# RealTalk: Real-time and Realistic Audio-driven Face Generation with 3D Facial Prior-guided Identity Alignment Network

Xiaozhong Ji[1], Chuming Lin[1], Zhonggan Ding[1], Ying Tai[2], Jian Yang[2], Junwei Zhu[1], Xiaobin Hu[1], Jiangning Zhang[1], Donghao Luo[1], and Chengjie Wang[1]

[1] Youtu Lab, Tencent
xiaozhongji@tencent.com
[2] Nanjing University
yingtai@nju.edu.cn

**Abstract.** Person-generic audio-driven face generation is a challenging task in computer vision. Previous methods have achieved remarkable progress in audio-visual synchronization, but there is still a significant gap between current results and practical applications. The challenges are two-fold: 1) Preserving unique individual traits for achieving high-precision lip synchronization. 2) Generating high-quality facial renderings in real-time performance. In this paper, we propose a novel generalized audio-driven framework `RealTalk`, which consists of an audio-to-expression transformer and a high-fidelity expression-to-face renderer. In the first component, we consider both identity and intra-personal variation features related to speaking lip movements. By incorporating cross-modal attention on the enriched facial priors, we can effectively align lip movements with audio, thus attaining greater precision in expression prediction. In the second component, we design a lightweight facial identity alignment (FIA) module which includes a lip-shape control structure and a face texture reference structure. This novel design allows us to generate fine details in real-time, without depending on sophisticated and inefficient feature alignment modules. Our experimental results, both quantitative and qualitative, on public datasets demonstrate the clear advantages of our method in terms of lip-speech synchronization and generation quality. Furthermore, our method is efficient and requires fewer computational resources, making it well-suited to meet the needs of practical applications.

**Keywords:** Audio-driven Face Generation · Real-time · 3D Facial Prior

## 1 Introduction

Audio-driven face generation has received much attention in recent years due to its great potential in real-world applications. The main challenges in generating realistic and expressive talking faces include: 1) Ensuring lip-speech synchronization that matches the audio input and the lip movements. 2) Achieving photo-realistic visual quality that preserves the details and textures of the face. 3)
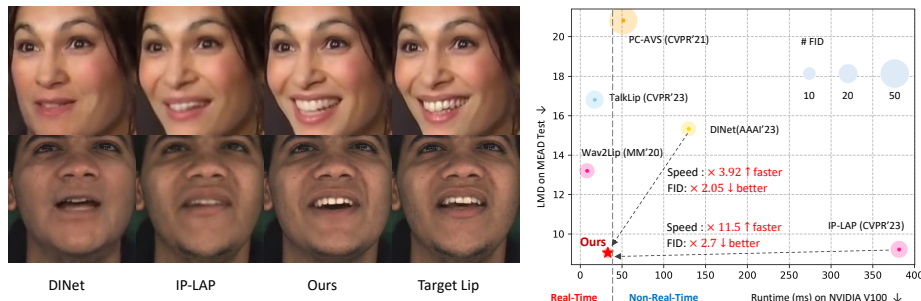
Fig. 1: Left: Visual Comparison on lip sync and generation quality with IP-LAP [38] and DINet [36]. Our method achieves precise lip-synced talking faces, closer to the target lip, with higher visual quality. **Right: Speed, LMD and FID comparisons**. Our method generates talking faces at 30 FPS on NVIDIA V100, showcasing the *best* LMD and FID scores while maintaining the *real-time* speed.

Maintaining identity preservation that keeps the expressions and facial features consistent with the original individual. 4) Enhancing efficiency that enables fast and robust face generation, especially for applications in real-world scenarios.

Existing person-generic talking face generation methods can be roughly divided into two categories: realtime-based and non realtime-based methods as listed in Table 1. In the first group, Wav2Lip [21] employs a sync-expert to improve the lip-synchronization performance with accurate lip motion. Talk-Lip [29] proposes an efficient framework that tackles the reading intelligibility problem by leveraging a lipreading expert. These methods *have fast inference speed but usually suffer from the unsatisfactory generation effects*, *e.g.* blurry faces in Wav2Lip or facial artifacts in TalkLip.

For methods in the second category, PC-AVS [39] incorporates disentangle learning for identity, speech content, and poses in talking face generation. DINet [36] develops a deformation part and an inpainting part for accurate mouth movements and textual details. StyleTalk [18] utilize an implicit style code to control global head and facial movements. IP-LAP [38] utilizes the guidance of prior landmark and appearance information, and proposes a two-stage framework, consisting of an audio-to-landmark generator and a landmark-to-video generation model. IP-LAP *produces better visual results than the realtime methods but is time-consuming*, making it impractical for real-world applications.

To achieve realtime efficiency and high-fidelity talking face effects simultaneously, in this paper we propose a novel framework, termed `RealTalk`, which consists of an audio-to-expression transformer converting input audio into 3D expression coefficients, and an expression-to-face renderer generating high-fidelity talking face from the estimated 3D expression. Specifically, there are three key designs in `RealTalk` to improve the performance and efficiency:

1) *Improved facial prior with cross-modal attention* in audio-to-expression transformer. Previous work [3] observed and discussed that the appearance of a face is influenced by two factors: identity and intra-personal variation (*e.g.*, expression, pose, lighting). Inspired by this observation, we enrich the input facial

**Table 1: Method categories and complexity comparsions**. The primary distinction between `RealTalk` and existing methods lies in the incorporation of cross-modal attention on 3D priors, learnable mask and the FIA module. The methods highlighted in **bold** are capable of real-time performance. Furthermore, our method attains real-time performance via its compact structure and reduced dependency on reference frames.

| Method *vs.* Category | Audio to Face Translation | Facial Mask | Identity Alignment ($n$ frames) | Speed (FPS) |
|---|---|---|---|---|
| **Wav2Lip** [21] | Audio encoder | Lower-half | Concat (1) | 120 |
| **TalkLip** [29] | Audio encoder | Lower-half | Concat (1) | 57 |
| PC-AVS [39] | Audio encoder | Full face | Concat (1) | 17 |
| DINet [36] | Audio encoder | Rectangle | Deformation (5) | 8 |
| StyleTalk [18] | Style decoder | Full face | Flow-based (1) | 7 |
| IP-LAP [38] | landmark transformer | Lower-half | Flow-based (25) | 3 |
| **RealTalk (Ours)** | 3D cross-modal temporal attention | Learnable | FIA module (1) | 30 |

prior by performing cross-modal on the 3D shape and historical expression coefficients as 3D facial prior guidance besides the input audio queries. Here, the shape represents the identity, while expressions from historical frames capture intra-individual lip amplitude variations.

