# FaceDiffuser: Speech-Driven 3D Facial Animation Synthesis Using Diffusion

Stefan Stan
Utrecht University
Utrecht, The Netherlands
st.stan96@gmail.com

Kazi Injamamul Haque
Utrecht University
Utrecht, The Netherlands
k.i.haque@uu.nl

Zerrin Yumak
Utrecht University
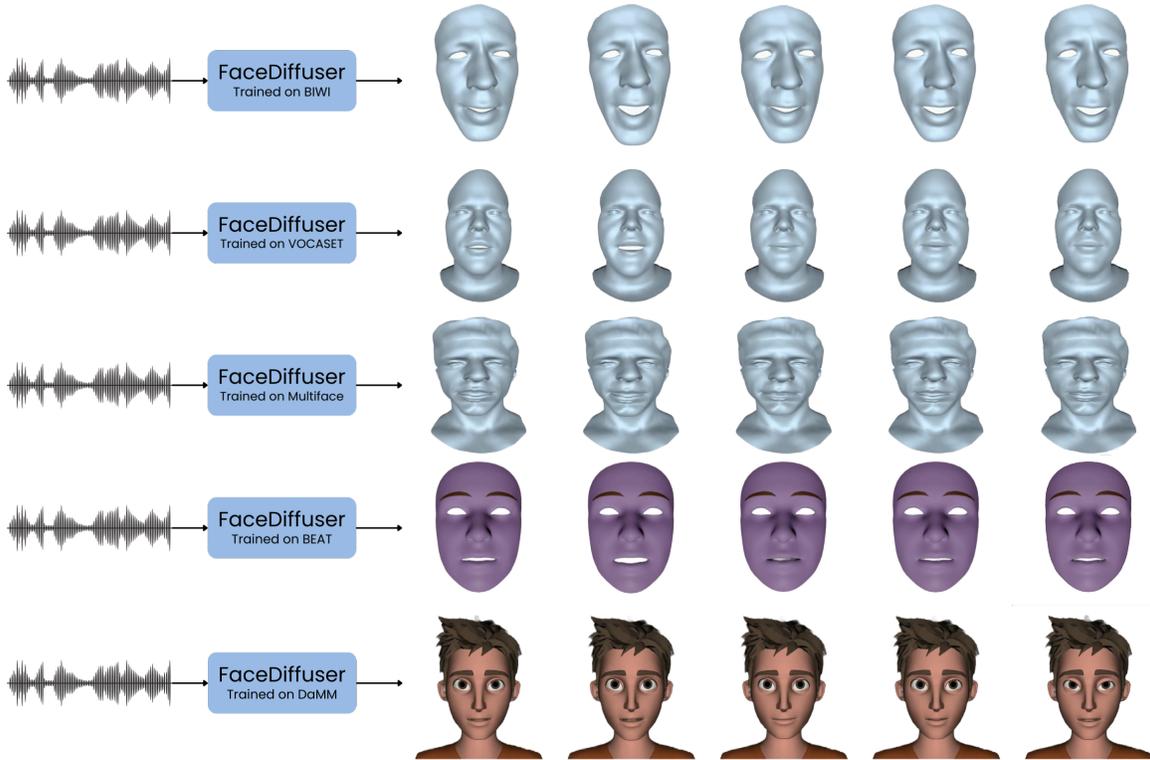Utrecht, The Netherlands
z.yumak@uu.nl

Figure 1: We present FaceDiffuser, an end-to-end non-deterministic neural network architecture for speech-driven 3D facial animation synthesis. Our proposed approach produces realistic and diverse animation sequences and is generalizable to both temporal 3D vertex based mesh animation datasets (top 3 rows) and temporal blendshape based datasets (bottom 2 rows).

## ABSTRACT

Speech-driven 3D facial animation synthesis has been a challenging task both in industry and research. Recent methods mostly focus on deterministic deep learning methods meaning that given a speech input, the output is always the same. However, in reality, the non-verbal facial cues that reside throughout the face are non-deterministic in nature. In addition, majority of the approaches focus on 3D vertex based datasets and methods that are compatible with existing facial animation pipelines with rigged characters is scarce. To eliminate these issues, we present FaceDiffuser, a non-deterministic deep learning model to generate speech-driven facial animations that is trained with both 3D vertex and blendshape based datasets. Our method is based on the diffusion technique and uses the pre-trained large speech representation model HuBERT to encode the audio input. To the best of our knowledge, we are the first to employ the diffusion method for the task of speech-driven 3D facial animation synthesis. We have run extensive objective and subjective analyses and show that our approach achieves better or comparable results in comparison to the state-of-the-art methods. We also introduce a new in-house dataset that is based on a blendshape based rigged character. We recommend watching the accompanying supplementary video. The code and the dataset will be released publicly[1].

## CCS CONCEPTS

• **Computing methodologies** → **Neural networks**; **Animation**;
• **Human-centered computing** → *User studies*.

## KEYWORDS

facial animation synthesis, deep learning, virtual humans, mesh animation, blendshape animation

---

[1]https://github.com/uuembodiedsocialai/FaceDiffuser

# 1 INTRODUCTION

3D facial animation is an important component in various applications such as gaming and XR. Generating facial animations is a tedious task, often done by experienced technical animators using keyframes or done by capturing and retargeting a performer's expression to a rigged model with blendshapes and facial controls. The former requires time and expertise to achieve while the latter requires specialised capture systems [11, 13, 16]. Recent pipelines such as Unreal Engine MetaHuman Animator [15] aims to provide more scalable facial animation capture solutions. Research on speech-driven 3D facial animation can be divided into phoneme-based procedural [6, 26] and data-driven approaches including recent deep learning based methods using intermediary representations of speech units [47, 58] and end-to-end deep learning [8, 17, 22, 29, 42, 54] eliminating the need for intermediary representations. Phoneme-based methods require explicit definition of co-articulation rules and manual work. Deep learning methods automatically discover patterns that can map new speech input to output animations. Our work is an end-to-end speech-driven 3D animation method that uses a diffusion based deep learning model to generate facial animations non-deterministically.

