NOISE-ROBUST AUDIO-VISUAL SPEECH-DRIVEN BODY LANGUAGE SYNTHESIS

Anonymous authors

Paper under double-blind review

ABSTRACT

With the continuous advancement of video generation, researchers have achieved speech-driven body language synthesis, such as co-speech gestures. However, due to the lack of paired data for visual speech (i.e., lip movements) and body languages, existing methods typically rely solely on audio-only speech, which struggles to correctly synthesize target results in noisy environments. To overcome this limitation, we propose an Audio-Visual Speech-Driven Synthesis (AV-SDS) method tailored for body language synthesis, aiming for robust synthesis even under noisy conditions. Given that each body language modality data has its corresponding audio speech, AV-SDS adopts a two-stage synthesis framework based on speech discrete units, consisting of the AV-S2UM and Unit2X modules. It uses speech discrete units as carriers to construct a direct mapping from audiovisual speech to each body language. Considering the distinct characteristics of different body languages, AV-SDS can be implemented based on semantic and acoustic discrete units, respectively, to achieve high-semantic and high-rhythm body language synthesis. Experimental results demonstrate that our AV-SDS achieves superior performance in synthesizing multiple body language modalities in noisy environments, delivering noise-robust body language synthesis. For samples and further information, please visit demo page at https://av-sds.github. io/.

028 029

031

004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

1 INTRODUCTION

Recent years have witnessed great advancements in video generation (Tian et al., 2024; Richard et al., 2021), and researchers have successfully achieved various forms of speech-driven body language synthesis (Liu et al., 2023), including talking head generation (Prajwal et al., 2020; 2022), co-gesture generation (Liu et al., 2022b;a; Yang et al., 2023b), 3D facial animation (Richard et al., 2021; Fan et al., 2022; Xing et al., 2023), etc. Despite significant progress in these fields, existing methods are limited to employing audio-only speech for synthesis, as shown in Figure 1. However, in some complex environments such as construction sites and plants (Ahmed & Gadelmoula, 2022), it is difficult to extract clear audio signal, especially where headsets are unavailable due to regulations of production safety.

041 To understand speech in noisy scenes, researchers often use visual speech assistance (Afouras et al., 042 2018b; Shi et al., 2022b). For example, (Wang & Zhu, 2021) proposed a vision-based human-043 machine communication framework through gesture recognition, and (Ray & Teizer, 2012) also 044 confirmed the feasibility of real-time posture analysis used in workers' ergonomics training. However, 045 in most speech-driven body language synthesis tasks, such as Cross-ID talking head synthesis (Wang et al., 2021; Liang et al., 2022), Cross-ID landmark generation (Hsu et al., 2022), 3D facial animation, 046 and speech gesture synthesis, there is a paucity of paired data between visual speech and body 047 languages (Daněček et al., 2022). This scarcity hinders the realization of audio-visual speech-driven 048 body language synthesis models. 049

Since the audio speech modality has a large amount of data supporting self-supervised learning (Kahn et al., 2020; Zen et al., 2019), many researchers have successfully used speech discrete units to represent speech information. Inspired by direct speech-to-speech translation (Lee et al., 2021), and considering that each body language (including visual speech, i.e., lip movements) is aligned with audio and has ample training data, can we use discrete speech units extracted from audio

055 056

064

065

066

067

068

087

880

089

090

091

092

093

094

096



Figure 1: Overview of speech-driven multimodal synthesis tasks under noise-free and noisy conditions. The audio-only speech-driven approach is inadequate for synthesis in noisy environments. Conversely, the audio-visual speech-driven method can effective enhance the robustness of speech-driven synthesis in noisy settings. This paper focuses primarily on four tasks: Cross-ID talking head generation, Cross-ID landmark generation, 3D facial animation, and co-speech gesture generation.

069 speech self-supervised learning (SSL) models (Hsu et al., 2021; Zeghidour et al., 2021) as a bridge 070 between visual speech and body language to address this challenge and achieve direct audiovisual 071 speech-driven multimodal body language synthesis?

072 In this paper, we introduce a two-stage audio-visual speech-driven body language synthesis model 073 (AV-SDS) based on speech discrete units, which includes two basic modules: AV-S2UM and Unit2X. 074 The AV-S2UM module consists of an audio-visual speech encoder (Shi et al., 2022a) and several 075 transposed convolutional layers, which can map the audiovisual speech to the corresponding discrete 076 speech units. Subsequently, these discrete speech units are input into the Unit2X module to 077 synthesize the corresponding multimodal body language data. In particular, to meet the characteristics of different body language modalities, the AV-SDS can be implemented based on different speech discrete units: an acoustic-centered model for co-speech gestures which focus more on the emotion 079 and rhythm of speech (Loehr, 2007), implemented based on acoustic discrete units; and a semantic-080 centered model for modalities that focus on the semantic information of speech (talking head, facial 081 landmark, and 3D face mesh), proposed based on semantic discrete units. Experiments show that the Unit2X module can successfully resynthesize various body languages from discrete speech units. 083 Furthermore, experiments on various speech-driven body language synthesis tasks under different 084 noise conditions demonstrate that our AV-SDS is capable of noise-robust speech-driven synthesis. 085 The main contributions are as follows: 086

- We propose a novel two-stage, noise-robust, audio-visual speech-driven body language synthesis model (AV-SDS) based on discrete speech units.
- We propose AV-S2UM module, which excels in retaining speech information in noisy environments.
- We propose Unit2X module to synthesize body language data from discrete speech units, demonstrating the sufficiency of speech information in these units for body language synthesis.
- · Our experiments confirm the robustness and effectiveness of AV-SDS across varying noise conditions, validating its potential for noise-robust speech-driven body language synthesis.

RELATED WORKS 2

SPEECH DRIVEN MULTI-MODAL BODY LANGUAGE SYNTHESIS 2.1 098

099 Body language (Liu et al., 2023) plays a pivotal role in facilitating effective communication and 100 enhancing social interactions. Over the years, numerous researchers (Ye et al., 2023; Zhou et al., 101 2020) have dedicated their efforts to speech-driven body language synthesis, aiming to create digital 102 avatars that seamlessly synchronize with spoken content. Leveraging advancements in generative 103 technology, researchers have made notable progress, culminating in the development of both 2D 104 talking head avatars (Zhou et al., 2020; Prajwal et al., 2020) and 3D facial animation (Richard et al., 105 2021; Xing et al., 2023), driven by audio and speech inputs. Moreover, recognizing the importance of comprehensive avatar technology, researchers (Yang et al., 2023b;c) have expanded their exploration 106 into co-speech gesture synthesis, harnessing the rhythm and cadence of speech to imbue digital 107 avatars with lifelike gestures.

However, existing methods can only achieve multimodal body language synthesis driven solely by audio speech, lacking robustness in noisy environments. To address this gap, we propose the first direct audio-visual speech driven multi-modal body language synthesis model, enabling speech-driven synthesis even in noisy conditions.

