# Co-speech Gesture Video Generation via Motion-Based Graph Retrieval

Yafei Song, Peng Zhang, Bang Zhang

Tongyi Lab, Alibaba Group

{huaizhang.syf, futian.zp, zhangbang.zb}@alibaba-inc.com

## Abstract

*Synthesizing synchronized and natural co-speech gesture videos remains a formidable challenge. Recent approaches have leveraged motion graphs to harness the potential of existing video data. To retrieve an appropriate trajectory from the graph, previous methods either utilize the distance between features extracted from the input audio and those associated with the motions in the graph or embed both the input audio and motion into a shared feature space. However, these techniques may not be optimal due to the many-to-many mapping nature between audio and gestures, which cannot be adequately addressed by one-to-one mapping. To alleviate this limitation, we propose a novel framework that initially employs a diffusion model to generate gesture motions. The diffusion model implicitly learns the joint distribution of audio and motion, enabling the generation of contextually appropriate gestures from input audio sequences. Furthermore, our method extracts both low-level and high-level features from the input audio to enrich the training process of the diffusion model. Subsequently, a meticulously designed motion-based retrieval algorithm is applied to identify the most suitable path within the graph by assessing both global and local similarities in motion. Given that not all nodes in the retrieved path are sequentially continuous, the final step involves seamlessly stitching together these segments to produce a coherent video output. Experimental results substantiate the efficacy of our proposed method, demonstrating a significant improvement over prior approaches in terms of synchronization accuracy and naturalness of generated gestures.*

## 1. Introduction

The synchronization of speech with corresponding co-speech gestures in video content is crucial for enhancing the expressiveness and effectiveness of verbal communication. Generating synchronized and natural-looking co-speech gesture videos remains a complex challenge. Recent advancements have seen various approaches aiming to address this issue, primarily falling into two categories:
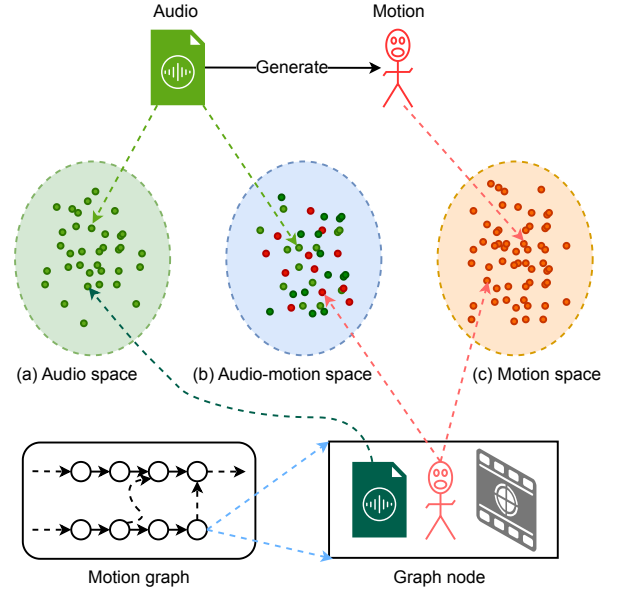


Figure 1. Our key idea is to generate the motion sequence conditioned on the input audio using a diffusion model and then retrieve its nearest trajectory from the pre-constructed motion graph.

those utilizing generative models, *e.g.* Generative Adversarial Networks (GANs) [2, 9, 12, 20, 21, 30, 31, 50] or diffusion models [8, 15, 17, 33, 41–43, 51], for gesture motion and video generation [7, 14, 19, 38, 46], and those employing motion graphs for retrieval and synthesis [24, 53].

One category of methodologies focuses on generating gesture motions adopting advanced generative models such as normalizing flow models [47], Variational AutoEncoder (VAE) [1, 22, 28, 36], GANs [13, 37, 49], diffusion models [5, 6, 32, 52, 54], and regression models [11, 25, 27, 48]. These methods are particularly adept at handling the mapping problem between spoken language and gestures. By learning intricate unidirectional mapping or joint distributions between audio and motion, these models can produce contextually appropriate gestures that are nuanced and expressive. However, while these techniques show promise, their application has often been limited to driving 3D mod-

els rather than generating real videos. Based on the generated gesture motions, a number of works further drive an image and generate corresponding real videos [7, 14, 19, 26, 38, 46]. However, due to the limited video generation capability, the results are still not satisfactory.

To obtain more natural video, another prominent approach involves leveraging motion graphs constructed from existing video data. These methods aim to retrieve and synthesize new gesture sequences by exploiting similarities between input audio features and those within the graph. While some studies [53] measure the distance in audio space via features extracted from the input audio and the audio from the graph node as shown in Fig. 1 (a). Others [24] embed both audio and motion into a shared feature space for retrieval as shown in Fig. 1 (b). Despite their effectiveness in utilizing real human motion data, these approaches struggle with the many-to-many mapping problem inherent in translating audio to gesture. The direct use of audio features or embedding audio and motion features into a common space often results in less natural transitions and mismatches due to the complexity of the mappings involved.

To address the limitations of existing methods, we propose a novel framework for Co-speech Gesture Video Generation via Motion-Based Graph Retrieval as shown in Fig. 1 (c). Our approach integrates the strengths of diffusion model-based gesture generation with sophisticated motion-based retrieval algorithms. Initially, we employ a diffusion model to generate motion sequences based on input audio, capturing the complex relationships between spoken language and body language. We enhance the training process by incorporating both low-level and high-level features extracted from the input audio, ensuring rich and accurate gesture generation.

Subsequently, we apply a meticulously designed motion-based retrieval algorithm to identify the most suitable path within the motion graph by assessing both global and local similarities in motion. This step ensures that the retrieved segments are not only contextually relevant but also seamlessly integrated. Given the potential discontinuity among nodes in the retrieved path, our final step involves stitching together these segments to produce a coherent and visually appealing video output.

