

Realistic Speech-to-Face Generation with Speech-Conditioned Latent Diffusion Model with Face Prior

Jinting Wang¹, Li Liu^{1*}, Jun Wang², Hei Victor Cheng³

¹The Hong Kong University of Science and Technology (Guangzhou)

²Tencent AI Lab

³Aarhus University

October 6, 2023

Abstract

Speech-to-face generation is an intriguing area of research that focuses on generating realistic facial images based on a speaker’s audio speech. However, state-of-the-art methods employing GAN-based architectures lack stability and cannot generate realistic face images. To fill this gap, we propose a novel speech-to-face generation framework, which leverages a **Speech-Conditioned Latent Diffusion Model**, called **SCLDM**. To the best of our knowledge, this is the first work to harness the exceptional modeling capabilities of diffusion models for speech-to-face generation. Preserving the shared identity information between speech and face is crucial in generating realistic results. Therefore, we employ contrastive pre-training for both the speech encoder and the face encoder. This pre-training strategy facilitates effective alignment between the attributes of speech, such as age and gender, and the corresponding facial characteristics in the face images. Furthermore, we tackle the challenge posed by excessive diversity in the synthesis process caused by the diffusion model. To overcome this challenge, we introduce the concept of residuals by integrating a statistical face prior to the diffusion process. This addition helps to eliminate the shared component across the faces and enhances the subtle variations captured by the speech condition. Extensive quantitative, qualitative, and user study experiments demonstrate that our method can produce more realistic face images while preserving the identity of the speaker

better than state-of-the-art methods. Highlighting the notable enhancements, our method demonstrates significant gains in all metrics on the AVSpeech dataset and Voxceleb dataset, particularly noteworthy are the improvements of 32.17 and 32.72 on the cosine distance metric for the two datasets, respectively.

Introduction

Generating realistic face portraits of speakers based on their audio speech has numerous applications, such as virtual anchors, teleconferencing, surveillance, and digital human animation. The field of speech-to-face generation aims to establish a meaningful mapping between audio and visual faces, resulting in high-quality and genuine face images. Previous studies have explored the relationship between human voice and facial structures, supporting the feasibility of speech-to-face generation [2, 10, 23]. Factors such as facial bone structure, joint configuration, and the tissues involved in sound production are closely intertwined with the shape and size of facial organs. Additionally, various factors including genetics, biology, and the environment impact both voice and face characteristics. For instance, gender, age, and ethnicity have been found to significantly influence both speech and facial appearance [10, 24, 12, 25].

Despite the potential of leveraging shared characteristics between face and voice for a speech-to-face generation, it remains a complex task due to the diverse nature

*Corresponding Author: avrillliu@hkust-gz.edu.cn.

