

# RADIO: Reference-Agnostic Dubbing Video Synthesis

Dongyeun Lee<sup>1,\*</sup> Chaewon Kim<sup>1,\*</sup> Sangjoon Yu<sup>1</sup> Jaejun Yoo<sup>2,†</sup> Gyeong-Moon Park<sup>3,†</sup>  
<sup>1</sup>Klleon AI Research <sup>2</sup>UNIST <sup>3</sup>Kyung Hee University

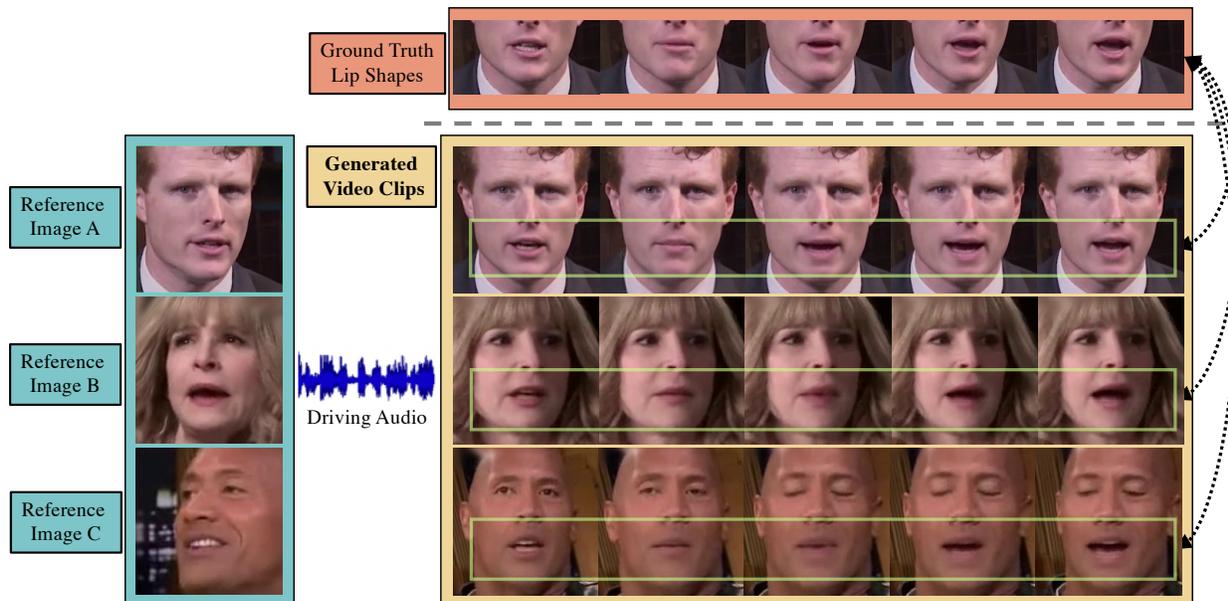


Figure 1. **Illustration of results generated by our RADIO framework.** Our method targets a one-shot audio-driven talking face generation, where synchronized mouth shapes are generated while holding on to the identity of a single reference frame. Even with diverse poses and expressions of reference frames, our method generates accurately synced lips robustly.

## Abstract

One of the most challenging problems in audio-driven talking head generation is achieving high-fidelity detail while ensuring precise synchronization. Given only a single reference image, extracting meaningful identity attributes becomes even more challenging, often causing the network to mirror the facial and lip structures too closely. To address these issues, we introduce RADIO, a framework engineered to yield high-quality dubbed videos regardless of the pose or expression in reference images. The key is to modulate the decoder layers using latent space composed of audio and reference features. Additionally, we incorporate ViT blocks into the decoder to emphasize high-fidelity details, especially in the lip region. Our experimental results demonstrate that RADIO displays high synchronization without the loss of fidelity. Especially in harsh scenarios where the reference frame deviates significantly from the ground truth, our method outperforms state-of-the-art methods, highlighting its robustness.

## 1. Introduction

Talking head generation [4, 13, 41, 53, 61] has become a focal point of research attention owing to its wide-ranging applications in the media industry, e.g. virtual human animation, audio-visual dubbing, and video content creation. Audio-driven talking face generation specifically aims to produce high-quality videos that exhibit precise synchronization with the driving audio. In particular, one-shot audio-driven methods are designed to generate talking faces of unseen speakers, given a single reference image.

However, it is challenging to consistently generate high-quality synced faces, due to the risk of over-fitting to the single image. In other words, previous methods face difficulties to generate mouth shapes and poses that deviate from the source image. We observed that this problem can be attributed to the susceptibility of previously proposed frameworks to the choice of reference frame. Early methods directly incorporate the information of reference image into the generator through skip-connections [9, 26, 40, 70]. These approaches constrain generated images to rarely di-

verge from the input image. Deformation-based methods [7, 52, 65, 66, 69, 72] aim to adjust facial alignment based on audio or target frames, but still struggles to generate realistic diverse poses. Significant changes in geometric priors like mesh and landmarks, or latent space also introduce artifacts and distortions in the images [44, 45, 60, 65, 71].

Despite various efforts made by previous works, we observed that there has been a limited exploration into scenarios where the dubbed video demands significantly different generated frames compared to the reference image in terms of pose and mouth expression. In fact, this scenario is the most frequently encountered in reality, as it is both inconvenient and impractical to manually select the appropriate reference frame. The key to consistently produce high-quality dubbed video is to effectively and distinctly extract the identity-related characteristics while eliminating undesired elements such as pose, facial expression, and mouth shape from the reference image.

To address this issue, we introduce **RADIO** - a method for **R**eference-**A**gnostic **D**ubbing **v**Id**e**O generation. **RADIO** aims to preserve high-fidelity details from the reference image and reduce sensitivity to the choice of the reference image, all within a unified framework. Specifically, we adopt the decoder structure from StyleGAN2 [23], and inject the reference frame information, *i.e.* style feature, through style-modulated convolution. Unlike previous methods with style-based generators as backbone, we do not inject the reference frame directly to the StyleGAN2 input [2, 60, 71]. While style modulation proves to be efficient in capturing identity-related features and diminishing structural reliance, it falls short in preserving high-fidelity details. To capture the fine texture and characteristic details of the source image, we introduce Vision Transformer (ViT) [12] within the intermediate decoder layers, preceding the style modulation. With this simple framework, we can produce talking faces in a more practical and challenging scenario where the reference face significantly differs from the target face.

Our main contributions are summarized as follows:

- We propose a simple yet effective architecture that extracts relevant information from a single reference image, thus able to create dubbing videos with improved lip synchronization that is robust from the reference pose or mouth shape.
- We improve fidelity preservation by incorporating carefully designed vision transformer blocks in the decoder, which specifically focus on lip-oriented details.
- We thoroughly evaluate **RADIO** with qualitative and quantitative experiments, and demonstrate its superiority over existing state-of-the-art methods.

## 2. Related Works

### 2.1. Audio-Driven Talking Head Generation

The task of audio-driven talking head generation learns to synthesize talking faces with lip movements synchronized with the driving audio. Early 3D-structure-based methods animate faces with 3D models such as meshes or vertex coordinates [20, 47, 73]. Unfortunately, the requirement of 3D model training data for individuals limits its application to animating general faces and struggles to reproduce teeth and hair details. In response to this challenge, recent research has shifted towards directly animating raw 2D images. 2D-based audio-driven works comprehensively fall into two categories: generating talking faces in a speaker-specific or a speaker-agnostic manner.

**Speaker-specific methods.** Personalized models generate faces in a speaker-specific manner and require re-training for an unseen identity [15, 27, 29, 32, 35, 42, 44, 48, 56, 58, 59]. Inspired by the development of neural rendering, recent methods model facial details implicitly by the hidden space of the neural radiance fields [36]. AD-NeRF [15] first proposes end-to-end audio-driven neural radiance fields for talking head generation. SSP-NeRF [32] introduces a semantic-aware dynamic ray sampling module, and DFA-NeRF [58] introduces two disentangled representations for improvement of realistic dynamics. DFRF [42] reduces the training speed via conditioning the face radiance field on 2D appearance images. Nevertheless, the need for additional training efforts and the capability of NeRF to generalize to unfamiliar identities heavily restricts its applicability.

**Speaker-agnostic methods.** Speaker-agnostic methods have gained popularity because they only require a single image of the target identity to animate the face with driving audio. Methods that generate the whole head either utilize warping techniques to drive the entire head movements [7, 18, 19, 31, 52, 64, 66, 71, 72], or generate inverted images via a well-trained encoder and a pre-trained face generator [2, 37, 60]. The former approach has controllability over head motions; however, it comes at the expense of fidelity degradation and artifacts due to the shifting of facial landmarks. The latter approach produces high-quality images but carries the risk of generating images biased by the pre-trained generators, leading to the potential leakage of the identity information.

