# 3DiFACE: Diffusion-based Speech-driven 3D Facial Animation and Editing

Balamurugan Thambiraja[1]    Sadegh Aliakbarian[3]    Darren Cosker[3]    Justus Thies[1,2]

[1] Max Planck Institute for Intelligent Systems, Tübingen, Germany
[2] Technical University of Darmstadt    [3] Microsoft Mixed Reality & AI Lab
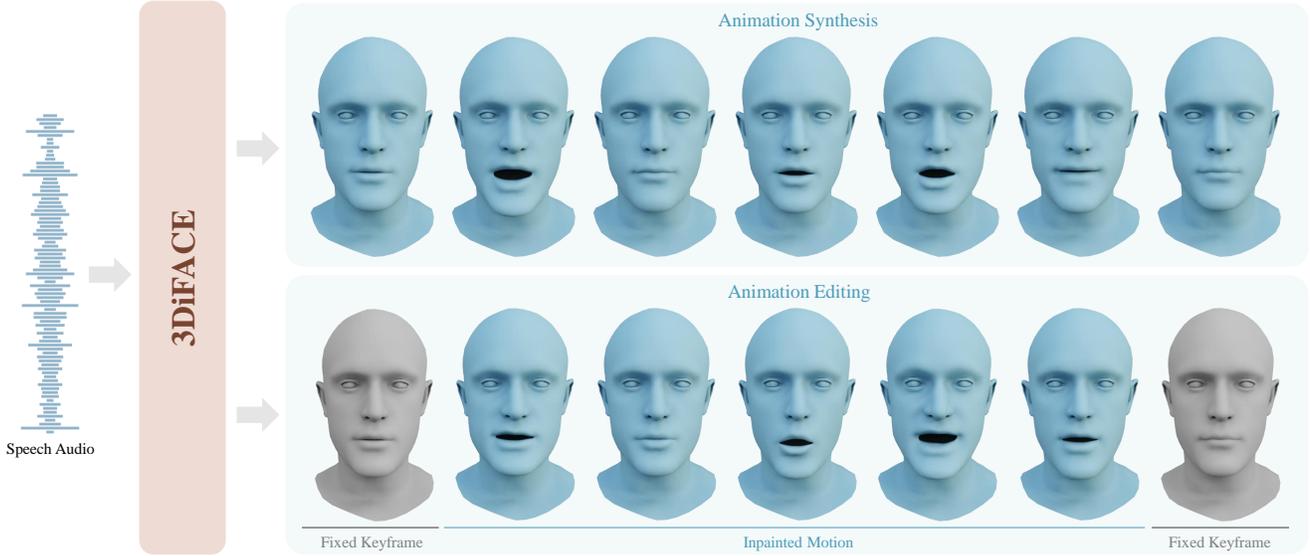https://balamuruganthambiraja.github.io/3DiFACE

Figure 1. *3DiFACE* is a novel diffusion-based method for speech-driven 3D facial animation. Given an audio sequence, our method can generate motion sequences with lip-sync and stochasticity. Additionally, 3DiFACE can be used for audio-consistent motion editing.

## Abstract

*We present 3DiFACE, a novel method for personalized speech-driven 3D facial animation and editing. While existing methods deterministically predict facial animations from speech, they overlook the inherent one-to-many relationship between speech and facial expressions, i.e., there are multiple reasonable facial expression animations matching an audio input. It is especially important in content creation to be able to modify generated motion or to specify keyframes. To enable stochasticity as well as motion editing, we propose a lightweight audio-conditioned diffusion model for 3D facial motion. This diffusion model can be trained on a small 3D motion dataset, maintaining expressive lip motion output. In addition, it can be fine-tuned for specific subjects, requiring only a short video of the person. Through quantitative and qualitative evaluations, we show that our method outperforms existing state-of-the-art techniques and yields speech-driven animations with greater fidelity and diversity.*

## 1. Introduction

Generating faithful 3D animations is crucial for realistic and immersive digital experience in games, movies, and other human-centric entertainment applications. Earlier works on 3D facial animation focused on animating faces based on procedural rules that map audio features to facial animation parameters [10, 17]. With advancements in machine learning, data-driven methods have been widely used to generate animations conditioned on audio input [11, 21, 42, 47, 69, 75]. Despite the recent advances in 3D facial animation, most of the existing methods learn a deterministic mapping between audio and facial animation, overlooking the inherent one-to-many relationship—one audio signal can match many different facial expression animations. This limits the diversity of synthesized animations.

1

In addition, most of the existing works focus solely on motion synthesis, but do not address motion editing, e.g. inbetweening facial motion between two keyframes. However, motion editing is an equally interesting and challenging problem with high relevance for practical content creation applications in the gaming and movie industry. In this work, we address this gap and propose a diffusion-based architecture to perform speech-driven motion synthesis and editing. In doing so, we face two challenges: (i) Diffusion models are known to require large training sets [48], yet the size of existing high-quality speech-to-3D-animation datasets is limited. (ii) Facial movements are highly person-specific. This requires the learning of person-specific speaking styles into the synthesis and editing pipeline. Especially for facial motion editing applications, without person-specific speaking-style, we observe an abrupt change in speaking style between the edited and unedited motion, resulting in unrealistic animations.

Recent works such as EMOTE [14] and Diff-PoseTalk [57] employ head trackers to annotate large-scale-video datasets with pseudo ground truth data and train their models on the resulting dataset. While effectively solving data scarcity, the synthesis fidelity is limited by the quality of the trackers and it is inferior to models that were trained on smaller datasets with higher quality [12]. Imitator [62] is a method that addresses speaker style adaptation, by learning a person-specific motion decoder, however, the motion is deterministic and does not allow for motion editing.

In this work, we propose a diffusion-based pipeline for speech-driven facial animation which can be efficiently trained on small datasets and enabling us to optimize person-specific speaking styles. In contrast to previous works [14, 20, 62, 68] that rely on transformer-based architectures [64], we propose an audio-conditioned diffusion model using a 1D convolutional backbone. In addition, we employ a concatenation-based convolution block over the widely used attention mechanism to inject audio conditioning signals into the 1D convolutional model and we propose a window-based training scheme. Based on our lightweight architecture, we take inspiration from Imitator [62] and propose a person-specific fine-tuning technique, that only requires a short video of the target actor. The diffusion architecture as well as the personalization allows us to generate and edit animations with person-specific speaking style. Since we regress vertex displacements directly, our method can capture more subtle speaking style variations than approaches that regress coefficients of a parametric head model [36]. Through qualitative and quantitative studies, we show that our proposed pipeline outperforms the state of the art in synthesizing and editing realistic 3D facial animations conditioned on input audio, while preserving the speaking style of a target subject. In summary, our contributions are:

- We introduce a novel diffusion-based method for speech-driven 3D facial animation that can be efficiently trained on small-scale high-quality datasets to generate diverse animations from a single audio source.
- Our method captures person-specific speaking styles from short reference videos, enabling personalized facial animations.
- We demonstrate exciting applications in editing facial motions, such as seamless motion interpolation, keyframing, and unconditional facial motion synthesis.

## 2. Related Work

Numerous prior investigations have delved into the realm of speech-driven generation, with a predominant emphasis on synthesizing 2D talking head videos. However, especially for 3D content creation in games, movies, and immersive telepresence, speech-driven 3D facial animation is of high interest to the research community. Our work generates 3D facial animations using a denoising diffusion probabilistic model, therefore we review motion diffusion work.

**Talking Head Videos** Approaches for talking head video generation can be mainly categorized into two groups: directly generating RGB videos from speech on the one hand, and utilizing a 3D Morphable Model (3DMM) for guided 2D or 3D rendering on the other. Suwajanakorn et al. [58] proposed a method belonging to the first category. It relies on recurrent models (LSTMs) to predict person-specific 2D lip landmarks to guide 2D image generation. Chung et al. [9] introduce a real-time approach for generating an RGB video of a talking face by directly mapping audio input to the video output space. Temporal generative adversarial networks (GANs) have been used for talking head generation [65, 80]. In particular, Vougioukas et al. [65] present an approach for generating facial animation from a single RGB image using a temporal GAN. While this approach directly utilizes speech information for talking head generation, MakeItTalk [80] disentangles content from style and speaker identity, facilitating speech-driven generation that can be applied to diverse types of realistic and hand-drawn head portraits. In the second category, an intermediate 3DMM [3, 18] is used to guide the 2D neural rendering of talking heads from audio [55, 63, 72, 78], concentrating on the facial expressions. Extending these approaches, Wang et al. [67] add the head movements of the speaker to the synthesis. Drawing inspiration from dynamic neural radiance fields [22], several works [26, 71] leverage dynamic neural radiance fields to learn personalized talking head models that can be rendered under novel views, controlled by audio inputs.

