

ExpGest: Expressive Speaker Generation Using Diffusion Model and Hybrid Audio-Text Guidance

Yongkang Cheng*
Tencent AILab
Shenzhen, China
ykcheng@tencent.com

Mingjiang Liang*
University of Technology Sydney
Sydney, Australia
mingjiang.liang@student.uts.edu.au

Shaoli Huang
Tencent AILab
Shenzhen, China
shaol.huang@gmail.com

Jifeng Ning[†]
Northwest A&F University
Xi'an, China
njf@nwsuaf.edu.cn

Wei Liu
University of Technology Sydney
Sydney, Australia
wei.liu@uts.edu.au

Abstract—Existing gesture generation methods primarily focus on upper body gestures based on audio features, neglecting speech content, emotion, and locomotion. These limitations result in stiff, mechanical gestures that fail to convey the true meaning of audio content. We introduce ExpGest, a novel framework leveraging synchronized text and audio information to generate expressive full-body gestures. Unlike AdaIN or one-hot encoding methods, we design a noise emotion classifier for optimizing adversarial direction noise, avoiding melody distortion and guiding results towards specified emotions. Moreover, aligning semantic and gestures in the latent space provides better generalization capabilities. ExpGest, a diffusion model-based gesture generation framework, is the first attempt to offer mixed generation modes, including audio-driven gestures and text-shaped motion. Experiments show that our framework effectively learns from combined text-driven motion and audio-induced gesture datasets, and preliminary results demonstrate that ExpGest achieves more expressive, natural, and controllable global motion in speakers compared to state-of-the-art models.

Index Terms—Gesture generation, multimodal learning, emotion guide, locomotion control

I. INTRODUCTION

In fields like virtual agents, movies, and human-computer interaction, virtual speakers convey information through co-speech gestures related to audio melody and content, and non-spontaneous movements. The emerging field of audio gesture generation has gained attention, with early research using rule-based methods [1], while data-driven techniques [2] improved diversity using statistical models. Deep models, such as VAE [3], flow models [4], and diffusion-based models, generate gestures directly from raw audio data. Methods combining audio melody and semantics [5] have advanced significantly. However, DiffStyleGesture [6] and Emog [7], which use emotion as guidance, perform poorly on the BEAT dataset [8].

In non-spontaneous motion generation, early methods [9]

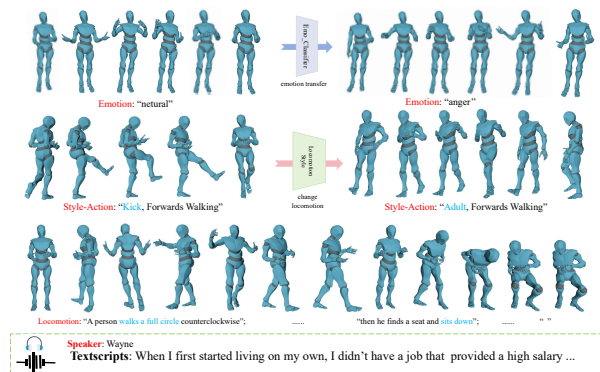


Fig. 1. **Our method demonstrates multimodal-driven effects.** The first row showcases the emotional control capability of gestures driven by audio alone, while the second row exhibits the motion style transferability driven by a combination of phrases and audio. The third row presents the results of long-frame textual descriptions and audio jointly driving the process.

considered deterministic signal-to-motion mapping using neural networks. The complexity and randomness of motion led to deeper research on generative models. TEMOS [10] pioneered action-to-action and text-to-action tasks. Recently, KIT-ML [11] and HumanML3D [12] datasets promoted text-guided motion generation. T2M-GPT [13] uses autoregressive models, improving results but with jitter issues. Diffusion model-based methods (e.g., MDM [14], MLD [15]) enhance motion data continuity through denoisers. Despite advancements, bottlenecks persist. No work has coherently integrated both motion categories, and different motion representations lead to dataset inconsistency, keeping the field in the single modality generation stage.

In this paper, we propose a diffusion model-based method aimed at using input text, audio, or a combination of both to guide the generation of expressive and diverse high-quality speakers. As shown in Figure 1, when only audio is input, our model aims to generate highly expressive and rich gestures.

