SYNCDIFF: Diffusion-based Talking Head Synthesis with Bottlenecked Temporal Visual Prior for Improved Synchronization

Xulin Fan^{1*} Heting Gao^{1*} Ziyi Chen² Peng Chang² Mei Han² Mark Hasegawa-Johnson¹ University of Illinois at Urbana-Champaign ²PAII Inc.

{xulinf2, hgao17, jhasegaw}@illinois.edu {chenziyi253, changpeng805, hanmei613}@paii-labs.com

Abstract

Talking head synthesis, also known as speech-to-lip synthesis, reconstructs the facial motions that align with the given audio tracks. The synthesized videos are evaluated on mainly two aspects, lip-speech synchronization and image fidelity. Recent studies demonstrate that GAN-based and diffusion-based models achieve state-of-the-art (SOTA) performance on this task, with diffusion-based models achieving superior image fidelity but experiencing lower synchronization compared to their GAN-based counterparts. To this end, we propose SYNCDIFF, a simple yet effective approach to improve diffusion-based models using a temporal pose frame with information bottleneck and facialinformative audio features extracted from AVHUBERT, as conditioning input into the diffusion process. We evaluate SYNCDIFF on two canonical talking head datasets, LRS2 and LRS3 for direct comparison with other SOTA models. Experiments on LRS2/LRS3 datasets show that SYNCDIFF achieves a synchronization score 27.7%/62.3% relatively higher than previous diffusion-based methods, while preserving their high-fidelity characteristics.

1. Introduction

Talking head synthesis has gained popularity as a research topic due to its expanding applications, which include virtual being creation, online conferences, audio dubbing, and video translation. Its primary objective is to generate lip-synced avatar videos based on given speech audio [5, 6, 19, 23, 41, 48].

Over recent years, many studies have been conducted to synthesize realistic talking head videos. Some approaches divide the entire generation process into two steps, utilizing facial landmarks or facial-model-based parameters as an intermediate feature. The first step is designed to predict precise intermediate features from speech audio, and the second step focuses on generating realistic images given the intermediate features [5, 6, 24, 25, 44, 48]. Recently, more works have focused on designing end-to-end video generation with advanced generative models [10, 12, 13, 28, 32, 39, 46, 47]. In addition to variations in the talking head generation pipeline, studies also vary in model generality. Some approaches are character-specific, necessitating the model to undergo training or adaptation using a substantial amount of data specific to the particular character [10, 13, 25, 44]. On the other hand, some models adopt a more general synthesis approach without character restrictions [12, 28, 32, 39, 46, 47]. These models are trained with large talking head datasets across various characters. This approach grants them the zero-shot ability to synthesize new character videos without additional character-specific training data and retraining the model.

A well-synchronized lip movement and high-fidelity generated image are the two most critical aspects of talking head synthesis. Several Generative Adversarial Network (GAN) based studies featuring a specific lip synchronization training design have demonstrated notable results in voice-lip synchronizations [12, 28, 39]. Nevertheless, these GAN-based approaches are beset by challenges such as unstable training processes and suboptimal visual quality. Conversely, certain diffusion-based models [10, 32] exhibit superior image quality, yet their generated videos may not attain the same degree of lip synchronization as observed in GAN-based models. More recently, diffusion-based models [38,43] focused on whole face generation were proposed with better image quality and temporal coherence. However, both models employ complex motion control modules and require large-scale internet-sourced data, which significantly hinders training and inference speed.

To address the limitations inherent in the aforementioned directions, we introduce SYNCDIFF, a diffusion model that incorporates temporal-augmented visual priors to condition the diffusion process [29], focused on lip region inpainting without the need for crowd-sourced data. Additionally, our model integrates a pre-trained audio-video self-supervised framework to enhance the synchronization of the generated videos.

^{*}Equal contribution.

Many studies leverage temporal information by employing a temporally aware audio encoder, yet they solely utilize a single image for identity information [10, 32, 39, 44, 46], thereby neglecting the temporal information present in the video modality. Due to the high similarity within short-term temporal images, directly incorporating neighborhood image information may lead to training shortcuts [28, 32]. In response to this concern, we propose a bottleneck layer to leverage temporal pose information and mitigate the risk of training shortcuts.

We conduct comprehensive experiments on two benchmark datasets, demonstrating that our proposed SYNCDIFF attains superior visual quality compared to all state-of-theart synthesis methods. Additionally, our approach significantly improves lip synchronization within the diffusion-based synthesis model, showing competitive performance compared to other advanced lip synchronization methods. In summary, our primary contributions are as follows:

- We present a novel approach by introducing a bottleneck layer to incorporate visual temporal pose information in talking head synthesis. This pioneering addition of visual temporal pose information results in a substantial improvement in lip synchronization.
- We leverage a self-supervised audio-visual pre-trained model, AVHUBERT [33, 34], within the diffusion model to facilitate talking head generation. This integration contributes to further enhancing lip synchronization.
- Our proposed SYNCDIFF conditional diffusion model excels in generating high-fidelity image quality and well-synchronized lip movement. Our experimental results demonstrate that SYNCDIFF surpasses other SOTA methods in terms of visual quality while significantly enhancing lip synchronization capabilities in diffusion models.

2. Related Work

2.1. 2D Audio-driven Talking Head Synthesis

Previous works on talking head synthesis can be categorized into two categories, 3D-based methods, which model in the 3D facial geometry space, and 2D based methods, which directly operate in the 2D pixel space, irrespective of the underlying 3D geometry.