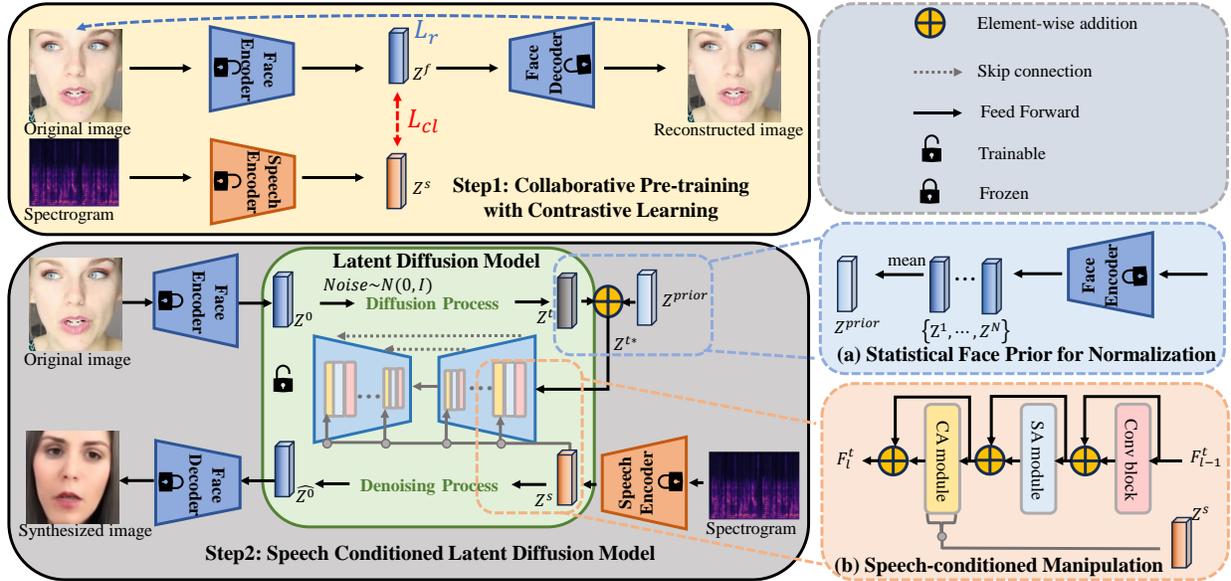


Figure 1: Details of the diffusion-based Speech-to-Face generation pipeline. It consists of two stages: **(1) Collaborative Pre-training with Contrastive Learning**. The speech encoder and face encoder take the audio and face image as input, respectively, to learn the representations, the captured embedding of face and speech are aligned by contrastive learning. **(2) SCLDM**. The aligned face embedding is gradually destroyed by Gaussian noise in the diffusion process within iterative time steps. Then it is combined with common face stuff provided by face prior and fed into LDM for face generation. The aligned speech embedding is utilized as conditioned variation to manipulate the face generation in the training process of LDM.

of human faces and the various speaking styles they encompass. There are two primary challenges: 1) How can a model learn the common characteristics shared by voice and face? 2) How can a model produce visually authentic facial images while preserving attribute details? To address the first challenge, some research utilizes extracted speech features to generate face images based on various existing generative models [17, 28, 6, 8]. Although delicate designs have been introduced to align the speech features and face features, a substantial gap persists between speech embedding and visual embedding, resulting in subpar face generation. Furthermore, to achieve speech-to-face conversion, GAN-based architectures are mainly employed in previous studies. However, methods based on GANs require simultaneous optimization of a generator and a discriminator, such a training process lacks stability and is prone to an effect known as mode

collapse [5]. Therefore, the generated speaker faces are of limited image quality.

Inspired by the Stable Diffusion [21], which employs latent diffusion models (LDMs) to achieve high-quality image generation, we propose a novel speech-to-face generation framework, which incorporates a **Speech-Conditioned LDM (SCLDM)**. Our framework is designed to address the limitations observed in previous works and to advance the progress of speech-to-face generation. Moreover, we aim to achieve state-of-the-art (SOTA) face generation quality that aligns seamlessly with the attributes of the spoken speech. To achieve this, we explore speech-conditional face manipulations with LDMs, which have not been explored before. Specifically, an LDM conditioned on the latent speech embedding is learned for face image generation. For accurate speech-conditioned face manipulations, we leverage the

face-aligned speech latent embedding space pre-trained by contrastive learning as the condition. Moreover, we introduce a statistical face prior to the LDM. This statistical prior serves as the shared component of the faces and assists the LDM to emphasize the subtle changes present in speech while alleviating its diversity for identity attribute preservation.

The main contributions of this work are as follows:

- A novel speech-conditioned diffusion-based network is proposed for speech-to-face generation. To the best of our knowledge, this is the first attempt to develop an LDM for this task.
- To achieve accurate and delicate speech-conditioned face manipulation, a contrastive pre-training is employed for aligning speech and face representations. Then the aligned speech latent embedding is utilized as a condition to generate a face image corresponding to the speech.
- To enhance the speech manipulation, a statistical face prior is introduced to the diffusion model, which allows for generating a more realistic face image conforming to the corresponding speech.
- Extensive experiments are conducted to validate the performance of our proposed method, the results show that SCLDM achieves the best among all speech-to-face generation methods.

Related Work

Diffusion Model

Diffusion models have demonstrated their SOTA generation performance in various tasks, including image generation [21, 5], speech generation [11, 13], and video generation [14, 3]. In diffusion-based frameworks, one of the major concerns is to improve the efficiency of diffusion models due to their iterative generation process in a high-dimensional data space. The latent diffusion model is one of the solutions that apply diffusion models in a small latent space [21], which was first presented in image generation [21]. An LDM model is trained as a flexible image generator with general conditioning inputs such as semantic map, text, and image. This approach inspired further

study on conditional LDMs in different domains, such as text-to-audio generation [9], text-to-image generation and editing [22], text-to-video generation and editing [3], conditional image-to-video generation [16], and audio-to-video generation [31]. In this work, our focus is on leveraging the LDM models and proposing a diffusion-based framework for speech-to-face generation.

