# SadTalker: Learning Realistic 3D Motion Coefficients for Stylized Audio-Driven Single Image Talking Face Animation

Wenxuan Zhang[*, 1]     Xiaodong Cun[*, 2]     Xuan Wang[3]     Yong Zhang[2]     Xi Shen[2]

Yu Guo[1]     Ying Shan[2]     Fei Wang[†, 1]

[1] Xi'an Jiaotong University     [2] Tencent AI Lab     [3] Ant Group

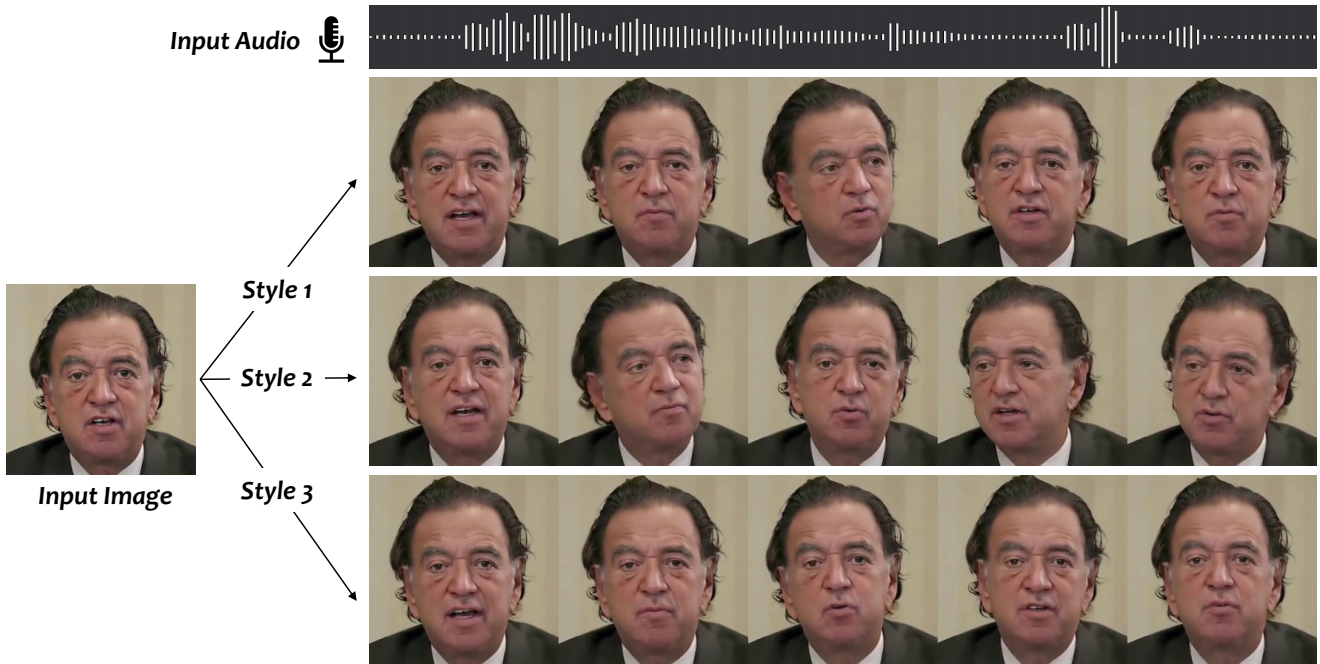https://sadtalker.github.io



Figure 1. The proposed SadTalker produces diverse, realistic, synchronized talking videos from an input audio and a single reference image.

## Abstract

*Generating talking head videos through a face image and a piece of speech audio still contains many challenges. i.e., unnatural head movement, distorted expression, and identity modification. We argue that these issues are mainly because of learning from the coupled 2D motion fields. On the other hand, explicitly using 3D information also suffers problems of stiff expression and incoherent video. We present SadTalker, which generates 3D motion coefficients (head pose, expression) of the 3DMM from audio and implicitly modulates a novel 3D-aware face render for talking head generation. To learn the realistic motion coefficients, we explicitly model the connections between audio and different types of motion coefficients individually. Precisely, we present ExpNet to learn the accurate facial expression from audio by distilling both coefficients and 3D-rendered faces. As for the head pose, we design PoseVAE via a conditional VAE to synthesize head motion in different styles. Finally, the generated 3D motion coefficients are mapped to the unsupervised 3D keypoints space of the proposed face render, and synthesize the final video. We conducted extensive experiments to demonstrate the superiority of our method in terms of motion and video quality.*

---

[*] Equal Contribution
[†] Corresponding Author

1

## 1. Introduction

Animating a static portrait image with speech audio is a challenging task and has many important applications in the fields of digital human creation, video conferences, *etc*. Previous works mainly focus on generating lip motion [2, 3, 30, 31, 51] since it has a strong connection with speech. Recent works also aim to generate a realistic talking face video containing other related motions, *e.g.*, head pose. Their methods mainly introduce 2D motion fields by landmarks [52] and latent warping [39, 40]. However, the quality of the generated videos is still unnatural and restricted by the preference pose [17, 51], month blur [30], identity modification [39, 40], and distorted face [39, 40, 49].

Generating a natural-looking talking head video contains many challenges since the connections between audio and different motions are different. *i.e.*, the lip movement has the strongest connection with audio, but audio can be talked via different head poses and eye blink. Thus, previous facial landmark-based methods [2, 52] and 2D flow-based audio to expression networks [39, 40] may generate the distorted face since the head motion and expression are not fully disentangled in their representation. Another popular type of method is the latent-based face animation [3, 17, 30, 51]. Their methods mainly focus on the specific kind of motions in talking face animation and struggle to synthesize high-quality video. Our observation is that the 3D facial model contains a highly decoupled representation and can be used to learn each type of motion individually. Although a similar observation has been discussed in [49], their methods also generate inaccurate expressions and unnatural motion sequences.

From the above observation, we propose SadTalker, a **S**tylized **A**udio-**D**riven **Talk**ing-head video generation system through implicit 3D coefficient modulation. To achieve this goal, we consider the motion coefficients of the 3DMM as the intermediate representation and divide our task into two major components. On the one hand, we aim to generate the realistic motion coefficients (*e.g.*, head pose, lip motion, and eye blink) from audio and learn each motion individually to reduce the uncertainty. For expression, we design a novel audio to expression coefficient network by distilling the coefficients from the lip motion only coefficients from [30] and the perceptual losses (lip-reading loss [1], facial landmark loss) on the reconstructed rendered 3d face [5]. For the stylized head pose, a conditional VAE [6] is used to model the diversity and life-like head motion by learning the residual of the given pose. After generating the realistic 3DMM coefficients, we drive the source image through a novel 3D-aware face render. Inspired by face-vid2vid [42], we learn a mapping between the explicit 3DMM coefficients and the domain of the unsupervised 3D keypoint. Then, the warping fields are generated through the unsupervised 3D keypoints of source and driving and it warps the reference image to generate the final videos. We train each sub-network

of expression generation, head poses generation and face renderer individually and our system can be inferred in an end-to-end style. As for the experiments, several metrics show the advantage of our method in terms of video and motion methods.

The main contribution of this paper can be summarized as:

- We present SadTalker, a novel system for a stylized audio-driven single image talking face animation using the generated realistic 3D motion coefficients.

- To learn the realistic 3D motion coefficient of the 3DMM model from audio, ExpNet and PoseVAE are presented individually.

- A novel semantic-disentangled and 3D-aware face render is proposed to produce a realistic talking head video.

- Experiments show that our method achieves state-of-the-art performance in terms of motion synchronization and video quality.

## 2. Related Work

**Audio-driven Single Image Talking Face Generation.** Early works [3, 30, 31] mainly focus on producing accurate lip motion with a perception discriminator. Since the real videos contain many different motions, ATVGnet [2] uses the facial landmark as the intermediate representation to generate the video frames. A similar approach has been proposed by MakeItTalk [52], differently, it disentangles the content and speaker information from the input audio signal. Since facial landmarks are still a highly coupled space, generating the talking head in the disentangled space is also popular recently. PC-AVS [51] disentangles the head pose and expression using implicit latent code. However, it can only produce low-resolution image and need the control signal from another video. Audio2Head [39] and Wang *et al*. [40] get inspiration from the video-driven method [36] to produce the talking-head face. However, these head movements are still not vivid and produce distorted faces with inaccurate identities. Although there are some previous works [33, 49] use 3DMMs as an intermediate representation, their method still faces the problem of inaccurate expressions [33] and obvious artifacts [49].