112 113

114

2.2 ROBUST AUDIO-VISUAL SPEECH LEARNING

Audio speech understanding (Baevski et al., 2020; Hsu et al., 2021) technology has advanced rapidly, 115 116 effectively conveying speech content information. However, in noisy environments like outdoors, audio-only models often lack robustness and struggle to resist environmental noise interference. 117 To address this, researchers (Afouras et al., 2018b) have begun exploring the use of visual speech 118 (lip movements) to enhance speech understanding capabilities under noisy conditions. Some re-119 searchers (Afouras et al., 2018a;b) initially collected paired audio-visual speech corpus from TED and 120 BBC for research. Subsequently, AV-HuBERT (Shi et al., 2022a) achieves noise-robust audio-visual 121 speech recognition as well as lip reading, successfully recognizing speech content in noisy condition. 122 Additionally, some (Gao & Grauman, 2021; Hsu et al., 2023) propose using visual speech to tackle 123 the cocktail party problem, effectively addressing the challenge of distinguishing the active speaker 124 track among multiple speakers.

However, in speech-driven multimodal generation tasks, many body language modalities (e.g., meshes and gestures) lack paired data with visual speech (Daněček et al., 2022), hindering robust audio-visual speech-driven synthesis. In this paper, we introduce a two-stage AV-SDS framework that utilizes discrete speech units as carriers to overcome this challenge of unpaired data.

129 130

131

2.3 SELF-SUPERVISED LEARNING IN SPEECH

132 Self-supervised learning methods (Baevski et al., 2020; Hsu et al., 2021; Zeghidour et al., 2021; Yang 133 et al., 2023a) leverage unlabeled audio speech data to significantly enhance speech representation and improve the performance of various speech-related tasks. For instance, some researchers (Baevski 134 et al., 2020; Hsu et al., 2021) employ a continuous alternation between unsupervised clustering 135 and mask prediction to augment the contextual semantic representation. Later, some (Zeghidour 136 et al., 2021; Yang et al., 2023a) integrate the RVQ (residual vector quantization) module to achieve 137 a fine-grained representation of speech acoustic information. Building upon these advancements, 138 researchers have proposed utilizing Speech SSL model to discretize speech, thereby exploring new 139 capabilities of speech models. Specifically, Lee et al. (2021) employs semantic discrete units as 140 a bridge to achieve direct speech-to-speech translation between speeches of different languages, 141 while Jiang et al. (2023) employs acoustic discrete units to represent acoustic information of speech, 142 enabling zero-shot TTS. 143

While speech discrete units hold promise in various speech-related tasks, their application in speechdriven multi-modal generation tasks remains largely unexplored. In this paper, we introduce a unitbased multi-modal generation module, Unit2X, which represents a pioneering effort in leveraging speech discrete unit information for multi-modal generation exploration.

1481493Audio-Visual Speech Driven Synthesis

- 150
- 151

3.1 OVERVIEW

152 Audio-visual speech-driven multimodal synthesis aims to generate multimodal content M =153 $\{M_1, M_2, \dots, M_N\}$ (e.g., meshes, talking heads, gestures, etc.) consistent with audio speech 154 $A = \{A_1, A_2, \cdots, A_N\}$ and visual speech $V = \{V_1, V_2, \cdots, V_N\}$, where N is the number of 155 audio-visual speech frames. Due to the lack of extensive paired data between visual speech and body 156 language modalities (e.g., meshes and co-speech gestures), it is challenging to train an audio-visual 157 speech-driven body language synthesis model. In this context, as shown in Figure 2, the AV-SDS 158 proposed in this paper aims to use speech self-supervised learning (Speech SSL) model to map 159 speech paired with various modalities into the corresponding unified speech discrete unit space $U = \{U_1, U_2, \cdots, U_T\}$, where T is the number of discrete speech units. By using discrete speech 160 units as carriers, we achieve direct audio-visual speech-driven multimodal synthesis. Specifically, as 161 described in Section 3.2, we adopt the AV-S2UM module to construct the mapping from audio-visual



176

191

201

177 Figure 2: Illustration of **AV-SDS**. While there is no paired data between visual speech (lip movements) 178 and most modalities (such as mesh and co-speech gestures), paired speech data exists for each 179 modality in the speech-driven synthesis task. In this context, we leverage Speech SSL model (Hsu et al., 2021; Yang et al., 2023a) to convert the audio speech data from various modal pairs for the 181 speech-driven generation task into corresponding speech discrete units, acting as a bridge between 182 audio-visual speech and different modalities. As detailed in Section 3.2, we employ AV-S2UM to translate audio-visual speech into the target discrete speech units, followed by the Unit2X module 183 introduced in Section 3.3 to synthesize the corresponding multi-modal data, thus achieving robust audio-visual speech-driven multimodal synthesis. 185

186 speech to unified speech discrete units. Subsequently, in Section 3.3, we introduce Unit2X, a unit-187 based multimodal generation framework, to reconstruct corresponding multimodal body language 188 data from unified speech discrete units. 189

190 3.2 AUDIO-VISUAL SPEECH-TO-UNIT MAPPING SYSTEM

Speech Discrete Units. The information conveyed by speech can be broadly classified into two 192 categories: semantic information and acoustic information. HuBERT (Hsu et al., 2021) employs 193 multiple iterations of mask prediction and \mathcal{K} -means clustering to continuously enhance its ability 194 to understand speech context, effectively extracting semantic discrete units that encapsulate speech 195 semantics, denoted as U^s . Encodec (Défossez et al., 2022) and SoundStream (Yang et al., 2023a; 196 Zeghidour et al., 2021) utilize the RVQ module for precise speech reconstruction, obtaining the 197 acoustic discrete unit, denoted as U^a . In this paper, we utilize the semantic discrete unit U^s for generating modalities closely associated with semantics, such as the talking head, facial landmarks, 199 and the mesh of a 3D avatar. Conversely, the acoustic discrete unit U^a is employed for generating 200 co-speech gestures, which are more closely related to emotion and rhythm.

202 **AV-S2UM** In Figure 2, each visual speech V is paired with its corresponding audio speech A. We implement the AV-S2U-Mapper model based on the AV-HuBERT model, pre-trained on a large dataset 203 of paired audio-visual speech utterances. While it's possible to use the AV-HuBERT representation 204 directly to obtain discrete units, we opt for the Speech SSL model, trained specifically on speech 205 data, to ensure unified speech discrete units across different modalities. We can obtain the acoustic 206 discrete units (U_{lip}^a) and semantic discrete units (U_{lip}^s) corresponding to the audio speech in (A, V). 207

208 By leveraging the pre-trained AV-HuBERT model (Shi et al., 2022a), we embed the audio-visual speech into robust speech features f, denoted as f = AV-HuBERT(A, V). Subsequently, we employ 209 210 n transposed convolutional layers to align these features with discrete unit sequences (n = 1 for semantic discrete units and n = 3 for acoustic discrete units). Each layer utilizes a kernel size (K) 211 of 4, a stride (S) of 2, padding (P) of 1, and output padding (Op) of 1. The output size (O) of each 212 transposed convolutional layer is calculated using the formula $O = ((I-1) \times S + K - 2 \times P) + Op$. 213