*These authors contributed equally and [†] is corresponding author.

When inputting a text-audio mixed modality, it generates co-speech gesture results with non-spontaneous motion. Specifically, for text-generated non-spontaneous motion, we adopt MDM as the baseline and address its lack of melody perception by encoding spectral features into text embeddings, thus alleviating the global motion issue in generated speakers unrelated to melody. For audio-generated gestures, we observe an objective fact: fingers and limbs exhibit different sensitivities to the two attributes (melody and semantics) in the audio. For instance, when calmly saying “one, two, three”, the arm tends to remain relatively still, while the fingers display the primary changes. In contrast, when the tone varies, significant changes occur in arm movements. Based on this observation, ExpGest is the first method to decouple fingers and limbs, assigning different weights to semantics and melody to guide the generation of gesture sequences that align with both speech content and melodic variations. Furthermore, we decouple emotions from the computation graph and optimize the noise classifier for emotion stylization, reasonably endowing gestures with emotional diversity without damaging the original semantic and melodic information. The main contributions of our work are: (1) We propose ExpGest, to the best of our knowledge, the first motion speaker generation framework under mixed control that combines audio-to-gesture and text-to-motion. (2) We decouple gesture components and introduce a semantic alignment module in the latent space. We separately assign stronger melody relevance and semantic relevance to arms and fingers, generating gesture poses that better express audio content. (3) We introduce a noise emotion classifier into the reverse diffusion process, controlling the emotional style of gestures by optimizing noise gestures through gradient backpropagation. (4) We demonstrate through extensive experiments that the naturalness and richness of generated speaker actions have been improved, surpassing existing methods in terms of motion quality.

II. APPROACH

Distinct from existing co-speech gesture generation methods [6], [16] that focus on upper-body or partial movements, **ExpGest** aims to create expressive talkers with full-body motion corresponding to textual descriptions and co-speech gestures aligned with audio. The overall network framework is shown in Figure 2. We first introduced a unified data representation in Section 2.1. Then, in Section 2.2, we presented a learning framework based on the diffusion model, effectively capturing expressive full-body gestures influenced by both speech and text information through combined text and audio guidance. Finally, in Sections 2.3 and 2.4, we proposed the semantic alignment module and emotion guidance, ensuring that our gesture results conform to semantic descriptions while maintaining emotional expressiveness.

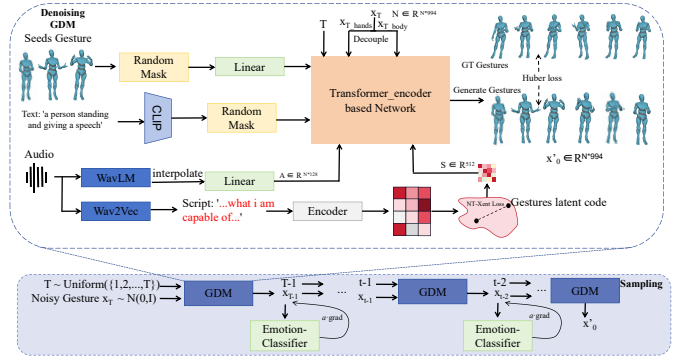


Fig. 2. **Architecture Diagram.** The upper part is the denoising model GDM. Noise step T , along with pure Gaussian noise and conditions (text description and audio), is fed into the model as input sequences. The lower part is the sampling step, where we predict x'_0 through the denoising process and add noise to x_{t-1} via the diffusion process. Subsequently, x_{t-1} is input into the noise emotion classifier to optimize the noise at that step, and the optimized noise is then passed back to the GDM. This cycle continues until $t=T$ becomes $t=0$.