**Audio-driven Video Portrait.** Our task is also related to visual dubbing, which aims to edit a portrait video through audio. Different from audio-driven single image talking face generation, this task is typically required to be trained and edited on the specific video. Following previous work of deep video portrait [19], these methods utilize 3DMM information for face reconstruction and animation. AudioDVP [45], NVP [38], AD-NeRF [11] learn to reenact the expression to
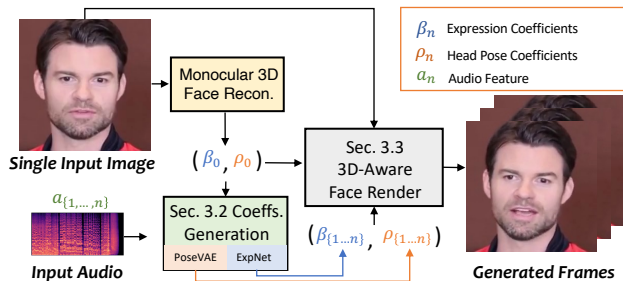
Figure 2. Main pipeline. Our method uses the coefficients of 3DMM as intermediate motion representation. To this end, we first generate realistic 3D motion coefficients (facial expression $\beta$, head pose $\rho$) from audio, then these coefficients are used to implicitly modulate the 3D-aware face render for final video generation.

edit the mouth shape. Beyond lip movement, *i.e.*, the head motions [23,48], emotional talking face [18] also get attention. The 3DMM-based method plays an important role in these tasks since it is practical to fit the 3DMM parameters from a video clip. Although these methods achieve satisfactory results in personalized video, their method can not be applied to arbitrary photos and in the wild audio.

**Video-Driven Single Image Talking Face Generation.** This task is also known as face reenactment or face animation, which aims to transfer the motion of the source image to the target person. It has been widely explored [14, 29, 33, 36, 37, 41, 42, 44, 47, 50] recently. Previous works also learn a shared intermediate motion representation from the source image and the target, which can be roughly divided into the landmark [41] and the unsupervised landmark-based methods [14,36,42,50], 3DMM based methods [7,33,47] and the latent animation [25, 44]. This task is much easier than our task since it contains the motion in the same domain. Our face render is also inspired by the method of unsupervised landmark-based method [42] and 3DMM-based method [33] to map the learned coefficient to generate the real video. However, they are not focused on generating realistic motion coefficients.

## 3. Method

As shown in Fig. 2, our system uses the 3D motion coefficients as the intermediate representation for talking head generation. We first extract the coefficients from the original image. Then, the realistic 3DMM motion coefficients are generated by ExpNet and PoseVAE individually. Finally, a 3D-aware face render is proposed to produce the talking head videos. Below, we give a brief introduction to the 3D face model as preliminaries in Sec. 3.1, the audio-driven motion coefficients generation and the coefficients-driven image animator we design in Sec. 3.2 and Sec. 3.3, respectively.

### 3.1. Preliminary of 3D Face Model

3D information is crucial to improve the realness of the generated video since the real video is captured in the 3D environment. However, previous works [30, 51, 52] have rarely been a consideration in 3D space since it is hard to obtain accurate 3D coefficients from a single image and the high-quality face render is also hard to design. Inspired by the recent single image deep 3D reconstruction method [5], we consider the space of the predicted 3D Morphable Models (3DMMs) as our intermediate representation. In 3DMM, the 3D face shape $\mathbf{S}$ can be decoupled as:

$$\mathbf{S} = \overline{\mathbf{S}} + \alpha \mathbf{U}_{id} + \beta \mathbf{U}_{exp}, \qquad (1)$$

where $\overline{\mathbf{S}}$ is the average shape of the 3D face, $\mathbf{U}_{id}$ and $\mathbf{U}_{exp}$ are the orthonormal basis of identity and expression of LSFM morphable model [1]. Coefficients $\alpha \in \mathbb{R}^{80}$ and $\beta \in \mathbb{R}^{64}$ describe the person identity and expression, respectively. To preserve pose variance, coefficients $\mathbf{r} \in SO(3)$ and $\mathbf{t} \in \mathbb{R}^3$ denote the head rotation and translation. To achieve identity irrelevant coefficients generation [33], we only model the parameters of motion as $\{\beta, \mathbf{r}, \mathbf{t}\}$. We learn the head pose $\rho = [\mathbf{r}, \mathbf{t}]$ and expression coefficients $\beta$ individually from the driving audio as introduced before. Then, these motion coefficients are used to implicitly modulate our face render for final video synthesis.

### 3.2. Motion Coefficients Generation through Audio

As introduced above, the 3D motion coefficients contain both head pose and expression where the head pose is a global motion and the expression is relatively local. To this end, learning everything altogether will cause huge uncertainty in the network since the head pose has a relatively weak relationship with audio while the lip motion is highly connected. We generate the motion of the head pose and expression using the proposed PoseVAE and ExpNet, respectively introduced below.

**ExpNet**   Learning a generic model which produces accurate expression coefficients from audio is extremely hard for two reasons: 1) audio-to-expression is not a one-to-one mapping task for different identities. 2) there are some audio-irrelevant motions in the expression coefficients and it will influence the prediction's accuracy. Our ExpNet is designed to reduce these uncertainties. As for the identity issue, we connect the expression motion to the specific person via the first frame's expression coefficients $\beta_0$. To reduce the motion weight of other facial components in natural talking, we use the *lip motion only* coefficients as the coefficient target through the pre-trained network of Wav2Lip [30] and deep 3D reconstruction [5]. Then, other minor facial motions (*e.g.*, eye blink) can be leveraged via the additional landmark loss on the rendered images.
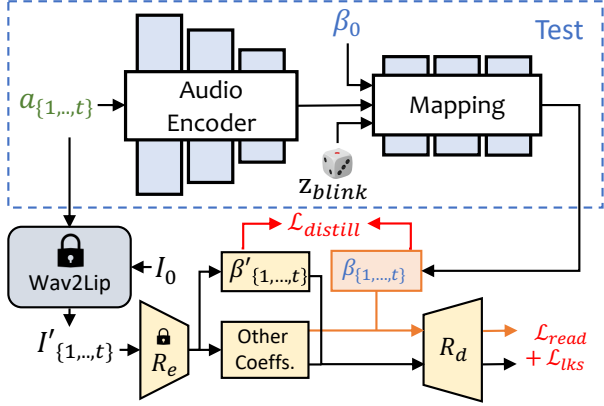
3

Figure 3. The structure of our ExpNet. We involve a monocular 3D face reconstruction model [5] ($R_e$ and $R_d$) to learn the realistic expression coefficients. Where $R_e$ is a pretrained 3DMM coefficients estimator and $R_d$ is a differentiable 3D face render without learnable parameters. We use the reference expression $\beta_0$ to reduce the uncertainty of identity and the generated frame from pre-trained Wav2Lip [30] and the first frame as target expression coefficients since it only contains the lip-related motions.