Methods that focus on mouth regions generate synchronized lip movements with the pose fixed by the target image. Inpainting-based methods [14, 39, 40, 65, 69] exhibit high accuracy in synchronization and identity preservation. However, in a one-shot scenario where only a single reference image is available, these models fail to preserve lip-oriented high-fidelity details. Furthermore, in harsh cases where the pose or expression of the reference image is significantly dissimilar to the target image, previous methods fail to ro-

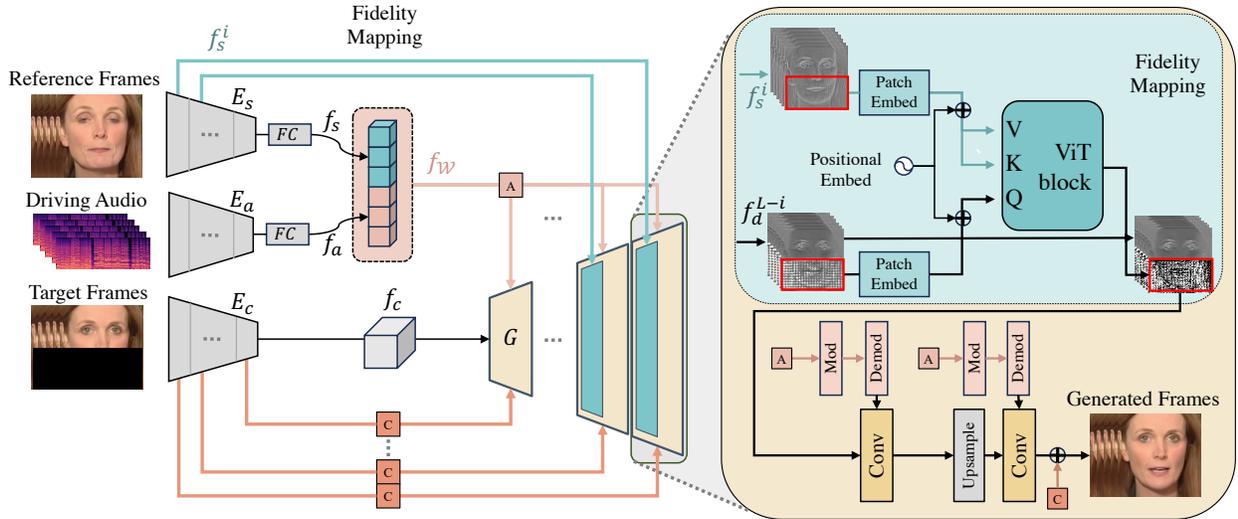


Figure 2. **Architecture of RADIO.** Our framework is composed of residual block encoders and a StyleGAN-based decoder. Lower half-masked target frames, reference frames, and mel-spectrograms are encoded by  $E_c$ ,  $E_s$ , and  $E_a$ , respectively. Basically, the generator  $G$  follows style modulation of StyleGAN2. The content feature  $f_c$  is fed to the generator  $G$  with residual mapping of the intermediate content features from each block of  $E_c$ .  $f_w$  is a concatenation of features  $f_s$  and  $f_a$ , and is mapped to each generator block as a  $W$  space for style modulation. Inserted into the last two blocks of the generator, the ViT block receives the lower half of the  $f_i$  (query), an intermediate feature flowing through the generator layers, and the lower half of the  $f_s^i$  (key, value), extracted from the front end of the style encoder  $E_s$ , with patch embedding and positional embedding. Finally, the output of ViT is re-concatenated with the upper half of the  $f_i$  and passes through the aforementioned style modulation layer. The final frames are the audio-driven high-fidelity results.

bustly generate accurate mouth shapes. Our method focuses on extracting the high-fidelity identity information robustly from a single image, without the guidance of additional geometric face priors.

## 2.2. Vision Transformer

The significant success of transformers [3, 51] in NLP has motivated numerous endeavors to extend their application to various vision tasks. Among them, Vision Transformer (ViT) [12] has shown remarkable performance across several discriminative tasks [6, 11, 25, 30, 33, 34, 38, 43, 50, 54, 57, 68]. Concurrently, recent efforts have also emerged to explore the integration of ViT into generative tasks. Several studies [28, 62, 67] have shown the competitive nature of ViT-based architectures when compared to CNN-based architectures [21–23] as the unconditional generator. Additionally, there have been attempts to utilize ViT in image-to-image translation [24, 49]. InstaFormer [24] leveraged ViT to capture the global consensus of a scene. UVCGAN [49] utilized ViT to learn pairwise relationships of low-frequency features. They commonly incorporate self-attention modules at low-resolution layers to discover the global information from a given image. On the other hand, our approach adopts ViT to generate high-fidelity results by capturing global relationships across features from

different images in high-resolution layers.

## 3. Method

In this section, we propose RADIO, an efficient one-shot audio-driven talking face architecture. Figure 2 shows the overview of the architecture design. During the training phase, RADIO receives consecutive series of target frames  $I_t \in \mathbb{R}^{T \times 3 \times H \times W}$ , randomly chosen reference frames  $I_r \in \mathbb{R}^{T \times 3 \times H \times W}$  of a target speaker, and an audio clip  $A$  aligned with the corresponding  $T$  frames as input. Our framework can create high-quality talking head videos  $I_{out} \in \mathbb{R}^{T \times 3 \times H \times W}$  where the target face speaks with high synchronization agnostic to the facial alignment of reference frame.  $T$  is the clip length and set to five, following the training strategy in [40] for the usage of sync discriminator [10]. Note that in the inference phase, only a single target and reference frame are used as inputs, *i.e.*  $T=1$ .  $H$  and  $W$  are the height and width of the frames, respectively.

### 3.1. Notation and Proposed Architecture

The proposed framework consists of four components: (i) a content encoder  $E_c$  for extracting the structural details of the target image, (ii) a style encoder  $E_s$  for capturing the visual attributes linked to the target identity, (iii) an audio encoder  $E_a$  for extracting the per-frame audio feature, and

at last (iv) a StyleGAN-based decoder  $G$  to generate images that exhibits the transferred style of the reference frames onto the target image.

**Encoder.** The content encoder  $E_c$  consists of  $L$  layers, which are constructed using  $L - 2$  residual down-blocks along with two additional convolution blocks. The encoder extracts the intermediate content features  $f_c^i$  from each of the layers,  $i \in 1, \dots, L$ , later used for residual mapping to the generator  $G$ . The final content feature  $f_c \in \mathbb{R}^{12 \times 12 \times 512}$  is used as an input of the decoder.

The structure of the style encoder  $E_s$  is similar to that of  $E_c$ . It is comprised of  $L$  layers, each producing intermediate style features  $f_s^i$ ,  $i \in 1, \dots, L$ . Note that  $E_s$  has an additional fully-connected layer to yield a sparse feature  $f_s \in \mathbb{R}^{512}$  of the reference image. We make the reference feature sparse, ensuring that the generator solely captures the broad attributes of the reference image while disregarding its finer structural details.

The audio encoder  $E_a$  receives mel-spectrogram  $A$  as input. We use the self-attentive pooling layer introduced in [5] to focus on important frame-level features. The final audio feature  $f_a \in \mathbb{R}^{512}$  is concatenated with the style feature  $f_s$  to formulate  $f_{\mathcal{W}} = \{f_s, f_a\} \in \mathbb{R}^{1024}$  as the  $\mathcal{W}$  space for the style mapping to the generator layers.

**Decoder.** The overall structure of the decoder follows the StyleGAN2 [23], with  $L$  hierarchical layers. With the content feature  $f_c$  as input and  $f_{\mathcal{W}}$  to modulate the convolution kernel weights of the generator, the decoder generates faces dubbed with the guidance of style and audio features upon the target image. Previous one-shot audio-driven works that utilize direct skip connections [40] have higher reliance to the structural information, like the poses and mouth shapes, of the reference image. That is, with a reference image with a dissimilar pose of the target image or ground truth mouth shape, the model struggles to generate high-fidelity results. Instead, we employ style modulation to convey the identity information, which eventually helps the robustness of distinct poses and mouth shapes from the reference images. We present empirical results in Section 4.4 to demonstrate that style modulation of reference image is more effective than skip-connections or direct input injections for high-quality lip-sync generation.

### 3.2. Design of Vision Transformer Blocks

With the sparse feature of the reference image delivered to the decoder, it is insufficient to reconstruct high-fidelity details of the target identity. As a solution, we incorporate Vision Transformer (ViT) [12] to restore these intricate details. We adopted the attention mechanism of ViT to understand the meaningful patterns and relationships between global image patches. With the aid of global attention, our framework is able to focus on lip-oriented regions even for misaligned reference frames compared to the target.

