

PLAYING FOR YOU: TEXT PROMPT-GUIDED JOINT AUDIO-VISUAL GENERATION FOR NARRATING FACES USING MULTI-ENTANGLED LATENT SPACE

Anonymous authors

Paper under double-blind review

ABSTRACT

We present a novel approach for generating realistic speaking and talking faces by synthesizing a person’s voice and facial movements from a static image, a voice profile, and a target text. The model encodes the prompt/driving text, the driving image, and the voice profile of an individual and then combines them to pass them to the multi-entangled latent space to foster key-value pairs and queries for the audio and video modality generation pipeline. The multi-entangled latent space is responsible for establishing the spatiotemporal person-specific features between the modalities. Further, entangled features are passed to the respective decoder of each modality for output audio and video generation. Our experiments and analysis through standard metrics demonstrate the effectiveness of our model. All model checkpoints, code, and the proposed dataset can be found at: <https://github.com/Playing-for-you>.

1 INTRODUCTION

AI-generated real-time audio-video multimedia communication by rendering realistic human talking faces has recently drawn massive attention^{1,2}. Such technology is promising in various applications such as digital communication, aiding communication with individuals with impairments, designing artificial instructors, and developing interactive healthcare (Xu et al., 2024b; Gan et al., 2023). In such applications, generating realistic and real-time speech and visual content simultaneously is a key requirement. Therefore, in an ideal scenario, given a prompt text along with a face image and the audio profile of an individual, a talking human face would be rendered as output with audio (generated speech) and visual narration according to the prompt text.

Generative AI has emerged as a key area of interest in the computer vision and learning representation community. Although existing approaches have made significant strides, they are constrained by their reliance on generating a single modality (Egger et al., 2020; Kim et al., 2021). For example, current text-to-speech models (TTSM) (Ao et al., 2022; Betker, 2022; Casanova et al., 2024) focus primarily on voice synthesis. Similarly, visual generation techniques i.e. talking face models (TFM) (Ren et al., 2021; Rombach et al., 2022; Siarohin et al., 2020; Zhang et al., 2023a; Xu et al., 2024b;c; Zhang et al., 2023b) aim at face video generation given a text or/and audio or/and image as a prompt. Hence both TTSM and TFM techniques are unsuitable for real-life audio-video multimedia communication scenarios such as audio-visual chatbots, as in such situations both realistic video and speech must be generated synchronously and simultaneously. Few efforts have been made in the literature to merge TTSM and TFM by cascading the pipeline (Wang et al., 2023; Zhang et al., 2022). Additionally, (Jang et al., 2024) made an effort to generate talking face and speaking audio jointly for a specific individual from a prompt text.

Further, these TFM (Chen et al., 2024; Zhang et al., 2019) depend on guidance from defined facial properties from the weakly supervised latent information from the reference modality. As a result,

¹https://www.business-standard.com/technology/tech-news/odisha-television-introduces-lisa-india-s-first-ai-news-presenter-123071000767_1.html

²<https://www.indiatoday.in/india/story/india-today-groups-ai-anchor-sana-wins-global-media-award-for-ai-led-newsroom-transformation-2532514-2024-04-27>

poor lip-synchronization and limited ability to tune an existing audio profile for personalizing the video content lead to generation that is far from being realistic. Moreover, expressiveness in facial dynamics along with subtle nuances for realistic facial behavior needs to simultaneously match with audio content temporally to produce realistic talking faces. Further, such synchronization also depends on individual traits, such as their speech intonation and other covariates. Although they are supposed to be important considerations for realistic speaking and taking faces models (STFM), However, this was not in the scope of existing work on STFM (Jang et al., 2024). Therefore, this gap in the literature motivates us to design a prompt text-guided audio-visual multimodal generative STFM that can jointly generate audio and video, given a reference image and reference audio along with the prompt text as input.

Consequently, in contrast to existing literature (See Figure 1), in this work we introduce a novel multi-modal framework designed to address these limitations by generating highly realistic speech and animations from a combination of prompt text, a driving image, and an audio profile as inputs. Specifically, our framework aims to synthesize videos of a talking human face where the person in the image appears to speak along with the generated voice from the provided text for the given identity. Our method enhances the capabilities of existing pre-trained models (Xu et al., 2024b) with an advanced parallel mechanism that leverages both visual and auditory data streams. This parallelism ensures that the synthesized videos not only align the subject’s facial movements with the spoken text but also synchronize with the generated personalized voice outputs that correspond to the subject’s appearance.

A person-agnostic generalized STFM model must encompass a large appearance and acoustic features variation. Furthermore, extracting such structure information along with the temporal synergy between the audio and video while preserving individual variance requires additional modules to model these complexities. Therefore, we introduce a parallel multiple entanglement in the latent space between the encoding and decoding of different modalities.

Our proposed architecture for STFM contains three main phases (See Figure 2). *Modality encoding phase*, at this stage a heterogeneous personal signature of the audio and video modality, and the driving feature from the text are extracted. The second stage is the *multi-entangled latent space* which glens the spatiotemporal relation and synchronization in the embeddings of the modalities, which further acts as the input to the *decoders phase* i.e the third stage of the proposed architecture. In the second stage, the exchange of information between the key and values (identity information from audio and video extracted from the individual encoders) and queries (driving features from encoded prompt text) are streamlined. To instrument this, an entanglement of the audio and text latent is performed which further entangles with video latent in transformers block and then to a diffusion block. The output of the diffusion block is passed to the video decoder. Similarly, an entanglement of the video and text latent is performed which further entangles with audio latent in a transformer space and passes to a text decoder block and then to the audio decoder. Such entanglements ensure to streamlining of the audio profile and the driving image by linear navigation in the latent space along with the encoded feature from the prompt text. Specifically, the temporal information for both the audio and video generation is constructed by linear displacement of codes in the latent space as per the encoded text prompt. In turn, the model also learns a set of orthogonal motion directions to simultaneously learn the audio and video temporal synergy, by exchanging their linear combination

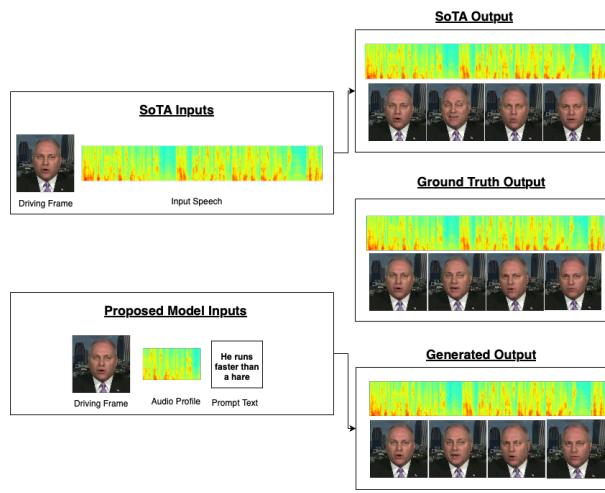


Figure 1: SOTA approaches of talking face generation use a face image as driving frame, with an audio prompt passed as input to the existing model such as Hallo (Xu et al., 2024b), VASA (Xu et al., 2024c) and the proposed model which generates a realistic audio-video synchronous multimodal talking face with face image and audio profile of an individual along with the prompt text.

