REALISTIC-GESTURE: CO-SPEECH GESTURE VIDEO GENERATION THROUGH CONTEXT-AWARE GESTURE REPRESENTATION

Anonymous authors

000

001

002

004

006

008 009

024

025

026 027 028

029

031

032

033

034

038

039

040

041

042

Paper under double-blind review



Long Sequence Generation

Video Gesture Editing

Gesture Pattern Transfer

Figure 1: Realistic-Gesture achieves various fine-grained control over video-level gesture motion.

Abstract

Co-speech gesture generation is crucial for creating lifelike avatars and enhancing human-computer interactions by synchronizing gestures with speech in computer vision. Despite recent advancements, existing methods often struggle with accurately aligning gesture motions with speech signals and achieving pixellevel realism. To address these challenges, we introduce Realistic-Gesture, a groundbreaking framework that transforms co-speech gesture video generation through three innovative components: (1) a speech-aware gesture tokenization that incorporate speech context into motion pattern representation, (2) a mask gesture generator that learns to map audio signals to gestures by predicting masked motion tokens, enabling bidirectional contextually relevant gesture synthesis and editing, and (3) a structure-aware refinement module that employs differentiable edge connection to link gesture keypoints to improve video generation. Our extensive experiments demonstrate that Realistic-Gesture not only produces highly realistic and speech-aligned gesture videos but also supports long-sequence generation and video gesture editing applications.

- 043 044 045 046
- 1 INTRODUCTION

In human communication, speech is often accompanied by gestures that enhance understanding and convey emotions De Ruiter et al. (2012). As these non-verbal cues play a vital role in effective interaction Burgoon et al. (1990), gesture generation a key component of natural human-computer interactions. As artificial intelligence advances, equipping virtual avatars with realistic gesture capabilities will become essential in creating immersive interactive experiences.

053 The relationships between the semantic and emotion content of speech context, the corresponding gestures, and the visual appearance of the speaker's performance are complex. As such, many

recent works Yi et al. (2023); Liu et al. (2023; 2022d;a) address a reduced form of this problem by generating a simplified representation of the 3D motion, consisting of joints and body parts, that plausibly accompanies a given speech sample, which can then be rendered using standard rendering pipelines. Such representations capture basic motion patterns, yet they neglect the importance of the speaker's visual appearance, resulting in a lack of realism that hinders effective communication.

Other works, e.g. ANGIE Liu et al. (2022c) and S2G-diffusion He et al. (2024), employ image-060 warping techniques, constrained by keypoints obtained from optical-flow-based deformations, for 061 co-speech video generation. However, such approaches encounter several critical issues. First, these 062 keypoints only define large-scale transformations, and thus miss subtle movements of specific body 063 parts (e.g. hands, fingers). Second, it is difficult to connect such broad and unconstrained motions 064 representation to speech content. This makes it difficult to conditionally generate gestures that are appropriately responsive to the audio, which inhibits the gestures' naturalism and expressivity. 065 Finally, the generated motion patterns are often unstructured and overly reactive to large motion, 066 resulting in noisy and imprecise renderings, especially in the hands and shoulders. Collectively, 067 these challenges significantly limit the overall quality and realism of the generated video content. 068

069 To address these challenges, we introduce *Realistic-Gesture*, a framework designed to generate speech-aligned gesture motions and high-fidelity speech video outputs. Our approach begins with 071 refined gesture motion representations using keypoints from pretrained human pose estimators, allowing for clearer disentanglement of human motions across the face, body, and hands. To uncover 072 the intrinsic temporal connections between gestures and speech, we employ contrastive learning 073 to align these two modalities. This joint representation captures the triggers of gesture patterns 074 influenced by speech. We incorporate speech-contextual features into the tokenization process 075 of gesture motions through knowledge distillation, aiming to infuse the gesture representations 076 with implicit intentions conveyed in the audio. This integration creates a clear linkage between 077 the gestures and the corresponding speech, enabling the conditional generation of gestures that accurately reflect the speaker's intended meaning based on the speech input. For latent motion 079 generation, inspired by Muse Chang et al. (2023) and MAGE Li et al. (2023), we introduce a masked gesture generator that refines the alignment of gesture motions with the speech signal 081 through bidirectional mask pretraining, enabling long sequence generation and editing capabilities. Finally, for uplifting the latent motion generation into 2D animations, we propose a structure-aware 082 image refinement module that generates heatmaps of edge connections from keypoints, providing 083 image-level supervision to improve the quality of body regions with large motion. Extensive 084 experiments demonstrate that our method outperforms the existing state-of-the-art approaches in 085 both quantitative and qualitative metrics.

In summary, our primary contributions are:

087

880

089

090 091

092 093

094

095 096

- 1. a *speech-aware gesture motion representation* obtained through knowledge distillation from the gesture-speech aligned features from contrastive learning;
- 2. a *masked gesture motion generator*, carefully designed to enable high-quality gesture motion generation with long sequence generation and edit-ability support; and
 - 3. a *pixel-level refinement module*, which uses a structure-aware edge heatmap as supervision to improve the final output fidelity.

2 RELATED WORK

097 **Co-speech Gesture generation** Most recent works on co-speech gesture generation employ 098 skeleton- or joint-level pose representations. Ginosar et al. (2019) use an adversarial framework to predict hand and arm poses from audio, and leverage conditional generation Chan et al. (2019) 100 based on pix2pixHD Wang et al. (2018) for videos. Some recent works Liu et al. (2022d); Deichler 101 et al. (2023); Xu et al. (2023) learns the hierarchical semantics or leverage contrastive learning to 102 obtain joint audio-gesture embeding to assist the gesture pose generation. Rhythmic gesticulator Ao 103 et al. (2022) construct high and low level audio-motion embedding based on lingustic theory for 104 gesture generation. TalkShow Yi et al. (2023) estimates SMPL Pavlakos et al. (2019) poses, and 105 models the body and hand motions for talk-show videos. CaMN Liu et al. (2022b) and EMAGE Liu et al. (2023) use large conversational and speech datasets for joint face and body modeling with 106 diverse style control. ANGIE Liu et al. (2022c) uses unsupervised 2D keypoints with image-107 warping features based on MRAA Siarohin et al. (2021) to model body motion. It leverages Vector Quantization van den Oord et al. (2018) to obtain common patterns, followed by a GPT-like network that outputs co-speech gesture videos. S2G-Diffusion He et al. (2024) uses TPS Zhao & Zhang (2022) and optical flow prediction to extract and refine latent motion features from videos. However, none of these works produce structure- and speech-aware motion patterns that are suitable for achieving natural and realistic gesture rendering.

