## **COCOGESTURE:** TOWARDS <u>COHERENT</u> <u>CO</u>-SPEECH 3D GESTURE GENERATION IN THE WILD

Anonymous authors

000

001

002 003 004

006

008

009 010 011

012

021

022 023

024 025

026

027

028

029

031

032

034

038

039

040

041

042

043

044





Figure 1: Our CoCoGesture framework pre-trained on the large-scale dataset can generate coherent and diverse 3D co-speech gestures corresponding with unseen zero-shot human audios.

#### ABSTRACT

Deriving co-speech 3D gestures has seen tremendous progress in virtual avatar animation. Yet, the existing methods often produce stiff and unreasonable gestures with unseen human speech inputs due to the limited 3D speech-gesture data. In this paper, we propose **CoCoGesture**, a novel framework enabling coherent and diverse gesture synthesis from unseen human speech prompts. Our key insight is built upon the custom-designed pretrain-fintune training paradigm. At the pretraining stage, we aim to formulate a large generalizable gesture diffusion model by learning the abundant postures manifold. Therefore, to alleviate the scarcity of 3D data, we first construct a large-scale co-speech 3D gesture dataset containing more than 40M meshed posture instances across 4.3K speakers, dubbed GES-X. Then, we scale up the large unconditional diffusion model to 1B parameters and pre-train it to be our gesture experts. At the finetune stage, we present the audio ControlNet that incorporates the human voice as condition prompts to guide the gesture generation. Here, we construct the audio ControlNet through a trainable copy of our pretrained diffusion model. Moreover, we design a novel Mixture-of-Gesture-Experts (MoGE) block to adaptively fuse the audio embedding from the human speech and the gesture features from the pre-trained gesture experts with a routing mechanism. Such an effective manner ensures audio embedding is temporal coordinated with motion features while preserving the vivid and diverse gesture generation. Extensive experiments demonstrate that our proposed CoCoGesture outperforms the state-of-the-art methods on the zero-shot speech-to-gesture generation. The dataset will be publicly available at: https://anonymous.4open.science/w/GES-X/.

045 046 047

048

1 INTRODUCTION

Co-speech gesture generation aims to synthesize vivid and diverse human postures coordinated with
the input speech audio. These non-verbal body languages greatly enhance the delivery of speech
content in daily conversations (Qi et al., 2024; 2023a; Liu et al., 2024a). Meanwhile, synthesizing
co-speech gestures of human avatars plays a significant role in wide applications like robotics (Farouk,
2022), virtual/augmented reality (AR/VR) (Fu et al., 2022), and human-machine interaction (Koppula & Saxena, 2013; Liu et al., 2023a).

Conventionally, recent researchers deal with speech-to-gesture tasks by modeling human upper-body dynamics with consistent speech voice (Liu et al., 2024a; 2022a; Yi et al., 2023; Chen et al., 2024; Liu et al., 2024b; Qi et al., 2024). Most of them address this task by conducting end-to-end mapping through the pre-defined corpus (Liu et al., 2022a; 2024a; Yi et al., 2023). However, they usually heavily rely on the paired audio-gesture data covering limited speaker identities, resulting in insufficient diversity of gestures. Moreover, the narrowed corpus data may lead to the model falling short of generalizing to unseen out-of-domain audio inputs, as shown in Figure 2(a).

In this work, we introduce the task of coher-061 ent and diverse co-speech 3D gesture genera-062 tion from in-the-wild human voices, depicted 063 in Figure 1. To achieve this goal, there are two 064 main challenges: 1) The existing 3D meshed 065 co-speech gesture datasets (Liu et al., 2024a; Yi 066 et al., 2023) are insufficient to train a general-067 izable model. Creating such a dataset through 068 accurate motion capture systems is extensively 069 labor-consuming. 2) Modeling the coherent and diverse co-speech gestures from unseen human audio in an end-to-end fashion is difficult, espe-071 cially in long sequences. 072

To overcome the issue of data scarcity, we first
newly construct a large-scale 3D meshed cospeech whole-body dataset that contains more
than 40M posture instances across about 4.3K
aligned speaker audios, dubbed GES-X. Specif-



Figure 2: Dataset statistical comparison between our GES-X and existing meshed co-speech gesture datasets (*i.e.*BEAT2 (Liu et al., 2022a), Talk-SHOW (Yi et al., 2023)). Our GES-X has a much larger word corpus and a more widely uniform distributed gesture motion.

ically, thanks to the advanced pose estimator (Zhang et al., 2023a), we can obtain high-quality 3D
postures (*i.e.*, SMPL-X (Pavlakos et al., 2019) and FLAME (Li et al., 2017)) from in-the-wild talk
show videos. Then, by employing WhisperX (Bain et al., 2023) for automatic speech recognition, we
ensure the acquired text transcript and phoneme consistency with speaker audio. In this fashion, our
GES-X provides the most comprehensive co-speech gestures with diverse modalities. As reported in
Figure 2 (b), the posture motion degree of the GES-X dataset displays a much more widely uniform
distribution against others, indicating our dataset contains more diverse gestures. Meanwhile, the
common mesh standards in our dataset also support other downstream human dynamics-related tasks, *e.g.*, talking head generation (Tian et al., 2024), human motion generation (Ao et al., 2023).

Along with this dataset, we propose **CoCoGesture**, a novel framework that enables the generation 087 of coherent human gestures from the unseen voice. Our key insight is built upon the custom-880 designed pretrain-fintune training paradigm. To ensure the generalization of the pre-trained model, 089 we leverage our large-scale co-speech gesture dataset GES-X as the source training set. Specifically, 090 we first conduct the pre-training phase based on the large unconditional diffusion transformer 091 backbone (Peebles & Xie, 2023). This diffusion model serves as a gesture expert and is scaled up to 092 1B parameters, thereby enabling the training model to build the sufficiently inherent motion manifold 093 from massive gesture dynamics. In this manner, our pre-trained model ensures the realism of the generated gestures while preserving vividness and diversity. 094

Moreover, to incorporate the human speech as the conditional prompt coordinately, inspired by 096 (Zhang et al., 2023b), we present the audio ControlNet for fine-tuning. Concretely, we refactor a trainable copy of our pre-trained unconditional large model for adapting various audio conditions. 098 Then, we propose a novel block, named Mixture-of-Gesture-Experts (MoGE), to fuse the audio embedding from the human voice and the gesture features from pre-trained gesture experts through a 099 routing mechanism. Here, the routing mechanism adaptively balances the input audio signal features 100 with the retained original motion clues. Meanwhile, the learned temporal-wise soft weight of the 101 routing mechanism greatly guarantees generated results to maintain the coherence rhythm with 102 input human speeches. Extensive experiments conducted on the out-of-domain datasets (Liu et al., 103 2024a; Yi et al., 2023) demonstrate our fine-tuned framework synthesizes vivid and diverse co-speech 104 gestures, outperforming the state-of-the-art counterparts. Our GES-X dataset will be open-sourced 105 soon to facilitate the research on the relevant community. 106

107 Overall, our contributions are summarized as follows:

We introduce the task of co-speech gesture generation from in-the-wild human speech incorporating the large 3D meshed whole-body human posture dataset, named GES-X. It includes more than 40M high-quality gesture instances with 4.3K speakers, significantly facilitating research on diverse gesture generation.