Speech-to-Face Generation

Audio-visual cross-modal learning, particularly speech-to-face generation, has gained significant attention in recent years. The goal of speech-to-face generation is to generate realistic facial appearances that correspond to the audio speech input. Many existing methods in this field employ GAN-based frameworks. One such framework is Wav2Pix [6], which proposes a speech-conditioned face generation framework. It consists of a speech encoder, a generator, and a discriminator. Wav2Pix is relatively simplistic and overlooks the preservation of identity information during the generation process. To address the preservation of identity information, Wen et al. [28] design a network capable of generating faces from voices by matching the identities of the generated faces with those of the speakers. This approach aims to preserve certain biometric characteristics of the speakers. Similarly, Fang et al. [8] incorporate speaker identity information to learn identity representations across different modalities. While explicitly modeling the identity relevance between audio and visual modalities is beneficial for ensuring the authenticity of generated face images, it has limitations when attempting to generate faces of different identities. On the other hand, Choi et al. [4] propose a two-stage framework comprising an inference stage for cross-modal matching and a generation stage for speech-conditioned face generation. This approach allows for more flexibility in generating faces with different identities. In the work of Oh et al. [17], a one-stage method is proposed, incorporating a multivariate mixing loss function to enforce consistency between speech and face features. This facilitates the learning of shared attribute information between the voice and the face. The existing methods for speech-to-face generation have shown some progress, but they still have limitations, particularly in terms of generation quality. Additionally, the alignment between speech and face modalities has not been adequately explored, which can

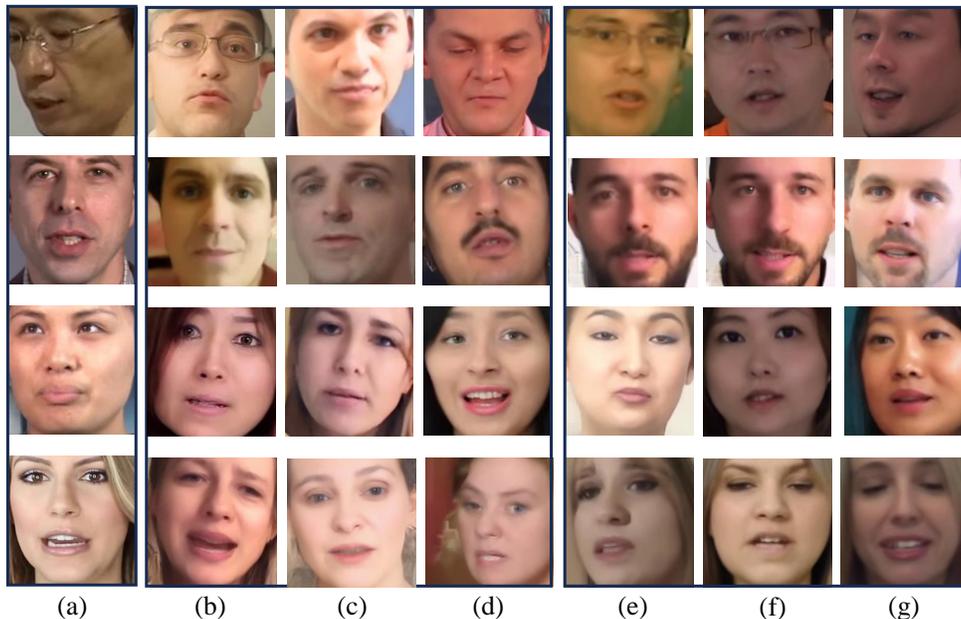


Figure 2: Qualitative comparison of the face prior. (a) Face images cropped from the video frame. (b-c) Top-3 generated faces of the same speech without prior normalization. (e-g) Top-3 generated faces of the same speech with prior normalization.

impact the overall coherence and accuracy of the generated results. In this work, we propose a speech-to-face generation network that employs LDM as the generation model, and speech-face alignment is conducted to enhance modality matching.

Proposed Method

In this work, we propose a diffusion-based framework for generating realistic face images from speech input, which exploits the SOTA LDM for speech-conditioned face synthesis. Considering shared information preservation between speech and face image modality is crucial for realistic image generation, we employ a pre-training step that utilizes a contrastive learning strategy to facilitate alignment between speech and face embedding. The aligned speech embedding is then used as a conditioning input to manipulate the face generation process of the LDM. Furthermore, we introduce a statistical face prior into the latent code of the LDM. This face prior provides a face

template that guides the model’s attention towards variations corresponding to the input speech.

Figure 1 provides detailed insights of the proposed diffusion-based speech-to-face generation framework, SCLDM, it can be seen that it consists of two stages. In the first stage, the face encoder, face decoder, and speech encoder are collaboratively pre-trained using a contrastive learning approach. This pre-training step ensures that the model learns to align speech and face embeddings effectively. In the second stage, the pre-trained face encoder and speech encoder are used to encode aligned embeddings for training the speech-conditioned LDM. The pre-trained face decoder is kept for the testing stage, where it is used to reconstruct the face images.

Collaborative Pre-training with Contrastive Learning

Inspired by the great success of contrastive pre-training in various cross-modal applications [1, 18], contrastive pre-

training is employed in this stage to facilitate speech-face alignment.