108 to represent any displacement in the latent space. To summarize, our key contributions are as follows:
 109

- 110 • To the best of our knowledge, the proposed architecture is the first person-agnostic STFM
 111 which fosters a text-driven multimodal realistic audio-video synthesis that can be generalized
 112 to any identity.
- 113 • We design a three-phase architecture which consists of the encoder, multi-entangled latent
 114 and decoder phase for audio and video pipeline. The multi-entangled latent space glens
 115 the spatiotemporal and synchronisation in the encoder embedding to exchange information
 116 between the modality and guided text and help to generate crucial visual and acoustic
 117 characteristics based on input profiles.
- 118 • With the comprehensive experiments, we demonstrate that the proposed method surpasses
 119 the state-of-the-art techniques available for STFM.

120 2 RELATED WORK

121 Text-to-speech (TTS) technology has seen remarkable progress in recent years, with the development
 122 of models that generate highly natural and expressive speech. Modern Text-To-Speech
 123 approaches(Casanova et al., 2024; Betker, 2022) leverage sequence-to-sequence architectures to
 124 map text directly to speech. Notable models among these are the Tacitron(Wang et al., 2017) and
 125 the newer Tacitron2(Shen et al., 2018). These models employ attention mechanisms to convert text
 126 sequences into mel-spectrograms. These spectrograms are then passed through neural vocoders
 127 like WaveNet(van den Oord et al., 2016) or HiFi-GAN(Kong et al., 2020) to generate high-quality
 128 audio waveforms. Other models, such as FastSpeech(Ren et al., 2019) and VITS(Kim et al., 2021),
 129 introduce optimizations to improve the speed of speech generation while maintaining or enhancing
 130 the naturalness and clarity of the output. Although models have advanced into more complex ar-
 131 chitectures, the underlying idea behind speech generation remains the same. TortoiseTTS(Betker,
 132 2022) is a modern, expressive TTS system with impressive voice cloning capabilities. This model
 133 incorporates a combination of the Auto-Regressive Model, followed by a Diffusion Model(Ho et al.,
 134 2020), to convert the input text into mel-spectrogram frames, via discrete acoustic tokens. This model
 135 also follows the standard of a vocoder(Univnet)(Jang et al., 2021) for generating the audio from the
 136 spectrogram frames. Only a few works have been made in the literature to attend STFM by cascading
 137 the pipeline (Wang et al., 2023; Zhang et al., 2022). In (Jang et al., 2024) advancements are made by
 138 generating a talking face and speaking audio jointly for a specific individual from a prompt text.
 139

140 2.1 FACE REENACTMENT AND LIP-SYNC MODELS

141 Recent advancements in face reenactment have enabled realistic video generation by synthesizing
 142 facial movements driven by audio inputs. Early models, such as SyncNet(Raina & Arora, 2022),
 143 focused on lip synchronization through facial key points and phoneme mapping but struggled
 144 with capturing detailed expressions and diverse facial structures. More recent models, such as
 145 LipGAN(K R et al., 2019) and Wav2Lip(Prajwal et al., 2020a), leverage GANs to improve lip-sync
 accuracy and generate more natural facial animations.

146 The multimodal synthesis of human videos, combining text, audio, and visual inputs, has advanced
 147 considerably in recent years. Early approaches focused on audio-driven models that primarily
 148 addressed lip-syncing, mapping speech inputs to corresponding facial movements. Models like
 149 SyncNet(Raina & Arora, 2022) played a crucial role in establishing baseline synchronization between
 150 audio and lip movements. However, these models often lacked expressive, natural face dynamics.
 151

152 2.2 DIFFUSION-BASED LIP-SYNC MODELS

153 Recent models have extended beyond simple lip-syncing to incorporate emotional expression and
 154 natural head motion. Audio2Head(Wang et al., 2021), for example, shifts from keypoint-based
 155 methods to a dense mapping of audio features onto facial expressions and head motion, resulting
 156 in a more fluid and expressive representation of speech-driven animations. Expressive Audio-
 157 driven Talking-heads (EAT)(Gan et al., 2023) enhances this by integrating text and audio as inputs,
 158 introducing more dynamic and natural facial expressions synchronized with speech.

159 The Hallo(Xu et al., 2024b) model builds on these advancements by using attention mechanisms
 160 to improve facial reenactment, ensuring smoother transitions and better coherence across diverse
 161 speakers. Furthermore, SadTalker(Zhang et al., 2023b) incorporates 3D facial representations,
 combining both speech and facial dynamics for more realistic head motions and expressive gestures.

162 FaceChain-ImagineID(Xu et al., 2024a) uses latent diffusion to generate talking faces directly from
 163 the only audio input, generating synthetic faces after disentangling the audio to extract aspects
 164 like expression, identity and emotion. Other notable works, such as Diffused Heads(Stypulkowski
 165 et al., 2023) and DreamTalk(Zhang et al., 2023a), have explored diffusion-based models for video
 166 generation, leveraging the success of image-to-video transformations in generating high-quality
 167 talking-head videos. These models focus on temporally consistent video generation, addressing
 168 fidelity and synchronization across frames.

169 3 METHODOLOGY

170 We propose a joint learning methodology for the audio, video, and natural language-based text prompts
 171 consisting of three main components – namely, (1) Encoding phase, (2) Entanglement of combined
 172 latent space, and (3) Decoding phase *i.e.*, Latent conditional generation of synthesized audio-video.
 173 Figure 2 illustrates detailed network architecture and roles of different model components to learn
 174 and dynamically synthesize audio video on a given source image.

175 3.1 MULTI-MODAL ENCODING PHASE.

176 We use HiFi-GAN (Kong et al., 2020) and Wav2Vec Encoder (Baevski et al., 2020) to extract
 177 high-dimensional embedding vectors from the reference audio. The HiFi-GAN generates a feature
 178 embedding \mathbf{f}_a that represents the audio waveform. At the same time, the Wav2Vec encoder produces
 179 a secondary set of embedding \mathbf{f}_s capturing semantic audio information. We treat the semantic audio
 180 embedding as a direct mapping of the speaker’s voice profile. Consequently, the combined features
 181 $\mathbf{f}_a \oplus \mathbf{f}_s$ provide a detailed audio profile necessary for driving the lip-sync and facial animations in the
 182 synthesized video. The input reference audio is represented as a 2-second MEL-spectrogram, encoder
 183 into a sequence of acoustic features per frame of 0.2 seconds duration with the shape of $\mathbb{R}^{5609 \times 512}$.

184 Our neural model’s newly inducted input text prompt undergoes Byte-Pair Encoding (BPE) and
 185 Tokenization (Zouhar et al., 2024) to convert textual information into a feature vector $\mathbf{f}_t \in \mathbb{R}^{512 \cdot T}$.
 186 This feature vector enables context-specific animations, allowing the synthesized video to align with
 187 the intended spoken words and expressions implied in the text. The purpose of concatenating \mathbf{f}_t
 188 with the combined feature of reference audio $\mathbf{f}_a \oplus \mathbf{f}_s$ is to obtain the speaker’s signature in the final
 189 flattened feature tokens of $\mathbf{f}_t \oplus \mathbf{f}_a \oplus \mathbf{f}_s \in \mathbb{R}^{5609+T \times 512}$.