The most relevant previous works are 2D-based methods with generative adversarial networks (GANs) emerging as a favored choice for talking head synthesis focused on lip region impainting [12, 28, 39]. WAV2LIP [28] proposes to leverage a cross-modality sync loss from a pretrained lip-sync expert to improve synchronization, whose synchronization performance is still among the highest nowa-

days. TALKLIP [39], one of the latest GAN-based methods, proposes to further incorporate lipreading loss which explicitly aligns the synthetic video with the text modality and further improves the video's lip-reading intelligibility, which is measured by the word error rate of lipreading the synthetic videos. STYLESYNC [12] leverages a STYLE-GAN [21] based architecture to solve the task by encoding the facial and audio information into the style space. In addition to GAN-based methods, diffusion-based methods emerge as an alternative approach to the task [10, 32]. DAETALKER [10] trains a speech-to-latent unit whose target is the visual embedding from the video stream. During inference, the latent from input speech is fed into a diffusion decoder to synthesize the video. However, this approach needs to train an identity-specific speech-tolatent encoder and does not generalize to unseen identities. DIFFTALK [32] is another work that proposes to use a combination of a masked frame and a reference frame to condition a diffusion model to synthesize the output. This approach generates high-fidelity videos and generalizes to unseen identities with a given reference image. However, it lags behind the GAN-based models in terms of synchronization, as measured by the SYNCNET [7] score and lipreading word error rate.

2.2. Audio Features for Talking Head Synthesis

As the audio features serve as the driving factor for video frame synthesis, it is crucial to enhance their quality in encoding temporal and lip movement enabling information. Instead of using the Mel-spectrogram features that have a limited context window as in WAV2LIP and SADTALKER [46], DIFFTALK chooses the DEEP-SPEECH [14] features, which uses bidirectional recurrent layers to encode speech spectrograms. With recent advances in self-supervised speech representations, audio features extracted from pretrained acoustic models such as WAV2VEC2 [3], HUBERT [16] have been shown to encode richer contextual information and have exhibited great performance advantages on various downstream tasks over traditional methods [42]. DAETALKER uses WAV2VEC2 features as speech representations for talking head synthesis. These aforementioned audio features are purely audio-driven. AVHUBERT pretrains a self-supervised HU-BERT-like Transformer-based model using masked multimodal cluster prediction objective on speech-video datasets to learn speech and video representations that have better correlations between the two modalities. TALKLIP [39] uses AVHUBERT features to drive the talking head generations and achieve the SOTA synchronization performance.

2.3. Diffusion Model

Denoising diffusion probabilistic models (DDPMs) have recently demonstrated great power in image synthesis [9,

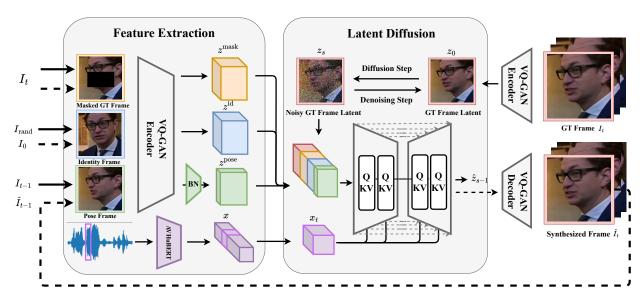


Figure 1. Architecture of the SYNCDIFF network. Solid lines denote inputs during training while dashed lines denote inputs during inference. I_t denotes the ground-truth frame at timestep t, \hat{I}_t denotes the synthesized frame at timestep t, and I_{rand} denotes randomly sampled frame from the groundtruth sequence. BN denotes the bottleneck layer which compresses the dimension of the pose prior.

15] and image impainting [26] tasks.

The main drawback of DDPM models is their high computational cost in both training and inference due to the pixel-based input and output space [31] and large number of diffusion steps [15]. To mitigate the cost incurred by operating on pixel space, latent diffusion models (LDMs) [29] are proposed, which use a VQGAN autoencoder [11] to reduce the high-resolution image to low-dimensional latent space for the training of the DDPMs. LDMs have been applied to many multimodal tasks including text-to-image translation [22, 30], text-to-audio generation [17] and speech-tovideo generation [32], etc. To reduce computational cost resulting from a large number of diffusion steps, denoising diffusion implicit models (DDIM) [35], consistency models [36] and EDM [20] are proposed. DDIM generalize the Markovian diffusion process assumed by DDPM into a class of non-Markovian diffusion process that lead to the same training objective as that of DDPM but with shorter Markov chains and thus fewer diffusion steps. DDIM follows the same training process as DDPM, it can be directly applied to accelerate the inference process of DDPM-based models [18, 32].

3. Method

3.1. Overview

A typical diffusion-based talking head synthesis model like DIFFTALK is conditioned on a given audio clip and a reference frame where the reference frame is used as a driving factor of the human identity in the generated videos. However, its deficiency in audio-visual synchronization is evident, suggested by its low SYNCNET scores and high lip-reading word error rate. In this paper, we propose to incorporate a pose frame that serves as temporal visual priors to further improve the synthesized videos' audio-visual synchronization. We additionally utilize a bottleneck layer over the pose frame to prevent the model from learning shortcuts. The overall architecture is presented in Figure 1.