The ViT blocks are strategically designed to focus on the lip regions, which in our scenario corresponds to the lower half of the image. ViT blocks are attached into the final two layers, namely the  $L - 1$  and  $L$ -th layers, of the decoder. We empirically found that attention in the final two layers were the most efficient and effective (see experimental results in the supplementary material). The lower half of the intermediate feature of the decoder  $f_d^{L-i}$ , is used as the query ( $Q$ ). The lower half of the intermediate style feature  $f_s^i$  is used as the key ( $K$ ) and value ( $V$ ). We first compute the attention output with the following equation :

$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) \cdot V, \quad (1)$$

where  $Q, K, V$  are each the output of layer normalization over corresponding features, and  $d_k$  is the dimension of the key vector. The attention output is then added to the the intermediate features. We used eight multi-heads to concat the attention score, and linearly transformed them with multi-level perceptron (MLP) layers to produce the final attention layer. Finally, The output of the ViT block is concatenated with the upper half of the intermediate feature  $f_d^{L-i}$ , then passed to the style modulation layer.

The patch size of ViT block is empirically set to  $\frac{H_i}{32} \times \frac{W_i}{32}$  for the  $i$ -th intermediate features layers. The ViT blocks consist of two attention layers, resulting in a total of four attention layers considering the entire architecture. We name each of these layers  $Att_{ij}$ , where  $i \in \{L - 1, L\}$  and  $j \in \{1, 2\}$ .

### 3.3. Loss Function

In the training phase, we use five consecutive target frames aligned with the audio clip, while the reference frames are randomly chosen . We use the following training objectives to enhance image quality and synchronization accuracy.

**Reconstruction Loss.** The reconstruction loss  $L_{rec}$  is composed of an  $L_1$  pixel loss and a perceptual loss:

$$\mathcal{L}_{rec} = \|I_t - I_{out}\|_1 + \sum_{i=1}^L \lambda_i \|\phi_i(I_t) - \phi_i(I_{out})\|_1, \quad (2)$$

where  $\phi_i$  is the  $i$ -th layer of the VGG network and  $L$  is the number of VGG layers. We use different weight  $\lambda_i$  for each layer, increasing for deeper layers.

**GAN Loss.** To maintain high fidelity of the generated image, we use adversarial loss  $\mathcal{L}_{adv}$  commonly used in generative networks:

$$\mathcal{L}_{adv} = \mathbb{E} [\log(1 + \exp(D(I_{out})) + \log(1 + \exp(-D(I_t)))], \quad (3)$$

where  $D$  is the StyleGAN2 [23] discriminator, trained jointly with our generator.

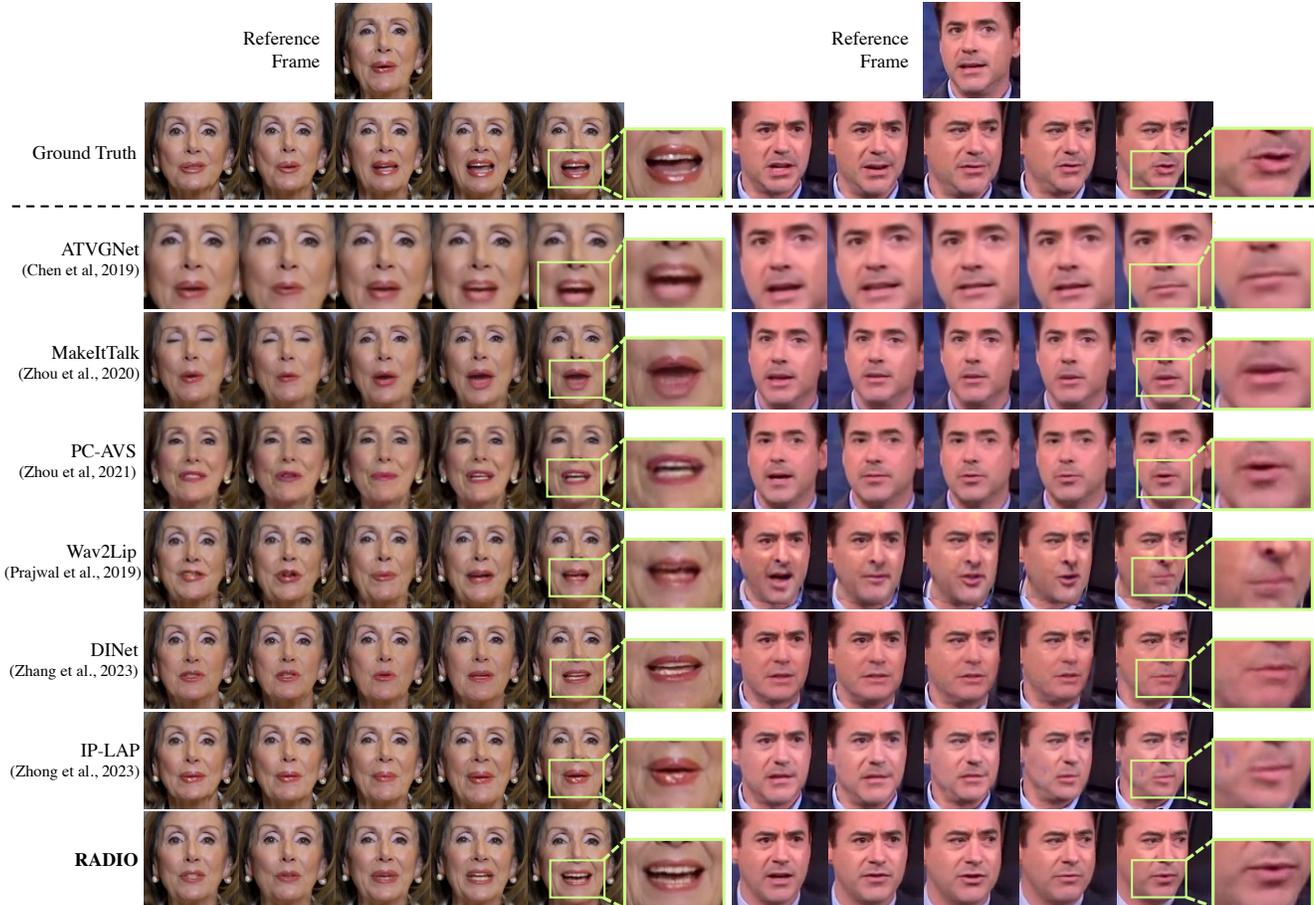


Figure 3. **Qualitative comparison with baselines.** We visualized the dubbed results for five adjacent frames from the HDTF (left) and VoxCeleb2 (right) dataset, and zoomed-in images of the mouth region for the closer inspection. The reference frame for HDTF clip was a frontalized face with closed mouth, and the reference frame for VoxCeleb2 was a face facing the left side. Our method showed the highest fidelity and accurately synced results, agnostic to the reference image.

**Sync Loss.** Following [40], we additionally train a grayscale sync discriminator  $S$ , consisting of a vision encoder  $S_v$  and audio encoder  $S_a$ . The encoder architecture follows [8] with self-attention pooling after ResNet layers. The details of our modified sync discriminator architecture can be found in the supplementary material. The sync discriminator is trained with a binary-cross entropy loss (eq. 5) to increase the cosine similarity (eq. 4) of the vision and audio features of five consecutive frames that are in-sync ( $y_i = 1$ ), while pursuing the opposite for frames that are off-sync ( $y_i = 0$ ).

$$p_i(I, A) = \frac{S_v(I_{i-2:i+2})^T \cdot S_a(A_{i-2:i+2})}{\|S_v(I_{i-2:i+2})\| \cdot \|S_a(A_{i-2:i+2})\|}, \quad (4)$$

$$\mathcal{L}_{sync} = -\mathbb{E}[y_i \log(p_i) + (1 - y_i) \log(1 - p_i)]. \quad (5)$$

During training the RADIO framework, we fix the weights of the pre-trained sync discriminator, and enhance the synchronization quality by the same sync-loss.

The overall training loss is formulated as follows:

$$\mathcal{L}_{total} = \mathcal{L}_{rec} + \lambda_{adv} \mathcal{L}_{adv} + \lambda_s \mathcal{L}_{sync}, \quad (6)$$

where  $\lambda_{adv}$  and  $\lambda_{sync}$  are balancing weights.