- We propose a novel framework named CoCoGesture that leverages the Mixture-of-Gesture-Experts (MoGE) blocks to adapt various unseen audio signals with pre-trained highly generalizable gesture experts effectively. The presented MoGE greatly enhances the temporal coherence between generated results and conditional prompts.
- Extensive experiments show that our CoCoGesture produces vivid and diverse co-speech gestures given unseen human voices, outperforming state-of-the-art counterparts.
- 2 RELATED WORK
- 119 120 121

112

113

114

115

116

117

118

Co-speech Gesture Generation. Generating vivid and diverse co-speech gestures has witnessed 122 impressive progress in recent years due to its practical value in wide-range applications (Qi et al., 123 2023c; Liu et al., 2023a; Zhu et al., 2023; Liang et al., 2024; Tian et al., 2024). Conventionally, 124 researchers utilize the rule-based workflow to bridge the gap between human speech and gestures 125 via the pre-defined corpus by linguistic experts (Marsella et al., 2013; Poggi et al., 2005). Other 126 works generate the results relying on mapping the audio signals to manually defined gesture features 127 through machine learning (Cassell et al., 1994; Huang & Mutlu, 2012). Nevertheless, these two approaches both need much more effort in preliminary dataset design, causing them to be limited by 128 the size and quality of the datasets. 129

130 Recently, thanks to the advanced deep learning methods and 3D human body modeling tech-131 niques (Loper et al., 2023; Zhang et al., 2023a; Pavlakos et al., 2019; Boukhayma et al., 2019; 132 Li et al., 2017), many works are proposed to generate the continuous 3D upper body postures. Speech-133 gesture-aligned datasets (Liu et al., 2022b; Yi et al., 2023; Yoon et al., 2020; Liu et al., 2024a; 2022a) are also proposed to address this challenging task. They involve multi-modality clues to promote 134 the generated gestures to be much more reasonable and diverse, like emotion (Liu et al., 2022a; Qi 135 et al., 2023a; 2024; Bhattacharya et al., 2021), identity (Yi et al., 2023; Liu et al., 2022b; 2024b), text 136 transcript (Liu et al., 2022b; Cheng et al., 2024). To be specific, Ao et. al (Ao et al., 2022)propose 137 a rhythm-based segmentation pipeline to boost the harmony between speech and gestures. Yang 138 et. al (Yang et al., 2023) leverage emotion as guidance to produce various stylized gestures with 139 the specifically designed diffusion model. Ahuja et. al (Ahuja et al., 2020) mix the disentangled 140 gesture styles as an ensemble to guide the gesture generation. However, they overlook that directly 141 generating the gesture from an in-the-wild human voice is much more practical in real-world scenes. 142 Considering the previous datasets are restricted to a limited scale, we thus propose a large-scale 143 meshed 3D co-speech dataset to facilitate the research on audio-driven gesture generation from 144 in-the-wild human speeches.

145

146 Zero-shot Human Motion Generation. Human motion generation strives to generate natural sequences of human poses. Recent advancements in motion data collection and generation methods 147 have sparked growing interest in this field. Existing research primarily revolves around generating 148 human motions using conditional signals like text (Tevet et al., 2022b; Chen et al., 2023; Dabral 149 et al., 2023), audio (Tseng et al., 2022; Ao et al., 2023; Zhu et al., 2023), and scene contexts (Araujo 150 et al., 2023; Huang et al., 2023). Currently, open-set human motion generation focuses on zero-shot 151 text-driven generation (Reed et al., 2016; Lin et al., 2023), which creates new content from text 152 prompts without relying on pre-defined data. MotionCLIP (Tevet et al., 2022a) enhances zero-shot 153 generation by employing a Transformer-based autoencoder to align the motion manifold with the 154 latent space of pre-trained vision-language model CLIP (Radford et al., 2021). However, without 155 sufficient high-quality 3D motion data, current approaches still face challenges in generating fine-156 grained motions from unseen audio prompts. Therefore, we propose a novel framework to generate 157 vivid and diverse gestures based on zero-shot human speech.

158

Mixture-of-Experts. Mixture-of-Experts (MoE) refers to combining the strengths of multiple
expert models to improve model generalization performance (Fedus et al., 2022; Jacobs et al., 1991;
Shazeer et al., 2017). Recently, MoE has been extensively applied to various research areas (Gale et al., 2023; Pini et al., 2023), demonstrating their versatility and effectiveness. In computer vision,

3

Table 1: Statistical comparison of our **GES-X** with existing ones. The dotted line separates whether 163 the posture in the dataset is built based on the mesh. Among meshed whole body co-speech gesture 164 datasets, the scale of our GES-X is  $15 \times$  larger than the existing ones (*i.e.*BEAT2). 165

-	Duration Attributes				Joint				
Dataset	(hours)	Speakers	Facial	Mesh	Phoneme	Text	Body	Hand	Annotation
Trinity (Ferstl & McDonnell, 2018)	4	1	×	x	×	1	24	38	mo-cap
TED (Yoon et al., $2020$ ) <sub>TOG</sub>	106.1	1,766	X	X	×	1	9	X	pseudo
SCG (Habibie et al., 2021) <sub>CVPR</sub>	33	6	X	X	×	X	14	24	pseudo
TED-Ex (Liu et al., 2022b) <sub>CVPR</sub>	100.8	1,764	X	X	×	1	13	30	pseudo
ZeroEGGS (Ghorbani et al., 2023)CGF	2	1	X	X	×	1	27	48	mo-cap
BEAT (Liu et al., 2022a) ECCV	35	30	1	X	1	1	27	48	mo-cap
TalkSHOW (Yi et al., 2023) CVPR	26.9				x	- X -	24 -	- 30 -	pseudo
BEAT2 (Liu et al., 2024a) <sub>CVPR</sub>	27	25	1	1	✓	1	24	30	mo-cap
GES-X (ours)	450	4,370	1	1	1	1	24	30	pseudo

175 176

177

178

179

181

182

183

184 185 186

187 188

189

162

> researchers employ the MoE paradigm to facilitate the multi-modal alignment tasks (Feng et al., 2023; Wang et al., 2023). Concretely, Shen et. al (Shen et al., 2023b) specifically investigates the scalability of MoE in vision-language models and showcases its potential to outperform dense models with equivalent computational cost. Regarding the human motion task, Liang et. al (Liang et al., 2024) propose a mixture-of-controllers mechanism that adaptively recognizes various ranges of the sub-motions with the text-token-specific experts, resulting in significant improvement on the text2motion research. Moreover, we notice that Mixture-of-Modality-Experts achieve promising performance in long-sequence modeling tasks (Liu et al., 2023b; Puigcerver et al., 2023; Shen et al., 2023a; Zhang et al., 2018). Motivated by this, we introduce Mixture-of-Gesture-Experts in our framework to enhance long-sequence gesture generation upon human speech guidance.

**PROPOSED METHOD** 3

#### PROBLEM FORMULATION 3.1

190 With the specifically designed generation framework, our goal is to synthesize vivid and diverse 3D 191 human gestures  $X = \{x_1, ..., x_N\}$  of the upper body through the given unseen continuous human 192 speech audio  $A = \{a_1, ..., a_N\}$  as input. Here, N denotes the number of the generated human postures coordinated with speech audio A. We leverage J joints with 3D representation to indicate 193 each pose  $x_i$ . Unlike the previous methods (Liu et al., 2024a; 2022a; Yi et al., 2023; Liu et al., 2024b) 194 that either utilize the text transcripts or speaker ID embedding as auxiliary input, our CoCoGesture 195 adopts only the human speech as model inputs. It should be noted this single modality input fashion 196 significantly facilitates the unseen speech-conditioned co-speech gesture generation. Our overall 197 workflow is displayed in Figure 3. 198

199 200

#### GESTURE DIFFUSION MODEL PRE-TRAINING 3.2

201 Large-scale Co-speech Gesture Dataset. To ensure the generalization of our pre-trained trans-202 former diffusion model, we newly collect a large-scale high-quality 3D meshed whole-body co-speech 203 gesture dataset, dubbed GES-X. In particular, we first leverage the advanced 3D pose estimator Pymaf-204 X (Zhang et al., 2023a) to obtain the meshed whole-body parameters upon SMPL-X (Pavlakos et al., 205 2019). The original raw data is collected from about 4.3K talk show videos including different stances 206 (*i.e.*, standing or sitting). After data processing<sup>1</sup>, our GES-X dataset contains more than 40M gesture 207 frames. To the best of our knowledge, this is the largest-scale whole-body meshed 3D co-speech gesture dataset, whose duration is 15x the current largest one, as reported in Table 1. 208

209 Specifically, the acquired human postures are represented as the unified standard SMPL (Loper 210 et al., 2023) body model accompanied by the MANO (Boukhayma et al., 2019) hand model. The 211 facial expression is presented in FLAME (Li et al., 2017) face model. Meanwhile, we leverage the 212 powerful speech recognition model WhisperX (Bain et al., 2023) to gain accurate word-level text 213 transcripts and linguistics phoneme (Studdert-Kennedy, 1987) aligned with the extracted motion 214 dynamics. In this manner, our GES-X not only facilitates the research on co-speech gesture generation

<sup>215</sup> 

<sup>&</sup>lt;sup>1</sup>Please refer to supplementary material for more details.