End-to-end deep learning based approaches effectively generate whole face animation producing promising results for accurate lip-sync and upper face animations. Karras et al. [29] proposes an end-to-end CNN based method mapping input waveforms to the 3D vertex coordinates. VOCA [8] proposes a CNN based approach that takes advantage of the pre-trained DeepSpeech [21] model including identity control. MeshTalk [42] learns a latent space representation of facial expressions by employing a U-Net autoencoder network and generates animations based on audio together with a template mesh as input. With the success of transformers, Fan et al. [17] proposed FaceFormer, a self-supervised representation learning model which employs Wav2Vec 2.0 as a speech encoder. Haque and Yumak [22] proposed FaceXHuBERT, an efficient end-to-end encoder-decoder model based on the large speech model HuBERT [25] including identity and emotion control. None of these methods take into account the non-deterministic nature of facial animations meaning given the same speech input, they produce the same results. However, the non-verbal facial cues that reside throughout the face are non-deterministic in nature [34].

There are a few works that employ non-deterministic methods. Learning2Listen [34] generates facial animation for the listening party in two-party dyadic interactions and uses a Vector Quantised Variational Auto Encoder (VQ-VAE) [51]. CodeTalker [54] incorporates self-supervised Wav2Vec 2.0 inspired by FaceFormer [17] together with the idea of having a latent codebook using VQ-VAE inspired by [34]. TalkShow [56] also employs VQ-VAE to generate both face, upper body and hand animations. Recent work on non-deterministic body motion synthesis such as [1, 49, 50] employs diffusion models. However, diffusion models have not yet been used in the speech-driven 3D facial animation domain and to our knowledge we are the first to do that. Our model employs a HuBERT [25] audio encoder together with a specialised diffusion model and produces facial animation both for 3D vertex based and rigged characters. Fig. 1 shows animations generated with our model using five different datasets. Our extensive objective and subjective analysis

shows that our model produces better or comparable results in comparison to state-of-the-art methods. The main contributions of our work are enumerated below:

- FaceDiffuser is the first to incorporate the diffusion mechanism for speech-driven 3D facial animation synthesis task.
- Our model performs better than the state-of-the-art methods in terms of objective metrics on multiple temporal high dimensional 3D vertex based mesh animation datasets.
- We extend our approach to show that the proposed model can generalise to lower dimensional blendshape and facial control based datasets. A new in-house built facial controls based facial animation dataset for rigged characters is also introduced.
- Extensive qualitative analysis and ablation studies were presented to demonstrate the importance of the diffusion mechanism and the ability to synthesise high-quality, diverse facial animation sequences with a discussion on the capabilities and limitations of deterministic and non-deterministic approaches.

# 2 RELATED WORK

Speech-driven facial animation can be classified into two categories: (i) neural rendering of 2D talking heads which resides in the pixel space [27, 32, 57] and (ii) 3D speech-driven animation synthesis using temporal 3D vertex data [8, 17, 22, 42, 54] and blendshape data for a rigged character [2, 58]. Another line of research focuses on 3D reconstruction of faces from 2D videos [9, 20, 44] however they are not speech-driven. For an extensive survey on 3D face reconstruction, tracking and morphable models, we refer to [14, 33].

In this paper, we focus on the 3D speech-driven facial animation synthesis problem using a diffusion model. Therefore in the following two sub-sections, we first present the state-of-the-art on speech-driven 3D facial animation and then motion synthesis using diffusion models.

## 2.1 Speech-driven 3D Facial Animation

3D speech-driven facial animation typically uses phoneme-based procedural approaches [6, 26]. Although these methods come with the advantage of animation control and easy integration to artist-friendly pipelines, they are not fully automatic and require defining explicit rules for co-articulation. Another line of research uses machine learning [48] or graph-based approaches [5] to learn speech-animation mappings from data. These methods rely on blending between speech units and cannot capture the complexity of the dynamics of visual speech [47]. Recent approaches on 3D speech animation synthesis effectively employ deep learning models [2, 8, 17, 22, 29, 42, 47, 58]. Taylor et al. [47] proposes a sliding window approach instead of an RNN focusing on capturing coarticulation effects. VisemeNet [58] builds upon the viseme-based JALI [26] model and combines this with an LSTM-based neural network. However, these two methods [47, 58] still rely on intermediary representations of phonemes and they focus on the mouth movement only. Most previous works do not include automatic tongue animation except [2]. Some methods use 3D face reconstruction methods from in-the-wild videos to generate their data, e.g. dyadic speech-driven facial animation [28, 34]. However, these methods are prone to 3D reconstruction errors. Most of the deep learning based approaches are based on 3D vertex based datasets

[8, 17, 22, 29, 42] and are not compatible with traditional animation pipelines with rigged characters except a few examples such as [2, 47, 58].

Closest to our work are [2, 8, 17, 22, 29, 42, 47, 58]. Karras et al. [29] proposed an end-to-end method using CNNs aiming to resolve the ambiguity in mapping between audio and face by introducing an additional emotion component to the network. However, the method is not trained on multiple speakers and cannot handle identity variations. Instead, Cudeiro et al. [8] presents the audio-driven facial animation method VOCA that generalizes to new speakers eliminating the need for retargeting. However, VOCA fails to realistically synthesise upper face motion and does not include emotional variations. Richard et al. [42] aims for audio-driven animation that can capture variations in multiple speakers including a large dataset of subjects. They address the problem of lack of upper face motions using a categorical latent space that disentangles audio-correlated and audio-uncorrelated information based on a cross-modality loss. Fan et al. [18] proposes an audio and text-driven facial animation method that incorporates the large language model GPT-2 [39]. FaceFormer [17] uses a self-supervised pretrained speech model that addresses the scarcity of available data in existing audio-visual datasets. The model uses a modified version of transformers to handle longer sequences of data. FaceXHuBERT [22] proposes a more efficient network and incorporates HuBERT [25] as the audio encoder as well as includes identity and emotion control.

However, these methods do not take into account the non-deterministic nature of facial animations. Learning2Listen [34] generates facial expressions non-deterministically in two-party dyadic interactions and uses a Vector Quantised Variational Auto Encoder (VQ-VAE) [51] to generate facial animation for the listening party. CodeTalker [54] incorporates self-supervised Wav2Vec 2.0 inspired by FaceFormer [17] and a modified version of VQ-VAE inspired by Learning2Listen [34]. TalkShow [56] also employs VQ-VAE to generate face, upper body and hand animations. None of these methods incorporate the diffusion process to generate a variety of 3D facial animations driven by speech input.