## 4. Experiments

### 4.1. Experimental Settings

**Datasets and Preprocessing.** We trained our RADIO and sync-discriminator with LRW [17] dataset, a commonly used audio-visual dataset with 1,000 utterances of 500 different words. For evaluation, we used 50 randomly selected videos from HDTF [66] and VoxCeleb2 [16] datasets. The HDTF dataset primarily consists of videos with frontalized

Methods	HDTF					VoxCeleb2				
	PSNR $\uparrow$	MS-SSIM $\uparrow$	LPIPS $\downarrow$	Sync-C $\uparrow$ /D $\downarrow$	LMD $\downarrow$	PSNR $\uparrow$	MS-SSIM $\uparrow$	LPIPS $\downarrow$	Sync-C $\uparrow$ /D $\downarrow$	LMD $\downarrow$
GT	100.0	1.0	0.0	8.716*/6.853*	0.0	100.0	1.0	0.0	6.363*/8.014*	0.0
ATVGNet [7]	30.101	0.780	0.170	6.422/8.640	9.893	29.555	0.716	0.194	5.203/9.087	11.454
MakeItTalk [72]	29.230	0.646	0.232	5.226/9.784	10.963	29.066	0.619	0.233	4.472/9.759	12.415
PC-AVS [71]	29.667	0.730	0.173	<u>9.060/6.461</u>	10.305	28.997	0.687	0.210	<u>7.179/7.413</u>	11.773
Wav2Lip [40]	30.645	0.818	0.135	<b>9.918/6.105</b>	7.882	30.487	0.791	0.146	<b>7.397/6.105</b>	<u>7.913</u>
DINet [65]	30.892	0.858	0.088	7.969/7.372	<u>7.632</u>	30.077	0.766	0.145	5.466/8.682	8.781
IP-LAP [69]	<u>32.111</u>	<u>0.902</u>	<u>0.068</u>	7.072/8.075	7.702	<u>31.395</u>	<u>0.844</u>	<u>0.106</u>	4.942/8.765	8.589
<b>RADIO</b>	<b>33.939</b>	<b>0.909</b>	<b>0.058</b>	8.310*/7.038*	<b>7.235</b>	<b>32.804</b>	<b>0.896</b>	<b>0.073</b>	6.671*/7.856*	<b>7.544</b>

Table 1. **Quantitative comparison of baselines.** We measured the perceptual and lip-sync quality for baselines in HDTF and VoxCeleb2 datasets. We denote the best scores in **bold**, second-best underlined, and the closest Sync-C/D scores to the ground truth with \* mark.

head poses, and VoxCeleb2 consists of videos with a wide range of head poses. All videos are resampled to a frame rate of 25 FPS.

For the image preprocessing, we first detected faces and cropped the images with FFHQ-alignment [22], then resized them to resolution  $192 \times 192$ . This alignment crop contains the whole lower half of faces, with both the mouth and nose vertically positioned at the center. We used  $L = 6$  layers for all encoders and decoder to match the resolution. For the audio preprocessing, we first converted audios to a sample rate of 16 kHz, then extracted mel-spectrograms using FFT window size 1,280, a hop length of 160, and 80 mel filter-banks.

**Baselines.** We compared our methods with person-agnostic talking face methods, including recent methods that claim to be state-of-the-art. The baselines include **ATVGNet** [7], **MakeItTalk** [72], **PC-AVS** [71], **Wav2Lip** [40], **DINet** [65], and **IP-LAP** [69]. Note that the first three methods generate the entire head driven by the audio, so the alignment of synthesized faces differs from the ground truth. The last three and RADIO correspond to inpainting-based methods, with slightly different masked regions around the mouth. DINet [65] and IP-LAP [69] authors used multiple reference frames for better grasp of identity. For a fair comparison, all baselines used only the first frame of the video as an identity reference.

## 4.2. Qualitative Evaluation

Figure 3 shows the qualitative comparisons for an example clip of five adjacent frames. Specifically, for HDTF [66] video clip on the left, generated images should resemble ground truth frames with widely opening mouths, given a source face image with a closed mouth. For VoxCeleb2 [16] video clip on the right, methods should generate realistic faces tilted rightwards, given a source face image facing the left. Except for ATVGNet [7], we aligned all generated frames using FFHQ alignment for clear comparison.

In comparison to other methods, our approach gener-

ated faces that closely resemble the ground truth in terms of visual fidelity and lip synchronization. Methods that synthesize the whole head [7, 71, 72] showed poor identity preservation and the alignment highly deviated from the ground truth due to the missing resource to drive the pose. Wav2Lip [40] generated blurry lower faces with indistinct mouth attributes. DINet [65] and IP-LAP [69] showed poor performance, with the generated lip regions more closely resembling the source image than the ground truth. IP-LAP failed to generate open lips throughout the whole video, and created artifacts for deviating pose alignments, *e.g.* seventh row. More qualitative comparisons for challenging scenarios are provided in the supplementary material.

**Experiments for Robustness.** In this section, we explore the robustness of our method by dubbing the same target frames with various reference images. Figure 4 shows generated images of RADIO using three different reference images. Specifically, each reference image are facing the front (A), the right with a closed lip (B), and the left with an open lip (C). An ideal scenario is to consistently generate accurate lips regardless of the varying pose and mouth shapes. Remarkably, our method demonstrated almost zero sensitivity to the reference image, as evidenced by the dubbed results. The synthesized images also closely resembled the ground truth, with nearly identical mouth shapes. With this guarantee of robustness, RADIO can be used without the need for additional reference frame selection processes.

## 4.3. Quantitative Evaluation

We evaluated methods with the following metrics to measure the reconstruction and lip synchronization quality. **PSNR**, **MS-SSIM** [55] and **LPIPS** [63] measure the pixel-wise and feature-wise similarity between generated and ground truth images. The SyncNet [10] confidence score **Sync-C** and distance score **Sync-D** measure the audio-visual synchronization quality. We used the officially released version of SyncNet\* [10] for a fair comparison. Lip

\*[https://github.com/joonson/syncnet\\_python](https://github.com/joonson/syncnet_python)

landmark distance (**LMD**) measures the normalized landmark distance of the mouth between generated and ground-truth images for lip synchronization evaluation.

Evaluating raw images with metrics designed for reconstruction quality would be unfair due to the varying sizes of the generation regions across different methods, *e.g.*, DNet [65] synthesizes only the small region around the mouth using its own cropping algorithm, Wav2Lip [40] inpaints the lower half of the tightly cropped face, and RADIO inpaints the lower half of the FFHQ aligned face which includes more background for generation. Since the final dubbed videos only need the tightly zoomed face to be attached, we employed a cropping method that zooms in on the faces with the same ratio and then resized them to the same resolution for evaluation. Please refer to the supplementary material for details on the cropping method.

Table 1 shows the quantitative comparison between competing methods. Our proposed method achieved the best performance on all visual quality metrics (PSNR, MS-SSIM, LPIPS). While IP-LAP [69] generated comparable visual quality synthesizing talking face to ours, it fell short in audio-visual synchronization, such as Sync-C/D and LMD. The reason for this is that IP-LAP focuses on the personalized facial traits and only uses landmarks priors for lip synchronization, which can be quite inaccurate from a single image. Regarding audio-visual synchronization, our method performed the best on LMD. Although Wav2Lip [40] performed the best on Sync-C/D, this result may be attributed its use of SyncNet, which was pre-trained on the same LRS2 [1] dataset as the official SyncNet used for our evaluation. PC-AVS [71] also achieved high Sync-C/D scores, but generated jerky lip movements to match the audio, which degraded the LMD scores. While our method demonstrated the third-best performance on Sync-C/D, it’s noteworthy that this score is the closest to the ground truth. RADIO generated the most synchronized natural lip movements, as supported by qualitative results. Compared to baselines, RADIO is the only method that can robustly deliver both fidelity preservation and synchronization.

#### 4.4. Ablation Study of ViT blocks

In this section, we first analyzed the ViT blocks via visualizing the attention map. Then, we demonstrated the effectiveness of our ViT design with a thorough ablation study.

**Attention Visualization.** Figure 5 displays the attention map, which highlights the important region in the reference image for each corresponding patch of the synthesized image. Arbitrary identity is presented on the first column in each row. For each row, the patch location is illustrated with green on the generated images in the upper half, and the reference image with the corresponding attention map in the lower half. In particular, we visualized the attention map of  $Att_{5,2}$ , which is located in the second layer of the



Figure 4. **Qualitative validation of the robustness of RADIO.** Our method consistently produced accurately dubbed videos, showcasing its robustness in generating lip-synchronized content regardless of the variations in the reference frames.

attention block in the fifth decoder layer. The attention map of ViT block in the last decoder layer is also visualized in the supplementary material.

