
GeneFace++: Generalized and Stable Real-Time Audio-Driven 3D Talking Face Generation

Zhenhui Ye^{*†}
Zhejiang University
zhenhuiye@zju.edu.cn

Jinzheng He^{*}
Zhejiang University
jinzhenghe@zju.edu.cn

Ziyue Jiang^{*}
Zhejiang University
ziyuejiang@zju.edu.cn

Rongjie Huang
Zhejiang University
rongjiehuang@zju.edu.cn

Jiawei Huang
Zhejiang University
jiawei Huang@zju.edu.cn

Jinglin Liu
Zhejiang University
jinglinliu@zju.edu.cn

Yi Ren
ByteDance
ren.yi@bytedance.com

Xiang Yin
ByteDance
yinxiang.stephen@bytedance.com

Zejun Ma
ByteDance
majejun@bytedance.com

Zhou Zhao[‡]
Zhejiang University
zhaozhou@zju.edu.cn

Abstract

Generating talking person portraits with arbitrary speech audio is a crucial problem in the field of digital human and metaverse. A modern talking face generation method is expected to achieve the goals of *generalized audio-lip synchronization*, *good video quality*, and *high system efficiency*. Recently, neural radiance field (NeRF) has become a popular rendering technique in this field since it could achieve high-fidelity and 3D-consistent talking face generation with a few-minute-long training video. However, there still exist several challenges for NeRF-based methods: 1) as for the *lip synchronization*, it is hard to generate a long facial motion sequence of high temporal consistency and audio-lip accuracy; 2) as for the *video quality*, due to the limited data used to train the renderer, it is vulnerable to out-of-domain input condition and produce bad rendering results occasionally; 3) as for the *system efficiency*, the slow training and inference speed of the vanilla NeRF severely obstruct its usage in real-world applications. In this paper, we propose GeneFace++ to handle these challenges by 1) utilizing the pitch contour as an auxiliary feature and introducing a temporal loss in the facial motion prediction process; 2) proposing a *landmark locally linear embedding* method to regulate the outliers in the predicted motion sequence to avoid robustness issues; 3) designing a computationally efficient NeRF-based motion-to-video renderer to achieves fast training and real-time inference. With these settings, GeneFace++ becomes the first NeRF-based method that achieves stable and real-time talking face generation with generalized audio-lip synchronization. Extensive experiments show that our method outperforms state-of-the-art baselines in terms of subjective and objective evaluation.⁴

^{*}Equal contribution.

[†]Work done while at an internship at ByteDance.

[‡]Corresponding author

⁴Video samples are available at <https://genefaceplusplus.github.io>

1 Introduction

Audio-driven talking face generation is a popular topic in the field of the digital person and metaverse[30, 18, 43], aiming at synthesizing an audio-lip-synchronized video of the target person given the input audio. Among recent works, some of them generate video frames with convolutional neural networks trained with adversarial objectives[38, 24, 35] and suffer from an unstable training process and 3D in-consistency artifacts, while others [12, 42] explore manipulating neural radiance field (NeRF)[25], a more stable and 3D-aware model, to render talking face videos. In general, modern talking face generation systems aim to achieve the following goals:

- *Generalized Audio-Lip Synchronization*: since people are sensitive to the slight misalignment between facial movements and speech audio [4], it is vital to maintain lip accuracy and temporal smoothness of the predicted facial motion. Considering the complex applications that may drive the system with out-of-domain (OOD) audio (such as cross-identity, cross-lingual, or singing audio), the model should also generalize well to various audios [42].
- *Good Video Quality*: the overall good video quality typically consists of high image fidelity, smooth transition between adjacent frames, and realistic 3D modeling of the talking avatar.
- *High System Efficiency*: to reduce the cost of computational resources and apply the model to real-world applications, the training process should be easy and the inference speed should be fast.

With the development of neural rendering techniques such as NeRF [25] and its variants [2], it is possible to build a high-fidelity 3D-aware talking face generation system with a several-minute-long video of the target person. Hence the second objective, i.e., *good video quality* has been partially achieved. However, as early NeRF-based methods [12][37] mainly train the model in an end-to-end manner, considering the limited amount of audio-lip pairs utilized in this process, it is hard for them to achieve the goal of *generalized audio-lip synchronization*. Recently, GeneFace[42] partially handles this problem by introducing a generalized audio-to-motion mapping learned from large-scale lip-reading datasets to predict accurate motion representations for the motion-conditioned NeRF-based renderer. Despite the community’s efforts to improve the practicality of NeRF-based talking face systems, there still exist several challenges that impede its usage in real-world applications:

- As for the *audio-lip synchronization*, it is still challenging to model the long audio sequence and generate time-consistent results. For instance, when a drawl or trill in the singing voice occurs, a single phoneme may last for more than 2 seconds (about 50 video frames), which requires the predicted lip motion to be consistent in the long term to seem perceptually natural.
- As for the *video quality*, it is hard to render diverse facial motions (such as an extra-big mouth) for a neural renderer that is only trained with consecutive frames from a few-minute-long video. GeneFace first notices this problem and learns a domain adaptative (DA) Postnet with adversarial training to map the predicted facial motion into the narrow input space of the renderer. However, it cannot guarantee the successful transformation of each frame and bad cases still occurs occasionally.
- As for the *system efficiency*, the slow training and inference speed caused by the expensive computation cost of vanilla NeRF has severely obstructed the usage in real-time applications.

In this work, we propose **GeneFace++** to address the three challenges in NeRF-based methods and achieve the goals of modern talking face generation. Specifically, 1) to improve long-term temporal consistency and naturalness of the predicted facial landmark sequence, we propose a *Pitch-Aware Audio-to-Motion* module. Specifically, we introduce the pitch contour as an auxiliary feature of the audio-to-motion mapping. We show that pitch could act as a helpful facial motion indicator to improve lip-sync quality. We also introduce a temporal smoothing loss to improve the overall temporal stability of the generated landmark sequence. 2) To improve the robustness of the system to diverse facial motions, we propose a manifold projection-based method named *Landmark Locally Linear Embedding* to post-process the predicted landmark, which solves a least square equation to reconstruct the input landmark with a linear combination of several nearest ground truth (GT) data points. It could map the landmark closer to the GT target person dataset and thus significantly reduce bad cases in rendering. 3) To improve the efficiency of the renderer, we propose an efficient dynamic NeRF named *Instant Motion-to-Video* module, which utilizes grid-based embedders to ease the training process and adopts deformable slicing surfaces to model the dynamics of the facial geometry conditioned on the 3D facial landmark. Compared with vanilla NeRF, the proposed renderer could be trained more efficiently and infer in real-time.

To summarize, GeneFace++ mitigates the difficulties in NeRF-based talking face video synthesis and achieves the goals of *generalized audio-lip synchronization*, *good video quality*, and *high system efficiency*. GeneFace++ could transform various OOD voices into accurate and time-consistent facial landmarks, and render high-fidelity 3D-aware human portraits in a stable and efficient manner. Extensive experiments show that GeneFace++ outperforms other state-of-the-art audio-driven talking face generation methods from the perspective of objective and subjective metrics. Ablation studies prove the necessity of each component.

2 Related Works

Talking face generation can be divided into two sequential steps: an audio-to-motion process that predicts facial motion given the input audio, and a motion-to-video process that renders the human portrait image given the input facial motion. We discuss related works about the lip-synchronized audio-to-motion transform and the human portrait rendering techniques, respectively.

2.1 Lip-Synchronized Audio-to-Motion

To achieve lip-synchronized motion prediction, there are two challenges faced by the community. The first challenge is the so-called one-to-many mapping problem, which means the same input audio may have several reasonable corresponding facial motions. Some early work [49, 47, 6] directly learn a deterministic model with a regression loss (e.g., L2 error), and suffer from over-smoothed lip results. Wav2Lip [30] first utilizes a discriminative sync expert to achieve a more sharp and accurate lip motion, which is followed by latter works[48, 45, 22, 19, 34]. MemFace[36] introduces memory retrieval in audio-to-motion to alleviate the one-to-many problem. The second challenge is to generate time consistent and stablized motion sequence given the long input audio. Some work [24] adopt auto-regressive structure to model the temporal sequence, but suffer from slow inference and error accumulation. Other works [41, 12] use parallel structure (such as 1D Convolution) with a sliding window, which partially address the shortage of auto-regressive methods. Recently, several works such as [7] and GeneFace[42] use feedforward structures (self-attention and convolution) to process the whole audio sequence in parallel. This framework enjoys high efficiency and capability to model the long term information, but is challenged in keeping temporal consistency and stability in the generated motion sequence.