214 The target distribution $p(U_t|\{U_i\}_{i=1}^{t-1}, (A, V))$ can be obtained with the robust speech feature f: 215

ŗ

$$p(U_t|\{U_i\}_{i=1}^{t-1}, (A, V)) = AV - S2UM(\mathbf{f}),$$
(1)

and the AV-S2UM module is trained using the cross-entropy loss:

$$L_{\text{AV-S2UM}} = -\sum_{t=1}^{N} \log p(U_t | \{U_i\}_{i=1}^{t-1}, (A, V)).$$
⁽²⁾

219 220 221

222 223

224

225

226

227

218

3.3 UNIT-BASED MULTI-MODAL SYNTHESIZER

With the Speech SSL model, the speech discrete unit U_m corresponding to the audio speech A paired with the multi-modal data M can be derived as either the corresponding semantic discrete unit (U_m^s) or the acoustic discrete unit (U_m^a) . In this subsection, we will elucidate how the Unit2X module synthesizes corresponding multi-modal data from semantic discrete units or acoustic discrete units.

228 Synthesizer based on Semantic Units. Drawing inspiration from Polyak et al. (2021), we employ 229 a lookup table (LUT) to map these discrete units $U_m^s = \{U_1^s, \dots, U_T^s\}$ to the corresponding speech 230 representation $\mathbf{f}^a = \{\mathbf{f}_1^a, \dots, \mathbf{f}_T^a\} = \text{LUT}(U_m^s)$. Subsequently, these speech representations \mathbf{f}^a are 231 inputted into various models to generate corresponding body language data.

For the acoustic discrete unit, we utilize 8 distinct units to represent various granular speech features at each speech frame, denoted as $U_m^a = \{U_{1,1}^a, \dots, U_{1,8}^a, \dots, U_{T,1}^a, \dots, U_{T,8}^a\}$. We obtain the audio speech representation \mathbf{f}^a by employing the pre-trained RVQ (Residual Vector Quantization) module within the Speech SSL model (Yang et al., 2023a) for acoustic units: $\mathbf{f}^a = \{\mathbf{f}_1^a, \dots, \mathbf{f}_T^a\} = \mathbb{RVQ}(U_m^a)$. Note that we keep the parameters of the RVQ module frozen during training.

237

244 245

254 255

256

257 258

259 260

261

Unit2X: Unit-Based multi-modal synthesizer. In the traditional speech-driven body language synthesis task, the audio encoder is utilized to encode the raw audio into the corresponding audio speech feature f^a . However, in this work, we substitute the original audio encoder with the unit-based encoder to obtain the audio speech feature. Once we obtain the audio speech representation f^a , we can apply the model used in the traditional speech-driven multi-modal synthesis task to generate corresponding multi-modal data M:

$$M = \text{Unit2X}(\mathbf{f}^a). \tag{3}$$

246 Here, we illustrate this process using the task of talking head generation as an example. The audio speech feature f^a is first input into the face decoder, where it is upsampled and combined with f^s , 247 the latter being extracted from randomly selected speaker reference frames and pose prior frames. 248 This combination generates the final talking head. The discriminator D consists of a series of 249 convolutional blocks and is trained alternately with the generator G. The loss function employed 250 during model training remains consistent with traditional speech-driven multi-modal generation 251 methods. Specifically, for the task of talking head generation, we use GAN loss L_G , lip reconstruction 252 loss L_{lip} , and synchronization loss L_{sync} as the training objectives: 253

$$L_{\text{unit}2x\,(\text{head})} = (1 - \lambda_{sync} - \lambda_{gen})L_{lip} + \lambda_{sync}L_{sync} + \lambda_{gen}L_G, \tag{4}$$

where $L_{unit2x (head)}$ is the training objective of Unit2X for talking head generation, and $\lambda_{sync} = 0.03$ and $\lambda_{gen} = 0.07$ as proposed by Prajwal et al. (2020).

4 EXPERIMENTS

4.1 DATASETS

262 The AV-S2UM module is trained on the LRS3 dataset (Afouras et al., 2018b). For the Unit2X 263 model, we utilize the most commonly employed datasets for each modality: 29h training split of 264 LRS2 (Petridis et al., 2018) for talking head generation, LRS3 (Afouras et al., 2018b) for facial 265 landmark synthesis, VOCASET (Cudeiro et al., 2019) for 3D facial animation, and TED-GESTURE 266 (Yoon et al., 2019) for co-speech gesture synthesis. The pre-trained Speech SSL models used in 267 this paper to obtain speech discrete units are trained on Libri-Light (Kahn et al., 2020) and the TTS Corpus (Yang et al., 2023a), respectively. Following the methodology described in Shi et al. (2022a), 268 we introduce noise into the audio speech by incorporating samples from the MUSAN dataset (Snyder 269 et al., 2015).

Table 1: Comparison of synthesis quality for various body language modalities. S2X denotes direct
 synthesis from real speech, Unit2S+S2X indicates synthesizing speech from units followed by
 additional synthesis, and Unit2X refers to the direct body language generation from discrete units.

Method	Ta	lking Head		M	esh	Landmark	Gesture				
	LSE-C↑	LSE-D↓	FID↓	SYNC.↑	REAL.↑	LMD↓	FGD↓				
Synthesize X-n	Synthesize X-modality data from real speech.										
S2X	7.50	7.14	5.08	45.01	39.72	4.287	4.133				
Synthesize X-n	Synthesize X-modality data from speech discrete units.										
Unit2S+S2X	5.74	7.96	5.52	41.37	35.29	5.180	4.228				
Unit2X(Ours)	7.34	7.54	5.14	46.83	43.28	4.718	3.976				

4.2 IMPLEMENTATION DETAILS

To ensure consistency across different body language modalities, we resample the audio speech in all datasets to 16kHz in this paper. This allows us to extract unified speech discrete units to represent semantic or acoustic information. Specifically, we use the HuBERT BASE model (Hsu et al., 2021) to extract semantic discrete units and the 16kHz version of the hificodec model (Yang et al., 2023a) to extract acoustic discrete units.

290 To facilitate effective research and encourage broad adoption, we intentionally selected fundamental, 291 widely applicable implementations for each modality. Specifically, we use Wav2Lip (Prajwal et al., 292 2020) for talking heads, GeneFace (Ye et al., 2023) for 3D landmarks, CodeTalker (Xing et al., 293 2023) for face meshes, and Tri-Modal (Yoon et al., 2020) for co-speech gestures to implement the 294 corresponding speech-driven body language synthesis tasks. We believe that experiments on basic implementations of different body language modalities are sufficient to demonstrate the effectiveness 295 of our approach, which can be integrated with any speech-driven body language synthesis model in 296 the future. For additional details and evaluation metrics, please refer to Appendix B. 297

298 299

281 282 283

284

4.3 UNIT2X: UNIT-BASED BODY LANGUAGE SYNTHESIS

Polyak et al. (2021) demonstrated the feasibility of reconstructing corresponding speech from discrete units. Building on this, we attempt to synthesize different body language data from corresponding speech discrete units. As shown in Table 1, we compare the performance of various synthesis methods across different body language modalities. This demonstrates that speech discrete units can effectively replace the original speech as input for synthesizing corresponding body language data. For additional qualitative comparisons, please refer to Appendix B.