3.2. Audio Feature Extraction

We use AVHUBERT to extract facial-informative audio features, which is a self-supervised multimodal Transformer pretrained on audio-visual datasets. Specifically, we first extract the Mel-spectrogram with a feature rate of 100 Hz using a moving window that has a window size of 25 ms and a hop-size of 10 ms. Every four consecutive Mel-spectrogram features are grouped as one feature to be input to a linear layer, resulting in the grouped Melspectrogram features with a rate of 25 Hz, the same rate as the video frames. The grouped features are then input to 12 layers of AVHUBERT Transformer layers to extract contextualized audio features. The first 9 layers of AVHU-BERT are frozen during training, and the remaining 3 layers are jointly trained with the diffusion model. To compare AVHUBERT features with other audio features, we additionally extract DEEPSPEECH features following the preprocessing in DIFFTALK. The extracted features are fed to a learnable temporal filtering network [37] for improved inter-frame consistency.

The audio features corresponding to the target frame are used as the conditioning inputs to the diffusion model.

3.3. Temporal-Augmented Visual Priors

Previous works [10,28,32,39] generally use the masked ground-truth (GT) frame, a randomly sampled reference frame and the audio features as priors to condition the GAN-based or diffusion-based generator during training, where the masked frame provides the head pose information and the reference frame provides the appearance identity of the lip region. During testing, the given first frame [28,39] or the synthesized previous frame [32] replaces the randomly sampled frame as the reference frame to provide identity information.

We argue that the model trained with a randomly sampled reference frame fails to learn the temporal transition between consecutive frames, which leads to interframe incoherence. Such incoherence is especially evident in the diffusion-based model as discussed in DIFFTALK. DIFFTALK mitigate the incoherence by training with random reference but inferring with generated previous frames. However, doing so introduces a train-test distribution mismatch, which can negatively influence the model's synchronization performance. One naive approach is to use the previous frames immediately preceding the target frame, instead of randomly sampled ones, as reference frames to provide the temporal information during training. However, this approach incurs significant model degradation because the neighborhood frames are usually so similar that the model can achieve a low loss by directly copying the reference frame. Such degradation is referred to as training shortcut in previous works [28, 32].

We discover that applying a simple bottleneck layer over the temporal pose frame resolves the shortcut issue. The bottleneck layer is a 2D convolutional layer and is trained to compress the reference features so that they contain mainly temporal pose priors. Furthermore, we propose to incorporate two reference frames, in addition to the masked GT frame, as priors to provide identity and temporal pose information separately. During training, one random sampled frame serves as the identity reference frame and one previous frame with bottleneck serves as the temporal pose frame, while during inference, the given first frame and the generated previous frame serve as the identity and pose frame.

The proposed triple-prior approach has two advantages compared to the traditional double-prior approach. Firstly, the explicit separation of speaker identity and temporal lip pose information into two sources can lead to a distribution easier to learn. Our empirical results in Sec 5.4 suggest that the triple-prior model can more effectively disentangle identity and temporal pose information while remaining invariant to the change of other aspects, e.g. lighting condition and camera position. Secondly, compared to DIFFTALK approach, our triple prior approach has unified reference frame distribution during training and inference.

3.4. Conditional Latent Diffusion

We model the generation process as a denoising process using a LDM. To project between the pixel space and the latent space, we adopt the VQGAN-based autoencoder [11, 29, 32], which consisting of a pair of encoder $\mathcal E$ and decoder $\mathcal D$ with a downsampling and upsampling factor of 4, respectively. As discussed in the previous section, three image frames that act as visual priors will go through the encoder $\mathcal E$ whose output space is $h \times w \times 3$. The temporal pose prior will go through one additional bottleneck convolution layer to further compress its dimension before it's concatenated with the other two priors. The visual feature extraction step can be formulated as:

$$E_m, E_i, E_p = \mathcal{E}(I_m, I_i, I_p) \tag{1}$$

$$E_v = \text{Concat}(E_m, E_i, BN(E_p))$$
 (2)

where I_m, I_i, I_p are respectively the masked frame, identity frame, and pose frame and BN is the bottleneck layer.

The extracted visual and audio priors E_v , E_a are then fed into a denoising UNET [31] \mathcal{M} together with a sampled timestep d and the respective z_d which denotes the latent representation of the target at the d timestep of the forward diffusion process. The architecture of the entire SyncDiff method is shown in 1. The loss function can then be formulated as

$$\mathcal{L}_{LDM} := \mathbb{E}_{z,\epsilon \sim \mathcal{N}(0,1), E_v, E_a, d}[\|\epsilon - \mathcal{M}(z_d, d, E_v, E_a)\|^2]$$
(3)

Our conditional UNET utilizes a cross-attention mechanism to learn across visual and audio modalities. The visual prior $E_v \in \mathcal{R}^{h \times w \times 7}$ is further concatenated with the noisy latent z_d and form an embedding $E_q \in \mathcal{R}^{h \times w \times 10}$ which is directly fed into the network and used to calculate the query of the cross-attention. The audio feature E_a is fed into the intermediate layers and used to calculate the key and value.

Denote the estimation of the latent at timestep d as \tilde{z}_d . During inference, the denoised latent at the final step d=0 can be decoded back to the pixel space using the decoder \mathcal{D} .

$$I_{pred} := \mathcal{D}(\tilde{z}_0) \tag{4}$$

4. Experiments

4.1. Datasets

We conduct our experiments using two canonical datasets, LRS2 [1] for training and evaluation and LRS3 [2] for evaluation only. LRS2 dataset contains around 29 hours of talking-head videos with the ground-truth transcripts of the speakers' speech. We follow the standard train/valid/test split in LRS2 dataset. LRS3 dataset contains thousands of spoken sentences from TED and TEDx videos. We use only the test set for evaluation, which contains 1321 video samples with a total duration of 51 minutes.