As shown in Figure 3, we generate the $t$-frame expression coefficients from an audio window $a_{\{1,..,t\}}$, where the audio feature of each frame is a 0.2s mel-spectrogram. For training, we first design a ResNet-based audio encoder $\Phi_A$ [12, 30] to embed the audio feature to a latent space. Then, a linear layer is added as the mapping network $\Phi_M$ to decode the expression coefficients. Here, we also add the reference expression $\beta_0$ from the reference image to reduce the identity uncertainty as discussed above. Since we use the lip-only coefficients as ground truth in the training, we explicitly add a blinking control signal $z_{blink} \in [0, 1]$ and the corresponding eye landmark loss to generate the controllable eye blinks. Formally, the network can be written as:

$$\beta_{\{1,...,t\}} = \Phi_M(\Phi_A(a_{\{1,...,t\}}), z_{blink}, \beta_0) \qquad (2)$$

As for the loss function, we first use $\mathcal{L}_{distill}$ to evaluate the differences between the lip only expression coefficients $R_e(\texttt{Wav2Lip}(I_0, a_{\{1,...,t\}}))$ and the generated $\beta_{\{1,...,t\}}$. Notice that, we only use the first frame $I_0$ of the wav2lip to generate the lip-sync video which reduces the influence of the pose variant and other facial expressions apart from lip movement. Besides, we also involve the differentiable 3D face render $R_d$ to calculate the additional perceptual losses in explicit facial motions space. As shown in Figure 3, we calculate the landmark loss $\mathcal{L}_{lks}$ to measure the range of eye blink and the overall expression accuracy. A pretrained lip reading network $\Phi_{reader}$ is also used as temporal lip reading loss $\mathcal{L}_{read}$ to keep the perceptual lip qualities [9, 30]. We provide more training details in the supplementary materials.
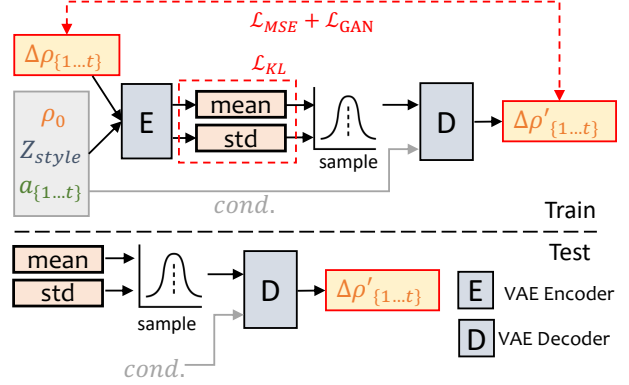


Figure 4. The pipeline of the proposed PoseVAE. We learn the residual of the input head pose $\rho_0$ via a conditional VAE structure. Given the conditions: first frame $\rho_0$, style identity $Z_{style}$ and the audio clip $a_{\{1,...,t\}}$, our method learns a distribution of the residual head pose $\Delta\rho_{\{1,...,t\}} = \rho_{\{1,...,t\}} - \rho_0$. After training, we can generate the stylized results through the pose decoder and the conditions ($cond.$) only.

**PoseVAE** As shown in Figure 4, a VAE [21] based model is designed to learn the realistic and identity-aware stylized head movement $\rho \in \mathbb{R}^6$ of the real talking video. In training, the pose VAE is trained on fixed $n$ frames using an encoder-decoder-based structure. Both the encoder and decoder are two-layer MLPs, where the inputs contain a sequential $t$-frame head poses and we embed it to a Gaussian distribution. In the decoder, the network is learned to generate the $t$-frame poses from the sampled distribution. Instead of generating the pose directly, our PoseVAE learns the *residual* of the condition pose $\rho_0$ of the first frame, which enables our method to generate longer, stable, and continuous head motion in testing under the condition of the first frame. Besides, according to CVAE [6], we add the corresponding audio feature $a_{\{1,...,t\}}$ and style identity $Z_{style}$ as conditions for rhythm awareness and identity style. The KL-divergence $\mathcal{L}_{KL}$ is used to measure the distribution of the generated motions. The mean square loss $\mathcal{L}_{MSE}$ and adversarial loss $\mathcal{L}_{GAN}$ are used to ensure the generated quality. We provide more details about the loss function in the supplementary materials.

### 3.3. 3D-aware Face Render

After generating the realistic 3D motion coefficients, we render the final video through a well-designed 3D-aware image animator. We get inspiration from the recent image animation method face-vid2vid [42] because it implicitly learns the 3D information from a single image. However, a real video is required as the motions driving signal in their method. Our face render makes it drivable through 3DMM coefficients. As shown in Figure 5, we propose mappingNet to learn the relationship between the explicit 3DMM motion coefficients (head pose and expression) and the implicit unsu-
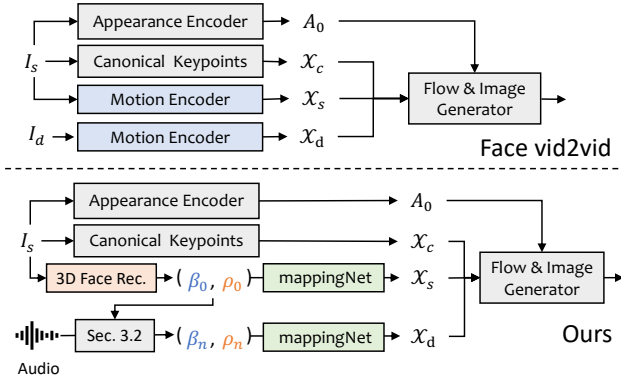
4

Figure 5. The proposed FaceRender and comparison with face-vid2vid [42]. Given source image $I_s$ and driving image $I_d$, face-vid2vid generates the motions in a unsupervised 3D keypoint spaces of $\mathcal{X}_c$, $\mathcal{X}_s$ and $\mathcal{X}_d$. Then, the image can be generated via the appearance $A_0$ and the keypoints. Since we do not have driving image, we use the explicit disentangled 3DMM coefficients as proxy and map it to the unsupervised 3D keypoints space.

pervised 3D keypoints. Our mappingNet is built via several 1D convolutional layers. We use the temporal coefficients from a time window for smoothing as PIRenderer [33]. Differently, we find the face alignment motion coefficients in PIRenderer will hugely influence the motion naturalness of audio-driven video generation and provide an experiment in Sec. 4.4. We only use the coefficients of expression and head pose.

As for training, our method contains two steps. Firstly, we train face-vid2vid [42] in a self-supervised fashion as in the original paper. In the second step, we freeze all the parameters of the appearance encoder, canonical keypoints estimator, and image generator for tuning. Then, we train the mapping net on the 3DMM coefficients of the ground truth video in a reconstruction style. We give the supervision in the domain of unsupervised keypoints using $\mathcal{L}_1$ loss and the final generated video following their original implementation. More details can be founded in the supplementary materials.

## 4. Experiments

### 4.1. Implementation Details and Metrics

**Datasets**   We use VoxCeleb [26] dataset for training which contains over 100k videos of 1251 subjects. We crop the original videos following previous image animation methods [36] and resize the video to 256×256. After preprocessing, the data is used to train our FaceRender. Since some videos and audios are not aligned in VoxCeleb, we select 1890 aligned videos and audios of 46 subjects to train our PoseVAE and ExpNet. The input audios are down-sampled to 16kHz and transformed to mel-spectrograms with the same setting as Wav2lip [30]. To test our method, we use the

346 videos' first 8-second video (around 70k frames in total) from HDTF dataset [49] since it contains high resolution and in-the-wild talking head videos. These videos are also cropped and processed following [36] and resized to 256 ×256 for evaluation. We use the first frame of each video as the reference image to generate videos.

**Implementation Details**   All of ExpNet, PoseVAE, and FaceRender are trained separately and we employ Adam optimizer [20] for all experiments. After training, our method can be inferred in an end-to-end fashion without any manual intervention. All the 3DMM parameters are extracted through pre-trained deep 3D face reconstruction method [5]. We perform all the experiments on 8 A100 GPUs. ExpNet, PoseVAE, and FaceRender are trained with a learning rate of $2e^{-5}$, $1e^{-4}$, and $2e^{-4}$, respectively. As for the temporal consideration, ExpNet uses continuous 5 frames to learn. PoseVAE is learned via continuous 32 frames. The frames in FaceRender are generated frame-by-frame with the coefficients of 5 continuous frames for stability.

**Evaluation Metrics**   We demonstrate the superiority of our method on multiple metrics that have been widely used in previous studies. We employ Frechet Inception Distance (FID) [13, 35] and cumulative probability blur detection (CPBD) [27] to evaluate the quality of the images, in which FID is for the realism of generated frames and CPBD is for the sharpness of generated frames. To evaluate identity preservation, we calculate the cosine similarity (CSIM) of identity embedding between the source images and the generated frames, in which we use ArcFace [4] to extract identity embedding of images. To evaluate lip synchronization and mouth shape, we evaluate the perceptual differences of the mouth shape from Wav2Lip [30], including the distance score (LSE-D) and confidence score (LSE-C). We also conduct some metrics to evaluate the head motions of generated frames. For the diversity of the generated head motions, a standard deviation of the head motion feature embeddings extracted from the generated frames using Hopenet [28] is calculated. For the alignment of the audio and generated head motions, we compute Beat Align Score as in Bailando [22].