190 Next, we process the input source image through a Variational Auto-Encoder (VAE) (Kingma &
 191 Welling, 2022) and a Landmarks Detection model (Zhang et al., 2020). The VAE generates an
 192 image embedding \mathbf{f}_i , representing the visual style and identity of the person in the source image.
 193 Concurrently, the landmarks detection network extracts structural features – face mask feature
 194 \mathbf{f}_{fm} and lip mask feature \mathbf{f}_{lm} , which are combined with the image embedding vectors to create a
 195 fused visual feature representation $\mathbf{f}_i \oplus \mathbf{f}_{lm} \oplus \mathbf{f}_{fm} \in \mathbb{R}^{3136 \times 512}$. The straightforward tendency of
 196 traditional methods is either to introduce prior 3D morphable models faces (Zhang et al., 2023b),
 197 motion priors of the facial parts (Jang et al., 2024), or guiding video frames (Wang et al., 2022) to
 198 learn nuances of facial articulation in relation to the audio in combined latent space. In contrast, we
 199 show that the entanglement of multiple latent spaces of text-audio-video using Transformer encoders
 200 (Vaswani et al., 2023) can eliminate the dependency on strong motion priors. As a result, we are able
 201 to use text prompt features as a set of anchoring tokens to both the Transformer encoders.

202 3.2 ENTANGLEMENT OF COMBINED TEXT-AUDIO-VIDEO LATENT SPACE.

203 As illustrated in Figure. 2, a smooth synergy between the text-audio latent embedding and the
 204 text-image latent embedding is established by two Transformer encoders followed by latent diffusion-
 205 guided (Xu et al., 2024b) synthesizer of visual nuances and decoder-only GPT-2 (Casanova et al.,
 206 2024) model for synthesizing text-conditioned audio latent.

207 The first Transformer encoder spatially contextualizes the audio MEL-spectrogram tokens using
 208 a dual-stream cross-modal attention mechanism with the flattened version, denoted by $\mathbf{L}(\cdot)$, of
 209 *categorically fixed speaker* embedding tokens merged with varying text embedding tokens, *i.e.*,
 210 $\mathbf{Q}_a = \mathbf{L}(\mathbf{f}_a \oplus \mathbf{f}_s)$, as

$$212 \text{Cross-Attention}(\mathbf{Q}_a, \mathbf{K}_{ti}, \mathbf{V}_{ti}) = \text{SoftMax}\left(\frac{\mathbf{Q}_a \mathbf{K}_{ti}^\top}{\sqrt{d_k}}\right) \mathbf{V}_{ti}, \quad (1)$$

213 where the query vector \mathbf{Q}_a is of dimension $\mathbb{R}^{5609 \times 512}$ and the key-value paring ($\mathbf{K}_{ti}, \mathbf{V}_{ti}$) between
 214 the tokens of $\mathbf{L}(\mathbf{f}_t \oplus \mathbf{f}_i \oplus \mathbf{f}_{lm} \oplus \mathbf{f}_{fm})$ has a variable spatial length (padded up-to a max length)

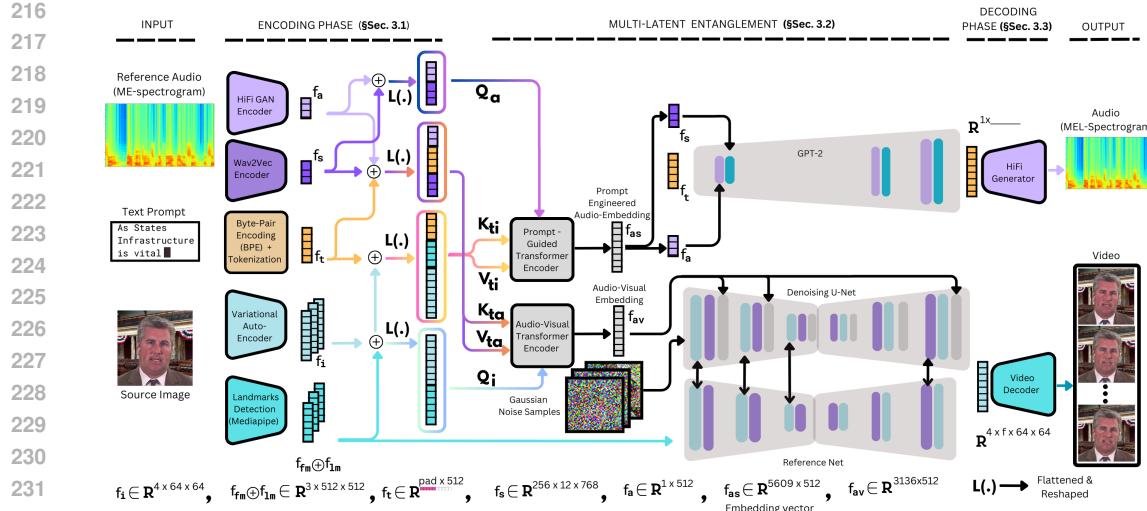


Figure 2: **Our Network Architecture:** Text Prompt-guided joint audio-visual learning representations using dual stream Transformer Encoders and Denoising Diffusion model. The model architecture can be divided into three phases – namely *Encoding Phase*, *Multi-Latent Entanglement*, and *Decoding Phase*. As an output, an audio-visual animation is generated from a single source image, reference audio, and a short text prompt.

with a fixed channel length of 512. Merging the varying text tokens serves two purposes – (1) first, querying audio tokens as well as the speaker tokens has been implicitly prompt-engineered by the text tokens, (2) second when the resulting prompt-engineered latent embedding vectors f_{as} are split into its respective constituents, they become proxy weights of text-image embedding vectors.

Similar to the previous encoder block, the second Transformer encoder spatially contextualizes the input masked-image embedding vectors $L(f_i \oplus f_{fm} \oplus f_{lm})$ using cross-modal attention with the key-value pairs (K_{ta}, V_{ta}) of merged text-audio embedding tokens $L(f_t \oplus f_a \oplus f_s)$ similar to the equation 1 as

$$\text{Cross-Attention}(Q_i, K_{ta}, V_{ta}) = \text{SoftMax}\left(\frac{Q_i K_{ta}^T}{\sqrt{d_k}}\right) V_{ta}. \quad (2)$$

As a result, the output latent embedding on audio-visual features f_{av} can serve as a compact and compressed representation of facial animation sequences in the high-dimensional space. Therefore, our next step is to learn a synthesizer *i.e.*, a hierarchical latent diffusion model Xu et al. (2024b) for video generation and a corresponding MEL-spectrogram synthesizer based on the X-Text-to-Speech (XTTS) model Casanova et al. (2024).

Latent Text Conditioned Spectrogram Synthesizer: The GPT-2 encoder is based on the TTS model (Casanova et al., 2023) and (Shen et al., 2018). This part is composed of a decoder-only transformer module that is conditioned by the audio and speaker embedding vectors f_a, f_s disentangled from the prompt-engineered audio embedding vector f_{av} , and the auto-regressive generation of spectrogram tokens is fully driven by the input text tokens from f_{av} .