| | LRS2 | | | | LRS3 | | | | | | | |
|-----------|-------|-------|--------|--------|--------|--------|-------|-------|--------|--------|--------|--------|
| Метнор | PSNR↑ | SSIM↑ | LPIPS↓ | LSE-C↑ | LSE-D↓ | WER↓ | PSNR↑ | SSIM↑ | LPIPS↓ | LSE-C↑ | LSE-D↓ | WER↓ |
| GT | N/A | 1.000 | 0.000 | 8.25 | 6.26 | 23.82 | N/A | 1.000 | 0.000 | 7.63 | 6.89 | 41.04 |
| GT-REC | 35.30 | 0.956 | 0.022 | 8.02 | 6.35 | 28.47 | 36.06 | 0.958 | 0.015 | 7.48 | 6.95 | 44.44 |
| WAV2LIP | 31.07 | 0.859 | 0.079 | 7.83 | 6.65 | 90.67 | 31.04 | 0.846 | 0.079 | 8.10 | 6.68 | 91.58 |
| SADTALKER | 29.27 | 0.543 | 0.109 | 5.45 | 8.34 | 110.64 | 29.08 | 0.641 | 0.144 | 5.27 | 8.65 | 107.13 |
| TALKLIP | 31.16 | 0.850 | 0.084 | 8.53 | 5.70 | 23.43 | 31.10 | 0.852 | 0.084 | 8.11 | 6.41 | 22.01 |
| DIFFTALK | 32.36 | 0.873 | 0.056 | 6.11 | 7.93 | 114.34 | 31.12 | 0.779 | 0.076 | 4.80 | 9.20 | 114.16 |
| SYNCDIFF | 32.39 | 0.874 | 0.049 | 7.80 | 6.58 | 67.40 | 32.18 | 0.846 | 0.058 | 6.82 | 7.58 | 83.94 |

Table 1. Comparison of visual quality and lip synchronization scores on LRS2 and LRS3 datasets. GT-REC are ground-truth videos encoded and reconstructed using our image autoencoder and serve as a topline for our SYNCDIFF models. The best results are highlighted in bold **black** and the second best results are highlighted in bold **violet**.

4.2. Data Preprocessing

Videos in both datasets have a frame rate of 25 frames per second (fps) and 3 RGB channels. The LRS2 videos have a dimension of $160 \times 160 \times 3$ and LRS3 videos have a dimension of $224 \times 224 \times 3$. During training, all videos are resized to a resolution of $256 \times 256 \times 3$. Following WAV2LIP [28] and TALKLIP [39], we use FACE-ALIGNMENT [4] to detect face region and use the lower half of the face region as lip region.

Audios are resampled to 16000 Hz. We extract Melspectrogram features with 10-ms hop sizes of 25-ms window lengths and 26-filter filter bank, resulting in spectrogram features with 26 bins and a feature rate of 100 Hz. We then group four consecutive features together to create 25-Hz spectrogram features with a dimensionality of 104. To replicate DIFFTALK [32] results, we additionally extract 25-Hz DEEPSPEECH [14] features with dimensions of 16×29 following the practice in VOCA [8]. We apply a window of 16 consecutive features centered at each timestamp to create 25-Hz DEEPSPEECH features with a dimensionality of $16 \times 16 \times 29$.

4.3. Implementation Details

We use a VQGAN-based image autoencoder [29] with a downsampling rate of 4 to compress the input images from the resized dimension of $256 \times 256 \times 3$ to a latent dimension of $64 \times 64 \times 3$, which is the input dimension of the LDM. We apply a 1×1 2D convolution layer on the channel dimension of the latent embedding of the previous frame to further compress its dimension to $64 \times 64 \times 1$. The conditioning input of the LDM has the same dimension as the audio features. In the case of AVHUBERT features, the audio features have a dimension of 768. In the case of DEEPSPEECH features, the audio features have a dimension as in DIFFTALK. The number of diffusion steps of LDM is set to 1000 during training and the number of steps of the DDIM sampler is set to 200 during

inference. The entire SYNCDIFF model is trained for 200 epochs on 6 GPUs with a batch size of 8 on each GPU.

4.4. Metrics

Following recent works on talking head generation [10, 27,28,32,39,40], we use peak signal-to-noise ratio (PSNR), structural similarity index measure (SSIM) and learned perceptual image patch similarity (LPIPS) [45] scores to measure visual quality, use "lip sync error - confidence" (LSE-C) and "lip sync error - distance" (LSE-D) scores generated by a pretrained SYNCNET [7] to measure lip-speech synchronization and use word error rate (WER) to measure lip reading intelligibility. We use AVHUBERT-LARGE, a lip-reading model finetuned on LRS2 dataset, to transcribe the generated video, and calculate WER based on the resulting transcripts.

4.5. Comparison Methods

We compare SYNCDIFF against four SOTA methods: WAV2LIP [28], SADTALKER [46], TALKLIP [39] and DIFFTALK [32]. WAV2LIP is a GAN-based model with its discriminator being a pretrained SYNCNET-based lip sync model and achieves high lip synchronization scores. SADTALKER is a 3D-based talking face model conditioned on generated realistic 3D pose coefficients. TALKLIP improves on WAV2LIP by replacing Mel-spectrogram features with pretrained AVHUBERT features and replacing SYNCNET lip-syncing model with AVHUBERT lipreading model. DIFFTALK uses diffusion models to generate face images and achieves the SOTA visual quality on talking head generation.