**Speech-Driven 3D Facial Animation** Speech-driven 3D facial animation is a long-standing research question in

computer graphics and animation. Traditional methods are procedural techniques [15, 17, 19, 32] wherein the goal is to animate pre-defined facial rigs through procedural rules. With the rise of deep learning, these methods have been extended by learning-based approaches [5, 12, 20, 33, 47, 59, 62, 63], where viseme patterns are directly learned from data. A common theme for procedural techniques was to use Hierarchical Hidden Markov Models (HMM) as the basis for generating visemes from input text or audio, and subsequent facial animations were generated either through viseme-dependent co-articulation models [15, 17] or by blending facial templates [32]. Unlike these procedural approaches, data-driven approaches learn to generate 3D facial animation from data. These approaches typically leverage pretrained speech models [1, 27, 51] to generate an abstract and generalized representation of the input audio, which then serves as the input to a convolutional or auto-regressive model, mapping to either a 3D Morphable Model (3DMM) parameter space or directly to 3D meshes. For instance, Karras et al. [33] exemplify learning a 3D facial animation model from small-scale but high-quality actor-specific 3D data. While demonstrating a strong baseline, this method lacks generalization to new subjects and exhibits improper lip movements. In VOCA [12], a model is trained on 3D data of multiple subjects, enabling the animation of corresponding identities from the input audio. While this approach improves generalization over [33], the generalization remains limited as it requires one-hot encoding of identities at inference time. MeshTalk [47] adopts a generalized approach by learning a categorical representation for facial expressions, auto-regressively sampling from this categorical space to animate a given 3D facial template mesh based on audio inputs. In FaceFormer [20] a pretrained Wav2Vec [1] audio representation and a transformer-based decoder to regress displacements onto a template mesh is used. Analogous to VOCA, FaceFormer incorporates a speaker identification code into the decoder, offering the flexibility to choose from talking styles present in the training set. CodeTalker [69] trains a Vector Quantized Variational Autoencoder (VQ-VAE) as motion prior. Once the codebooks are trained, CodeTalker uses a transformer to learn the conditional distribution of codes given a speech signal. More recently, EMOTE [14] utilizes the MEAD dataset [66] and generates pseudo-ground truth meshes using EMOCA [13], a state-of-the-art 3D face reconstruction method. Since MEAD contains different emotion and intensity labels, the trained model can generate speech-driven facial animation with various emotions, at the cost of slightly inferior lip-sync and realism compared to other approaches.

**Motion Diffusion** Diffusion models [29, 53] have become a popular choice for generative tasks, among them the generation of images [16, 24, 48, 50, 54], videos [4, 30, 52, 79], 3D objects [7, 37, 44, 73], audio [6, 35], and human motion [39, 46, 61, 74, 77]. At their core, diffusion models are trained to iteratively denoise samples such that at inference time, new samples can be created from white-noise input. The stochastic nature of this process makes diffusion models highly suitable for modeling complex distributions. In contrast to generative adversarial networks (GAN)[16], they show higher diversity and quality. Typically, the denoising process is conditioned on additional modalities such as text [48], audio [57], or depth maps [76]. For an in-depth state-of-the-art report on diffusion models, we kindly refer the reader to the excellent surveys of Po et al. [43] and Yang et al. [70]. In terms of methods that use diffusion for motion synthesis, MDM [60] and Mofusion [39] are the closest to our work, producing impressive results for the task of body motion synthesis and editing from input text and music. In contrast to both methods, we use a 1D-convolution architecture with concatenation-based condition injection to train on small-scale datasets more efficiently.

Despite the rapidly growing popularity of diffusion models, to the best of our knowledge, only two concurrent works apply them to the task of speech-driven face animation: FaceDiffuser [56] uses a pretrained speech representation model to convert the audio signal into sequences of latent feature vectors which condition a diffusion model that is based on recurrent GRU layers [8]. Similarly, DiffPoseTalk [57] uses a diffusion model for speech-to-motion, however, they employ a transformer-based architecture and further propose a style encoder to personalize the synthesis based on short reference sequences. In contrast to these two methods, we use a convolution-based architecture which allows us to train our model on a small high-quality dataset and obtain higher synthesis quality and diversity. Furthermore, FaceDiffuser and DiffPoseTalk both perform auto-regressive inference, which restricts their applicability to animation editing tasks, e.g., motion between two keyframes cannot be synthesized consistently. Lastly, FaceDiffuser does not allow for personalization, and DiffPoseTalk's personalization capabilities are restricted within the space of the coefficients of a pretrained head model. In contrast, our method regresses vertex displacements directly and allows for personalization with higher fidelity.

## 3. Preliminaries

**Denoising Diffusion Probabilistic Models** Our method is based on the diffusion framework of Sohl et al. [53]. During the diffusion process, a data sample $x_0$ from the training distribution is iteratively disturbed by Gaussian noise for $T$ steps, resulting in a transition of the sample to white noise.
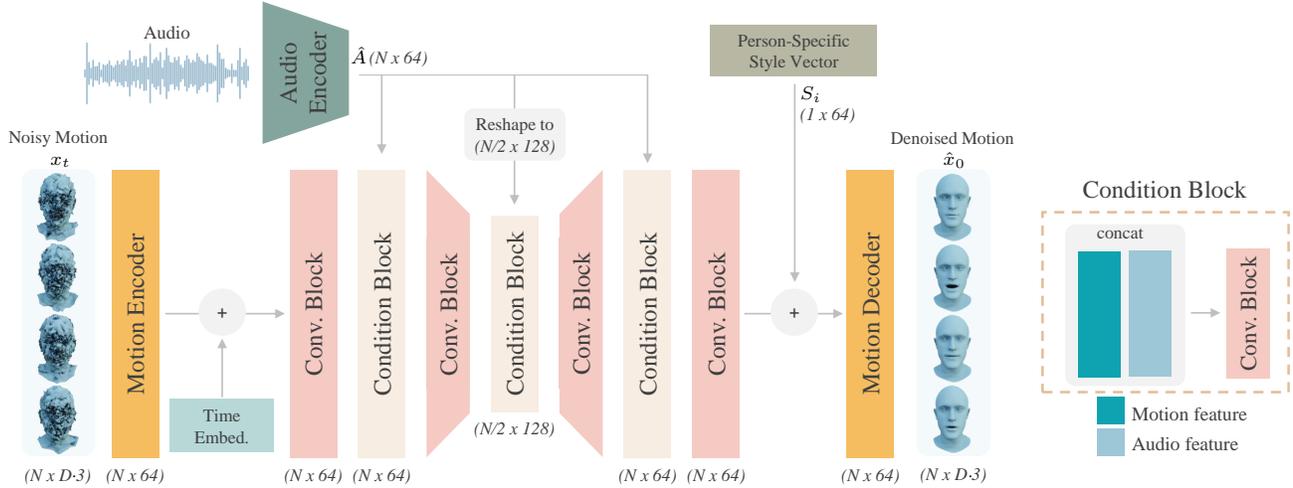
Figure 2. Our method takes noised vertex displacements, denoted as $x_t$, and the diffusion time step embedding as inputs to predict a denoised sample $\hat{x}_0$, leveraging both the audio conditioning signal $\hat{A}$ and a person-specific feature vector $S_i$. Our approach employs wav2vec2.0 [1] for extracting audio features from the raw audio signal. The audio condition is injected into the network by concatenation through a series of convolutional blocks. Note that $N$ corresponds to the frame count of the sequence and $D$ to the number of vertices.

The forward diffusion step $t$ is defined as:

$$x_t \sim q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1-\beta_t}x_{t-1}, \beta_t I), t = 1...T, \tag{1}$$

where $\beta_t$ is following a predefined variance schedule.