Audio-Based vs. Unit-Based. The audio-based method (S2X) directly extracts the corresponding 307 embedding from the mel spectrum for body language modality synthesis, while the unit-based method 308 (Unit2X) synthesizes from discrete units. As demonstrated in the experiments, across various body 309 language modalities, the unit-based method achieves performance comparable to that of the audio-310 based method, affirming that speech discrete units can efficiently represent the speech information 311 necessary to generate the corresponding body language data. Additionally, in scenarios with limited 312 training data (e.g., mesh) or where cross-modal mapping is challenging to build (e.g., co-speech 313 gesture), unit-based methods prove to be more effective than methods using raw speech as input. For 314 instance, the FGD of Unit 2X is 3.976, while that of S2X is 4.133. This demonstrates that speech 315 discrete units can distinctly capture various semantic and acoustic information from speech, leading to superior performance in these challenging contexts. 316

U2S+S2X vs. Unit2X. To synthesize the body language data corresponding to the speech discrete unit, the simplest approach is to first re-synthesize the audio speech corresponding to the discrete unit, and then synthesize the corresponding body language data from the audio speech (i.e., U2S+S2X).
 However, this cascade method tends to accumulate errors, resulting in a significant decline in the correlation between the generated body language modality data and the audio speech compared to direct synthesis from the unit (the LSE-C of Unit2X is 7.34, while the LSE-C of U2S+S2X is only 5.74). This further underscores the importance of the unit-based body language synthesizer (Unit2X).

Table 2: Comparison of speech-driven synthesis performance across four different modalities (talking head, co-speech gesture, facial landmark, and mesh) under varying noise conditions. We present performance comparisons across different signal-to-noise ratios (SNRs) of SNR = {15, 5, -5, -15}.
 The cascade method with a dagger ([†]) employs the AV-S2UM+U2S+S2X cascade method. The mesh modality results are assessed against the generated mesh of clean audio speech-driven synthesis.

(a) Comparison of Talking Head Generation on LRS3. LSE-C and LSE-D in this table are evaluated between the generated talking head video and the clean audio speech.

Method	LSE-C↑			LSE-D↓				FID↓				
Methou	15	5	-5	-15	15	5	-5	-15	15	5	-5	-15
Wav2Lip	5.63	4.65	3.15	2.05	8.19	8.49	8.93	9.35	5.88	6.42	7.45	8.73
Cascade [†]	5.45	5.23	<u>5.09</u>	<u>4.83</u>	8.47	8.67	8.94	<u>9.12</u>	6.12	<u>6.26</u>	<u>6.41</u>	<u>6.67</u>
AV-SDS	5.72	5.59	5.41	5.12	<u>8.32</u>	8.45	8.51	8.68	<u>5.97</u>	6.12	6.23	6.37

(b) Comparison of 3D Mesh of Talking Head Generation on LRS3. The SYNC. and REAL. in this table are expressed as preference ratios compared to the results generated based on clean audio speech.

Met	hod		SYN	I C. ↑			REA	L. ↓		
	nou	15	5	-5	-15	15	5	-5	-15	-
Cod	eTalker	45.13	38.12	29.22	20.31	42.62	36.44	26.83	18.28	-
Case	cade [†]	38.78	35.94	32.13	28.89	36.38	35.04	33.71	31.66)
AV-S	SDS	<u>40.29</u>	39.75	38.32	34.49	43.13	41.26	39.85	36.72	
(c) Comparis	son of Lan	1000000000000000000000000000000000000			(d) Co	ompariso	FGD ↓			
Method			D↓		Met	hod		FGI)↓	
Method	15	5	D↓ -5	-15	Met	hod	15	FGI 5)↓ -5	-15
Method GeneFace	<u>15</u> 5.115	5.357	•••••••••••••••••••••••••••••••••••••	-15 6.356	Met Tri-r	hod nodal	15 4.322	FGI 5 4.726	•↓ -5 5.217	-15 5.910

353 354

329

330

339

355 356

4.4 ROBUST AUDIO-VISUAL SPEECH DRIVEN MULTI-MODAL SYNTHESIS

357 To assess the performance of different speech-driven body language synthesis methods in noisy 358 environments, we conducted validation across various body language modalities: (1) For talking head and facial landmarks, since there are no ground truth results in the cross-identity audio-visual 359 speech-driven talking head or facial landmark generation task, we follow the experimental setting of 360 Prajwal et al. (2020) and use talking heads and facial landmarks paired with audio speech as targets 361 for audio-visual speech-driven synthesis. Talking heads are evaluated on the LRS2 dataset, while 362 facial landmarks are evaluated on the LRS3 dataset. To further enhance the convincingness, we 363 also present cross-identity audio-visual speech-driven talking head generation results in Section 4.6, 364 validated solely through audio-visual synchronization (Prajwal et al., 2020). (2) For 3D face mesh, due to the absence of paired mesh and visual speech data for testing, we utilized the speech-driven 366 method to generate 3D face mesh corresponding to clean audio speech on the LRS3 dataset, using 367 these as reference videos for qualitative comparison. (3) For co-speech gesture, we curated a test set 368 comprising audio-visual speech and co-speech gesture paired data by re-collecting and processing the original video clips from the TED-GESTURE dataset as described by Afouras et al. (2018b). This 369 Audio-Visual TED-GESTURE (AV-GES) test dataset included 565 utterances in total. In Table 2, we 370 present the experimental results across various tasks, demonstrating the effectiveness of our proposed 371 AV-SDS method under noisy conditions. 372

Noise-Robust Speech-Driven Synthesis Traditional audio-only methods exhibit significant performance degradation in noisy environments. For instance, in Table 2d, the audio-only speech-driven method (Tri-modal) experiences a notable decrease in performance, dropping by 1.588 from 4.322 to 5.910 as SNR decreases from 15 to -15. In contrast, our proposed AV-SDS, driven by audio-visual speech, only experiences a slight decrease of 0.409 from 4.434 to 4.843. This demonstrates that AV-SDS achieves noise-robust speech-driven body language synthesis and is more robust to noise.

Method		$ $ WER(%) \downarrow		LSE-C↑			LSE-D↓				MOS ↑		
memou	15	5	-5	-15	15	5	-5	-15	15	5	-5	-15	Avg.
Inp.Audio	7.8	17.8	63.9	87.9	6.72	6.68	3.38	2.05	7.58	8.12	10.82	11.93	2.88±0.17
Resynthesis	10.2	19.9	83.5	97.7	6.64	6.53	3.19	1.82	7.65	8.29	11.04	12.14	2.81 ± 0.10
Demucs	6.9	15.1	48.0	81.3	6.98	6.89	3.87	2.21	7.38	7.83	10.12	11.79	3.22 ± 0.13
VisualVoice	6.6	8.8	23.4	58.0	7.03	6.75	5.70	4.81	7.34	7.61	8.57	9.79	3.52 ± 0.12
ReVISE	9.4	9.7	11.7	20.5	6.63	6.59	6.41	5.79	7.49	7.56	7.78	8.46	3.40 ± 0.1

378 Table 3: Comparison of speech enhancement performance under different noise conditions on LRS3. 379 We re-evaluated the LSE metrics following Prajwal et al. (2020), and reproduced the results of 380 ReVISE based on the AV-S2UM (U^s) module and unit-based vocoder (Polyak et al., 2021).