113

114 Conditional Video Generation Conditional Video Generation has undergo significant progress 115 for various modalities, like text Blattmann et al. (2023), pose Karras et al. (2023); Wang et al. 116 (2023b), and audio Ruan et al. (2023). Diffusion Models Ho et al. (2020) improve generation 117 qualities. AnimateDiff Guo et al. (2024) presents an efficient motion adaptation module based on low-rank adaptation Hu et al. (2022) (LoRA) to adapt image diffusion model for video motion 118 generation. AnimateAnyone Hu et al. (2023) construct referencenet for fine-grained control based 119 on skeleton. Make-Your-Anchor Huang et al. (2024) improves avatar video generation through 120 disentangled face and body based on SMPL-X conditions. Champ Zhu et al. (2024) introduces 121 human SMPL models for guidance. EMO Tian et al. (2024) leverages audio as control signal for 122 talking head generation. However, these methods are based on large amount of training data and 123 slow in inference speed. None of them focus on the speech-gesture pixel-level video generation.

124 125

Masked Representation Learning for Generation Masked Representation Learning has been 126 demonstated an effective representation learning for various modalities. Devlin (2018); He et al. 127 (2022) Some works explored the generation capabilities using this paradigm. MAGE Li et al. (2023) 128 achieves high-quality image generation through iterative remasking. Muse Chang et al. (2023) 129 extends this idea to leverage language with region masking for image editing and achieve fine-130 grained control. Recent Masking Models Pinyoanuntapong et al. (2024); Wang (2023); Mao et al. 131 (2024) bring this strategy to the motion and gesture domain and improves the motion generation 132 speed, quality, and editing capability. Inspired by these work, we propose the masked gesture 133 generation conditioned the audio to learn the gesture-speech correspondence during generation.

134 135

3 PRELIMINARY

136 137

Warping-Based Image Animation. Warping-based image animation methods have risen to
 prominence recently Siarohin et al. (2021; 2019); Zhao & Zhang (2022). They leverage keypoint
 predictor to identify pairwise corresponding keypoints between a source image and a driving image.
 This information is used to warp the source image to match the driving image, thereby producing a
 deformation that aligns with the driving scene. Following this, pixel-level optical flow and occlusion
 masks are estimated from the deformed images to capture global motion and handle occlusions for
 achieving driving image reconstruction. We defer additional details to the Appendix.

145

146 Image-Animation Based Co-Speech Gesture Video Synthesis. In the context of co-speech 147 gesture video synthesis, recent advancements have employed warping-based image animation 148 techniques to derive motion patterns and learn the correspondence between these patterns and audio, facilitating speech-driven generation. Given a video clip $V = \{I_0, I_1, \ldots, I_N\}$ and 149 an accompanying audio sequence $A = \{a_1, a_2, \ldots, a_N\}$, the objective is to predict motion 150 151 representations M based on the initial frame I_0 and the audio input. The image animation module 152 reconstructs all video frames I_1 through unsupervised learning to derive motion representations and transformations. The audio sequence serves as guidance for reconstructing motion patterns across 153 the entire sequence of frames following the initial frame. 154

However, this approach faces three significant challenges: (1) keypoints derived from global optical-flow-based transformations, learned through unsupervised methods, often fail to capture subtle movements of specific body parts; (2) the motion representations do not include contextual information from the speech, making it difficult to generate gestures that are conditionally responsive to audio; and (3) the lack of structural awareness in the motion representations leads to blurry and noisy predictions, particularly affecting the hands and shoulders, while also rendering the system sensitive to large motion patterns. To address these challenges, we propose the following methods to enhance control over co-speech gesture video generation.

¹⁶² 4 REALISTIC-GESTURE

163 164

As shown in Fig. 2, our framework targets at gen-165 erating realistic gesture videos. To achieve this 166 goal, we first learn gesture-speech alignment to build 167 speech-aware gesture motion representation through 168 contrastive learning. (Sec. 4.1) To achieve finegrained control over the gesture motion generation, 169 170 we propose a Masking-based Gesture Generator, with long sequence and editing capabilities. (Sec. 4.2). 171 To improve the noisy hand and shoulder movement 172 during the uplifting of latent motion to pixel space, we 173 propose a structure-aware image refinement through 174 differentiable edge heatmaps for guidance. (Sec. 4.3).



Figure 2: An overview of our framework.

175 176

177

189

200 201 202

4.1 SPEECH-AWARE GESTURE MOTION REPRESENTATION

178 Unlike ANGIE Liu et al. (2022c) and S2G-Diffusion He et al. (2024), which rely on unsupervised 179 learning of keypoints for warping-based on optical flow, we utilize 2D poses extracted from 180 images. While using 2D poses for image warping may slightly decrease fidelity, it significantly 181 enhances the perceptual quality of the generated video, shown in Sec. 5.4. With poses, gestures 182 can be decomposed into facial and body movements. We represent a gesture motion sequence as 183 $G = [F;B] = [f_t; b_t]_{t=1}^T$, where T denotes the length of the motion, f represents the 2D facial landmarks, and b denotes the 2D body landmarks. Further details on gesture representations can be 184 found in the appendix. For speech representation, we extract audio embeddings from WavLM Chen 185 et al. (2022). In addition, we extract Mel spectrogram features Rabiner & Schafer (2010) and beat information using librosa McFee et al. (2015). These features are concatenated to form the speech 187 representation. For image-warping transformation, we select TPS Zhao & Zhang (2022). 188

Speech-Gesture Alignment. To align gesture motion patterns with the content of speech and 190 beats, we draw inspiration from image-language contrastive learning Radford et al. (2021). We first 191 project both speech and gesture modalities into a shared embedding space to enhance the speech 192 content awareness of gesture features. As illustrated in Fig. 3, we separately train two gesture 193 content encoders, \mathcal{E}_f for face motion and \mathcal{E}_b for body motion, alongside two speech encoders, \mathcal{E}_{S_f} 194 and \mathcal{E}_{S_b} , to map face and body movements and speech signals into this joint embedding space. For 195 simplicity, we represent the general gesture motion sequence as G. We then apply mean pooling to 196 aggregate content-relevant information from each feature sequence, resulting in the embeddings z^s 197 and z^{g} for speech and gestures, respectively. We leverage CLIP-style contrastive learning to train these content encoders. Given a batch of paired embeddings $\mathcal{B} = \{(z_i^t, z_j^g)\}_{i=1}^{B}$, we optimize the 198 following loss with τ as the temperature: 199

$$\mathcal{L}_{\text{contrast}} = -\frac{1}{B} \sum_{i=1}^{B} \left(\log \frac{\exp(z_i^g \cdot z_i^t/\tau)}{\sum_{j=1}^{B} \exp(z_i^g \cdot z_j^t/\tau)} + \log \frac{\exp(z_i^t \cdot z_i^g/\tau)}{\sum_{j=1}^{B} \exp(z_i^t \cdot z_j^g/\tau)} \right)$$
(1)

Unlike previous methods Ao et al. (2022); Liu et al. (2022d); Deichler et al. (2023), which primarily
capture sequence-level alignment and may overlook local temporal dynamics, to mitigate this
limitation, we randomly mask 30% of segments from both speech and gesture sequences within
the same temporal regions during training. Furthermore, we apply a linear classifier on the gesture
embedding to predict speech beats, enhancing the temporal alignment between gestures and speech.
We defer additional details of temporal-level improvement by our strategy in the Appendix.