2) *Learnable facial mask* as the bridge connecting the two networks. Different from the previous methods that occlude half of the face or adopt a fixed position black square, our method leverages the learned 3D expressions from the audio-to-expression transformer, and converts them into an adaptive facial mask that better estimates the output facial structure given the input audio, leading to better performance in facial contour generation and lip motion accuracy.

3) *Efficient and effective network design* in expression-to-face renderer. We highlight the advantages of our FIA module in inference speed, which dominates the overall runtime. Unlike recent methods [36,38] that require *time-consuming* feature extraction from multiple reference images to enhance visual quality, our FIA module is meticulously crafted to achieve high-quality texture generation in *real-time* using only 1 image. Specifically, different from the methods that use optical flow [38] or deformation module [36], our method designs a novel Facial Identity Alignment (FIA) module to achieve high-fidelity talking face synthesis.

Overall, our contributions are summarized as follows:

– The proposed `RealTalk` makes full use of the improved 3D facial prior by applying cross-modal attention to shape and variation to help predict more accurate facial expressions.
– The proposed FIA module exhibits strong control over lip movements and texture referencing capabilities, thereby producing high-quality facial images without sacrificing efficiency.
– To our best knowledge, the proposed `RealTalk` is the best choice considering both accuracy and efficiency (*i.e.*, 30FPS) for talking face generation as shown in Fig. 1.

## 2   Related Work

**Audio-driven Talking Face Generation.** Existing audio-driven face generation can be primarily divided into two categories of methods: person-specific and person-generic approaches. Person-specific methods  [9, 16, 17, 24, 28, 32, 33] require training or fine-tuning on specific individuals before inference, whereas person-generic methods [4, 5, 8, 13, 14, 18, 20, 21, 25, 27, 29, 31, 36, 38–40] enable the direct generation of talking face videos for unseen person. To address the challenge of audio-visual synchronization in person-generic methods, Wav2Lip [21] introduces a lip synchronization discriminator using SyncNet [6]. In contrast, TalkLip [29] utilizes a lip reading network to enhance the comprehensibility of the lip region. Furthermore, several approaches [4, 14, 18, 38] focus on modeling the mapping from audio to facial expressions, which simplifies the process of lip synchronization. To enhance the visual quality, DiffTalk [25] employs a diffusion model and StyleSync [8] utilizes StyleGAN [15] to provide high-fidelity facial priors. Additionally, certain methods [14, 36, 38] employ identity reference alignment to preserve facial identity and texture. However, achieving a balance between efficiency, visual quality, and accuracy of lip movements is a formidable challenge for the aforementioned person-generic methods.

**Audio to Facial Expressions Modeling.** Modeling the integration of audio into facial expressions in a general context enables more efficient and accurate learning of lip movements. IP-LAP [38] predicts facial landmarks of the target face based on the audio input, while EAMM [14] utilizes the unsupervised FOMM [26] to extract the keypoints of the target face. Although facial landmarks or keypoints are relatively easy to obtain, they suffer from sparsity, making it challenging to represent complex facial movements adequately, such as the actions of puckering or pursing the lips. To address these limitations, we propose using 3DMM [1] to extract more accurate decoupled information about facial identity, pose, and expression, and learn the mapping from audio to expression coefficients. Compared to 2D facial landmarks, 3DMM provides denser keypoints, allowing for a more comprehensive representation of intricate facial region movements.

**Identity Reference Alignment.** Previous methods commonly employ encoder-decoder architectures to directly fuse reference identity frames, but they often fail to effectively preserve identity features. In contrast, identity reference alignment ensures a stronger resemblance between the generated results and the identity. For example, IP-LAP [38] relies on optical flow [11] to align multiple identity reference features by warping them onto the target frame. DINet [36] employs adaptive affine transformation [35] to process multiple reference frames and generate deformed features that enhance identity information. These methods employ intricate alignment strategies and *multiple* identity references, slowing down the inference speed. In contrast, we propose an efficient facial identity alignment module by utilizing *single* frame of identity reference.
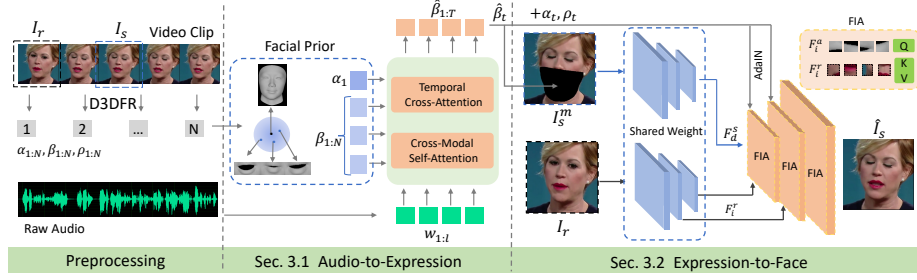
**Fig. 2: Framework of our approach.** Our network is divided into two parts: *Audio-to-expression Transformer*, and *Expression-to-face Renderer*. The preprocessing is to extract 3D shapes $\alpha_{1:N}$, expressions $\beta_{1:N}$, poses $\rho_{1:N}$, and audio feature $w_{1:l}$. In the first part, the shape $\alpha_1$ and historical expressions $\beta_{1:N}$ are utilized as **Improved Facial Prior** to predict $\hat{\beta}_{1:N}$ while preserving identity and intra-personal lip amplitude variations. In the second part, the predicted expressions are injected into the proposed **Facial Identity Alignment (FIA)** module to inpaint the masked source frame $I_s^m$ the target lip through cross-attention with the identity reference $I_r$.

## 3   Method

**Overview.** Our goal is to generate a video that synchronizes the lip movements with a target audio clip while maintaining the consistency of the facial identity from the original video. Fig. 2 illustrates our method, which consists of two stages. In the first stage, we utilize shape and historical expressions as conditions to map the audio to 3D expression coefficients with the proposed audio-to-expression transformer. In the second stage, we design an lightweight face renderer including a facial identity alignment module to generate the target lip based on the predicted expression coefficients and reference frame.