390 391 392

393

394

395

397

407

408 409

381 382

Particularly under harsh noise conditions (SNR= $\{5, -5, -15\}$), AV-SDS achieves superior performance in speech-driven body language synthesis, preserving more information from the speech and yielding more reliable body language synthesis. It's worth noting that since the AV-S2UM model is trained solely on LRS3, the experiments on co-speech gesture in Table 2d represent zero-shot scenar-396 ios for the AV-S2UM module. Even in such conditions, AV-SDS demonstrates better performance under noisy conditions, underscoring the significance of audio-visual speech-driven synthesis.

 $AV-S2UM(U^a)$ 5.8 6.5 13.6 46.0 7.25 7.14 7.06 5.85 7.15 7.32 7.51 8.30 3.79±0.12

Direct System vs. Cascade System. While the cascade approach using the three modules 399 AV-S2UM+U2S+S2X can achieve a certain degree of noise-resistant audio-visual speech-driven 400 synthesis, the accumulation of errors from multiple module cascades hampers its synthesis perfor-401 mance under varying noise conditions compared to AV-SDS. Particularly, robust audio-visual speech understanding in noisy environments is significantly challenging. The errors generated by each 402 module under noisy conditions are non-negligible and significantly impact the final outcome. For 403 instance, in the results for SNR = -15 in Table 2d, the FGD of the cascade method is 0.524 lower 404 than that of AV-SDS. Therefore, minimizing the number of cascade layers is crucial for speech-driven 405 synthesis tasks in noisy environments. 406

4.5 AV-S2UM: PRESERVATION OF SPEECH INFORMATION IN NOISY AUDIO.

Based on visual speech, AV-S2UM effectively preserves speech information in noisy environments 410 and can reconstruct corresponding audio speech using speech discrete units. As shown in Table 411 3, we compared the speech enhancement performance to evaluate the ability of various models to 412 retain speech information under different noise conditions. Among these models, Demucs (Defossez 413 et al., 2020) represents an audio-only approach, VisualVoice (Gao & Grauman, 2021) represents 414 a naive audio-visual method, ReVISE (Hsu et al., 2023) relies on semantic discrete units (i.e., 415 $AV-S2UM(U^s)$), and $AV-S2UM(U^a)$ refers to the model that relies on acoustic discrete units.

416 Audio-Only vs. Audio-Visual. In a noisy environment, models that rely solely on audio-speech 417 are unable to resist noise interference and suffer significant loss of speech information. The Word 418 Error Rate (WER) for the audio-only method (Demucs) at SNR=-15 is 81.3%, showing only a 419 6.6% improvement from 87.9% for Inp.Audio. However, audio-visual speech-based methods use 420 visual speech as auxiliary information to help models resist noise interference. VisualVoice, for 421 example, maintains a WER of 58.0% at SNR=-15, which is 23.3% better than the audio-only method, 422 demonstrating the importance of visual speech for audio understanding in noisy environments.

423 Naive Audio-Visual Method vs. Unit-Based Audio-Visual Method. Traditional visually guided 424 speech enhancement methods only rely on a limited amount of audio-visual speech pairing data 425 for training, making it difficult to achieve high-fidelity audio reconstruction. However, methods 426 based on discrete units, such as ReVISE and our AV-S2UM (U^a), train the AV-S2UM module and 427 the corresponding vocoder separately on large-scale audio-visual speech and massive audio speech, 428 achieving more effective audio information retention. Under the conditions of SNR = -5 and SNR =-15, the WER of unit-based ReVISE is 12% and 37.5% better than that of end-to-end VisualVoice, 429 respectively. This experiment demonstrates that the two-stage method, which has more training data, 430 can retain more original audio information compared to the end-to-end model, which can only use 431 extremely limited paired data.

432 Semantic Discrete Unit vs Acoustic Discrete Unit. ReVISE employs semantic discrete units to 433 retain semantic speech information while disregarding acoustic elements such as timbre, emotion, 434 and rhythm (the speaker's timbre in the ReVISE speech output differs from the original timbre). In 435 contrast, $AV-S2UM(U^a)$ advances by utilizing acoustic discrete units as intermediaries to connect 436 visual speech with high-fidelity audio speech. AV-S2UM (U^a) not only preserves speech semantic information but also retains fine-grained acoustic details, including emotion, timbre, and rhythm. It 437 excels in audio-visual synchronization (LSE-C and LSE-D are the best among all noise conditions), 438 showcasing its applicability to body language modalities closely related to emotion and rhythm, such 439 as co-speech gestures. Notably, although ReVISE cannot retain the acoustic speech information, 440 it maintains an excellent WER even under high noise conditions (WER of 20.5% at SNR=-15). 441 It demonstrates that the AV-S2UM (U^s) module based on semantic discrete units can effectively 442 preserve semantic speech information in noisy environments, making it well-suited for synthesizing 443 body language modalities related to speech semantics. 444

445 446

4.6 AUDIO-VISUAL SPEECH-DRIVEN VS. VIDEO-DRIVEN.

447 Although we use visual speech (i.e., lip movements) 448 to facilitate audio-visual speech-driven body language 449 synthesis, our AV-SDS differs significantly from the 450 video-driven body language synthesis methods (Pang et al., 2023). Typically, video-driven methods gen-451 erate the target individual's expressions and gestures 452 based on a reference video. In contrast, our approach 453 relies solely on speech information extracted from au-454 diovisual speech to ensure that the generated results 455 are consistent with the driving speech, without con-456 sidering other visual information from the reference 457 video. Given the difficulty in distinguishing speaker 458 identity details from facial movements, results gener-459 ated by video-driven methods (Yin et al., 2022) often 460 retain speaker identity attributes (such as face shape 461 and makeup) from the driving video. Additionally, we present lip sync metrics for cross-identity driven talking 462

Table 4: Comparison of cross-id talking head generation results using different modality driving methods on LRS3 at SNR=-5. The LSE metrics (LSE-C and LSE-D) are evaluated between the generated talking head video and the clean audio speech. AO-Speech: Audio-Only Speech.

Method	Driven	LSE-C↑	LSE-D↓
Wav2Lip	AO-Speech	2.467	9.466
DPE	Video	3.620	9.824
AV-SDS	AV-Speech	4.447	9.242

head synthesis in Table 4. Due to the lack of speech-related supervision, video-driven methods struggle to maintain a high level of lip sync during the synthesis process, only simulating the expressions of the driving video to a limited extent. Consequently, the generated video cannot effectively convey the corresponding speech content. Notably, even under extremely strong noise interference conditions (SNR=-5), AV-SDS outperforms video-driven method (DPE), demonstrating the relevance of our method for the task of audio-visual speech-driven body language synthesis.