2.2 Human Portrait Rendering

The modern techniques for dynamic human portrait synthesis could be categorized into three classes: 1) 2D-based, 2) 3D Morphable Model [29] (3DMM)-based, and 3) neural rendering-based. The earliest works typically belong to the 2D-based methods[39, 35, 30, 49, 46, 48], which either adopt GANs [10] or image-to-image translation [17] as the image renderer. Although achieves good image fidelity, these methods fail to generate pose-controllable videos due a lack of 3D geometry modeling. The 3DMM-based methods[41, 38, 44] inject the 3D prior knowledge by using 3DMM coefficients as auxiliary conditions, but using 3DMM as the intermediate is known to cause an information loss, which degrades the performance. Recently, the neural rendering-based methods[3, 9, 31, 15, 50] adopt NeRF[25] or its variants[2] to enjoy a realistic 3D modeling of the human head. AD-NeRF[12] is the first NeRF-based talking face method, which presents an end-to-end audio-to-video NeRF renderer to generate face images conditioned on audio features. Then some works explore to improve the sample efficiency [23] and achieve few-shot synthesis [33]. To achieve good lip synchronization, GeneFace[42] introduces an independent audio-to-motion module for the NeRF-based renderer. To improve the system efficiency, RAD-NeRF [37] introduces discrete learnable grids [26] in AD-NeRF, which accelerates training and inference.

As discussed above, with the rapid advances in the audio-to-motion and motion-to-video process, a modern talking face system that simultaneously achieves the goals of "generalized audio-lip synchronization", "good video quality" and "high system efficiency" is dawning. This observation motivates the design of GeneFace++. As for the audio-to-motion phase, we solve the long-term consistency problem with pitch information and a temporal smoothing loss. A manifold projection-based postprocessing method is also proposed to improve the robustness of the system. As for the motion-to-video phase, we utilize grid encoders and deformable slicing surfaces to achieve a high-quality and efficient motion-conditioned human portrait renderer.

3 Preliminaries: GeneFace

Our proposed GeneFace++ follows the two-stage paradigm in GeneFace [42]. Therefore, in this section, we introduce preliminary knowledge about the *audio-to-motion* and *motion-to-video* stage in GeneFace.

Audio-to-Motion At the *audio-to-motion* stage, GeneFace first learns a conditional variational auto-encoder (VAE) [20] model in a large-scale lip-reading dataset to achieve generalized and accurate facial landmark prediction given various audios. Specifically, the training loss of the VAE is:

$$\mathcal{L}_{\text{VAE}} = \mathbb{E}[\|\mathbf{l} - \hat{\mathbf{l}}\|_2^2 + KL(z|\hat{z}) + l_{\text{sync}}(\mathbf{a}, \hat{\mathbf{l}})], \quad s.t. \quad z \sim N(0, 1), \quad (1)$$

where \mathbf{a} is the input audio and \mathbf{l} is the corresponding GT facial landmark. *Enc* and *Dec* denote the encoder and decoder in VAE, respectively. $\hat{z} = \text{Enc}(\mathbf{l}, \mathbf{a})$ is the latent encoding of GT landmark and $\hat{\mathbf{l}} = \text{Dec}(\hat{z}, \mathbf{a})$ is the predicted facial landmark. *KL* denotes KL divergence and l_{sync} is a perceptual loss that measures the audio-visual synchronization, which is provided by a pretrained sync expert [30]. After training, the encoder is discarded and only the decoder is needed to predict the facial landmark given the input audio.

To overcome the significant domain gap between the large-scale lip-reading dataset and the target person video, GeneFace adopts adversarial domain adaptation to learn a domain adaptative (DA) Postnet that projects the predicted facial motion into the target person domain. The training loss of DA Postnet is:

$$\mathcal{L}_{\text{Postnet}} = \mathbb{E}[l_{\text{Adv}}(\bar{\mathbf{l}})] + l_{\text{sync}}(\mathbf{a}, \bar{\mathbf{l}}), \quad s.t. \quad \bar{\mathbf{l}} = PN(\hat{\mathbf{l}}), \quad (2)$$

where $\bar{\mathbf{l}}$ represents the refined landmark generated by the DA Postnet and *PN* is the Postnet (a shallow 1D convolutional neural network) to be trained. $\hat{\mathbf{l}} = \text{Dec}(z, \mathbf{a})$ is the facial landmark generated by the VAE decoder and z is the noise sampled from normalized gaussian distribution. l_{Adv} is the LSGAN-styled[?] adversarial loss, whose objective is to minimize the distance between the refined landmark distribution and the GT target person landmark set.

Once the training of the DA Postnet is done, the audio-to-motion module, which consists of a VAE decoder followed by a DA Postnet, could generate lip-sync and personalized facial landmarks given the input audio:

$$\bar{\mathbf{l}} = PN(\hat{\mathbf{l}}) = PN(\text{Dec}(z, \mathbf{a})). \quad (3)$$

Motion-to-Video As for the *motion-to-video* stage, GeneFace designs a landmark-conditioned dynamic NeRF network to render the human portrait given the input facial landmark. Specifically, it learns an implicit function F which can be formulated as follows:

$$F : (x, d, \mathbf{l}) \rightarrow (c, \sigma), \quad (4)$$

where x and d are positions and view direction. \mathbf{l} is the facial landmark that morphs the human head, which is implicitly modeled in the function F . c and σ denote the predicted color and density in the 3D radiance field. We can conveniently render an image from this radiance field via a differentiable volume rendering equation that aggregates the color c along the ray r :

$$C(r, \mathbf{l}) = \int_{t_n}^{t_f} \sigma(r(t), \mathbf{l}) \cdot c(r(t), \mathbf{l}, d) \cdot T(t) dt, \quad (5)$$

where C is the RGB value of the pixel that corresponds to the ray r emitted in the 3D space. t_n and t_f are the near bound and far bound of the ray. $r(t)$ is a shorthand of the position x and direction d at the sampled point t of the ray. $T(t) = \exp(-\int_{t_n}^t \sigma(r(\tau), \mathbf{l}) d\tau)$ is the accumulated transmittance along the ray from t_n to t . Now that the image could be rendered, the training objective of NeRF is to reduce the L2 error between the rendered and ground-truth images:

$$\mathcal{L}_{\text{NeRF}} = \mathbb{E}[\|C(r, \mathbf{l}) - C_g\|_2^2]. \quad (6)$$

During inference, by cascading through the *audio-to-motion* module and the *motion-to-video* module, GeneFace achieves to generate lip-sync and high-fidelity talking face video to various OOD audios.

4 GeneFace++

In this section, we introduce the architecture of GeneFace++, which aims to improve GeneFace to achieve more natural *audio-lip synchronization*, more robust *video quality*, and higher *system efficiency*. As shown in Fig. 1(a), GeneFace++ is composed of three parts: 1) a *pitch-aware audio-to-motion* module that transforms audio features into facial motion; 2) a *landmark locally linear embedding* method to post-process the predicted motion; and 3) an *instant motion-to-video* module that could render the final talking face video efficiently. We describe the designs and the training process of these three parts in detail in the following subsections. We provide detailed network structures in Appendix A.1.