Experimental results demonstrate that our proposed method significantly improves upon previous approaches, offering more natural and accurately synchronized co-speech gesture videos. By combining the advantages of diffusion models and motion graph-based retrieval, we provide a robust solution to the challenges posed by the many-to-many mapping between audio and gestures.

## 2. Related Work

Co-speech gesture generation has gained significant attention, with two primary approaches emerging: direct audio-to-gesture synthesis and motion graph-based methods that retrieve and assemble gestures from existing videos.

**Direct Generation Approaches**

Direct generation methods aim to convert spoken language directly into corresponding gesture animations without relying on pre-existing video datasets. Early efforts in this domain utilized rule-based systems, where predefined rules mapped specific words or phrases to corresponding gestures [4, 18]. While these methods were straightforward, they suffered from a lack of flexibility and naturalness, often resulting in robotic and unnatural movements.

More recently, machine learning techniques, particularly deep learning models, have been employed for direct gesture generation. For instance, VAE [1, 22, 28, 36], GANs [13, 37, 49], regression models [11, 25, 27, 48], have been used to model the temporal dependencies between speech and gestures. Although these models improved upon the naturalness and synchronization of generated gestures, they struggled with capturing the full complexity of human movement due to limitations in modeling long-range dependencies and multimodal data.

Diffusion models represent a more recent advancement in this category [5, 6, 32, 52, 54]. By implicitly constructing the joint distribution between audio and motion, diffusion models can generate highly contextually appropriate gestures. However, the effectiveness of these models heavily relies on the quality and diversity of training data, as well as the complexity of the feature extraction process. Despite these challenges, diffusion models offer significant improvements in generating nuanced and contextually accurate gestures compared to previous methods.

**Motion Graph-Based Approaches**

Motion graph-based methods utilize pre-recorded human videos to construct a graph structure, enabling retrieval and synthesis of new gesture sequences [24, 53]. Each graph node consists of a piece of audio and one or more continuous video frames. These approaches typically involve creating a motion graph by connecting similar poses between different nodes, allowing for smooth animation generation.

One common strategy involves measuring the distance between features extracted from input audio and those associated with nodes in the graph [53]. This method ensures that the retrieved motions are closely aligned with the input audio but may fail to account for the many-to-many mapping between audio and gestures, potentially resulting in less natural or contextually inappropriate gestures.

Another approach embeds both the input audio and motion into a shared feature space, facilitating the retrieval of motions that are most similar to the input [24]. While this technique can improve the contextual relevance of the retrieved motions, it also faces challenges related to the accuracy of embeddings and the inherent limitations of one-to-one mappings between audio and gestures.
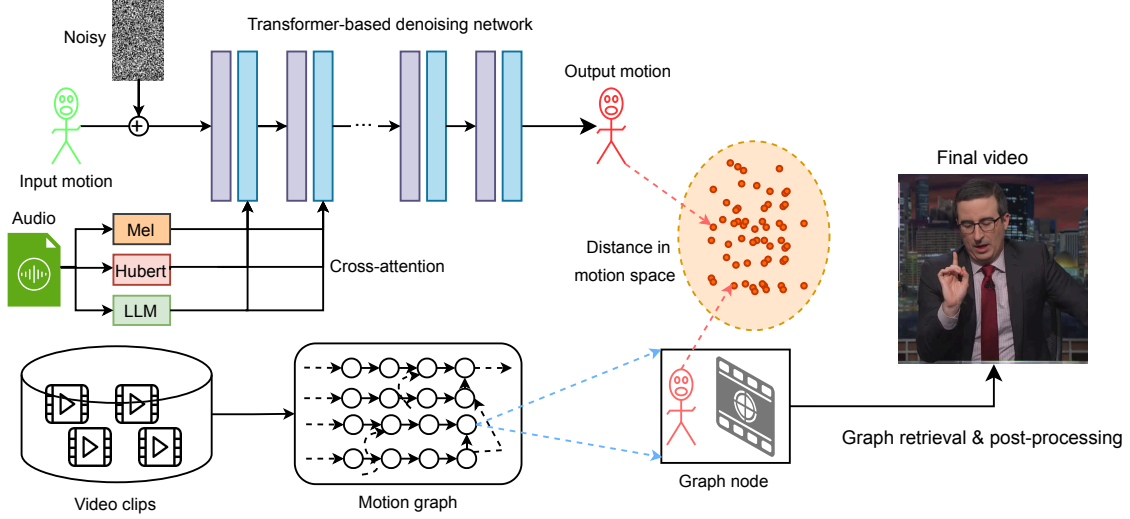
Figure 2. The pipeline of our method. We first train a transformer-based denoising network to generate motion conditioned on the audio, then construct a motion graph using the existing video data. By defining a hybrid motion similarity metric, we could retrieve the optimal trajectory from the motion graph. Combining all nodes in the trajectory, we could get the final video.

## 3. Method

In this section, we elaborate on the proposed framework for co-speech gesture video generation via motion-based graph retrieval. As shown in Fig. 2, we first train a diffusion model to learn the joint distribution of audio and gesture. Then, we construct a Motion Graph from a collection of human action videos. At inference time, we generate 3D motion sequences based on input audio using the trained diffusion model and then retrieve matching video segments from the motion graph by aligning generated 3D motions. At last, these retrieved video segments are stitched to form the final coherent video output. Our approach effectively combines the strengths of diffusion models in capturing complex audio-motion relationships with sophisticated motion-based retrieval algorithms, thereby addressing the many-to-many mapping problem between audio and gestures.

### 3.1. Generating Gestures Using Diffusion Model

Our motion generation framework employs a denoising diffusion implicit model (DDIM) [43] with a transformer architecture to learn the joint distribution of audio and motion. Given an input audio sequence $A = \{a_t\}_{t=1}^T$, we aim to generate corresponding motion parameters $X = \{x_t\}_{t=1}^T$, where $x_t = \{r\}_{j=1}^J$, $r \in SO(3)$ is the 3D rotation of a joint, and $J$ denotes the number of joints in SMPL-X [34].