### 4.2. Compare with other state-of-the-art methods

We compare several state-of-the-art methods for the talking head video generations (MakeItTalk [52], Audio2Head [39] and Wang *et al*. [40] [1]) and audio to expression generations (Wav2Lip [30], PC-AVS [51]). The evaluation is performed on their publicly available checkpoint directly. As shown in Table 1, the proposed method shows better overall video qualities and head pose diversity

---

[1]This method needs to extract the phoneme from audio, which only works on the specific language.

| Method | Eye Blink | Lip Synchronization | | Learned Head Motion | | Video Quality | | |
|---|---|---|---|---|---|---|---|---|
| | | LSE-C↑ | LSE-D↓ | Diversity↑ | Beat Align↑ | FID↓ | CPBD↑ | CSIM↑ |
| Real Video | N./A. | 8.211 | 6.982 | 0.259 | 0.271 | 0.000 | 0.428 | 1.000 |
| Wav2Lip* [30] | N./A. | 10.221 | 5.535 | N./A. | N./A. | 21.725 | 0.368 | 0.849 |
| PC-AVS** [51] | from ref. | 9.053 | 6.355 | N./A. | N./A. | 69.127 | 0.206 | 0.683 |
| MakeItTalk [52] | automatic | 5.110 | 10.059 | 0.257 | 0.268 | 28.243 | 0.283 | 0.838 |
| Audio2Head [39] | automatic | 7.357 | 7.535 | 0.181 | 0.267 | 24.392 | 0.281 | 0.823 |
| Wang *et al.* [40] | automatic | 4.932 | 10.055 | 0.226 | 0.268 | 22.432 | 0.295 | 0.811 |
| Ours | controllable | 7.290 | 7.772 | **0.278** | **0.293** | **22.057** | **0.335** | **0.843** |

Table 1. Comparison with the state-of-the-art method on HDTF dataset. We evaluate Wav2Lip [30] and PC-AVS [51] in the one-shot settings. Wav2Lip* achieves the best video quality since it only animates the lip region while other regions are the same as the original frame. PC-AVS** is evaluated using the fixed reference pose and fails in some samples.



Figure 6. We compare our method with several state-of-the-art methods for single image audio-driven talking head generation. Our method produces much higher quality results in terms of lip synchronization, identity preservation, head motion and image quality. *We give the target image above for both lip shape and identity reference.* Please refer our supplementary video for better comparison.

and also shows comparable performance with other fully talking-head generation methods in terms of the lip synchronization metrics. We argue that these lip synchronization metrics are too sensitive to the audio where the unnatural lip movement may get a better score. However, our method

achieves a similar score to the real videos, which demonstrates our advantages. We also illustrate the visual results of different methods in Figure 6. Here, we give the lip reference to visualize the lip synchronization of our method. From the figure, our method has a very similar visual quality to the
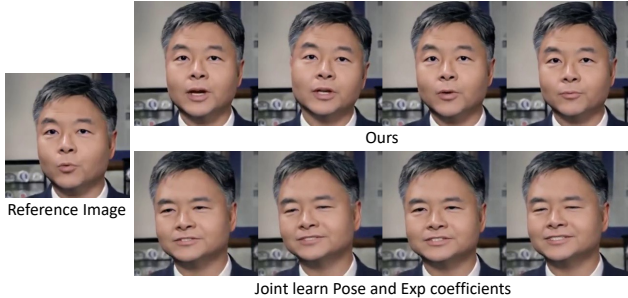
Figure 7. We compare our method with a baseline method which learn all the coefficients from a single network without any condition (from Speech2Gesture [10]). Our method shows clear head movements, identity preservation and diverse expressions.

original target video and with different head poses as we expected. Compared with other methods, Wav2Lip [30] produces blur half-face. PC-AVS [51] and Audio2Head [39] are struggling for identity preservation. Audio2Head can only generate the front talking face. Besides, MakeItTalk [52] and Audio2Head [39] generate the distorted face video due to 2D warping. We give the video comparison in the supp. to show more clear comparisons.

### 4.3. User Studies

We conduct user studies to evaluate the performance of all the methods. We generate overall 20 videos as our testing. These samples contain almost equal genders with different ages, poses and expressions to show the robustness of our method. We invert 20 participants and let them choose the best method in terms of video sharpness, lip synchronization, the diversity and naturalness of the head motion, and overall quality. The results are shown in Table 2, where the participants like our method mostly because of the video and motion quality. We also find that 38% of the participants think our methods show better lip synchronization than other methods, which is inconsistent with Table 1. We think it might be because most of the participants focus on the overall quality of the video, where the blurry and still face videos [30, 51] influence their opinions.

| Method | Lip Sync. | Motion Diversity | Video Sharpness | Overall Naturalness |
|---|---|---|---|---|
| Wav2Lip [30] | 15.6% | 3.1% | 2.0% | 2.8% |
| PC-AVS [51] | 18.1% | 9.6% | 3.4% | 9.1% |
| MakeItTalk [52] | 5.6% | 5.3% | 5.7% | 6.9% |
| Wang *et al.* [40] | 12.5% | 12.1% | 16.3% | 11.6% |
| Audio2Head [39] | 9.5% | 12.1% | 9.7% | 14.7% |
| Ours | **38.7%** | **57.9%** | **62.8%** | **54.8%** |

Table 2. User study.

### 4.4. Ablation Studies

**Ablation of ExpNet** For ExpNet, we mainly evaluate the necessity of each component via the lip synchronization

| Method | LSE-C $\uparrow$ | LSE-D $\downarrow$ |
|---|---|---|
| Speech2Gesture [10] | 0.878 | 13.889 |
| OursFull (Lip coeffs. + $\beta_0$ + $\mathcal{L}_{read}$) | **7.290** | **7.772** |
| w/o $\beta_0$ & $\mathcal{L}_{read}$ | 5.241 | 9.532 |
| w/o $\mathcal{L}_{read}$ | 6.993 | 7.841 |
| w/ real coeffs. | 6.567 | 8.061 |

Table 3. Ablation for ExpNet. Both the initial expression $\beta_0$, lip reading loss $\mathcal{L}_{read}$ improve the performance a lot. However, the lip synchronization metric drops a lot when using the real coefficients.



Figure 8. The ablation of ExpNet. We choose four frames from the generated video as comparison. Our method largely reduces the uncertainty of audio to expression generation. The reference $\beta_0$ is used to provide the identity information while the lip only coefficients generate better lip synchronization. *Notice that, the target image is provided as the identity and lip motion reference.*

metrics. Since there are no disentangled methods before, we consider a baseline (Speech2Gesture [10], which is an audio to keypoint generation network) to learn the head pose and expression coefficients jointly. As shown in Table 3 and Figure 7, learning all the motion coefficients altogether is hard to generate truth-worthy talking head videos. We then consider the variants of the proposed ExpNet, both the initial expression $\beta_0$, lip reading loss $\mathcal{L}_{read}$ and the necessity of lip-only coefficients are critical. The visual comparison is shown in Figure 8, where our method w/o the initial expression $\beta_0$ shows huge identity changes as expected. Also, if we use the real coefficients to replace the lip-only coefficients we use, the performance drops a lot in lip synchronization.