## 2.2 Diffusion for Motion Synthesis

The concept of diffusion was introduced in 2015 by Sohl-Dickstein et al. [45] and is based on a concept in non-equilibrium thermodynamics. The idea is that a sample from the data distribution is gradually noised by the diffusion process and then a neural model learns the reverse process of gradually denoising the sample [49]. It is widely used in the computer vision domain by denoising images noised e.g. with a Gaussian noise and a neural network is trained to reverse the diffusion process [23, 46]. It was used successfully in text-to-image generation tasks leading to examples such as DALL-E2 [40] and StableDiffusion [43]. For a survey of diffusion models applied in the image domain, we refer to [7].

3D body motion generation and 2D talking face generation are the closest work we found in the literature to our work with respect to the use of diffusion process. In the domain of 3D body motion synthesis, Tevet et al. [49] proposed Human Motion Diffusion Model (MDM), a model that can generate body animations

based on text descriptions. They employ a transformer-based architecture and introduce the noised ground truth motions as an additional input to the network. By doing this, they succeed in generating non-deterministic animations at inference time. With the success of MDM, other works generate body animations given music and audio as input [1, 50, 55]. Tseng et al. [50] train a model that can generate dance animations conditioned on music, while Yang et al. [55] use a similar approach for speech-driven gesture motion synthesis. Alexanderson et al. [1] apply diffusion both for co-speech gesture and dance motion generation. In the domain of 2D talking faces, a speech driven video editing method is proposed by Bigioi et al. [4]. By taking a template video as input along with a new speech segment, the model generates new lip motions that follows the target speech sequence. The model is capable of generalising across different speaker identities. DAE-Talker [12] makes use of diffusion for generating talking head animation with a 2-stage learning process. They employ a diffusion autoencoder approach initially introduced by Preechakul et al. [36] on images and extend it to generating videos. They first train a diffusion autoencoder that learns the latent space of facial expressions from the training data. The first stage has no temporal awareness and only learns to reconstruct an image from its encoding and its noised representation. In the second stage, a transformer-based encoder-decoder architecture is used to encode the audio input and output frame embeddings. To our knowledge, no work in the literature apply diffusion models for the 3D speech-driven facial animation task.

## 3 PROBLEM FORMULATION

Let $A$ be an audio input associated with a sequence of ground truth frames $x_0^{1:N} = (x_0^1, x_0^2, ..., x_0^N)$, where $N$ is the number of visual frames sampled at a certain FPS specific to datasets. Each frame in the sequence $x_0^n$ is represented as an array of vertex positions with the length $V$ x 3, where $V$ is the number of mesh vertices in the topology, and 3 is the number of spatial axes. In the case of blendshape or facial control datasets, $x_0^n$ represents a vector of rig controls or blendshape values, having the length $C$, the number of controls or blendshapes driving the facial rig. Based on audio input $A$, the goal of our architecture is to predict an animation sequence, $\hat{x}_0^{1:N}$ that resemble the ground truth frames $x_0^{1:N}$. Additionally, the predictions will be guided by a style $S$ in the form of a one-hot vector with a length equal to the number of training subjects (for vertex based datasets only in our experiments) and noise $x_T^{1:N}$ sampled from the normal distribution $\mathcal{N}(0, 1)$ and with the same shape as the ground truth sequence $x_0^{1:N}$. It is to be noted that, models trained on vertex data generates animations in the form of vertex displacements with respect to neutral face templates. Whereas, models trained on blendshape or facial control does not require this step as the neutral face of the rigged face is not a variable in our setting.

The abstraction of the formulated problem is presented with the following equation:

$$\hat{x}_0 = \text{FaceDiffuser}(\mathcal{A}, x_t, t) \tag{1}$$

Where, $\hat{x}_0$ is the predicted animation sequence, $\mathcal{A}$ is the input audio sequence and $x_t$ is the sequence $x_0$ after t diffusion steps. Here, $t = T$, $x_t = \sigma^{1:N}$ drawn from $\mathcal{N}(0, 1)$.
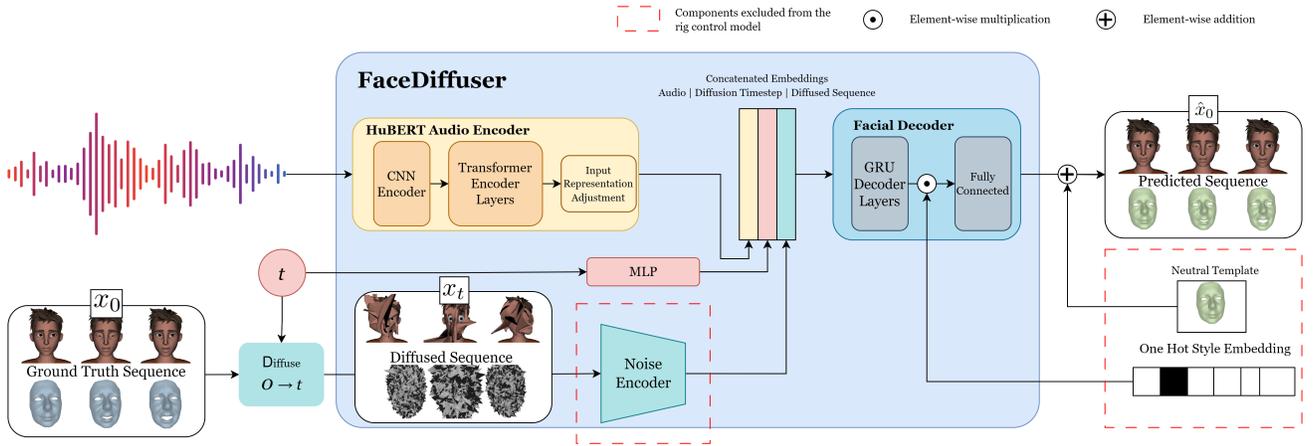
Figure 2: FaceDiffuser learns to denoise facial expressions and generate animations based on input speech. Audio speech embeddings from pre-trained HuBERT model combined with embeddings from the noised ground truth animation sequence are used to train the Facial Decoder. The Facial Decoder is comprised of a sequence of GRU layers followed by a fully connected layer and learns to predict (i) vertex displacements or (ii) rig control (blendshape) values. The predicted sequence $\hat{x}_0$ is compared with the ground truth sequence $x_0$ by computing the loss, which is then backpropagated to update the model parameters.

## 4 PROPOSED APPROACH

### 4.1 Training

We propose a general model that can be employed for both vertex based and blendshape based datasets with slight modifications in terms of hyperparameters. We refer to the vertex based model configuration as V-FaceDiffuser, and to the blendshape based model as B-FaceDiffuser. The main difference being the additional *Noise Encoder* as it can be seen in Fig. 2 enclosed in dashed red box. The noise encoder helps to project high dimensional vertex data into low dimensional latent representation. The diffusion noising process takes in $x_0^{1:N}$ to compute noised $x_t^{1:N}$, retaining its original shape.