### 3.1   Audio-to-expression Transformer

In this stage, our objective is to generate precise and stylistically consistent lip movements. Leveraging 3D face reconstruction technology, we can effectively control dense facial regions with a reduced number of coefficients. Given a series of facial images, we use the D3DFR model [7] to extract 3D coefficients. These coefficients consist of three components: shape $\alpha$, expression $\beta$, and pose $\rho$. As shown in Fig. 3, the audio-to-expression network takes three inputs, including driving audio features, the 3D shape feature, and historical expressions.

**Shape and Expression Prior.** The uniqueness of each individual's facial and mouth structures results in variations in how audio and lip movements align. Here, we enrich two types of personalized facial priors, *i.e.* shape and historical expressions.

*Shape* signifies identity, which is typically related to the natural face size and mouth proportions. On the other hand, *historical expressions* capture the
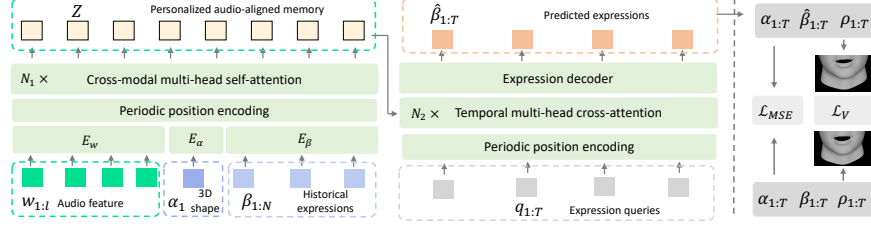
**Fig. 3: Architecture of the Audio-to-expression Transformer.** (Left) The audio, shape, and historical expressions are processed through an encoder to obtain memory features, which are then combined with expression queries in decoder to generate predictions. (Right) The predicted expressions and GT expressions are optimized using reconstruction and vertex losses.

individual's unique lip amplitude, which supply the individual variations from the standard shape.

Firstly, we designate the first frame to obtain the default shape coefficient, denoted as $\alpha_1$ and $N$ frames of historical expressions, represented as $\beta_1, \cdots, \beta_N$. We utilize Hubert [12] to extract audio features. The audio features can be represented as $w_1, \cdots, w_l$, where $l$ denotes the length of audio feature. Specifically, the shape, expression coefficients and audio features are passed through fully connected networks to obtain embeddings, respectively. These embeddings are then concatenated in sequence order, resulting in a total of $l + N + 1$ tokens, which are input into a periodic position encoding layer and $N_1$ cascaded Cross-Modal multi-head Self-Attention (CMSA) to obtain the personalized audio-aligned memory:

$$Z = \mathrm{CMSA}(w_1, \cdots, w_l, \alpha_1, \beta_1, \cdots, \beta_N). \tag{1}$$

The expression queries $q_t$ (initially set to zero), are combined with the memory and fed into $N_2$ cascaded Temporal multi-head Cross-Attention (TCA) network and the linear decoder to get expression prediction:

$$\hat{\beta}_1, \hat{\beta}_2, \cdots, \hat{\beta}_T = \mathrm{TCA}(q_1, q_2, \cdots, q_T, Z). \tag{2}$$

**Loss Function.** The loss function primarily consists of Mean Squared Error (MSE) and 3D vertex loss. MSE calculates the error between the predicted expression coefficients and the Ground Truth (GT):

$$\mathcal{L}_{MSE} = \frac{1}{T}\Sigma_{t=1}^{T}\|\beta_t - \hat{\beta}_t\|_2^2. \tag{3}$$

The vertex loss calculates 3D vertices by combining the predicted expression coefficients $\hat{\beta}$ with shape $\alpha_1$ and pose $\rho$, and chooses points from the mouth region to evaluate the distance. These loss components optimize the process by minimizing the gap between the predicted and GT expression coefficients and ensure accurate alignment of the generated 3D vertices with the mouth keypoints.
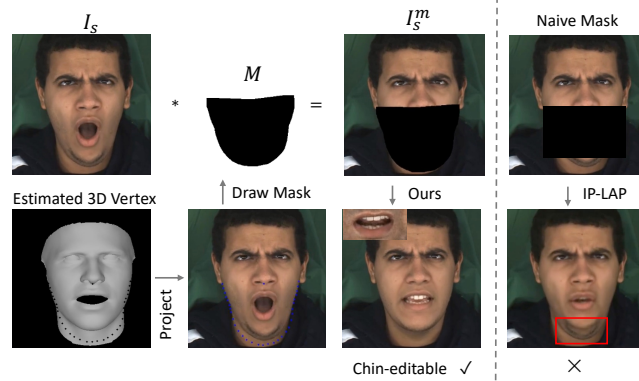
**Fig. 4: Illustration of the learnable mask** based on predicted facial expressions. 1) The estimated 3D vertex allows us to select points with fixed positions relative to the face. We choose points that emphasizes a larger facial contour to accommodate diverse lip movements. 2) Comparisons between our learnable mask and the naive mask (IP-LAP), with the left-top in our result representing the target lip, show that our method effectively adjusts the face shape based on the spoken content, while IP-LAP yields unnatural results, *e.g.* a double chin.

Let V represent the 3D vertex computed from the coefficients. The vertex loss can be described as follows:

$$\mathcal{L}_V = \frac{1}{T}\Sigma_{t=1}^{T}\| \, \mathrm{V}(\alpha_t, \beta_t, \rho_t) - \mathrm{V}(\alpha_t, \hat{\beta}_t, \rho_t)\|_2^2. \tag{4}$$

The overall loss is $\mathcal{L}_{a2e} = \mathcal{L}_{MSE} + 0.1 * \mathcal{L}_V$.

### 3.2 Expression-to-face Renderer

We design a lightweight network for generating facial images with edited lips. The face renderer takes a masked source image $I_s^m$, a reference image $I_r$, and 3D coefficients $\{\alpha_t, \hat{\beta}_t, \rho_t\}$ as inputs. Our network adopts an encoder-decoder architecture, where both images are processed via a shared-weight encoder to extract multi-scale features, subsequently integrated with 3D coefficients within the decoder. Next, we describe the detailed process below.