A. Unified Data Representation

Our method requires the correct preservation of features from different motion datasets, necessitating the unification of various data representations. First, we extract the Euler angles from motion capture data (BVH format) and convert them into rot6D representations corresponding to the 55 joints in SMPL-X [17]. For 3D locations, we randomly extract a skeleton as a reference and align other data to this skeleton to unify the 3D coordinates. We then scale the root node’s displacement based on the extracted skeleton, adjusting the orientation to ensure consistency across all data. Finally, we combine the extracted rot6D rotation representation ($J \times 6$), 3D location ($J \times 3$), linear velocity ($J \times 3$), angular velocity ($J \times 6$), and ground contact signals (4) into a kinematic feature representation, where J represents the number of joints. Each frame of the unified motion sequence has a 994-dimensional feature representation, denoted as $x_0 \in \mathbb{R}^{N \times 994}$, where N is the number of motion sequence frames. Due to the lack of mixed-modality data, we artificially synthesized some data. Specifically, we considered the lower body (including the root) below the third spine as the locomotion-involved part, and the remaining upper body as the gesture-involved part. We concatenated the two parts, resulting in 20K artificially generated text-audio-motion matching pairs. We used a ratio of (0.6, 0.4) to mix gesture data and artificial data as the training dataset.

B. Diffusion Model for Generating Motion Speakers

ExpGest builds upon the diffusion model [18] to generate co-speech gestures. As illustrated in Figure 2, our approach involves learning how to progressively denoise from pure Gaussian noise to obtain expressive full-body gestures. This process encompasses a forward diffusion step, where noise is added to the original gestures, and a subsequent reverse denoising procedure.

Denoising Diffusion Probabilistic Model. We sample clean

gestures x_0 from the distribution of real gestures $x_0 \sim q(x_0)$, where $q(x_0)$ represents the distribution of the original data. Subsequently, Gaussian noise is incrementally added to x_0 . When the number of noise steps T is sufficiently large, the final noisy gesture x_T converges to pure Gaussian noise. Each step is modeled as a forward process denoted by q . The forward process diffuses the data samples through Gaussian transitions parameterized with a Markov process:

$$\begin{aligned} q(x_t|x_{t-1}) &= \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I), \\ &= \mathcal{N}\left(\sqrt{\frac{\alpha_t}{\alpha_{t-1}}}x_{t-1}, \left(1 - \frac{\alpha_t}{\alpha_{t-1}}\right)I\right), \end{aligned} \quad (1)$$

where $\{\beta_t\}_{t=1}^T$ is the variance schedule and $\alpha_t = \prod_{s=1}^t (1 - \beta_s)$. Then the reverse process becomes $p_\theta(x_{0:T}) := p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)$, starting from $x_T \sim \mathcal{N}(0, I)$ with noise predictor ϵ_t^θ :

$$x_{t-1} = \frac{1}{\sqrt{1 - \beta_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \alpha_t}} \epsilon_t^\theta(x_t) \right) + \sigma_t z_t, \quad (2)$$

where $z_t \sim \mathcal{N}(0, I)$ and $\sigma_t^2 = \beta_t$ means the variance schedule stays constant.

Expressive Gesture Diffusion Generator. The objective of ExpGest is to generate motion-capable talking avatars of arbitrary lengths, accompanied by semantically relevant and expressively rich synchronized gestures. However, unlike the original DDPM [18], we take into account the extreme physical constraints inherent in 3D human bodies and replace the predicted noise with the original human representation, deviating from image generation. As a result, in each step of our denoising process, we reconstruct the original representation from pure Gaussian noise, ultimately producing the final generated result through a cyclic process of noise addition and denoising:

$$\begin{aligned} \hat{x}_0 &= \frac{x_t - \sqrt{1 - \alpha_t} \epsilon_t^\theta(x_t|c)}{\sqrt{\alpha_t}}, \\ x_{t-1} &= \sqrt{\alpha_{t-1}} \hat{x}_0 + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \epsilon_t^\theta(x_t|c) + \sigma_t z_t, \end{aligned} \quad (3)$$

where c is the conditions. As shown in the upper part of Figure 2, our conditions include noise step t , seed posture, motion text, audio information, and semantic latent code. The noise step t and seed posture are projected to the same dimensions through MLP and Linear layers, respectively, and then added together. The audio is encoded using WavLM, and to match the number of gesture frames, it is interpolated in the time dimension. The text description is directly encoded into the CLIP space, while the semantic features are obtained through our semantic alignment module. Finally, we perform Huber loss on \hat{x}_0 reconstructed at each step and x_0 sampled in the real distribution to optimize ExpGest:

$$\mathcal{L}_{Huber} = E_{x_0 \sim q(x_0|c), t \sim [1, T]} [HuberLoss(x_0 - \hat{x}_0)], \quad (4)$$

It is worth noting that during the denoising process, if there is only a single modality input, the other modality is directly masked and added to the conditions in the form of \emptyset .