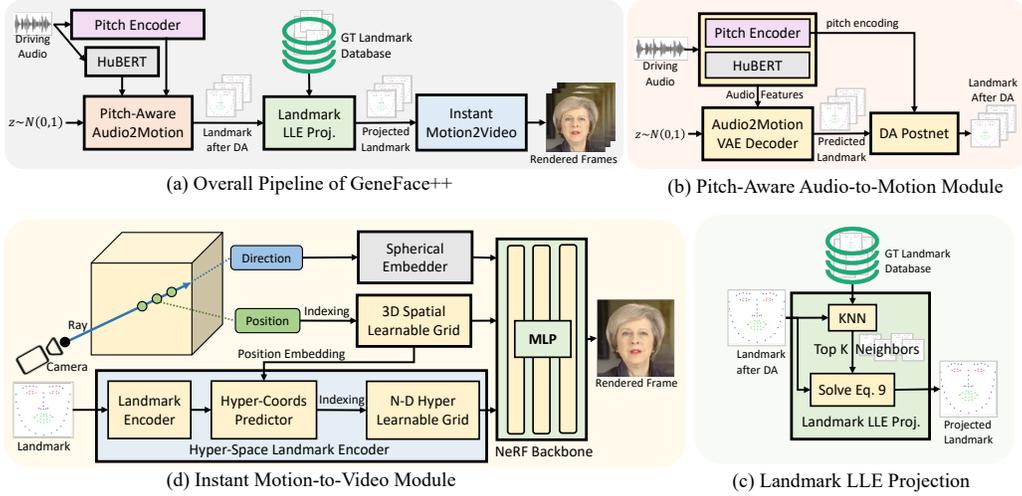


Figure 1: The inference process of GeneFace++. In subfigure (a), we show the overall three-stage pipeline. In subfigure (b), "DA Postnet" denotes the Domain Adaptative Postnet proposed in GeneFace. In subfigure (c), "KNN" denotes finding the K-nearest neighbors of the input landmark, and "Landmark LLE Proj." denotes Landmark Locally Linear Embedding Projection method proposed in Section 4.2. In subfigure (d), as for the *indexing* operation, we perform bi-linear interpolation to query the continuous coordinates in the discrete (spatial/hyper) grids.

4.1 Pitch-Aware Audio-to-Motion Transform

The motivation for considering pitch information in the audio-to-motion mapping is that pitch is known to highly correlated to facial expressions [40]. For instance, a high and steady pitch contour may correlate to a large and steady lip motion. Then we further found two advantages of using pitch in the audio-to-motion module: (1) We notice that previous NeRF-based methods [12][42] utilize phonetic posteriorgrams features [13][16] as the audio feature to ease the training and improve cross-lingual cross-identity generalizability. However, since PPGs is intended for ASR usage, it ignores the acoustic information in the waveform, which is necessary to achieve stable and expressive facial motion prediction. Under this circumstance, auxiliary acoustic features such as pitch contours are helpful to improve the expressiveness and temporary consistency of the predicted facial motion. (2) The second reason to introduce pitch information is the observation of the unstable performance of the DA Postnet in GeneFace: the Postnet is only provided with the predicted facial motion and is requested to project it into the target domain. It is hard for the Postnet to model this domain shift in its implicit space without any condition, which leads to unstable training and occasional bad cases. We suggest that the pitch information could act as a lightweight and helpful hint for the Postnet to better process the input facial landmark.

Pitch Encoder As shown in Fig. 2(c), we design a pitch encoder to effectively extract information from the pitch encoder. The key idea to designing the pitch encoder is to keep it lightweight and efficient. To be specific, we first discretize the continuous pitch (fundamental frequency) value into several discrete tokens to ensure temporary smoothness and ease the training of the pitch encoder.

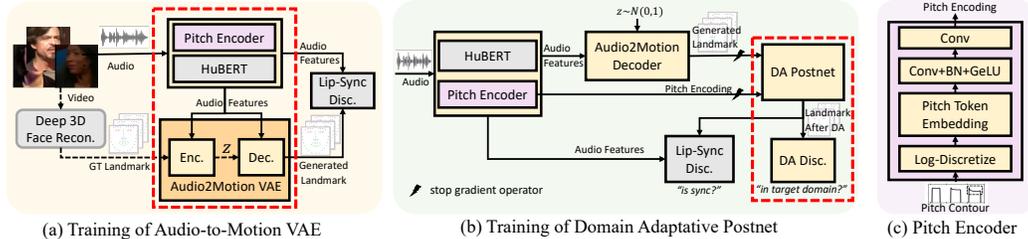


Figure 2: The training process of the Pitch-Aware Audio-to-Motion module. Learnable models are marked with red dotted rectangles and parameter-frozen models are colored in gray. In subfigure (a), Enc. and Dec. denote Encoder and Decoder in VAE, respectively. Disc. means Discriminator. In subfigure (b), the thunder-like symbol represents the "stop gradient" operator. In subfigure(c), "Log-Discretize" denotes the operation that quantizes the log-scale continuous pitch value into discrete tokens.

Note that we discretize the bins in log-scale to adhere to human perception. A group of pitch embedding corresponding to the discrete pitch tokens is learned from scratch during the training of the *pitch-aware VAE* (discussed latter). The pitch embedding is fed into a shallow convolutional network, which consists of a 1D convolutional layer with batch normalization and GeLU activation, and an extra 1D convolutional layer to produce the final pitch encoding.

Pitch-Aware VAE As shown in Fig. 2(a), we plug the pitch encoder as an auxiliary condition encoder into the audio-to-motion VAE. We train the pitch encoder and VAE on a large-scale lip-reading dataset following Equation 1. During inference, the pitch encoder and VAE decoder are used to predict the facial landmark:

$$\hat{\mathbf{I}} = Dec(z, \mathbf{h}, PE(\mathbf{p})), \quad s.t. \quad z \sim N(0, 1), \quad (7)$$

where \mathbf{h} and \mathbf{p} are the HuBERT feature and pitch value extracted from the input audio. PE represents the pitch encoder.

Pitch-Aware DA Postnet As shown in Fig. 2(b), the pretrained pitch encoder provides the auxiliary pitch encoding to the DA Postnet. We train the DA Postnet following Equation 2. Since the supervision signal of the adversarial training process is known to be unstable, we fix the parameters of the pitch encoder and VAE decoder to prevent deteriorated performance. During inference, as shown in Figure 1(a), the pitch-aware decoder and DA Postnet could generate personalized facial landmark $\bar{\mathbf{I}}$ given the input audio:

$$\bar{\mathbf{I}} = PN(\hat{\mathbf{I}}, PE(\mathbf{p})), \quad (8)$$

where $\hat{\mathbf{I}}$ is the facial landmark predicted by decoder in Equation 7.

Training We follow the major setting of GeneFace to train the VAE and Postnet with two modifications: (1) To remove the jitter in the predicted landmark, inspired by FACIAL[47], we encourage the VAE to optimize the velocity of the landmark sequence. Specifically, a temporal smoothing term is added to the training objective of VAE. (2) To stabilize the adversarial training, we adopt the gradient penalty in WGAN-GP[11] when updating the discriminator of Postnet. Due to spatial limitation, we illustrate the modified training loss in Appendix A.2.

4.2 Landmark Locally Linear Embedding

By introducing pitch information into the audio-to-motion module, we have improved the temporal consistency and naturalness of the predicted landmark. However, improving the quality of predicted landmarks is not adequate to achieve *good video quality*, as it requires the NeRF-based motion-to-video module to accurately render the human portrait corresponding to the assigned facial motion. This NeRF-based renderer, however, is typically learned from a very small dataset (a few-minute-long video), hence are expected to only work well on a narrow input space of facial landmark. When confronted with OOD landmarks, the renderer may produce inaccurate facial motion or even crashed rendering results. GeneFace utilize adversarial domain adaptation to train a Postnet to map all landmarks into the NeRF's narrow input space. However, due to the instability of adversarial training, it is not theoretically guaranteed to correctly project every frame into the target domain, and bad case occurs occasionally, which raises robust challenges to real-world applications.