**Learnable Mask.** Unlike previous methods, as shown in Fig. 4, we selectively mask the source image $I_s$ allowing for more accurate control over the modified areas. The 3D vertices are estimated according to the predicted expression co-efficients and then projected onto the image. Based on predetermined sectional points, we draw and fill the mouth and neck regions. The process of generating the mask is as follows:

$$\begin{aligned} V_{xy} &= \mathrm{P}(\mathrm{V}(\alpha_t, \hat{\beta}_t, \rho_t), \tau_t), \\ M &= \mathrm{C}(V_{xy}), (x, y) \in S, \\ I_s^m &= M * I_s, \end{aligned} \tag{5}$$
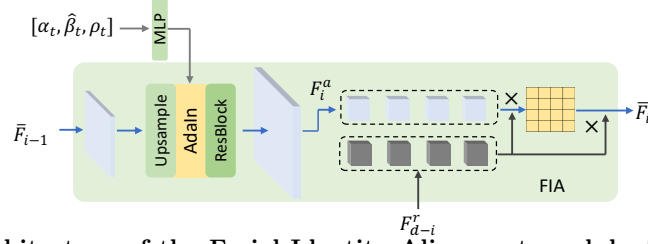
**Fig. 5: Architecture of the Facial Identity Alignment module.** The inputs to the FIA module include the predicted facial expression coefficients $\hat{\beta}_t$, the known shape $\alpha_t$ and pose $\rho_t$, the feature from last module $\bar{F}_{i-1}$, and the identity reference image feature $F_{d-i}^r$. The 3D coefficients are injected through AdaIN, enabling control over the lip. Multi-scale reference features interact with the current features through cross-attention, facilitating effective texture transfer.

where P is the project function and $\tau_t$ indicates the translation matrix of $t$. C is the contour of points in the face and neck area of interest, and the set of these points is $S$. $I_s^m$ represents the image operated with the learned mask $M$.

**Shared Encoder.** The encoder consists of $d$ stages of stacked residual blocks, where the resolution is reduced by half and the feature dimension is increased at each stage. Since the encoder weights are shared, we parallelize the computation of the source and reference along the batch dimension. The encoder extracts multi-scale features as:

$$F_i^s, F_i^r = \mathrm{SE}_i([F_{i-1}^s, F_{i-1}^r]),\tag{6}$$

where $\mathrm{SE}_i$ represents $i$-th level of the shared encoder. Then we get $[F_1^s, \cdots, F_d^s]$ and $[F_1^r, \cdots, F_d^r]$.

**Facial Identity Alignment Module (FIA).** In decoder, each layer's input features $\bar{F}_{i-1}$, output from last stage (initially $F_d^s$), are first upsampled to increase the resolution, denoted as $F_i^u$. As depicted in Fig. 5, the 3D coefficients $\{\alpha_t, \hat{\beta}_t, \rho_t\}$ undergo dimension mapping through a three-layer MLP before being injected into the network through an AdaIN module. The feature is modulated by the injected 3D coefficients to match the specific lip shapes. We then incorporate multiple additional residual blocks (2 in our final model) to further enhance these features. Finally, features from the reference image $F_{d-i}^r$ are aligned at the same resolution to control face generation and aggregate textures:

$$\bar{F}_i = \mathrm{FIA}_i([\alpha_t, \hat{\beta}_t, \rho_t], \bar{F}_{i-1}, F_{d-i}^r).\tag{7}$$

To avoid generating unnecessary background, we adopt a blending strategy by merging the results of the final layer with the input to obtain the final outcome:

$$\hat{I}_s = M * I_s + (1 - M) * \bar{F}_d.\tag{8}$$

**Loss Function.** We employ multiple loss functions to constrain the accuracy of lip-sync and visual quality, including pixel loss, perceptual loss, adversarial loss,

and local pixel loss to enhance the details of teeth:

$$\begin{aligned}
\mathcal{L}_1 &= \Sigma \|\hat{I}_s - I_s\|_1, \\
\mathcal{L}_2 &= \Sigma \| \text{VGG}(\hat{I}_s) - \text{VGG}(I_s) \|_1, \\
\mathcal{L}_3 &= \mathbb{E}_{I_s}[\log D(I_s) + \log(1 - D(\hat{I}_s))], \\
\mathcal{L}_4 &= \Sigma \| M' * \hat{I}_s - M' * I_s \|_1,
\end{aligned} \tag{9}$$

where $D(\cdot)$ is the discriminator and $M'$ represents the binary mask of the teeth area. The overall loss function is obtained by weighting the above losses according to their respective weights.

$$\mathcal{L}_{e2f} = \lambda_1 * \mathcal{L}_1 + \lambda_2 * \mathcal{L}_2 + \lambda_3 * \mathcal{L}_3 + \lambda_4 * \mathcal{L}_4, \tag{10}$$

where $\lambda_1 = 1, \lambda_2 = 1, \lambda_3 = 0.1, \lambda_4 = 1$.

**Discussion with IP-LAP and DINet on Efficient Design.** To reduce computational burden, we employ several optimization strategies. Firstly, we utilize highly compressed 3D coefficients as the context for audio-to-face conversion. This approach is computationally more efficient than IP-LAP [38]'s use of 2D landmarks images, thereby circumventing the need to process high-dimensional features. Secondly, we implement a shared encoder to concurrently extract multi-scale features from both masked source and unmasked reference images. Thirdly, we employ only 1 frame for texture transfer, in contrast to the 5 and 25 frames used by DINet [36] and IP-LAP respectively. For example, IP-LAP increases the computational load of the alignment module by warping each reference frame to the current image via optical flow. Conversely, DINet accomplishes mouth inpainting by extracting deformation features from reference frames. Lastly, our FIA module is designed for efficiency. Its cross-attention mechanism can adaptively query similar texture features and can be flexibly embedded at various scales within the network. It should be noted that we only perform cross-attention on resolutions of $\frac{1}{8}$ and $\frac{1}{16}$.

## 4 Experiments

### 4.1 Experiments Setting

**Implementation Details.** In our experiments, $N = 16$, $T = 16$, $l = 32$. The $N$ historical expressions are randomly selected during training and inference. We train the expression-to-face renderer at $256 \times 256$ resolution. The identity reference is set to the first frame of a video clip. The number of encoder and decoder stages $d$ is 4, and each stage has 2 stacked residual blocks. More details are included in the supplementary material.

**Datasets.** We conducted experiments on three popular datasets: VoxCeleb1 [19], MEAD [30], HDTF [37]. VoxCeleb1 comprises over $100,000$ utterances from $1,251$ celebrities, extracted from videos uploaded to YouTube. We utilized these utterances that are available from Internet (about 10% of total) for training and randomly selected 50 utterances for evaluation. MEAD is a talking-face video

**Table 2: Quantitative results with SOTA methods on benchmark datasets.** '↑' and '↓' mean higher and lower are desired. The $\text{Sync}_{conf}$ are marked in gray for its weak reflection of audio-visual synchronization. The runtime is evaluated on V100.