Speech-Pattern Learning Through Knowledge Distillation. For gesture motion tokenization, we utilize Residual Vector Quantization Lee et al. (2022) (RVQ) to capture the high diversity and complexity of facial and body motions. To construct context-aware motion representations, we directly encode alignment information into the gesture motion codebook. This allows the semantics and contextual triggers from speech (e.g., pronouns like "this" or "they") to be fused into the motion embedding, and enables the generator to easily identify the corresponding motion representation in response to speech triggers. To achieve this goal, we leverage gesture content encoder as the teacher and distill knowledge to codebook latent representation. We aim to maximize the cosine similarity



Figure 3: Left: Contrastive Learning for gesture-speech alignment. We distill the joint speech contextualaware feature into latent codebook. Right: We use speech for generating gesture motion tokens with Mask Gesture Generator. We apply random mask for reconstruction during training and iterative remask based on probability for inference. Residual Gesture Generator based on the base VQ-tokens to predict the residuals.

over time between the final RVQ quantization output and the representation from the gesture content encoder, formulated as follows:

$$\mathcal{L}_{\text{distill}} = \sum_{t=1}^{T} \cos\left(p(Q_R)^t, Es(G)^t\right)$$
(2)

where p denotes a linear projection layer, Q_R is the final quantization output from the RVQ-VAE, Es(G) represents the output from the gesture content encoder, and T is the total time frames. The overall training objective for the RVQ-VAE is defined as:

$$\mathcal{L}_{\text{rvq}} = \mathbb{E}_{x \sim p(x)} \left[\left\| x - \hat{x} \right\|^2 \right] + \alpha \sum_{r=1}^{R} \mathbb{E}_{z_r \sim q(z_r|x)} \left[\left\| e_r - \text{sg} \left(z_r - e_r \right) \right\|^2 \right] + \beta \mathcal{L}_{\text{distill}}$$
(3)

where \mathcal{L}_{rvq} combines a motion reconstruction loss, a commitment loss van den Oord et al. (2018) for each layer of quantizer with a distillation loss, with α and β weighting the contributions.

4.2 Speech-conditioned Gesture Motion Generation

To enhance the generation of gesture motions across different layers of the quantized codebooks,
 we draw inspiration from VALL-E Wang et al. (2023a) to design Masked Gesture Generator for
 jointly decoding facial and body motions for the base-layer outputs of quantizers, and Residual
 Gesture Generator for the face and body tokens from the subsequent *R* residual quantization layers.

Masked Gesture Generator. As shown in Fig. 3, during training, we derive motion tokens 255 by processing raw gesture sequences through both body and face tokenizers. The motion token 256 corresponding to the source image acts as the conditioning for all subsequent frames. For speech 257 control, we initialize the audio content encoder from alignment pre-training as described in Sec. 4.1. 258 This pre-alignment of gesture tokens with audio encoder features enhances the coherence of gesture 259 generation. We employ cross-attention, using the audio input as keys and values while the gesture 260 representation serves as the query, integrating audio information with gesture feature. To refine 261 control over gesture patterns, we apply Adaptive Instance Normalization (AdaIN) Huang & Belongie 262 (2017) after the feed-forward layers, enabling diverse gesture styles based on the speaker's identity.

Residual Gesture Generator. The Residual Gesture Generator shares a similar architecture with the Masked Gesture Generator, but it includes R separate embedding layers corresponding to each RVQ residual layer. During training, we randomly select a quantizer layer $j \in [1, R]$ for learning. All tokens from the preceding layers $t^{0:j-1}$ are embedded and summed to form the token embedding input. After generating the base layer predictions of discrete tokens from the Masked Gesture Generator, these tokens are fed into the Residual Gesture Generator. This module iteratively predicts the tokens from the base layers, ultimately producing the final quantized output.

241

246 247 248

249

254

263

270 Inference. While existing works Liu et al. (2023); Yi et al. (2023); Chen et al. (2024) leverage 271 auto-regressive next-token prediction or diffusion-based generation process, these strategies hinder 272 the fast synthesis for real-time applications. To resolve this problem, as in Fig. 3, we employ an 273 iterative mask prediction strategy to decode motion tokens during inference. Initially, all tokens are 274 masked except for the first token from the source frame. Conditioned on the audio input, the Mask Gesture Generator predicts probabilities for the masked tokens. In the l-th iteration, the tokens with 275 the lowest confidence are re-masked, while the remaining tokens stay unchanged for subsequent 276 iterations. This updated sequence continues to inform predictions until the final iteration, when 277 the base-layer tokens are fully generated. Upon completion, the Residual Gesture Generator uses 278 the predicted base-layer tokens to progressively generate sequences for the remaining quantization 279 layers. Finally, all tokens are transformed back into motion sequences via the RVQ-VAE decoder. 280

Training Objective. To train our gesture generation models, \mathcal{L}_{mask} , and \mathcal{L}_{res} functions for two generators respectively by minimizing the categorical cross-entropy loss, as illustrated below:

$$\mathcal{L}_{mask} = \sum_{i=1}^{T} -\log p_{\phi}(t_i | Es(S), \text{MASK}), \quad \mathcal{L}_{res} = \sum_{j=1}^{V} \sum_{i=1}^{T} -\log p_{\phi}(t_i^j | t_i^{1:j-1}, Es(S), j).$$
(4)

In this formulation, \mathcal{L}_{mask} predicts the masked motion tokens t_i at each time step *i* based on the input audio and the special [MASK] token. Conversely, \mathcal{L}_{res} focuses on learning from multiple quantization layers, where t_i^j represents the motion token from quantizer layer *j* and $t_i^{1:j-1}$ includes the tokens from preceding layers. We also feed the predicted tokens into the RVQ decoder for gesture reconstructions, with velocity and acceleration losses Tevet et al. (2022); Siyao et al. (2022).