**Ablation of PoseVAE** We evaluate the proposed PoseVAE in terms of motion diversities and audio beat alignments. As shown in Table. 4, the baseline Speech2Gesture [10] also performs worse in pose evaluation. As for our variants, since

| Method | Diversity↑ | Beat Align↑ |
|---|---|---|
| Speech2Gesture [10] | 0.1574 | 0.274 |
| OurFull (Single Fixed Style) | 0.2735 | 0.287 |
| w/o $\mathcal{L}_{gan}$ | 0.2500 | 0.271 |
| w/o initial pose | 0.2725 | 0.278 |
| w/o audio | 0.2566 | 0.274 |
| w/o all conditions | 0.2631 | 0.279 |
| OursFull (Mixed Style) | **0.2778** | **0.293** |

Table 4. Ablation the diversity and audio alignment of the proposed PoseVAE. Each component or conditional contribute largely to generate realistic head motions.
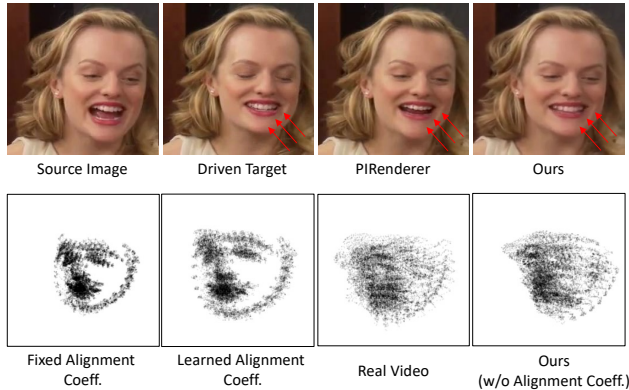


Figure 9. Ablation studies of face render. In the first row, we directly compare our method with PIRenderer [33] for face animation and our method shows better expression modeling. The second row is the trace map of the generated facial landmarks from the same motion coefficients. Using additional face alignment coefficients as part of the motion coefficients [33] will generate unrealistic aligned head video.

our method contains several identity style labels, to better evaluate other components, we first consider the perform the ablation studies on a fixed one-hot style of our full method for evaluation (OurFull, Single Fixed Style). Each condition in our settings benefits the overall motion quality in terms of diversity and beat alignment. We further report the results of the mixed style of our full method, which uses the randomly-selected identity label as style and shows a better diversity performance also. Since the pose differences are hard to be shown in the figure, please refer to our supplementary materials for better comparison.

**Ablation of Face Render**　We conduct the ablation study on the proposed face render from two aspects. On the one hand, we show the reconstruction quality of our method with the PIRenderer [33], since both methods use 3DMM as an intermediate representation. As shown in the first row of Fig. 9, the proposed face render shows better expression reconstruction qualities thanks to the mapping of sparse unsupervised keypoints. Where accurate expression mapping is also the key to achieving lip synchronization. Besides, we
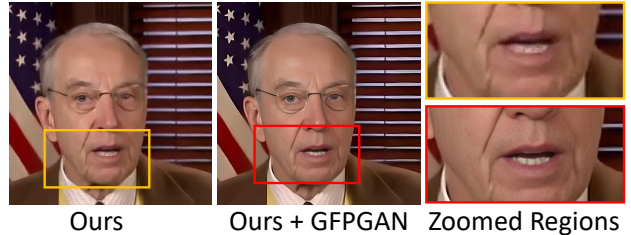


Figure 10. Limitation. Our method may show some teeth artifacts in the lip region in some examples, it can be improved via the face restoration network, *i.e.*, GFPGAN [43].

evaluate the pose unnaturalness caused by the additional alignment coefficients used in PIRenderer [33]. As shown in the second row of Fig. 9, we plot the trace map of the landmarks from the generated video with the same head pose and expression coefficients. Using the fixed or learning-able crop coefficients (as part of pose coefficients in our poseVAE) will generate the face-aligned video, which is strange as a natural video. We remove it and directly use the head pose and expression as modulation parameters showing a more realistic result.

## 4.5. Limitation

Although our method generates realistic video from a single image and audio, there still have some limitations in our system. Since 3DMMs do not model the variation of eyes and teeth, the mappingNet in our Face Render will also struggle to synthesize the realistic teeth in some cases. This limitation can be improved via the blind face restoration networks [43] as shown in Fig. 10. Another limitation of our work is that we only concern the lip motion and eye blinking other than the other facial expressions, *e.g.*, emotion and gaze direction. Thus, the generated video has a fixed emotion, which also reduces the realism of generated content. We consider it as future work.

## 5. Conclusion

In this paper, we present a new system for stylized audio-driven talking head video generation. We use the motion coefficients from 3DMM as an intermediate representation and learn the relationships. To generate realistic 3D coefficients from audio, we propose ExpNet and PoseVAE for realistic expressions and diverse head poses. To model the relationships between 3DMM motion coefficients and the real video, we propose a novel 3D-aware face render inspired by the image animation method [42]. The experiments demonstrate the superiority of our entire framework. Since we predict the realistic 3D facial coefficients, our method can also be used in other modalities directly, *i.e.*, personalized 2D visual dubbing [45], 2D Cartoon animation [52], 3D face animation [8, 46] and NeRF-based 4D talking-head generation [15].

**Ethical Considerations** We consider the misuse of the proposed method since it can generate very realistic video from a single face image. Both visible and invisible video watermarks will be inserted into the produced video for generated content identification similar to Dall-E [32] and Imagen [34]. We also hope our method can provide new research samples in the area of forgery detection.

# References

[1] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces. In *ACM SIGGRAPH*, 1999. 2, 3

[2] Lele Chen, Ross K Maddox, Zhiyao Duan, and Chenliang Xu. Hierarchical cross-modal talking face generation with dynamic pixel-wise loss. In *CVPR*, 2019. 2

[3] Kun Cheng, Xiaodong Cun, Yong Zhang, Menghan Xia, Fei Yin, Mingrui Zhu, Xuan Wang, Jue Wang, and Nannan Wang. Videoretalking: Audio-based lip synchronization for talking head video editing in the wild. In *SIGGRAPH Asia 2022 Conference Papers*, pages 1–9, 2022. 2

[4] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *CVPR*, 2019. 5

[5] Yu Deng, Jiaolong Yang, Sicheng Xu, Dong Chen, Yunde Jia, and Xin Tong. Accurate 3d face reconstruction with weakly-supervised learning: From single image to image set. In *CVPR Workshops*, 2019. 2, 3, 4, 5, 12, 13, 14

[6] Carl Doersch. Tutorial on variational autoencoders. *arXiv preprint arXiv:1606.05908*, 2016. 2, 4

[7] Michail Christos Doukas, Stefanos Zafeiriou, and Viktoriia Sharmanska. Headgan: One-shot neural head synthesis and editing. In *ICCV*, 2021. 3

[8] Yingruo Fan, Zhaojiang Lin, Jun Saito, Wenping Wang, and Taku Komura. Faceformer: Speech-driven 3d facial animation with transformers. In *CVPR*, 2022. 8

[9] Panagiotis P. Filntisis, George Retsinas, Foivos Paraperas-Papantoniou, Athanasios Katsamanis, Anastasios Roussos, and Petros Maragos. Visual speech-aware perceptual 3d facial expression reconstruction from videos. *arXiv preprint arXiv:2207.11094*, 2022. 4, 13

[10] Shiry Ginosar, Amir Bar, Gefen Kohavi, Caroline Chan, Andrew Owens, and Jitendra Malik. Learning individual styles of conversational gesture. In *CVPR*, 2019. 7, 8, 14

[11] Yudong Guo, Keyu Chen, Sen Liang, Yong-Jin Liu, Hujun Bao, and Juyong Zhang. Ad-nerf: Audio driven neural radiance fields for talking head synthesis. In *ICCV*, 2021. 2

[12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 4

[13] 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. In *NeurIPS*, 2017. 5

[14] Fa-Ting Hong, Longhao Zhang, Li Shen, and Dan Xu. Depth-aware generative adversarial network for talking head video generation. In *CVPR*, 2022. 3

[15] Yang Hong, Bo Peng, Haiyao Xiao, Ligang Liu, and Juyong Zhang. Headnerf: A real-time nerf-based parametric head model. In *CVPR*, 2022. 8

[16] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. *CVPR*, 2017. 14

[17] Xinya Ji, Hang Zhou, Kaisiyuan Wang, Qianyi Wu, Wayne Wu, Feng Xu, and Xun Cao. Eamm: One-shot emotional talking face via audio-based emotion-aware motion model. In *ACM SIGGRAPH*, 2022. 2

[18] Xinya Ji, Hang Zhou, Kaisiyuan Wang, Wayne Wu, Chen Change Loy, Xun Cao, and Feng Xu. Audio-driven emotional video portraits. In *CVPR*, 2021. 3