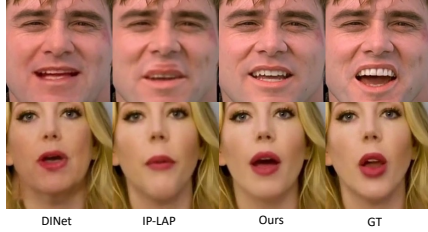| Method | Dataset | Reconstruction | | | | | | | Cross Audio | | | Runtime (ms) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | LMD↓ | M-LMD↓ | F-LMD↓ | FID↓ | LPIPS↓ | SSIM↑ | $\text{Sync}_{conf}$ | FID↓ | CSIM↑ | $\text{Sync}_{conf}$ | |
| GT | | - | - | - | - | - | - | 6.54 | - | - | 6.54 | - |
| Wav2Lip [21] | | 8.78 | 15.87 | 5.82 | 17.58 | 0.1097 | 0.9351 | 7.23 | 20.87 | 0.9329 | 6.73 | 8.3 |
| PC-AVS [39] | | 23.90 | 22.08 | 24.66 | 65.21 | 0.3281 | 0.6960 | 7.17 | 69.41 | 0.7621 | 7.20 | 51.7 |
| TalkLip [29] | VoxCeleb1 [19] | 19.29 | 31.83 | 14.06 | 34.70 | 0.1557 | 0.8987 | 6.92 | 22.27 | 0.9238 | 4.90 | 17.4 |
| DINet [36] | | 18.25 | 24.57 | 15.61 | 23.83 | 0.1235 | 0.9091 | 5.52 | 27.56 | 0.8385 | 4.66 | 129.9 |
| IP-LAP [38] | | 8.69 | 14.44 | 6.29 | 16.84 | 0.1196 | 0.9279 | 6.05 | 23.81 | 0.9287 | 4.20 | 381.5 |
| Ours | | **6.72** | **11.02** | **4.92** | **12.73** | **0.0916** | **0.9361** | 6.21 | **17.52** | **0.9434** | 5.00 | 33.1 |
| GT | | - | - | - | - | - | - | 4.65 | - | - | 4.65 | - |
| Wav2Lip [21] | | 13.20 | 29.02 | 6.61 | 24.97 | 0.1346 | 0.9219 | 6.87 | 23.68 | 0.9307 | 6.73 | 8.3 |
| PC-AVS [39] | | 20.82 | 26.73 | 18.35 | 86.02 | 0.3457 | 0.7553 | 7.26 | 90.81 | 0.7550 | 7.69 | 51.7 |
| TalkLip [29] | MEAD [30] | 16.80 | 34.10 | 9.59 | 35.64 | 0.1622 | 0.9073 | 6.75 | 29.34 | 0.9316 | 4.94 | 17.4 |
| DINet [36] | | 15.33 | 38.14 | 5.82 | 23.90 | 0.1131 | 0.9236 | 5.14 | 24.16 | 0.8099 | 4.70 | 129.9 |
| IP-LAP [38] | | 9.22 | 18.68 | 5.27 | 31.57 | 0.1441 | **0.9285** | 5.77 | 36.68 | 0.9472 | 4.01 | 381.5 |
| Ours | | **9.04** | **18.65** | **5.02** | **11.68** | **0.0958** | 0.9251 | 4.00 | **13.22** | **0.9638** | 3.84 | 33.1 |
| DINet [36] | | 8.0255 | 15.36 | 5.185 | 12.94 | 0.0975 | 0.9327 | - | - | - | - | 129.9 |
| IP-LAP [38] | HDTF [37] | 6.076 | 10.20 | 4.658 | 9.490 | 0.1101 | 0.9416 | - | - | - | - | 381.5 |
| Ours | | **6.011** | **9.966** | **4.207** | **6.065** | **0.0820** | **0.9418** | - | - | - | - | 33.1 |



**Fig. 6: Visual comparisons with state-of-the-art competitors.** Our method achieves the best lip-speech sync and visual quality.

corpus that features 60 actors and actresses expressing 8 different emotions at 3 distinct intensity levels. We exclusively use the front view videos, selecting 40 actors for training and 3 actors for evaluation. 20 video clips in HDTF testset are used only for evaluation under reconstruction setting.

**Evaluation Metrics.** We employ facial Landmarks Distance (LMD) [2] to assess the accuracy of lip-sync, M- (Mouth) and F- (Face) separately for better evaluation, FID (Fréchet Inception Distance) [10], LPIPS (Learned Perceptual Image Patch Similarity) [34], and SSIM (Structural Similarity Index Measure)

**Table 3:** Mean Opinion Score (MOS) on benchmark datasets.

| MOS / Method | Wav2Lip | DINet | TalkLip | IP-LAP | Ours |
|---|---|---|---|---|---|
| Visual Quality | 1.53 | 1.72 | 1.55 | 2.84 | **3.77** (33%↑) |
| Lip Sync | 2.15 | 2.50 | 2.06 | 2.58 | **3.72** (44%↑) |

**Table 4:** Metric and runtime comparison with StyleTalk.

| Metric | M-LMD ↓ | F-LMD ↓ | FID ↓ | LPIPS ↓ | SSIM ↑ | Runtime (ms) ↓ |
|---|---|---|---|---|---|---|
| StyleTalk [18] | 7.910 | 5.220 | 15.42 | 0.1486 | 0.8016 | 141.0 |
| Ours | **4.387** | **2.215** | **12.25** | **0.0908** | **0.9093** | **33.1** |



DINet          IP-LAP          Ours          GT

**Fig. 7:** Visual results on HDTF. Please zoom in for more detail.



StyleTalk          Ours          GT

**Fig. 8:** Visual comparison with StyleTalk on its official demos.

to evaluate the quality of generated images. We also assess the identity preservation of generated faces by measuring CSIM (cosine similarity between identity features) [23]. LMD assesses lip-audio sync accuracy among different methods, with a lower value indicating closer alignment with GT and thus better synchronization with the audio. A lower FID signifies that the generated image quality more closely resembles the original video, reflecting image clarity and naturalness. Moreover, a higher CSIM indicates higher facial similarity, suggesting superior identity preservation by the corresponding method. Additionally, we employ the SyncNet [6] metric to evaluate audio-visual consistency. However, it is crucial to emphasize that a higher SyncNet score doesn't necessarily indicate better audio-visual sync, as discussed in [8].