**Data Preprocessing**: Motion parameters are estimated from raw videos using the SMPL-X parametric body model [3, 3, 34, 48]. For each 3-second video clip (sampled at 30 fps), we obtain motion sequence $X \in r^{90 \times J}$ and corresponding audio features.

**Conditional Diffusion Process**: Following the DDIM

framework [43], we define the forward process as a Markov chain gradually adding Gaussian noise to the motion data as

$$q(x_{1:K}|x_0) := \prod_{k=1}^K q(x_k|x_{k-1}), \quad q(x_k|x_{k-1})$$
$$:= \mathcal{N}(x_k; \sqrt{\alpha_k}x_{k-1}, (1-\alpha_k)\mathbf{I}) \quad (1)$$

where $\{\alpha_k\}_{k=1}^K$ defines a monotonically decreasing noise schedule. Our denoising transformer $\epsilon_\theta$ predicts the noise component at each step conditioned on multiple audio features

$$\epsilon_\theta(x_k, k, F) = \text{Transformer}(x_k \oplus E_k \oplus F), \quad (2)$$

where $E_k$ is the sinusoidal position encoding for diffusion step $k$, and $F$ represents the fused conditional features.

**Multi-Modal Conditioning**: We extract three complementary audio features: 1) Mel-Spectrogram features: $F_{\text{mel}} = \text{STFT}(A) \in \mathbb{R}^{T \times 128}$. 2) HuBERT embeddings [16]: $F_{\text{hubert}} = \text{HuBERT}(A) \in \mathbb{R}^{T \times 1024}$. 3) LLM (Large Language Model) semantic features: Through automatic speech recognition (ASR) [39] and QWen2-7B [45], we obtain text tokens $S = \{s_i\}_{i=1}^M$ with embeddings $F_{\text{token}} \in \mathbb{R}^{M \times 3584}$. These are temporally aligned to motion frames via nearest-neighbor interpolation:

$$F_{\text{LLM}}[t] = F_{\text{token}}[\arg\min_i |t - \frac{T}{M}i|]. \quad (3)$$

Note that the timestamps obtained by ASR are word-based, not token-based. We suppose each character in one word has the same time duration and get the timestamp of each token via combining all the characters in it.

**Algorithm 1** Motion Generation with Overlapping Diffusion

---
**Require:** Input audio $A$, trained model $\epsilon_\theta$
**Ensure:** Generated motion sequence $X$
 1: Segment $A$ into $\{A^{(i)}\}$ with 0.2s overlap
 2: **for** each clip $A^{(i)}$ **do**
 3:     **if** $i = 1$ **then**
 4:         Sample $X^{(1)} \sim p_\theta(X|A^{(1)})$
 5:     **else**
 6:         Initialize $X^{(i)}_{1:6} \leftarrow X^{(i-1)}_{85:90}$
 7:         Inpaint $X^{(i)}_{7:90}$ using $\epsilon_\theta$
 8:     **end if**
 9:     Blend $X^{(i-1)}$ and $X^{(i)}$ in overlap region
10: **end for**
11: **return** Concatenated motion $X$

---

The final conditioning vector combines these features through adaptive weights

$$F = W_{\text{mel}}F_{\text{mel}} + W_{\text{hubert}}F_{\text{hubert}} + W_{\text{LLM}}F_{\text{LLM}}, \quad (4)$$

where $\{W.\}$ are learnable projection matrices.

**Training Objective**: The model is trained via minimizing the noise prediction error

$$\mathcal{L} = \mathbb{E}_{k,x_0,\epsilon}\left[\|\epsilon - \epsilon_\theta(x_k, k, F)\|_2^2\right]. \quad (5)$$

**Inference with Temporal Consistency**: For inference on long audio inputs, we follow the in-painting algorithms with diffusion model [5, 29] as: 1) Segment input audio into 3-second clips with 0.2-second overlap, 2) Generate initial clip $X^{1:90}$ using ancestral sampling, 3) For subsequent clips $X^{(i)}$, fix the overlapping 6 frames (0.2s) and perform inpainting:

$$\hat{X}^{(i)}_{7:90} = \epsilon_\theta(X^{(i)}_{7:90}, F^{(i)}) \quad \text{s.t.} \quad X^{(i)}_{1:6} = X^{(i-1)}_{85:90}, \quad (6)$$

4) Blend overlapping regions using linear interpolation. This approach ensures temporal continuity while allowing error correction through overlapping generations. The complete algorithm is summarized in Algorithm 1.

This formulation enables the generation of long, coherent motion sequences while maintaining synchronization with the input audio through multiple complementary conditioning signals. The combination of low-level acoustic features and high-level semantic embeddings allows the model to capture both prosodic and linguistic aspects of gesture generation.

### 3.2. Constructing the Motion Graph

The motion graph serves as a structured representation of motion continuity and transition possibilities within the available video data. Our construction process consists of three key phases: node creation, continuous edge establishment, and transition edge identification.

**Node Representation**: For each video frame $i$, we create a graph node $n_i$ containing the following components:
- *Visual Data*: Raw RGB frame $I_i \in \mathbb{R}^{H \times W \times 3}$.
- *SMPL-X Parameters*: Body model parameters $\Theta_i = (\beta_i, R_i, t_i)$ where $\beta_i \in \mathbb{R}^{10}$ denotes shape parameters, $R_i \in \mathrm{SO}(3)^{J+1}$ contains joint rotations (including root orientation), $t_i \in \mathbb{R}^3$ represents root translation.
- *Joint Positions*: Camera-space coordinates $P_i \in \mathbb{R}^{J \times 3}$ computed via forward kinematics

$$P_i = \mathcal{F}_k(\Theta_i, \pi_i), \quad (7)$$

where $\mathcal{F}_k$ denotes the SMPL-X kinematic function and $\pi_i$ represents camera parameters, *i.e.* focal length $f_i \in \mathbb{R}$.
- *Joint Velocities*: Instantaneous velocity $V_i \in \mathbb{R}^{J \times 3}$ calculated through central differencing

$$V_i^j = \frac{P_{i+1}^j - P_{i-1}^j}{2\Delta t}, \quad (8)$$

where $\Delta t = 1/\text{fps}$ and $j$ indexes joints.

