 2020.
- [7] D. Rezende and S. Mohamed, "Variational inference with normalizing flows," in *International Conference on Machine Learning*. PMLR, 2015, pp. 1530–1538.
- [8] L. Dinh, D. Krueger, and Y. Bengio, "Nice: Non-linear independent components estimation," *arXiv preprint arXiv:1410.8516*, 2014.
- [9] L. Dinh, J. Sohl-Dickstein, and S. Bengio, "Density estimation using real nvp," *arXiv preprint arXiv:1605.08803*, 2016.
- [10] A. Simon, H. G. Eje, K. Taras, and B. Jonas, "Style-Controllable Speech-Driven Gesture Synthesis Using Normalising Flows," in *Computer Graphics Forum*, vol. 39. Wiley Online Library, 2020, pp. 487–496, issue: 2.
- [11] L. Jing, K. Di, P. Wenjie, Z. Xuefei, Z. Ying, H. Zhenyu, and B. Linchao, "Audio2gestures: Generating diverse gestures from speech audio with conditional variational autoencoders," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 11293–11302.
- [12] T. Sarah, W. Jonathan, G. David, and M. Iain, "Speech-Driven Conversational Agents using Conditional Flow-VAEs," in *European Conference on Visual Media Production*, 2021, pp. 1–9.
- [13] A. G. Alias Parth Goyal, N. R. Ke, S. Ganguli, and Y. Bengio, "Variational walkback: Learning a transition operator as a stochastic recurrent net," *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [14] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *Advances in Neural Information Processing Systems*, vol. 33, pp. 6840–6851, 2020.
- [15] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, "Deep unsupervised learning using nonequilibrium thermodynamics," in *International Conference on Machine Learning*. PMLR, 2015, pp. 2256–2265.
- [16] K. Rasul, A.-S. Sheikh, I. Schuster, U. Bergmann, and R. Vollgraf, "Multivariate probabilistic time series forecasting via conditioned normalizing flows," *arXiv preprint arXiv:2002.06103*, 2020.
- [17] Y. Song and S. Ermon, "Generative modeling by estimating gradients of the data distribution," *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [18] S. Alexanderson, R. Nagy, J. Beskow, and G. E. Henter, "Listen, denoise, action! audio-driven motion synthesis with diffusion models."
- [19] A. Gulati, J. Qin, C.-C. Chiu, N. Parmar, Y. Zhang, J. Yu, W. Han, S. Wang, Z. Zhang, and Y. Wu, "Conformer: Convolution-augmented transformer for speech recognition."
- [20] T. Ao, Z. Zhang, and L. Liu, "Gesturediffclip: Gesture diffusion model with clip latents."
- [21] S. Ghorbani, Y. Ferstl, D. Holden, N. F. Troje, and M.-A. Carbonneau, "Zeroeggs: Zero-shot example-based gesture generation from speech."
- [22] S. Nyatsanga, T. Kucherenko, C. Ahuja, G. E. Henter, and M. Neff, "A comprehensive review of data-driven co-speech gesture generation."
- [23] J. Windle, D. Greenwood, and S. Taylor, "Uea digital humans entry to the geneva challenge 2022," in *GENEA: Generation and Evaluation of Non-Verbal Behaviour for Embodied Agents Challenge \$\{\\$\backslash\\$backslash\\$&*.
- [24] Alexanderson Simon, Henter Gustav Eje, Kucherenko Taras, and Beskow Jonas, "Style-controllable speech-driven gesture synthesis using normalising flows," in *Computer Graphics Forum*, vol. 39, no. 2. Wiley Online Library, pp. 487–496.
- [25] Taylor Sarah, Windle Jonathan, Greenwood David, and Matthews Iain, "Speech-driven conversational agents using conditional flow-vaes," in *European Conference on Visual Media Production*, pp. 1–9.
- [26] S. Ghorbani, Y. Ferstl, D. Holden, N. F. Troje, and M.-A. Carbonneau, "Zeroeggs: Zero-shot example-based gesture generation from speech," in *Computer Graphics Forum*, vol. 42, no. 1. Wiley Online Library, pp. 206–216.