Text-Anchored Audio-Video Latent Conditioned Denoising Diffusion: The Denoising Diffusion model aims to reverse a diffusion process(Ho et al., 2020; Song et al., 2022) that progressively adds random Gaussian noise to data. Inspired by the Hallo method (Xu et al., 2024b), we employ an additional augmentation of the text-anchored latent embedding vector learned to combine the audio and motion nuances on a single image inside the Denoising U-Net (Ronneberger et al., 2015) model of Hallo. The model is initialized with pre-trained weights and fine-tuned during the training step.

Throughout each step of the diffusion process, we introduce embedding cross-attention, which incorporates the combined latent space embedding, particularly our f_{av} , into each diffusion step. This cross-attention mechanism allows the diffusion models to leverage the shared information across modalities, ensuring that the generated outputs (audio and video) are consistent with the input

embedding. The inclusion of cross-attention helps to maintain coherence between the synthesized motion across all the pixels of the source image.

Additionally, diffusion cross-attention facilitates mutual information exchange between the audio and video diffusion blocks. This cross-attention mechanism enables the audio and video models to synchronize their outputs, ensuring that the generated audio and video components are temporally aligned. By integrating this cross-attention, our framework effectively coordinates the diffusion processes, leading to synchronized and coherent multimedia output.

3.3 DECODING PHASE FOR AUDIO-VIDEO GENERATION

The outputs of the previous steps are processed by their respective final decoders. For audio generation, similar to the XTTS method (Casanova et al., 2024), the synthesized spectrogram is passed through a Vocoder component of HiFi Generator module to obtain the final audio signal. For video, the Denoising UNet generates f number of frames of dimension $\mathbb{R}^{4 \times f \times 64 \times 64}$, which are decoded by a pre-trained decoder component of (Kingma & Welling, 2019) to produce the complete video.

3.4 LOSS FUNCTIONS

To train our model, we use –

(1) Video Loss as the Pixel-wise L1 Loss *i.e.*, sum of the N number of pixel intensities between the ground truth image frame $\mathcal{I}_{\text{gt}}^f$ and the generated frame $\mathcal{I}_{\text{gen}}^f$ for all the f number of frames as $\mathcal{L}_{\text{video}} = \sum_f \sum_{i=1}^N \|(\mathcal{I}_{\text{gt}}^f)^i - (\mathcal{I}_{\text{gen}}^f)^i\|$, (2) Audio Loss as the Spectrogram MSE loss at the spectrogram S domain as mean squared error between the ground-truth magnitudes and generated magnitudes at different of time step t as $\mathcal{S}_{\text{gt}}^t$ and the generated frame $\mathcal{S}_{\text{gen}}^t$ as $\mathcal{L}_{\text{audio}} = \frac{1}{T} \sum_{t \in T} \|(\mathcal{I}_{\text{gt}}^f)^i - (\mathcal{I}_{\text{gen}}^f)^i\|^2$. Total loss as $\mathcal{L}_{\text{Total}} = \lambda \mathcal{L}_{\text{audio}} + \mathcal{L}_{\text{video}}$ with balancing factor $\lambda = 0.1$.

4 EXPERIMENTAL RESULTS

4.1 DATASETS, PREPROCESSING, IMPLEMENTATION DETAILS AND EVALUATION MATRICES

Datasets: We have primarily conducted our experiments on 4 datasets. Our model training was done on a combination of **VoxCeleb** Dataset (Nagrani et al., 2019), **FakeAVCeleb** dataset (Khalid et al., 2022), **HDTF** (Zhang et al., 2021) and the **CelebV-HQ** dataset (Zhu et al., 2022). VoxCeleb is an audio-visual dataset consisting of short clips of human speech, extracted from interview videos uploaded to YouTube. FakeAVCeleb is a novel audio-video multimodal deepfake dataset. We only considered the non-deepfake part of the dataset. CelebV-HQ is a large-scale video facial attributes dataset demonstrating a diverse quality of data, which is important to test the robustness of our model. HDTF is a large in-the-wild high resolution audio-visual dataset built for talking face generation.

Preprocessing: Our preprocessing involved resizing the videos to 512x512 and then cropping each video sample to the first 20 seconds (at 25FPS which equates to 500 frames). We then separated the audio from the video using ffmpeg, and then ran the OpenAI’s Whisper model(Radford et al., 2022) to transcribe the audio speeches.

Implementation details: The optimizer used for our model is AdamW with a learning rate of 1e-4 and weight decay of 1e-2, and the scheduler has a step-wise learning rate with a step size of 1000 and gamma of 0.5. The weight decay regularizes the model, preventing any overfitting. We have used Nvidia 1xA6000s GPU for training each model, **and the model inference requires 12GB of VRAM**. The total parameter size of the model comes to 1,575,936 and performs 5.39 GFLOPs (Giga Floating Point Operations) per generation. We have trained the models for 10 epochs, with a batch size of 8. The Hifi-Gan, Wav2Vec Encoders, the Variational Autoencoder, Diffusion Models, and the GPT2 Decoder are pre-trained, which were further trained with the rest of the entire proposed network

Evaluation Metrics: Following are the evaluation matrices employed.

Video Metrics: *Fréchet Video Distance (FVD)*: A measure of the quality of generated videos, comparing them to real videos based on spatio-temporal features. Lower values indicate better performance (Unterthiner et al., 2019). *FID (Fréchet Inception Distance)*: Evaluates the visual quality of individual frames by comparing the distributions of generated and real images. Lower scores represent better visual quality (Heusel et al., 2018). *Fréchet Video Motion Distance(FVMD)*: Measures the quality of motion in generated videos, capturing the difference between real and generated motion trajectories. Lower values signify a more realistic motion.(Liu et al., 2024).

324 Audio Metrics: *Fréchet Audio Distance (FAD)*: Assesses the similarity between generated and real
325 audio samples, with lower scores indicating closer resemblance. *Short-Time Objective Intelligibility*
326 (*STOI*): Measures the intelligibility of the generated speech. Higher values represent more intelligible
327 speech (Kilgour et al., 2019). *Mel Cepstral Distortion(MCD)*: A metric used to evaluate the quality
328 of speech synthesis by comparing the spectral features of generated and reference audio. Lower
329 scores imply better audio quality (Zezario et al., 2020).

330 Audio-visual (AV) synchronisation: We used two metrics proposed in Wav2Lip Prajwal et al.
331 (2020b) to find the audio-visual synchronisation. The first is the average error measure calculated in
332 terms of the distance between the lip and audio representations, “LSE-D” (“Lip Sync Error Distance”).
333 A lower LSE-D denotes a higher audio-visual match, i.e., the speech and lip movements are in
334 synchronization. The second metric is the average confidence score, “LSE-C” (Lip Sync Error
335 Confidence). The higher the confidence, the better the audio-video correlation.

336 Training and Testing: Our primary training dataset is the VoxCeleb dataset(Nagrani et al., 2019),
337 where our training set comprised of approximately 36000 videos. We chose this training set by
338 filtering out individuals whose speech was in English. We tested on more than 200 samples from
339 each of the four datasets (VoxCeleb, FakeAVCeleb, CelebV-Hq and HDTF.), resulting in a test set of
340 over 800 unseen samples.