### 4.2 Comparison with SOTA Methods

**Comparison Methods.** We compare the proposed method with state-of-the-art person-generic audio-driven face generation methods, including Wav2Lip [21], PC-AVS [39], DINet [36], TalkLip [29], and IP-LAP [38]. Among these methods, Wav2Lip excels in reconstruction performance. PC-AVS stands out for editing lip motions and poses. DINet employs deformable convolution to construct its reconstruction network. TalkLip, structurally similar to Wav2Lip, introduces a new lip-reading loss function. IP-LAP utilizes 2D landmarks as intermediate information, serving as our primary comparison method.

**Reconstruction.** In this setting, the video clip is driven by the corresponding original audio. Since the video clip illustrates the target lip motions, it serves as the ground truth for calculating metrics requiring paired data.

As shown in Table 2, our method achieves the best results among most metrics on all three datasets. Our method exhibits a significant advantage in the FID metric, surpassing the second-best method by 51% and 36% on MEAD and HDTF, respectively. IP-LAP, employing multiple reference frames (*i.e.*, 25 frames), achieves comparable results on LMD but severely impacts efficiency (*i.e.*, 10× slower than `RealTalk`). Although Wav2Lip achieves noteworthy out-

comes on metrics, our examination reveals its inadequate visual quality, as corroborated by the user study results in Table 3.

The top two rows of Fig. 6 and 7 display the qualitative results under the reconstruction setting, with the rightmost column representing the ground truth. By incorporating the improved facial prior, our proposed audio-to-expression transformer can precisely predict lip shapes according to the individual's movement amplitude, resulting in outcomes closer to the GT. Furthermore, our method captures similar textures resembling the original face with a single image, demonstrating the efficacy and efficiency of our FIA module. Concerning the Sync-Net metric, although both Wav2Lip and PC-AVS scored exceptionally high, the audio-visual synchronization did not improve correspondingly, creating inconsistency with the visual effects.

**Dubbing with Cross Audio.** In this setting, the video clip is driven by another audio segment, providing a more representative depiction of real-world scenarios. Our method achieves the best FID on both datasets, indicating more natural and realistic results in cross-audio settings. The bottom two rows of Fig. 6 depict qualitative results under cross-audio testing. The rightmost column displays the lip motion from the video aligned to the cross audio, serving as a pseudo GT for the accurate lip shape. Our results exhibit closer lip shapes to the pseudo GT, indicating better lip-speech sync against the SOTA competitors.

We further conduct a user study to evaluate the generation quality and lip synchronization of different methods. 10 videos were selected, and scores were collected from 15 participants, ranging from 1 (worst) to 5 (best). As shown in Table 3, our method excels in both generation quality and lip synchronization, outperforming the second-best IP-LAP by a significant margin (*i.e.*, 33% and 44% improvements).

**Runtime Analysis.** Our approach outperforms the SOTA methods, speeding up $3.92\times$ and $11.5\times$ faster than DINet and IP-LAP, respectively. Specifically, the runtime of the audio-to-expression transformer and the expression-to-face renderer are 2.3ms and 30.8ms, respectively. By utilizing 3D priors, our method facilitates the precise expression generation with reduced computational requirements. Our FIA module efficiently executes reference texture transfer, thereby eliminating the necessity for multi-encoder and multi-reference feature alignment processes. Moreover, our method offers a distinct quality advantage over the real-time methods (e.g., TalkLip, Wav2Lip). As evidenced by the reconstruction metrics in Table 2, despite its rapid processing speed, TalkLip's reconstruction capability is subpar. Visual inspection reveals that Wav2Lip, the fastest method, presents significant artifacts in cross-audio scenarios, leading to a decrease in generalization performance. Overall, our method constitutes the optimal solution for achieving a balance between effectiveness and efficiency.

**Comparison with One-shot Method.** We emphasize the main differentials compared with one-shot talking head methods, *e.g.*, StyleTalk [18]. In `RealTalk`, we explicitly introduce historical expressions and focus on local lip movements by vertex loss in Equation (4). In StyleTalk, the style is implicit and is used to control global head and facial movements. As a lip movement specialist,

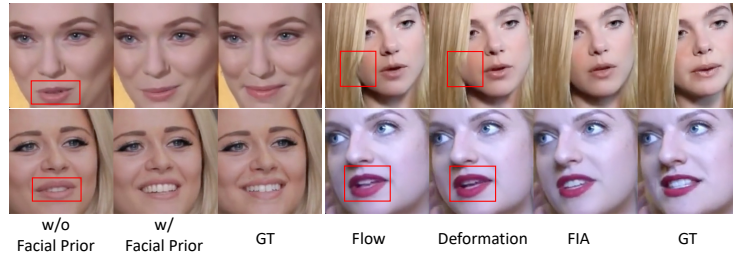**Table 5:** Ablation study on different facial priors.

| Facial | Shape | × | ✓ | × | ✓ |
|---|---|---|---|---|---|
| Prior | Historical Expression | × | × | ✓ | ✓ |
| Expression Error | | 0.2680 | 0.1342 | 0.1186 | **0.1128** |

**Table 6:** Reconstructions with different masks on VoxCeleb1.

| Mask | LMD↓ | FID ↓ | LPIPS ↓ | SSIM ↑ |
|---|---|---|---|---|
| Naive | 7.08 | 13.31 | 0.1133 | 0.9088 |
| Learnable | **6.72** | **12.73** | **0.0916** | **0.9361** |

**Table 7:** Reconstruction performance of FIA with different reference module and across different scales on VoxCeleb1. The runtime only reflects the face renderer.

| Configuration | FID↓ | LPIPS↓ | SSIM↑ | Runtime (ms) | Param. (M) |
|---|---|---|---|---|---|
| Flow | 13.68 | 0.0963 | 0.9329 | 30.48 | 82.94 |
| Deformation | 13.38 | 0.0948 | 0.9332 | 31.20 | 98.79 |
| Blocks=1 | 13.42 | 0.0959 | 0.9326 | 24.28 | 52.53 |
| Blocks=3 | **12.11** | 0.0924 | 0.9352 | 38.65 | 85.96 |
| Final (FIA) | 12.73 | **0.0916** | **0.9361** | 30.82 | 69.24 |



| w/o Facial Prior | w/ Facial Prior | GT | Flow | Deformation | FIA | GT |

**Fig. 9: Left**: Effects on with or without improved facial priors. **Right**: Effects on FIA paired with different reference module.