From Fig. 2, we can identify the following main components that are included in both versions of the model:

**Audio Encoder:** We use the pretrained large speech model, HuBERT as the audio encoder similar to [22] and it is kept the same in both versions of the architecture. We employ a pre-trained version of the HuBERT architecture and use the released *hubert-base-ls960* version of it, which was trained on 960 hours of LibriSpeech[35] dataset.

**Diffusion Process:** Let $x_0^{1:N}$ be a sequence of ground truth visual frames from the dataset with shape $(N, C)$, where $C$ is either the number of vertices, multiplied by the 3 (for 3 spatial axes), or the number of rig facial controls (or blendshape values). During training, we randomly sample an integer timestep $t$ from $[1, T]$, indicating the number of noising steps to be applied to $x_0^{1:N}$ and to obtain $x_t^{1:N}$ with the formula:

$$x_t^{1:N} = q(x_t^{1:N}|x_{t-1}^{1:N}) = \mathcal{N}(\sqrt{1-\beta_t}x_{t-1}^{1:N}, (\beta_t)\mathcal{I}) \quad (2)$$

Where, $N$ is the number of visual frames in the sequence, $t$ is the diffusion timestep and $\beta_t$ is the constant noise at timestep $t$ such that $0 < \beta_1 < \beta_2 < ... < \beta_T < 1$.

After the forward noising process, ideally, we want to be able to compute the reverse process and go backwards from $x_T^{1:N} \sim \mathcal{N}(0, 1)$ to $x_0^{1:N}$. Therefore the conditional distribution function $p(x_{t-1}^{1:N}|x_t^{1:N})$ needs to be known beforehand. *Ho et al.* [24] proposes to achieve that by learning the latent representation variance of the dataset. The training objective is defined as learning to predict the noise $\epsilon$ that was added to the input $x_0$. However, we deviate from [24] and follow MDM [49] and EDGE [50], choosing for our model to learn to predict actual animation data instead of the noise level in the data. We consider this to be more suitable for our task since the results are also conditioned on the input audio. Furthermore, by choosing this approach, our model is able to predict acceptable results even from the first denoising steps of the inference process, allowing for faster sampling. However following the full inference process would give the best results.

We employ a simple loss for training similar to [49] and [50]. More thorough experimentation was conducted by *Ho et al.* [24], who also claim that utilising the simple loss for learning the variational bound proved to be both easier to implement and also advantageous for the quality of the sampled results. The loss is defined as:

$$\mathcal{L} = E_{x_0 \sim q(x_0|c), t \sim [1,T]}[\|x_0 - \hat{x}_0\|] \quad (3)$$

**Facial Decoder:** The facial decoder is responsible for producing the final animation frames from latent representation of the encoded audio and noise. It is comprised of multiple GRU layers followed by a final fully connected layer that predicts the output sequence. During the decoding step, a style embedding can also be added in the form of an element-wise product between a learned style embedding vector and the hidden states output. We explain the choice of the GRU decoder in the ablation section.
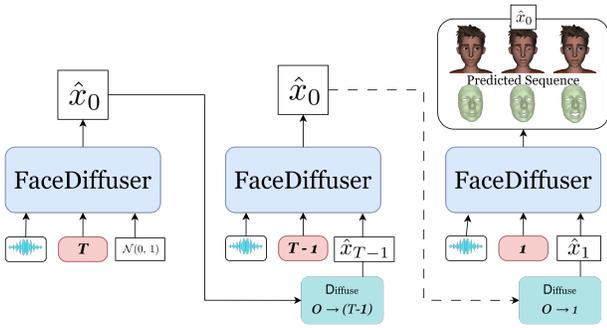
Figure 3: FaceDiffuser inference is an iterative process from T decreasing to 1. The initial noise being represented by actual noise from the normal distribution $\mathcal{N}(0, 1)$. At each step, we provide the network with the audio and noised animation input. The predicted motion is then diffused again and fed to the next step of the iteration.

## 4.2 Inference

The inference is an iterative process, during which we go through all of the diffusion timesteps backwards from $T$ to 1, gradually improving the prediction at each inference step. During inference, the ground truth noised sequence as input to the model is missing compared to training time. Therefore, at inference time $T$, a randomly sampled noise is provided to the model from the uniform distribution $\mathcal{N}(0, 1)$. Fig. 3 depicts the inference process.

## 5 EXPERIMENTAL SETUP

For our experiments, we trained our model on 3D vertex datasets. We use BIWI [19] as our primary vertex based dataset for comparisons against the state-of-the-art methods such as VOCA [8], FaceFormer [17] and CodeTalker [54]. The dataset contains audio-4D scan pairs of 14 subjects, each uttering 40 sentences twice- with neutral emotion and with emotional expressions. We adopt the exact same dataset split for BIWI as done in [17, 54] and only use the emotional sequences subset. The split results into training set *BIWI-Train* that contains 192 sentences and validation set *BIWI-Val* that contains 24 sentences from 6 training subjects. There are two test sets: BIWI-Test-A, containing 24 sentences from seen subjects and BIWI-Test-B, containing 32 sentences from the 8 remaining unseen subjects. The former test set is used for computing objective metrics while both facilitate the qualitative analysis and the perception study. BIWI-Test-B is further used to compute the diversity metric that we define in the next section. We also train all the methods on VOCASET, adopting the VOCASET-Train and VOCASET-Test split following [8, 17, 54] and similar to these works, we use the test set to generate animations for the perceptual user study. In addition, we also employ the Multiface dataset [53] which was also used in MeshTalk [42] to demonstrate the generalisability capability of our proposed method and compare our results to two state-of-the-art methods. Among these three datasets, only Multiface has proper eye-blinks while for VOCASET, there are no examples of eye-blinks and for BIWI, the face topology does not contain proper eye-lids. In our experiments, we used BIWI and MultiFace for objective analysis and BIWI and VOCASET for perceptual studies with users. For the

| Hyperparameter | V-FaceDiffuser | B-FaceDiffuser |
|---|---|---|
| Optimizer | Adam | Adam |
| Learning Rate | $1e^{-4}$ | $1e^{-4}$ |
| Number of epochs | 50 | 100 |
| Diffusion Steps | 500 | 1000 |
| $\beta$ schedule | linear | linear |
| Input Embedding Dim | 512(B); 256(V,M) | N/A |
| Number of GRU Layers | 2 | 4(U);2(B) |
| GRU hidden size | 512 | 1024(U); 256(B) |
| GRU dropout | 0.3 | 0.3 |

Table 1: Hyperparameter values of our proposed approach. For V-FaceDiffuser, since the vertex data is high dimensional, we embed it to a latent dimension. The input embedding dimension is 512 for BIWI (represented as 512(B)) while for VOCASET and Multiface, it is 256 (i.e.- 256(V,M)). B-FaceDiffuser does not need this projection as the data is low dimensional. Different number of GRU layers and hidden sizes were used for BEAT and UUDaMM, denoted by (B) and (U) respectively.

qualitative analysis, all datasets were used. For clarity, we defined the models trained on vertex data in V-FaceDiffuser configuration.