A denoising model is trained to reverse the diffusion process, hence to estimate $q(x_{t-1}|x_t)$. Following recent work [57, 60], we train a neural network $\theta$ to estimate $x_0$ from its noised version $x_t$: $\hat{x}_0 = \theta(x_t, t, C)$ with $C$ denoting additional conditions. Following [29], the inverse diffusion process is then given through:

$$q(x_{t-1}|x_t) = \mathcal{N}\left(x_{t-1}; \sqrt{\bar{\alpha}_{t-1}}\theta(x_t, t, C), (1-\bar{\alpha}_{t-1})I\right), \tag{2}$$

with $\alpha_t := 1 - \beta_t$ and $\bar{\alpha}_t := \prod_{k=1}^{t} \alpha_k$. For generating new samples, we randomly sample $x_T$ from a Gaussian distribution and iteratively denoise it until $t = 0$ is reached.

To add diversity, we employ Classifier-Free Guidance (CFG) [28] and calculate the output as a weighted sum of the conditional and unconditional prediction:

$$\theta_s(x_t, t, C) := \theta(x_t, t, \emptyset) + s \cdot [\theta(x_t, t, C) - \theta(x_t, t, \emptyset)], \tag{3}$$

where $s$ is the guidance scale and $\theta(x_t, t, \emptyset)$ denotes the unconditional prediction in which we set the audio conditions to zero. Note that while CFG is typically used with a guidance scale $> 1$ to enhance alignment with the condition, we set it to values $< 1$ (0.5 unless specified otherwise) to increase diversity.

**Audio Encoding** Similar to other state-of-the-art methods [12, 20, 62, 68], we adopt the pretrained

Wav2Vec2.0 [1] model to generate audio features from the raw audio signal. Wav2Vec2.0 uses a self-supervised learning approach to map audio to quantized feature vectors with 768 channels. We resample the output of Wav2Vec2.0 via linear interpolation to match the sampling rate of the motion sequences (30fps for VOCAset [12]). A trainable linear layer is applied to project the feature vectors to 64 channels, resulting in a speech representation $\hat{A} \in \mathbb{R}^{N \times 64}$ for $N$ frames.

## 4. Method

Our goal is to synthesize and edit facial animations given speech audio as a conditioning input. We represent facial animations as a sequence of 3D vertex displacements that can be applied on top of a template mesh. To generate those displacements from audio, we employ a diffusion-based model that is trained to iteratively denoise the displacement sequences. This architecture does not only produce stochastic outputs, it also allows us to edit the animation sequence by defining keyframing. As facial motions are person-specific, we design the architecture such that it can be adapted and fine-tuned to specific subjects, only requiring a short video sequence of the actor. In Figure 2, we show an overview of our method, which is detailed in the following.

As written above, we represent facial animations as sequences of 3D vertex displacements w.r.t. a template mesh. Let $x_0 \in \mathbb{R}^{N \times D \cdot 3}$ denote such a sequence where $N$ is the sequence length and $D$ is the number of vertices in the template mesh. The input to our diffusion model $\theta$ is a noised

4

vertex displacement sequence $x_t \in \mathbb{R}^{N \times D \cdot 3}$ and we aim to predict its noise-free counterpart: $\hat{x}_0 = \theta(x_t, t, C)$ given diffusion step $t$ and conditions $C$. As a first step, we employ a single fully connected layer as *Motion Encoder* to project $x_t$ to a 64-dimensional latent space. We positionally encode the diffusion step $t$ [53], map it to the latent space with a linear layer, and add it to the encoded $x_t$. We apply a series of 1D-convolution blocks to first reduce the temporal dimension of the activations, followed by an up-sampling convolution block to restore the original temporal dimension. Each convolution block is followed by a condition block to incorporate the audio features. The condition blocks concatenate the input features with the audio conditions and apply a dimension-preserving convolution. We add a person-specific feature vector $S_i \in \mathbb{R}^{1 \times 64}$ to the output of the convolutional layers prior to applying to the *Motion Decoder*. This produces the final noise-free sample $\hat{x}_0$. Similar to the *Motion Encoder*, the *Motion Decoder* is a single fully connected layer. Note that in our formulation, the condition $C$ represents the set of both the per-frame audio features $\hat{A}$ and the person-specific feature vector $S_i$.

In contrast to state-of-the-art methods on 3D facial animation synthesis that use transformer architectures [20, 57, 62, 68], we take inspiration from Pavllo et al. [41] and adopt a 1D-convolutional network as our backbone. Specifically, instead of infusing the condition through an attention mechanism, we use feature concatenation. We found that these architecture changes are crucial for efficiently training the model on the small available VOCA training set [12] (see Table 2). Note that while other methods train transformer architectures on bigger datasets with pseudo-ground-truth annotations, we show that our architecture changes allow us to achieve superior results while training on a smaller, yet high-quality dataset. In particular, the fully convolutional nature of our network allows us to randomly crop the sequences to 30 frames for training and generalize to sequences of arbitrary length at inference time. We find this data augmentation strategy to be vital for improved generalization and convergence in the unconditional setting. Note that this data augmentation strategy is not possible for transformer-based architectures since they rely on a consistent positional encoding, which prevents them from generalizing to longer sequences. While auto-regressive motion synthesis could in theory resolve this limitation, it would make crucial animation editing tasks impossible, such as in-betweening distant motion frames. Further, we empirically find that for the unconditional case, in which the audio conditions are set to 0, transformer-based architectures do not converge on the small VOCA training set due to its limited size. However, as outlined in Section 3, unconditional synthesis is crucial for synthesis diversity.

## 4.1. Training

Similar to [57, 61], we train our diffusion model to predict the ground truth vertex displacements $x_0$ from their noised counterparts $x_t$:

$$\mathcal{L}_{\text{simple}} = ||x_0 - \theta(x_t, t, C)||^2. \tag{4}$$

In comparison to predicting the applied noise which is common practice in related work [39, 48, 77], we empirically found that predicting the ground truth displacements yields better convergence in the unconditional and person-specific fine-tuning setup. Furthermore, we take inspiration from [12, 62] and add a velocity loss $\mathcal{L}_{\text{vel}}$ to improve temporal smoothness:

$$\mathcal{L}_{\text{vel}} = \frac{1}{N-1} \sum_{n=1}^{N} ||(x_{0,n} - x_{0,n-1}) - (\hat{x}_{0,n} - \hat{x}_{0,n-1})||^2, \tag{5}$$

where $x_{0,n}$ denotes the ground truth vertex displacements in frame $n$. Our final training objective is formulated as:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{simple}} + \lambda_{\text{vel}} \cdot \mathcal{L}_{\text{vel}}. \tag{6}$$

We empirically set $\lambda_{\text{vel}} = 10.0$ unless specified otherwise. Note that during training, we randomly set the audio condition $C$ to 0 in 10% of the cases in order to enable unconditional synthesis at inference time.

## 4.2. Person-Specific Fine-tuning

For capturing the speaking style of a subject that is not part of the training set, we require a short reference talking head video. The facial movements are extracted with the state-of-the-art monocular face tracker MICA [82]. We use the tracked meshes as pseudo ground truth and fine-tune the entire model to fit the expression distribution of the target subject using the training objective from Eq. (6).

## 5. Dataset

We train our model on the VOCAset [12] since it provides high-quality, speech-aligned 3D face scan sequences. It consists of 12 actors (6 female and 6 male) with 40 sequences each with a length of 3-5 seconds, resampled at 30fps. Following previous work [62], we use the train/val/test set split of 8, 2, 2 actors. All 40 sequences of the training actors are used during training. However, for the test and validation, only 20 sequences without overlap with the speech scripts of the training sequences are used. For the style adaption experiment, we split the 40 sequences of the test actors to 18, 2, 20 for train/val/test sets. The test sequences of the experiments w/ and w/o style adaptation are identical, allowing a direct comparison of the scores in Table 1.