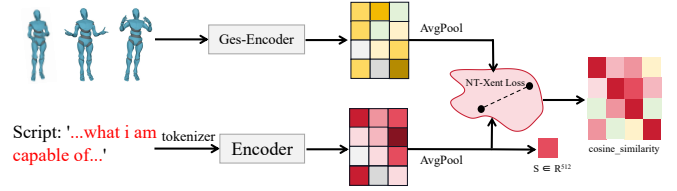


Fig. 3. **Semantic Alignment Module.** We employ contrastive learning to encode gestures and audio semantics into a shared latent space and achieve alignment in the latent space.

C. Semantic Alignment in Latent Space

In previous gesture generation methods, the crucial aspect of semantic alignment was often overlooked, posing a significant challenge for accurately generating semantically aligned actions due to the many-to-many mapping relationship between audio content and gesture sequences. To tackle this issue, we learn a joint embedding space for gestures and audio transcriptions, allowing their alignment in latent space and uncovering semantic associations between both modalities. As depicted in Figure 3, we initially train gesture and transcription encoders using a motion-VAE structure and a BERT tokenizer [19], respectively. We parameterize transcriptions into tokenized word embedding sequences and linearly map them to a space with the same dimensions as the gesture latent code. We employ global average pooling to extract global information from both modalities and utilize CLIP-style contrastive learning to fine-tune the encoders. The NT-Xent [20] Loss serves as the objective to maximize the similarity of transcription-gesture matching pairs in latent space while minimizing the similarity of non-matching pairs. Formally, the loss function is as follows:

$$\mathcal{L}(s, g) = -\log \frac{\exp(\text{sim}(z_s, z_g)/\tau)}{\sum_{k \in K} \exp(\text{sim}(z_s, z_k)/\tau)}, \quad (5)$$

where, z_s and z_g are the latent space representations of a matching transcription-gesture pair. sim is similarity score between two latent codes, K is a set containing one positive sample transcription and a group of negative sample gestures, and τ is the temperature parameter used to adjust the sensitivity of the function. Finally, we freeze the trained semantic alignment module and deploy only the transcription encoder into the GDM, ensuring that the final generated results accurately capture the semantic content.

D. Noise-based Emotion Guided Classifier

Previous methods transformed emotions into a set of one-hot encodings embedded directly into the conditions, controlling emotions by appending identity. However, this approach fails to capture the underlying relationships and continuity between emotions, making it difficult for the model to learn smooth transitions between different emotions. To address this issue, we introduce a noise emotion classifier that optimizes the denoised noise gesture at each sampling step, moving it towards a human-specified emotional direction. This method is decoupled from the diffusion computation graph, preserving the original semantic content and making it applicable to different structures in the same representation. Specifically, we randomly select emotion-gesture matching pairs and add Gaussian noise of random steps to the gestures, turning them into noisy gestures. We then train the noise emotion

classifier based on the noisy gesture-emotion matching pair data. After training, we graft the classifier onto the GDM, input the denoised x_t into the GDM, compute the gradient with the specified emotion, and then backpropagate it to x_t for optimization. The formula is as follows:

$$\hat{x}_t = x_t + \alpha \cdot \nabla_{x_t} \mathcal{L}(\text{Emo_Classifier}(x_t), y), \quad (6)$$

where α is gradient update weight and y is the real emotion label.

III. EXPERIMENTS

In this section, we discuss ExpGest’s technical details, evaluation, comparison with other frameworks, potential applications, and design choices validation through ablation studies.

Implementation Details. We first downsample all data to 20 FPS to unify the frame rate. For text-motion data, we randomly truncate long-frame data to a motion length of 180 frames, and for short-frame data, we select a length of 60 to 180 frames, padding with zeros to extend to a uniform 180-frame length. We set the maximum text length for CLIP encoding to 20. For audio-gesture data, we randomly truncate the audio and corresponding gesture sequence to 180 frames. For artificially synthesized data, we only use motions with a large displacement for the lower body (including running, walking, jumping, standing, sitting, etc.) and replace the gesture data below the third spine with motion data, removing some unnatural sequences. For the semantic alignment module, we align transcription texts and gesture motions in a shared embedding space using a VAE encoder and a BERT tokenizer. For the noise emotion classifier module, we project the noise gesture motion representation and predict 8-bit one-hot encoded emotion categories. We construct a diffusion model using a 12-layer Transformer encoder and set the diffusion steps to 1000. We train the model for approximately 72 hours on 4 NVIDIA Tesla V100 GPUs, totaling 800K steps.