`RealTalk` exhibits superior lip-synchronization than StyleTalk in comparison on their official demos, as shown in Tab. 4 and Fig. 8. In summary, `RealTalk` outperforms in all 5 metrics and is 4.27× faster than StyleTalk, which adopts PIRender [22] in the second stage.

### 4.3 Ablation Study

**Effectiveness of the Improved Facial Prior.** To validate the effectiveness of our improved facial priors, we separately removed the shape and historical expressions, and evaluate on VoxCeleb1 dataset under reconstruction setting. Since the GT expression coefficients are known, we directly quantified the mean squared error between the predicted expression coefficients of different models and the GT coefficients. In Table 5, introducing both shape and historical expressions positively impacts lip synchronization prediction. Compared to the model without shape and expression prior, the full model improves prediction accuracy by 57.9%. The visual results depicted in Fig. 9-Left demonstrate the efficacy of the enhanced facial prior in maintaining intra-personal expressions.

**Effectiveness of the Learnable Mask.** To validate the benefits of our proposed learnable mask in talking face generation, we replace it with a naive mask obscuring the lower half of the image and conduct experiment on VoxCeleb1 dataset under reconstruction setting. As shown in Table 6, the performance drops

when using the naive mask. The naive mask lacks information about the target face shape and includes irrelevant background in the area that the network has to generate, posing increased learning difficulty. Conversely, the learnable mask is intrinsically associated with the target audio, remaining impervious to the original facial contour, thereby guaranteeing enhanced accuracy of lip movements and the naturalness of facial expressions.

**Comparison with Deformation and Flow-based Module.** To validate our FIA module, we replace its cross-attention component with other common alignment modules: flow-based warping and deformation convolution. Table 7 shows the impact of building FIA with different modules on generated image quality, along with comparisons in terms of runtime and parameters. FIA (Final in Table 7) achieves superior visual quality with fewer parameters over the deformation and flow-based structure. Fig. 9-Right illustrates the proposed FIA paired with cross-attention, highlighting its strong ability to restore textures, such as hair and teeth. Unlike flow-based and deformable convolution methods overly relying on the reference, cross-attention allows for a weighted fusion of features from different regions, enabling a more flexible generation of facial textures.

**Complexity.** Table 7 also displays the impact of adjusting the number of residual blocks within the FIA module on generated image quality, assessing the performance of `RealTalk` across different scales. Configurations with 1, 2 (final model), and 3 residual blocks are explored. Setting num=1 accelerates speed but diminishes visual quality, while num=3 improves FID results at the expense of real-time performance. Consequently, num=2 strikes a good balance between visual quality and speed.

## 5   Conclusion

We propose a novel audio-driven framework `RealTalk`, incorporating an audio-to-expression transformer and a high-fidelity expression-to-face renderer. Our improved facial prior adeptly adjusts speech content while maintaining identity through a cross-modal attention on both identity and intra-persona variation features. A specialized learnable mask tackles challenges associated with altering facial structures. Our FIA module, combining AdaIN and cross-attention structures, facilitates precise lip-shape control using 3D coefficients and efficient facial texture transfer from a single frame. Experimental results on benchmarks affirm our method's superiority in lip-speech sync and generation quality, emphasizing its efficiency and applicability.

**Limitation and Social Impacts.** Our approach encounters limitations with facial obstructions, such as microphones or hand movements in front of the face. This is expected, given our primary focus on facial generation, particularly in modeling mouth shapes, without an additional segmentation model to predict facial obstructions. Efficient talking face technology can be used for digital human live streaming and interaction, but it also carries risks in illicit industries, including the manipulation of spoken content for deceptive purposes. To prevent misuse, generated videos should be clearly marked. Ongoing research should also be dedicated to identifying AI-generated videos, evolving alongside advancements in generative models.

# References

1. Blanz, V., Vetter, T.: A morphable model for the synthesis of 3d faces. In: Proceedings of the 26th annual conference on Computer graphics and interactive techniques. pp. 187–194 (1999) 4
2. Bulat, A., Tzimiropoulos, G.: 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) 10
3. Chen, D., Cao, X., Wang, L., Wen, F., Sun, J.: Bayesian face revisited: A joint formulation. In: European Conference on Computer Vision (2012) 2
4. Chen, L., Maddox, R.K., Duan, Z., Xu, C.: Hierarchical cross-modal talking face generation with dynamic pixel-wise loss. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 7832–7841 (2019) 4
5. Cheng, K., Cun, X., Zhang, Y., Xia, M., Yin, F., Zhu, M., Wang, X., Wang, J., Wang, N.: Videoretalking: Audio-based lip synchronization for talking head video editing in the wild. In: SIGGRAPH Asia 2022 Conference Papers. pp. 1–9 (2022) 4
6. Chung, J.S., Zisserman, A.: Out of time: automated lip sync in the wild. In: Computer Vision–ACCV 2016 Workshops: ACCV 2016 International Workshops, Taipei, Taiwan, November 20-24, 2016, Revised Selected Papers, Part II 13. pp. 251–263. Springer (2017) 4, 11
7. Deng, Y., Yang, J., Xu, S., Chen, D., Jia, Y., Tong, X.: Accurate 3d face reconstruction with weakly-supervised learning: From single image to image set. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops. pp. 0–0 (2019) 5
8. Guan, J., Zhang, Z., Zhou, H., Hu, T., Wang, K., He, D., Feng, H., Liu, J., Ding, E., Liu, Z., et al.: 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. pp. 1505–1515 (2023) 4, 11
9. Guo, Y., Chen, K., Liang, S., Liu, Y.J., Bao, H., Zhang, J.: Ad-nerf: Audio driven neural radiance fields for talking head synthesis. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5784–5794 (2021) 4
10. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: Gans trained by a two time-scale update rule converge to a local nash equilibrium. Advances in neural information processing systems **30** (2017) 10
11. Horn, B.K., Schunck, B.G.: Determining optical flow. Artificial intelligence **17**(1-3), 185–203 (1981) 4
12. Hsu, W.N., Bolte, B., Tsai, Y.H.H., Lakhotia, K., Salakhutdinov, R., Mohamed, A.: 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) 6
13. Hu, X., Ren, W., LaMaster, J., Cao, X., Li, X., Li, Z., Menze, B., Liu, W.: Face super-resolution guided by 3d facial priors. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IV 16. pp. 763–780. Springer (2020) 4
14. Ji, X., Zhou, H., Wang, K., Wu, Q., Wu, W., Xu, F., Cao, X.: Eamm: One-shot emotional talking face via audio-based emotion-aware motion model. In: ACM SIGGRAPH 2022 Conference Proceedings. pp. 1–10 (2022) 4
15. Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., Aila, T.: Analyzing and improving the image quality of stylegan. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 8110–8119 (2020) 4