5. Results

5.1. Main Results

Table 1 compares SYNCDIFF against ground truth videos, ground truth videos reconstructed using the im-

age autoencoder, and four state-of-the-art models quantitatively. The performances are measured using visual quality (PSNR, SSIM and LPIPS), lip-speech synchronization (LSE-C and LSE-D), and lip reading intelligibility (WER) scores. on the in-domain LRS2 dataset, we make three key observations. First, our SYNCDIFF can generate talking head videos with the best image quality compared to GAN-based WAV2LIP and TALKLIP and diffusion-based DIFFTALK. Second, SYNCDIFF, trained with AVHU-BERT features and additional previous frame, achieves significantly better synchronization scores than DIFFTALK, and the LSE-C and LSD-D scores of SYNCDIFF approach to the reconstructed ground truth topline. Finally, if we exclude TALKLIP, which explicitly trains using paired video and text data, SYNCDIFF outperforms other methods that do not use additional text information on lip-reading intelligibility.

The same observation generally holds on the out-of-domain LRS3 dataset. However, we do notice a more evident drop in synchronization using diffusion-based DIFFTALK and SYNCDIFF compared to GAN-based WAV2LIP and TALKLIP, which, we hypothesize, is due to SYNCDIFF reconstruct lip movement as well as the surrounding texture while WAV2LIP and TALKLIP focus less on texture but more on lip movement. Such behavior may limit the generalizability of SYNCDIFF to unseen identities.

We should make a side note that we use the officially released SADTALKER checkpoint, which is not pretrained on LRS2. This choice possibly leads to its suboptimal performance in our evaluation.

5.2. Perceptual Evaluation

We provide uniformly sampled snapshots of two generated talking face videos in Figure 2 to give a better visual demonstration of the advantage of SYNCDIFF compared against other SOTA methods. We can see that compared to GAN-based WAV2LIP and TALKLIP, Diffusion-based DIFFTALK and SYNCDIFF generate the talking face frames with sharper textures and more details. Also, because the DIFFTALK and SYNCDIFF are trained to generate the entire image instead of the face region, we do not see an obvious boundary mismatch at the face region that appears frequently in videos generated by WAV2LIP and TALKLIP. Comparing DIFFTALK and SYNCDIFF, we observe that video frames from the latter exhibit lip shapes closer to the ground truth frames than those from the former, for example, at frames for "A", "V", "G".

5.3. Comparison on Audio Features

We compare the performance of using DEEPSPEECH and AVHUBERT audio features. The results are shown in Table 2. We observe that the synchronization and intelligibility scores are greatly improved using AVHUBERT fea-

| FEAT | PSNR↑ | SSIM↑ | LPIPS↓ | LSE-C↑ | LSE-D | WER↓ |
|------|-------|-------|--------|--------|-------|------------------------|
| DS | 32.36 | 0.873 | 0.056 | 6.11 | 7.93 | 114.34 |
| AVH | 32.40 | 0.877 | 0.055 | 7.05 | 7.11 | 114.34 90.24 |

Table 2. Comparison between DEEPSPEECH (DS) and AVHUBERT (AVH) speech features with randomly sampled video frame as visual prior during training.

| REF FRAME | PSNR↑ | SSIM↑ | LPIPS↓ | LSE-C↑ | LSE-D↓ | WER↓ |
|-------------|-------|-------|--------|--------|--------|-------|
| RND | 32.40 | 0.877 | 0.055 | 7.05 | 7.11 | 90.24 |
| PRV | 31.61 | 0.851 | 0.057 | 7.38 | 6.91 | 76.16 |
| PRV^- | | | 0.053 | | 6.55 | |
| $RND+PRV^-$ | 32.39 | 0.874 | 0.049 | 7.80 | 6.58 | 67.40 |

Table 3. Comparison between using different video frames as input to predict the next frame during training. RND, PRV and PRV denotes using a randomly selected frame, previous frame and the bottlenecked previous frame, respectively. The best results are highlighted in bold **black** and the second best results are highlighted in bold **violet**.

tures. We believe this is because RNN-based DEEPSPEECH audio representations have three drawbacks compared to Transformer-based AVHUBERT. First, DEEPSPEECH contains fewer parameters than AVHUBERT and therefore lower model capacities. Second, DEEPSPEECH uses bidirectional RNN layers, that do not sufficiently model long-term context compared to self-attention layers, resulting in better audio representations and better overall video quality. Finally, DEEPSPEECH is trained on pure audio while AVHUBERT is trained to match the audio representation with corresponding lip shapes. With additional facial priors, the audio features extracted from the latter are more conducive to lip movement generation and thus higher synchronization and intelligibility scores.

5.4. Effect of Reference Frames

We conduct experiments to compare the effect of using different reference frames to condition the diffusion generator: a randomly sampled video frame (RND), the previous frame of the target frame (PREV), the bottlenecked previous frame (PREV⁻), and the combination of the random and bottlenecked previous frame (RND+PREV⁻). The objective results are shown in Table 3 and the snapshots of the generated videos are shown in Figure 3. By comparing the RND and the PREV results, we verify the existence of the short-cut issue that leads to the grey-out of the lip region. By comparing PREV and PREV⁻, we observe that the bottleneck layer after the previous frame resolves this issue. By comparing RND and PREV⁻, we observe that using the bottlenecked previous frame greatly improves lip synchronization with a slight sacrifice of image quality. The best