As shown in Figure 1, given an audio clip of a speaker and the corresponding face image, a speech encoder $E_{Speech}(\cdot)$ and a face encoder $E_{Face}(\cdot)$ are employed to extract the speech embedding $Z^s \in \mathcal{R}^d$ and face embedding $Z^f \in \mathcal{R}^d$, respectively, where d is the dimension of the embedding vectors. To align the speech embedding and face embedding, a symmetric cross-entropy loss [20] is applied, leveraging contrastive learning techniques. Specifically, we use VGGFace [19] as the face encoder, and a model combined with CNN (the speech encoder architecture in Speech2Face [17]) and CBAM module [29] as the speech encoder. Since we apply the diffusion model in a latent space, a face decoder is required to reconstruct face images from the latent representation. A CNN-based model symmetrical to the VGGFace is designed as the face decoder $D_{Face}(\cdot)$. Finally, a combination of MAE loss and LPIPS loss [30] is used as the reconstruction loss.

The objective function of the collaborative pre-training L_C is defined as:

$$L_C = L_{cl} + L_r, \quad (1)$$

where L_{cl} and L_r denote the contrastive loss and reconstruction loss, respectively. After training, the speech representation Z^s of a random audio sample is used for providing shared-attributes information.

Explicit Face Prior as Normalization

Motivations. While the speech-conditioned LDM can generate positive results, we observe that it occasionally falls short in producing realistic face images that accurately match the speaker attributes. This can be seen in Figure 2, where an LDM trained with speech condition is employed to create speaker portraits. However, the outputs originating from the same speech clip exhibit diverse characteristics, rendering them distinct from one another. This diversity has a detrimental impact on the visual quality, hindering their use in real-world applications. This shortcoming is attributed to the inherent diversity present in the latent space learned by the LDM, which is a result of the complex data distribution. During the denoising process for face synthesis, the latent is randomly sam-

pled from the latent space, resulting in the variance observed in the generated images. This observation suggests that relying solely on speech conditions to exert control is insufficient for generating the desired face images that accurately correspond to the speech clips. This highlights the necessity for supplementary mechanisms to guide the generation process.

Latent Diffusion with Face Prior as Normalization.

To address the aforementioned issue and generate more realistic face images conforming to the speech, we propose to explicitly introduce a statistical face prior to the LDM for normalization purposes. This approach exploits the concept of residuals to eliminate the shared component found in the faces and effectively emphasizes the subtle variations that are captured by the speech embedding.

As illustrated in Figure 1, the statistical face prior Z^{prior} is combined with the latent code as a weighted sum:

$$Z^{t*} = Z^t + \beta Z^{prior}. \quad (2)$$

Incorporating the face prior maintains the fundamental mechanisms of the LDM. The inclusion of the prior enriches the LDM’s learning process by guiding it to learn the difference between the input face images and the introduced face prior within the diffusion process. In the denoising process, we effectively alter the distribution of the latent code that we sample from, shifting from $\mathcal{N}(\mathbf{0}, \mathbf{I})$ to $\mathcal{N}(\beta Z^{prior}, \mathbf{I})$. The choice of weight $\beta = 0.01$ is based on empirical knowledge, and we will provide an ablation study of the weight β in the experiment part.

The face prior is constructed by averaging the features extracted from the pre-trained face encoder E_{Face} on given gender-balance data:

$$Z^{prior} = \frac{1}{N} \sum_1^N E_{Face}(f), \quad (3)$$

where f denotes a face image and N is the total number of face images. Through experiments we observe that when N gradually increases, the face prior tends to converge, suggesting that the calculated prior becomes representative of the shared characteristics. As indicated in Table 1, we take $N = 10000$ in this work.

With the incorporation of the face prior into LDM as normalization, SCLDM demonstrates improved precision. It is capable of synthesizing face images that better

N_1	N_2	$L1(N_1, N_2)$
100	500	11.86
500	1000	5.21
1000	5000	3.35
5000	10000	1.28
10000	15000	0.67

Table 1: $L1$ distance between face prior calculated with different number N of face image.

preserve the identity of the speaker, as depicted in Figure 2.

Speech-Conditioned Latent Diffusion Model with Face Prior

In speech-to-face generation, the speaker portrait is synthesized given the speech clip as the condition. With speech-conditioned LDM, we are interested in generating face images conforming to the speech. Specifically, in the training stage, the face image is embedded into latent representation Z^f (Z^0 in diffusion process) by the trained face encoder E_{Face} , and then the face embedding is destroyed into a noised vector Z^t with Gaussian distribution in the diffusion process after t time steps, which is denoted as:

$$Z^t := \alpha^t * Z^0 + (1 - \alpha^t) * \epsilon, \quad (4)$$

where $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ denotes the injected noise, and α^t represents the noise level at the t time step. With the prior normalization mentioned above, Z^t would be transformed into Z^{t*} by introducing the face prior Z^{prior} . For noise estimation, a denoising model comprised of the UNet backbone and attention mechanism is employed in the denoising process. The condition vector Z^s captured by the pre-trained speech encoder E_{Speech} would be mapped to the intermediate layers of the UNet via a cross-attention mechanism. As shown in Figure 1, for a specific layer l , the speech-conditioned manipulation is that the layer input F_{l-1}^t successively processed by a convolution block and a self-attention module, and then interacted with the condition vector Z^s via a cross-attention module. The speech-conditioned operation is defined as:

$$CA(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V, \quad (5)$$

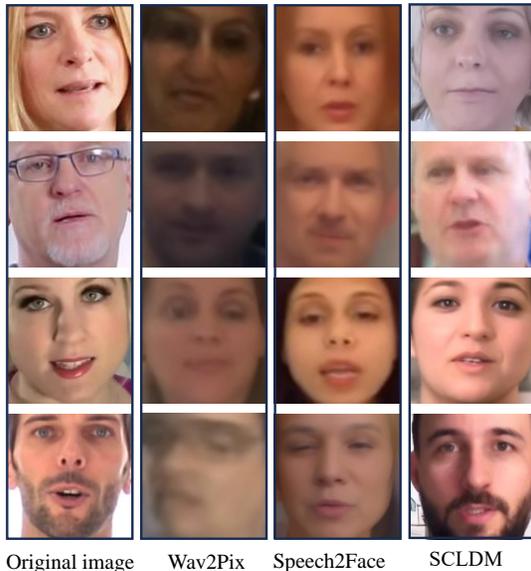


Figure 3: Qualitative comparison of our model and SOTA methods on the AVSpeech dataset.

where CA is the abbreviation of cross attention module, Q is the linear transformation of the output of self-attention module (SA). While K and V are transformed from condition vector Z^s by two linear layers. Based on speech-face pairs, the conditional LDM is trained via

$$L_{LDM} := \mathbb{E}_{\epsilon, Z^{t*}, t} \left[\|\epsilon - \epsilon_{\theta}(Z^s, Z^{t*}, t)\|^2 \right], \quad (6)$$

where ϵ_{θ} denotes the optimized denoising model, and θ denotes its parameters.

In the generation process, starting from injecting face prior Z^{prior} into Gaussian noise $Z^t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, the trained denoising model gradually generates latent face \hat{Z}^0 with the speech condition, and then it is fed into the trained face decoder and the face reconstruction is obtained. Details of the training and generation process are provided in the supplementary materials.

Method	Feature Similarity			Identity Preservation	
	L1 ↓	L2 ↓	cos ↓	gender (%) ↑	age (%) ↑
Wav2Pix	144.72	24.32	82.51	67.4	41.3
Speech2Face	67.18	3.94	46.97	95.6	65.2
SCLDM (Ours)	35.01	1.48	12.81	98.8	84.5

Table 2: Comparison results on AVSeech dataset.

Method	Feature Similarity			Identity Preservation	
	L1 ↓	L2 ↓	cos ↓	gender (%) ↑	age (%) ↑
Wav2Pix	137.58	22.19	79.36	74.5	49.6
Speech2Face	66.46	2.77	44.38	96.1	69.4
Wen <i>et al.</i>	59.82	2.41	42.54	97.4	72.5
SCLDM (Ours)	27.10	1.09	11.54	99.4	88.6

Table 3: Comparison results on VoxCeleb dataset.

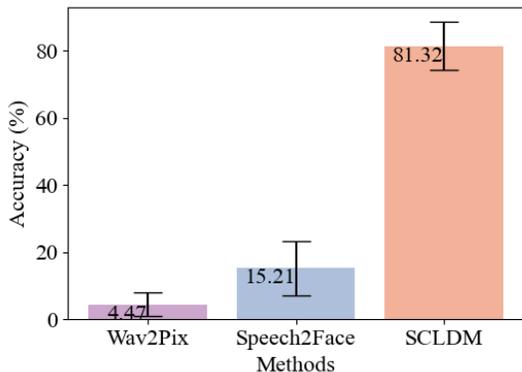


Figure 4: Results of the user study. Among the three methods, our method achieves the highest accuracy for users evaluation, in terms of image quality and identity preservation.

Experiments

Dataset

Throughout the experiments, two datasets are used. **AVSpeech.** AVSpeech [7] is a large-scale audio-visual dataset collected from YouTube, which consists of 2.8m video clips. It is the significant diversity present in the included face images which were extracted from videos captured “in the wild”.

VoxCeleb. VoxCeleb [15] contains 1251 speakers, which span a wide range of ethnicities, accents, professions, and ages. The nationality and gender of each speaker are also provided in the dataset.