Even with differently synthesized mouth shapes compared to the reference image, patches near the mouth successfully attended to important features like cheeks and similar locations of the mouth, as seen in the second and third columns. In contrast, patches unrelated to mouth details primarily focused their attention on corresponding regions on the reference face, with less attention directed towards the mouth, as evident in the fourth and fifth columns. This phenomenon showed that our proposed ViT block leverages its global context understanding and semantic knowledge to successfully focus on lip-oriented details. With the capacity of ViT blocks to substantially guide global attention, RADIO can generate high-fidelity talking faces, even with misaligned reference frames.

**Ablation Study.** We further conducted ablation studies to validate the effectiveness of our proposed framework. In Table 2, we quantitatively compared the performance while changing each component we proposed, evaluated on the LRW [17] validation dataset. We used PSNR and LPIPS to assess the perceptual quality compared to the ground truth. For assessing lip-sync accuracy, we calculated the similarity score between audio and visual SyncNet features (eq. 4), employing our SyncNet pre-trained on the LRW train dataset. In this comparison, all models were trained for 210K iterations with a batch size of 16, with resolution scaled down to  $96 \times 96$  for the sake of resource efficiency.

We added up different parts of our model starting from a **baseline (A)**, which directly injects the reference frame in-



Figure 5. **Analysis of ViT Blocks.** For three different identities, we visualized the green patches on generated frames (upper half) with the attention map on reference frames (lower half). Our well-trained attention layer consistently focused on the globally relevant region of the reference frame for each local patch.

formation to the decoder consisting of no ViT blocks. In this configuration, the model takes concatenated target and reference frames as input to the visual encoder, following the input design of Wav2Lip [40]. Instead of learning identity via style modulated convolution, the baseline uses the reference feature as input to the decoder, with decoder layers only modulated by the audio feature, *i.e.*  $f_{\mathcal{W}} = f_a \in \mathbb{R}^{512}$ . Quantitative results show that the generated images become more susceptible to the influence of the reference frame, thereby compromising the synchronization quality.

**Method B** separates the reference and target frames using individual encoders, specifically the style encoder and content encoder. Then it utilizes the style feature  $f_s$  to compose the latent space for style modulation, along with the audio feature, *i.e.*  $f_{\mathcal{W}} = \{f_a, f_s\} \in \mathbb{R}^{1024}$ . In this regime, the model captures less structural information from the reference frame, which enables it to improve synchronization quality by reducing its dependence on the source mouth shape. **Method C** builds upon method B by incorporating a fidelity mapping via a straightforward addition adding the lower half of reference frame features to the intermediate decoder layers, *i.e.*,  $f_d^{L-i} + f_s^i$ . Even with this naïve skip-connection of reference frame, we observe improvements across perceptual evaluation metrics. However, this simple integration eventually degraded the synchronization quality, because of its sensitivity to the reference frame.

**Method D**, our RADIO framework, attached ViT blocks into the decoder layers to selectively extract high-fidelity details necessary for generating synchronized mouth movements. With the combination of modulated convolution via style features and fidelity mapping via ViT, our framework earned the most gain compared to the baseline model. This configuration demonstrated the best quantitative per-

Method	PSNR $\uparrow$	LPIPS $\downarrow$	Sync $\uparrow$
A Baseline	32.672	0.040	0.520
B + Style modulation	33.089	0.072	0.576
C + Fidelity mapping w/o ViT	34.493	0.037	0.554
D + Fidelity mapping w/ ViT	<b>34.938</b>	<b>0.031</b>	<b>0.609</b>

Table 2. **Ablation study with quantitative evaluation on LRW [17].** We varied the latent space for modulated convolution and the method for fidelity mapping. Our framework (D) achieved the most improvement compared to the baseline.

formance by simultaneously learning high-fidelity details and maintaining high synchronization.

## 5. Conclusion

In this paper, we presented an efficient framework that generates accurately dubbed faces with lip-oriented details preserved from a single source image. Our work especially stood out in the challenging yet under-investigated scenario where the face orientation and lip shape between the source and target frames are significantly different. RADIO adapts to the identity of a person via StyleGAN2 style modulation, whilst reducing the reliance on facial alignment. With the aid of ViT blocks, RADIO is finally able to synthesize faces with high-fidelity details by focusing on the important facial attributes of the reference image. Through extensive experimentations, our method has demonstrated its unique capability to consistently generate high-fidelity videos while maintaining precise lip synchronization. This achievement establishes it as the new state-of-the-art in the field of one-shot audio-driven talking face generation. Considering the simplicity and practical applicability of our framework, we look forward to wide usage in future work.

## 6. Acknowledgements

This work was supported by the National Research Foundation of Korea (NRF) grant funded by the Korea government Ministry of Science and ICT (MSIT) (No. 2021R1G1A1094379), and in part by MSIT under the Information Technology Research Center (ITRC) support program (IITP-2023-RS-2023-00258649) supervised by the Institute for Information & Communications Technology Planning & Evaluation (IITP), and in part by the IITP grant funded by MSIT (Artificial Intelligence Innovation Hub) under Grant 2021-0-02068, and in part by IITP grant funded by MSIT (No.2020-0-01336, Artificial Intelligence Graduate School Program (UNIST), No.2022-0-00264, Comprehensive Video Understanding and Generation with Knowledge-based Deep Logic Neural Network), and NRF grant funded by MSIT (No. 2.220574.01).

## References

- [1] T. Afouras, J. S. Chung, A. Senior, O. Vinyals, and A. Zisserman. Deep audio-visual speech recognition. In *arXiv:1809.02108*, 2018. 7
- [2] Mohammed M. Alghamdi, He Wang, Andrew J. Bulpitt, and David C. Hogg. Talking head from speech audio using a pre-trained image generator. In *Proceedings of the 30th ACM International Conference on Multimedia*, 2022. 2
- [3] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020. 3
- [4] Egor Burkov, Igor Pasechnik, Artur Grigorev, and Victor Lempitsky. Neural head reenactment with latent pose descriptors. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 1
- [5] Weicheng Cai, Jinkun Chen, and Ming Li. Exploring the encoding layer and loss function in end-to-end speaker and language recognition system. In *The Speaker and Language Recognition Workshop*, 2018. 4
- [6] Hu Cao, Yueyue Wang, Joy Chen, Dongsheng Jiang, Xiaopeng Zhang, Qi Tian, and Manning Wang. Swin-unet: Unet-like pure transformer for medical image segmentation. In *European conference on computer vision*, pages 205–218. Springer, 2022. 3
- [7] 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 Conference on Computer Vision and Pattern Recognition*, pages 7832–7841, 2019. 2, 6, 12
- [8] Joon Son Chung, Jaesung Huh, Seongkyu Mun, Minjae Lee, Hee Soo Heo, Soyeon Choe, Chiheon Ham, Sunghwan Jung, Bong-Jin Lee, and Icksang Han. In defence of metric learning for speaker recognition. In *Proc. Interspeech*, 2020. 5, 12
- [9] Joon Son Chung, Amir Jamaludin, and Andrew Zisserman. You said that? In *British Machine Vision Conference (BMVC)*, 2017. 1
- [10] J. S. Chung and A. Zisserman. Out of time: automated lip sync in the wild. In *Workshop on Multi-view Lip-reading, ACCV*, 2016. 3, 6
- [11] Zhigang Dai, Bolun Cai, Yugeng Lin, and Junying Chen. Up-detr: Unsupervised pre-training for object detection with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1601–1610, 2021. 3
- [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR*, 2021. 2, 3, 4
- [13] Michail Christos Doukas, Stefanos Zafeiriou, and Viktoriia Sharmanska. Headgan: One-shot neural head synthesis and editing. In *Proceedings of the IEEE/CVF International conference on Computer Vision*, pages 14398–14407, 2021. 1
- [14] 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)*, 2023. 2
- [15] 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 *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 2
- [16] A. Nagrani, J. S. Chung and A. Zisserman. Voxceleb2: Deep speaker recognition. In *INTERSPEECH*, 2018. 5, 6
- [17] A. Zisserman, J. S. Chung. Lip reading in the wild. In *Asian Conference on Computer Vision (ACCV)*, 2016. 5, 7, 8, 12, 13
- [18] Xinya Ji, Hang Zhou, Kaisiyuan Wang, Qianyi Wu, Wayne Wu, Feng Xu, and Xun Cao. Eamm: One-shot emotional talking face via audio-based emotion-aware motion model. *arXiv preprint arXiv:2205.15278*, 2022. 2
- [19] Xinya Ji, Hang Zhou, Kaisiyuan Wang, Wayne Wu, Chen Change Loy, Xun Cao, and Feng Xu. Audio-driven emotional video portraits. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021. 2
- [20] Tero Karras, Timo Aila, Samuli Laine, Antti Herva, and Jaakko Lehtinen. Audio-driven facial animation by joint end-to-end learning of pose and emotion. *ACM Trans. Graph.*, 36(4), 2017. 2
- [21] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. *Advances in Neural Information Processing Systems*, 34:852–863, 2021. 3
- [22] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (CVPR)*, pages 4401–4410, 2019. 3, 6
- [23] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2, 3, 4, 12, 13
- [24] Soohyun Kim, Jongbeom Baek, Jihye Park, Gyeongnyeon Kim, and Seungryong Kim. Instaformer: Instance-aware image-to-image translation with transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18321–18331, 2022. 3
- [25] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023. 3
- [26] Prajwal KR, Rudrabha Mukhopadhyay, Jerin Philip, Abhishek Jha, Vinay Namboodiri, and CV Jawahar. Towards automatic face-to-face translation. In *Proceedings of the 27th ACM international conference on multimedia*, pages 1428–1436, 2019. 1