**Continuous Edges**: We first establish temporal continuity by connecting consecutive frames within original videos:

$$E_{\text{cont}} = \{(n_i, n_{i+1})|\forall i < T_v\}, \quad (9)$$

where $T_v$ denotes the frame count of video $v$. These edges receive a special "continuous" flag.

**Transition Edges**: To enable non-sequential transitions while preserving motion plausibility, we compute potential edges between non-consecutive nodes using a dual-threshold criterion. For any node pair $(n_s, n_t)$ where $s \neq t \pm 1$, we calculate *positional discrepancy* for joint $j$:

$$d_p(s,t,j) = \|P_s^j - P_t^j\|_2. \quad (10)$$

and *velocity discrepancy* for joint $j$:

$$d_v(s,t,j) = \|V_s^j - V_t^j\|_2. \quad (11)$$

The adaptive thresholds $\tau_p(s)$ and $\tau_v(s)$ are determined based on local motion characteristics:

$$\tau_p(s) = \lambda_p \cdot \frac{1}{2}\left(\|P_s - P_{s+1}\|_F + \|P_s - P_{s-1}\|_F\right), \quad (12)$$

$$\tau_v(s) = \lambda_v \cdot \frac{1}{2}\left(\|V_s - V_{s+1}\|_F + \|V_s - V_{s-1}\|_F\right), \quad (13)$$

where we set $\lambda_p = \lambda_v = 1.3$, and $\|\cdot\|_F$ denotes the Frobenius norm across all joints. A transition edge $(n_s, n_t)$ is established if

$$\frac{1}{J}\sum_{j=1}^{J}\mathbb{I}\left[d_p(s,t,j) \leq \tau_p(s) \wedge d_v(s,t,j) \leq \tau_v(s)\right] \geq \text{Th},$$

$$(14)$$

4

where $\mathbb{I}[\cdot]$ is the indicator function. This ensures that Th = 95% of joints maintain position and velocity continuity within local motion patterns.

To enable infinite-length video generation without dead-ends, we prune the motion graph to retain only its largest strongly connected component (SCC). With this manner, we identify nodes that form cyclic paths where every node is reachable from all others. The remaining graph preserves original continuous edges while maintaining valid transition edges within the SCC, ensuring any retrieved path can theoretically continue indefinitely by cycling through connected components. The complete motion graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ thus contains node set $\mathcal{V} = \{n_i\}_{i=1}^N$ where $N$ is the total node count and edge set $\mathcal{E} = E_{\text{cont}} \cup E_{\text{trans}}$.

This construction methodology ensures that the motion graph preserves original video continuity while enabling plausible transitions between semantically similar motion segments. The adaptive thresholding mechanism accounts for natural motion variability, and the velocity constraints maintain dynamic consistency during transitions. The SCC pruning further guarantees the graph's capacity for infinite generation by eliminating terminal nodes.

### 3.3. Retrieving Matching Video Segments

Given the generated upper-body motion sequence $X^{\text{gen}} = \{x_t^{\text{gen}}\}_{t=1}^T$ from Algorithm 1, our objective is to retrieve the optimal path $\mathcal{P} = (n_1, n_2, ..., n_L)$ through the motion graph $\mathcal{G}$ that best matches $X^{\text{gen}}$. Considering the co-speech gesture generation task, we focus exclusively on upper-body joints, excluding lower-body joints and global body rotations/translations in similarity computations. Our retrieval process consists of two key components: a hybrid motion similarity metric and a pruned tree search algorithm.

**Motion Similarity Metric**: We propose a composite distance measure combining both rotational and positional discrepancies between generated motions and graph nodes. Let $x_t^{\text{gen}}$ and $n_i$ denote a generated motion frame and a graph node respectively. We compute their distance as:

$$D(x_t^{\text{gen}}, n_i) = \lambda_r D_r(x_t^{\text{gen}}, n_i) + \lambda_p D_p(x_t^{\text{gen}}, n_i), \quad (15)$$

where $\lambda_r$ and $\lambda_p$ are balancing weights, with $D_r$ and $D_p$ defined as follows:

1) *Rotational Distance* $D_r$: For each joint $j$ in the upper body ($j \in \mathcal{J}_{\text{upper}}$), we compute the quaternion angular distance between generated rotations $r_{t,j}^{\text{gen}}$ and node rotations $r_{i,j}$:

$$D_r(x_t^{\text{gen}}, n_i) = \frac{1}{|\mathcal{J}_{\text{upper}}|} \sum_{j \in \mathcal{J}_{\text{upper}}} 2 \arccos\left(|q_{t,j}^{\text{gen}} \cdot q_{i,j}|\right),$$
$$(16)$$

where $q$ denotes rotation quaternions.

2) *Positional Distance* $D_p$: We compute the Euclidean distance between joint positions derived through forward kinematics:

$$D_p(x_t^{\text{gen}}, n_i) = \frac{1}{|\mathcal{J}_{\text{upper}}|} \sum_{j \in \mathcal{J}_{\text{upper}}} \|\mathcal{F}_k(r_t^{\text{gen}})_j - P_i^j\|_2, \quad (17)$$

where $\mathcal{F}_k(\cdot)$ denotes the forward kinematic function given rotation and $\mathcal{F}_k(\cdot)_j$ denotes the result for joint $j$, and $P_i^j$ is the 3D position of joint $j$ in node $n_i$.