Evaluation Metrics

Feature Similarity. Following [17], we measure the cosine, L1, and L2 distance between the features of the true face image and generated face image extracted by VG-GFace [19], a trained face recognition network.

Identity Preservation. We employ Face++¹, a commercial API for face attribute recognition, to evaluate face attributes, including age and gender. Note that the age accuracy is computed within a 10-years range.

Implementation Details

Following Speech2Face [17], we use 6 seconds of audio taken from the video clip and then process it into a spectrogram by STFT. The face images are resized to 256×256 pixels. For collaborative pre-training, we set a learning rate for the face encoder and face decoder of 0.0001, and 0.001 for the speech encoder. We adopt LDM [21] as a generator since it achieves a good balance between quality and speed. In the LDM training stage, the

¹<https://www.faceplusplus.com/attributes>.

Method			Feature Similarity			Identity Preservation	
base	CP	PN	L1 ↓	L2 ↓	cos ↓	gender (%) ↑	age (%) ↑
✓			56.39	3.30	29.83	95.9	74.7
✓	✓		47.61	2.69	21.79	97.2	82.6
✓		✓	44.27	2.38	20.41	96.4	80.3
✓	✓	✓	35.01	1.48	12.81	98.8	84.5

Table 4: Ablation results on AVSpeech dataset. Abbreviations “CP” and “PN” denote Collaborative Pre-training and Prior Normalization, respectively.

face encoder and speech encoder are frozen, and optimization is performed with a learning rate of $2e-5$.

Comparative Study

We compare our proposed method with three SOTA speech-to-face generation methods, i.e., Wav2Pix [6], Speech2Face [17], and the work of Wen *et al.* [28]. We conduct experiments using the default settings and official implementations for Wav2Pix [6] and the work of Wen *et al.* [28]. Since the code of Speech2Face [17] is not available, we reproduce it according to the paper. We only compare with Wen *et al.* on the Voxceleb dataset since the identity information of speakers is lacking in AVSpeech dataset.

Quantitative Comparison. The comparison results on AVSpeech and VoxCeleb datasets are reported in Table 2 and Table 3, respectively. Our method outperforms all the competitors in all metrics. Specifically, the cosine distance of our method achieves 12.81 on AVSpeech test set and 13.54 on VoxCeleb test set. The gender recognition accuracy achieves 98.8 and 98.4 on the two dataset. These results verify our effectiveness in speech-conditioned quality.

Qualitative Comparison. The qualitative comparison is presented in Figure 3. It is observed that our framework is capable of synthesizing realistic outputs consistent with the attributes of speaker compared with Wav2Pix [6] and Speech2Face [17].

User Study. We conduct a user study with 20 human evaluators to measure the effectiveness of the methods perceptually. We randomly sample 50 speech clips in the AVSpeech test set, then synthesize the speaker’s face images given the speech. The evaluators are provided with the true face and the generated face images. They are

asked to choose the best image based on 1) image realism, 2) identity preservation. Figure 4 showed the mean and standard deviation of the results, one can see that our framework significantly outperforms the existing STOA methods, which verifies the effectiveness in both image quality and identity preservation.

Ablation Study

Ablation Experiment on Model Components. We conduct ablation studies on the AVSpeech dataset to validate the effectiveness of different components. The comparison results of different versions are listed in Table 4. It can be seen that accuracy gains a lot in both gender and age attributes with collaborative pre-training, which indicates that the identity information shared between face and speech is aligned and preserved by contrastive learning. With face prior normalization, the feature distances between generated images and original faces are lower, which implies that the synthesized results present similar appearance as original faces.

Additionally, we provide visual examples in Figure 5 to illustrate the generated face images. It is evident that with the collaborative pre-training and prior normalization, the generated face images exhibit a similar appearance and attributes to the speaker in the corresponding speech.

Ablation Experiment on Prior Normalization Weight. We conduct experiments to evaluate the impact of varying the weight of the face prior, β , during both the training and inference stages. The results of these experiments are shown in Figure 6. In this analysis, we calculate the mean and variance of the cosine distance between the top-3 generated faces and the original images within the AVSpeech test set. This finding suggests that incorporating the face prior with a weight of 0.01 during training and inference