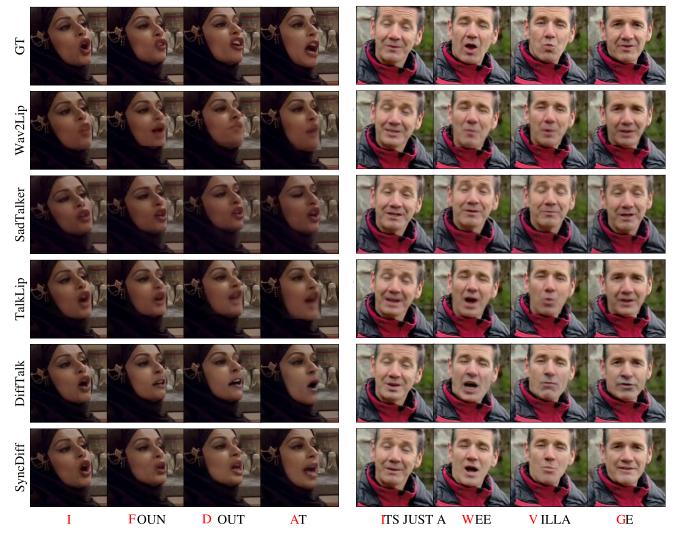


Figure 2. Visual comparison with SOTA talking head generation methods. The letter that each frame corresponds to is marked in red.



Figure 3. Visual comparison of using different reference frames.

strategy is RND+PREV⁻, which combines the advantages of the two to achieve both high image quality and high synchronization.

To further study the behavior of SYNCDIFF with the additional input of the previous frame, we measure the variance of latent representations encoded by the UNET of the diffusion model, across different identity/pose frames that share the same identity/pose. Specifically, Two sets of experiments, namely IDENT and POSE, are conducted on two models, RND and RND+PREV⁻. In IDENT, we randomly

select 10 videos from LRS2. For each video, we choose the second frame to be the target and fix the pose frame to be the first frame. The subsequent frames in the same video serve as the identity frames sharing the same speaker identity. One latent representation per identity frame is extracted from the diffusion model. We compute the variances in each dimension and average them across dimensions and videos. In POSE, we choose 5 videos from LRS2 featuring the same speaker uttering the same sentence: "Thanks for watching". Frames corresponding to the phone "wa" are manually labeled and treated as pose frames sharing the same pose. We set the identity frames as the first frame of the first video, select the frame after each pose frame as the GT frame, and vary the masked GT frame and the pose frame across 5 videos. Similar to IDENT experiment, we compute the averaged variance across latent dimensions. The results are shown in Table 4. In both experiments, we observe RND+PREV⁻ has notably lower variance. This ob-

| | | Ident | Pose | | |
|------|-------------|-------|-------|---------|--|
| REF | RND RND+PRV | | RND | RND+PRV | |
| VAR↓ | 0.340 | 0.075 | 0.530 | 0.328 | |

Table 4. Comparison between the average variance over latent dimensions when varying identity frame or pose frame while keeping the other fixed. IDENT denotes varying the identity frames and fixing the pose frames and POSE denotes the reverse. VAR denotes the averaged variance.

| NP | PSNR↑ | SSIM↑ | LPIPS↓ | LSE-C↑ | LSE-D↓ | WER↓ |
|----|-------|-------|--------|-----------------------|--------|-------|
| 1 | 32.39 | 0.874 | 0.049 | 7.80 | 6.58 | |
| 5 | 32.34 | 0.869 | 0.052 | 7.46 | 6.85 | 84.13 |
| 10 | 32.45 | 0.872 | 0.051 | 7.80 7.46 7.48 | 6.82 | 86.20 |

Table 5. Comparison on using different number previous frames (NP) as input to LDM.

servation implies that the latent features with RND+PREV⁻ more effectively capture both identity and pose-related information, remaining invariant to changes in other aspects.

5.5. Ablation Study on Number of Previous Frames

We experiment on providing more temporal information to SYNCDIFF by increasing the number of previous frames fed to the LDM. The results are shown in Table 5. Unfortunately, we observe no improvement in lip-synchronization; Using one previous frame as a reference yields the best performance. We believe this is because the LDM architecture lacks recurrent or self-attention layers over the input sequence to capture temporal dependencies of lip movement, which we leave as future works to explore.

Table 6. Comparison on freezing different number of layers (NF) in AVHUBERT.

| NF | PSNR↑ | SSIM↑ | LPIPS↓ | LSE-C↑ | LSE-D↓ | WER_{\downarrow} |
|----|-------|-------|--|--------|--------|--------------------|
| 3 | 32.35 | 0.872 | 0.050 | 7.56 | 6.75 | 80.90 |
| 6 | 32.37 | 0.874 | 0.049 | 7.64 | 6.67 | 72.33 |
| 9 | 32.39 | 0.874 | 0.049 | 7.80 | 6.58 | 67.40 |
| 12 | 31.97 | 0.855 | 0.050 0.049 0.049 0.057 | 7.46 | 6.87 | 75.92 |

5.6. Ablation Study on Freezing AVHuBERT

AVHUBERT It is known that the features from different layers of HUBERT based model capture different levels of linguistic information [33]. We therefore experiment on freezing the first few layers of the pretrained 12-layer AVHUBERT model during the training. The results are shown in Table 6. We observe a U-shape curve in performance with the increase of number of frozen layers and

freezing the first nine layers of AVHUBERT yields the best overall performance.

6. Limitation

Our model mainly suffers from three limitations. Firstly, the inference of the diffusion model is slow. Although we apply DDIM sampler to speed up the inference, it still takes a significant amount of time to generate one single video compared to GAN-based methods. Diffusion models with fewer diffusion steps such as consistency models and EDM can be explored to further speed up the inference. Secondly, although SYNCDIFF achieves significant improvement in lip synchronization, we observe an evident performance degradation on out-of-distribution testing data, compared to the GAN-based methods. Adding explicit crossmodality contrastive loss for better audio-visual representation as in WAV2LIP and TALKLIP might be one promising solution. However, it is not easy to directly apply such loss in diffusion models because of the heavy noise in the intermediate diffusion steps. We leave this as future works to explore. Last but not least, our ablation studies show that one previous frame is a better temporal prior than 5 or 10 previous frames, suggesting SYNCDIFF does not leverage long-term temporal patterns efficiently. Self-attention or recurrent layers can be experimented to extract better temporal features.