[19] Hyeongwoo Kim, Pablo Garrido, Ayush Tewari, Weipeng Xu, Justus Thies, Matthias Niessner, Patrick Pérez, Christian Richardt, Michael Zollhöfer, and Christian Theobalt. Deep video portraits. *ACM Transactions on Graphics (TOG)*, 2018. 2

[20] Diederik P Kingma and Jimmy Ba. A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5

[21] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *CoRR*, abs/1312.6114, 2014. 4

[22] Siyao Li, Yu Weijiang, Gu Tianpei, Lin Chunze, Wang Quan, Qian Chen, Loy Chen Change, and Liu Ziwei. Bailando: 3d dance generation via actor-critic gpt with choreographic memory. In *CVPR*, 2022. 5

[23] Yuanxun Lu, Jinxiang Chai, and Xun Cao. Live speech portraits: real-time photorealistic talking-head animation. *ACM Transactions on Graphics (TOG)*, 2021. 3

[24] Pingchuan Ma, Yujiang Wang, Stavros Petridis, Jie Shen, and Maja Pantic. Training strategies for improved lip-reading. In *ICASSP*, 2022. 13

[25] Arun Mallya, Ting-Chun Wang, and Ming-Yu Liu. Implicit Warping for Animation with Image Sets. In *NeurIPS*, 2022. 3

[26] Arsha Nagrani, Joon Son Chung, and Andrew Zisserman. Voxceleb: a large-scale speaker identification dataset. In *INTERSPEECH*, 2017. 5, 11, 12

[27] Niranjan D. Narvekar and Lina J. Karam. A no-reference image blur metric based on the cumulative probability of blur detection (cpbd). *TIP*, 2011. 5

[28] Ruiz Nataniel, Eunji Chong, and Rehg James M. Fine-grained head pose estimation without keypoints. In *CVPR Workshops*, 2018. 5

[29] 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. *arXiv preprint arXiv:2301.06281*, 2023. 3

[30] K R Prajwal, Rudrabha Mukhopadhyay, Vinay P.Namboodiri, and C.V.Jawahar. A lip sync expert is all you need for speech to lip generation in the wild. In *ACM MM*, 2020. 2, 3, 4, 5, 6, 7, 11, 12

[31] Prajwal K R, Rudrabha Mukhopadhyay, Jerin Philip, Abhishek Jha, Vinay Namboodiri, and C V Jawahar. Towards automatic face-to-face translation. In *ACM MM*, 2019. 2

[32] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 9

[33] Yurui Ren, Ge Li, Yuanqi Chen, Thomas H Li, and Shan Liu. Pirenderer: Controllable portrait image generation via semantic neural rendering. In *ICCV*, 2021. 2, 3, 5, 8, 11, 14

[34] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022. 9

[35] Maximilian Seitzer. pytorch-fid: FID Score for PyTorch. https://github.com/mseitzer/pytorch-fid, August 2020. Version 0.2.1. 5

[36] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. First order motion model for image animation. In *NeurIPS*, 2019. 2, 3, 5, 11

[37] Aliaksandr Siarohin, Oliver Woodford, Jian Ren, Menglei Chai, and Sergey Tulyakov. Motion representations for articulated animation. In *CVPR*, 2021. 3

[38] Justus Thies, Mohamed Elgharib, Ayush Tewari, Christian Theobalt, and Matthias Nießner. Neural voice puppetry: Audio-driven facial reenactment. In *ECCV*, 2020. 2

[39] Suzhen Wang, Lincheng Li, Yu Ding, Changjie Fan, and Xin Yu. Audio2head: Audio-driven one-shot talking-head generation with natural head motion. In *IJCAI*, 2021. 2, 5, 6, 7, 12

[40] Suzhen Wang, Lincheng Li, Yu Ding, and Xin Yu. One-shot talking face generation from single-speaker audio-visual correlation learning. In *AAAI*, 2022. 2, 5, 6, 7, 12

[41] Ting-Chun Wang, Ming-Yu Liu, Andrew Tao, Guilin Liu, Jan Kautz, and Bryan Catanzaro. Few-shot video-to-video synthesis. In *NeurIPS*, 2019. 3

[42] Ting-Chun Wang, Arun Mallya, and Ming-Yu Liu. One-shot free-view neural talking-head synthesis for video conferencing. In *CVPR*, 2021. 2, 3, 4, 5, 8, 11, 14

[43] Xintao Wang, Yu Li, Honglun Zhang, and Ying Shan. Towards real-world blind face restoration with generative facial prior. In *CVPR*, 2021. 8

[44] Yaohui Wang, Di Yang, Francois Bremond, and Antitza Dantcheva. Latent image animator: Learning to animate images via latent space navigation. *arXiv preprint arXiv:2203.09043*, 2022. 3

[45] Xin Wen, Miao Wang, Christian Richardt, Ze-Yin Chen, and Shi-Min Hu. Photorealistic audio-driven video portraits. *IEEE Transactions on Visualization and Computer Graphics*, 26(12):3457–3466, 2020. 2, 8

[46] Jinbo Xing, Menghan Xia, Yuechen Zhang, Xiaodong Cun, Jue Wang, and Tien-Tsin Wong. Codetalker: Speech-driven 3d facial animation with discrete motion prior. *arXiv preprint arXiv:2301.02379*, 2023. 8

[47] 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 *ECCV*, 2022. 3

[48] Chenxu Zhang, Yifan Zhao, Yifei Huang, Ming Zeng, Saifeng Ni, Madhukar Budagavi, and Xiaohu Guo. Facial: Synthesizing dynamic talking face with implicit attribute learning. In *ICCV*, 2021. 3

[49] Zhimeng Zhang, Lincheng Li, Yu Ding, and Changjie Fan. Flow-guided one-shot talking face generation with a high-resolution audio-visual dataset. In *CVPR*, 2021. 2, 5, 11, 12

[50] Jian Zhao and Hui Zhang. Thin-plate spline motion model for image animation. In *CVPR*, 2022. 3

[51] Hang Zhou, Yasheng Sun, Wayne Wu, Chen Change Loy, Xiaogang Wang, and Ziwei Liu. Pose-controllable talking face generation by implicitly modularized audio-visual representation. In *CVPR*, 2021. 2, 3, 5, 6, 7, 12

[52] Yang Zhou, Xintong Han, Eli Shechtman, Jose Echevarria, Evangelos Kalogerakis, and Dingzeyu Li. Makelttalk: speaker-aware talking-head animation. *ACM Transactions on Graphics (TOG)*, 2020. 2, 3, 5, 6, 7, 8, 12

# A. Additional Experiments

## A.1. PIRenderer *v.s* Our FaceRender for Face Reenactment

We compare our FaceRender and PIRenderer [33] on the task of video-driven face reenactment. We have already shown the visual comparisons in Fig 9 of the main paper, here, we give the numerical comparison results on the HDTF dataset. We evaluate these two methods using cross-identity settings and the results are conducted over 354 videos. As shown in Tab A5, the proposed method shows a much better visual quality in terms of FID and CSIM, which demonstrate the advantage of the proposed methods for audio-driven talking-head generation. For more differences between the proposed method and PIRenderer for this task, we also discuss the influence of the face alignment coefficient in Sec. B.4.

| Method | FID $\downarrow$ | CPBD $\uparrow$ | CSIM $\uparrow$ |
|---|---|---|---|
| PIRenderer [33] | 26.521 | **0.363** | 0.857 |
| Our face render | **19.646** | 0.334 | **0.880** |

Table A5. Face render evaluation.

## A.2. Cross-ID Settings and More Test Datasets

We conduct the *same-identity* experiment in the main paper, where the first frame of the test video is regarded as the reference image and the corresponding audio is regarded as the driving signal, generating a video that has the synchronized expression but diverse head poses. Differently, in the cross-identity experiment, the driving audio comes from another video. This kind of setting is also widely used in the comparison of the video-driven face reenactment [36]. In the cross-identity experiment, the reference pose of PC-AVS comes from the audio's corresponding video.