291 292

293

300

308

312 313 314

322

283 284 285

4.3 STRUCTURE-AWARE IMAGE REFINEMENT

To transfer gesture generation to pixel-level video synthesis, we leverage TPS Zhao & Zhang (2022) to achieve portrait animation based on gesture pattern keypoints from Sec. 4.2 through image warping. To address the uncertainties by optical-flow-based deformation, particularly in large motion regions such as the hands and shoulders, we propose a Semantic-Aware Generator. Auto-Link He et al. (2023) demonstrates that the learning of keypoint connections for image reconstruction aids the model in understanding image semantics. Based on this, we leverage keypoint connections as semantic guidance for image refinement.

Learnable Edge Heatmaps. Using the gesture motion keypoints, we establish linkages between them to provide structural information. To optimize computational efficiency, we limit the number of keypoint connections to those defined by body joint relationships Wan et al. (2017), rather than considering all potential connections in He et al. (2023).

For two keypoints k_i and k_j within predefined connection groups, we create a differentiable edge map S_{ij} . This edge is modeled as a Gaussian function extending along the line connecting the keypoints. Formally, the edge map S_{ij} for keypoints (k_i, k_j) is defined as:

$$\mathbf{S}_{ij}(\mathbf{p}) = \exp\left(v_{ij}(\mathbf{p})d_{ij}^2(\mathbf{p})/\sigma^2\right),\tag{5}$$

where σ is a learnable parameter controlling the edge thickness, and $d_{ij}(p)$ is the L_2 distance between the pixel p and the edge defined by keypoints k_i and k_j :

$$\boldsymbol{d}_{ij}(\boldsymbol{p}) = \begin{cases} \|\boldsymbol{p} - \boldsymbol{k}_i\|_2 & \text{if } t \leq 0, \\ \|\boldsymbol{p} - ((1-t)\boldsymbol{k}_i + t\boldsymbol{k}_j)\|_2 & \text{if } 0 < t < 1, \\ \|\boldsymbol{p} - \boldsymbol{k}_j\|_2 & \text{if } t \geq 1, \end{cases} \text{ where } t = \frac{(\boldsymbol{p} - \boldsymbol{k}_i) \cdot (\boldsymbol{k}_j - \boldsymbol{k}_i)}{\|\boldsymbol{k}_i - \boldsymbol{k}_j\|_2^2}.$$
(6)

Here, t denotes the normalized distance between k_i and the projection of p onto the edge.

To derive the edge map $S \in \mathbb{R}^{H \times W}$, we take the maximum value at each pixel across all heatmaps: $S(p) = \max_{ij} S_{ij}(p).$ (7)

We generate heatmaps at various resolutions. Inspired by SPADE Park et al. (2019), we treat these structural heatmaps as semantic guidance for image generation. A U-Net with residual blocks utilizes spatial semantic control from the edge heatmaps to refine the final video output.

Training Objective. We employ a conditional adversarial loss Mirza & Osindero (2014), along with perceptual similarity loss Johnson et al. (2016) and L1 loss for image refinement. The



Figure 4: Visual comparisons. Our method generates high-quality hand and shoulder motions, and presents metaphoric gestures when saying "90 joules," and "in each case."

discriminator utilizes the edge heatmap as a condition to compare generation against ground truth:

 $\mathcal{L}_{cGAN}(G, D) = \mathbb{E}_{I_{gt}, map} \left[\log D\left(I_{gt}, map\right) \right] + \mathbb{E}_{map} \left[\log \left(1 - D\left(map, I_{gan}\right)\right) \right], \tag{8}$ where *map* denotes the edge heatmap, and (\cdot, \cdot) indicates concatenation.

5 EXPERIMENTS

Since our work focuses on joint gesture motion and video generation, the main experiments primarily compare our approach with existing methods that also address joint generation. Comparisons specifically for gesture motion generation and avatar video rendering are deferred to the Appendix, where they are treated as separate, disentangled modules and compared with relevant works.

360 361 5.1 Experimental Settings

Dataset and preprocessing. We utilize PATS Ginosar et al. (2019); Ahuja et al. (2020) for the experiments. It contains 84,000 clips from 25 speakers with a mean length of 10.7s, 251 hours in total. For a fair comparison, following the literature Liu et al. (2022c); He et al. (2024) and replace the missing subject, with 4 speakers are selected (*Noah, Kubinec, Oliver*, and *Seth*). All video clips are cropped with square bounding boxes, centering speaks, resized to 256 × 256. We defer the additional details in the Appendix. After filtering, we obtain around 1000 clips for each speaker, randomly divided into 90% for training and 10% for evaluation, 4,000 in total.

Baseline Methods. We benchmark Realistic-Gesture against several co-speech gesture video generation methods: (1) ANGIE Liu et al. (2022c), a work in co-speech gesture video synthesis; (2) MM-Diffusion Ruan et al. (2023), an audio-video generation model demonstrated on the AIST++ dataset Li et al. (2021) that produces audio-driven human motion videos; and (3) S2G-Diffusion He et al. (2024), the most recent advancement in this domain. Notably, due to MM-Diffusion's fixed generation of 34 frames, we segment the audio accordingly for each generation.

375 376

348

349 350 351

352

353

354

355 356

357

358

359

5.2 QUANTITATIVE EVALUATION

Evaluation Metrics. We evaluate gesture motion and pixel-level video quality separately. For gesture motion metrics, we use **Fréchet Gesture Distance (FGD)** Yoon et al. (2020) to measure

Nama	Gesture-Motion Evaluation				Video Quality Assessment			
Indiffe	$FGD\downarrow$	Div. ↑	BAS \uparrow	PCM ↑	$FVD\downarrow$	$VQA_A \uparrow$	$VQA_T \uparrow$	
Ground Truth (GT)	0.0	14.012	1.00	1.000	0.0	95.694	5.329	
ANGIE	67.524	6.674	0.783	0.372	526.247	88.144	4.729	
MM-Diffusion	137.62	3.212	0.646	0.112	-	79.561	4.238	
S2G-Diffusion	23.646	10.848	0.974	0.447	486.134	93.553	5.401	
Ours	1.303	13.260	0.996	0.572	476.120	96.326	6.081	

Table 1: Quantitative results on the test set. Bold indicates the best performance. Our method performs betterin terms of both gesture motions and video generation quality.

the distribution gap between real and generated gestures in feature space, Diversity (Div.) Lee et al.
(2019) to calculate the average feature distance between generated gestures, Beat Alignment Score
(BAS) following Li et al. (2021), and Percent of Correct Motion parameters (PCM), difference of
generation deviate from ground-truth following Chen et al. (2024). We extract 2D human poses for
face and body using MMPose OpenMMLab (2020), which differs from S2G-Diffusion that focuses
solely on body poses. The first two metrics utilize an auto-encoder trained on PATS poses.