**Pruned Tree Search Algorithm**: To find the optimal path through the motion graph, we employ a beam search strategy with adaptive pruning as following:

1) *State Representation*: Each search state $s = (n_c, \tau, \mathcal{C})$ consists of: - Current node $n_c \in \mathcal{V}$ - Time alignment $\tau \in \{1, ..., T\}$ indicating correspondence to $x_\tau^{\text{gen}}$ - Accumulated cost $\mathcal{C} = \sum_{t=1}^\tau D(x_t^{\text{gen}}, n_{\pi(t)})$, where $\pi(t)$ maps generated frames to path nodes.

2) *Search Initialization*: Populate the frontier with all nodes having initial cost:

$$\mathcal{C}_0(n_i) = D(x_1^{\text{gen}}, n_i). \quad (18)$$

3) *Path Expansion*: For each state $s$ in current frontier:
- Generate successor states by following all outgoing edges $(n_c, n') \in \mathcal{E}$
- Update time alignment: $\tau' = \tau + 1$
- Recode node mapping: $\pi(\tau') = n'$
- Compute incremental cost: $\Delta\mathcal{C} = D(x_{\tau'}^{\text{gen}}, n')$
- Update accumulated cost: $\mathcal{C}' = \mathcal{C} + \Delta\mathcal{C} + \beta \cdot \mathbb{I}_{\text{transition}}$, where $\mathbb{I}_{\text{transition}} = 0$ for continuous edges and $\mathbb{I}_{\text{transition}} = 1$ for transition edges, and $\beta = 0.1$ penalizes transition edges to prefer original video continuity.

4) *Pruning Strategy*: Maintain only the top-$K$ states (beam width $K = 200$) with minimal accumulated cost at each time step $\tau$. Additionally, prune states where:

$$\mathcal{C} > \mathcal{C}_{\text{min}} + \gamma \cdot (\tau/T), \quad (19)$$

where $\mathcal{C}_{\text{min}}$ is the current minimum cost and $\gamma = 1.5$ controls the pruning threshold.

The complete algorithm terminates when all states reach $\tau = T$, returning the path with minimal final cost $\mathcal{C}_T$. This beam search with adaptive pruning balances exploration of diverse motion possibilities with computational efficiency, ensuring retrieval of contextually appropriate and temporally coherent gesture sequences.

### 3.4. Stitching Retrieved Video Segments

The retrieved path $\mathcal{P}$ from the motion graph may contain transition edges that introduce discontinuities between adjacent nodes. To ensure temporal coherence and visual continuity in the final video output, we implement a two-stage stitching process combining frame interpolation and lip synchronization.

**Frame Interpolation for Smooth Transitions**: For each transition edge $(n_s, n_t)$ where $n_t$ is not the immediate successor of $n_s$ in the original video, we apply the

FILM frame interpolation model [40] to generate intermediate frames. Given the two frames before the transition $(I_{s-1}, I_s)$ and two frames after the transition $(I_t, I_{t+1})$, we synthesize two intermediate frames $\{\hat{I}_1, \hat{I}_2\}$ that bridge the motion gap. The interpolation is formulated as:

$$\hat{I}_k = \mathcal{F}_{\text{FILM}}(I_{s-1}, I_{t-1}), \quad k = 1, 2, \tag{20}$$

where $\mathcal{F}_{\text{FILM}}$ denotes the interpolation network.

**Lip Synchronization**: To enhance audio-visual consistency, we employ the pre-trained Wav2Lip model [35] to refine the mouth region of stitched frames. For each interpolated frame $I_k$ aligned with audio segment $A^{(k)}$, we generate lip masks $M_k$ and blend the synthesized mouth regions:

$$I_k^{\text{final}} = \mathcal{F}_{\text{Wav2Lip}}(I_k, A^{(k)}), \tag{21}$$

where $\odot$ denotes element-wise multiplication.

Therefore, we improve the temporal consistency to eliminate flickering artifacts, ensuring smooth transitions between stitched frames and lip synchronization.

# 4. Experiments

In this section, we introduce the dataset and experimental settings, then compare our results with several previous methods quantitatively and qualitatively. The results demonstrate that the proposed method outperforms these methods with a notable margin. Finally, an ablation study is conducted to verify the effectiveness of several key settings of the proposed method.

## 4.1. Dataset and experimental settings

Following many previous works, such as TANGO [24] and SDT [38], we also conduct experiments on the Oliver subset of the SHOW dataset [48], which contains 121 episodes of talk show recordings. The anchor in this show had diverse gestures which is suitable for co-speech gesture generation. Each episode comprises multiple video clips with synchronized speech and upper-body gestures, totaling 28.7 hours of footage. The dataset provides SMPL-X [34] parameters estimated through elaborate optimization, with joint rotations, global translation and orientation, shape parameters, and camera parameters.

We adopt a stratified splitting strategy to ensure content diversity:
- Training set: 97 episodes (5,521 clips), which is used to train the diffusion model.
- Validation set: 12 episodes (757 clips), which is used to validate the training and tune the hyper-parameters.
- Test set: 12 episodes (833 clips), which is used to test the generation results.

For the test episodes, we further partition each episode into:

- Motion graph construction: 80% of clips (606 clips), which is used to construct the motion graph. Because the host wears different outfits in each episode, we construct a motion graph for each episode's data.
- Test queries: 20% of clips (137 3-10 seconds clips (after removing shorter clips than 3 seconds), which is used for final generation and comparisons. The audio is used as input fore each method, and the video is taken as the ground truth.

Our diffusion model uses a 8-layer transformer with 512 hidden dimensions, trained for 3000 epochs using AdamW optimizer (learning rate 2e-4 and reduced on plateau till 1e-5, batch size 512). We extract audio features using:
- Mel-spectrograms: 128 bins, 33ms window, 33ms hop.
- HuBERT [16]: Large model fine-tuned on LibriSpeech.
- QWen2-7B [45]: The features output from the last layer. The model does not output this feature by default, we modify the model to get it.