To this end, besides the HDTF [49] dataset in the main paper (we evaluate the results under the same-identity experiment), we also evaluate our method on HDTF and Vox-Celeb2 [26] datasets in ***cross-identity setting*** as in Table A6 and Table A7. VoxCeleb2 contains over 1 million utterances of 6112 speakers, of which there are 36k utterances of 118 speakers in the test set. We randomly select 3 videos for each speaker, obtaining 354 videos as for evaluation. The evaluation metrics are the same as those the same-identity experiment on the HDTF dataset. We directly evaluate the pretrained model of all the models on this dataset. We crop the videos in the same way used in [36] and resize the frames to 256×256. As shown in the Tables, the proposed method shows better lip synchronization in this kind of setting on both datasets in most metrics. The same trend is also observed in the head motion and visual quality of the final videos.

# B. More Implementation Details

We provide the detailed audio pre-processing, network structures, loss function, the discussion on alignment coefficients in Sec. B.1, Sec. B.2, Sec. B.3 and Sec. B.4.

## B.1. Audio Pre-processing Details

We follow Wav2Lip [30] to pre-process the audio. Specifically, we pre-process all the audio to 16k Hz. Then, we convert it to the mel-spectrograms with FFT window size 800, hop length 200 and 80 Mel filter banks. Thus, for each frame, we have 0.2s mel-spectrogram feature with the shape of 16×80.

## B.2. Network Structure Details

**ExpNet** Our ExpNet is built via an audio encoder $\Phi_A$ and a linear layer $\Phi_M$. We use the parameters from the pre-trained Wav2Lip to initialize the audio encoder. As discussed in the main paper, we first encode the audio feature into an audio embedding. And then, we generate the expression coefficients $\beta^g_{\{1,...,t\}}$ with additional coefficients of the first frame $\beta_0$ and blink control signal $z_{blink}$. As shown in Fig. B11 (c), the audio encoder $\Phi_A$ is built via a four stages ResBlock-C as in Fig. B11 (a). we only use a single linear layer $\Phi_M$ as in Fig. B11 (b).

**PoseVAE** As shown in Fig. B11 (e), both the encoder and decoder of our PoseVAE contain several linear layers. For the conditions, the encoder $\mu$ and $\sum$ is mapped through the concatenation of the $\Delta\rho_{\{1:T\}}$, the 46 dimensional one-hot vector $Z_{style}$ (our training dataset contains 46 identities) and the encoded features from audio encoder $\Phi_A$. As for the decoder, we first add the re-parameterized feature and the style embedding. Then, we concatenate the audio feature similar to the encoder.

**FaceRender** As discussed in the main paper of Fig. 5, our face render is inspired by the motion transfer method face-vid2vid [42]. We introduce a mappingNet to remap the learned 3DMM motion coefficients to the space of unsupervised 3D keypoints. As shown in Figure B11 (d), the mapping network contains the $t$-frames ($[t-2:t+2]$) motion coefficients of pose $\rho_{[t-2:t+2]}$ and expression $\beta_{[t-2:t+2]}$ to generate the motions representation of face-vid2vid [42] ( yaw, pitch, roll, tr, and $\delta$) in frame $t$. Other networks in our FaceRender have the same structures in [42]. Please refer to face-vid2vid [42] for more network details on the FaceRender.

## B.3. Loss Function Details

**ExpNet** As described in the main paper, we use the expression coefficients which are generated from the pre-trained

| Method | Lip Synchronization | | Learned Head Motion | | Video Quality | | |
|---|---|---|---|---|---|---|---|
| | LSE-C↑ | LSE-D↓ | Diversity↑ | Beat Align↑ | FID↓ | CPBD↑ | CSIM↑ |
| Real Video | 8.211 | 6.982 | 0.259 | 0.271 | 0 | 0.428 | 1.000 |
| Wav2Lip* [30] | 9.641 | 6.035 | N./A. | N./A. | 21.727 | 0.368 | 0.846 |
| PC-AVS** [51] | 8.959 | 6.435 | N./A. | N./A. | 99.098 | 0.201 | 0.648 |
| MakeItTalk [52] | 4.937 | 10.231 | 0.2553 | 0.276 | 26.829 | 0.333 | 0.834 |
| Audio2Head [39] | 7.237 | **7.648** | 0.1783 | 0.260 | 24.404 | 0.282 | 0.818 |
| Wang *et al.* [40] | 4.634 | 10.457 | 0.2260 | 0.265 | 22.302 | 0.294 | 0.805 |
| Ours | **7.343** | 7.709 | **0.2759** | **0.284** | **20.886** | 0.334 | **0.846** |

Table A6. Comparison with the state-of-the-art method on HDTF dataset [49] with *cross-identity* setting. Wav2Lip* achieves the best video quality since it only animates the lip region while other regions are the same as the original frame. In cross-identity setting, PC-AVS** is evaluated using the reference pose from the driving video and fails in some samples.

| Method | Lip Synchronization | | Learned Head Motion | | Video Quality | | |
|---|---|---|---|---|---|---|---|
| | LSE-C↑ | LSE-D↓ | Diversity↑ | Beat Align↑ | FID↓ | CPBD↑ | CSIM↑ |
| Real Video | 6.209 | 7.911 | 0.4879 | 0.266 | 0 | 0.099 | 1.000 |
| Wav2Lip* [30] | 7.640 | 7.099 | N./A. | N./A. | 19.293 | 0.107 | 0.936 |
| PC-AVS** [51] | 7.168 | 7.443 | N./A. | N./A. | 111.043 | 0.074 | 0.494 |
| MakeItTalk [52] | 3.756 | 10.222 | **0.5230** | 0.275 | 23.501 | 0.063 | 0.883 |
| Audio2Head [39] | 5.266 | 8.788 | 0.2064 | 0.273 | 54.694 | 0.098 | 0.602 |
| Wang *et al.* [40] | 3.441 | 10.519 | 0.2547 | 0.272 | 42.092 | **0.136** | 0.750 |
| Ours | **5.571** | **8.503** | 0.5211 | **0.277** | **22.738** | 0.081 | **0.893** |

Table A7. Comparison with the state-of-the-art method on VoxCeleb2 [26] dataset under *cross-identity* setting. Wav2Lip* achieves the best video quality since it only animates the lip region while other regions are the same as the original frame. In cross-identity setting, PC-AVS** is evaluated using the reference pose from the driving video and fails in some samples.

wav2lip [30] and then perform 3D face capture [5] as guidance (lip-only expression coefficients for short). Basically, we calculate $\mathcal{L}_{distill}$ through the Mean-Squared loss between lip-only expression coefficients and the generated expression coefficients in training. Formally, a $T$-frames $\mathcal{L}_{distill}$ can be written as:

$$\mathcal{L}_{distill} = \frac{1}{T}\sum_{t=1}^{T}\left(\beta_t^g - \beta_t^{lip}\right)^2 \quad (3)$$

Where $\beta_t^{lip}$ and $\beta_t^g$ are the lip only and the generated expression coefficients, respectively.

We also calculate the loss function on the projected 2D landmarks of the rendered 3D face. In detail, as shown in Fig. B12, the height and width of the eye area in the $t$-th frame are defined as follows:

$$E_t^w = \frac{\left\|P_t^{39} - P_t^{36}\right\|_2 + \left\|P_t^{45} - P_t^{42}\right\|_2}{2} \quad (4)$$

$$E_t^h = \frac{\left\|P_t^{37} + P_t^{38} - P_t^{40} - P_t^{41}\right\|_2}{2} \quad (5)$$

$$+ \frac{\left\|P_t^{43} + P_t^{44} - P_t^{46} - P_t^{47}\right\|_2}{2}. \quad (6)$$

Where $E_t^w$ is the width of the eye area in frame $t$, $E_t^h$ is the width of the eye area in frame $t$, $P_t^i$ is the $i$-th landmark in frame $t$. we define $R_t = \frac{E_t^h}{E_t^w}$ as the predicted and calculate eye loss as follows:

$$\mathcal{L}_{eye} = \sum_{t=1}^{T}\left\|R_t - Z_t^{blink}\right\|_1, \quad (7)$$

where $Z_t^{blink}$ is the eye blinking control signal the of $t$-th frame which is generated uniformly and randomly. To eliminate the effects of $\mathcal{L}_{eye}$ on other facial expression, we also constrain the minimal modification in the other landmarks.