469

5 CONCLUSION

470 471

Speech-driven body language synthesis aims to create intelligent digital humans that align with audio 472 speech. However, due to the lack of paired data for visual speech and body language modalities, 473 existing methods can rely on audio-only speech, which struggles to produce accurate results under 474 noisy conditions. To address this issue, we propose the first direct audio-visual speech-driven 475 multi-modal synthesis framework, AV-SDS. This framework employs speech discrete units as an 476 intermediate carrier in a two-stage approach to bridge audio-visual speech and various body language 477 modalities. Firstly, AV-S2UM maps audio-visual speech to unified discrete units. Then, Unit2X 478 synthesizes various multi-modal body language data from these units. We introduce Unit2X, the 479 first multi-modal body language synthesis model based on speech discrete units, and demonstrate the 480 feasibility of using speech discrete units instead of raw audio speech for body language synthesis. Additionally, we propose two different implementations based on semantic discrete units and acoustic 481 discrete units for semantically related and rhythm-related body language modalities, respectively. In 482 various speech-driven multi-modal body language synthesis tasks, our AV-SDS achieves state-of-the-483 art performance under different noise conditions, confirming its effectiveness in noisy environments. 484

485

486 REFERENCES

516

523

527

528

529

Triantafyllos Afouras, Joon Son Chung, Andrew Senior, Oriol Vinyals, and Andrew Zisserman. Deep audio-visual speech recognition. *IEEE transactions on pattern analysis and machine intelligence*, 44(12):8717–8727, 2018a.

- Triantafyllos Afouras, Joon Son Chung, and Andrew Zisserman. Lrs3-ted: a large-scale dataset for visual speech recognition. *arXiv preprint arXiv:1809.00496*, 2018b.
- 494 SS Ahmed and AM Gadelmoula. Industrial noise monitoring using noise mapping technique: a case
 495 study on a concrete block-making factory. *International Journal of Environmental Science and* 496 *Technology*, 19(2):851–862, 2022.
- Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460, 2020.
- Lele Chen, Zhiheng Li, Ross K Maddox, Zhiyao Duan, and Chenliang Xu. Lip movements generation at a glance. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 520–535, 2018.
- Daniel Cudeiro, Timo Bolkart, Cassidy Laidlaw, Anurag Ranjan, and Michael J Black. Capture, learning, and synthesis of 3d speaking styles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10101–10111, 2019.
- Radek Daněček, Michael J Black, and Timo Bolkart. Emoca: Emotion driven monocular face capture
 and animation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 20311–20322, 2022.
- Alexandre Defossez, Gabriel Synnaeve, and Yossi Adi. Real time speech enhancement in the waveform domain. *arXiv preprint arXiv:2006.12847*, 2020.
- Alexandre Défossez, Jade Copet, Gabriel Synnaeve, and Yossi Adi. High fidelity neural audio
 compression. *arXiv preprint arXiv:2210.13438*, 2022.
- Yingruo Fan, Zhaojiang Lin, Jun Saito, Wenping Wang, and Taku Komura. Faceformer: Speech driven 3d facial animation with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18770–18780, 2022.
- Ruohan Gao and Kristen Grauman. Visualvoice: Audio-visual speech separation with cross-modal consistency. In 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 15490–15500. IEEE, 2021.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. Advances in neural information processing systems, 30, 2017.
 - Gee-Sern Hsu, Chun-Hung Tsai, and Hung-Yi Wu. Dual-generator face reenactment. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 642–650, 2022.
- Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:3451–3460, 2021.
- Wei-Ning Hsu, Tal Remez, Bowen Shi, Jacob Donley, and Yossi Adi. Revise: Self-supervised speech resynthesis with visual input for universal and generalized speech regeneration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18795–18805, 2023.
- Ziyue Jiang, Yi Ren, Zhenhui Ye, Jinglin Liu, Chen Zhang, Qian Yang, Shengpeng Ji, Rongjie
 Huang, Chunfeng Wang, Xiang Yin, et al. Mega-tts: Zero-shot text-to-speech at scale with intrinsic inductive bias. *arXiv preprint arXiv:2306.03509*, 2023.

540 541 542 543 544	Jacob Kahn, Morgane Rivière, Weiyi Zheng, Evgeny Kharitonov, Qiantong Xu, Pierre-Emmanuel Mazaré, Julien Karadayi, Vitaliy Liptchinsky, Ronan Collobert, Christian Fuegen, et al. Libri-light: A benchmark for asr with limited or no supervision. In <i>ICASSP 2020-2020 IEEE International</i> <i>Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pp. 7669–7673. IEEE, 2020.
545 546 547	Ann Lee, Peng-Jen Chen, Changhan Wang, Jiatao Gu, Sravya Popuri, Xutai Ma, Adam Polyak, Yossi Adi, Qing He, Yun Tang, et al. Direct speech-to-speech translation with discrete units. <i>arXiv</i> preprint arXiv:2107.05604, 2021.
548 549 550	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.
551 552 553 554	Borong Liang, Yan Pan, Zhizhi Guo, Hang Zhou, Zhibin Hong, Xiaoguang Han, Junyu Han, Jingtuo Liu, Errui Ding, and Jingdong Wang. Expressive talking head generation with granular audio-visual control. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 3387–3396, 2022.
555 556 557	Li Liu, Lufei Gao, Wentao Lei, Fengji Ma, Xiaotian Lin, and Jinting Wang. A survey on deep multi- modal learning for body language recognition and generation. <i>arXiv preprint arXiv:2308.08849</i> , 2023.
558 559 560 561	Xian Liu, Qianyi Wu, Hang Zhou, Yuanqi Du, Wayne Wu, Dahua Lin, and Ziwei Liu. Audio-driven co-speech gesture video generation. <i>Advances in Neural Information Processing Systems</i> , 35: 21386–21399, 2022a.
562 563 564 565	Xian Liu, Qianyi Wu, Hang Zhou, Yinghao Xu, Rui Qian, Xinyi Lin, Xiaowei Zhou, Wayne Wu, Bo Dai, and Bolei Zhou. Learning hierarchical cross-modal association for co-speech gesture generation. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 10462–10472, 2022b.
566 567	Daniel Loehr. Aspects of rhythm in gesture and speech. Gesture, 7(2):179-214, 2007.
568 569 570	Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. fairseq: A fast, extensible toolkit for sequence modeling. <i>arXiv preprint arXiv:1904.01038</i> , 2019.
572 573 574	Youxin Pang, Yong Zhang, Weize Quan, Yanbo Fan, Xiaodong Cun, Ying Shan, and Dong-ming Yan. Dpe: Disentanglement of pose and expression for general video portrait editing. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 427–436, 2023.
575 576 577	Stavros Petridis, Themos Stafylakis, Pingchuan Ma, Georgios Tzimiropoulos, and Maja Pantic. Audio-visual speech recognition with a hybrid ctc/attention architecture. In 2018 IEEE Spoken Language Technology Workshop (SLT), pp. 513–520. IEEE, 2018.
579 580 581	Adam Polyak, Yossi Adi, Jade Copet, Eugene Kharitonov, Kushal Lakhotia, Wei-Ning Hsu, Ab- delrahman Mohamed, and Emmanuel Dupoux. Speech resynthesis from discrete disentangled self-supervised representations. <i>arXiv preprint arXiv:2104.00355</i> , 2021.
582 583 584	KR Prajwal, Rudrabha Mukhopadhyay, Vinay P Namboodiri, and CV Jawahar. A lip sync expert is all you need for speech to lip generation in the wild. In <i>Proceedings of the 28th ACM international conference on multimedia</i> , pp. 484–492, 2020.
586 587 588	KR Prajwal, Triantafyllos Afouras, and Andrew Zisserman. Sub-word level lip reading with visual attention. In <i>Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition</i> , pp. 5162–5172, 2022.
589 590	Soumitry J Ray and Jochen Teizer. Real-time construction worker posture analysis for ergonomics training. <i>Advanced Engineering Informatics</i> , 26(2):439–455, 2012.
591 592 593	Alexander Richard, Michael Zollhöfer, Yandong Wen, Fernando De la Torre, and Yaser Sheikh. Meshtalk: 3d face animation from speech using cross-modality disentanglement. In <i>Proceedings</i> of the IEEE/CVF International Conference on Computer Vision, pp. 1173–1182, 2021.