Motion graph construction employs adaptive thresholds $\lambda_p = 1.3$, $\lambda_v = 1.3$ with 95% joint consensus. The retrieval beam search uses $K = 200$, $\gamma = 1.5$, and $\beta = 0.1$. Training is conducted on 4×A100 GPUs with PyTorch 2.0 and testing on one A100.

## 4.2. Quantitative results and comparisons

For each test query, we use the trained diffusion model to generate the corresponding motion sequence and take the sequence to retrieve and synthesise the final video based on its corresponding motion graph.

For comparison, we evaluate against three state-of-the-art methods: graph-based method TANGO [24], diffusion-based method MDDiffusion [14], and GAN-based method SDT [38]. For TANGO [24], we use its Audio-Motion CLIP to retrieve from the motion graph that is identical to ours for a fairer comparison. For MDDiffusion [14] and SDT [38], these two methods are designed to drive an image, therefore, we select a clear frame from the corresponding motion graph as the source image.

To quantitatively evaluate the quality of generated videos, we use the content debiased fréchet video distance (CD-FVD) [10] instead of FVD [44], as FVD is not sensitive to temporal quality. To quantitatively evaluate the generated motion, we use SMPLer-X [3] to reconstuct the motion from the generated videos and adopt fréchet motion distance (FMD) [49] to measure the motion result. The auto-encoder model used by FMD to compute latent features is trained on the whole Oliver dataset. To measure the beat aligned with the audio, we adopt the beat consistency (BC) [23]. We also report the motion diversity that is described in [22].

As presented in Tab. 1, our method achieves the best performance on CD-FVD, FMD, and BC. Compared to the second method, the improvements on CD-FVD (170.1

| Method | CD-FVD↓ | FMD↓ | Diversity↑ | BC↑ |
|---|---|---|---|---|
| SDT [38] | 445.1 | 38.23 | 1.150 | 0.823 |
| MDDiffusion [14] | 563.6 | 40.00 | 1.044 | 0.788 |
| TANGO [24] | 228.3 | 29.48 | 1.370 | 0.924 |
| Ours | 170.1 | 21.75 | 1.177 | 0.931 |
| GT | 0 | 0 | 1.340 | 0.941 |

Table 1. Our method achieves the best performance on CD-FVD, FMD, and BC. TANGO exhibits higher diversity, which may stem from its direct reuse of existing motions.

vs 228.3) and FMD (21.75 vs 29.48) are both remarkable, demonstrating enhanced visual quality and motion fidelity. Our method also achieves a slightly better BC as the diffusion-based motion prior enables better alignment with speech context compared to pure retrieval approaches. While TANGO exhibits higher diversity (1.37 vs 1.177), this may stem from its direct reuse of existing motions, which may preserve dataset biases. Our method maintains competitive beat consistency (0.931 vs 0.924) while generating novel motions. Without a doubt, compared with single-image-based methods [14, 38], our method could achieve obviously better results on all aspects.

### 4.3. Ablation studies

*Impact of LLM Features*: As shown in Tab. 2, removing LLM semantic features degrades CD-FVD by 26.5% and FMD by 31.5%, validating that high-level semantic cues help generate contextually appropriate gestures. The marginal BC improvement (0.935 vs 0.931) suggests LLM features primarily enhance macro-level motion semantics rather than micro-rhythmic patterns. This result verifies that co-speech gestures not only align to beat but also have a high correlation to the semantical information embedded in the speech.

*Retrieval Strategy*: Replacing our hybrid similarity without rotation distance degrades BC by 9.3%, while without positional distance, it increases FMD by 14.2%. This experiment indicates that the combined metric optimally balances local rotation and global spatial constraints, leading to overall better performance.

| Method | CD-FVD↓ | FMD↓ | Div.↑ | BC↑ |
|---|---|---|---|---|
| Ours w/o LLM | 215.1 | 28.60 | 1.124 | 0.935 |
| Ours w/o $D_r$ | 174.4 | 22.37 | 1.182 | 0.844 |
| Ours w/o $D_p$ | 178.2 | 24.84 | 1.143 | 0.893 |
| Ours | 170.1 | 21.75 | 1.177 | 0.931 |

Table 2. Ablation experiments verify the effectiveness of LLM features, the rotational distance $D_r$ and the position distance $D_p$.

### 4.4. Qualitative analysis

We present some qualitative result comparisons in Fig. 3 and the videos in the supplement. From the videos, our method produces smoother and higher-fidelity co-speech gesture videos than the compared methods. While TANGO [24] suffers from jitter frames, which may be due to its audio-motion CLIP retrieves more non-continuous segments. MDDiffusion [14] generates blurry and jittery frames, which indicates that current diffusion models have limited ability to generate quick-moving content. SDT [38] suffers from obvious warping artifacts, which are usually observed in GAN-based methods.

## 5. Conclusion

This paper presents a novel framework for co-speech gesture video generation that integrates diffusion-based motion synthesis with motion graph retrieval. Our key innovation lies in decoupling the problem into two stages: (1) a diffusion model conditioned on multi-modal audio features (Mel-spectrograms, HuBERT embeddings, and LLM-derived semantics) to generate plausible gesture sequences, and (2) a hybrid motion similarity metric for retrieving and stitching video segments from a pre-constructed motion graph. This approach effectively addresses the many-to-many mapping challenge between speech and gestures while leveraging real motion data for natural video synthesis. Extensive experiments demonstrate state-of-the-art performance in synchronization accuracy and visual quality, outperforming existing methods across multiple metrics, including CD-FVD, FVD, and BeatConsistency.

The main limitations stem from dependency on motion graph construction from video data and occasional transition artifacts. Future work will explore few-shot graph adaptation and enhanced transition modeling. Our framework advances co-speech gesture generation by combining the strengths of learned motion priors and structured motion retrieval, offering practical value for virtual agent animation and human-computer interaction systems.

Transcript: just a couple of things there. first, Britain was already independent.