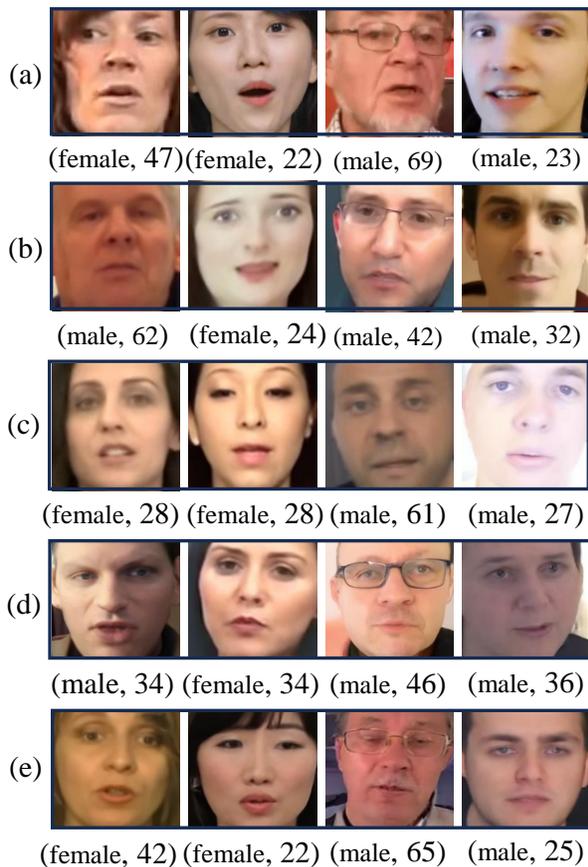


Figure 5: Qualitative comparison of ablation studies. (a) Original face images cropped from the video frame. (b) Generated face by speech-conditioned LDM (base). (c) Generated face by speech-conditioned LDM with collaborative pre-training. (d) Generated face by speech-conditioned LDM with face prior normalization. (e) Generated face by conditional LDM with collaborative pre-training and face prior normalization (SCLDM).

achieves the best quality for the generated faces.

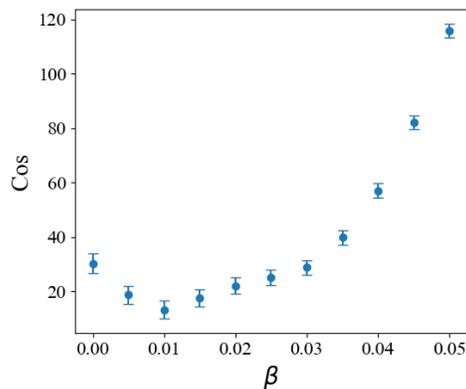


Figure 6: Ablation study on the weight β of face prior normalization.

Conclusion and Discussion

In this work, we propose a speech-conditioned latent diffusion model for speech-to-face generation, with the aim of generating realistic and identity-preserving face images given speech clips in real-world scenarios. With the assistance of contrastive pre-training, the speech encoder and face encoder provide aligned embedding that preserves the shared identity information, which enables effective manipulation of the generated faces. In addition, we utilize a statistical face prior in combination with a residual trick, which serves as normalization, removing the common components shared among faces. With the face prior, the LDM can better focus on capturing the subtle variations in the speech condition, resulting in more accurate and realistic face generation. Extensive experiments demonstrate that our method achieves new SOTA performance on speech-to-face generation. We believe that our idea of utilizing latent diffusion model with a statistical face prior would be a good inspiration for future works in different face generation tasks.

Limitations and Future Work. Our framework focuses on exploiting diffusion models for speech-to-face generation, therefore, the face encoder and speech encoder are both built following the structures in previous work. However, the representation ability of the embeddings, partic-

ularly the speech embedding poses a high influence on conditioned generation. Future work should focus on finding better embeddings aligning speech and face features. A possible direction is to employ a large speech model to find audio signal embedding. On the other end, the face prior introduced in our work is combined with the latent code of diffusion model with equal weights across samples, which may affect the diversity when the speech condition has limited variation. The ways of constructing and employing the face prior can be further explored for a more realistic generation.