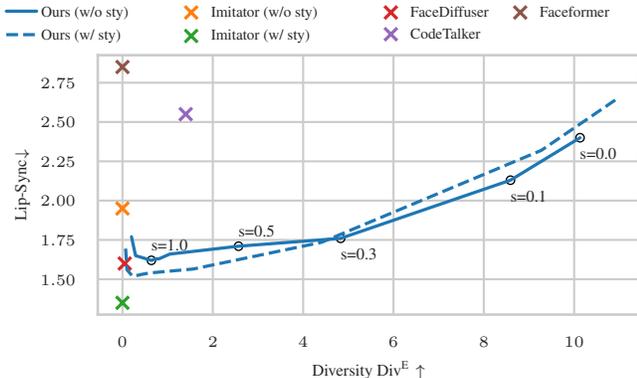We evaluate person-specific fine-tuning for in-the-wild sequences on the video sequences from Imitator [62]. The

5

Figure 3. We investigate the impact of the classifier-free-guidance scale $s$ [28] using the 'Lip-sync' and $Div^E$ metrics. Lower guidance values yield animations with significantly more diverse motion but inferior lip-sync quality. Conversely, higher guidance values result in high-quality animation with reduced diversity. For our experiments, we maintain a fixed guidance value of 0.5, a sweet spot that balances high diversity with excellent lip-sync quality.

provided videos are 2 minutes long which we divide into 60/30/30 seconds for train/val/test respectively.

## 6. Results

We evaluate and we compare our method against state-of-the-art methods: VOCA [12], Faceformer [20], CodeTalker [68], EMOTE [14] , FaceDiffuser [56] and Imitator [62]. Figure 4 shows a qualitative comparison to the baselines on a test sequence from the VOCAset. We find that only our method and Imitator produce expressive facial animations that match the speaking style of the target subject. For additional qualitative results, we refer to the suppl. video.

**Quantitative Comparison**  In Table 1, we present a quantitative evaluation based on the following metrics: *Lip-Sync* measures the lip synchronization using Dynamic Time Warping to compute the temporal similarity [62]. *Lip-max* [47] reports the mean of the maximal per-frame lip distances. $L_2^{\text{lip}}$ and $L_2^{\text{face}}$ correspond to the mean $L_2$ vertex errors for the lip region and the entire face respectively. Additionally, we adopt the diversity metric $Div^E$ proposed by Ren et al. [46] to assess the diversity of animations generated from the same audio.

Note that only Imitator [62] and our method allow for optimizing person-specific speaking styles. For a fair comparison, we report results both with and without person-specific style optimization for these two models. For the non-personalized synthesis, we find that our method with the default guidance scale $s = 0.5$, improves synthesis diversity by over 80% compared to the closest competitor (see

| | Method | $Div^E \uparrow$ | Lip-Sync $\downarrow$ | Lip-max $\downarrow$ | $L_2^{\text{lip}} \downarrow$ | $L_2^{\text{face}} \downarrow$ |
|---|---|---|---|---|---|---|
| | | Non-Personalized Synthesis | | | | |
| 1 | VOCA [12] | – | 5.30 | 7.06 | 0.20 | 0.94 |
| 2 | Faceformer [20] | – | 2.85 | 5.41 | 0.14 | **0.80** |
| 3 | Imitator [62] | – | 1.95 | **4.95** | **0.12** | 0.85 |
| 4 | CodeTalker [68] | 1.40 | 2.55 | 5.02 | 0.14 | 0.88 |
| 5 | FaceDiffuser [56] | 0.05 | **1.60** | 5.20 | 0.16 | 0.89 |
| 6 | Ours$_{s=0.5}$ (w/o sty) | **2.57** | 1.71 | 5.20 | 0.15 | 0.86 |
| 7 | Ours$_{s=1.0}$ (w/o sty) | 0.64 | 1.62 | 5.13 | 0.15 | 0.84 |
| | | Personalized Synthesis | | | | |
| 8 | Imitator (w/ sty) | – | **1.35** | **3.43** | **0.09** | **0.76** |
| 9 | Ours (w/ sty) | **1.57** | 1.56 | 4.01 | 0.11 | 0.78 |

Table 1. The quantitative results from the VOCAset [12] demonstrate that our method outperforms the baseline in generating diverse motions. It achieves comparable performance to the state-of-the-art regression method in terms of lip-sync accuracy, both in personalized and non-personalized setups.

row 4 and 6 of Table 1). Furthermore, the same model performs second-best in terms of *Lip-Sync* and is competitive on all other metrics. This demonstrates that our method with guidance scale $s = 0.5$ offers significantly improved synthesis diversity while still ensuring plausible lip synchronization. Note that we can use the guidance scale parameter to freely trade synthesis diversity for lip-sync accuracy. As we increase the guidance scale to $s = 1.0$ (row 7 of Table 1), we are able to match the *Lip-Sync* score of the top-performing baseline. We visualize the trade-off between lip-sync accuracy and synthesis diversity in Figure 3. We find that the guidance scale $s$ is an effective tool to increase synthesis diversity beyond all baselines with only a small loss of lip-sync accuracy for $0.3 \le s \le 1.0$.

When personalizing our model, we find that while the synthesis diversity decreases, all other scores improve. The resulting model now consistently outperforms all non-personalized models on all metrics. Note that we only require $\sim 100s$ of video to personalize our method to an unseen identity. The moderate decline in synthesis diversity during personalization is an expected behavior and even is an indicator of successful personalization. During personalization, the model learns to suppress movements that do not align with the target identity and as a natural consequence, the synthesis diversity is reduced.

In comparison to Imitator [62] after personalization, we achieve higher synthesis diversity and comparable accuracy scores, yet we are not able to outperform this baseline. However, note that Imitator is a deterministic model that does not allow for stochastic synthesis. Also, in contrast to Imitator and all other baselines, our method is the only one that enables animation editing like motion inbetweening.

**User Study**  We conducted an A/B user study to assess our method's perceptual performance. We sample 20 sequences combined from the VOCAset test set and the in-the-wild se-

**Get a calico cat to keep**

GT

Ours — *Accurate Motion*

FaceDiffuser[56] — *Improper Articulations*

Imitator[62] — *Accurate Motion*

CodeTalker[68] — *Muted Expressions*

FaceFormer[20] — *Muted Expressions*
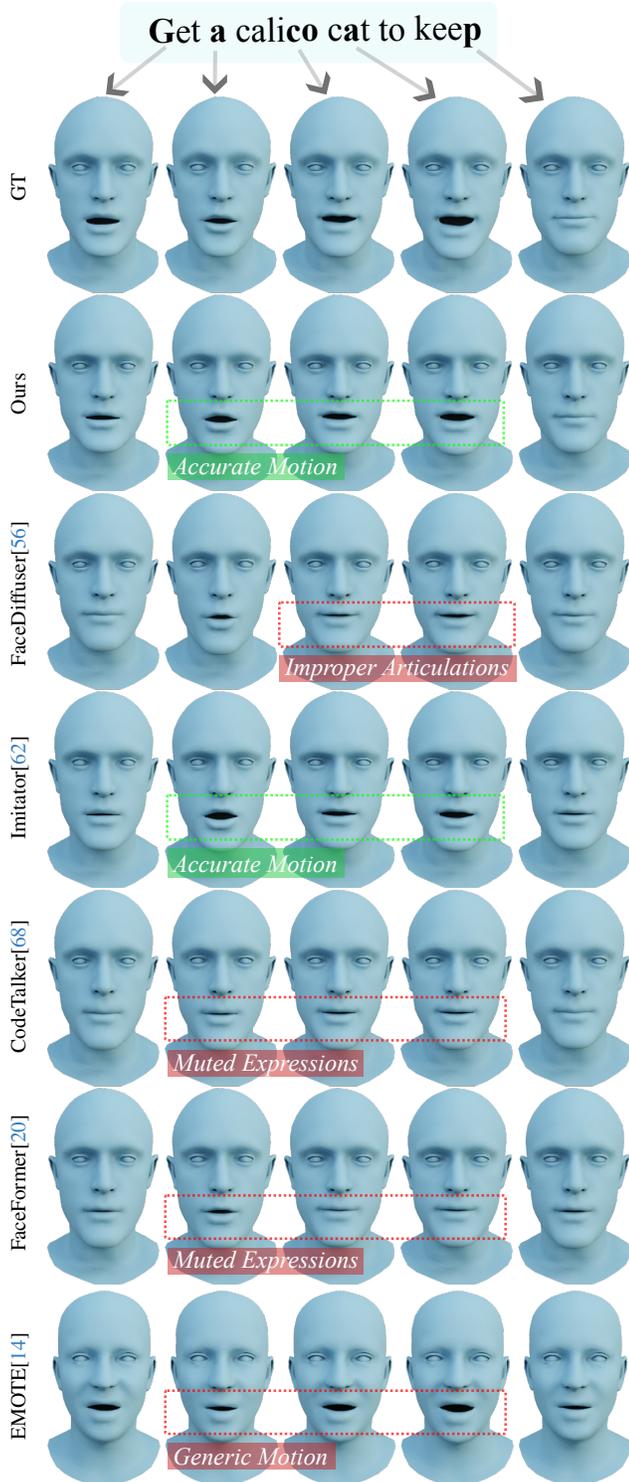
EMOTE[14] — *Generic Motion*