594 595 596	Bowen Shi, Wei-Ning Hsu, Kushal Lakhotia, and Abdelrahman Mohamed. Learning audio-visual speech representation by masked multimodal cluster prediction. <i>arXiv preprint arXiv:2201.02184</i> , 2022a.
597 598	Bowen Shi, Wei-Ning Hsu, and Abdelrahman Mohamed. Robust self-supervised audio-visual speech
599	recognition. arXiv preprint arXiv:2201.01763, 2022b.
600 601	David Snyder, Guoguo Chen, and Daniel Povey. Musan: A music, speech, and noise corpus. <i>arXiv</i> preprint arXiv:1510.08484, 2015.
602	
603 604	sive portrait videos with audio2video diffusion model under weak conditions. <i>arXiv preprint</i>
605	arXiv:2402.17485, 2024.
606 607	Christophe Veaux, Junichi Yamagishi, Kirsten MacDonald, et al. Superseded-cstr vctk corpus: English multi-speaker corpus for cstr voice cloning toolkit. 2016.
608 609 610 611	Ting-Chun Wang, Arun Mallya, and Ming-Yu Liu. One-shot free-view neural talking-head synthesis for video conferencing. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 10039–10049, 2021.
612 613	Xin Wang and Zhenhua Zhu. Vision–based framework for automatic interpretation of construction workers' hand gestures. <i>Automation in Construction</i> , 130:103872, 2021.
614	The Min Market Min V. Walter Zhao, Min Lee Charles Westernet The This Was
615	Codetalker: Speech-driven 3d facial animation with discrete motion prior. In <i>Proceedings of</i>
616	the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 12780–12790, 2023.
619	
619	Dongchao Yang, Songxiang Liu, Rongjie Huang, Jinchuan Tian, Chao Weng, and Yuexian Zou.
620 621	arXiv:2305.02765, 2023a.
622	Sicheng Yang, Zhiyong Wu, Minglei Li, Zhensong Zhang, Lei Hao, Weihong Bao, Ming Cheng, and
623 624	Long Xiao. Diffusestylegesture: Stylized audio-driven co-speech gesture generation with diffusion models. <i>arXiv preprint arXiv:2305.04919</i> , 2023b.
625 626 627 628	Sicheng Yang, Zhiyong Wu, Minglei Li, Zhensong Zhang, Lei Hao, Weihong Bao, and Haolin Zhuang. Qpgesture: Quantization-based and phase-guided motion matching for natural speech- driven gesture generation. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and</i> <i>Pattern Recognition</i> , pp. 2321–2330, 2023c.
629 630	Zhenhui Ye, Ziyue Jiang, Yi Ren, Jinglin Liu, JinZheng He, and Zhou Zhao. Geneface: Generalized and high-fidelity audio-driven 3d talking face synthesis. <i>arXiv preprint arXiv:2301.13430</i> , 2023.
632 633 634	Fei Yin, Yong Zhang, Xiaodong Cun, Mingdeng Cao, Yanbo Fan, Xuan Wang, Qingyan Bai, Baoyuan Wu, Jue Wang, and Yujiu Yang. Styleheat: One-shot high-resolution editable talking face generation via pre-trained stylegan. In <i>European conference on computer vision</i> , pp. 85–101. Springer, 2022.
635 636 637 638	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.
639 640 641	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.
642 643 644 645	Neil Zeghidour, Alejandro Luebs, Ahmed Omran, Jan Skoglund, and Marco Tagliasacchi. Sound- stream: An end-to-end neural audio codec. <i>IEEE/ACM Transactions on Audio, Speech, and</i> <i>Language Processing</i> , 30:495–507, 2021.
646 647	Heiga Zen, Viet Dang, Rob Clark, Yu Zhang, Ron J Weiss, Ye Jia, Zhifeng Chen, and Yonghui Wu. Libritts: A corpus derived from librispeech for text-to-speech. arXiv preprint arXiv:1904.02882, 2019.

Yang Zhou, Xintong Han, Eli Shechtman, Jose Echevarria, Evangelos Kalogerakis, and Dingzeyu Li. Makelttalk: speaker-aware talking-head animation. *ACM Transactions On Graphics (TOG)*, 39(6): 1–15, 2020.

A DATASETS

TTS Corpus (Yang et al., 2023a). The TTS Corpus integrates various public datasets, including LibriTTS (Zen et al., 2019) and VCTK (Veaux et al., 2016), encompassing 1,000 hours of high-fidelity English speeches. Importantly, all speech segments within the corpus have been meticulously confirmed to be free of discernible background noise. The 16kHz audio speech codec (Yang et al., 2023a) employed in this work is pre-trained on this corpus. This codec, rooted in the high-fidelity characteristics of the corpus, plays a crucial role in achieving the desired performance and fidelity in various speech-related tasks undertaken in this work.

LRS2 (Afouras et al., 2018a) and LRS3 (Afouras et al., 2018b). LRS2 and LRS3 stand out as the
 most expansive publicly accessible lip-reading dataset at the sentence level, boasting over 229/443
 hours of video content sourced from BBC and TEDx talks. In our experiments, we harnessed the
 train set to extract facial landmarks, following the methodology outlined by Ye et al. (2023).

In this paper, LRS3 serves as a crucial resource for evaluating the performance across various audio visual speech-driven synthesis tasks for many body language modalities, such as talking head, facial
 landmarks and 3d mesh.