In addition, we trained our proposed method on blendshape based datasets such as BEAT [31] and our in-house dataset, UUDaMM (Utrecht University Dyadic Multimodal Motion Capture Dataset) for a rigged character. While both datasets include both face and body animations, we use the facial data together with the synchonously captured audio data only. While BEAT facial animations are based on Apple ARKIT 52 blendshape standard, UUDaMM dataset includes AutoDesk Maya facial controls. These are much lower dimensional datasets compared to the vertex based ones. Since there has not yet been any speech-driven facial animation work done with these two datasets in the literature to our knowledge, we compare the results of our approach with a baseline method that is identical to the proposed method without the diffusion component. Moreover, because the compared state-of-the-art methods are designed and proposed for vertex based datasets specifically, employing those on the blendshape datasets for direct comparison is not applicable. For BEAT, we use a subset ($\approx$16 hour data by native English speakers) of the full ($\approx$76 hours) dataset. UU-DaMM ($\approx$ 10 hours) consists of multimodal motion capture data of 2 actors interacting in a natural dyadic setting. The dataset contains full body motion, captured with Vicon [52], facial performance capture with Dynamixyz [13], synchronised audio recording and reference videos. The facial performance data was solved and retargeted to a publicly available model - Ray[41]. We then export the temporal facial control values (where the facial controls drive the artist generated blendshapes) to form the training dataset. We defined the models trained on lower dimensional blendshape based datasets as B-FaceDiffuser. Eye-blinks are present in both these datasets while UUDaMM also includes eye gaze in the training data. More details on these 5 datasets we used are available in the supplementary material.

**Implementation Details:** All the model training was done on a shared compute cluster running Linux with AMD EPYC 7313 CPU,

| Dataset | Method | MVE ↓ x10$^{-3}$ mm | LVE ↓ x10$^{-4}$ mm | FDD ↓ x10$^{-5}$ mm | Diversity↑ x10$^{-3}$ mm |
|---------|--------|------|------|------|------|
| BIWI | VOCA | 8.3606 | 6.7155 | 7.5320 | 7.8507 |
| | FaceFormer | 7.1658 | 4.9847 | 5.0972 | 5.9201 |
| | CodeTalker | 7.3980 | 4.7914 | 4.1170 | 0.0003 |
| | V-FaceDiffuser | **6.8088** | **4.2985** | **3.9101** | **9.2459** |
| Multiface | VOCA | 15.782 | 25.067 | 14.253 | 0.5292 |
| | FaceFormer | 7.6132 | 7.0770 | **5.0127** | 14.745 |
| | CodeTalker | 12.170 | 20.392 | 6.6857 | 13.423 |
| | V-FaceDiffuser | **7.0004** | **6.2295** | 5.9020 | **15.500** |

**Table 2: Objective results computed over the temporal 3D vertex datasets. Our approach achieves the best results in all four objective metrics for the BIWI dataset. For the Multiface dataset we score best on all metrics except the FDD metric.**

| Dataset | Method | MBE ↓ | LBE ↓ | FDD ↓ |
|---------|--------|-------|-------|-------|
| BEAT | w/o Diffusion | **0.4170** | **0.1077** | 0.1482 |
| | B-FaceDiffuser | 0.5152 | 0.1358 | **0.1471** |
| UUDaMM | w/o Diffusion | **2.6963** | **1.5924** | 2.0553 |
| | B-FaceDiffuser | 3.4479 | 1.6671 | **1.7752** |

**Table 3: Objective evaluation results computed over the blend-shape and controller based datasets. Best results are marked as bold. Our model generates slightly higher error based on frame-level low dimensional GT values while achieving better FDD.**

Nvidia A16 GPU, 1TB RAM. Tab. 1 shows the hyperparameters we use for the proposed approaches.

## 6 RESULTS

We evaluate our proposed approach quantitatively, qualitatively and with a perceptual user study and compared our results to the state-of-the-art methods. In the following subsections, we will present and discuss the results.

### 6.1 Quantitative Evaluation

Following [17, 22, 54], we employ the lip vertex error (LVE for V-FaceDiffuser, LBE for B-FaceDiffuser) in order to measure the lip-sync error. Similar to [54], we also adopt the FDD metric that gives an indication of the upper face motion variation in terms of statistics and how close it is to the variation observed in the ground truth. Additionally, we also compute mean full-face vertex error (MVE for V-FaceDiffuser, MBE for B-FaceDiffuser) as we are interested in not only the lip-sync but also the motion that resides throughout the entire face. We use the exact same set of lip and upper face vertices to compute the mentioned metrics as provided in code repositories of [54] and [42] for BIWI and Multiface respectively. In order to demonstrate the diversity capability of our proposed model, we introduce a novel diversity metric that is subject to animation generated conditioned on different training subjects for the vertex based datasets. For the blendshape based datasets, we manually select the blendshapes related to lip and upper face movements for LBE and FDD respectively.

***MVE and MBE***. Mean vertex (or blendshape) error measures the deviation of all the face vertices (or all the blendshapes) of a sequence with respect to the ground truth by computing the maximal L2 error for each frame and by taking a mean over all corresponding frame pairs.

***LVE and LBE***. Lip vertex (or blendshape) error is identical to MVE (or MBE). We only consider the lip vertices (or blendshapes related to lips) for computing the metrics.

***FDD***. Introduced in [54], it measures the variation of facial dynamics for a motion sequence against ground truth. It gives an

indication of how close the standard deviation (or upper face motion variation) of a generated sequence is compared to the observed variation in ground truth.