GT

Our

TANGO

MDDiff

SDT

Transcript: even though, somewhere in a dimly lit room, Paul Krugman worked very hard to make it. Online! Online!

GT

Our

TANGO

MDDiff

SDT

Figure 3. Qualitative comparison showing our method's ability to generate context-specific gestures (*e.g.*, raised arms during emphatic speech). As images hardly represent temporal information, please refer to the videos in the supplement for better observation.

# References

[1] Tenglong Ao, Qingzhe Gao, Yuke Lou, Baoquan Chen, and Libin Liu. Rhythmic gesticulator: Rhythm-aware co-speech gesture synthesis with hierarchical neural embeddings. *ACM TOG*, 41(6):1–19, 2022. 1, 2

[2] Guha Balakrishnan, Amy Zhao, Adrian V Dalca, Fredo Durand, and John Guttag. Synthesizing images of humans in unseen poses. In *CVPR*, pages 8340–8348, 2018. 1

[3] Zhongang Cai, Wanqi Yin, Ailing Zeng, Chen Wei, Qingping Sun, Wang Yanjun, Hui En Pang, Haiyi Mei, Mingyuan Zhang, Lei Zhang, Chen Change Loy, Lei Yang, and Ziwei Liu. SMPLer-X: Scaling up expressive human pose and shape estimation. In *NeurIPS*, 2023. 3, 6

[4] Justine Cassell, Hannes Högni Vilhjálmsson, and Timothy Bickmore. Beat: the behavior expression animation toolkit. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, pages 477–486, 2001. 2

[5] Junming Chen, Yunfei Liu, Jianan Wang, Ailing Zeng, Yu Li, and Qifeng Chen. Diffsheg: A diffusion-based approach for real-time speech-driven holistic 3d expression and gesture generation. In *CVPR*, pages 7352–7361, 2024. 1, 2, 4

[6] Qingrong Cheng, Xu Li, and Xinghui Fu. Siggesture: Generalized co-speech gesture synthesis via semantic injection with large-scale pre-training diffusion models. In *SIGGRAPH Asia 2024 Conference Papers*, pages 1–11, 2024. 1, 2

[7] Enric Corona, Andrei Zanfir, Eduard Gabriel Bazavan, Nikos Kolotouros, Thiemo Alldieck, and Cristian Sminchisescu. Vlogger: Multimodal diffusion for embodied avatar synthesis. *arXiv preprint arXiv:2403.08764*, 2024. 1, 2

[8] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 1

[9] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *CVPR*, pages 12873–12883, 2021. 1

[10] Songwei Ge, Aniruddha Mahapatra, Gaurav Parmar, Jun-Yan Zhu, and Jia-Bin Huang. On the content bias in fréchet video distance. In *CVPR*, pages 7277–7288, 2024. 6

[11] Shiry Ginosar, Amir Bar, Gefen Kohavi, Caroline Chan, Andrew Owens, and Jitendra Malik. Learning individual styles of conversational gesture. In *CVPR*, pages 3497–3506, 2019. 1, 2

[12] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *NeurIPS*, 27, 2014. 1

[13] Ikhsanul Habibie, Weipeng Xu, Dushyant Mehta, Lingjie Liu, Hans-Peter Seidel, Gerard Pons-Moll, Mohamed Elgharib, and Christian Theobalt. Learning speech-driven 3d conversational gestures from video. In *Proceedings of the 21st ACM International Conference on Intelligent Virtual Agents*, pages 101–108, 2021. 1, 2

[14] Xu He, Qiaochu Huang, Zhensong Zhang, Zhiwei Lin, Zhiyong Wu, Sicheng Yang, Minglei Li, Zhiyi Chen, Songcen Xu, and Xiaofei Wu. Co-speech gesture video generation via motion-decoupled diffusion model. In *CVPR*, pages 2263–2273, 2024. 1, 2, 6, 7

[15] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NeurIPS*, 33:6840–6851, 2020. 1

[16] Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. *IEEE/ACM Trans. Audio, Speech and Lang. Proc.*, 29:3451–3460, 2021. 3, 6

[17] Li Hu. Animate anyone: Consistent and controllable image-to-video synthesis for character animation. In *CVPR*, pages 8153–8163, 2024. 1

[18] Chien-Ming Huang and Bilge Mutlu. Robot behavior toolkit: generating effective social behaviors for robots. In *Proceedings of the seventh annual ACM/IEEE international conference on Human-Robot Interaction*, pages 25–32, 2012. 2

[19] Ziyao Huang, Fan Tang, Yong Zhang, Xiaodong Cun, Juan Cao, Jintao Li, and Tong-Yee Lee. Make-your-anchor: A diffusion-based 2d avatar generation framework. In *CVPR*, pages 6997–7006, 2024. 1, 2

[20] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, pages 4401–4410, 2019. 1

[21] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Training generative adversarial networks with limited data. *NeurIPS*, 33:12104–12114, 2020. 1

[22] Jing Li, Di Kang, Wenjie Pei, Xuefei Zhe, Ying Zhang, Zhenyu He, and Linchao Bao. Audio2gestures: Generating diverse gestures from speech audio with conditional variational autoencoders. In *ICCV*, pages 11293–11302, 2021. 1, 2, 6

[23] Ruilong Li, Shan Yang, David A Ross, and Angjoo Kanazawa. Ai choreographer: Music conditioned 3d dance generation with aist++. In *ICCV*, pages 13401–13412, 2021. 6

[24] Haiyang Liu, Xingchao Yang, Tomoya Akiyama, Yuantian Huang, Qiaoge Li, Shigeru Kuriyama, and Takafumi Taketomi. Tango: Co-speech gesture video reenactment with hierarchical audio motion embedding and diffusion interpolation. *arXiv preprint arXiv:2410.04221*, 2024. 1, 2, 6, 7