672 VOCASET (Cudeiro et al., 2019). VOCASET consists of 480 paired audio-visual sequences
673 recorded from 12 subjects. The facial motion is captured at 60fps, lasting approximately 4 seconds
674 each. Each 3D face mesh is registered to the FLAME (Li et al., 2017) topology, featuring 5023 vertices.
675 To ensure fair comparisons, we utilize the same training (VOCA-Train), validation (VOCA-Val), and
676 testing (VOCA-Test) splits as VOCA (Cudeiro et al., 2019).

TED GESTURE (Yoon et al., 2019). The TED Gesture Dataset encompasses a substantial volume
 of paired audio-visual sequences derived from TED talks, offering both a sizable dataset for inves tigating the intricate relationship between speech and gestures. Covering a diverse range of topics,
 TED talks feature thousands of unique speakers sharing their individual ideas and stories, capturing a
 broad spectrum of speech content.

MUSAN (Snyder et al., 2015). In this paper, we randomly selected audio samples from MUSAN datasets to introduce background noise to the speech content. MUSAN (Snyder et al., 2015) consists of music, speech, and babble noise. Following the approach of Shi et al. (2022a), we used the audio samples from MUSAN to add noise to the speech.

B MORE IMPLEMENTATION DETAILS

692 B.1 TRAINING DETAILS

For the training of AV-S2UM module, we loaded the publicly available pretrained weights of AV-HuBERT (Shi et al., 2022a) and fine-tuned the model over a total of 45,000 steps. During the first 5,000 steps, we exclusively trained the decoder by freezing the encoder. Afterward, we unfroze the encoder and trained the entire model together. The learning rate was adjusted using a tri-stage LR scheduling strategy with specific phases set at (10%, 20%, 70%) and a peak learning rate of 6e-5. The training was conducted on one V100 GPUs. We utilized the Adam optimizer with parameters set to (0.9, 0.98).

For the training of Unit2X module, we strictly follow the training details of each body language modality for training, and all models are trained on a single V100 GPU.

Method	Talkin	g Head	Landmark	Gesture	
	Qual.	Sync.	Sync.	Sync.	
S2X	4.12 ± 0.09	4.09 ± 0.13	3.76 ± 0.18	3.88±0.15	
Ground Truth	4.33 ± 0.12	4.15 ± 0.09	3.95 ± 0.15	4.03±0.12	
U2S+S2X	4.06±0.15	3.87±0.15	3.65±0.18	3.76±0.18	
Unit2X(ours)	4.09±0.12	3.98±0.16	3.68±0.16	3.92±0.16	

Table 5: Qualitative comparison of body language synthesis performance across different methods.

Table 6: Qualitative comparison for speech enhancement under various noise conditions.

Method	SNR=15	SNR=5	SNR=-5	SNR=-15	Mean
Inp.Audio Resynthesis Demucs	3.82 ± 0.11 3.74 ± 0.09 3.97 ± 0.10	3.53 ± 0.14 3.46 ± 0.17 3.62 ± 0.13	2.33 ± 0.20 2.26 ± 0.19 2.92 ± 0.17	$\begin{array}{c} 1.83{\pm}0.21\\ 1.78{\pm}0.20\\ 2.35{\pm}0.18\end{array}$	$\begin{array}{c} 2.88 {\pm} 0.17 \\ 2.81 {\pm} 0.16 \\ 3.22 {\pm} 0.15 \end{array}$
VisualVoiceReVISEAV-S2UM (U^a)	$\begin{array}{c} 4.01{\pm}0.08\\ 3.52{\pm}0.09\\ \textbf{4.25}{\pm}\textbf{0.07}\end{array}$	3.77±0.10 3.42±0.12 4.05±0.09	3.42±0.13 3.37±0.15 3.68±0.12	2.87±0.16 3.29±0.15 3.18±0.18	3.52±0.12 3.40±0.13 3.79±0.12

B.2 METRICS

Unit-Based Body-language Synthesis. (1) For lip movements, we employ LSE-C and LSE-D (Prajwal et al., 2020) as evaluation metrics to assess the synchronization between audio speech and lip movements. In the context of talking heads, we used FID (Heusel et al., 2017) to assess the dissimilarity between the generated images and the real images. (2) For facial landmarks, the facial landmark distance (LMD) (Chen et al., 2018) is used to measure the distance between the generated landmarks and ground truth landmarks. (3) For mesh, we conduct A/B testing to evaluate the authenticity (Real.) and synchronicity (Sync.) of various mesh synthesis methods. Similar to Xing et al. (2023), the evaluation is based on the percentage of samples with a higher user preference than ground truth (GT) videos. (4) For co-speech gesture, we employ FGD (Fused Gaussian Distance) as metrics. FGD (Yoon et al., 2020) measures the distribution disparity between generated output and ground truth with a pre-trained autoencoder.

Speech Enhancement. For the speech enhancement task, we employ the ASR model (Ott et al., 2019) to transcribe the denoised speech, using Word Error Rate (WER) as a metric to assess content accuracy. Additionally, to evaluate the synchronization between the denoised speech and the talking head video, LSE-C and LSE-D are adopted to assess lip synchronization.

C QUALITATIVE EXPERIMENTS

We performed a manual evaluation of all generated results, appraising qualitative outcomes through the Mean Opinion Score (MOS) methodology. Each sample was randomly presented to 15 participants for scoring. The composite MOS was subsequently calculated by averaging scores across the relevant dimensions. Each dimension was independently rated on a scale of 1 (lowest) to 5 (highest). Kindly visit the demo page (https://av-sds.github.io/) to view the corresponding generation results. The comprehensive MOS evaluation details for each task are outlined below:

Unit-Based Body Language Synthesis. In Table 5, we present a qualitative comparison among multiple methods on different body language modalities. For talking heads, we evaluated image quality (Qual.) and lip synchronization (Sync.). In the case of facial landmarks and co-speech gestures, the focus was on evaluating the synchronization (Sync.) of audio speech and bodyity data.

As the evaluation metric for the 3D mesh modality depends on manual assessment, we refrained from conducting further experiments on this particular modality.

Speech Enhancement. For the visually guided high-fidelity speech denoising task, we evaluated denoising outcomes across various noise conditions (SNR= $\{15, 5, -5, -15\}$), as depicted in Table 6. It is important to highlight that ReVISE consistently receives low subjective scores due to its inability to reconstruct the timbre of the corresponding speech. Notably, when SNR= -15, owing to its robust semantic reconstruction capability, ReVISE obtained the highest MOS. However, in other instances, AV-S2UM(U^a) demonstrated superior high-fidelity speech noise reduction results, earning top ratings.

D LIMITATION

This paper verifies only a limited range of body language modalities (talking head, mesh, co-speech gesture, and 3D landmark). However, we believe these modalities are sufficient to demonstrate the effectiveness of our method. In the future, we will also validate it on additional modalities, such as listener responses and others.

E ETHICAL DISCUSSION

The task studied in this paper involves the field of virtual human synthesis, which carries a certain risk of video forgery. However, since the focus of this paper is not on the authenticity of the synthesis but on the robustness of the voice-driven body movement synthesis task in a noisy environment, this concern is not particularly serious.