341 We benchmarked the video outputs for the unseen samples against SoTA Portrait Animation models,
342 like Hallo(Xu et al., 2024b), Sadtalker(Zhang et al., 2023b), EAT(Gan et al., 2023) and Au-
343 dio2Head(Wang et al., 2021). We also benchmarked the audio outputs for the unseen samples against
344 SoTA Speech generation models, like Tortoise(Betker, 2022), Your_TTS(Casanova et al., 2023),
345 XTTS_v2(Casanova et al., 2024) and GlowTTS(Kim et al., 2020).

346 4.2 RESULT ANALYSIS

348 Video Results: From Table 1, we can observe that our model shows superior performance across
349 all three metrics FID, FVD, and FVMD on VoxCeleb, CelebV-Hq and HDTF. This indicates high
350 fidelity and minimal discrepancies are attended by the proposed model. On the FakeAVCeleb, the
351 performance is slightly poorer but can be comparable, it still maintains strong visual consistency and
352 realism on visual inspection. For the CelebV-HQ our model excels again, demonstrating its capability
353 to produce high-quality video outputs. On HDTF our model shows incredible performance in the FID
354 and FVD metrics, beating all the other models, while our model is admirably performing considering
355 FVMD when compared to Hallo.

356 Table 1: Video pipeline evaluation scores across datasets.

Dataset	Model	FID Score (↓)	FVD Score (↓)	FVMD Value (↓)
VoxCeleb	Audio2Head	81.00	90.12	5100.92
	Hallo	67.28	70.69	5703.44
	EAT	85.16	80.38	4878.36
	SadTalker	119.36	112.77	6352.19
	Our Model	42.88	49.78	4192.07
FakeAVCeleb	Audio2Head	93.59	97.85	1329.23
	Hallo	26.88	39.42	2351.20
	EAT	94.34	98.49	1324.91
	SadTalker	81.77	77.10	4158.18
	Our Model	47.24	49.15	2263.54
CelebV-HQ	Audio2Head	90.22	102.76	2939.49
	Hallo	42.76	56.10	2816.68
	EAT	47.88	56.21	2894.31
	SadTalker	52.60	52.55	2789.19
	Our Model	34.01	43.67	2743.29
HDTF	Audio2Head	37.78	32.69	2633.04
	Hallo	20.54	25.81	1290.57
	EAT	29.57	29.34	2573.05
	SadTalker	22.34	23.57	2410.89
	Our Model	11.72	15.58	1784.16

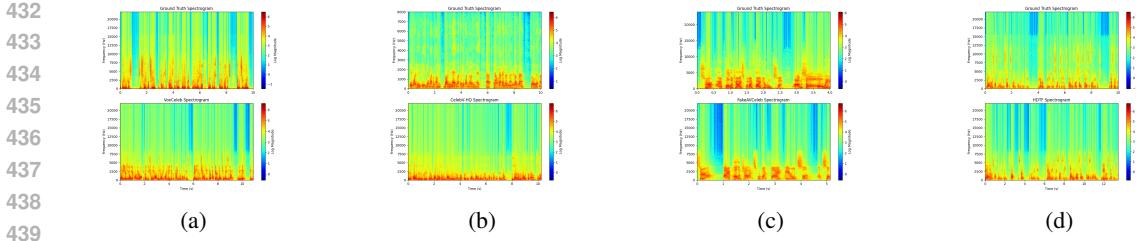
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
Table 2: Audio pipeline evaluation scores across datasets.

Dataset	Model	FAD Score (↓)	MCD Score (↓)	STOI Score (↑)
VoxCeleb	Tortoise	258.54	82.37	0.10
	Your_TTS	199.52	111.79	0.19
	XTTS_v2	249.17	100.80	0.13
	GlowTTS	329.21	103.94	0.15
	Our Model	241.75	75.39	0.17
FakeAVCeleb	Tortoise	871.14	82.12	0.10
	Your_TTS	445.38	65.60	0.21
	XTTS_v2	184.39	77.88	0.11
	GlowTTS	482.04	87.11	0.18
	Our Model	171.52	55.12	0.19
CelebV-HQ	Tortoise	529.06	113.18	0.09
	Your_TTS	520.01	137.58	0.16
	XTTS_v2	509.90	124.61	0.07
	GlowTTS	549.18	139.81	0.22
	Our Model	244.83	85.76	0.18
HDTF	Tortoise	425.30	67.15	0.11
	Your_TTS	467.42	49.38	0.15
	XTTS_v2	135.11	49.65	0.14
	GlowTTS	510.61	66.42	0.12
	Our Model	106.43	44.05	0.15



420
421
422
423
424
425
Figure 3: The figures in each row show frames from the videos generated by each technique in the
order: Ground Truth, Our proposed Model, Audio2Head (Wang et al., 2021), EAT (Gan et al., 2023),
Hallo (Xu et al., 2024b), and SadTalker (Zhang et al., 2023b) on the VoxCeleb Dataset. A frame
in each column for both videos corresponds to the same time-stamp (frames were sampled at equal
intervals of 25 seconds across the videos).

426
427
428
429
430
431
Based on the results, we observed that for some datasets certain models work slightly better than the
proposed model, and the reason behind this is that those models try to memorize certain properties
from individual datasets. Whereas our model is a more generalized version that can perform
consistently on cross datasets having varying resolution, and video quality. The visualization from
Figure 3 also concludes that our model can generate video very close to the ground truth and better
than any model. From Figure 4 it can be concluded that our model can generate nearby results for
HDTF, FakeAVCeleb and CelbV-HQ when compared to ground truth.



440 Figure 5: Ground Truth vs. Generated Audio Spectrograms for (a) VoxCeleb, (b) CelebV-HQ, (c)
441 FakeAVCeleb and (d) HDTF datasets

442
443
444 **Audio Results:** We can infer from [Table 2](#) that
445 our model consistently performs the best in the
446 MCD Score metric, which suggests that it min-
447 imizes distortion between the spectral features
448 of synthetic and reference speech. While con-
449 sidering the FAD scores, our model also per-
450 formed on par state-of-the-art, except on Vox-
451 Celeb where Your_TTS is better, these show-
452 case that the proposed model can generate con-
453 stantly similar audio compared to the ground
454 truth. Considering the STOI metric, the per-
455 formance of our model is similar to or slightly
456 lower than Your_TTS. The analysis of all the
457 measures showcases that our model is more gen-
458 eralized and realistic as it can minimize distor-
459 tion and also generate accurate distributions, and
460 maintain intelligibility of the speech consistently
461 better than any other models. The visualization
462 from Figure 5 also concludes that our model can
463 generate audio very close to the ground truth.

464 **AV synchronization results:** From Table 3 we
465 can conclude that our proposed model has per-
466 formed better audio-video synchronization than
467 SOTA and is close to the ground truth. The pro-
468 posed model has the lowest LSE-D, i.e. better
469 audio-visual match, i.e. and LSE-C i.e. better
470 audio-video correlation. We have also analyzed
471 the model with varying accents, blurred audio
472 profiles, and audio profiles of a kid with a source
473 image of an adult and vice versa, and the results
474 were found to be effective, no bias was found in
475 any aspect. Models fail in a few scenarios where a very noisy audio profile is used, output audio is
476 feeble or for source images with closed eyes face dynamics get affected (details are in supplementary).