***Diversity***. We introduce an additional metric that measures the model's ability to produce diverse animations. With the introduction of diffusion, our goal is to develop a model able to generate a great range of motions and non-deterministic expressions for the regions of the face that are uncorrelated or loosely correlated to speech. Based on similar metrics used for diffusion models in literature [49], we introduce a novel diversity metric for the face. The diversity metric could be applied to different iterations of the inference algorithm, however, to be compatible with deterministic state-of-the-art baselines, we define the diversity across different subjects. Therefore, we define the diversity metric as follows. Let $\hat{x}_0^s$ be a generated sequence, conditioned on subject $s \in S$, where $S$ is the list of training subjects. We compute the mean vertex difference between every 2 sequences conditioned on different subjects. We then take the mean of these differences and define the diversity of a sequence as follows:

$$Diversity = \frac{\sum_{i=1}^{|S|-1} \sum_{j=i+1}^{|S|} \|\hat{x}_0^i - \hat{x}_0^j\|}{\frac{(|S|-1)\cdot|S|}{2}} \qquad (4)$$

where $S$ denotes the list of the training subjects and $\hat{x}_0^i$ is the predicted animation sequence conditioned on the $i$th subject from $S$. In our training setting, the diversity metric is only suitable for the vertex based datasets that allow using neutral face template meshes that are different across subjects. This ensures different subjects have different neutral facial physiognomy. Whereas for the blend-shape based datasets, the neutral configuration of the blendshape values remains identical for different actors/speakers. Hence, B-FaceDiffuser is trained without the subject conditioning and we evaluate the diversity in terms of animation graphs.

***Discussion***. Tab. 2 and Tab. 3 show the objective results for V-FaceDiffuser and B-FaceDiffuser respectively. Our approach performs better than all the other methods on all the objective metrics for the BIWI dataset. For the Multiface dataset, ours perform the best on all objective metrics but FDD, for which FaceFormer performs slightly better. For the blendshape based datasets, we cannot compare our method with state-of-the-art methods as mentioned earlier. Instead, we compare the diffusion MBE and LBE are higher than the baseline model as our approach encourages randomness
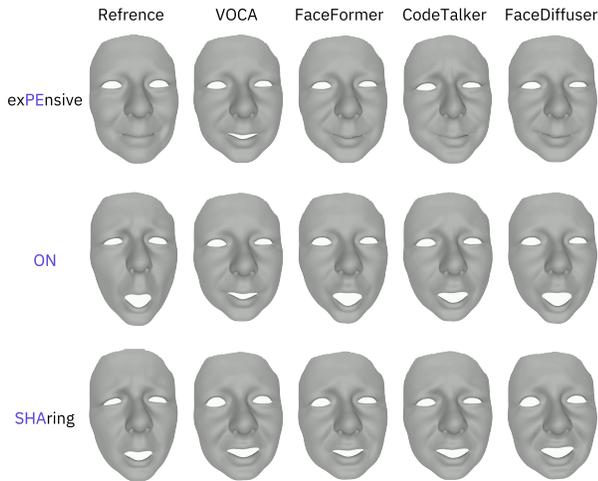
**Figure 4: Visual comparison of frames from synthesised facial animation sequences generated by different methods together with reference frames from GT sequence. The highlighted utterances are represented in visual frames. Our method generates lip shapes that are close to the references while encouraging diversity in the upper face.**

and non-determinism while more resembling the upper face variation observed in the ground truth with lower FDD value. Furthermore, the diversity for B-FaceDiffuser is evaluated qualitatively in Section 6.2.

## 6.2 Qualitative Evaluation

We carried out extensive qualitative analysis of the generated animation sequences and compared them to both ground truth and other methods. Our method generates accurate lip shapes resembling ground truth lip motions while generating diverse upper-face motions. A visual comparison can be observed in Fig. 4. Moreover, our approach is generalisable to unseen speakers, noisy audio input, multiple overlapping speakers in audio, speech in different languages and text-to-speech audio.

In order to qualitatively demonstrate the diversity metric presented in Section 6.1, we sample animation sequences from BIWI-Test-B generated with the same audio but conditioned on different training subjects. Using the generated sequences, we plot the mean and standard deviation of the motion and present it with a heatmap visualisation as shown in Fig. 5. Because there is no subject conditioning for B-FaceDiffuser, we demonstrate the diversity of the generated sequences with animation graphs. We sample our model multiple times using the same audio input and plot the animation graphs of some key facial controls. This can be seen in Fig. 6 where we can observe that upper face controls which are uncorrelated or loosely correlated to speech, do not follow the ground truth whereas the speech correlated lip controls resemble more to the ground truth, especially the peaks in the graph. Furthermore, a high diversity across different results can be observed, especially in the case of eye blinks. For visual judgement, we refer to our supplementary video.
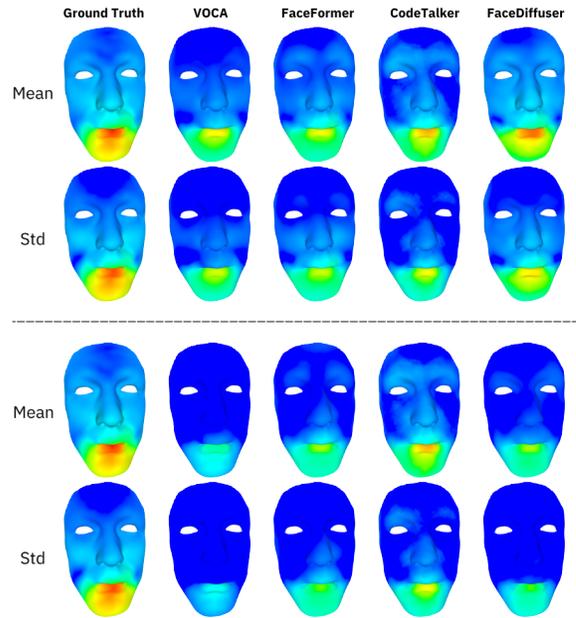


**Figure 5: Mean Motion comparison: Animation sequences were sampled from BIWI-Test-B conditioned on different training subjects. One set of two rows separated by the dashed line depicts motion statistics of the same inference. Here we present two sets of inferences generated using the same audio but conditioned on two different training subjects. We notice that conditioning of different subjects produces diversity in generated animations. Whereas, the other methods are much less diverse as seen in the mean motion heatmaps, where dark blue means less observed motion and bright red means more observed motion.**

## 6.3 Perception Study

In addition to quantitative and qualitative analyses, we also conducted a series of user studies to evaluate the user perception in terms of realism, lip-sync and appropriateness. We adopt a similar A/B testing strategy for the user study as done in the previous state of the arts-[17, 54]. We conducted three separate user studies for the compared models trained on three datasets - BIWI, VOCASET and UUDaMM. As visual renders of Multiface resemble VOCASET while BEAT resembles BIWI, these two were dropped for our perception study. For the studies, we used the rendered videos of the generated animation sequences by different methods and presented randomly ordered pairs of videos in a side-by-side manner. The participants were asked to judge three subjective aspects- realism, lip-sync and appropriateness of the rendered animations with respect to the audio. The user studies were hosted on Qualtrics[38] and participants were recruited using Prolific[37], ensuring proper remuneration for their time. Some additional participants were recruited as well who did the studies voluntarily. In total, there were 83 survey responses for the three studies where 31 participated in the study conducted on BIWI, 29 participated in the study on VOCASET and finally, 23 people did the study on UUDaMM. More

(a) Lower lip facial control.