A. Dataset and Evaluation Metrics

Datasets. Our audio data comes from the large-scale, high-quality **BEAT** dataset [8], featuring 76 hours of multimodal speech data from 30 speakers in 8 emotions and 4 languages. We selected the **English** data from both presentations and conversations for training. We also used **AMASS** [21] and **100-STYLE** [22] for locomotion training. From **AMASS**, we selected motion pairs with text descriptions, containing only the body’s SMPL representation and zero-filled hand representations. **100-STYLE**, a multi-style motion dataset, includes various actions each with a BVH file and short phrase description. We processed these data into the motion representation described in Sec 2.1 for unified training.

Evaluation Metrics. We use three metrics to evaluate the quality, emotion alignment, and semantic matching of gestures generated by ExpGest. Gesture quality is assessed using Fréchet Gesture Distance (FGD) [23], which calculates the distance between latent feature distributions of generated and real gestures, evaluating gesture quality. Lower FGD values indicate higher motion quality. Additionally, we propose a

TABLE I
QUANTITATIVE ASSESSMENT OUTCOMES. BOLDFACE DENOTES THE TOP-PERFORMING METHOD FOR EACH EVALUATION CRITERION. ALL EXPERIMENTS ARE CONDUCTED AND APPRAISED ON THE BEAT DATASET.

Method	<i>FGD(Raw)</i>	<i>FGD(Feature)</i>	<i>EA</i>	<i>SA</i>
Ground Truth	0.0	0.0	0.97	0.91
CAMN [8]	52.4	263.9	-	-
Trimodal [23]	47.6	212.7	-	-
DiffStyleGesture [6]	33.7	133.9	0.60	0.11
ExpGest	11.7	76.6	0.91	0.61
ExpGest w/ hybrid-guided	25.4	129.3	-	-

TABLE II
EMOTION-GUIDED RESULTS. THE BOLD VALUES REPRESENT THE BEST RESULTS. ALL METHODS ARE TRAINED ON THE SAME DATASET.

Method	<i>EA</i>	<i>EC</i>	<i>EA_{hands}</i>	<i>EC_{hands}</i>
Ground Truth	0.97	-	0.95	-
DiffStyleGesture [6]	0.60	0.27	0.49	0.19
DiffStyleGesture w/ EG	0.83	0.69	0.70	0.63
ExpGest	0.91	0.83	0.81	0.70

Semantic Alignment (SA) metric to assess the semantic consistency between generated gestures and audio. We use the trained VAE encoder to encode transcribed audio text and gestures into a shared latent space and evaluate their pre-established similarity as the assessment score. The calculation formula is as follows:

$$SA = \cos(\text{avg_pool}(V_g(G_{pred})), \text{avg_pool}(V_s(S))), \quad (7)$$

where G represents the gestures generated by the model, and S denotes the hidden states encoded by the BERT [8] model after tokenizing the transcribed text, serving as a representation of the semantics. We also introduce Emotion Alignment (EA) and Emotion Control Success Rate (EC) to evaluate the model’s stylization capabilities. Unlike in Section 2.4, we pre-train a gesture emotion classifier, which doesn’t require adding noise to the gesture motion but directly classifies emotions in the original gestures. EA represents consistency with the true emotions of the audio, while EC denotes the success rate of style transfer for specified emotions. Higher values for both EA and EC indicate more accurate emotion guidance.

$$EC_{|A} = \frac{\text{Emo_Classifier}(G_{pred}, \text{Emotions})}{\text{Sum}(\text{Emotions})}, \quad (8)$$