References

- [1] M. Afham, I. Dissanayake, D. Dissanayake, A. Dharmasiri, K. Thilakarathna, and R. Rodrigo. Crosspoint: Self-supervised cross-modal contrastive learning for 3d point cloud understanding. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9902–9912, 2022.
- [2] P. Belin, S. Fecteau, and C. Bedard. Thinking the voice: neural correlates of voice perception. *Trends in Cognitive Sciences*, 8(3):129–135, 2004.
- [3] A. Blattmann, R. Rombach, H. Ling, T. Dockhorn, S. W. Kim, S. Fidler, and K. Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22563–22575, 2023.
- [4] H.-S. Choi, C. Park, and K. Lee. From inference to generation: End-to-end fully self-supervised generation of human face from speech. In *International Conference on Learning Representations*, 2019.
- [5] P. Dhariwal and A. Nichol. Diffusion models beat gans on image synthesis. In *Advances in Neural Information Processing Systems*, pages 8780–8794, 2021.
- [6] A. Duarte, F. Roldan, M. Tubau, J. Escur, S. Pascual, A. Salvador, E. Moledano, K. McGuinness, J. Torres, and X. Giro-i Nieto. Wav2pix: Speech-conditioned face generation using generative adversarial networks. In *2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8633–8637, 2019.
- [7] A. Ephrat, I. Mosseri, O. Lang, T. Dekel, K. Wilson, A. Hassidim, W. T. Freeman, and M. Rubinstein. Looking to listen at the cocktail party: A speaker-independent audio-visual model for speech separation. *arXiv preprint arXiv:1804.03619*, 2018.
- [8] Z. Fang, Z. Liu, T. Liu, C.-C. Hung, J. Xiao, and G. Feng. Facial expression gan for voice-driven face generation. *The Visual Computer*, 38(3):1151–1164, 2022.
- [9] R. Huang, Z. Zhao, H. Liu, J. Liu, C. Cui, and Y. Ren. Prodiff: Progressive fast diffusion model for high-quality text-to-speech. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 2595–2605, 2022.
- [10] M. Kamachi, H. Hill, K. Lander, and E. Vatikiotis-Bateson. Putting the face to the voice’: Matching identity across modality. *Current Biology*, 13(19):1709–1714, 2003.
- [11] H. Kim, S. Kim, and S. Yoon. Guided-tts: A diffusion model for text-to-speech via classifier guidance. In *International Conference on Machine Learning*, pages 11119–11133, 2022.
- [12] D. Kwasny and D. Hemmerling. Gender and age estimation methods based on speech using deep neural networks. *Sensors*, 21(14):4785, 2021.
- [13] J. Lee, J. S. Chung, and S.-W. Chung. Imaginary voice: Face-styled diffusion model for text-to-speech. In *2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5, 2023.
- [14] Z. Luo, D. Chen, Y. Zhang, Y. Huang, L. Wang, Y. Shen, D. Zhao, J. Zhou, and T. Tan. Videofusion: Decomposed diffusion models for high-quality video generation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10209–10218, 2023.
- [15] A. Nagrani, J. S. Chung, and A. Zisserman. Voxceleb: a large-scale speaker identification dataset. *arXiv preprint arXiv:1706.08612*, 2017.

- [16] H. Ni, C. Shi, K. Li, S. X. Huang, and M. R. Min. Conditional image-to-video generation with latent flow diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18444–18455, 2023.
- [17] T.-H. Oh, T. Dekel, C. Kim, I. Mosseri, W. T. Freeman, M. Rubinstein, and W. Matusik. Speech2face: Learning the face behind a voice. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7539–7548, 2019.
- [18] M. Parelli, A. Delitzas, N. Hars, G. Vlassis, S. Anagnostidis, G. Bachmann, and T. Hofmann. Clip-guided vision-language pre-training for question answering in 3d scenes. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5606–5611, 2023.
- [19] Z. Qawaqneh, A. A. Mallouh, and B. D. Barkana. Deep convolutional neural network for age estimation based on vgg-face model. *arXiv preprint arXiv:1709.01664*, 2017.
- [20] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763, 2021.
- [21] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis with latent diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022.
- [22] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. L. Denton, K. Ghasemipour, R. Gontijo Lopes, B. Karagol Ayan, T. Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. In *Advances in Neural Information Processing Systems*, pages 36479–36494, 2022.
- [23] J. Wang, X. Hu, L. Liu, W. Liu, M. Yu, and T. Xu. Attention-based residual speech portrait model for speech to face generation. *arXiv preprint arXiv:2007.04536*, 2020.
- [24] J. Wang, J. Liu, L. Zhao, S. Wang, R. Yu, and L. Liu. Acoustic-to-articulatory inversion based on speech decomposition and auxiliary feature. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, page 4808–4812. IEEE, 2022.
- [25] J. Wang, Z. Wang, X. Hu, X. Li, Q. Fang, and L. Liu. Residual-guided personalized speech synthesis based on face image. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, page 4743–4747. IEEE, 2022.
- [26] J. Wang, Y. Zhao, H. Fan, T. Xu, Q. Li, S. Li, and L. Liu. Memory-augmented contrastive learning for talking head generation. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, page 1–5. IEEE, 2023.
- [27] J. Wang, Y. Zhao, L. Liu, T. Xu, Q. Li, and S. Li. Emotional talking head generation based on memory-sharing and attention-augmented networks. *arXiv preprint arXiv:2306.03594*, 2023.
- [28] Y. Wen, R. Singh, and B. Raj. Face reconstruction from voice using generative adversarial networks. In *the 33rd International Conference on Neural Information Processing Systems*, pages 5265–5274, 2019.
- [29] S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon. Cbam: Convolutional block attention module. In *European Conference on Computer Vision (ECCV)*, pages 3–19, 2018.
- [30] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 586–595, 2018.
- [31] L. Zhu, X. Liu, X. Liu, R. Qian, Z. Liu, and L. Yu. Taming diffusion models for audio-driven co-speech gesture generation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10544–10553, 2023.