(b) Eye brow facial control.



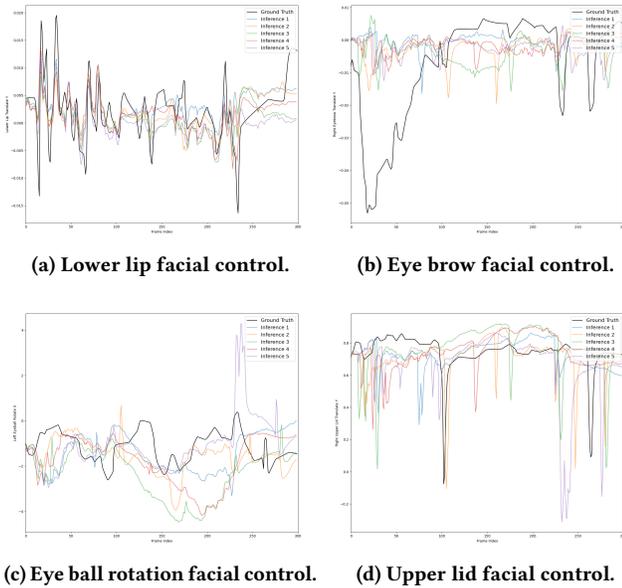(c) Eye ball rotation facial control.



(d) Upper lid facial control.

**Figure 6: Animation graphs of some facial controls (i.e.- low-erlip, eyebrow, gaze, upperlid) of the UUDaMM dataset. We synthesise animation data multiple times using the same audio and plot the graphs together with the ground truth. The black plots depict the ground truth whereas different coloured plots depict different inferences. It is evident that our approach produces lip control values similar to the ground truth as seen in Fig. 6a while encouraging diversity for the other facial controls as seen in Fig. 6b, Fig. 6c and Fig. 6d.**

| Dataset- BIWI | | | |
|---|---|---|---|
| Competitor | Realism | LipSync | Appropriateness |
| VOCA | 77.27 % | 69.32 % | 79.55% |
| FaceFormer | 31.03 % | 34.48 % | 37.93 % |
| CodeTalker | 44.94 % | 44.94 % | 41.57 % |
| GT | 38.71 % | 34.41 % | 36.56 % |
| Dataset- VOCASET | | | |
| VOCA | 76.83 % | 78.05 % | 78.05 % |
| FaceFormer | 49.38 % | 46.91 % | 51.85 % |
| CodeTalker | 27.16 % | 26.40 % | 27.16 % |
| GT | 23.81 % | 33.33 % | 27.38 % |
| Dataset- UUDaMM | | | |
| w/o Diffusion | 63.77 % | 66.67 % | 65.94 % |
| GT | 18.48 % | 27.17 % | 20.65 % |

**Table 4: User study results. We conduct A/B testing and report the percentage of responses where A (i.e. ours) was preferred over B (i.e. competitor) in terms of realism, lip-sync and appropriateness of the rendered animations.**

details about the user studies can be found in the supplementary material.

Tab. 4 shows the results of the 3 surveys we conducted. Our model clearly outperforms VOCA in all three aspects for both BIWI and VOCASET. However, renders of FaceFormer and CodeTalker were perceived better than ours for BIWI whereas for VOCASET ours perform similarly to FaceFormer while worse than CodeTalker. As reported in Tab. 2, unlike CodeTalker which produces the same motion conditioned on different training subjects, our model produces diverse animation. This might have affected the user study results. For example, a generated sequence conditioned on a less expressive training subject will have less motion for our model while for CodeTalker, the motion is not affected by the subject condition, see Fig. 5. Resulting in more motion and expressiveness in the rendered videos for CodeTalker that might have affected the perceived subjective aspects when presented side-by-side for this kind of instances. For UUDaMM, there is a clear preference for the results generated by our final model with diffusion than the baseline model without diffusion. For all three datasets, ground truth was perceived as superior, which is expected.

## 6.4 Ablation Study

To analyse the effects of the different components of our proposed architecture, we experimented with different configurations of it,

by either removing or changing different modules. We conducted ablation studies in terms of (i) diffusion mechanism, (ii) audio encoder and (iii) facial decoder. For the ablation study experiments, we only employ BIWI for vertex based dataset, and UUDaMM for blendshape based dataset.

**Ablation on Diffusion Process:** To understand the contribution of the diffusion mechanism in our proposed model, we train a similar model without diffusion and compare the results. We experimented with both BIWI and UUDaMM datasets. The first segment of Tab. 5 shows that the model without diffusion achieves slightly worse results in terms of the objective metrics. Fig. 7 shows that models with diffusion produce more motion throughout the face, resembling mean motion observed in ground truth.

**Ablation on Audio Encoder:** We use HuBERT as our audio encoder in the proposed model. Following [17], a lot of recent works employed Wav2Vec 2.0[3] as the audio encoder. In order to justify our choice of using HuBERT, we trained our model with Wav2Vec 2.0 as our audio encoder. The second segment in Tab. 5 shows a clear improvement of using HuBERT over Wav2Vec 2.0 to encode audio information for speech-driven downstream tasks of 3D animation synthesis.