Figure 3: The overview of our CoCoGesture. In the **Pre-training**, we first pre-train a large unconditional diffusion model upon our large-scale GES-X dataset as the gesture expert. The **Finetuning** stage incorporates audio signal as gesture generation guidance. In the **Inference** stage, our CoCoGesture can generate vivid and diverse 3D co-speech gestures from unseen zero-shot human speeches.

but also supports various other human avatar creation tasks, *e.g.*, talking face (Tian et al., 2024), human behavior analysis (Qi et al., 2023b). Along with this large-scale dataset, the pretraining of the unconditional diffusion model is greatly enhanced with generalization and vividness.

240 Model Scaling-up & Pre-training Inspired by (Guo et al., 2022; Liang et al., 2024), we formulate 241 the popular diffusion transformer (DiT (Peebles & Xie, 2023)) as our model backbone owing to the 242 scalability and excellent compatibility of large-scale training data. Here, similar to the foundation model stable diffusion (Rombach et al., 2022), we scale up the original DiT from 120M to 1B with 243 different layers and latent dimensions, enabling learning massive gesture features so as to apply to 244 different downstream applications. During training, we enforce our denoiser to produce continuous 245 human motions given the diffusion time step t and noised postures  $x^{t}$ . The denoising processing is 246 constrained by the simple objective: 247

$$\mathcal{L}_{simple} = \mathbb{E}_{\mathbf{x},t,\epsilon} \left[ \left\| \mathbf{x} - \mathcal{D}_u(\mathbf{x}^t, t) \right\|_2^2 \right],\tag{1}$$

where  $\mathcal{D}_u$  is our unconditional denoiser,  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  is the added random Gaussian noise,  $\mathbf{x}^t = \mathbf{x} + \sigma_t \epsilon$  is the gradually noise adding process at step t.  $\sigma_t \in (0, 1)$  is the constant hper-parameter. Moreover, we follow the setting of (Tevet et al., 2022b; Guo et al., 2022) to leverage the velocity loss  $\mathcal{L}_{vel}$  and foot contact loss  $\mathcal{L}_{foot}$  for improving generated results more smoothness and physically reasonable. To this end, the overall objective is

$$\mathcal{L}_{total} = \lambda_{simple} \mathcal{L}_{simple} + \mathcal{L}_{vel} + \mathcal{L}_{foot}, \tag{2}$$

where  $\lambda_{simple}$  is trade-off weight coefficients.

#### 258 259 3.3 AUDIO CONTROLNET FINETUNE

232

233

234

235 236

237

238

239

248

249

255 256

257

267

In the finetuning phase, we intend to incorporate the audio condition *A* into the pre-trained gesture model. Inspired by text2image ControlNet (Zhang et al., 2023b), we introduce an audio ControlNet consisting of the trainable copy of the unconditional diffusion model and a novel proposed Mixtureof-Gesture-Experts (**MoGE**) block, as shown in Figure 4. The frozen pre-trained model serves as a strong gesture expert and the MoGE blocks follow a trainable copy to produce the temporally coordinated joint embedding of the audio signal and gesture features. Then the joint embedding is adaptively added to the denoised motion features of the next layer through a novel routing mechanism.

Mixture-of-Gesture-Experts. Inspired by MoE (Zhu et al., 2024; Yu et al., 2024; Shazeer et al., 2017), the key insight of the MoGE is adaptively fusing the information from the gesture expert (*i.e.*, pre-trained model) and the speech audio expert (*i.e.*, audio encoder), thereby the generated gestures

preserving temporal consistent with speech rhythms. To enhance the sequence-aware correspondence of the fused features, we first leverage the audio embedding  $f^a$  as the query Q to match the key feature K and values features V belonging motion embeddings  $f_l^{x''}$  via cross-attention mechanism:

$$Q_l = \mathbf{f}^a \mathbf{W}_l, K_l = \mathbf{f}_l^{x''} \mathbf{W}_l, V_l = \mathbf{f}_l^{x''} \mathbf{W}_l.$$
(3)

Here, l represents the index of each attention layer, and W denotes the projection matrix. Once we obtain these fused trainable features

277  $\mathbf{f}^{train}$ , we adopt an adaptive instance normal-278 ization (Ada-IN) layer conditioned on audio fea-279 tures to further boost  $\mathbf{f}^{train}$ . Then, we utilize a  $f_{l}^{x}$ 280 learnable routing adaptor to combine the output 281 of the gesture expert and trainable copy branch. 282 To be specific, we leverage the output of the frozen original last layer as motion guidance 283 representation to indicate the soft weight. By 284 doing so, we derive the blending process as fol-285 lows 286

$$\begin{aligned} \mathbf{f}_{l+1}^{x} &= \mathbf{R}_{l} \odot \mathbf{f}_{l}^{x'} + (1 - \mathbf{R}_{l}) \odot \mathbf{f}_{l}^{train}, \\ \mathbf{R}_{l} &= Softmax(\mathbf{W}_{R,l} \otimes \mathbf{f}_{l}^{x}), \end{aligned}$$
(4)

where **R** is the learnable router,  $\mathbf{W}_{R,l}$  denotes the weight matrix,  $\odot$  indicates the Hadamard product and  $\otimes$  indicates matrix multiplication. Afterward, we exploit the zero-initialized convolution layers to ensure the audio condition in



Figure 4: Details of our proposed Mixture-of-Gesture-Experts (MoGE) block. The pre-trained transformer layer is frozen and serves as the gesture expert, while the audio embedding is extracted from the audio expert.

the trainable copy branch cannot be impacted by the harmful noise.

**Training and Inference.** During the training, we leverage the same loss function in Eq. 2 to constrain the trainable conditional denoiser parameters. In the inference, we utilize the classifier-free guidance unconditional denoiser and audio-conditioned one  $D_a$ :

$$\hat{\mathbf{x}}^{(0)} = s \cdot \mathcal{D}_a(\mathbf{x}^{(t)}, t, a) + (1 - s) \cdot \mathcal{D}_u(\mathbf{x}^{(t)}, t),$$
(5)

where  $\hat{x}^{(0)}$  denotes the denoised gesture motions, and s is the set as 4.0 in practice.

#### 4 EXPERIMENTS

#### 4.1 EXPERIMENTAL SETTING AND DATASETS

**Implementation Details.** In the pretraining phase, we set  $\lambda_{simple} = 10$ , empirically. The total diffusion time step is 1,000 with the cosine noisy schedule (Nichol & Dhariwal, 2021). The initial learning rate is set as  $1 \times 10^{-4}$  with AdamW optimizer. Our model is trained on 8 NVIDIA H800 GPUs with a batch size of 256. The total training process takes 100 epochs, accounting for one week of the largest model version within 1B parameters. We provide three-version models with different architectures and parameters to explore the dependence of performance on model size.