B. Comparison to State-of-the-art Methods

We compare our method with three cutting-edge approaches: Trimodal [23], CaMN [8], and DiffStyleGesture [6]. We remove the text guidance describing motion and use an audio-only mode for comparison. Quantitative results in Table I show our method consistently outperforms all other methods across all evaluations. Specifically, in terms of Fréchet Gesture Distance, we achieve state-of-the-art performance in the audio-only mode. Our method improves the feature space by 57.3 (42.7%) and the raw space by 22 (65.2%) in evaluations with only upper-body limbs compared to our baseline. Although there is a slight decrease in the mixed mode, our method still outperforms the baseline. Thanks to

TABLE III

95% CONFIDENCE INTERVAL FOR USER STUDY AVERAGE SCORE. SINCE THE GT DATA FROM THE BEAT DATASET LACKS GLOBAL LOCOMOTION, WE ONLY COMPARE GESTURE SEQUENCES GENERATED BY EXPGEST IN AUDIO-ONLY MODE WITH GT AND DIFFSTYLEGESTURE, USING USER RATINGS. FOR THE GLOBAL-COHERENCE METRIC, WE COMBINE ACTION OR TEXT DESCRIPTIONS WITH AUDIO AS INPUT AND HAVE USERS EVALUATE THE OVERALL COORDINATION OF THE GENERATED SPEAKER MOTION.

Method	Human-likeness	Gesture-appropriateness	Emotion-compatibility	Global-coherence
Ground Truth	4.61 ± 0.17	4.72 ± 0.20	3.62 ± 0.27	-
DiffStyleGesture [6]	3.46 ± 0.16	3.61 ± 0.11	1.22 ± 0.19	-
ExpGest (audio-only)	4.17 ± 0.13	4.20 ± 0.18	3.88 ± 0.21	-
ExpGest (audio-action)	3.76 ± 0.08	3.83 ± 0.11	2.63 ± 0.09	3.61 ± 0.17
ExpGest (audio-text)	3.72 ± 0.10	3.91 ± 0.15	2.27 ± 0.06	3.47 ± 0.16

the semantic-gesture joint embedding space and using aligned semantic features as guidance conditions, our model achieves improvements in SA. Furthermore, by optimizing the noise gradient and guiding generation towards specified emotion styles, our method significantly outperforms existing methods in EA and EC metrics.

C. Is the emotion-guided classifier helpful?

To evaluate our emotion classifier’s guidance capability, we randomly select two to four emotions from the eight available in the BEAT [8] dataset and assign them to each audio segment in the validation set. We feed the generated results into a pre-trained emotion classifier to determine the emotion category and compare the output with the ground truth emotion labels. As shown in Table II, our baseline method significantly outperforms DiffStyleGesture [6] in emotion-alignment and emotion-controllability metrics. We also provide an emotion evaluation metric for fingers, observing noticeable changes in some emotions. For instance, under the angry emotion, the virtual character often displays pointing gestures while speaking. To substantiate our emotion classifier’s effectiveness, we incorporate it into the DiffStyleGesture method’s inference stage for evaluation. Using the same data representation ensures a fair comparison with consistent evaluation logic and classifier model across both methods. Table II shows that, with the emotion classifier, DiffStyleGesture achieves higher emotion alignment compared to the original approach (relying on one-hot encoding for emotion representation).

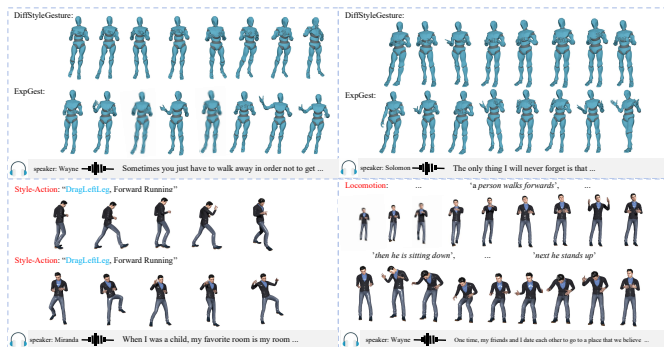


Fig. 4. The first row compares our method’s generative performance in an audio-guided scenario with state-of-the-art techniques, yielding more expressive gestures. The second row demonstrates a speaker’s generation with rich movements, guided by audio combined with action or text. See **Supp** for more videos.

TABLE IV

ABLATION STUDY. DEMONSTRATING THE EFFECTIVENESS OF SEMANTIC ALIGNMENT MODULE AND EMOTION GUIDANCE MODULE.