4.3 ABLATION STUDY

477 [Table 4](#) shows the ablation study of our proposed model.
478 We have 3 main sub-networks that define the output
479 of our model. The **Transformer Encoder Block(TE)**
480 ([Vaswani et al., 2023](#)) with two variations shared-TE
481 (**STE**) where both audio and video pipeline shares a
482 transformer block and explicit-TE (**ETE**) where audio
483 and video pipeline has explicit or separate transformer
484 block. **Diffusion**([Song et al., 2022](#)) **Cross Atten-**
485 **tion(DC)**, and the **Embedding Cross Attention(EC)**.
486 From our results, it is understandably explainable that
487 the transformer encoder block, which encodes our in-



487 Figure 4: Results of our model on FakeAVCeleb,
488 CelebV-HQ and HDTF datasets.

489 Table 3: Evaluation of audio-visual syn-
490 chronization

	LSE-C(\uparrow)	LSE-D(\downarrow)
Groundtruth	5.45	8.52
Hallo	3.03	8.71
Audio2Head	2.51	10.34
EAT	4.39	9.35
SadTalkert	5.44	10.09
STE	5.71	8.41
ETE/ Proposed	5.74	8.38

486 puts into a common latent space, is the most important modality of our network, with its removal
 487 drastically reducing our metric values. Our experiments also show that the cross-attention blocks
 488 between the diffusion models are more important than the embedding cross-attention since our
 489 metric values drop more when we remove the diffusion cross-attention, probably since the diffusion
 490 cross-attention already syncs the modalities during the parallel learning stage. Another important
 491 aspect of ablation is the encoding latent in the individual transformer i.e. ETE is much better than
 492 STE. This infers that it is important to encode the latent for each modality separately while sharing
 493 information among the generated modalities. **Table 5 shows our ablation study on the encoders.**
 494 "Only Visual Tokens Attended" involves eliminating the audio prompt-guided transformer. Similarly,
 495 the "Only Audio Tokens Attended" involves using only the audio prompt-guided transformer encoder.
 496 "No Hifi-GAN" and "No Wav2Vec" are results obtained by eliminating the encoding process of the
 497 Hifi-GAN and Wav2Vec Models respectively. "No Visual Attention in AV-Transformer" involves not
 498 attending the visual tokens to the Audio Tokens in the Audio-Visual Transformer Encoder. These
 499 ablations quantify the importance of each of the components.

500 Table 4: Ablation study of the transformers.

501 ETE	502 STE	503 DC	504 EC	505 FID (\downarrow)	506 FVD (\downarrow)	507 FVMD (\downarrow)	508 FAD Score (\downarrow)	509 MCD (\downarrow)	510 STOI (\uparrow)
		✓	✓	86.70	80.88	5275.89	328.27	95.44	0.07
		✓		68.83	74.19	4412.74	260.91	87.51	0.11
		✓	✓	63.68	71.38	4298.30	250.12	83.96	0.14
		✓	✓	61.44	69.15	2720.41	241.77	81.60	0.17
		✓	✓	42.88	49.78	4192.07	241.75	75.39	0.17

511 Table 5: Ablation study of the encoders.

512 Ablation	513 FID (\downarrow)	514 FVD (\downarrow)	515 FVMD (\downarrow)	516 FAD Score (\downarrow)	517 MCD (\downarrow)	518 STOI (\uparrow)
Only Visual Tokens Attended	68.31	78.42	5747.04	304.98	81.17	0.13
Only Audio Tokens Attended	69.02	79.35	6576.85	301.49	80.65	0.13
No Hifi-GAN	85.25	94.28	7483.40	498.33	87.51	0.09
No Wav2Vec	70.10	80.96	5926.64	309.95	89.58	0.11
No Visual token in Prompt Guided transformer	54.38	62.02	5481.36	221.07	63.25	0.12
Proposed Model	42.88	49.78	4192.07	241.75	75.39	0.17

519 4.4 SOCIAL RISKS AND MITIGATIONS

520 There are social risks with technology development for text-driven audio video talking face generation.
 521 The foremost risk is the ethical implications of creating highly realistic talking faces, it can be used
 522 for malicious purposes, such as deepfakes. To mitigate such risk, ethical guidance for the use of such
 523 generation techniques is required. Also, concerns regarding privacy and consent are implicit in such
 524 work. Transparent data usage policies by consent, and safeguarding the privacy of individuals can
 525 mitigate such concerns. By addressing these we aim to promote responsible and produce ethical
 526 generative technology.

527 5 CONCLUSION

528 This paper introduces a novel method for realistic speaking and talking faces by joint multimodal
 529 video and audio generation. We provide a holistic architecture where the information is exchanged
 530 between the modalities via the proposed multi-entangled latent space. A source image of an individual
 531 as a driving frame, reference audio which can be referred to as the audio profile of the individual
 532 and a driving or prompt text is passed as an input. The model encodes the input driving image,
 533 prompt/driving text, and the voice profile which are further combined and passed to the proposed
 534 multi-entangled latent space consisting of two separate transformers and diffusion block for video
 535 and text decoder for audio pipeline to foster key-vale and query representation for each modality.
 536 By this spatiotemporal person-specific featuring between the modalities is also established. The
 537 entangled-based learning representation is further passed to the respective decoder of audio and
 538 video modality for respective outputs. Conducted experiments and ablation studies prove that the
 539 proposed multi-entangled latent-based learning representation has helped our model obtain superior
 540 results on both video and audio outputs as compared to state-of-the-art models. While there is always
 541 scope for improvement in the future, we believe that our model has shown promising new learning
 542 representation for realistic speaking and talking face generation models.