313 During the finetuning stage, the audio signal is processed to mel-spectrograms with FFT window size 314 1,024, and hop length 512. Similar to (Liu et al., 2022b; Qi et al., 2023a; 2024), we take an advanced 315 speech recognizer (Chung et al., 2020) as the audio encoder. We train the audio ControlNet with a 316 batch size of 128 for 100 epochs. The initial learning rate is set as  $1 \times 10^{-5}$ . We take the DDIM (Song 317 et al., 2020) sampling strategy within 25 denoising timesteps during inference. Temporally, our 318 CoCoGesture synthesizes the 10-second gesture motions including 43 upper joints (*i.e.*13 body joints 319 + 30 hand joints) in practice. Each joint is converted to the 6D rotation representation (Zhou et al.) for better modeling in the experiments. 320

321

274

287 288 289

296

297

298

299 300 301

302 303

304 305

306

**GES-X Dataset.** We newly propose a large-scale co-speech gesture dataset, dubbed GES-X, to train our unconditional diffusion model. Firstly, we leverage 16 NVIDIA RTX 4090 GPUs to extract the 3D human poses from downloaded in-the-wild 4, 370 talk show videos. This process takes *more than* 

		BEAT	2 (Liu et al., 2	024a)	TalkSHOW (Yi et al., 2023) (zero-shot)			
	Methods	$FGD\downarrow$	Diversity $\uparrow$	$\mathrm{BA}\uparrow$	$FGD\downarrow$	Diversity $\uparrow$	$\mathbf{BA}\uparrow$	
1	Trimodal (Yoon et al., $2020$ ) <sub>TOG</sub>	13.05	33.54	0.75	-	-	-	
	HA2G (Liu et al., $2022b$ ) <sub>CVPR</sub>	9.37	45.81	0.76	15.25	58.41	0.65	
	CAMN (Liu et al., 2022a) <sub>ECCV</sub>	7.12	44.02	0.82	-	-	-	
	TalkSHOW (Yi et al., 2023) <sub>CVPR</sub>	10.59	45.23	0.79	16.41	57.30	0.64	
	DiffuGesture (Zhu et al., 2023) <sub>CVPR</sub>	11.82	48.53	0.81	17.03	50.52	0.72	
	ProbTalk (Liu et al., 2024b) <sub>CVPR</sub>	6.06	66.03	0.82	11.18	65.95	0.78	
	EMAGE (Liu et al., $2024a$ ) <sub>CVPR</sub>	4.09	69.70	0.85	-	-	-	
	CoCoGesture (ours)	3.92	70.47	0.87	9.62	69.10	0.83	

Table 2: Comparison with the state-of-the-art counterparts on BEAT2 and TalkSHOW datasets. ↑
means the higher the better, and ↓ indicates the lower the better. "-" denotes that the method cannot
be applied to the TalkSHOW dataset due to the lack of text transcripts. The term "zero-shot" implies
that the dataset contains unseen human voices.

one month, acquiring more than 88 million raw frames. After filtering the unreasonable gestures,
 we obtain 40 million high-quality postures. Then, we resample the FPS as 15, thereby the total
 generated gesture frames are 150 in a sequence. Finally, we obtain the 100, 162 motion clips with
 corresponding audio/text transcripts/phonemes.

**BEAT2 and TalkSHOW Datasets.** To fully verify the generalization and effectiveness of our 344 pr-trained model, we adopt two meshed datasets BEAT2 (Liu et al., 2024a) and TalkSHOW (Yi et al., 345 2023) in the evaluation phases. BEAT2 contains 3D meshed whole-body postures with multi-modality 346 information such as speaker ID and text transcripts. The content of the speech is based on 25 speakers' 347 answers to predefined questions. All the instances in BEAT2 are standing postures collected by the 348 motion-capture system. In the TalkSHOW dataset, only sitting postures with 4 speakers are collected 349 by 3D pose estimator from in-the-wild talk show videos. It is noted that the TalkSHOW dataset does 350 not provide text transcript annotation. 351

**Evaluation Metrics.** To fully evaluate the realism and diversity of the generated co-speech gestures, we introduce various metrics:

- **FGD**: Fréchet Gesture Distance (FGD) (Yoon et al., 2020) is leveraged to measure the distribution distance between the motions of real ones and generated ones.
- **BA**: Beat Alignment Score (BA) (Liu et al., 2022a;b) measures whether the generated human motions are rhythmically aligned with the speech beat.
- **Diversity**: Similar to (Liu et al., 2022b; Zhu et al., 2023; Qi et al., 2024), the same feature extractor is exploited to acquire feature embeddings of the synthesized gestures. We leverage the average distance between 500 randomly assembled pairs to indicate the diversity score.
- 361 362 363

324

338

343

352

353

354

355

356

357

358

359

360

#### 4.2 QUANTITATIVE RESULTS

364 **Comparisons with the State-of-the-art.** To fully verify the effectiveness of our method, we 365 compare our CoCoGesture framework with various state-of-the-art counterparts: Trimodal (Yoon 366 et al., 2020), HA2G (Liu et al., 2022b), CAMN (Liu et al., 2022a), TalkSHOW (Yi et al., 2023), 367 DiffuGesture (Zhu et al., 2023), ProbTalk (Liu et al., 2024b) and EMAGE (Liu et al., 2024a). For 368 a fair comparison, all the models are implemented by the source code released by the authors. We adopt GES-X in the finetuning stage to train our audio ControlNet. Then, we exploit both BEAT2 369 and TalkSHOW as testing sets. As for all the other counterparts, we adopt only the BEAT2 as the 370 training set. The TalkSHOW serves as the out-of-domain testing dataset, measuring the comparison 371 of the zero-shot ability. Since the TalkSHOW dataset does not provide the text transcript, it cannot be 372 used by some competitors (Yoon et al., 2020; Liu et al., 2022a; 2024a) that rely on text. 373

As reported in Table 2, our framework achieves the best results on both datasets. We observe that
both EMAGE and ours generate high-quality results in the FGD metric on the BEAT2 dataset.
However, different from EMAGE trained on BEAT2, our CoCoGesture is directly tested on this
dataset. Meanwhile, since our method only depends on the audio signal input, we can easily apply it
to another dataset. In terms of diversity score, our classifier-free inference strategy enables diverse

			-	-	<u> </u>		-	
Model	niquers	dmodel	nheads	dheads	Parms	BEAT	2 (Liu et al., 20	)24a)
	-lugers	-mouel	- neuus	-neuus		$\mathrm{FGD}\downarrow$	Diversity $\uparrow$	$\mathrm{BA}\uparrow$
CoCoGesture-Base	25	512	8	128	120M	6.00	52.73	0.81
CoCoGesture-Medium	25	1024	16	128	480M	4.96	57.75	0.83
CoCoGesture-Large ‡	50	1024	16	128	1B	4.30	68.33	0.85
CoCoGesture-Large	50	1024	16	128	1B	3.92	70.47	0.87

Table 3: Ablation study on model scale and pre-training setting. ‡ denotes without pre-training stage.

gestures while preserving the authority and vividness of the results. Considering the zero-shot inference, our approach outperforms all the counterparts by a large margin. Remarkably, on the TalkSHOW dataset, our CoCoGesture reduces FGD by a significant amount of 16.22% over the sub-optimal counterparts. The better performance demonstrates our model's superior generalization ability, verifying our insight on pre-training and finetune strategy.

Ablation Study. To further evaluate the effectiveness of our proposed framework, we conduct a series of ablation studies of different components and training strategies as variations.

**Effects on Model scale & Pre-training:** To investigate the impact of the model scale and pre-training stage, we conduct the ablation study on the BEAT2 dataset, as reported in Table 3. We design three model variants with different architectures. Here,  $n_{layers}$  is the total transformer layers,  $d_{model}$ denotes dimension of latent vectors,  $n_{heads}$  means number of attention heads,  $d_{heads}$  indicates the dimension of each attention head. It is observed that our model performance is gradually improved with model scaling up. This aligns our insight on larger models to learn massive gesture manifold. It is noticed that without pre-training, the model achieves lower performance. This suggests that pre-training on our GES-X dataset is effective in improving model generalization ability.

Effects of the MoGE Block: To fully analyze the effectiveness of our proposed Mixture-of Gesture-Experts (MoGE) blocks, we conduct the ablation study through detailed components.
 As reported in Table 4, we demon-

405	As reported in Table 4, we demon-
406	strate the exclusion of cross-attention
407	and routing mechanisms respectively
408	from our full large model version
409	leads to performance degradation. To
410	be specific, the cross-attention module
/11	effectively models the dependency of
410	audio signals with generated results,
412	thus implementation without it leads
413	to worse performance in all the met-

|--|

	BEAT2 (Liu et al., 2024a)						
Methods	$FGD\downarrow$	Diversity $\uparrow$	$BA\uparrow$				
w/o Cross-attn	4.79	62.48	0.86				
w/o Routing	4.28	67.14	0.79				
CoCoGesture (full)	3.92	70.47	0.87				

rics. Meanwhile, the exclusion of the routing mechanisms results in an obvious decrease in the BA
score. This demonstrates that our routing mechanism significantly enhances the temporal coherency
between the audio embeddings *w.r.t.* gesture features, thus producing vivid and coherency gestures.