For pixel-level video quality, we assess **Fréchet Video Distance (FVD)** Unterthiner et al. (2018) for the overall quality of gesture videos, **VQA**_A for aesthetics and **VQA**_T for technical quality based on Dover Wu et al. (2023), pretrained on a large-scale dataset with labels ranked by real users.

397 **Evaluation Results.** We present quantitative evaluations in Tab. 1. Our approach significantly 398 outperforms existing methods in gesture motion metrics, achieving an FGD of 1.303 and a Diversity 399 score of 13.260. These results indicate that our generated gesture patterns closely resemble those from ground-truth videos, exhibiting both high naturalness and a broader range of motion patterns. 400 For video quality assessment, we use the FVD metric to evaluate the similarity of the generated video 401 distribution to the ground-truth videos. Our model achieves the lowest FVD among the compared 402 methods, demonstrating superior performance. The VQA_A and VQA_T metrics measure perceived 403 user preferences for video generation content. Notably, our approach yields a VQA_A of **96.326** 404 and a VQA_T of **6.081**, surpassing the ground-truth videos. This success can be attributed to our 405 structure-aware image enhancement design. In contrast, MM-Diffusion produces limited gesture 406 patterns due to its design, which generates only a few continuous frames and struggles to learn 407 diverse motion patterns from speech audio. ANGIE leverages MRAA Siarohin et al. (2021) for 408 regional coarse motion patterns but lacks the precision necessary for motion control aligned with 409 speech, resulting in low diversity and beat alignment. S2G-Diffusion performs better than ANGIE 410 but still fails to generate fine-grained gesture patterns, as it relies on image optical flows without adequately focusing on the nuances of human facial and body movements. 411

412

380 381 382

384 385 386

413 5.3 QUALITATIVE EVALUATION

Table 2:Subjective evaluation results areshown as Mean Opinion Scores (MOS).

Evaluation Results. We provide qualitative evaluations of video generation in Fig. 4. MM-Diffusion generates unrealistic shoulder movements. ANGIE produces misaligned gesture motions with the accompanying speech. Although S2G-Diffusion shows improvement over ANGIE, it struggles with local

	1			/				
Methods	MOS_1	MOS_2	MOS ₃	MOS_4				
Wiethous	User Study							
GT	4.7	4.7	4.7	4.65				
MM-Diffusion	1.35	1.65	1.4	1.55				
ANGIE	1.95	3.25	1.9	2.25				
S2G-Diffusion	3.0	3.6	3.15	3.0				
Ours	3.35	3.05	3.35	3.25				

regions, such as the hands, due to its reliance on unsupervised keypoints for global transformations, which neglects local deformations. In contrast, our method demonstrates high-quality video generation, particularly in the facial and body areas. The alignment between gesture and speech is notably enhanced through our speech-content-aware gesture latent representation. For example, when the actor says "90 *joules*," he points to the screen, and he emphasizes phrases like "*so two ways*" and "*in each case*" by raising his hands as metaphoric gestures. This coordination exemplifies our approach's capability to produce contextually relevant and expressive gestures.

User Study. We conducted a user study to evaluate the visual quality of our method. We sampled
80 videos from each method and ground-truth, and invited 20 participants to conduct Mean Opinion
Scores (MOS) evaluations. The rating ranges from 1 (poorest) to 5 (highest). Participants rated the
videos on: (1) MOS₁: "*How realistic does the video appear*?", (2) MOS₂: "*How diverse does the gesture pattern present*?", (3) MOS₃: "*Are speech and gesture synchronized in this video*?" and (4)
MOS₄: "What is your overall evaluation of the video". The videos were presented in random order

Kp Repr.	FVD↓	LPIPS↓	PSNR↑	G-Repr.	FGD↓	Div.↑	PCM↑	G-0	Gen.	FGD↓	Div.↑	PCM↑
Unsup-kp	387.05	0.05	27.41	baseline	262.675	18.142	0.279	w/o	res	3.372	11.359	0.513
2D-pose	272.18	0.05	27.20	+KVQ	34.940	0./13	0.327	con	cat	3.415	11.314	0.514
+ nex kp	3/7.14	0.06	25.30	+ separat	e 21.4/3	10.536	0.412	w/o a	align	8.382	11.452	0.373
ruii-modei	225.77	0.04	27.17	+ distill	1.303	13.260	0.582	full-n	nodel	1.303	13.260	0.582
(a) Configu design.	urations	for keyp	ooint	(b) Gest tions.	ure motio	on repre	esenta-	(c) Ge	enera	tor arch	iecture	design.
Refine	VQA _A	↑ VQA _T	↑ FVD,	\downarrow <i>M-Ra</i>	<i>ıtio</i> FGD↓	Div.↑	PCM↑	· ·	iter.	FGD↓	Div.↑	PCM↑
w/o refine	91.248	5.381	492.34	41 Uni ()-1 3.348	14.312	0.513		5	1.303	13.260	0.582
+ UNet	93.958	5.479	484.32	23 Uni.	3-1 3.232	12.58	0.512		10	1.642	13.40	0.575
+ skeleton	95.902	5.479	475.63	6 Uni.	5-1 1.303	13.260	0.582		15	1.828	13.49	0.573
+ heatmap	96.326	6.081	476.12	20 Uni .	7-1 1.790	13.49	0.572		20	1.881	13.49	0.572
(d) Image-	refineme	ent strateg	gies.	(e) n	ask-ratio o	during tr	aining.		(f) N	Aask de	coding	steps.
	Not and Belleving Provide the State Provide the	۲	*		1 ⁴⁸ (S 	3		*		tandare Ze-2 Ze-2 Afra - 2 Afra - 2 Afr	
unsupervised	keypoints		- RVQ +	Distill		-RVQ base la	iyer (No resid	uai) + Distili		w/t	Tenne	
	Note and Construction and un-operation Note of the Note of the Not		۲. ۱	۲. ۱۹		· · ·	E	3	141 141		Sandar (27-1) 27-1)	
2D poses +RVQ			2		+RV	/Q + Distill			W	refine		

Table 3: Ablations of our method. We exam the keypoint design, gesture representation, gesture generator
 architecture, training & inference strategy and image-refinement. Bold indicates the best performance.

Figure 5: Ablation visualizations. Left: motion by unsupervised keypoints or 2d poses; Middle: RVQ-based gesture representation and generation; Right: image-refinement helps hand generation.

to capture participants' initial impressions. As shown in Tab. 2, our method outperformed others across realness, synchronization, and overall quality, but lower in diversity than S2G-Diffusion. We defer additional details of the user study in the appendix.