[25] Haiyang Liu, Zihao Zhu, Giorgio Becherini, Yichen Peng, Mingyang Su, You Zhou, Xuefei Zhe, Naoya Iwamoto, Bo Zheng, and Michael J Black. Emage: Towards unified holistic co-speech gesture generation via expressive masked audio gesture modeling. In *CVPR*, pages 1144–1154, 2024. 1, 2

[26] Xian Liu, Qianyi Wu, Hang Zhou, Yuanqi Du, Wayne Wu, Dahua Lin, and Ziwei Liu. Audio-driven co-speech gesture video generation. *NeurIPS*, 35:21386–21399, 2022. 2

[27] Xian Liu, Qianyi Wu, Hang Zhou, Yinghao Xu, Rui Qian, Xinyi Lin, Xiaowei Zhou, Wayne Wu, Bo Dai, and Bolei Zhou. Learning hierarchical cross-modal association for co-speech gesture generation. In *CVPR*, pages 10462–10472, 2022. 1, 2

[28] Yifei Liu, Qiong Cao, Yandong Wen, Huaiguang Jiang, and Changxing Ding. Towards variable and coordinated holistic

co-speech motion generation. In *CVPR*, pages 1566–1576, 2024. 1, 2

[29] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *CVPR*, pages 11461–11471, 2022. 4

[30] Liqian Ma, Xu Jia, Qianru Sun, Bernt Schiele, Tinne Tuytelaars, and Luc Van Gool. Pose guided person image generation. *NeurIPS*, 30, 2017. 1

[31] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014. 1

[32] Muhammad Hamza Mughal, Rishabh Dabral, Ikhsanul Habibie, Lucia Donatelli, Marc Habermann, and Christian Theobalt. Convofusion: Multi-modal conversational diffusion for co-speech gesture synthesis. In *CVPR*, pages 1388–1398, 2024. 1, 2

[33] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International conference on machine learning*, pages 8162–8171. PMLR, 2021. 1

[34] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed AA Osman, Dimitrios Tzionas, and Michael J Black. Expressive body capture: 3d hands, face, and body from a single image. In *CVPR*, pages 10975–10985, 2019. 3, 6

[35] K R Prajwal, Rudrabha Mukhopadhyay, Vinay P. Namboodiri, and C.V. Jawahar. A lip sync expert is all you need for speech to lip generation in the wild. In *ACM MM*, page 484–492, 2020. 6

[36] Xingqun Qi, Chen Liu, Lincheng Li, Jie Hou, Haoran Xin, and Xin Yu. Emotiongesture: Audio-driven diverse emotional co-speech 3d gesture generation. *IEEE TMM*, 2024. 1, 2

[37] Xingqun Qi, Jiahao Pan, Peng Li, Ruibin Yuan, Xiaowei Chi, Mengfei Li, Wenhan Luo, Wei Xue, Shanghang Zhang, Qifeng Liu, et al. Weakly-supervised emotion transition learning for diverse 3d co-speech gesture generation. In *CVPR*, pages 10424–10434, 2024. 1, 2

[38] Shenhan Qian, Zhi Tu, Yihao Zhi, Wen Liu, and Shenghua Gao. Speech drives templates: Co-speech gesture synthesis with learned templates. In *ICCV*, pages 11077–11086, 2021. 1, 2, 6, 7

[39] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *International conference on machine learning*, pages 28492–28518. PMLR, 2023. 3

[40] Fitsum Reda, Janne Kontkanen, Eric Tabellion, Deqing Sun, Caroline Pantofaru, and Brian Curless. Film: Frame interpolation for large motion. In *ECCV*, 2022. 6

[41] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, pages 10684–10695, 2022. 1

[42] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. pmlr, 2015.

[43] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2021. 1, 3

[44] Thomas Unterthiner, Sjoerd Van Steenkiste, Karol Kurach, Raphaël Marinier, Marcin Michalski, and Sylvain Gelly. Fvd: A new metric for video generation. In *ICLR Workshop DeepGenStruct*, 2019. 6

[45] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024. 3, 6

[46] Quanwei Yang, Jiazhi Guan, Kaisiyuan Wang, Lingyun Yu, Wenqing Chu, Hang Zhou, ZhiQiang Feng, Haocheng Feng, Errui Ding, Jingdong Wang, et al. Showmaker: Creating high-fidelity 2d human video via fine-grained diffusion modeling. In *NeurIPS*, pages 51039–51062, 2024. 1, 2

[47] 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 *ECCV*, pages 712–728. Springer, 2022. 1

[48] Hongwei Yi, Hualin Liang, Yifei Liu, Qiong Cao, Yandong Wen, Timo Bolkart, Dacheng Tao, and Michael J Black. Generating holistic 3d human motion from speech. In *CVPR*, pages 469–480, 2023. 1, 2, 3, 6

[49] 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. *ACM TOG*, 39(6):1–16, 2020. 1, 2, 6

[50] Lantao Yu, Weinan Zhang, Jun Wang, and Yong Yu. Seqgan: Sequence generative adversarial nets with policy gradient. In *AAAI*, 2017. 1

[51] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *ICCV*, pages 3836–3847, 2023. 1

[52] Mingyuan Zhang, Zhongang Cai, Liang Pan, Fangzhou Hong, Xinying Guo, Lei Yang, and Ziwei Liu. Motiondiffuse: Text-driven human motion generation with diffusion model. *IEEE TPAMI*, 46(6):4115–4128, 2024. 1, 2

[53] Yang Zhou, Jimei Yang, Dingzeyu Li, Jun Saito, Deepali Aneja, and Evangelos Kalogerakis. Audio-driven neural gesture reenactment with video motion graphs. In *CVPR*, pages 3418–3428, 2022. 1, 2

[54] Lingting Zhu, Xian Liu, Xuanyu Liu, Rui Qian, Ziwei Liu, and Lequan Yu. Taming diffusion models for audio-driven co-speech gesture generation. In *CVPR*, pages 10544–10553, 2023. 1, 2