417

419

378

388

389

390

391

392

#### 418 4.3 QUALITATIVE EVALUATION

**Visualization:** To fully demonstrate the superior performance of our CoCoGesture framework, we 420 show the visualized key frames synthesized by ours compared with various counterparts on BEAT2 421 and TalkSHOW datasets, respectively. As shown in Figure 5, our method displays vivid and diverse 422 gestures against others. In particular, we observe the Trimodal tends to synthesize unreasonable 423 and stiff results (e.g., the red rectangle in the BEAT2 dataset). Although the HA2G and EMAGE 424 can generate the natural upper body postures, we find that their body movements are of limited 425 dynamics (e.g., the blue rectangle in the BEAT2 dataset). In terms of the zero-shot inference in the 426 TalkSHOW dataset, both DiffuGesture and our method produce reasonable gestures. However, the 427 results generated by DiffuGesture are misaligned with the input audio. This may be caused by the 428 limited word corpus of the BEAT2 dataset restricting the generalization of the model. In contrast, our 429 method can synthesize the vivid and synchronous co-speech gestures (e.g., the arms become lifting while the hands stretch out). This highly aligns with our motivation about the model generalization 430 improved by pre-training on our large-scale dataset GES-X. For more demo results please refer to our 431 anonymous website: https://anonymous.4open.science/w/GES-X/.



Figure 5: Visualization of our generated 3D co-speech gestures against various state-of-the-art methods. The samples on the left are from BEAT2, and the samples on the right are from TalkSHOW.

**User Study:** To further analyze the quality of results synthesized by various counterparts and ours, we conduct a user study by inviting 15 volunteers. The statistical mean results are reported in Figure 6. All the volunteers are recruited anonymously from schools with different majors. Each participant is required to rate the randomly selected visualization videos from 0 (worst) to 5 (best) in terms of naturalness, smoothness, and speech-gesture coherency. Our CoCoGesture framework demonstrates the best performance among all the competitors.

Especially, in terms of smoothness and speechgesture coherency, our method outperforms others with noticeable improvements, verifying the effectiveness of our Mixture-of-Gesture-Expert.

#### 5 CONCULSION

454

455 456

457

458

459

460

461

462

463

464

465

466

467

468

469 In this paper, we propose CoCoGesture to gen-470 erate vivid and diverse co-speech 3D gestures 471 from in-the-wild zero-shot human speech. To fulfill this goal, we first newly collect a large-472 scale dataset that contains more than 40M high-473 quality 3D meshed postures across 4.3K speak-474 ers from in-the-wild talk show videos. Along 475 with this dataset, we pre-train a large generaliz-476 able diffusion model to be our gesture expert in



Figure 6: User study on gesture naturalness, motion smoothness, and speech-gesture coherency.

the first stage. To incorporate human speech as guidance, we further propose a novel audio ControlNet
 that adaptively fuses the audio embeddings and the motion clues from the pre-trained gesture expert.
 Extensive experiments conducted on two out-of-domain datasets show the superiority of our model.

Limitation: Our framework only takes the audio signal as model input to generate gestures. It might be possible that our model produces emotionally insensitive cases (*e.g.*, moving faster or more intensely when angry or happy). Meanwhile, the automated pose extraction and speech techniques may have an impact on the datasets we newly collect, despite the huge effort we put into data clean filtering and processing. In future works, we will incorporate our model with emotional conditions and investigate more stable data processing techniques to improve the quality of generated gestures.

# 486 REFERENCES

487	
488	Chaitanya Ahuja, Dong Won Lee, Yukiko I Nakano, and Louis-Philippe Morency. Style transfer
489	for co-speech gesture animation: A multi-speaker conditional-mixture approach. In Computer
490	Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings,
491	Part XVIII 16, pp. 248–265. Springer, 2020.
492	Tenglong Ao Oingzhe Gao Yuke Lou Baoguan Chen and Libin Liu. Rhythmic gesticulator:
493	Rhythm-aware co-speech gesture synthesis with hierarchical neural embeddings. ACM Transac-
494	tions on Graphics (TOG), 41(6):1–19, 2022.
495	······································
496	Tenglong Ao, Zeyi Zhang, and Libin Liu. Gesturediffuclip: Gesture diffusion model with clip latents,
497	2023.
498	Jose Dedre Areuie, Jiemen Li, Kerthik Vetrivel, Dishi Agerwal, Deepek Conjugth, Jiejun Wu
499	Alexander Clegg, and C. Karen Liu. Circle: Canture in rich contextual environments 2023
500	Alexander Clegg, and C. Karen Eld. Chele. Capture in nen contextuar environments, 2025.
501	Max Bain, Jaesung Huh, Tengda Han, and Andrew Zisserman. Whisperx: Time-accurate speech
502	transcription of long-form audio. INTERSPEECH 2023, 2023.
503	
504	Uttaran Bhattacharya, Elizabeth Childs, Nicholas Rewkowski, and Dinesh Manocha.
505	Speech 2 affective gestures: Synthesizing co-speech gestures with generative adversarial affec-
506	np 2027_2036 2021
507	pp. 2027–2050, 2021.
508	Adnane Boukhayma, Rodrigo de Bem, and Philip HS Torr. 3d hand shape and pose from images in the
509	wild. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition,
510	pp. 10843–10852, 2019.
511	Henri Dardin annante andie 2.1 angelen dissingtion gineling, grinziale hensekanade and graine. In
512	24th INTEPSDEECH Conference (INTEPSDEECH 2022) pp 1082 1087 ISCA 2022
513	24 <i>in not EKSI ELCH Conjetence (not EKSI ELCH 2023)</i> , pp. 1703–1707. ISCA, 2023.
514	Justine Cassell, Catherine Pelachaud, Norman Badler, Mark Steedman, Brett Achorn, Tripp Becket,
515	Brett Douville, Scott Prevost, and Matthew Stone. Animated conversation: rule-based generation of
516	facial expression, gesture & spoken intonation for multiple conversational agents. In Proceedings
517	of the 21st annual conference on Computer graphics and interactive techniques, pp. 413–420,
518	1994.
519	Junming Chen Yunfei Liu, Jianan Wang, Ailing Zeng, Yu Li, and Oifeng Chen, Diffsheg: A
520	diffusion-based approach for real-time speech-driven holistic 3d expression and gesture generation.
521	In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2024.
522	
523	Xin Chen, Biao Jiang, Wen Liu, Zilong Huang, Bin Fu, Tao Chen, Jingyi Yu, and Gang Yu. Executing
524	your commands via motion diffusion in latent space, 2023.
525	Oingrong Cheng, Xu Li, and Xinghui Fu. Siggesture: Generalized co-speech gesture synthesis via
526	semantic injection with large-scale pre-training diffusion models. <i>arXiv preprint arXiv:2405.13336</i> .
527	2024.
528	
529	Joon Son Chung, Jaesung Huh, Seongkyu Mun, Minjae Lee, Hee-Soo Heo, Soyeon Choe, Chiheon
03U	Ham, Sunghwan Jung, Bong-Jin Lee, and Icksang Han. In defence of metric learning for speaker
500	recognition. Interspeech 2020, 2020.
502	Rishabh Dabral, Muhammad Hamza Mughal, Vladislav Golvanik, and Christian Theobalt, Mofusion:
504	A framework for denoising-diffusion-based motion synthesis, 2023.
534	
535	Maged Farouk. Studying human robot interaction and its characteristics. <i>International Journal of</i>
537	Computations, Information and Manufacturing (IJCIM), 2(1), 2022.
538	William Fedus Barret Zonh and Noam Shazeer Switch transformers: Scaling to trillion parameter
539	models with simple and efficient sparsity. <i>Journal of Machine Learning Research</i> . 23(120):1–39.
565	2022.