7. Conclusion

In this work, we present SYNCDIFF, a diffusion-based talking head synthesis model that simultaneously achieves both high visual quality and good synchronization. Our contribution can be summarized in mainly three aspects. First, we introduce a bottleneck layer to incorporate visual temporal pose information in the diffusion-based talking head synthesis model, which resolves the learning short-cut issue of using previous frames during training and leads to a substantial improvement in lip synchronization. This modification opens the possibility of incorporating further temporal video information in the training of existing GANbased and diffusion-based talking head generators. Second, we leverage a self-supervised audio-visual pretrained model, AVHUBERT, to facilitate diffusion-based talking head generation which contributes to further enhanced lip synchronization. Finally, We conduct extensive experiments to demonstrate the advantage of SYNCDIFF over other state-of-the-art methods in terms of visual quality and lip synchronization.

References

[1] Triantafyllos Afouras, Joon Son Chung, Andrew Senior, Oriol Vinyals, and Andrew Zisserman. Deep audio-visual

- speech recognition. *IEEE transactions on pattern analysis and machine intelligence*, 44(12):8717–8727, 2018. 4
- [2] Triantafyllos Afouras, Joon Son Chung, and Andrew Zisserman. Lrs3-ted: a large-scale dataset for visual speech recognition. ArXiv, abs/1809.00496, 2018. 4
- [3] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. Advances in neural information processing systems, 33:12449–12460, 2020. 2
- [4] Adrian Bulat and Georgios Tzimiropoulos. How far are we from solving the 2d & 3d face alignment problem? (and a dataset of 230,000 3d facial landmarks). In *International Conference on Computer Vision*, 2017. 5
- [5] Lele Chen, Guofeng Cui, Celong Liu, Zhong Li, Ziyi Kou, Yi Xu, and Chenliang Xu. Talking-head generation with rhythmic head motion. *arXiv preprint arXiv:2007.08547*, 2020.
- [6] Lele Chen, Ross K Maddox, Zhiyao Duan, and Chenliang Xu. Hierarchical cross-modal talking face generation with dynamic pixel-wise loss. In *Proceedings of the IEEE Con*ference on Computer Vision and Pattern Recognition, pages 7832–7841, 2019.
- [7] J. S. Chung and A. Zisserman. Out of time: automated lip sync in the wild. In Workshop on Multi-view Lip-reading, ACCV, 2016. 2, 5
- [8] Daniel Cudeiro, Timo Bolkart, Cassidy Laidlaw, Anurag Ranjan, and Michael J. Black. Capture, learning, and synthesis of 3d speaking styles. 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 10093–10103, 2019. 5
- [9] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 8780–8794. Curran Associates, Inc., 2021.
- [10] Chenpng Du, Qi Chen, Tianyu He, Xuejiao Tan, Xie Chen, K. Yu, Sheng Zhao, and Jiang Bian. Dae-talker: High fidelity speech-driven talking face generation with diffusion autoencoder. Proceedings of the 31st ACM International Conference on Multimedia, 2023. 1, 2, 4, 5
- [11] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12873–12883, 2021. 3, 4
- [12] Jiazhi Guan, Zhanwang Zhang, Hang Zhou, Tianshu Hu, Kaisiyuan Wang, Dongliang He, Haocheng Feng, Jingtuo Liu, Errui Ding, Ziwei Liu, and Jingdong Wang. Stylesync: High-fidelity generalized and personalized lip sync in style-based generator. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 1505–1515, June 2023. 1, 2
- [13] Yudong Guo, Keyu Chen, Sen Liang, Yongjin Liu, Hujun Bao, and Juyong Zhang. Ad-nerf: Audio driven neural radiance fields for talking head synthesis. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 1
- [14] Awni Y. Hannun, Carl Case, Jared Casper, Bryan Catanzaro, Gregory Frederick Diamos, Erich Elsen, Ryan J. Prenger,

- Sanjeev Satheesh, Shubho Sengupta, Adam Coates, and A. Ng. Deep speech: Scaling up end-to-end speech recognition. *ArXiv*, abs/1412.5567, 2014. 2, 5
- [15] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Advances in Neural Information Processing Systems, 33:6840–6851, 2020. 2, 3
- [16] Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:3451–3460, 2021. 2
- [17] Rongjie Huang, Jiawei Huang, Dongchao Yang, Yi Ren, Luping Liu, Mingze Li, Zhenhui Ye, Jinglin Liu, Xiang Yin, and Zhou Zhao. Make-an-audio: Text-to-audio generation with prompt-enhanced diffusion models. arXiv preprint arXiv:2301.12661, 2023. 3
- [18] Rongjie Huang, Zhou Zhao, Huadai Liu, Jinglin Liu, Chenye Cui, and Yi 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. 3
- [19] Amir Jamaludin, Joon Son Chung, and Andrew Zisserman. You said that?: Synthesising talking faces from audio. *International Journal of Computer Vision*, 127:1767–1779, 2019.
- [20] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. ArXiv, abs/2206.00364, 2022. 3
- [21] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 8107–8116, 2019. 2
- [22] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 6007–6017, 2023. 3
- [23] 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), 37(4):1–14, 2018.
- [24] Avisek Lahiri, Vivek Kwatra, Christian Frueh, John Lewis, and Chris Bregler. Lipsync3d: Data-efficient learning of personalized 3d talking faces from video using pose and lighting normalization. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 2755–2764, June 2021.
- [25] Yuanxun Lu, Jinxiang Chai, and Xun Cao. Live Speech Portraits: Real-time photorealistic talking-head animation. ACM Transactions on Graphics, 40(6), 2021.
- [26] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11461–11471, 2022. 3