Figure 4. Qualitative comparison. Only our method and Imitator produce expressive motions that match the target speaking style. While Imitator synthesizes similarly convincing animations, its outputs are not diverse (see Table 1) and it cannot be used for animation editing.
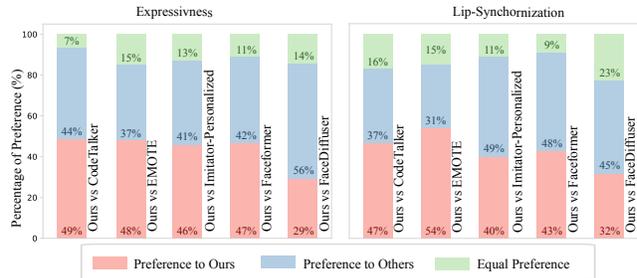


Figure 5. User-study results on the VOCAset [12]. Overall, our performance matches or surpasses the state-of-the-art (SOTA) in terms of expressiveness and lip-sync, except for FaceDiffuser [56]. FaceDiffuser generates animations with less diversity and high lip-sync. Additionally, we evaluate the person-specific speaking-style modeling by comparing against Imitator-Personalized [62]. Users were presented with a reference video and asked to select the synthesized method that closely resembled it. Notably, 40% of users favored our method for better style-similarity, while 17% found both methods equally comparable.

quences from Imitator, resulting in 100 A/B comparisons across five baselines. On Amazon Mechanical Turk(AMT), we divided the A/B comparisons into 5 HITs (Human Intelligence Task), each with 25 individual assignments. For each HIT, users select their preference for a method based on expressiveness and lip-synchronization. The results in Figure 5 show that our method outperforms most of the baselines in terms of lip-synchronization and expressiveness, except for the concurrent work FaceDiffuser. The difference in performance w.r.t. FaceDiffuser is an inherent trade-off between diversity and lip-sync quality. Our method with a guidance value of 1.0 matches the performance of FaceDiffuser with better diversity and additionally, our method allows for person-specific fine-tuning and motion editing.

We, further, conduct a second user study to evaluate the speaking style preservation of our personalized model in comparison to Imitator. To this end, the AMT users rated the similarity based on a reference video and the synthesized videos of the VOCA test set. 40% of the users preferred our method, 43% preferred Imitator, 17% voted for tie. This demonstrates that our method for the first time allows to synthesize diverse personalized face animations and enables animation editing without significantly reducing speaking style faithfulness.

**Motion Editing** We show motion editing using keyframes in Figure 6. In this application, we selectively replace the predicted denoised vertex-displacement sequences $\hat{x}_0$ with ground truth values during the denoising process. This is similar in spirit to well-established diffusion-based image inpainting methods [38]. We additionally show unconditional motion synthesis and editing

| | Method | $Div^E$ ↑ | Lip-Sync ↓ | Lip-max ↓ | $L_2^{lip}$ ↓ | $L_2^{face}$ ↓ |
|---|---|---|---|---|---|---|
| | | (a) Architecture ablation | | | | |
| 1 | Ours (w/ conv attn) | 0 | 1.68 | 5.17 | 0.15 | 0.88 |
| 2 | Ours (w/ FF arch) | 0 | 3.49 | 6.2 | 0.19 | 0.96 |
| 3 | Ours | 1.57 | 1.56 | 4.01 | 0.11 | 0.78 |
| | | (b) Person-specific Fine-tuning | | | | |
| 4 | Ours ($\sim$ 5s) | 29.95 | 4.89 | 13.05 | 0.37 | 1.95 |
| 5 | Ours ($\sim$ 30s) | 0.18 | 1.81 | 4.90 | 0.13 | 0.86 |
| 6 | Ours ($\sim$ 60s) | 0.67 | 1.69 | 4.18 | 0.12 | 0.83 |
| 7 | Ours ($\sim$100s) | 1.57 | 1.56 | 4.01 | 0.11 | 0.78 |
| | | (c) Audio noise ablation | | | | |
| 8 | Ours (high noise) | 6.41 | 2.56 | 6.68 | 0.19 | 0.94 |
| 9 | Ours (med. noise) | 2.54 | 1.97 | 4.93 | 0.13 | 0.84 |
| 10 | Ours (low noise) | 1.85 | 1.78 | 4.4 | 0.12 | 0.77 |

Table 2. We ablate our method with respect to (a) architecture, (b) person-specific fine-tuning dataset size, and (c) audio noise level robustness. Our architecture outperforms attention-based conditioning mechanisms (row 1), and the transformer backbone of Faceformer [20] (row 2). Further, we show that 30s of video suffice to perform person-specific fine-tuning while 100s further improve all scores (row 4-7). Our method is robust wrt. medium and low audio noise levels (row 9-10).

results in the supplemental material. As can be seen in Figure 6, the personalization of the motion synthesis is important to match the talking style, preventing an abrupt style change.

**Ablation** We evaluate the benefits of our 1D-convolutional architecture over the transformer-based architectures used in the baseline methods, and the attention-based convolutional approach in MoFusion [39]. To this end, we compare our proposed architecture against two variants. (i) We replace the conditional convolution blocks with attention layers such as those used in Mofusion (row 1 in Table 2). (ii) We replace the entire backbone with the transformer-based architecture from Faceformer (row 2 in Table 2). We observe that both changes to our method significantly worsen all scores. This confirms the improved effectiveness of our proposed architecture when training on small datasets (see discussion in Section 4).

*Training data for person-specific fine-tuning:* We evaluate the impact of the data set size for person-specific fine-tuning in Table 2 row 4-7. To this end, we perform fine-tuning on the VOCAset test set by varying the dataset size to 5/30/60/100s. We observe that while the fine-tuning on 5 s diverges, 30s and 60s suffice to achieve acceptable results. For 100s of data, our model is able to synthesize motion with a better *Lip-sync* and diversity $Div^E$.

*Sensitivity study:* Additionally, we conducted a noise sensitivity experiment similar to [12, 62], where we added white noise to the input audio with a negative gain of 36db (low), 24db (medium), and 12db (high). As reported in Table 2 rows 8-10, our method produces robust high-quality facial animations for low and medium noise levels.
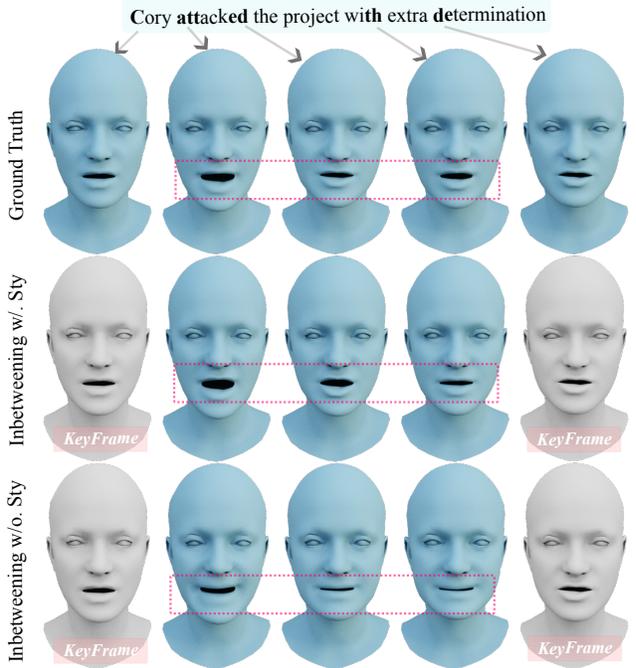


Figure 6. Qualitative evaluation of the importance of person-specific finetuning for motion editing. As highlighted in purple, without finetuning we observe an abrupt change in speaking style between the keyframes and the generated motion, thus rendering the results unrealistic.