- [27] A. Lahiri, V. Kwatra, C. Frueh, J. Lewis, and C. Bregler. Lipsync3d: Data-efficient learning of personalized 3d talking faces from video using pose and lighting normalization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2754–2763, 2021. 2
- [28] Kwonjoon Lee, Huiwen Chang, Lu Jiang, Han Zhang, Zhuowen Tu, and Ce Liu. Vitgan: Training gans with vision transformers. *arXiv preprint arXiv:2107.04589*, 2021. 3
- [29] Jiahe Li, Jiawei Zhang, Xiao Bai, Jun Zhou, and Lin Gu. Efficient region-aware neural radiance fields for high-fidelity talking portrait synthesis. *arXiv preprint arXiv:2307.09323*, 2023. 2
- [30] Yanghao Li, Chao-Yuan Wu, Haoqi Fan, Karttikeya Mangalam, Bo Xiong, Jitendra Malik, and Christoph Feichtenhofer. Mvitv2: Improved multiscale vision transformers for classification and detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4804–4814, 2022. 3
- [31] Borong Liang, Yan Pan, Zhizhi Guo, Hang Zhou, Zhibin Hong, Xiaoguang Han, Junyu Han, Jingtuo Liu, Errui Ding, and Jingdong Wang. Expressive talking head generation with granular audio-visual control. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3387–3396, June 2022. 2
- [32] Xian Liu, Yinghao Xu, Qianyi Wu, Hang Zhou, Wayne Wu, and Bolei Zhou. Semantic-aware implicit neural audio-driven video portrait generation, 2022. 2
- [33] Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, et al. Swin transformer v2: Scaling up capacity and resolution. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12009–12019, 2022. 3
- [34] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 10012–10022, 2021. 3
- [35] Yuanxun Lu, Jinxiang Chai, and Xun Cao. Live speech portraits: real-time photorealistic talking-head animation. *ACM Transactions on Graphics (TOG)*, 40(6):1–17, 2021. 2
- [36] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European Conference on Computer Vision*, 2020. 2
- [37] Dongchan Min, Minyoung Song, and Sung Ju Hwang. Styletalker: One-shot style-based audio-driven talking head video generation, 2022. 2
- [38] Ishan Misra, Rohit Girdhar, and Armand Joulin. An end-to-end transformer model for 3d object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2906–2917, 2021. 3
- [39] 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. 2
- [40] K R Prajwal, Rudrabha Mukhopadhyay, Vinay P. Nambodiri, and C.V. Jawahar. 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*, page 484–492, 2020. 1, 2, 3, 4, 5, 6, 7, 8, 13, 14
- [41] Yurui Ren, Ge Li, Yuanqi Chen, Thomas H. Li, and Shan Liu. Pirenderer: Controllable portrait image generation via semantic neural rendering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 1
- [42] Shuai Shen, Wanhua Li, Zheng Zhu, Yueqi Duan, Jie Zhou, and Jiwen Lu. Learning dynamic facial radiance fields for few-shot talking head synthesis. In *European conference on computer vision*, 2022. 2
- [43] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 7262–7272, 2021. 3
- [44] Supasorn Suwajanakorn, Steven M. Seitz, and Ira Kemelmacher-Shlizerman. Synthesizing obama: Learning lip sync from audio. *ACM Trans. Graph.*, 36(4), 2017. 2
- [45] Anni Tang, Tianyu He, Xu Tan, Jun Ling, Runnan Li, Sheng Zhao, Li Song, and Jiang Bian. Memories are one-to-many mapping alleviators in talking face generation, 2022. 2
- [46] Xuansong Xie Tao Yang, Peiran Ren and Lei Zhang. Gan prior embedded network for blind face restoration in the wild. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 14
- [47] Sarah L. Taylor, Taehwan Kim, Yisong Yue, Moshe Mahler, James Krahe, Anastasio Garcia Rodriguez, Jessica K. Hodgins, and Iain A. Matthews. A deep learning approach for generalized speech animation. *ACM Trans. Graph.*, 36(4):93:1–93:11, 2017. 2
- [48] Justus Thies, Mohamed Elgharib, Ayush Tewari, Christian Theobalt, and Matthias Nießner. Neural voice puppetry: Audio-driven facial reenactment. In *European Conference on Computer Vision*, 2020. 2
- [49] Dmitrii Torbunov, Yi Huang, Haiwang Yu, Jin Huang, Shinjae Yoo, Meifeng Lin, Brett Viren, and Yihui Ren. Uvcgan: Unet vision transformer cycle-consistent gan for unpaired image-to-image translation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 702–712, 2023. 3
- [50] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International conference on machine learning*, pages 10347–10357. PMLR, 2021. 3
- [51] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 3
- [52] Suzhen Wang, Lincheng Li, Yu Ding, and Xin Yu. One-shot talking face generation from single-speaker audio-visual correlation learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 2531–2539, 2022. 2

- [53] 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*, 2021. 1
- [54] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 568–578, 2021. 3
- [55] Z. Wang, E.P. Simoncelli, and A.C. Bovik. Multiscale structural similarity for image quality assessment. In *The Thirty-Seventh Asilomar Conference on Signals, Systems & Computers, 2003*, volume 2, pages 1398–1402 Vol.2, 2003. 6
- [56] Xin Wen, Miao Wang, Christian Richardt, Ze-Yin Chen, and Shi-Min Hu. Photorealistic audio-driven video portraits. *IEEE Transactions on Visualization and Computer Graphics*, 26(12):3457–3466, 2020. 2
- [57] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in Neural Information Processing Systems*, 34:12077–12090, 2021. 3
- [58] Shunyu Yao, RuiZhe Zhong, Yichao Yan, Guangtao Zhai, and Xiaokang Yang. Dfa-nerf: Personalized talking head generation via disentangled face attributes neural rendering. *arXiv preprint arXiv:2201.00791*, 2022. 2
- [59] Ran Yi, Zipeng Ye, Juyong Zhang, Hujun Bao, and Yongjin Liu. Audio-driven talking face video generation with learning-based personalized head pose. *arXiv preprint arXiv:2002.10137*, 2020. 2
- [60] Fei Yin, Yong Zhang, Xiaodong Cun, Mingdeng Cao, Yanbo Fan, Xuan Wang, Qingyan Bai, Baoyuan Wu, Jue Wang, and Yujiu Yang. Styleheat: One-shot high-resolution editable talking face generation via pre-trained stylegan. *arxiv:2203.04036*, 2022. 2
- [61] Egor Zakharov, Aliaksandra Shysheya, Egor Burkov, and Victor Lempitsky. Few-shot adversarial learning of realistic neural talking head models, 2019. 1
- [62] Bowen Zhang, Shuyang Gu, Bo Zhang, Jianmin Bao, Dong Chen, Fang Wen, Yong Wang, and Baining Guo. Styleswin: Transformer-based gan for high-resolution image generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11304–11314, 2022. 3
- [63] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018. 6
- [64] Wenxuan Zhang, Xiaodong Cun, Xuan Wang, Yong Zhang, Xi Shen, Yu Guo, Ying Shan, and Fei Wang. Sadtalker: Learning realistic 3d motion coefficients for stylized audio-driven single image talking face animation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8652–8661, 2023. 2
- [65] Zhimeng Zhang, Zhipeng Hu, Wenjin Deng, Changjie Fan, Tangjie Lv, and Yu Ding. Dinet: Deformation inpainting network for realistic face visually dubbing on high resolution video. In *AAAI*, 2023. 2, 6, 7, 12, 13
- [66] Zhimeng Zhang, Lincheng Li, Yu Ding, and Changjie Fan. Flow-guided one-shot talking face generation with a high-resolution audio-visual dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3661–3670, 2021. 2, 5, 6
- [67] Long Zhao, Zizhao Zhang, Ting Chen, Dimitris Metaxas, and Han Zhang. Improved transformer for high-resolution gans. *Advances in Neural Information Processing Systems*, 34:18367–18380, 2021. 3
- [68] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao Xiang, Philip HS Torr, et al. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6881–6890, 2021. 3
- [69] Weizhi Zhong, Chaowei Fang, Yinqi Cai, Pengxu Wei, Gangming Zhao, Liang Lin, and Guanbin Li. Identity-preserving talking face generation with landmark and appearance priors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9729–9738, 2023. 2, 6, 7, 13, 14
- [70] Hang Zhou, Yu Liu, Ziwei Liu, Ping Luo, and Xiaogang Wang. Talking face generation by adversarially disentangled audio-visual representation. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2019. 1
- [71] Hang Zhou, Yasheng Sun, Wayne Wu, Chen Change Loy, Xiaogang Wang, and Ziwei Liu. Pose-controllable talking face generation by implicitly modularized audio-visual representation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2, 6, 7
- [72] Yang Zhou, Xintong Han, Eli Shechtman, Jose Echevarria, Evangelos Kalogerakis, and Dingzeyu Li. Makeittalk: speaker-aware talking-head animation. *ACM Transactions on Graphics (TOG)*, 39(6):1–15, 2020. 2, 6, 12
- [73] Yang Zhou, Zhan Xu, Chris Landreth, Evangelos Kalogerakis, Subhansu Maji, and Karan Singh. Visemenet: Audio-driven animator-centric speech animation. *ACM Transactions on Graphics (TOG)*, 2017. 2