16. Liu, X., Xu, Y., Wu, Q., Zhou, H., Wu, W., Zhou, B.: Semantic-aware implicit neural audio-driven video portrait generation. In: European Conference on Computer Vision. pp. 106–125. Springer (2022) 4

17. Lu, Y., Chai, J., Cao, X.: Live Speech Portraits: Real-time photorealistic talking-head animation. ACM Transactions on Graphics **40**(6) (December 2021). https://doi.org/10.1145/3478513.3480484 4

18. Ma, Y., Wang, S., Hu, Z., Fan, C., Lv, T., Ding, Y., Deng, Z., Yu, X.: Styletalk: One-shot talking head generation with controllable speaking styles. arXiv preprint arXiv:2301.01081 (2023) 2, 3, 4, 11, 12

19. Nagrani, A., Chung, J.S., Zisserman, A.: Voxceleb: a large-scale speaker identification dataset. In: INTERSPEECH (2017) 9, 10

20. Park, S.J., Kim, M., Hong, J., Choi, J., Ro, Y.M.: Synctalkface: Talking face generation with precise lip-syncing via audio-lip memory. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 36, pp. 2062–2070 (2022) 4

21. Prajwal, K., Mukhopadhyay, R., Namboodiri, V.P., Jawahar, C.: A lip sync expert is all you need for speech to lip generation in the wild. In: Proceedings of the 28th ACM international conference on multimedia. pp. 484–492 (2020) 2, 3, 4, 10, 11

22. Ren, Y., Li, G., Chen, Y., Li, T.H., Liu, S.: Pirenderer: Controllable portrait image generation via semantic neural rendering. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 13759–13768 (2021) 13

23. Schroff, F., Kalenichenko, D., Philbin, J.: Facenet: A unified embedding for face recognition and clustering. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 815–823 (2015) 11

24. Shen, S., Li, W., Zhu, Z., Duan, Y., Zhou, J., Lu, J.: Learning dynamic facial radiance fields for few-shot talking head synthesis. In: European Conference on Computer Vision. pp. 666–682. Springer (2022) 4

25. Shen, S., Zhao, W., Meng, Z., Li, W., Zhu, Z., Zhou, J., Lu, J.: Difftalk: Crafting diffusion models for generalized audio-driven portraits animation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1982–1991 (2023) 4

26. Siarohin, A., Lathuilière, S., Tulyakov, S., Ricci, E., Sebe, N.: First order motion model for image animation. Advances in neural information processing systems **32** (2019) 4

27. Sun, Y., Zhou, H., Wang, K., Wu, Q., Hong, Z., Liu, J., Ding, E., Wang, J., Liu, Z., Hideki, K.: Masked lip-sync prediction by audio-visual contextual exploitation in transformers. In: SIGGRAPH Asia 2022 Conference Papers. pp. 1–9 (2022) 4

28. Thies, J., Elgharib, M., Tewari, A., Theobalt, C., Nießner, M.: 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. pp. 716–731. Springer (2020) 4

29. Wang, J., Qian, X., Zhang, M., Tan, R.T., Li, H.: Seeing what you said: Talking face generation guided by a lip reading expert. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 14653–14662 (2023) 2, 3, 4, 10, 11

30. Wang, K., Wu, Q., Song, L., Yang, Z., Wu, W., Qian, C., He, R., Qiao, Y., Loy, C.C.: Mead: A large-scale audio-visual dataset for emotional talking-face generation. In: ECCV (Augest 2020) 9, 10

31. Xu, C., Zhu, J., Zhang, J., Han, Y., Chu, W., Tai, Y., Wang, C., Xie, Z., Liu, Y.: High-fidelity generalized emotional talking face generation with multi-modal emotion space learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2023) 4

32. Ye, Z., Jiang, Z., Ren, Y., Liu, J., He, J., Zhao, Z.: Geneface: Generalized and high-fidelity audio-driven 3d talking face synthesis. arXiv preprint arXiv:2301.13430 (2023) 4

33. Zhang, C., Zhao, Y., Huang, Y., Zeng, M., Ni, S., Budagavi, M., Guo, X.: Facial: Synthesizing dynamic talking face with implicit attribute learning. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 3867–3876 (2021) 4

34. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 586–595 (2018) 10

35. Zhang, Z., Ding, Y.: Adaptive affine transformation: A simple and effective operation for spatial misaligned image generation. In: Proceedings of the 30th ACM International Conference on Multimedia. pp. 1167–1176 (2022) 4

36. Zhang, Z., Hu, Z., Deng, W., Fan, C., Lv, T., Ding, Y.: Dinet: Deformation inpainting network for realistic face visually dubbing on high resolution video. arXiv preprint arXiv:2303.03988 (2023) 2, 3, 4, 9, 10, 11

37. Zhang, Z., Li, L., Ding, Y., Fan, C.: Flow-guided one-shot talking face generation with a high-resolution audio-visual dataset. IEEE (2021) 9, 10

38. Zhong, W., Fang, C., Cai, Y., Wei, P., Zhao, G., Lin, L., Li, G.: Identity-preserving talking face generation with landmark and appearance priors. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 9729–9738 (2023) 2, 3, 4, 9, 10, 11

39. Zhou, H., Sun, Y., Wu, W., Loy, C.C., Wang, X., Liu, Z.: Pose-controllable talking face generation by implicitly modularized audio-visual representation. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 4176–4186 (2021) 2, 3, 4, 10, 11

40. Zhu, F., Zhu, J., Chu, W., Tai, Y., Xie, Z., Huang, X.H., Wang, C.: Hifihead: One-shot high fidelity neural head synthesis with 3d control. In: International Joint Conference on Artificial Intelligence (2022) 4