In this context, as shown in Figure 1(c), we propose *landmark locally linear embedding* (Landmark LLE), a manifold projection-based post-processing method that guarantees each predicted landmark is successfully mapped into (the vicinity of) the input space of the landmark-conditioned renderer. In other words, with the help of Landmark LLE, each predicted landmark is dragged closer to the GT landmark set that is used to train the renderer. We follow the main idea of the classic locally linear embedding (LLE) algorithm [32] on the facial representation manifold: each facial landmark and its neighbors are locally-linear on the manifold. Motivated by the success of 3DMM, which could reconstruct arbitrary human face with a linear combination of around 144 template meshes, we assume that the facial landmark data points themselves are locally-linear, and there is no need to project them into a higher dimension manifold as vanilla LLE does. To be specific, given a predicted 3D facial landmark $\bar{\mathbf{l}} \in \mathbb{R}^{68 \times 3}$, the goal of Landmark LLE is to compute the reconstructed facial landmark $\bar{\mathbf{l}}' \in \mathbb{R}^{68 \times 3}$ on each dimension. As is illustrated in Figure 1(c), we first find the K nearest landmark of the input landmark $\bar{\mathbf{l}}$ in the GT landmark database $\mathcal{D} \in \mathbb{R}^{N \times 68 \times 3}$ by computing the Euclidean distance, where N is the number of landmark data points used to train the NeRF-based renderer. After obtaining the K nearest neighbors $\{\mathbf{l}_1, \dots, \mathbf{l}_K\} \subseteq \mathcal{D}$, we then seek a linear combination of these neighbors to reconstruct $\bar{\mathbf{l}}'$ by minimizing the reconstruction error $\|\bar{\mathbf{l}} - \bar{\mathbf{l}}'\|_2^2$, which could be formulated as the following least-squared optimization problem:

$$\min \|\bar{\mathbf{l}} - \sum_{k=1}^K w_k \cdot \mathbf{l}_k\|_2^2, \quad s.t. \sum_{k=1}^K w_k = 1, \quad (9)$$

where w_k is the barycentric weight of the k -th nearest landmark \mathbf{l}_k . The optimal weights $\mathbf{w}^* = \{w_1^*, \dots, w_K^*\} \in \mathbb{R}^K$ can be obtained by solving Equation 9. The hyper-parameter K is chosen as 20 in our experiment via grid searching. Then we could compute the reconstructed facial landmark $\bar{\mathbf{l}}' = \sum_{k=1}^K w_k^* \cdot \mathbf{l}_k$. Ideally, the reconstructed landmark $\bar{\mathbf{l}}'$ can be regarded as an in-domain landmark data point that also possesses the semantic facial motion (such as eye blinking, laughing) in the input landmark $\bar{\mathbf{l}}$, though with some information loss. In practice, during inference, we use a linear combination of originally predicted landmark $\bar{\mathbf{l}}$ and reconstructed landmark $\bar{\mathbf{l}}'$ as the final motion representation of the NeRF-based renderer:

$$\bar{\mathbf{l}}'' = \bar{\mathbf{l}}' \cdot \alpha + \bar{\mathbf{l}} \cdot (1 - \alpha) = (\sum_{k=1}^K w_k^* \cdot \mathbf{l}_k) \cdot \alpha + \bar{\mathbf{l}} \cdot (1 - \alpha), \quad (10)$$

where $\alpha \in [0, 1]$ is a temperature hyper-parameter that tunes the trade-off between image quality and facial motion expressiveness: Intuitively, a large α would drag the $\bar{\mathbf{l}}$ closer to the GT dataset distribution and hence lead to better image quality and less rendering bad cases, while a smaller α would keep more details in the originally predicted landmark, which denotes more diverse and expressive facial motion. Our experiment results in Table 4 show that the proposed method helps improve our system’s robustness to predicted OOD landmarks. We provide a pseudo code of Landmark LLE in Appendix A.3.

4.3 Instant Motion-to-Video Rendering

In the previous sections, we obtain an expressive, time-consistent, and robust audio-to-motion mapping through the *pitch-aware audio-to-motion* module and *landmark LLE* postprocessing method. Next, we design an *instant motion-to-video* module to efficiently render video frames conditioned on the predicted 3D landmarks, which is shown in Figure 1(d).

Grid-based NeRF Renderer Recent progress in grid-based NeRF [26] proposes to encode 3D spatial information with a learnable feature grid. Compared with vanilla NeRF which obtains the spatial features with dense MLP forwarding, this new paradigm could directly query the features in the continuous 3D space via linear interpolation in the discrete feature grid, which is more efficient in both the training and inference stage. Following [26], we utilize a learnable 3D grid to encode the queried position. Besides, an occupancy grid, which is a 3D grid that records the density value σ estimated by NeRF, is maintained during training to prune the ray marching path. Since the utilization of a 3D feature grid eases the burden of NeRF to model the continuous spatial space, a lightweight grid-based NeRF could achieve comparable rendering quality with a deeper vanilla NeRF.

Hyper-Space Landmark Encoder To utilize the facial landmark to morph the human head in NeRF, some works [27] adopt a conditional deformation field to warp the spatial points in the canonical space, which cannot model the non-rigid transform of the human head; GeneFace adopts a modulation-based method that directly concatenates the facial landmark and spatial features as the input to learn a

landmark-conditioned NeRF. However, this method requires deep MLP forwarding (along with the input spatial feature) to accurately morph the head geometry, which is computationally expensive and no longer feasible in a shallow grid-based NeRF. The motivation is to keep the accuracy of the input landmark to morph the head geometry while using a lightweight structure.

Inspired by Hyper-NeRF [28], as shown in Figure 1(d), we project the input facial landmark into an N-dimensional ambient coordinate conditioned on the grid-based spatial features, which allows an efficient fusion of spatial information and landmark condition. Once the ambient coordinate is obtained, instead of querying the landmark features with a dense MLP, we use an extra N-D learnable grid to improve efficiency. We empirically set $N = 3$ via grid search, as a trade-off of performance and efficiency. The queried spatial features and landmark features are concatenated and fed into the NeRF backbone (a shallow MLP) to generate the density and color. Specifically, the implicit function can be formulated as:

$$F : (f_x, f_l, d) \rightarrow c, \sigma \quad (11)$$

where f_x and f_l are the spatial/landmark feature queried from the grid, respectively. d is the view direction. Then we could obtain a rendered image following the volume rendering technique illustrated in Equation 5 and train the renderer with Equation 6.

5 Experiments

5.1 Experimental Setup

Data Preparation and Preprocessing To learn a generalized audio-to-motion module, we use a clean subset of LRS3-TED [1] to provide 190 hours of high-quality audio-motion pairs. To learn NeRF-based person-specific renderers, we adopt the dataset collected by [24] and [12], which consist of 5 videos of an average length of 6,000 frames in 25 FPS. During the data preprocessing phase, HuBERT [16] features and pitch contours are extracted from the audio track; head pose and 3D landmark are extracted from the video frames following [42]. To train the NeRF, the target person videos are cropped into 512x512 resolution and each frame is processed with the help of an automatic parsing method [21] for segmenting the head and torso part and extracting a clean background.

Compared Baselines We compare our GeneFace++ with several remarkable works: 1) Wav2Lip [30], which pretrains a sync-expert to improve the lip-synchronization performance; 2) MakeItTalk [49], which also utilizes 3D landmark as the action representation; 3) LiveSpeechPortrait (LSP) [24], which achieves photorealistic results at over 30 FPS; 4) AD-NeRF [12], which first utilize NeRF to achieve talking head generation. 5) RAD-NeRF [37], which adopts grid-based encoders in AD-NeRF to achieve real-time inference. 6) GeneFace [42], which achieves accurate lip-synchronization to OOD audios in NeRF-based talking face generation.

Model Configurations GeneFace++ consists of a *pitch-aware audio-to-motion* and an *instant motion-to-video* module. The pitch-aware audio-to-motion module follows the major settings from GeneFace, with an additional pitch encoder. In the instant motion-to-video module, the network consists of two 3D learnable grid encoders to store the spatial and landmark information, respectively. The spatial and landmark information is then concatenated and fed into a shallow MLP backbone. For a fair comparison, the hyper-parameters of all NeRF-based baselines are adjusted to be coherent with our model. We provide detailed hyper-parameters of GeneFace++ in Appendix B.1.

Training Details We train the GeneFace++ on 1 NVIDIA RTX 3090 GPU. For VAE and Postnet in the pitch-aware audio-to-motion module, it takes about 40k and 10k steps to converge (about 12 hours). For the instant motion-to-video renderer, we train each model for 400k iterations (200k for the head and 200k for the torso, respectively), which takes about 10 hours.