5.4 ABLATION STUDY

In this section, we present ablation study of keypoint design for image warping, gesture pattern
 representation exploraton, gesture generator architecture design, and varios comparisons of image refinement. We defer additional experiments in the Appendix.

Motion Keypoint Design We evaluate three keypoint representations for image-warping: (1) unsupervised keypoints for global optical-flow transformation (as in ANGIE and S2G-Diffusion), (2) 2D human poses, and (3) 2D human poses augmented with flexible learnable points. Each design is assessed using TPS Zhao & Zhang (2022) transformation, with self-reconstruction based on these keypoints for evaluation. We compare the full-model reconstruction with refinement against the first three designs without refinement. As shown in Tab. 3a, learnable keypoints lead to a significant decrease in FVD, highlighting their inadequacy for motion control. The 2D landmarkbased keypoints yield slightly lower SSIM scores, likely due to their limited capacity to represent global transformations. The inclusion of flexible keypoints does not enhance the image-warping outcomes. Consequently, we opt to utilize 2D pose landmarks exclusively for our study.

Motion Representation. We evaluate several configurations: (1) baseline: no motion representation, relying solely on the generator to synthesize raw 2D landmarks; (2) + RVQ: utilizing Residual VQ (RVQ) to encode joint face-body keypoints; (3) + separate motion: employing two RVQs for independent face and body motions; (4) + distill: learning joint embeddings for speech and gesture in both face and body motions, followed by distillation for RVQ tokenization. We discover RVQ significantly improve the precise pose location while distillation leads to natural movements.



Figure 6: Our model supports multiple video gesture generation end editing applications.

494 Generator Design. We explore various designs for the gesture generator: (1) w/o res: no residual 495 gesture decoder; (2) concat: instead of using cross-attention for audio control, we concatenate 496 the audio features with gesture latent features element-wise during generation; (3) w/o align: the 497 audio encoder is randomly initialized rather than initialized from face and body contrastive learning. 498 Our findings indicate that the Residual Gesture Generator significantly enhances finger motion 499 generation. The cross-attention design outperforms element-wise concatenation, while the pre-500 alignment of the audio encoder notably improves FGD. We attribute this improvement to the shared similarities in the codebook and audio encoders during contrastive alignment. 501

502 **Image Refinement.** We examine various network designs for motion generation, specifically: (1) 503 w/o refine: no image refinement, relying solely on image warping; (2) + UNet: employing a standard 504 UNet; (3) + pose skeleton: integrating connected skeleton maps as in the diffusion ReferenceNet Hu 505 et al. (2023); (4) + edge heatmap: substituting the previous design with our learnable edge heatmap. Our experiments reveal that the edge heatmap outperforms skeleton maps, likely due to the learnable 506 thickness of connections, which provides better semantic-aware generation guidance. 507

508 Training and Inference Strategy. We evaluate the mask ratio during training and the number of 509 inference steps during decoding. As shown in Tab. 3f, our model requires only 5 inference steps, in 510 contrast to over 50 or 100 steps in diffusion-based models. Furthermore, a uniform masking ratio 511 between 0.5 and 1 during training yields optimal performance.

512 513

514

493

5.5 APPLICATION

515 **Long Sequence Generation.** Shown on the left of Fig. 6, to generate long sequences, we start 516 with the initial frame and the corresponding target audio, which we segment into smaller windows. 517 After generating the first segment, we use the last few frames of the generated output as the new starting frame conditions for the next segment of audio, allowing for a iterative outpainting. 518

519

520 **Video Gesture Editing.** For gesture editing and inpainting, we first extract the keypoints from a given video sequence and tokenize the face and body movements into motion tokens. Thanks to the 521 model's bidirectional decoding capability, we insert [MASK] tokens wherever edits are needed. 522 Trained on temporal masking, the model generate coherent gestures in the masked areas. By 523 incorporating different speech input and speaker embeddings, we can create new gesture patterns 524 and re-render the video based on the edited latent motion tokens. 525

526 Gesture Pattern Transfer. Given the design of our framework, with different identity embedding, the model can generate different gesture patterns given the input identity embedding control given 528 the same audio. Please see the demo videos in our Appendix for more details. 529

530 531

532

527

6 CONCLUSION

533 We present **Realistic-Gesture**, a framework for generating realistic co-speech gesture videos. To 534 ensure the gestures cohere well with speech, we propose speech-content aware gesture motion 535 representation though knowledge distillation from the gesture-speech aligned features obtained 536 through contrastive learning. Our masked gesture motion generator enables the creation and editing 537 of long, high-quality gesture motion sequences. Our pixel-level refinement module further improves the transformation of inferred gesture motions into realistic animations for large-scale body motion. 538 We believe this work will encourage further exploration of the relationship between gesture patterns and speech context for more compelling gesture video generations in the future.