# 7. Discussion

Our proposed method excels in synthesizing and editing diverse 3D facial animations based on speech. Its architecture is carefully designed for efficient training on a small, yet high-quality dataset. Nonetheless, we believe the performance can be further improved by increasing the amount of training data. In the user study, our method outperforms most baselines and is on par with the concurrent work FaceDiffuser [56]. At the same time, our method synthesizes animations with higher diversity and allows the editing of existing animations such as inserting new motion frames.

Currently, we do not consider head motion (e.g. neck rotation) in our method, due to the lack of this data in the VOCAset. However, we believe that it can be extended to head motion, given a suitable dataset.

# 8. Conclusion

With 3DiFACE we present the first method that can both generate and edit diverse 3D facial animations from speech input. Employing classifier-free guidance provides us with an effective tool to balance synthesis diversity and accuracy allowing us to generate animations with unprecedented diversity while outperforming or matching all baselines with respect to synthesis accuracy. Through personalization,

we can extract person-specific speaking styles from short ($\sim 100s$) videos which significantly improves performance. Further, our architecture allows us to edit animations by using keyframes. We are convinced that these properties make 3DiFACE a powerful tool for content creators and are excited for future applications.

## 9. Acknowledgements

## 10. Additional Applications

### 10.1. Unconditional motion synthesis and editing:

While unconditional motion synthesis has been extensively applied in the motion synthesis domain [45, 61], to the best of our knowledge, its application in 3D facial animation synthesis remains widely unexplored. The significance of an unconstrained facial motion synthesis method cannot be overstated. It holds substantial potential for various applications, such as animating background characters like NPCs in movies and games. Additionally, it enables targeted editing of specific facial elements—such as eye blinks and eyebrow motions—since these non-verbal facial expressions often exhibit weak or no correlation with audio features. Moreover, an unconditional model serves as a valuable motion prior for downstream tasks, extending its utility beyond synthesis and editing applications. Our demonstration of unconditional synthesis and editing are showcased in Figure 7, underscoring the potential and versatility of such unconstrained models for 3D facial animation synthesis.

### 10.2. Motion Keyframing/Inbetweeing:

We evaluate the performance of motion inbetweening with respect to the input data. To this end, we preserve 5%, 10%, 20%, 50% of the starting and ending frames, and then perform inbetweening for the intermediate motion sequences. Furthermore, we assess the robustness of the inbetweening by randomly inserting keyframes at different rates: 1KF/sec, 2KF/sec, and 3KF/sec. These evaluations are conducted for all sequences of the test subject 024 from the VOCAset [12], and the resulting metrics are presented in Table 3. These metrics demonstrate the efficacy of our method in robustly editing motion.
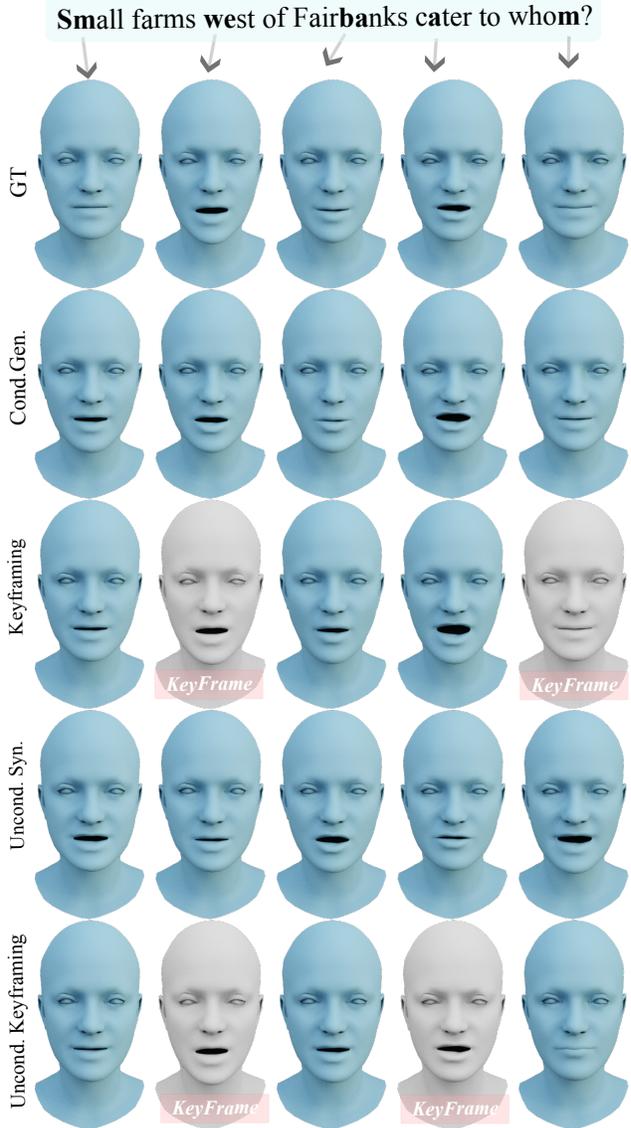


Figure 7. Qualitative illustration of motion inbetweening using our conditional and unconditional model. In row 2 and 3, we showcase a sequence synthesized conditionally and subsequently refined using keyframes. In Row 4 (Uncond. Syn.), we present our unconditional synthesis results. As observed from the results, our model can unconditionally synthesize facial animations that appear plausible. Further in row 5 (Uncond. Keyframing), we see that our method can unconditionally inbetween facial animation while preserving the speaking-style of the target actor. This progression demonstrates our model's capabilities: from conditional synthesis and keyframe-based editing to unconditional synthesis and editing, while maintaining the actor's speaking-style.

| | Method | Div$^{\text{E}}$ ↑ | Lip-Sync ↓ | Lip-max ↓ | L$_2^{\text{lip}}$ ↓ | L$_2^{\text{face}}$ ↓ |
|---|---|---|---|---|---|---|
| 1 | Ours (synthesis) | 1.35 | 1.4 | 3.41 | 0.09 | 0.69 |
| 2 | Ours (Ip 5%) | 1.27 | 1.17 | 3.19 | 0.09 | 0.65 |
| 3 | Ours (Ip 10%) | 1.24 | 1.15 | 2.98 | 0.08 | 0.61 |
| 4 | Ours (Ip 20%) | 1.15 | 1.01 | 2.66 | 0.07 | 0.54 |
| 5 | Ours (Ip 50%) | 0.9 | 0.68 | 1.74 | 0.05 | 0.33 |
| 6 | Ours (1KF/sec) | 1.26 | 1.28 | 3.17 | 0.09 | 0.65 |
| 7 | Ours (2KF/sec) | 1.14 | 1.2 | 2.98 | 0.08 | 0.62 |
| 8 | Ours (3KF/sec) | 1.05 | 1.1 | 2.77 | 0.07 | 0.59 |

Table 3. We quantitatively evaluate our motion editing capability on all the test sequences of the subject 024 in the VOCAset [12]. To this end, first we preserve 5%, 10%, 20%, 50% of the starting and ending frames, and then perform inbetweening for the intermediate motion sequences. In addition, we assess the robustness of the inbetweening by randomly inserting keyframes (KF): 1KF/sec, 2KF/sec, and 3KF/sec. From the metrics, we can see that the synthesis quality increases significantly with adding more keyframes, which is a clear indication that the model matches the ground truth and produces realistic motion. For animators and artists, this means that they can insert any number of keyframes they want and have fine-grained control over the motion synthesis. Note that keyframes could also stem from previously generated motion sequences using our method (iterative refinement).