Method	FGD (raw)	FGD (feature)	SA	EA
DiffStyleGesture [6]	33.7	133.9	0.11	0.60
Baseline	25.9	95.5	0.14	0.57
Baseline w/ EG	28.2	115.3	0.10	0.89
Baseline w/ SA	18.6	89.1	0.63	0.62
ExpGes	11.7	76.6	0.61	0.91

D. User Study and Results

User Study. To evaluate ExpGest’s real-world visual performance, we conducted a user study comparing generated gesture sequences under various guidance modes, DiffStyleGesture, and ground truth data. We curated 40 audio clips from the BEAT test set, featuring diverse voices, and divided them into short (8-16 seconds) and long (60-90 seconds) segments. About 120 participants evaluated the gesture sequences rendered into dynamic virtual speakers using Blender. We employed four evaluation metrics: “Human-likeness”, “Gesture-appropriateness”, “Emotion-compatibility”, and “Global-coherence”, each scored on a 1-5 scale. Table III presents the average scores, showing ExpGest outperforms previous methods in visual performance and gains higher acceptance from participants. The mixed guidance approach demonstrates remarkable synergy in generated gestures.

Results. Figure 4 (top) demonstrates a comparison of ExpGest with prior state-of-the-art methods under single audio guidance. We uniformly selected serene recordings from Wayne and Solomon, two speakers in the BEAT dataset, as test audio. The generation results of DiffStyleGesture, which relies solely on melody features, are relatively static with little variation. In contrast, ExpGest, benefiting from semantic alignment in the latent space, reveals highly expressive gesture results. Furthermore, our proposed noise emotion guide can naturally and smoothly perform emotional transitions, thereby enhancing the expressiveness of the generated results for reference). Figure 4 (bottom) showcases the superior performance of ExpGest under mixed control. By integrating audio with action, ExpGest can generate a variety of motion sequences, such as “dragging left leg, running forward”. Simultaneously, using existing large language models to edit the action text description for a speech (“he looks for a place to sit while talking”), combined with

the audio itself, a vivid speaking character can be generated.

E. Ablation Studies

we designed the following ablation experiments. The results are detailed in Table IV. Our baseline method has already achieved significant improvements compared to DiffStyleGesture. Subsequently, we incorporated an Emotion-Guided Classifier (EG) during the reverse diffusion process, optimizing noise using gradient backpropagation. Although the performance in FGD has decreased, it still outperforms previous methods and significantly improves emotion alignment. Then, we introduced the Semantic Alignment (SA) module. Compared to the baseline, FGD has increased significantly, indicating the importance of semantic alignment in gesture generation, and since there are no additional optimization operations in reverse diffusion, the generation speed is also close to DiffStyleGesture, taking approximately 20 seconds to generate a 180-frame gesture slice. In this experiment, we conducted ablation studies only on the proposed baseline, without involving mixed guidance. All evaluation results were generated using pure audio guidance.

IV. CONCLUSIONS

In this paper, we propose ExpGest, a novel framework using the diffusion generation model to create motion speakers from mixed audio-text guidance. Training on heterogeneous gesture-audio and text-motion pairs allows effective capture of human motion nuances. Our approach generates natural-looking speaker gesture sequences, laying the groundwork for large-scale motion speaker generation with potential applications in virtual agents, movies, and human-computer interaction.