## A. Architectural Details

In this section, we provide a detailed description of the RADIO encoders and the sync discriminator.

**Encoders.** The two encoders ( $E_c$ ,  $E_s$ ) mentioned in the main paper have a similar structure consisting of residual-down blocks. The  $E_c$  and  $E_s$  receive RGB frames with  $192 \times 192$  resolution and pass four numbers of a residual-down block along with two additional convolution blocks. A single residual-down block involves two convolutional layers (kernel size=3) and LeakyReLU activation with a skip connection that adds intermediate features passed through the additional convolutional layer (kernel size=1). The  $f_c \in \mathbb{R}^{12 \times 12 \times 512}$ , which is an output of  $E_c$ , is fed to generator  $G$ . The only difference between  $E_c$  and  $E_s$  is that  $E_s$  extracts  $f_s \in \mathbb{R}^{1 \times 1 \times 512}$  via a spatial dimensional global average pooling operation and a fully-connected layer after four residual-down blocks.

The audio encoder  $E_a$  receives mel-spectrogram as inputs and encodes to  $f_a \in \mathbb{R}^{1 \times 1 \times 512}$  through 2D convolutional layers. The  $E_a$  is implemented to follow [8] structure. [8] is considered one of the most effective architectures for speaker recognition using audio inputs and comprises SE layers and self-attention pooling with ResNet layers. We modified the activation function as LeakyReLU and normalization as instance normalization.

**Sync Discriminator.** The sync discriminator consists of an encoder that receives mel-spectrogram as input and an encoder that receives facial images as input. The audio encoder features follow exactly the structure of [8], while the visual encoder utilizes channel-attention and spatial-attention operations instead of self-attention of [8]. This is because it is important to concentrate on the mouth’s shape within the face image or the local area around it. Finally, the sync discriminator is pre-trained with a loss function (eq. 4 in main paper) to increase the cosine similarity (eq. 5 in main paper) of the vision and audio features so that we can provide superior audio-video synchronization errors during the RADIO training scheme.

We pre-trained the sync discriminator on the LRW [17] dataset using tightly zoomed face images with a resolution of  $144 \times 144$ . This approach allowed the discriminator to focus specifically on the lip shape of the synthesized facial image. In addition, the images were converted to grayscale, making the sync discriminator color-agnostic and enabling it to focus solely on learning the sync accuracy of mouth shapes.

## B. Pre-processing algorithm for evaluation.

Including our proposed method, baselines generated different sizes of images with different alignments and different target regions for synthesis. In order to evaluate quantitative metrics fairly, we applied a pre-processing algo-

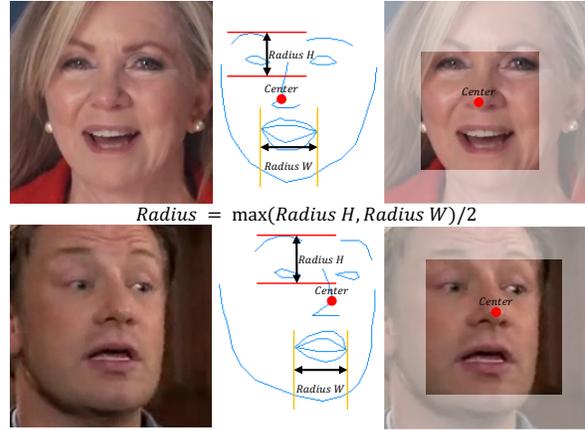


Figure 6. **Cropping method used for evaluation.** We applied this cropping method to all generated baseline results to evaluate the quantitative metrics with the same ground truth images.

rithm to all generated images before comparison. First, we aligned all baseline methods with FFHQ alignment [23]. Then, we applied a face cropping method based on the DInet [65] masking algorithm. Last, we resized tightly cropped faces to the same resolution.

Figure 6 depicts the face cropping method. We assumed that the tip of the nose (the thirty-fourth point of the facial landmark) is in the center of the human face. Then, we computed two radius values:  $Radius H$ , which measures the distance between the highest point and the thirtieth point along the y-axis of the facial landmark, and  $Radius W$ , which measures the distance between the fifty-fifth point and the forty-ninth point along the x-axis of the facial landmark. We then set the final  $Radius$  value as the maximum value between these two distances. Finally, we cropped the attached facial image with this  $Radius$  value, starting from the thirty-fourth point on the facial landmark.

## C. Baseline Models

In this section, we describe additional details about baselines mentioned in Section 4.1 of the main paper.

**ATVGnet.** ATVGnet [7] proposes constructing high-level representation (facial landmarks) from the audio signal and generating talking head videos conditioned on the facial landmark. ATVGnet leverages the pixel-wise loss with attention mechanisms to ensure temporal consistency and utilizes a regression-based discriminator to generate accurate facial shapes and realistic-looking images in the training scheme. Finally, ATVGnet can only generate  $128 \times 128$  resolution videos and cannot keep up with the head motion of the target frames.

**MakeItTalk.** MakeItTalk [72] also proposes an audio-to-landmarks approach for controlling the motion of lips

while determining the specifics of facial expressions and the rest of the talking-head dynamics from an audio signal. After that, MakeItTalk generates talking head animations ( $256 \times 256$  resolution) with a single image (cartoon or natural human) and predicted landmarks using image-to-image translation. MakeItTalk animates talking head videos based on facial landmarks extracted from audio signals. However, these facial landmarks are too sparse to describe lip motion details and do not represent significant head motion.

**Wav2Lip.** To the best of our knowledge, Wav2Lip [40] is the first approach to utilize a pre-trained Sync Discriminator in a training scheme and generate the lower half masked of the target frame. This method guarantees high audio-visual synchronization. However, it is highly dependent on the reference frame by feeding with concatenating the reference frame and masked input frame, and generates blurry results.

**DINet.** DINet [65] proposes a deformation inpainting network, which performs spatial deformation on feature maps of reference images to synthesize high-fidelity dubbing videos. DINet uses five reference facial images to create deformed features in order to align head poses and driving audio to preserve high-frequency details. In addition, DINet develops its masking algorithm around the lip, resulting in efficient inpainting synthesis of mouth shapes. Finally, DINet can generate high-resolution ( $416 \times 320$ ) videos. Although DINet utilizes multiple reference images, the synthesized results vary sensitively depending on the selected reference images. Also, their framework is restricted to frontalized head poses, and generates artifacts when the mouth region covers the background.

**IP-LAP.** IP-LAP [69] follows a two-stage training scheme like any other method of utilizing facial landmarks. IP-LAP is implemented to leverage facial sketch maps rather than face landmark coordinates so the framework can learn the driving face shapes clearly. Finally, IP-LAP aligns the twenty-five reference images using a warping-based alignment module and utilizes them to generate  $128 \times 128$  resolutions while preserving the target head pose and expression. Despite the usage of a large number of reference images, the accuracy of the lip shape is insufficient to learn the audio-visual synchronization only with facial landmarks.

## D. Ablation Study of ViT design

**Attention Visualization of last ViT block.** Figure 7 visualizes the attention map of  $Att_{6,2}$ , located in the second layer of the attention block within the last ( $L = 6$ ) decoder layer. The upper half of the figure displays the green patches on the generated frames, while the lower half presents the corresponding attention map on the reference image. In contrast to Figure 5 in the main paper, each patch attended to the entire image for the last ViT block.