## 11. Implementation:

### 11.1. Baselines:

For VOCA [12], Faceformer [20], Imitator [62] and FaceDiffuser [56], we use the pre-trained model provided in the offical repositories. For CodeTalker [68], we adapt the official implementation to add the functionality of generating diverse motion. Especially, we re-train the audio conditioned codebook sampling (stage 02) to randomly sample a code from top 'm' closest codes instead of always using the closest code. This process is in spirit close to training the language-based models, where a new diverse text sequence is generated by sampling the 2nd or 3rd closest language token over the token with maximum probability. By adapting this method, we ensure that CodeTalker could generate diverse samples for a given audio input. For EMOTE [14], we request the authors to run their method on the VOCAset [12] and use it for the qualitative and perceptual user-study.

### 11.2. Training Details:

We train our method using ADAM [34] with a learning rate of *1e-4* for 140K iterations with a batch size of 64. Our diffusion framework is based on the Gaussian diffusion from Nichol *et al.* [40], we set the diffusion step to 500 for our experiments. During training, we randomly crop the sequences to length of 30 frames. Our light-weight architecture enables to train our model on a single Nvidia quadro P6000 $32GB$ within 30 hours. The lightweight architecture is also critical for person-specific style-adaption with a short reference video. For person-specific speaking-style, we use

the same training setup as from the generalized setting, except that we only train it for $30K$ iterations. For evaluating the best checkpoint, we fix the guidance scale $s = 0.99$ and evaluate all the saved checkpoints on the validation set. Further, we fix the best checkpoint and vary the guidance scale from *s=0, 0.1 ... to 1.0* with increment of $0.1$ and find the best guidance factor. From our experiment, we found the guidance scale of $0.5$ balances the lip-synchornization and diversity and provides best results.

### 11.3. Inference

Our method takes 3.15sec to produce 1sec (30 frame) of facial animation on a single Nvidia GeForce RTX 3090 $24GB$, compared to 5.78sec for the concurrent method FaceDiffuser [56].

## 12. Broader Impact

We introduce a method for realistic facial animation synthesis and editing that matches the speaking-style of any given target actor. These animations hold promise for driving virtual avatars in AR or VR settings, especially, in immersive communication technologies. Yet, it is essential to acknowledge the potential pitfalls of such advancements, notably in the realm of 'DeepFakes.' By employing voice cloning techniques, our method can generate 3D facial animations that drive digital avatar methods like [2, 23, 25, 31, 81], which could be abused for identity theft, cyberbullying, and various criminal activities. Advocating for transparent research practices, we strive to illuminate the risks associated with technology misuse. Sharing our implementation aims to foster research in digital multimedia forensics, particularly in developing synthesis methods crucial for training data utilized in spotting forgeries [49].

## References

[1] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. 3, 4

[2] Shrisha Bharadwaj, Yufeng Zheng, Otmar Hilliges, Michael J. Black, and Victoria Fernandez Abrevaya. FLARE: Fast learning of animatable and relightable mesh avatars. *ACM Transactions on Graphics*, 42:15, 2023. 10

[3] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces. In *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*, pages 187–194, 1999. 2

[4] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 3

[5] Yong Cao, Wen C. Tien, Petros Faloutsos, and Frédéric Pighin. Expressive speech-driven facial animation. *ACM Trans. Graph.*, 24(4):1283–1302, 2005. 3

[6] Nanxin Chen, Yu Zhang, Heiga Zen, Ron J. Weiss, Mohammad Norouzi, and William Chan. Wavegrad: Estimating gradients for waveform generation, 2020. 3

[7] Rui Chen, Yongwei Chen, Ningxin Jiao, and Kui Jia. Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation. *arXiv preprint arXiv:2303.13873*, 2023. 3

[8] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling, 2014. 3

[9] Joon Son Chung, Amir Jamaludin, and Andrew Zisserman. You said that? *arXiv preprint arXiv:1705.02966*, 2017. 2

[10] Michael M. Cohen, Rashid Clark, and Dominic W. Massaro. Animated speech: research progress and applications. In *Proc. Auditory-Visual Speech Processing*, page 200, 2001. 1

[11] Daniel Cudeiro, Timo Bolkart, Cassidy Laidlaw, Anurag Ranjan, and Michael J Black. Capture, learning, and synthesis of 3d speaking styles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10101–10111, 2019. 1

[12] Daniel Cudeiro, Timo Bolkart, Cassidy Laidlaw, Anurag Ranjan, and Michael J. Black. Capture, Learning, and Synthesis of 3D Speaking Styles. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10093–10103, Long Beach, CA, USA, 2019. IEEE. 2, 3, 4, 5, 6, 7, 8, 9, 10

[13] Radek Danecek, Michael J. Black, and Timo Bolkart. EMOCA: Emotion driven monocular face capture and animation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 20311–20322, 2022. 3

[14] Radek Daněček, Kiran Chhatre, Shashank Tripathi, Yandong Wen, Michael Black, and Timo Bolkart. Emotional speech-driven animation with content-emotion disentanglement. ACM, 2023. 2, 3, 6, 7, 10

[15] José Mario De Martino, Léo Pini Magalhães, and Fábio Violaro. Facial animation based on context-dependent visemes. *Computers & Graphics*, 30(6):971–980, 2006. 3

[16] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021. 3

[17] Pif Edwards, Chris Landreth, Eugene Fiume, and Karan Singh. Jali: an animator-centric viseme model for expressive lip synchronization. *ACM Transactions on graphics (TOG)*, 35(4):1–11, 2016. 1, 3

[18] Bernhard Egger, William AP Smith, Ayush Tewari, Stefanie Wuhrer, Michael Zollhoefer, Thabo Beeler, Florian Bernard, Timo Bolkart, Adam Kortylewski, Sami Romdhani, et al. 3d morphable face models—past, present, and future. *ACM Transactions on Graphics (TOG)*, 39(5):1–38, 2020. 2

[19] T. Ezzat and T. Poggio. MikeTalk: a talking facial display based on morphing visemes. In *Proceedings Computer Animation '98 (Cat. No.98EX169)*, pages 96–102, Philadelphia, PA, USA, 1998. IEEE Comput. Soc. 3

[20] Yingruo Fan, Zhaojiang Lin, Jun Saito, Wenping Wang, and Taku Komura. Faceformer: Speech-driven 3d facial animation with transformers. *CoRR, abs/2112.05329*, 2021. 2, 3, 4, 5, 6, 7, 8, 10

[21] Yingruo Fan, Zhaojiang Lin, Jun Saito, Wenping Wang, and Taku Komura. Faceformer: Speech-driven 3d facial animation with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18770–18780, 2022. 1

[22] Guy Gafni, Justus Thies, Michael Zollhöfer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. *CoRR, abs/2012.03065*, 2020. 2

[23] Guy Gafni, Justus Thies, Michael Zollhöfer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8649–8658, 2021. 10

[24] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H. Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion, 2022. 3

[25] Philip-William Grassal, Malte Prinzler, Titus Leistner, Carsten Rother, Matthias Nießner, and Justus Thies. Neural head avatars from monocular rgb videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18653–18664, 2022. 10

[26] Yudong Guo, Keyu Chen, Sen Liang, Yongjin Liu, Hujun Bao, and Juyong Zhang. Ad-nerf: Audio driven neural radiance fields for talking head synthesis. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 2

[27] Awni Hannun, Carl Case, Jared Casper, Bryan Catanzaro, Greg Diamos, Erich Elsen, Ryan Prenger, Sanjeev Satheesh, Shubho Sengupta, Adam Coates, and Andrew Y. Ng. DeepSpeech: Scaling up end-to-end speech recognition. 2014. 3

[28] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance, 2022. 4, 6

[29] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. 3, 4

[30] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P. Kingma, Ben Poole, Mohammad Norouzi, David J. Fleet, and Tim Salimans. Imagen video: High definition video generation with diffusion models, 2022. 3

[31] Berna Kabadayi, Wojciech Zielonka, Bharat Lal Bhatnagar, Gerard Pons-Moll, and Justus Thies. Gan-avatar: Controllable personalized gan-based human head avatar. In *International Conference on 3D Vision (3DV)*, 2024. 10

[32] G.A. Kalberer and L. Van Gool. Face animation based on observed 3D speech dynamics. In *Proceedings Computer Animation 2001. Fourteenth Conference on Computer Animation (Cat. No.01TH8596)*, pages 20–251, Seoul, South Korea, 2001. IEEE Comput. Soc. 3