540 REFERENCES
541

- 542 Junyi Ao, Rui Wang, Long Zhou, Chengyi Wang, Shuo Ren, Yu Wu, Shujie Liu, Tom Ko, Qing Li,
543 Yu Zhang, Zhihua Wei, Yao Qian, Jinyu Li, and Furu Wei. Speecht5: Unified-modal encoder-
544 decoder pre-training for spoken language processing, 2022. URL <https://arxiv.org/abs/2110.07205>. 1
- 545 Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework
546 for self-supervised learning of speech representations, 2020. URL <https://arxiv.org/abs/2006.11477>. 4
- 547 James Betker. Tortoise text-to-speech, 2022. URL <https://github.com/neonbjb/tortoise-tts>. Accessed: [date you accessed the repository]. 1, 3, 7
- 548 Edresson Casanova, Julian Weber, Christopher Shulby, Arnaldo Candido Junior, Eren Gölge, and
549 Moacir Antonelli Ponti. Yourtts: Towards zero-shot multi-speaker tts and zero-shot voice conver-
550 sion for everyone, 2023. URL <https://arxiv.org/abs/2112.02418>. 5, 7
- 551 Edresson Casanova, Kelly Davis, Eren Gölge, Görkem Göknar, Iulian Gulea, Logan Hart, Aya Aljafari,
552 Joshua Meyer, Reuben Morais, Samuel Olayemi, and Julian Weber. Xtts: a massively multilingual
553 zero-shot text-to-speech model, 2024. URL <https://arxiv.org/abs/2406.04904>. 1, 3,
554 4, 5, 6, 7
- 555 Ken Chen, Sachith Seneviratne, Wei Wang, Dongting Hu, Sanjay Saha, Md. Tarek Hasan, Sanka
556 Rasnayaka, Tamasha Malepathirana, Mingming Gong, and Saman Halgamuge. Anifacediff:
557 High-fidelity face reenactment via facial parametric conditioned diffusion models, 2024. URL
558 <https://arxiv.org/abs/2406.13272>. 1
- 559 Bernhard Egger, William A. P. Smith, Ayush Tewari, Stefanie Wuhrer, Michael Zollhoefer, Thabo
560 Beeler, Florian Bernard, Timo Bolkart, Adam Kortylewski, Sami Romdhani, Christian Theobalt,
561 Volker Blanz, and Thomas Vetter. 3d morphable face models – past, present and future, 2020. URL
562 <https://arxiv.org/abs/1909.01815>. 1
- 563 Yuan Gan, Zongxin Yang, Xihang Yue, Lingyun Sun, and Yi Yang. Efficient emotional adaptation for
564 audio-driven talking-head generation, 2023. URL <https://arxiv.org/abs/2309.04946>.
565 1, 3, 7, 8
- 566 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.
567 Gans trained by a two time-scale update rule converge to a local nash equilibrium, 2018. URL
568 <https://arxiv.org/abs/1706.08500>. 6
- 569 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models, 2020. URL
570 <https://arxiv.org/abs/2006.11239>. 3, 5
- 571 Won Jang, Dan Lim, Jaesam Yoon, Bongwan Kim, and Juntae Kim. Univnet: A neural vocoder with
572 multi-resolution spectrogram discriminators for high-fidelity waveform generation, 2021. URL
573 <https://arxiv.org/abs/2106.07889>. 3
- 574 Youngjoon Jang, Ji-Hoon Kim, Junseok Ahn, Doyeop Kwak, Hong-Sun Yang, Yoon-Cheol Ju,
575 Il-Hwan Kim, Byeong-Yeol Kim, and Joon Son Chung. Faces that speak: Jointly synthesising
576 talking face and speech from text. In *Proceedings of the IEEE/CVF Conference on Computer
577 Vision and Pattern Recognition*, pp. 8818–8828, 2024. 1, 2, 3, 4
- 578 Prajwal K R, Rudrabha Mukhopadhyay, Jerin Philip, Abhishek Jha, Vinay Namboodiri, and C V
579 Jawahar. Towards automatic face-to-face translation. In *Proceedings of the 27th ACM International
580 Conference on Multimedia*, MM ’19. ACM, October 2019. doi: 10.1145/3343031.3351066. URL
581 <http://dx.doi.org/10.1145/3343031.3351066>. 3
- 582 Hasam Khalid, Shahroz Tariq, Minha Kim, and Simon S. Woo. Fakeavceleb: A novel audio-video
583 multimodal deepfake dataset, 2022. URL <https://arxiv.org/abs/2108.05080>. 6
- 584 Kevin Kilgour, Mauricio Zuluaga, Dominik Roblek, and Matthew Sharifi. Fréchet audio distance:
585 A metric for evaluating music enhancement algorithms, 2019. URL <https://arxiv.org/abs/1812.08466>. 7

- 594 Jaehyeon Kim, Sungwon Kim, Jungil Kong, and Sungroh Yoon. Glow-tts: A generative flow for text-
 595 to-speech via monotonic alignment search, 2020. URL <https://arxiv.org/abs/2005.11129>. 7
 596
- 597 Jaehyeon Kim, Jungil Kong, and Juhee Son. Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech, 2021. URL <https://arxiv.org/abs/2106.06103>. 1, 3
 598
- 601 Diederik P. Kingma and Max Welling. An introduction to variational autoencoders. *Foundations and Trends® in Machine Learning*, 12(4):307–392, 2019. ISSN 1935-8245. doi: 10.1561/2200000056.
 602 URL <http://dx.doi.org/10.1561/2200000056>. 6
 603
- 605 Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2022. URL <https://arxiv.org/abs/1312.6114>. 4
 606
- 607 Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. Hifi-gan: Generative adversarial networks for
 608 efficient and high fidelity speech synthesis, 2020. URL <https://arxiv.org/abs/2010.05646>. 3, 4
 609
- 611 Jiahe Liu, Youran Qu, Qi Yan, Xiaohui Zeng, Lele Wang, and Renjie Liao. Fréchet video motion
 612 distance: A metric for evaluating motion consistency in videos, 2024. URL <https://arxiv.org/abs/2407.16124>. 6
 613
- 615 Arsha Nagrani, Joon Son Chung, Weidi Xie, and Andrew Zisserman. Voxceleb: Large-scale speaker
 616 verification in the wild. *Computer Science and Language*, 2019. 6, 7
 617
- 618 K R Prajwal, Rudrabha Mukhopadhyay, Vinay P. Namboodiri, and C.V. Jawahar. A lip sync expert is
 619 all you need for speech to lip generation in the wild. In *Proceedings of the 28th ACM International
 620 Conference on Multimedia*, MM ’20. ACM, October 2020a. doi: 10.1145/3394171.3413532. URL
<http://dx.doi.org/10.1145/3394171.3413532>. 3
 621
- 622 KR Prajwal, Rudrabha Mukhopadhyay, Vinay P Namboodiri, and CV Jawahar. A lip sync expert is
 623 all you need for speech to lip generation in the wild. In *Proceedings of the 28th ACM international
 624 conference on multimedia*, pp. 484–492, 2020b. 7
 625
- 626 Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever.
 627 Robust speech recognition via large-scale weak supervision, 2022. URL <https://arxiv.org/abs/2212.04356>. 6
 628
- 629 Akshay Raina and Vipul Arora. Syncnet: Using causal convolutions and correlating objective for time
 630 delay estimation in audio signals, 2022. URL <https://arxiv.org/abs/2203.14639>. 3
 631
- 632 Yi Ren, Yangjun Ruan, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu. Fastspeech: Fast,
 633 robust and controllable text to speech, 2019. URL <https://arxiv.org/abs/1905.09263>.
 634 3
 635 Yurui Ren, Ge Li, Yuanqi Chen, Thomas H. Li, and Shan Liu. Pirenderer: Controllable portrait
 636 image generation via semantic neural rendering, 2021. URL <https://arxiv.org/abs/2109.08379>. 1
 637
- 638 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
 639 resolution image synthesis with latent diffusion models, 2022. URL <https://arxiv.org/abs/2112.10752>. 1
 640
- 642 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical
 643 image segmentation, 2015. URL <https://arxiv.org/abs/1505.04597>. 5
 644
- 645 Jonathan Shen, Ruoming Pang, Ron J. Weiss, Mike Schuster, Navdeep Jaitly, Zongheng Yang, Zhifeng
 646 Chen, Yu Zhang, Yuxuan Wang, RJ Skerry-Ryan, Rif A. Saurous, Yannis Agiomyrgiannakis, and
 647 Yonghui Wu. Natural tts synthesis by conditioning wavenet on mel spectrogram predictions, 2018.
 URL <https://arxiv.org/abs/1712.05884>. 3, 5