540 REFERENCES

Chaitanya Ahuja, Dong Won Lee, Ryo Ishii, and Louis-Philippe Morency. No Gestures Left Behind: 542 Learning Relationships between Spoken Language and Freeform Gestures. In Proceedings of the 543 2020 Conference on Empirical Methods in Natural Language Processing: Findings, pp. 1884– 544 1895, 2020. 546 Tenglong Ao, Qingzhe Gao, Yuke Lou, Baoquan Chen, and Libin Liu. Rhythmic gesticulator: 547 Rhythm-aware co-speech gesture synthesis with hierarchical neural embeddings. ACM Transac-548 tions on Graphics (TOG), 41(6):1-19, 2022. 549 Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik 550 Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, Varun Jampani, and Robin 551 Rombach. Stable Video Diffusion: Scaling Latent Video Diffusion Models to Large Datasets, 552 2023. URL https://arxiv.org/abs/2311.15127. 553 554 Judee K Burgoon, Thomas Birk, and Michael Pfau. Nonverbal Behaviors, Persuasion, and Credibility. Human communication research, 17(1):140-169, 1990. 555 556 Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei A Efros. Everybody Dance Now. In 557 Proceedings of the IEEE International Conference on Computer Vision, 2019. 558 559 Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Murphy, William T Freeman, Michael Rubinstein, et al. Muse: Text-to-Image 560 Generation Via Masked Generative Transformers. arXiv preprint arXiv:2301.00704, 2023. 561 562 Junming Chen, Yunfei Liu, Jianan Wang, Ailing Zeng, Yu Li, and Qifeng Chen. DiffSHEG: A 563 Diffusion-Based Approach for Real-Time Speech-driven Holistic 3D Expression and Gesture 564 Generation, 2024. URL https://arxiv.org/abs/2401.04747. 565 Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki 566 Kanda, Takuya Yoshioka, Xiong Xiao, et al. WavLM: Large-Scale Self-Supervised Pre-Training 567 for Full Stack Speech Processing. IEEE Journal of Selected Topics in Signal Processing, 16(6): 568 1505-1518, 2022. 569 570 Jan P De Ruiter, Adrian Bangerter, and Paula Dings. The Interplay Between Gesture and Speech 571 in the Production of Referring Expressions: Investigating the Tradeoff Hypothesis. Topics in 572 cognitive science, 4(2):232-248, 2012. 573 Anna Deichler, Shivam Mehta, Simon Alexanderson, and Jonas Beskow. Diffusion-based co-574 speech gesture generation using joint text and audio representation. In INTERNATIONAL 575 CONFERENCE ON MULTIMODAL INTERACTION, ICMI '23. ACM, October 2023. doi: 576 10.1145/3577190.3616117. URL http://dx.doi.org/10.1145/3577190.3616117. 577 578 Jacob Devlin. BERT: Pre-Training of Deep Bidirectional Transformers for Language Understanding. 579 arXiv preprint arXiv:1810.04805, 2018. 580 S. Ginosar, A. Bar, G. Kohavi, C. Chan, A. Owens, and J. Malik. Learning Individual Styles of 581 Conversational Gesture. In Proceedings of the IEEE Conference on Computer Vision and Pattern 582 Recognition. IEEE, June 2019. 583 584 Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh Agrawala, Dahua Lin, and Bo Dai. AnimateDiff: Animate Your Personalized Text-to-Image 585 Diffusion Models without Specific Tuning. Proceedings of the International Conference on 586 Machine Learning, 2024. 588 Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked 589 Autoencoders Are Scalable Vision Learners. In Proceedings of the IEEE Conference on Computer 590 Vision and Pattern Recognition, pp. 16000–16009, 2022. 591 Xingzhe He, Bastian Wandt, and Helge Rhodin. AutoLink: Self-Supervised Learning of Human 592 Skeletons and Object Outlines by Linking Keypoints, 2023. URL https://arxiv.org/ abs/2205.10636.

614

630

- 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 MotionDecoupled Diffusion Model. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2263–2273, 2024.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising Diffusion Probabilistic Models. Advances in neural information processing systems, 33:6840–6851, 2020.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
 and Weizhu Chen. LoRA: Low-Rank Adaptation of Large Language Models. In *Proceedings* of the International Conference on Machine Learning, 2022. URL https://openreview.
 net/forum?id=nZeVKeeFYf9.
- Li Hu, Xin Gao, Peng Zhang, Ke Sun, Bang Zhang, and Liefeng Bo. Animate Anyone: Consistent and Controllable Image-to-Video Synthesis for Character Animation. *arXiv preprint arXiv:2311.17117*, 2023.
- 609Xun Huang and Serge Belongie. Arbitrary Style Transfer in Real-time with Adaptive Instance610Normalization, 2017. URL https://arxiv.org/abs/1703.06868.
- ⁶¹¹ Ziyao Huang, Fan Tang, Yong Zhang, Xiaodong Cun, Juan Cao, Jintao Li, and Tong-Yee
 ⁶¹² Lee. Make-your-anchor: A diffusion-based 2d avatar generation framework. *arXiv preprint* 613 *arXiv:2403.16510*, 2024.
- Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual Losses for Real-Time Style Transfer and Super-Resolution, 2016. URL https://arxiv.org/abs/1603.08155.
- Johanna Karras, Aleksander Holynski, Ting-Chun Wang, and Ira Kemelmacher-Shlizerman. Dream Pose: Fashion Image-to-Video Synthesis via Stable Diffusion. *arXiv preprint arXiv:2304.06025*, 2023.
- Doyup Lee, Chiheon Kim, Saehoon Kim, Minsu Cho, and Wook-Shin Han. Autoregressive Image
 Generation Using Residual Quantization, 2022. URL https://arxiv.org/abs/2203.
 01941.
- Hsin-Ying Lee, Xiaodong Yang, Ming-Yu Liu, Ting-Chun Wang, Yu-Ding Lu, Ming-Hsuan Yang, and Jan Kautz. Dancing to Music. *Proceedings of the Neural Information Processing Systems Conference*, 32, 2019.
- Ruilong Li, Shan Yang, David A Ross, and Angjoo Kanazawa. AI Choreographer: Music
 Conditioned 3D Dance Generation with AIST++. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 13401–13412, 2021.
- Tianhong Li, Huiwen Chang, Shlok Mishra, Han Zhang, Dina Katabi, and Dilip Krishnan.
 MAGE: Masked Generative Encoder to Unify Representation Learning and Image Synthesis. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 2142–2152, 2023.
- Haiyang Liu, Naoya Iwamoto, Zihao Zhu, Zhengqing Li, You Zhou, Elif Bozkurt, and Bo Zheng.
 DisCo: Disentangled Implicit Content and Rhythm Learning for Diverse Co-Speech Gestures
 Synthesis. In *Proceedings of the 30th ACM International Conference on Multimedia*, pp. 3764–3773, 2022a.
- 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. *arXiv preprint arXiv:2203.05297*, 2022b.
- Haiyang Liu, Zihao Zhu, Giorgio Becherini, Yichen Peng, Mingyang Su, You Zhou, Naoya Iwamoto, Bo Zheng, and Michael J Black. EMAGE: Towards Unified Holistic Co-Speech Gesture Generation via Masked Audio Gesture Modeling. *arXiv preprint arXiv:2401.00374*, 2023.
- Kian Liu, Qianyi Wu, Hang Zhou, Yuanqi Du, Wayne Wu, Dahua Lin, and Ziwei Liu. Audio-Driven
 Co-Speech Gesture Video Generation. *Proceedings of the Neural Information Processing Systems Conference*, 35:21386–21399, 2022c.