- [27] Se Jin Park, Minsu Kim, Joanna Hong, Jeongsoo Choi, and Yong Man Ro. Synctalkface: Talking face generation with precise lip-syncing via audio-lip memory. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 2062–2070, 2022. 5
- [28] Prajwal K R, 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. *Proceedings of the 28th* ACM International Conference on Multimedia, 2020. 1, 2, 4, 5
- [29] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 1, 3, 4, 5
- [30] Robin Rombach, Andreas Blattmann, and Björn Ommer. Text-guided synthesis of artistic images with retrieval-augmented diffusion models. arXiv preprint arXiv:2207.13038, 2022. 3
- [31] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. Unet: Convolutional networks for biomedical image segmentation. In Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18, pages 234–241. Springer, 2015. 3, 4
- [32] Shuai Shen, Wenliang Zhao, Zibin Meng, Wanhua Li, Zhengbiao Zhu, Jie Zhou, and Jiwen Lu. Difftalk: Crafting diffusion models for generalized audio-driven portraits animation. 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 1982–1991, 2023.
- [33] Bowen Shi, Wei-Ning Hsu, Kushal Lakhotia, and Abdelrahman Mohamed. Learning audio-visual speech representation by masked multimodal cluster prediction. arXiv preprint arXiv:2201.02184, 2022. 2, 8
- [34] Bowen Shi, Wei-Ning Hsu, and Abdelrahman Mohamed. Robust self-supervised audio-visual speech recognition. arXiv preprint arXiv:2201.01763, 2022. 2
- [35] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 3
- [36] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. In *International Conference* on Machine Learning, 2023. 3
- [37] Justus Thies, Mohamed Elgharib, Ayush Tewari, Christian Theobalt, and Matthias Nießner. Neural voice puppetry: Audio-driven facial reenactment. In *Computer Vision–ECCV* 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVI 16, pages 716–731. Springer, 2020. 3
- [38] Linrui Tian, Qi Wang, Bang Zhang, and Liefeng Bo. Emo: Emote portrait alive-generating expressive portrait videos with audio2video diffusion model under weak conditions. arXiv preprint arXiv:2402.17485, 2024. 1
- [39] Jiadong Wang, Xinyuan Qian, Malu Zhang, Robby T. Tan, and Haizhou Li. Seeing what you said: Talking face generation guided by a lip reading expert. 2023 IEEE/CVF Confer-

- ence on Computer Vision and Pattern Recognition (CVPR), pages 14653–14662, 2023. 1, 2, 4, 5
- [40] Jiayu Wang, Kang Zhao, Shiwei Zhang, Yingya Zhang, Yujun Shen, Deli Zhao, and Jingren Zhou. Lipformer: High-fidelity and generalizable talking face generation with a prelearned facial codebook. 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 13844–13853, 2023. 5
- [41] Ting-Chun Wang, Arun Mallya, and Ming-Yu Liu. One-shot free-view neural talking-head synthesis for video conferencing. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 10039–10049, 2021. 1
- [42] Shu wen Yang, Po-Han Chi, Yung-Sung Chuang, Cheng-I Lai, Kushal Lakhotia, Yist Y. Lin, Andy T. Liu, Jiatong Shi, Xuankai Chang, Guan-Ting Lin, Tzu hsien Huang, Wei-Cheng Tseng, Ko tik Lee, Da-Rong Liu, Zili Huang, Shuyan Dong, Shang-Wen Li, Shinji Watanabe, Abdel rahman Mohamed, and Hung yi Lee. Superb: Speech processing universal performance benchmark. In *Interspeech*, 2021. 2
- [43] Mingwang Xu, Hui Li, Qingkun Su, Hanlin Shang, Li-wei Zhang, Ce Liu, Jingdong Wang, Luc Van Gool, Yao Yao, and Siyu Zhu. Hallo: Hierarchical audio-driven visual synthesis for portrait image animation. arXiv preprint arXiv:2406.08801, 2024.
- [44] Zhenhui Ye, Ziyue Jiang, Yi Ren, Jinglin Liu, Jinzheng He, and Zhou Zhao. Geneface: Generalized and high-fidelity audio-driven 3d talking face synthesis. *arXiv preprint arXiv:2301.13430*, 2023. 1, 2
- [45] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 5
- [46] Wenxuan Zhang, Xiaodong Cun, Xuan Wang, Yong Zhang, Xiaodong Shen, Yu Guo, Ying Shan, and Fei Wang. Sadtalker: Learning realistic 3d motion coefficients for stylized audio-driven single image talking face animation. 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 8652–8661, 2022. 1, 2, 5
- [47] Hang Zhou, Yu Liu, Ziwei Liu, Ping Luo, and Xiaogang Wang. Talking face generation by adversarially disentangled audio-visual representation. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 9299–9306, 2019. 1
- [48] Yang Zhou, Xintong Han, Eli Shechtman, Jose Echevarria, Evangelos Kalogerakis, and Dingzeyu Li. Makelttalk: speaker-aware talking-head animation. *ACM Transactions On Graphics (TOG)*, 39(6):1–15, 2020. 1