- 648 Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. First order
 649 motion model for image animation, 2020. URL <https://arxiv.org/abs/2003.00196>.
 650 1
- 651 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models, 2022. URL
 652 <https://arxiv.org/abs/2010.02502>. 5, 9
- 653 Michał Stypułkowski, Konstantinos Vougioukas, Sen He, Maciej Zięba, Stavros Petridis, and Maja
 655 Pantic. Diffused heads: Diffusion models beat gans on talking-face generation, 2023. URL
 656 <https://arxiv.org/abs/2301.03396>. 4
- 657 Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, and
 658 Sylvain Gelly. Towards accurate generative models of video: A new metric challenges, 2019.
 659 URL <https://arxiv.org/abs/1812.01717>. 6
- 660 Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves,
 661 Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw
 662 audio, 2016. URL <https://arxiv.org/abs/1609.03499>. 3
- 663 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz
 664 Kaiser, and Illia Polosukhin. Attention is all you need, 2023. URL <https://arxiv.org/abs/1706.03762>. 4, 9
- 665 Suzhen Wang, Lincheng Li, Yu Ding, Changjie Fan, and Xin Yu. Audio2head: Audio-driven one-shot
 666 talking-head generation with natural head motion, 2021. URL <https://arxiv.org/abs/2107.09293>. 3, 7, 8
- 667 Yaohui Wang, Di Yang, Francois Bremond, and Antitza Dantcheva. Latent image animator: Learning
 668 to animate images via latent space navigation. *arXiv preprint arXiv:2203.09043*, 2022. 4
- 669 Yuxuan Wang, RJ Skerry-Ryan, Daisy Stanton, Yonghui Wu, Ron J. Weiss, Navdeep Jaitly, Zongheng
 670 Yang, Ying Xiao, Zhifeng Chen, Samy Bengio, Quoc Le, Yannis Agiomyrgiannakis, Rob Clark,
 671 and Rif A. Saurous. Tacotron: Towards end-to-end speech synthesis, 2017. URL <https://arxiv.org/abs/1703.10135>. 3
- 672 Zhichao Wang, Mengyu Dai, and Keld Lundgaard. Text-to-video: a two-stage framework for
 673 zero-shot identity-agnostic talking-head generation. *arXiv preprint arXiv:2308.06457*, 2023. 1, 3
- 674 Chao Xu, Yang Liu, Jiazheng Xing, Weida Wang, Mingze Sun, Jun Dan, Tianxin Huang, Siyuan
 675 Li, Zhi-Qi Cheng, Ying Tai, and Baigui Sun. Facechain-imagineid: Freely crafting high-fidelity
 676 diverse talking faces from disentangled audio, 2024a. URL <https://arxiv.org/abs/2403.01901>. 4
- 677 Mingwang Xu, Hui Li, Qingkun Su, Hanlin Shang, Liwei Zhang, Ce Liu, Jingdong Wang, Yao Yao,
 678 and Siyu Zhu. Hallo: Hierarchical audio-driven visual synthesis for portrait image animation,
 679 2024b. URL <https://arxiv.org/abs/2406.08801>. 1, 2, 3, 4, 5, 7, 8
- 680 Sicheng Xu, Guojun Chen, Yu-Xiao Guo, Jiaolong Yang, Chong Li, Zhenyu Zang, Yizhong Zhang,
 681 Xin Tong, and Baining Guo. Vasa-1: Lifelike audio-driven talking faces generated in real time,
 682 2024c. URL <https://arxiv.org/abs/2404.10667>. 1, 2
- 683 Ryandhimas E. Zezario, Szu-Wei Fu, Chiou-Shann Fuh, Yu Tsao, and Hsin-Min Wang. Stoi-
 684 net: A deep learning based non-intrusive speech intelligibility assessment model, 2020. URL
 685 <https://arxiv.org/abs/2011.04292>. 7
- 686 Chenxu Zhang, Chao Wang, Jianfeng Zhang, Hongyi Xu, Guoxian Song, You Xie, Linjie Luo,
 687 Yapeng Tian, Xiaohu Guo, and Jiashi Feng. Dream-talk: Diffusion-based realistic emotional
 688 audio-driven method for single image talking face generation, 2023a. URL <https://arxiv.org/abs/2312.13578>. 1, 4
- 689 Fan Zhang, Valentin Bazarevsky, Andrey Vakunov, Andrei Tkachenka, George Sung, Chuo-Ling
 690 Chang, and Matthias Grundmann. Mediapipe hands: On-device real-time hand tracking, 2020.
 691 URL <https://arxiv.org/abs/2006.10214>. 4

- 702 Sibo Zhang, Jiahong Yuan, Miao Liao, and Liangjun Zhang. Text2video: Text-driven talking-
703 head video synthesis with personalized phoneme-pose dictionary. In *ICASSP 2022-2022 IEEE*
704 *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 2659–2663.
705 IEEE, 2022. 1, 3
- 706 Wenzuan Zhang, Xiaodong Cun, Xuan Wang, Yong Zhang, Xi Shen, Yu Guo, Ying Shan, and Fei
707 Wang. Sadtalker: Learning realistic 3d motion coefficients for stylized audio-driven single image
708 talking face animation, 2023b. URL <https://arxiv.org/abs/2211.12194>. 1, 3, 4, 7, 8
- 710 Yunxuan Zhang, Siwei Zhang, Yue He, Cheng Li, Chen Change Loy, and Ziwei Liu. One-shot face
711 reenactment, 2019. URL <https://arxiv.org/abs/1908.03251>. 1
- 712 Zhimeng Zhang, Lincheng Li, Yu Ding, and Changjie Fan. Flow-guided one-shot talking face
713 generation with a high-resolution audio-visual dataset. In *Proceedings of the IEEE/CVF Conference*
714 *on Computer Vision and Pattern Recognition*, pp. 3661–3670, 2021. 6
- 716 Hao Zhu, Wayne Wu, Wentao Zhu, Liming Jiang, Siwei Tang, Li Zhang, Ziwei Liu, and Chen Change
717 Loy. CelebV-hq: A large-scale video facial attributes dataset, 2022. URL <https://arxiv.org/>
718 [abs/2207.12393](https://arxiv.org/abs/2207.12393). 6
- 719 Vilém Zouhar, Clara Meister, Juan Luis Gastaldi, Li Du, Tim Vieira, Mrinmaya Sachan, and Ryan
720 Cotterell. A formal perspective on byte-pair encoding, 2024. URL <https://arxiv.org/>
721 [abs/2306.16837](https://arxiv.org/abs/2306.16837). 4
- 722
- 723
- 724
- 725
- 726
- 727
- 728
- 729
- 730
- 731
- 732
- 733
- 734
- 735
- 736
- 737
- 738
- 739
- 740
- 741
- 742
- 743
- 744
- 745
- 746
- 747
- 748
- 749
- 750
- 751
- 752
- 753
- 754
- 755