5.2 Quantitative Evaluation

Evaluation Metrics We employ the FID score [14] to measure image quality. We utilize the landmark distance (LMD)[5] and syncnet confidence score (Sync score) [30] to evaluate audio-lip synchronization. To evaluate the inference speed, we implement a frame-per-second (FPS) assessment on a single NVIDIA 3090Ti GPU.

Evaluation Results The results are shown in Table 1. We have the following observations. (1) GeneFace++ achieves good *audio-lip synchronization*, as it outperforms other baselines in terms of

Method	LMD↓	Sync ↑	PSNR ↑	FID ↓	FPS ↑
GT	0.000	6.735	0.00	0.000	N/A
Wav2Lip	3.902	7.618	29.16	71.963	11.95
MakeItTalk	4.838	4.383	27.94	50.087	22.03
LSP	4.431	5.094	29.58	33.993	25.35
AD-NeRF	4.038	4.685	30.89	33.340	0.069
RAD-NeRF	3.984	4.886	31.02	33.029	28.37
GeneFace	3.827	5.596	31.08	29.677	0.064
GeneFace++	3.776	6.114	31.22	29.147	23.55

Table 1: Quantitative evaluation with different methods. Best results are in **bold**.

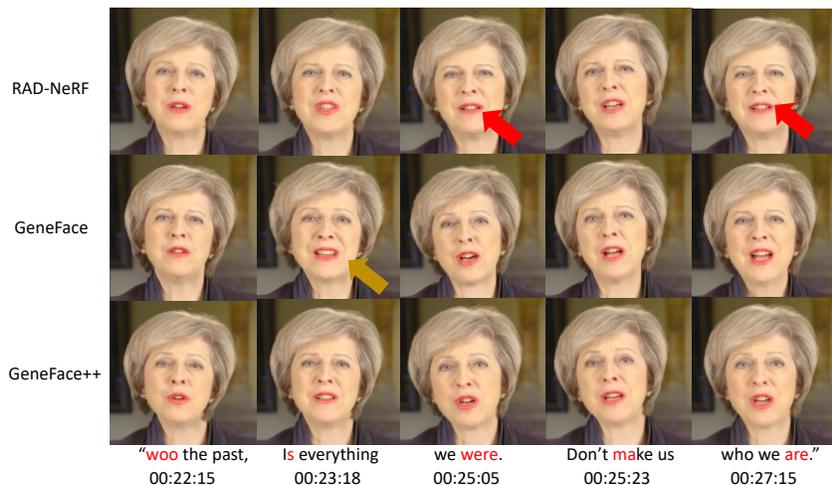


Figure 3: The comparison of generated key frame results. We show the speaking word and time step in the demo video. We mark the un-sync and bad rendering quality results with the red and brown arrows, respectively. **Please zoom in for better visualization.**

LMD and Sync score.⁵ (2) GeneFace++ achieves high *video quality*, as it has the best performance in terms of PSNR and FID. (3) As for the *system efficiency*, GeneFace++ could render talking face video in nearly real-time (about 23.55 FPS), which is a large acceleration from the GeneFace of 0.064 FPS. Both RAD-NeRF and GeneFace++ adopt the acceleration techniques from Instant-NGP to design the renderer, but RAD-NeRF achieves higher FPS with the comparable model scale of the renderer. This is due to the fact that GeneFace++ needs to execute an additional audio-to-motion and landmark LLE process to obtain lip-synchronized facial landmarks as the condition of the NeRF-based renderer. Although the audio-to-motion module and landmark LLE process are computationally cheap compared with the neural rendering process, it still brings slight latency in the overall system, which can be regarded as the sacrifices made to achieve a more accurate audio-lip synchronization.

5.3 Qualitative Evaluation

To make a qualitative comparison of each method, we provide a demo video⁶ in which each method is driven by a three-minute-long English song as a hard case. We recommend that readers watch this video for a better comparison. We also show the keyframes of this demo video in Figure 3. Due to space limitations, we only compare our GeneFace++ with the two most competitive baselines (GeneFace and RAD-NeRF) here and provide a comparison with all baselines in Appendix C.1. We observe that GeneFace++ manages to handle several problems in previous NeRF-based methods:

⁵An exception is Wav2Lip in terms of Sync Score, which is possibly due to the fact that Wav2Lip is jointly trained with SyncNet, it obtains a high sync score that is even higher than the ground truth video.

⁶The video URL is https://genefaceplusplus.github.io/GeneFace++/dream_it_possible.mp4

Methods	Wav2Lip	MakeItTalk	LSP	AD-NeRF	RAD-NeRF	GeneFace	GeneFace++
Audio-Lip Sync	3.67±0.25	3.12±0.16	3.42±0.21	2.97±0.18	3.05±0.26	3.62±0.18	3.82±0.18
Image Quality	2.84±0.25	2.67±0.32	3.55±0.23	3.31±0.20	3.43±0.24	3.39±0.21	3.58±0.18
Video Realness	3.07±0.27	2.39±0.30	3.34±0.28	3.15±0.23	3.31±0.24	3.39±0.22	3.54±0.26

Table 2: User study with different methods. The error bars are 95% confidence interval.

Setting	L2 Error ↓	TE	LMD↓	Sync ↑
GeneFace++	0.0108	0.0031	3.776	6.114
w/o P-VAE	0.0133	0.0042	3.783	5.975
w/o P-Postnet	0.0192	0.0047	3.810	5.732
GeneFace	0.0237	0.0058	3.827	5.596

Table 3: Ablation study results on pitch-aware audio-to-motion module. Temporal error is the L2 error on the velocity of the landmark sequence.

(1) while RAD-NeRF produces unsynchronized lip motion (red arrows in Figure 3) due to the weak generalizability to OOD audios, GeneFace++ achieves accurate lip motion. (2) GeneFace++ handles the occasionally occurred rendering bad cases in GeneFace (brown arrows in Figure.3), which are caused by OOD landmarks. (3) We also show that GeneFace++ improves the time consistency in the demo video.

User Study We conduct user studies to test the quality of audio-driven portraits. Specifically, we sample 5 audio clips from English, French, and Italian (including a three-minute-long English song as a hard case) for all methods to generate the videos, and then involve 20 attendees for user studies. We adopt the Mean Opinion score (MOS) rating protocol for evaluation, which is scaled from 1 to 5. The attendees are required to rate the videos based on three aspects: (1) *audio-lip synchronization*; (2) *image quality*; (3) *video realness*, which mainly measures the time consistency and 3D realness.

We compute the average score for each method, and the results are shown in Table 2. We have the following observations: (1) Our GeneFace++ achieves the best *audio-lip synchronization*, which shows the performance of our proposed pitch-aware audio-to-motion module. (2) As for the *image quality*, with the delicately designed instant motion-to-video module, we found GeneFace++ achieves better image quality than previous person-specific NeRF-based (AD-NeRF, RAD-NeRF, and GeneFace) and GAN-based (LSP) renderers. We also observe that the one-shot talking face generation methods (Wav2Lip and MakeItTalk) obtains low scores in terms of image quality. (3) GeneFace also achieves the best *video realness* among the tested methods. We attribute the pitch-aware audio-to-motion module and landmark LLE method to good video realness, which could predict more accurate and time-consistent landmarks and automatically refine the outliers.

5.4 Ablation Studies

In this section, we perform ablation studies to prove the necessity of each component in GeneFace++.

Pitch-Aware Audio-to-Motion We test three settings on the audio-to-motion module: (1) removing pitch information in the VAE (w/o P-VAE); (2) removing pitch information in the Postnet (w/o P-Postnet); and (3) removing pitch information in VAE and Postnet, which is equivalent to the vanilla audio-to-motion module in GeneFace. We use the L2 error of the predicted landmark sequence to measure the lip accuracy and adopt a temporal error (which is the L2 error of the velocity of the landmark sequence) to measure the temporal consistency. We abbreviate temporal error as TE. The LMD and Sync score of the downstream rendered frames is also measured. The results are shown in Table 3. We observe that removing pitch information leads to a significant degradation in lip accuracy and time consistency.