540 541 542 543	Zhida Feng, Zhenyu Zhang, Xintong Yu, Yewei Fang, Lanxin Li, Xuyi Chen, Yuxiang Lu, Jiaxi- ang Liu, Weichong Yin, Shikun Feng, et al. Ernie-vilg 2.0: Improving text-to-image diffusion model with knowledge-enhanced mixture-of-denoising-experts. In <i>Proceedings of the IEEE/CVF</i> <i>Conference on Computer Vision and Pattern Recognition</i> , pp. 10135–10145, 2023
544 545 546 547	<ul><li>Ylva Ferstl and Rachel McDonnell. Investigating the use of recurrent motion modelling for speech gesture generation. In <i>Proceedings of the 18th International Conference on Intelligent Virtual Agents</i>, pp. 93–98, 2018.</li></ul>
548 549 550	Yu Fu, Yan Hu, and Veronica Sundstedt. A systematic literature review of virtual, augmented, and mixed reality game applications in healthcare. <i>ACM Transactions on Computing for Healthcare (HEALTH)</i> , 3(2):1–27, 2022.
551 552 553	Trevor Gale, Deepak Narayanan, Cliff Young, and Matei Zaharia. Megablocks: Efficient sparse training with mixture-of-experts. <i>Proceedings of Machine Learning and Systems</i> , 5, 2023.
554 555 556	Zhifu Gao, Zerui Li, Jiaming Wang, Haoneng Luo, Xian Shi, Mengzhe Chen, Yabin Li, Lingyun Zuo, Zhihao Du, Zhangyu Xiao, et al. Funasr: A fundamental end-to-end speech recognition toolkit. <i>arXiv preprint arXiv:2305.11013</i> , 2023.
557 558 559	Saeed Ghorbani, Ylva Fe2rstl, Daniel Holden, Nikolaus F Troje, and Marc-André Carbonneau. Zeroeggs: Zero-shot example-based gesture generation from speech. In <i>Computer Graphics</i> <i>Forum</i> , volume 42, pp. 206–216. Wiley Online Library, 2023.
560 561 562 563	Chuan Guo, Shihao Zou, Xinxin Zuo, Sen Wang, Wei Ji, Xingyu Li, and Li Cheng. Generating diverse and natural 3d human motions from text. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 5152–5161, 2022.
564 565 566 567	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 <i>Proceedings of the 21st ACM International Conference on Intelligent Virtual Agents</i> , pp. 101–108, 2021.
568 569 570	Chien-Ming Huang and Bilge Mutlu. Robot behavior toolkit: generating effective social behaviors for robots. In <i>Proceedings of the seventh annual ACM/IEEE international conference on Human-Robot Interaction</i> , pp. 25–32, 2012.
571 572 573	Siyuan Huang, Zan Wang, Puhao Li, Baoxiong Jia, Tengyu Liu, Yixin Zhu, Wei Liang, and Song- Chun Zhu. Diffusion-based generation, optimization, and planning in 3d scenes, 2023.
574 575	Robert A Jacobs, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. Adaptive mixtures of local experts. <i>Neural computation</i> , 3(1):79–87, 1991.
576 577 578	Hema Swetha Koppula and Ashutosh Saxena. Anticipating human activities for reactive robotic response. In <i>IROS</i> , pp. 2071. Tokyo, 2013.
579 580	Tianye Li, Timo Bolkart, Michael J Black, Hao Li, and Javier Romero. Learning a model of facial shape and expression from 4d scans. <i>ACM Trans. Graph.</i> , 36(6):194–1, 2017.
582 583 584	Zhihao Li, Jianzhuang Liu, Zhensong Zhang, Songcen Xu, and Youliang Yan. Cliff: Carrying location information in full frames into human pose and shape estimation. In <i>European Conference on Computer Vision</i> , pp. 590–606. Springer, 2022.
585 586 587	Han Liang, Jiacheng Bao, Ruichi Zhang, Sihan Ren, Yuecheng Xu, Sibei Yang, Xin Chen, Jingyi Yu, and Lan Xu. Omg: Towards open-vocabulary motion generation via mixture of controllers. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , 2024.
588 589 590 591	<pre>Junfan Lin, Jianlong Chang, Lingbo Liu, Guanbin Li, Liang Lin, Qi Tian, and Chang Wen Chen. OhMG: Zero-shot open-vocabulary human motion generation, 2023. URL https: //openreview.net/forum?id=41GL_ruft.</pre>
592 593	Chen Liu, Peike Patrick Li, Xingqun Qi, Hu Zhang, Lincheng Li, Dadong Wang, and Xin Yu. Audio- visual segmentation by exploring cross-modal mutual semantics. In <i>Proceedings of the 31st ACM</i> <i>International Conference on Multimedia</i> , pp. 7590–7598, 2023a.

616

623

594	Haiyang Liu, Zihao Zhu, Naoya Iwamoto, Yichen Peng, Zhengqing Li, You Zhou, Elif Bozkurt, and
595	Bo Zheng. Beat: A large-scale semantic and emotional multi-modal dataset for conversational
596	gestures synthesis. In European Conference on Computer Vision, pp. 612–630. Springer, 2022a.
597	

- 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. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024a.
- Juncai Liu, Jessie Hui Wang, and Yimin Jiang. Janus: A unified distributed training framework for
   sparse mixture-of-experts models. In *Proceedings of the ACM SIGCOMM 2023 Conference*, pp.
   486–498, 2023b.
- Kian 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10462–10472, 2022b.
- Yifei Liu, Qiong Cao, Yandong Wen, Huaiguang Jiang, and Changxing Ding. Towards variable and
   coordinated holistic co-speech motion generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024b.
- Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multi-person linear model. In *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*, pp. 851–866. 2023.
- Stacy Marsella, Yuyu Xu, Margaux Lhommet, Andrew Feng, Stefan Scherer, and Ari Shapiro. Virtual
   character performance from speech. In *Proceedings of the 12th ACM SIGGRAPH/Eurographics symposium on computer animation*, pp. 25–35, 2013.
- Michael McAuliffe, Michaela Socolof, Sarah Mihuc, Michael Wagner, and Morgan Sonderegger.
   Montreal forced aligner: Trainable text-speech alignment using kaldi. In *Interspeech*, volume 2017, pp. 498–502, 2017.
- Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models.
   In *International conference on machine learning*, pp. 8162–8171. PMLR, 2021.
- 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 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,
   pp. 10975–10985, 2019.
- William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4195–4205, 2023.
- Stefano Pini, Christian S Perone, Aayush Ahuja, Ana Sofia Rufino Ferreira, Moritz Niendorf, and
   Sergey Zagoruyko. Safe real-world autonomous driving by learning to predict and plan with a
   mixture of experts. In 2023 IEEE International Conference on Robotics and Automation (ICRA),
   pp. 10069–10075. IEEE, 2023.
- Isabella Poggi, Catherine Pelachaud, Fiorella de Rosis, Valeria Carofiglio, and Berardina De Carolis.
   Greta. a believable embodied conversational agent. *Multimodal intelligent information presentation*, pp. 3–25, 2005.
- Joan Puigcerver, Carlos Riquelme Ruiz, Basil Mustafa, and Neil Houlsby. From sparse to soft mixtures of experts. In *The Twelfth International Conference on Learning Representations*, 2023.
- Kingqun Qi, Chen Liu, Lincheng Li, Jie Hou, Haoran Xin, and Xin Yu. Emotiongesture: Audio-driven diverse emotional co-speech 3d gesture generation. *arXiv preprint arXiv:2305.18891*, 2023a.
- Kingqun Qi, Chen Liu, Muyi Sun, Lincheng Li, Changjie Fan, and Xin Yu. Diverse 3d hand gesture prediction from body dynamics by bilateral hand disentanglement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4616–4626, June 2023b.