- 648 Xian Liu, Qianyi Wu, Hang Zhou, Yinghao Xu, Rui Qian, Xinyi Lin, Xiaowei Zhou, Wayne Wu, 649 Bo Dai, and Bolei Zhou. Learning Hierarchical Cross-Modal Association for Co-Speech Gesture 650 Generation. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 651 pp. 10462–10472, 2022d. 652 Xiaofeng Mao, Zhengkai Jiang, Qilin Wang, Chencan Fu, Jiangning Zhang, Jiafu Wu, Yabiao Wang, 653 Chengjie Wang, Wei Li, and Mingmin Chi. Mdt-a2g: Exploring masked diffusion transformers 654 for co-speech gesture generation. In Proceedings of the 32nd ACM International Conference on 655 Multimedia, MM '24, pp. 3266–3274. ACM, October 2024. doi: 10.1145/3664647.3680684. 656 URL http://dx.doi.org/10.1145/3664647.3680684. 657 658 Brian McFee, Colin Raffel, Dawen Liang, Daniel PW Ellis, Matt McVicar, Eric Battenberg, and 659 Oriol Nieto. librosa: Audio and Music Signal Analysis in Python. In Proceedings of the 14th Python in Science Conference, volume 8, 2015. 660 661 Mehdi Mirza and Simon Osindero. Conditional Generative Adversarial Nets, 2014. URL https: 662 //arxiv.org/abs/1411.1784. 663 664 OpenMMLab. OpenMMLab Pose Estimation Toolbox and Benchmark. https://github. 665 com/open-mmlab/mmpose, 2020. 666 Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic Image Synthesis with 667 Spatially-Adaptive Normalization. In Proceedings of the IEEE Conference on Computer Vision 668 and Pattern Recognition, 2019. 669 670 Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman, 671 Dimitrios Tzionas, and Michael J. Black. Expressive Body Capture: 3D Hands, Face, and Body 672 from a Single Image. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019. 673 674 Ekkasit Pinyoanuntapong, Pu Wang, Minwoo Lee, and Chen Chen. MMM: Generative Masked 675 Motion Model. In Proceedings of the IEEE Conference on Computer Vision and Pattern 676 Recognition, pp. 1546–1555, 2024. 677 678 Lawrence Rabiner and Ronald Schafer. Theory and Applications of Digital Speech Processing. Prentice Hall Press, 2010. 679 680 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, 681 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya 682 Sutskever. Learning Transferable Visual Models From Natural Language Supervision, 2021. URL 683 https://arxiv.org/abs/2103.00020. 684 685 Ludan Ruan, Yiyang Ma, Huan Yang, Huiguo He, Bei Liu, Jianlong Fu, Nicholas Jing Yuan, Qin Jin, and Baining Guo. MM-Diffusion: Learning Multi-Modal Diffusion Models for Joint Audio 686 and Video Generation. In Proceedings of the IEEE Conference on Computer Vision and Pattern 687 Recognition, 2023. 688 689 Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. First 690 Order Motion Model for Image Animation. In Proceedings of the Neural Information Processing 691 Systems Conference, 2019. 692 Aliaksandr Siarohin, Oliver Woodford, Jian Ren, Menglei Chai, and Sergey Tulyakov. Motion 693 Representations for Articulated Animation. In Proceedings of the IEEE Conference on Computer 694 Vision and Pattern Recognition, 2021. 695 696 Li Siyao, Weijiang Yu, Tianpei Gu, Chunze Lin, Quan Wang, Chen Qian, Chen Change Loy, and 697 Ziwei Liu. Bailando: 3D Dance Generation by Actor-Critic GPT with Choreographic Memory. 698 In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 11050– 699 11059, 2022. 700
- ⁷⁰¹ Guy Tevet, Sigal Raab, Brian Gordon, Yonatan Shafir, Daniel Cohen-Or, and Amit H Bermano. Human Motion Diffusion Model. *arXiv preprint arXiv:2209.14916*, 2022.

702 703 704	Linrui Tian, Qi Wang, Bang Zhang, and Liefeng Bo. EMO: Emote Portrait Alive-Generating Expressive Portrait Videos with Audio2Video Diffusion Model Under Weak Conditions. <i>arXiv</i> preprint arXiv:2402.17485, 2024.
705 706 707 708	Thomas Unterthiner, Sjoerd Van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, and Sylvain Gelly. Towards Accurate Generative Models of Video: A New Metric & Challenges. <i>arXiv preprint arXiv:1812.01717</i> , 2018.
709 710	Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural Discrete Representation Learning, 2018. URL https://arxiv.org/abs/1711.00937.
712 713	Qingfu Wan, Wei Zhang, and Xiangyang Xue. DeepSkeleton: Skeleton Map for 3D Human Pose Regression, 2017. URL https://arxiv.org/abs/1711.10796.
714 715 716	Chengyi Wang, Sanyuan Chen, Yu Wu, Ziqiang Zhang, Long Zhou, Shujie Liu, Zhuo Chen, Yanqing Liu, Huaming Wang, Jinyu Li, et al. Neural Codec Language Models Are Zero-Shot Text to Speech Synthesizers. <i>arXiv preprint arXiv:2301.02111</i> , 2023a.
717 718 719	Congyi Wang. T2m-hifigpt: Generating high quality human motion from textual descriptions with residual discrete representations, 2023. URL https://arxiv.org/abs/2312.10628.
720 721 722	Tan Wang, Linjie Li, Kevin Lin, Chung-Ching Lin, Zhengyuan Yang, Hanwang Zhang, Zicheng Liu, and Lijuan Wang. DisCo: Disentangled Control for Referring Human Dance Generation in Real World. arXiv preprint arXiv:2307.00040, 2023b.
723 724 725	Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High- Resolution Image Synthesis and Semantic Manipulation with Conditional GANs. In <i>CVPR</i> , 2018.
726 727 728 729	Haoning Wu, Erli Zhang, Liang Liao, Chaofeng Chen, Jingwen Hou Hou, Annan Wang, Wenxiu Sun Sun, Qiong Yan, and Weisi Lin. Exploring Video Quality Assessment on User Generated Contents from Aesthetic and Technical Perspectives. In <i>Proceedings of the IEEE International Conference on Computer Vision</i> , 2023.
730 731 732	Zunnan Xu, Yachao Zhang, Sicheng Yang, Ronghui Li, and Xiu Li. Chain of generation: Multi- modal gesture synthesis via cascaded conditional control, 2023. URL https://arxiv.org/ abs/2312.15900.
734 735 736	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 <i>Proceedings of the</i> <i>IEEE Conference on Computer Vision and Pattern Recognition</i> , 2023.
737 738 739	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. <i>ACM Transactions on Graphics</i> , 39(6), 2020.
740 741 742	Jian Zhao and Hui Zhang. Thin-Plate Spline Motion Model for Image Animation. In <i>Proceedings</i> of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 3657–3666, 2022.
743 744 745 746	Shenhao Zhu, Junming Leo Chen, Zuozhuo Dai, Yinghui Xu, Xun Cao, Yao Yao, Hao Zhu, and Siyu Zhu. Champ: Controllable and Consistent Human Image Animation with 3D Parametric Guidance. <i>arXiv preprint arXiv:2403.14781</i> , 2024.
747 748 749	
750 751	
752 753	
754	