Landmark LLE We test three settings of Landmark LLE by tuning the hyper-parameter α defined in Equation 10 from $\{0, 0.5, 1.0\}$. Note that $\alpha = 0$ equals to not using Landmark LLE to post-process the predicted landmark and $\alpha = 1.0$ means we use the LLE-reconstructed landmark l' as the final

Setting	LMD ↓	Sync ↑	BadCase(%) ↓
$\alpha = 0.0$	3.708	6.153	0.513
$\alpha = 0.5$	3.776	6.114	0.164
$\alpha = 1.0$	3.825	6.039	0.102

Table 4: Ablation study results on landmark LLE method. α is the hyper-parameter defined in Equation 10. BadCase% denotes the percentage of bad cases in the generated frames.

Setting	PSNR(Face) ↑	PSNR(Lip) ↑	FPS ↑
dim=2	31.13	29.87	26.91
dim=3	31.22	30.19	23.55
dim=4	31.20	30.21	18.31

Table 5: Ablation study results on number of hyper coordinate dimensions of landmark.

input of the renderer. We are interested in two questions: (1) whether the additional postprocessing process would degrade the lip synchronization of the predicted landmark; (2) whether the proposed Landmark LLE method could effectively handle the bad cases caused by OOD landmarks. Therefore, we use LMD and Sync score to reflect the lip synchronization and compute the percentage of the bad case caused by OOD landmark in the generated videos, represented as BadCase%. The results are shown in Table 4. We observe that adopting Landmark LLE significantly decreases the bad case rate, while slightly degrading the lip synchronization. As a trade-off, we use $\alpha = 0.5$ in our experiments.

To further investigate the efficacy of Landmark LLE, we utilize T-SNE to visualize the landmarks at different stages in Appendix C.2. We can see that Landmark LLE manages to drag outliers into the target person domain, hence improving the stability of the rendering results.

Instant-Motion-to-Video We ablate the number of hyper coordinate dimensions of the landmark encoder in the instant motion-to-video module. As a small coordinate dimension may limit the model’s capacity and harms the rendering quality while a large coordinate requires more computation costs, we are interested in a suitable hyper coordinate dimension to efficiently model the landmark information. To this end, we test the number of dimensions within $\{2, 3, 4\}$. We take PSNR to measure the rendering quality and report the FPS to compare the inference speed. The results are shown in Table 5. An interesting finding is that increasing the hyper coordinate dimension from 2 to 3 significantly improves the rendering quality on the Lip area. We suggest it is because the lip is a high-frequency dynamic part in a talking human head and could benefit from a larger model capacity. Therefore, we use a 3D hyper coordinate as a trade-off between rendering quality and inference speed.

6 Conclusion

In this paper, we propose GeneFace++, which aims to achieve three goals in modern talking face generation: *generalized audio-lip synchronization*, *good video quality*, and *high system efficiency*. GeneFace++ takes a big step forward from GeneFace. A pitch-aware audio-to-motion module is proposed to predict the facial landmark of generalized lip synchronization and high time consistency. A Landmark LLE method is introduced to automatically refine the predicted landmark and avoid the potential bad cases of rendering. A delicate instant motion-to-video renderer is designed to generate high-quality video efficiently. Extensive experiments show that our method achieves the three goals of the modern talking face system and outperforms existing methods. Due to space limitations, we discuss the limitations and future works in Appendix D.

References

- [1] Triantafyllos Afouras, Joon Son Chung, and Andrew Zisserman. Lrs3-ted: a large-scale dataset for visual speech recognition. *arXiv preprint arXiv:1809.00496*, 2018.

- [2] Eric R. Chan, Connor Z. Lin, Matthew A. Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J. Guibas, Jonathan Tremblay, Sameh Khamis, Tero Karras, and Gordon Wetzstein. Efficient geometry-aware 3d generative adversarial networks. In *CVPR*, pages 16123–16133, June 2022.
- [3] Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *CVPR*, 2021.
- [4] Lele Chen, Guofeng Cui, Ziyi Kou, Haitian Zheng, and Chenliang Xu. What comprises a good talking-head video generation?: A survey and benchmark. *arXiv preprint arXiv:2005.03201*, 2020.
- [5] Lele Chen, Zhiheng Li, Ross K Maddox, Zhiyao Duan, and Chenliang Xu. Lip movements generation at a glance. In *ECCV*, pages 520–535, 2018.
- [6] Lele Chen, Ross K Maddox, Zhiyao Duan, and Chenliang Xu. Hierarchical cross-modal talking face generation with dynamic pixel-wise loss. In *CVPR*, pages 7832–7841, 2019.
- [7] Liyang Chen, Zhiyong Wu, Jun Ling, Runnan Li, Xu Tan, and Sheng Zhao. Transformer-s2a: Robust and efficient speech-to-animation. In *ICASSP*, 2022.
- [8] Shu-Yu Chen, Wanchao Su, Lin Gao, Shihong Xia, and Hongbo Fu. Deepfacedrawing: Deep generation of face images from sketches. *ACM Transactions on Graphics (TOG)*, 39(4):72–1, 2020.
- [9] Guy Gafni, Justus Thies, Michael Zollhofer, and Matthias Niessner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In *CVPR*, pages 8649–8658, June 2021.
- [10] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.
- [11] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. *NIPS*, 2017.
- [12] Yudong Guo, Keyu Chen, Sen Liang, Yong-Jin Liu, Hujun Bao, and Juyong Zhang. Ad-nerf: Audio driven neural radiance fields for talking head synthesis. In *ICCV*, pages 5784–5794, 2021.
- [13] Awni Hannun, Carl Case, Jared Casper, Bryan Catanzaro, Greg Diamos, Erich Elsen, Ryan Prenger, Sanjeev Satheesh, Shubho Sengupta, Adam Coates, et al. Deep speech: Scaling up end-to-end speech recognition. *arXiv preprint arXiv:1412.5567*, 2014.
- [14] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.
- [15] Yang Hong, Bo Peng, Haiyao Xiao, Ligang Liu, and Juyong Zhang. Headnerf: A real-time nerf-based parametric head model. In *CVPR*, pages 20374–20384, June 2022.
- [16] Wei-Ning Hsu, Yao-Hung Hubert Tsai, Benjamin Bolte, Ruslan Salakhutdinov, and Abdelrahman Mohamed. Hubert: How much can a bad teacher benefit asr pre-training? In *ICASSP*, 2021.
- [17] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, 2017.
- [18] Amir Jamaludin, Joon Son Chung, and Andrew Zisserman. You said that?: Synthesising talking faces from audio. *International Journal of Computer Vision*, 127(11):1767–1779, 2019.
- [19] Xinya Ji, Hang Zhou, Kaisiyuan Wang, Qianyi Wu, Wayne Wu, Feng Xu, and Xun Cao. Eamm: One-shot emotional talking face via audio-based emotion-aware motion model. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–10, 2022.

- [20] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [21] Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive facial image manipulation. In *CVPR*, pages 5549–5558, 2020.
- [22] 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 *CVPR*, 2022.
- [23] Xian Liu, Yinghao Xu, Qianyi Wu, Hang Zhou, Wayne Wu, and Bolei Zhou. Semantic-aware implicit neural audio-driven video portrait generation. *arXiv preprint arXiv:2201.07786*, 2022.
- [24] Yuanxun Lu, Jinxiang Chai, and Xun Cao. Live speech portraits: real-time photorealistic talking-head animation. *ACM Transactions on Graphics*, 40(6):1–17, 2021.
- [25] 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 *ECCV*, pages 405–421. Springer, 2020.
- [26] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022.
- [27] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *ICCV*, 2021.
- [28] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M Seitz. Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. *arXiv preprint arXiv:2106.13228*, 2021.
- [29] Pascal Paysan, Reinhard Knothe, Brian Amberg, Sami Romdhani, and Thomas Vetter. A 3d face model for pose and illumination invariant face recognition. In *2009 sixth IEEE international conference on advanced video and signal based surveillance*, pages 296–301. Ieee, 2009.
- [30] KR Prajwal, Rudrabha Mukhopadhyay, Vinay P Namboodiri, and CV Jawahar. A lip sync expert is all you need for speech to lip generation in the wild. In *ACM MM*, pages 484–492, 2020.
- [31] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *CVPR*, pages 10318–10327, 2021.
- [32] Sam T Roweis and Lawrence K Saul. Nonlinear dimensionality reduction by locally linear embedding. *science*, 290(5500):2323–2326, 2000.
- [33] 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 *ECCV*, 2022.
- [34] Yasheng Sun, Hang Zhou, Ziwei Liu, and Hideki Koike. Speech2talking-face: Inferring and driving a face with synchronized audio-visual representation. In *IJCAI*, 2021.
- [35] Supasorn Suwajanakorn, Steven M Seitz, and Ira Kemelmacher-Shlizerman. Synthesizing obama: learning lip sync from audio. *ACM Transactions on Graphics (ToG)*, 36(4):1–13, 2017.
- [36] 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. *arXiv preprint arXiv:2212.05005*, 2022.
- [37] Jiaxiang Tang, Kaisiyuan Wang, Hang Zhou, Xiaokang Chen, Dongliang He, Tianshu Hu, Jingtuo Liu, Gang Zeng, and Jingdong Wang. Real-time neural radiance talking portrait synthesis via audio-spatial decomposition. *arXiv preprint arXiv:2211.12368*, 2022.