**(a) ResBlock-C**

$3\times3$-s1,1-p1-Conv-C, BN

ReLU

**(b) Linear Layer $\Phi_M$**

$(a_t^{emb}, \beta_0, Z_t^{blink})$

$(512+64+1)$

Linear-64

$(64)$

Expression coefficients

**(c) Audio encoder $\Phi_A$**

$a_t$   $(1, 80, 16)$

$3\times3$-s1,1-p1-Conv-32, BN

ResBlock-32   $\times2$

$3\times3$-s3,1-p1-Conv-64, BN

ResBlock-64   $\times2$

$3\times3$-s3,3-p1-Conv-128, BN

ResBlock-128   $\times2$

$3\times3$-s3,2-p1-Conv-256, BN

ResBlock-256   $\times2$

$3\times3$-s1,1-p0-Conv-512, BN

$1\times1$-s1,1-p0-Conv-512, BN

Reshape $(512,1,1)\rightarrow(512)$

$(512)$

Audio embedding $a_t^{emb}$

**(d) MappingNet**

$(\beta_{t-2:t+2}, \rho_{t-2:t+2})$   $(70, 5)$

3-s1-p0-Conv1D-256

Leaky-ReLu

3-s1-p0-Conv1D-256

Reshape C256$\times$D1$\rightarrow$C256

fc66 | fc66 | fc66 | fc3 | fc60

softmax | softmax | softmax

$yaw_t$   $pitch_t$   $roll_t$   $tr_t$   $\delta_t$

**(e) PoseVAE**

$\Delta\rho_{1:T}$ $(T,6)$   $Z_{style}$ $(46)$   $a_{1:T}$ $(T,1,80,16)$   $Z_{style}$ $(46)$   $(\mu, \Sigma)$

Linear-6 | Linear-64 | Audio encoder $\Phi_A$ | Linear-64 | reparametrization

Linear-6

Reshape $(T,6)\rightarrow(T\times6)$

Linear-192

Linear-128

Linear-64 | Linear-64

$\mu$ $(64)$   $\Sigma$ $(64)$

Encoder

Linear-128

Linear-192

Linear-64

$\Delta\rho'_{1:T}$ $(T,6)$

Decoder

Figure B11. The architectures of the networks in our model. Here, '$3\times3$-s1,1-p1-Conv-32' means a convolutional layer with the kernel size $3\times3$, the stride size (1,1), padding size (1,1) and the output channel is 32.

Thus,

$$\mathcal{L}_{lks} = \lambda_{eye}\mathcal{L}_{eye} + \frac{1}{T}\frac{1}{N}\sum_{t=1}^{T}\sum_{i=1}^{M}\left\|P_t^i - P_t^{i\prime}\right\|_2^2, \quad (8)$$

where $\lambda_{eye}$ is set to 200, $P_t^{i\prime}$ is the landmarks predicted by the lip-only expression coefficients, $\mathcal{M}$ is a set of landmarks other than the eye areas.

Besides, we use pretrained lip-reading models proposed in [24] to calculate lip reading loss $\mathcal{L}_{read}$ inspired by [9]. We use the pretrained video-based lip-reading model where the input is a sequence (5 frames in our case) of the cropped interesting region around the mouth (as shown in Fig B13) and the target is a series of the character sequence. So we employ a differentiable 3D face render [5] in ExpNet to render the images through the generated expression coefficients, then, we crop the mouth area of the rendered images using the bounding box of the mouth landmarks, obtaining the logit of the character sequences $\mathbf{C_P}$. As for the supervision, we
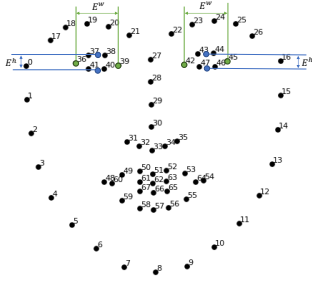
13

Figure B12. Face landmarks visualization.



Figure B13. Example of the cropped interesting region in the 3D rendered face sequences to calculate the lip-reading loss.

generate the logit of the character sequences $\mathbf{C_{gt}}$ from the ground truth audio using the audio-driven lip-reading model. Thus, our goal is to minimize the difference between $\mathbf{C_p}$ and $\mathbf{C_{gt}}$. In other words,

$$\mathcal{L}_{read} = \text{CrossEntropy}(\mathbf{C_{gt}}, \mathbf{C_p}) \qquad (9)$$

Overall, the final loss of ExpNet is given by :

$$\mathcal{L}_{exp} = \lambda_{distill}\mathcal{L}_{distill} + \lambda_{read}\mathcal{L}_{read} + \lambda_{lks}\mathcal{L}_{lks} \qquad (10)$$

Where $\lambda_{distill}$, $\lambda_{read}$, $\lambda_{lks}$ are set to 2, 0.01, and 0.01, respectively.

**PoseVAE**    We first calculate the reconstruction loss by applying Mean-Squared loss between the generated $\Delta\rho'_{\{1...T\}}$ and the original $\Delta\rho_{\{1...T\}}$:

$$\mathcal{L}_{MSE} = \frac{1}{T}\sum_{t=1}^{T}(\Delta\rho'_t - \Delta\rho_t)^2 \qquad (11)$$

Meanwhile, we encourage the similarity of the latent space distribution and the Gaussian distribution with the mean vector $\mu$ and covariance matrix $\sum$. So we define $\mathcal{L}_{KL}$ as the Kullback–Leibler (KL) divergence between the latent space distribution and the Gaussian distribution. We also employ a discriminator $D$ based on the PatchGAN [16] to perform 1D convolution on the head motion sequence as Speech2Gesture [10]. We define the adversarial loss $\mathcal{L}_{GAN}$:

$$\mathcal{L}_{GAN} = \arg\min_{G}\max_{D}(G, D) \qquad (12)$$

Where $G$ is proposed PoseVAE. The total loss of PoseVAE can be summarized as follows.

$$\mathcal{L}_{pose} = \lambda_{MSE}\mathcal{L}_{MSE} + \lambda_{KL}\mathcal{L}_{KL} + \lambda_{GAN}\mathcal{L}_{GAN} \qquad (13)$$

where $\lambda_{MSE}$, $\lambda_{KL}$ and $\mathcal{L}_{GAN}$ are set to 1, 1, and 0.7, respectively.

**FaceRender**    We add a MappingNet to map the explicit 3DMM coefficients to the space of the face-vid2vid [42], to training, apart from the loss functions used in face-vid2vid [42], we add $L_1$ regularization on the domain of unsupervised keypoints:

$$\mathcal{L}_1 = \frac{1}{N}\sum_{n=1}^{N}||K'_n - K_n||_1, \qquad (14)$$

where $K'_n$ and $K_n$ are the $n$-th keypoint generated by our MappingNet and the motion generator of the original face-vid2vid, respectively. The weight of $\mathcal{L}_1$ is set to 20, and the weights of the other loss functions keep the same as in face-vid2vid and they are calculated on the final generated image. Please refer to face-vid2vid [42] for more details.

### B.4. More Details about the Alignment Coefficients.

In the main paper, we show the effect of the alignment coefficients in Fig 9. Here, we give more details about the alignment coefficients. Generally, the alignment coefficients are the transformation parameters (translation and scaling) to transform and crop the arbitrary video to the aligned face video for deep 3D face reconstruction [5]. The implicit modulation of PIRenderer [33] contains 73 dimensional motion coefficients, including the expression (64), head pose (6) and the alignment coefficients (3). Since their method focuses on video driving animation, the alignment coefficients are known in testing. However, it is hard to *learn* from the audio since there is no relationship between the alignment and the audio. We also try to learn and use the alignment coefficients of the first frame in our method (as shown in Fig. 9 and the supp. video), however, the produced head motion is aligned and unnatural. We discard it to obtain more natural video results. Thus, our motion coefficients (70) only contains the expression (64) and head pose (6).

## C. Supplementary Video

We provide a supplementary video to include all the video results of our method and other related methods as comparisons, the ablation study of each component, and more results of our method in different languages.