Figure 7. **Visualization of attention map in the last ViT block.** We visualized the green patches on generated frames (upper half) with the attention map on reference frames (lower half).

Configuration	PSNR $\uparrow$	LPIPS $\downarrow$	Sync $\uparrow$
Baseline	33.089	0.072	0.576
+ ViT(1, 1)	34.637	0.033	0.557
+ ViT(1, 2)	34.757	0.032	0.559
+ ViT(2, 2)	<b>34.938</b>	<b>0.031</b>	<b>0.609</b>

Table 3. **Ablation for the different number of attention layers.**

We relate this phenomenon to the hierarchical nature of StyleGAN2 [23], which generates course-to-fine information for low-to-higher layers. While the intermediate attention layer, *i.e.*,  $Att_{5,2}$ , focused on the globally relevant features for each local patch, the last attention layer, *i.e.*,  $Att_{5,2}$ , captured the fine-grained textures and colors across the entire image.

**Design of ViT attention layers.** We additionally present quantitative results, evaluated on the LRW [17] validation dataset, for an ablation study of RADIO with varying numbers of ViT layers. Our evaluation metrics include PSNR, LPIPS, and the similarity score between audio and visual features. To obtain the similarity score, both the audio and visual features are encoded using the encoders of our sync discriminator described in Section A. We specifically used our pre-trained sync discriminator, which was trained with the LRW [17] training dataset, for accurate evaluation. For this experiment, all models were trained for 210K iterations with a batch size of 16, with resolution scaled down to  $96 \times 96$ , like the main ablation experiment. We only conducted the experiment with ViT on the last two layers of the decoder, because patches for earlier layers were too small to deliver semantically interpretable results. For example, the feature resolution is  $12 \times 12$  for the fourth decoder layer, which is too small to divide into patches.

In Table 3, the baseline refers to the framework that generates audio-driven images by decoder layers modulated with style features, without additional components for fidelity mapping (method B in Table. 2 of the main paper). We denote  $ViT(n, m)$  as our RADIO framework with  $n$  ViT blocks, consisting of  $m$  attention layers. Note that in our main experiment, we applied ViT block to the last two decoder layers, with each block comprising two attention lay-

ers, *i.e.*, ViT(2, 2). The results indicated that having two attention layers in a single ViT block was better than using only one layer. Additionally, employing ViT blocks in two decoder layers was more effective than placing them in a single decoder layer. Finally, ViT(2, 2) achieved the best PSNR, LPIPS, and sync similarity scores compared to the baseline.

## E. Additional Experimental Results

In this section, we show the additional experimental results of RADIO in Figure 8 and Figure 9. Throughout all examples, our results consistently generated the most natural and realistic mouth shapes, with high synchronization accuracy compared to the ground truth. ATVGNet, MakeItTalk, and PC-AVS commonly failed to generate identity-preserving details. Wav2Lip consistently created blurry images and generated artifacts for extreme face poses. Especially in harsh scenarios, IP-LAP and DInet struggled to generate realistic-looking mouth shapes, due to the significant distortion caused by warping and deformation. The mouth shapes generated by these methods were similar across all time steps, which also led to a degradation in synchronization quality. Especially, DInet failed to generate realistic faces with extreme poses, as their framework is limited to generate frontalized faces.

## F. Limitation and Broader Societal Impact

While our model excels in producing high-quality images around the mouth region, it struggles to generate a natural-looking background. During our evaluation, we observed that frames significantly misaligned with the reference frame exhibited artifacts in the background. This limitation is observed across all baseline models [40, 69], but is more conspicuous for ours due to the alignment method that includes a larger portion of the background for generation. This issue can be easily fixed by borrowing a face-parsing model [46] to attach only the face region to the original video, thus improving the overall video quality.

Previous one-shot audio-driven frameworks have struggled to consistently generate realistic, high-fidelity frames. These challenges arise because they heavily rely on the reference image, which typically requires a frontalized pose with a neutral facial expression. In contrast, our reference-agnostic framework demonstrates exceptional capabilities in generating high-quality dubbed videos, even in the most challenging scenarios. This makes it suitable for a wide range of real-world industrial applications where diverse poses and expressions are encountered. We look forward to the application of our framework to generate realistic audio-driven faces for unseen speakers in real-time. Looking ahead, we aspire to enhance and extend our RADIO framework to support higher resolutions in the near future.

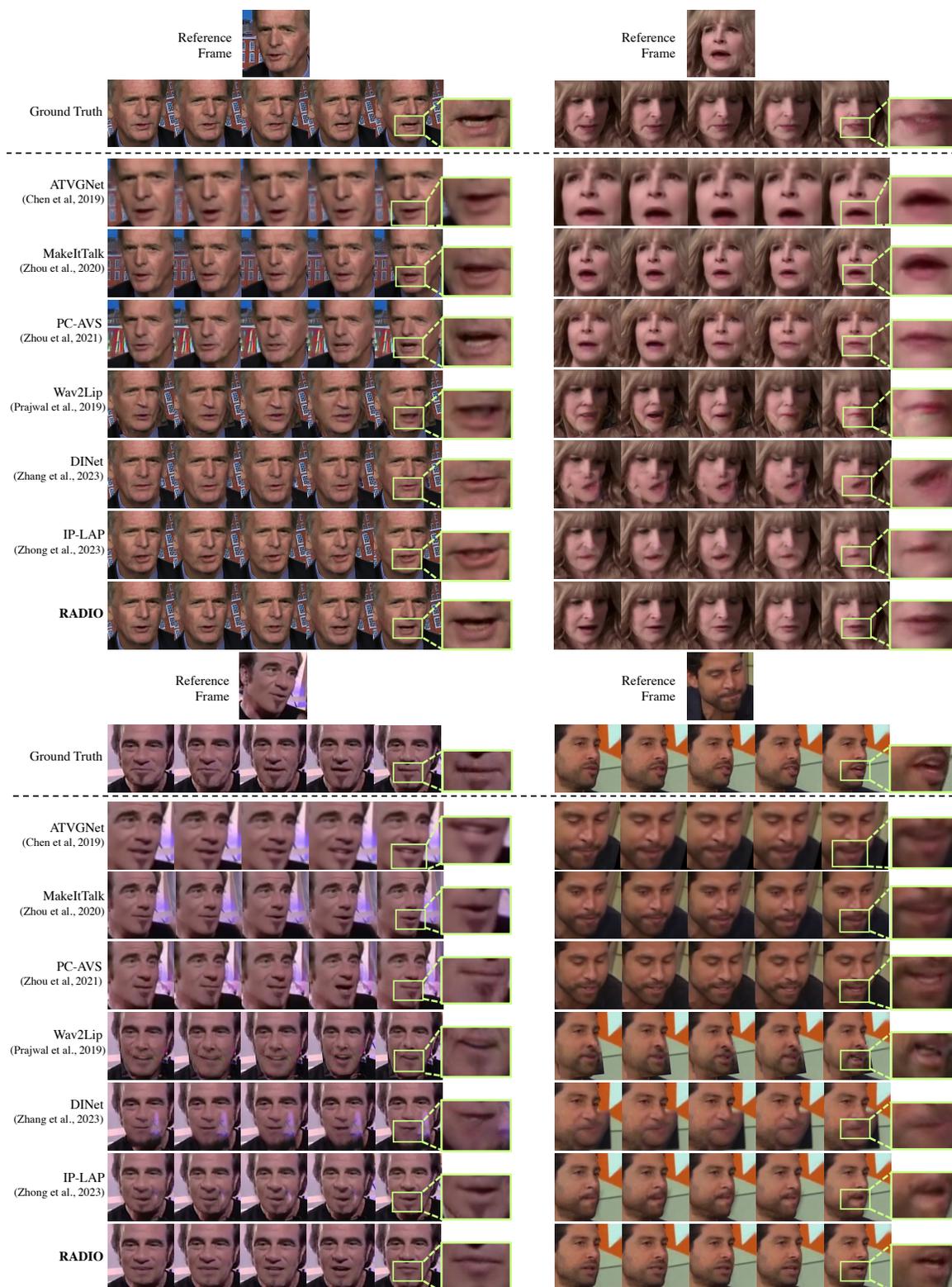


Figure 8. **Qualitative comparison with baselines.** We visualized the dubbed results for challenging scenarios where the ground truth pose and expression significantly differ from the reference frame.



Figure 9. **Qualitative comparison with baselines.** We visualized the dubbed results for challenging scenarios where the ground truth pose and expression significantly differ from the reference frame.