648 649 650	Xingqun Qi, Chen Liu, Muyi Sun, Lincheng Li, Changjie Fan, and Xin Yu. Diverse 3d hand gesture prediction from body dynamics by bilateral hand disentanglement. In <i>Proceedings of the IEEE/CVF</i> Conference on Computer Vision and Pattern Recognition, pp. 4616–4626, 2023c
651 652	Xingqun Qi, Jiahao Pan, Peng Li, Ruibin Yuan, Xiaowei Chi, Mengfei Li, Wenhan Luo, Wei Xue,
653	Shanghang Zhang, Qifeng Liu, and Yike Guo. Weakly-supervised emotion transition learning
654	for diverse 3d co-speech gesture generation. In <i>Proceedings of the IEEE/CVF Conference on</i> Computer Vision and Pattern Recognition 2024
655	Computer vision and Fallern Recognition, 2024.
656 657	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever.
658	Learning transferable visual models from natural language supervision, 2021.
659	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever.
660 661	Robust speech recognition via large-scale weak supervision. In <i>International Conference on Machine Learning</i> , pp. 28492–28518. PMLR, 2023.
662	South Dood Zayman Alexa, Vinchan Van Lajanugan Lagaguaran Damt Sabiala and Hanglak Lag
663 664	Generative adversarial text to image synthesis, 2016.
665	Robin Rombach Andreas Blattmann Dominik Lorenz Patrick Esser and Biorn Ommer High-
666	resolution image synthesis with latent diffusion models. In <i>Proceedings of the IEEE/CVF confer-</i>
667	ence on computer vision and pattern recognition, pp. 10684–10695, 2022.
668	New Glassia Asti Mitasia' Kasa (CM, 's Ast De 's O ast Cooffee History of
669	Noam Snazeer, Azalia Mirnoseini, Krzysztor Maziarz, Andy Davis, Quoc Le, Geomrey Hinton, and Leff Dean. Outrageously large neural networks: The sporsely geted mixture of experts layer. <i>arXiv</i>
670	prenrint arXiv:1701.06538, 2017
671	preprint drXiv.1701.00550, 2017.
672	Sheng Shen, Le Hou, Yanqi Zhou, Nan Du, Shayne Longpre, Jason Wei, Hyung Won Chung, Bar
673	Zoph, William Fedus, Xinyun Chen, et al. Mixture-of-experts meets instruction tuning: A winning
674	combination for large language models. In <i>The Twelfth International Conference on Learning</i>
675	<i>Representations</i> , 2023a.
676	Sheng Shen, Zhewei Yao, Chunyuan Li, Trevor Darrell, Kurt Keutzer, and Yuxiong He. Scaling
677	vision-language models with sparse mixture of experts. In The 2023 Conference on Empirical
670	Methods in Natural Language Processing, 2023b.
620	Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. arXiv
681	preprint arXiv:2010.02502, 2020.
682	Michael Studdert-Kennedy. The phoneme as a perceptuomotor structure. <i>Haskins Laboratories:</i>
683	Status Report on Speech Research, SR, 91:45–57, 1987.
695	Guy Tevet, Brian Gordon, Amir Hertz, Amit H. Bermano, and Daniel Cohen-Or. Motionclip:
686 686	Exposing human motion generation to clip space, 2022a.
687	Guy Tevet, Sigal Raab, Brian Gordon, Yoni Shafir, Daniel Cohen-or, and Amit Haim Bermano.
688	Human motion diffusion model. In The Eleventh International Conference on Learning Represen-
689	tations, 2022b.
690	Lingui Tian Oi Wang Bang Zhang and Liefeng Bo. Emo: Emote portrait alive - generating
691	expressive portrait videos with audio2video diffusion model under weak conditions, 2024.
692	
093	Jonathan Tseng, Rodrigo Castellon, and C. Karen Liu. Edge: Editable dance generation from music,
094	2022.
606	Mingjie Wang, Hao Cai, Yong Dai, and Minglun Gong. Dynamic mixture of counter network
607	for location-agnostic crowd counting. In Proceedings of the IEEE/CVF winter conference on
608	applications of computer vision, pp. 167–177, 2023.
699	Sicheng Yang, Zhiyong Wu, Minglei Li, Zhensong Zhang, Lei Hao, Weihong Rao, Ming Chang
700	and Long Xiao. Diffusestylegesture: stylized audio-driven co-speech gesture generation with
701	diffusion models. In <i>Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence</i> , pp. 5860–5868, 2023.

702 703 704 705	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 IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 469–480, 2023.
706 707 708	Youngwoo Yoon, Woo-Ri Ko, Minsu Jang, Jaeyeon Lee, Jaehong Kim, and Geehyuk Lee. Robots learn social skills: End-to-end learning of co-speech gesture generation for humanoid robots. In 2019 International Conference on Robotics and Automation (ICRA), pp. 4303–4309. IEEE, 2019.
709 710 711	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 (TOG)</i> , 39(6):1–16, 2020.
712 713 714 715	Jiazuo Yu, Yunzhi Zhuge, Lu Zhang, Ping Hu, Dong Wang, Huchuan Lu, and You He. Boosting continual learning of vision-language models via mixture-of-experts adapters. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 23219–23230, 2024.
716 717 718	Ailing Zeng, Lei Yang, Xuan Ju, Jiefeng Li, Jianyi Wang, and Qiang Xu. Smoothnet: A plug-and-play network for refining human poses in videos. In <i>European Conference on Computer Vision</i> , pp. 625–642. Springer, 2022.
719 720 721	He Zhang, Sebastian Starke, Taku Komura, and Jun Saito. Mode-adaptive neural networks for quadruped motion control. <i>ACM Transactions on Graphics (TOG)</i> , 37(4):1–11, 2018.
722 723 724	Hongwen Zhang, Yating Tian, Yuxiang Zhang, Mengcheng Li, Liang An, Zhenan Sun, and Yebin Liu. Pymaf-x: Towards well-aligned full-body model regression from monocular images. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 2023a.
725 726 727	Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pp. 3836–3847, 2023b.
728 729 730 731	Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation representations in neural networks. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 5745–5753.
732 733 734	Lingting Zhu, Xian Liu, Xuanyu Liu, Rui Qian, Ziwei Liu, and Lequan Yu. Taming diffusion models for audio-driven co-speech gesture generation. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 10544–10553, 2023.
735 736 737 738	Pengfei Zhu, Yang Sun, Bing Cao, and Qinghua Hu. Task-customized mixture of adapters for general image fusion. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 7099–7108, 2024.
739	
740	
741 742	
743	
744	
745	
746	
747	
748	
749	
750	
751	
752	
753	
754	
755	



Figure 7: The overall workflow of our dataset construction. The talk show videos are processed to obtain high-quality postures through advanced automatic technologies and expert proofreading.

## A SUPPLEMENTARY MATERIAL

To demonstrate the effectiveness of our data construction techniques and the proposed method of coherent co-speech gesture generation, we further elaborate on the detailed data synthesis and vision perception in the supplementary material.

A.1 DATASET

776

777 778 779

781

782

783

784 785

786

### 787 A.1.1 CONSTRUCTION OF OUR GES-X

In this section, we detail the overall pipeline for creating GES-X, a large-scale dataset that contains over 40M co-speech gesture frames. The whole procedure consists of four folds: internet video collection, motion annotation, post-processing, and manual inspection, as summarised in Figure 7.

**Internet Videos Collection (Step 1**&2): Acquiring the paired speech-gesture 3D data via motion 792 capture system is expensive and labor-consuming. Consequently, some previous works (Liu et al., 793 2022b; Yoon et al., 2020; 2019; Qi et al., 2023a; 2024; Yi et al., 2023) leverage in-the-wild talk show 794 videos as the source to extract 3D postures via advanced pose estimator. Following this fashion, we 795 intend to obtain large-scale co-speech 3D gestures from YouTube talk show videos covering diverse 796 topics and speaker styles. We obtain 4,370 videos and their corresponding text transcripts. Given the 797 substantial volume of our video data, we employ PySceneDetect to segment lengthy videos into clips. 798 YOLOv8 is also used for human detection, discarding clips that do not show a person within the first 799 30 frames. These processes allow us to obtain potential clips containing speakers, with an average 800 duration of 9.85 seconds of each.