- [27] Yoon Youngwoo, Ko Woo-Ri, Jang Minsu, Lee Jaeyeon, Kim Jaehong, and Lee Geehyuk, "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)*, Howard A., Althoefer K., Arai F., Arrichiello F., Caputo B., Castellanos J., Hauser K., Isler V., Kim J., Liu H., Oh P., Santos V., Scaramuzza D., Ude A., Voyles R., Yamane K., and Okamura A., Eds., pp. 4303–4309.
- [28] Kucherenko Taras, Jonell Patrik, van Waveren Sanne, Henter Gustav Eje, Alexandersson Simon, Leite Iolanda, and Kjellström Hedvig, "Gesticulator: A framework for semantically-aware speech-driven gesture generation," in *Proceedings of the 2020 International Conference on Multimodal Interaction*, pp. 242–250.
- [29] Devlin Jacob, Chang Ming-Wei, Lee Kenton, and Toutanova Kristina, "Bert: Pre-training of deep bidirectional transformers for language understanding."
- [30] Bhattacharya Uttaran, Rewkowski Nicholas, Banerjee Abhishek, Guhan Pooja, Bera Aniket, and Manocha Dinesh, "Text2gestures: A transformer-based network for generating emotive body gestures for virtual agents," in *2021 IEEE Virtual Reality and 3D User Interfaces (VR)*. IEEE, pp. 1–10.
- [31] Pennington Jeffrey, Socher Richard, and Manning Christopher D., "Glove: Global vectors for word representation," in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1532–1543.
- [32] Mikolov Tomas, Sutskever Ilya, Chen Kai, Corrado Greg S., and Dean Jeff, "Distributed representations of words and phrases and their compositionality," in *Advances in Neural Information Processing Systems*, pp. 3111–3119.
- [33] Bojanowski Piotr, Grave Edouard, Joulin Armand, and Mikolov Tomas, "Enriching word vectors with subword information," vol. 5, pp. 135–146.
- [34] S. Chen, C. Wang, Z. Chen, Y. Wu, S. Liu, Z. Chen, J. Li, N. Kanda, T. Yoshioka, and X. Xiao, "Wavlm: Large-scale self-supervised pre-training for full stack speech processing," vol. 16, no. 6, pp. 1505–1518.
- [35] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, "Deep unsupervised learning using nonequilibrium thermodynamics," in *International Conference on Machine Learning*. PMLR, pp. 2256–2265.

- 
- [36] L. Paul, “sur la théorie du mouvement brownien,” *C. R. Acad. Sci.*, vol. 65, no. 11, pp. 146,530–533, 1908, publisher: American Association of Physics Teachers.
  - [37] K. Rasul, C. Seward, I. Schuster, and R. Vollgraf, “Autoregressive denoising diffusion models for multivariate probabilistic time series forecasting,” in *International Conference on Machine Learning*, 2021, pp. 8857–8868.
  - [38] W. Peebles and S. Xie, “Scalable diffusion models with transformers.”
  - [39] Z. Li, F. Liu, W. Yang, S. Peng, and J. Zhou, “A survey of convolutional neural networks: Analysis, applications, and prospects.”
  - [40] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” vol. 60, no. 6, pp. 84–90.
  - [41] Y. Chen, “Convolutional neural network for sentence classification.”
  - [42] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778.
  - [43] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, “Going deeper with convolutions,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1–9.
  - [44] A. Brock, J. Donahue, and K. Simonyan, “Large scale gan training for high fidelity natural image synthesis,” *arXiv preprint arXiv:1809.11096*, 2018.
  - [45] P. Dhariwal and A. Nichol, “Diffusion models beat gans on image synthesis,” *Advances in neural information processing systems*, vol. 34, pp. 8780–8794, 2021.
  - [46] Ferstl Ylva and McDonnell Rachel, “Investigating the use of recurrent motion modelling for speech gesture generation,” in *Proceedings of the 18th International Conference on Intelligent Virtual Agents*, pp. 93–98.
  - [47] 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.”
  - [48] Grassia F. Sebastian, “Practical parameterization of rotations using the exponential map,” vol. 3, no. 3, pp. 29–48.
  - [49] 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.
  - [50] T. Kucherenko, P. Wolfert, Y. Yoon, C. Viegas, T. Nikolov, M. Tsakov, and G. E. Henter, “Evaluating gesture-generation in a large-scale open challenge: The genea challenge 2022.”