- [38] Samuli Laine Antti Herva Tero Karras, Timo Aila and Jaakko Lehtinen. Audio-driven facial animation by joint end-to-end learning of pose and emotion. *ACM Transactions on Graphics*, 36(4), 2017.
- [39] Justus Thies, Mohamed Elgharib, Ayush Tewari, Christian Theobalt, and Matthias Nießner. Neural voice puppetry: Audio-driven facial reenactment. In *ECCV*, pages 716–731. Springer, 2020.
- [40] William Forde Thompson, Frank A Russo, and Steven R Livingstone. Facial expressions of singers influence perceived pitch relations. *Psychonomic bulletin & review*, 17(3):317–322, 2010.
- [41] Haozhe Wu, Jia Jia, Haoyu Wang, Yishun Dou, Chao Duan, and Qingshan Deng. Imitating arbitrary talking style for realistic audio-driven talking face synthesis. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 1478–1486, 2021.
- [42] Zhenhui Ye, Ziyue Jiang, Yi Ren, Jinglin Liu, JinZheng He, and Zhou Zhao. Geneface: Generalized and high-fidelity audio-driven 3d talking face synthesis. In *ICLR*, 2023.
- [43] Zhenhui Ye, Zhou Zhao, Yi Ren, and Fei Wu. Syntaspeech: syntax-aware generative adversarial text-to-speech. *arXiv preprint arXiv:2204.11792*, 2022.
- [44] Ran Yi, Zipeng Ye, Juyong Zhang, Hujun Bao, and Yong-Jin Liu. Audio-driven talking face video generation with learning-based personalized head pose. *arXiv preprint arXiv:2002.10137*, 2020.
- [45] Fei Yin, Yong Zhang, Xiaodong Cun, Mingdeng Cao, Yanbo Fan, Xuan Wang, Qingyan Bai, Baoyuan Wu, Jue Wang, and Yujiu Yang. Styleheat: One-shot high-resolution editable talking face generation via pre-trained stylegan. In *ECCV*, 2022.
- [46] Lingyun Yu, Jun Yu, Mengyan Li, and Qiang Ling. Multimodal inputs driven talking face generation with spatial-temporal dependency. *IEEE Transactions on Circuits and Systems for Video Technology*, 31(1):203–216, 2020.
- [47] Chenxu Zhang, Yifan Zhao, Yifei Huang, Ming Zeng, Saifeng Ni, Madhukar Budagavi, and Xiaohu Guo. Facial: Synthesizing dynamic talking face with implicit attribute learning. In *ICCV*, 2021.
- [48] Hang Zhou, Yu Liu, Ziwei Liu, Ping Luo, and Xiaogang Wang. Talking face generation by adversarially disentangled audio-visual representation. In *AAAI*, 2019.
- [49] Yang Zhou, Xintong Han, Eli Shechtman, Jose Echevarria, Evangelos Kalogerakis, and Dingzeyu Li. Makelttalk: speaker-aware talking-head animation. *ACM Transactions on Graphics (TOG)*, 39(6):1–15, 2020.
- [50] Yiyu Zhuang, Hao Zhu, Xusen Sun, and Xun Cao. Mofanerf: Morphable facial neural radiance field. In *ECCV*, 2022.

A Details of Models

A.1 Detailed Network Structure

We provide the detailed network structure of audio-to-motion VAE and Postnet in Figure 4 and Figure 5, respectively. As shown in 4(b), in practice, we additionally train a flow-based model to predict the latent code during inference, which is adhere to GeneFace.

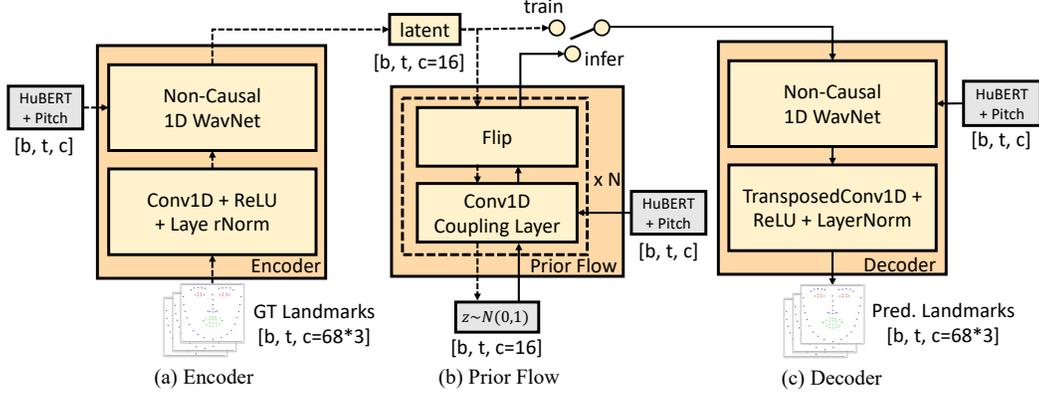


Figure 4: The detailed structure of VAE in the Pitch-Aware Audio-to-Motion module. In subfigure (a), (b), and (c), we show the structure of encoder, prior flow, and the decoder, respectively. Dotted arrows are only operated at the training phase. "Pred." is a shorthand of "predicted".

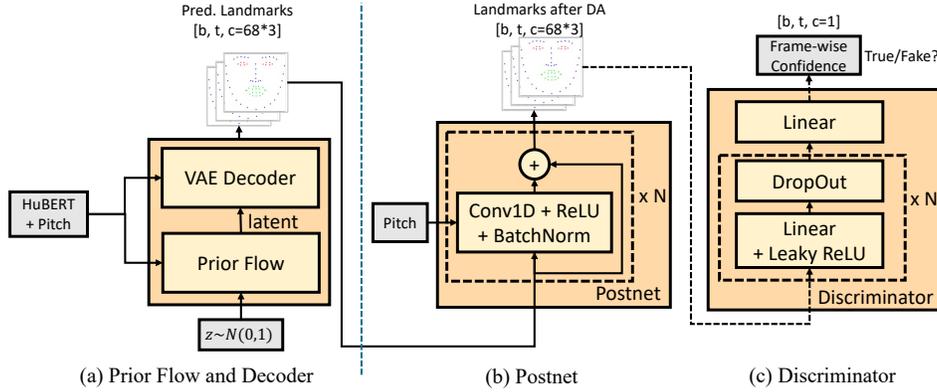


Figure 5: The detailed structure of Postnet and its discriminator in the Pitch-Aware Audio-to-Motion module. In subfigure (a), Prior Flow and Decoder are pretrained and fixed to predict raw landmark sequence. We separated them with a blue dotted line. In subfigure (b) and (c) we show the structure of Postnet and its discriminator, respectively. Dotted arrows are only operated at the training phase. "Pred." is a shorthand of "predicted".