801 Motion Annotation (Step 3&4): Here, we employ SMPL-X (Pavlakos et al., 2019) to represent 802 whole-body poses, a widely 3D human representation standard adopted in various downstream tasks. 803 Then, we exploit the advanced pose estimator PyMAF-X (Zhang et al., 2023a) to extract high-quality 804 3D postures including body poses, subtle fingers, shapes, and expressions of the speakers. For 805 audio processing, we use FunASR (Gao et al., 2023) with the Whisper-large-v3 model to generate 806 transcripts. We then apply eight criteria to filter the clips and motion annotations: *clips that are too* 807 short, contain multiple people, involve looking back or sideways, have missing joints, show small or static individuals, or briefly miss the speakers. Additionally, transcripts with fewer than five words 808 are discarded, though the corresponding video clips are retained to increase the data scale for certain 809 audio-to-gesture tasks.

810 **Post-Processing (Step 5**&**6):** Once we obtain a large amount of raw pose sequences, we conduct 811 the post-processing to boost the quality of our data. Specifically, we visualize the motion sequences 812 with render mesh vertices and observe there are some temporal jittering issues. These jitters usually 813 result from heavy occlusion, truncation, and motion blur caused by changes in camera angles and 814 large-scale human movements of speakers. To address this, similar to CLIFF (Li et al., 2022), we utilize SmoothNet (Zeng et al., 2022) for temporal smoothing and jitter motion refinement. In practice, 815 through manual review, we notice that SmoothNet effectively produces cleaner and more reliable 816 motion sequences without sacrificing the diversity of postures. Despite that, given the frequent 817 extreme variations in camera angles, speaker poses, and lighting in talk show videos, some inaccurate 818 pose estimations from PyMAF-X are inevitable. Therefore, we leverage an automatic abnormal pose 819 detection method to further improve the pose quality. By representing the arm poses as Euler angles 820 using the x, y, and z convention, based on findings from (Pavlakos et al., 2019), we focus particularly 821 on the poses of the wrists. Once the wrist poses exceed 150 degrees on any axis or if the pose changes 822 by more than 25 degrees between adjacent frames (at 15 fps), we discard these abnormal postures 823 surrounding 150 frames. 824

Manual Inspection (Step 7): Finally, we perform the manual review for the processed clips with 825 a uniform ratio of 10:1. In particular, we follow the order of scenecut and sample one clip from 826 every ten groups of clips. Since these 10 clips typically originate from the same video, making 827 this assumption reasonably valid. For all clips, we divide them into ten groups for ten inspectors to 828 manually review. These inspectors evaluate the visualizations based on obtained SMPL-X parameters 829 to determine whether they are smooth, jittering, or abnormal. If the motion sequences appear jittering 830 or abnormal, the entire group of ten clips from which the sample originated is discarded. Through 831 meticulous evaluation and significant effort, the quality of our GES-X is greatly ensured. 832

833 **Text Transcript and Phonme Alignment:** To acquire accurate semantic annotations from speech, we transcribe audio files to extract text, phonemes, and their corresponding timestamps. Specifically, 834 we utilize WhisperX (Bain et al., 2023) as our transcription tool, which employs pyannote (Bredin, 835 2023) for speaker diarization and the Whisper (Radford et al., 2023) model for automatic speech 836 recognition (ASR). This tool incorporates a VAD Cut & Merge strategy to address the issue of 837 inaccurate timestamp predictions in long audio. We configure the system to recognize only one 838 speaker and utilize the Whisper Large V3 model for ASR. This approach splits long audio into 839 segments, each with its corresponding text. Subsequently, all data and labels are manually reviewed 840 by skilled human annotators. Finally, we apply the verified transcriptions and segment results to 841 perform Forced Phoneme Alignment using the Montreal Forced Aligner (McAuliffe et al., 2017) to 842 accurately label all phonemes and their respective timestamps.

844 A.1.2 BEAT2 & TALKSHOW DATASETS

Similar to our GES-X, we first resample the BEAT2 and TalkSHOW datasets with the FPS 15. Then, we divide datasets into 10s clips. Finally, we obtain 35, 758 clips in BEAT2 and 9, 629 in TalkSHOW. We follow the convention of (Liu et al., 2024a) to split the train/validation/test with the proportion of 85%, 7.5%, and 7.5% of both datasets.

850 A.2 Additional Experiments851

# 852 A.2.1 METRIC CALCULATION DETAILS

Inspired by (Yoon et al., 2020; Liu et al., 2022b), we leverage the FGD to evaluate whether the generated gestures preserve realism with the ground truth in the perceptive of distribution. We first pre-train an auto-encoder as the feature extractor. Then the FGD is calculated among the latent vectors belonging to sequential prediction and ground truth, respectively. The dimension of the latent vector is 128, similar to (Yoon et al., 2020; Liu et al., 2022b).

859 860

843

A.2.2 DISCUSSION OF EXPERIMENTAL SETTING

In our experiments, we only take human audio as a condition to guide the gesture generation. Although
 current speech-to-text methods can provide high-quality results, it requires an additional module to
 obtain word-level transcripts with accurate timestamps before modeling gestures from human speech.
 Meanwhile, during our pretraining phases, there are more than 4.3k speaker identities. In this fashion,

it is difficult to model the speaker's characteristics. In contrast, our method directly generates the gestures from speech signals. In this universal manner, our model is more practical in real sence applications (e.g., outdoor background noise may have a serious impact on speech-to-text). Therefore, similar to (Yi et al., 2023; Liu et al., 2024b; Zhu et al., 2023), our setting of directly generating gestures from speech audios without textual information is one of the common methodology streams in the community.

#### A.2.3 ADDITIONAL ABLATION RESULTS

We further conduct experiments to train our CoCoGesture on the BEAT2 dataset (denoted as CoCoGesture\*). Our method attains the best performance against all the counterparts, which highly demonstrates the effectiveness of our proposed CoCoGesture framework. Although the FGD of our framework

pre-trained on the GES-X dataset (denoted by †) is slightly worse than the
one trained on BEAT2 due to crossdataset evaluation, it still achieves better results than other competitors.

Methods	BEAT2 (Liu et al., 2024a)		
	$FGD\downarrow$	Diversity ↑	$\mathrm{BA}\uparrow$
CoCoGesture*	3.66	71.08	0.87
CoCoGesture <sup>†</sup>	3.92	70.47	0.87

881

870

871 872

873

874

875

#### 882 A.2.4 USER STUDY DETAILS

B83 During the user study, we utilize eight

models to randomly generate demo

885 videos in each of the BEAT2 and TalkSHOW datasets. For each method, we randomly generate 886 two demo videos from two datasets. For those that can be performed on the Talkshow dataset, the 887 generated results are guaranteed to come from both datasets. Therefore, each participant needs to respond to 16 samples from eight methods. Then, all the volunteer students are requested to rate all videos without any hint about which model produces this video. The higher score means the 889 better results. 5 points means that the video meets the audience's requirements perfectly. 0 points 890 indicates that the video is totally unacceptable. To ensure fairness, each demo video is played on a 891 PPT slide with a blank background. When all students have completed the grading, their results will 892 be collected anonymously and the average score will be calculated and announced. For each sample, 893 the participants are allowed to rate only after watching the entire video. To ensure that participants 894 will not have biased results due to recency bias, we invite participants to take the test at different 895 periods and not strictly limit the test duration. Participants can watch each video repeatedly. We 896 double-check the rating results by randomly selecting 60% of participants to redo the same test one 897 week later, and there are no significant changes to the final results. 898

#### 899 A.2.5 ADDITIONAL VISUALIZATION RESULTS

Here, we provide more visualized results of our CoCoGesture framework and other counterparts in the anonymous website: *https://anonymous.4open.science/w/GES-X/*. Moreover, to fully demonstrate the effectiveness of our proposed components and different model scales, we visualize the key frames of the generated results in Figure 8 and Figure 9.

905

900

- 906
- 907
- 908
- 909
- 910
- 911
- 912 913
- 913
- 914
- 915 916
- 917



969

Figure 9: Visual comparisons of ablation study on BEAT2. We show the key frames of the generated motions given the human speech. Best viewed on screen.