[33] 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 Transactions on Graphics*, 36(4):1–12, 2017. 3

11

[34] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. 10

[35] Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile diffusion model for audio synthesis, 2021. 3

[36] Tianye Li, Timo Bolkart, Michael. J. Black, Hao Li, and Javier Romero. Learning a model of facial shape and expression from 4D scans. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)*, 36(6), 2017. 2

[37] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content creation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 3

[38] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models, 2022. 7

[39] Jianxin Ma, Shuai Bai, and Chang Zhou. Pretrained diffusion models for unified human motion synthesis. *arXiv preprint arXiv:2212.02837*, 2022. 3, 5, 8

[40] Alex Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models, 2021. 10

[41] Dario Pavllo, Christoph Feichtenhofer, David Grangier, and Michael Auli. 3d human pose estimation in video with temporal convolutions and semi-supervised training, 2019. 5

[42] Ziqiao Peng, Haoyu Wu, Zhenbo Song, Hao Xu, Xiangyu Zhu, Jun He, Hongyan Liu, and Zhaoxin Fan. Emotalk: Speech-driven emotional disentanglement for 3d face animation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 20687–20697, 2023. 1

[43] Ryan Po, Wang Yifan, Vladislav Golyanik, Kfir Aberman, Jonathan T Barron, Amit H Bermano, Eric Ryan Chan, Tali Dekel, Aleksander Holynski, Angjoo Kanazawa, et al. State of the art on diffusion models for visual computing. *arXiv preprint arXiv:2310.07204*, 2023. 3

[44] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv*, 2022. 3

[45] Sigal Raab, Inbal Leibovitch, Peizhuo Li, Kfir Aberman, Olga Sorkine-Hornung, and Daniel Cohen-Or. Modi: Unconditional motion synthesis from diverse data, 2022. 9

[46] Zhiyuan Ren, Zhihong Pan, Xin Zhou, and Le Kang. Diffusion motion: Generate text-guided 3d human motion by diffusion model, 2023. 3, 6

[47] Alexander Richard, Michael Zollhofer, Yandong Wen, Fernando de la Torre, and Yaser Sheikh. MeshTalk: 3D Face Animation from Speech using Cross-Modality Disentanglement. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1153–1162, Montreal, QC, Canada, 2021. IEEE. 1, 3, 6

[48] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. 2, 3, 5

[49] A. Rössler, D. Cozzolino, L. Verdoliva, C. Riess, J. Thies, and M. Nießner. Faceforensics++: Learning to detect manipulated facial images. *ICCV 2019*, 2019. 10

[50] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models, 2022. 3

[51] Steffen Schneider, Alexei Baevski, Ronan Collobert, and Michael Auli. wav2vec: Unsupervised pre-training for speech recognition. In *Interspeech 2019, 20th Annual Conference of the International Speech Communication Association, Graz, Austria, 15-19 September 2019*, pages 3465–3469. ISCA, 2019. 3

[52] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oran Ashual, Oran Gafni, Devi Parikh, Sonal Gupta, and Yaniv Taigman. Make-a-video: Text-to-video generation without text-video data, 2022. 3

[53] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015. 3, 5

[54] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 3

[55] Linsen Song, Wayne Wu, Chen Qian, Ran He, and Chen Change Loy. Everybody's talkin': Let me talk as you want. *IEEE Transactions on Information Forensics and Security*, 17:585–598, 2022. 2

[56] Stefan Stan, Kazi Injamamul Haque, and Zerrin Yumak. Facediffuser: Speech-driven 3d facial animation synthesis using diffusion. In *ACM SIGGRAPH Conference on Motion, Interaction and Games (MIG '23), November 15–17, 2023, Rennes, France*, New York, NY, USA, 2023. ACM. 3, 6, 7, 8, 10

[57] Zhiyao Sun, Tian Lv, Sheng Ye, Matthieu Gaetan Lin, Jenny Sheng, Yu-Hui Wen, Minjing Yu, and Yong jin Liu. Diffposetalk: Speech-driven stylistic 3d facial animation and head pose generation via diffusion models, 2023. 2, 3, 4, 5

[58] 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. 2

[59] 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. 3

[60] Guy Tevet, Sigal Raab, Brian Gordon, Yoni Shafir, Daniel Cohen-or, and Amit Haim Bermano. Human motion diffusion model. In *The Eleventh International Conference on Learning Representations*, 2023. 3, 4

[61] Guy Tevet, Sigal Raab, Brian Gordon, Yoni Shafir, Daniel Cohen-or, and Amit Haim Bermano. Human motion diffusion model. In *The Eleventh International Conference on Learning Representations*, 2023. 3, 5, 9

[62] Balamurugan Thambiraja, Ikhsanul Habibie, Sadegh Aliakbarian, Darren Cosker, Christian Theobalt, and Justus Thies. Imitator: Personalized speech-driven 3d facial animation. In *Proceedings of the IEEE/CVF International Conference on*

*Computer Vision (ICCV)*, pages 20621–20631, 2023. 2, 3, 4, 5, 6, 7, 8, 10

[63] Justus Thies, Mohamed Elgharib, Ayush Tewari, Christian Theobalt, and Matthias Nießner. Neural voice puppetry: Audio-driven facial reenactment. *ECCV 2020*, 2020. 2, 3

[64] 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. 2

[65] Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic. Realistic speech-driven facial animation with gans. *International Journal of Computer Vision*, 128(5):1398–1413, 2020. 2

[66] Kaisiyuan Wang, Qianyi Wu, Linsen Song, Zhuoqian Yang, Wayne Wu, Chen Qian, Ran He, Yu Qiao, and Chen Change Loy. Mead: A large-scale audio-visual dataset for emotional talking-face generation. In *ECCV*, 2020. 3

[67] S Wang, L Li, Y Ding, C Fan, and X Yu. Audio2head: Audio-driven one-shot talking-head generation with natural head motion. In *International Joint Conference on Artificial Intelligence*. IJCAI, 2021. 2

[68] Jinbo Xing, Menghan Xia, Yuechen Zhang, Xiaodong Cun, Jue Wang, and Tien-Tsin Wong. Codetalker: Speech-driven 3d facial animation with discrete motion prior. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12780–12790, 2023. 2, 4, 5, 6, 7, 10

[69] Jinbo Xing, Menghan Xia, Yuechen Zhang, Xiaodong Cun, Jue Wang, and Tien-Tsin Wong. Codetalker: Speech-driven 3d facial animation with discrete motion prior. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12780–12790, 2023. 1, 3

[70] Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Yingxia Shao, Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and applications. *arXiv preprint arXiv:2209.00796*, 2022. 3

[71] 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

[72] 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. 2

[73] Xin Yu, Yuan-Chen Guo, Yangguang Li, Ding Liang, Song-Hai Zhang, and Xiaojuan Qi. Text-to-3d with classifier score distillation, 2023. 3

[74] Ye Yuan, Jiaming Song, Umar Iqbal, Arash Vahdat, and Jan Kautz. Physdiff: Physics-guided human motion diffusion model. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023. 3

[75] Chenxu Zhang, Saifeng Ni, Zhipeng Fan, Hongbo Li, Ming Zeng, Madhukar Budagavi, and Xiaohu Guo. 3d talking face with personalized pose dynamics. *IEEE Transactions on Visualization and Computer Graphics*, 2021. 1

[76] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models, 2023. 3

[77] Mingyuan Zhang, Zhongang Cai, Liang Pan, Fangzhou Hong, Xinying Guo, Lei Yang, and Ziwei Liu. Motiondiffuse: Text-driven human motion generation with diffusion model. *arXiv preprint arXiv:2208.15001*, 2022. 3, 5

[78] 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

[79] Daquan Zhou, Weimin Wang, Hanshu Yan, Weiwei Lv, Yizhe Zhu, and Jiashi Feng. Magicvideo: Efficient video generation with latent diffusion models, 2023. 3

[80] 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. 2

[81] Wojciech Zielonka, Timo Bolkart, and Justus Thies. Instant volumetric head avatars. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4574–4584, 2022. 10

[82] Wojciech Zielonka, Timo Bolkart, and Justus Thies. Towards metrical reconstruction of human faces. ECCV, 2022. 5