A.2 Training Losses of Audio-to-Motion Module

Pitch-Aware VAE We add a temporal term l_{temp} in the training of VAE and Postnet, which is the L2 error of the velocity of the landmark sequence and could be represented as:

$$l_{temp}(v, \hat{v}) = \|v - \hat{v}\|_2^2 \quad (12)$$

where v and \hat{v} is the velocity of GT landmark l and predicted landmark \hat{l} , respectively. Now that the temporal term is defined, the training loss of Pitch-Aware VAE could be formulated on the top of Equation 1:

$$\mathcal{L}_{P-VAE} = \mathbb{E}[\|l - \hat{l}\|_2^2 + KL(z|\hat{z}) + l_{sync}(a, \hat{l}) + l_{temp}(v, \hat{v})]. \quad (13)$$

Pitch-Aware Postnet As for the training of Pitch-Aware Postnet, the Postnet follows the Equation 2 in GeneFace, and the discriminator is enhanced by introducing a gradient penalty term:

$$l_{gp}(l) = \nabla_l Disc(l), \quad (14)$$

where l is the GT landmark. This term encourages the discriminator to give a similar confidence score to a similar input landmark, which successfully prevents overfitting in the adversarial training. Now that the gradient penalty term is obtained, we formulate the loss of discriminator as:

$$\mathcal{L}_{\text{Disc}} = \mathbb{E}[l_{adv}(l, \hat{l}) + l_{gp}(l)], \quad (15)$$

where l_{adv} is the discriminator loss of LS-GAN.

A.3 Pseudo Code of Landmark LLE Method

We provide the pseudo code of the proposed Landmark LLE method in Algorithm 1.

Algorithm 1 Landmark Locally Linear Embedding Method

- 1: **Input:** GT landmark set \mathcal{D} ; Predicted landmark \mathbf{l} .
 - 2: **Output:** Projected landmark l' .
 - 3: Compute K nearest neighbors of \mathbf{l} in the GT landmark set \mathbf{L} .
 - 4: Solve the least squared problem defined in Equation 9, obtain the optimal weights of K neighbors $\mathbf{w}^* = \{w_1^*, \dots, w_K^*\}$.
 - 5: Compute the projected landmark following Equation 10.
-

B Detailed Experimental Settings

B.1 Model Configurations

We list the hyper-parameters of GeneFace++ in Table 6.

C Additional Experiments

C.1 Qualitative Results with All Baselines

We show the rendered keyframes of different talking face methods in Figure 6. We recommend the readers watch the demo video for a better comparison. The video URL is https://genefaceplusplus.github.io/GeneFace++/dream_it_possible.mp4

C.2 T-SNE Visualization of 3DMM Landmark

To validate the effectiveness of the proposed Landmark LLE method, we adopt T-SNE to visualize the predicted landmark at different stages of GeneFace++. The result is shown in Figure 7. The brown points denotes the ground truth landmarks in the target person video. Ideally the predicted landmarks should within this distribution so that the motion-to-video module could render images of good quality. The red points are landmarks predicted by the VAE, and the green points are the prediction refined by the Postnet. We can see that Postnet manages to project the majority of predicted landmark into the target person domain, yet there still exist some outliers out of the distribution (marked with a red circle) and will raise rendering bad cases. Then the blue points are the predicted landmark further postprocessed by the Landmark LLE method. We can see that LLE successfully drag the outliers into the target person video, hence significantly improves the stability of the rendering result.

D Limitation and Future Works

To validate the effectiveness of the proposed Landmark LLE method, we adopt T-SNE to visualize the predicted landmark at different stages of GeneFace++. The result is shown in Figure 7. The brown points denote the ground truth landmarks in the target person’s video. Ideally, the predicted landmarks should be within this distribution so that the motion-to-video module could render images of good quality. The red points are landmarks predicted by the VAE, and the green points are the prediction refined by the Postnet. We can see that Postnet manages to project the majority of predicted landmarks into the target person domain, yet there still exist some outliers out of the distribution

Table 6: Hyper-parameter list

	Hyper-parameter	Value
Pitch Encoder	Number of Pitch Bins	300
	Pitch Embedding Channel Size	64
	Number of Conv1D Layers	2
	Pitch Encoder Conv1D Kernel	3
	Pitch Encoder Conv1D Channel Size	64
Audio-to-Motion VAE	Encoder Layers	8
	Decoder Layers	4
	Encoder/Decoder Conv1D Kernel	5
	Encoder/Decoder Conv1D Channel Size	192
	Latent Size	16
	Prior Flow Layers	4
	Prior Flow Conv1D Kernel	3
	Prior Flow Conv1D Channel Size	64
	Sync-expert Layers	14
	Sync-expert Channel Size	512
Post-net and its DA-Discriminator	Post-net Layers	8
	Post-net Conv1D Kernel	3
	Post-net Conv1D Channel Size	256
	Discrimnator Layers	5
	Discrimnator Linear Hidden Size	256
	Discrimnator Dropout Rate	0.25
Landmark LLE	Number of Neighbors (K)	20
	Weights of LLE results to construct the final landmark (α)	0.5
Instant Motion-to-Video Renderer	Head/Torso NeRF Layers	5
	Head/Torso NeRF Hidden Size	64
	Head NeRF Spatial Grid Dimension	3
	Head NeRF Spatial Grid Channel Size per Resolution Revel	2
	Head NeRF Spatial Grid Number of Resolution Revel	16
	Landmark Encoder Layers	6
	Landmark Encoder Hidden Size	64
	Landmark Hyper Grid Dimension	3
	Landmark Hyper Grid Channel Size	2

(marked with a red cycle) and will raise rendering bad cases. Then the blue points are the predicted further post-processed by the Landmark LLE method. We can see that LLE successfully drags the outliers into the target person’s video, hence significantly improving the stability of the rendering result.

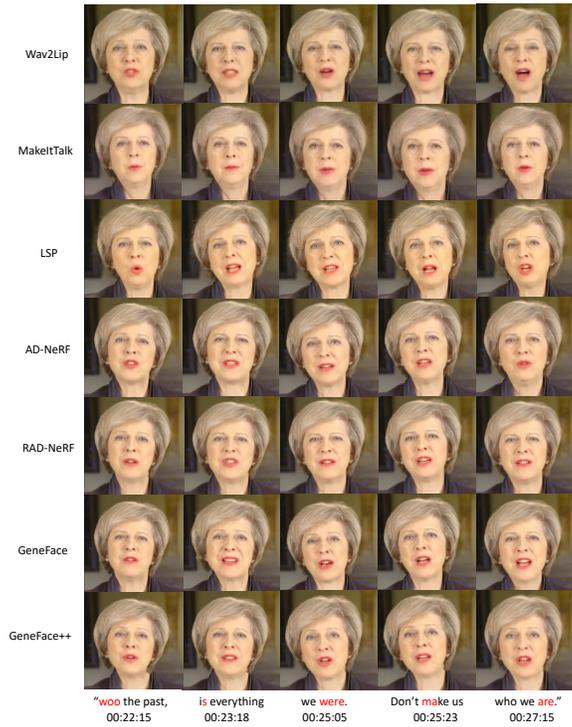


Figure 6: The comparison of generated key frame results. We show the speaking word and time step in the demo video. **Please zoom in for better visualization.**

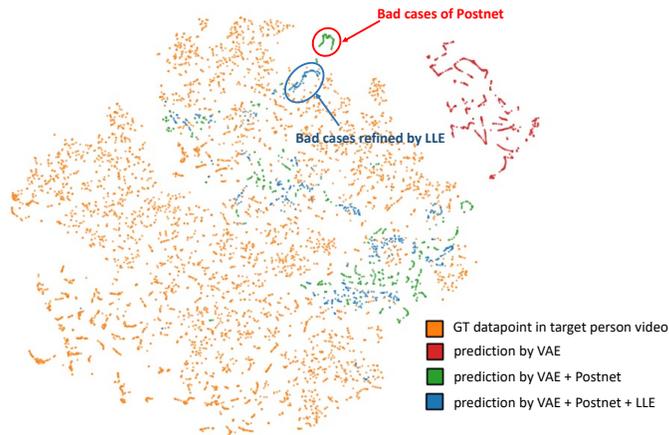


Figure 7: The T-SNE visualization of 3DMM landmarks in the different stages of GeneFace++.