REFERENCES

- [1] Michael Kipp, *Gesture generation by imitation: From human behavior to computer character animation*, UniversalPublishers, 2005.
- [2] Ikhsanul Habibie, Mohamed Elgharib, Kripasindhu Sarkar, Ahsan Abdullah, Simbarashe Nyatsanga, Michael Neff, and Christian Theobalt, "A motion matching-based framework for controllable gesture synthesis from speech," in *ACM SIGGRAPH 2022 Conference Proceedings*, 2022, pp. 1–9.
- [3] Diederik P Kingma and Max Welling, "Auto-encoding variational bayes," 2013.
- [4] Sheng Ye, Yu-Hui Wen, Yanan Sun, Ying He, Ziyang Zhang, Yaoyuan Wang, Weihua He, and Yong-Jin Liu, "Audio-driven stylized gesture generation with flow-based model," in *European Conference on Computer Vision*. Springer, 2022, pp. 712–728.
- [5] Tenglong Ao, Qingzhe Gao, Yuke Lou, Baoquan Chen, and Libin Liu, "Rhythmic gesticulator: Rhythm-aware co-speech gesture synthesis with hierarchical neural embeddings," 2022, vol. 41, pp. 1–19, ACM New York, NY, USA.
- [6] Sicheng Yang, Zhiyong Wu, Minglei Li, Zhensong Zhang, Lei Hao, Weihong Bao, Ming Cheng, and Long Xiao, "Diffusestylegesture: Stylized audio-driven co-speech gesture generation with diffusion models," 2023.
- [7] Lianying Yin, Yijun Wang, Tianyu He, Jinming Liu, Wei Zhao, Bohan Li, Xin Jin, and Jianxin Lin, "Emog: Synthesizing emotive co-speech 3d gesture with diffusion model," 2023.
- [8] Haiyang Liu, Zihao Zhu, Naoya Iwamoto, Yichen Peng, Zhengqing Li, You Zhou, Elif Bozkurt, and Bo Zheng, "Beat: A large-scale semantic and emotional multi-modal dataset for conversational gestures synthesis," in *European Conference on Computer Vision*. Springer, 2022, pp. 612–630.
- [9] Partha Ghosh, Jie Song, Emre Aksan, and Otmar Hilliges, "Learning human motion models for long-term predictions," in *2017 International Conference on 3D Vision (3DV)*. IEEE, 2017, pp. 458–466.
- [10] Mathis Petrovich, Michael J Black, "Temos: Generating diverse human motions from textual descriptions," in *European Conference on Computer Vision*. Springer, 2022, pp. 480–497.
- [11] Matthias Plappert, Christian Mandery, and Tamim Asfour, "The kit motion-language dataset," 2016, vol. 4, pp. 236–252, Mary Ann Liebert, Inc. 140 Huguenot Street, 3rd Floor New Rochelle, NY 10801 USA.
- [12] Chuan Guo, Shihao Zou, Xinxin Zuo, Sen Wang, Wei Ji, Xingyu Li, and Li Cheng, "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.
- [13] Jianrong Zhang, Yangsong Zhang, Xiaodong Cun, ShaoliHuang, Yong Zhang, Hongwei Zhao, Hongtao Lu, and Xi Shen, "T2m-gpt: Generating human motion from textual descriptions with discrete representations," 2023.
- [14] Guy Tevet, Sigal Raab, Brian Gordon, Yonatan Shafir, Daniel Cohen-Or, and Amit H Bermano, "Human motion diffusion model," 2022.
- [15] Xin Chen, Biao Jiang, Wen Liu, Zilong Huang, Bin Fu, Tao Chen, and Gang Yu, "Executing your commands via motion diffusion in latent space," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 18000–18010.
- [16] Tenglong Ao, Zeyi Zhang, and Libin Liu, "Gesturediffclip: Gesture diffusion model with clip latents," 2023.
- [17] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black, "Smpl: A skinned multiperson linear model," 2015, vol. 34, pp. 1–16, ACM New York, NY, USA.
- [18] Jonathan Ho, Ajay Jain, and Pieter Abbeel, "Denoising diffusion probabilistic models," 2020, vol. 33, pp. 6840–6851.
- [19] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," 2018.
- [20] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton, "A simple framework for contrastive learning of visual representations," in *International conference on machine learning*. PMLR, 2020, pp. 1597–1607.
- [21] Naureen Mahmood, Nima Ghorbani, Nikolaus F Troje, Gerard Pons-Moll, and Michael J Black, "Amass: Archive of motion capture as surface shapes," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 5442–5451.
- [22] Ian Mason, Sebastian Starke, He Zhang, Hakan Bilen, and Taku Komura, "Few-shot learning of homogeneous human lo-comotion styles," 2018, vol. 37, pp. 143–153, Wiley.
- [23] Youngwoo Yoon, Bok Cha, Joo-Haeng Lee, Minsu Jang, Jaeyeon Lee, Jaehong Kim, and Geehyuk Lee, "Speech gesture generation from the trimodal context of text, audio, and speaker identity," 2020, vol. 39, pp. 1–16, ACM New York, NY, USA.