**Ablation on Decoder:** In order to motivate our choice of the facial decoder in the proposed model, we carried out ablation study experiments by using different sequence modelling method for the decoder. We replaced the proposed GRU decoder with a simpler RNN decoder and more complex transformer decoder and report the objective results in Tab. 6 for both BIWI and UUDaMM. GRU decoder performs the best out of the tested configurations, resulting in lower error values on both datasets.
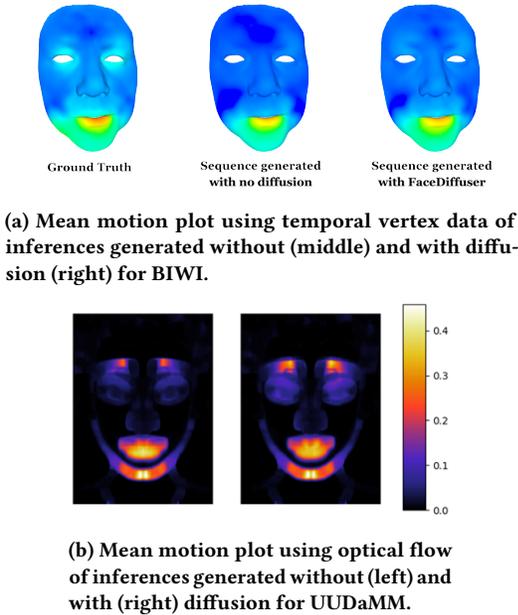
**(a) Mean motion plot using temporal vertex data of inferences generated without (middle) and with diffusion (right) for BIWI.**



**(b) Mean motion plot using optical flow of inferences generated without (left) and with (right) diffusion for UUDaMM.**

**Figure 7: Ablation on the diffusion process. It is evident that the diffusion mechanism encourages more non-verbal cues throughout the face.**

## 7 DISCUSSION AND LIMITATION

Our approach performs objectively better than state of the arts but the perceptual evaluation shows that there is still room for improvement in terms of subjective aspects. There is a limitation in terms of available datasets. We argue that the sequences are too short to capture the full range of expressions a person can have in [8, 19, 53], therefore our model cannot capture long-term context from the data. It would be interesting to see how our model performs against other methods when a larger dataset is used. Furthermore, a dataset comprising diverse captured sequences where for one textual content, subjects perform multiple times, would help to better analyse the diversity capability of non-deterministic models in addition to analysing the non-deterministic nature of non-verbal facial cues uncorrelated or loosely correlated with speech. Only in BIWI, this is available but only twice in terms of binary emotion conditions. Incorporating emotion information can be a future direction to explore, including categorical or continuous emotion models. Using BEAT dataset or creating synthetic datasets using vision-based 4D reconstruction and emotion recognition models such as EMOCA[9] similar to [10], can be a potential future direction to this end. Since our results have limitations in terms of perceived realism, employing a diffusion autoencoder similar to [12] can prove to be a potential future direction to achieve better quality 3D facial animation synthesis in terms of the subjective aspects. Moreover, it would be interesting to see how models proposed for vertex based datasets perform when trained on lower dimensional blendshape datasets. Furthermore, due to the iterative sampling process of the diffusion mechanism, our model's inference time is long and subject to the diffusion timesteps used during training. Therefore, our approach is not suitable for real-time applications.

| Ablation on Diffusion Process | | | | |
|---|---|---|---|---|
| Model | MVE ↓ x10$^{-3}$mm | LVE ↓ x10$^{-4}$mm | FDD ↓ x10$^{-5}$mm | Training Time (m) |
| w/o Diffusion | 6.8833 | 4.5870 | 4.6690 | ≈ 67 |
| FaceDiffuser | 6.8088 | 4.2985 | 3.9100 | ≈ 67 |
| Ablation on Audio Encoder | | | | |
| Model | MVE ↓ x10$^{-3}$mm | LVE ↓ x10$^{-4}$mm | FDD ↓ x10$^{-5}$mm | Training Time (m) |
| Wav2Vec2 | 7.4593 | 5.1590 | 4.1950 | ≈ 67 |
| HuBERT | **6.8088** | **4.2985** | **3.9100** | ≈ 67 |

**Table 5: Objective metrics computed over the test results of BIWI dataset for ablation experiments of the diffusion process and of different audio encoder.**

| Dataset- BIWI | | | | |
|---|---|---|---|---|
| Decoder Type | MVE ↓ x10$^{-3}$ mm | LVE ↓ x10$^{-4}$ mm | FDD ↓ x10$^{-5}$ mm | Training Time (h) |
| GRU | 6.8088 | 4.2985 | 3.9100 | ≈ 1.12 |
| RNN | 7.0833 | 4.7870 | 4.0690 | ≈ 1.12 |
| Transformer(TF) | 9.9767 | 10.1300 | 4.8587 | ≈ 1.34 |
| Transformer(AR) | 7.0213 | 4.6941 | 4.3272 | ≈ 5.00 |
| Dataset- UUDaMM | | | | |
| Decoder Type | MBE ↓ | LBE ↓ | FDD ↓ | Training Time (h) |
| GRU | **3.6791** | **3.5812** | 1.8862 | ≈ 4.5 |
| RNN | 3.8654 | 3.9196 | **1.8426** | ≈ 2.92 |
| Transformer(AR) | 3.6881 | 3.7105 | 1.8533 | ≈ 4.58 |

**Table 6: Objective metrics computed over the BIWI test results for ablation experiments of facial decoder. Here, (TF) and (AR) depict teacher-forcing scheme and autoregressive scheme respectively for Transformer decoders.**

## 8 CONCLUSION

We integrated the diffusion mechanism into a generative deep neural network trained to generate 3D facial animations conditioned on speech. The proposed approach is generalisable to both high dimensional temporal 3D vertex data as well as low dimensional blendshape data with minimal changes. The quantitative analysis shows that our approach performs better than the state of the arts. We showed that our model is able to produce higher diversity of motions between different style conditions than the competitors. Our approach also produces diverse animation sequences for rigged characters that can be observed in animation graphs of multiple generations conditioned on the same audio. Extensive ablation studies support our network architecture design choices, showing the benefits of different proposed components of the neural network architecture.

**Ethical Consideration** Face data can be used for generating content that can jeopardise privacy. We must act responsibly by considering the aspects pertaining to privacy and ethics.

# REFERENCES

[1] Simon Alexanderson, Rajmund Nagy, Jonas Beskow, and Gustav Eje Henter. 2023. Listen, Denoise, Action! Audio-Driven Motion Synthesis with Diffusion Models. *ACM Trans. Graph.* 42, 4 (2023), 1–20. https://doi.org/10.1145/3592458

[2] Mónica Villanueva Aylagas, Héctor Anadon Leon, Mattias Teye, and Konrad Tollmar. 2022. Voice2Face: Audio-driven Facial and Tongue Rig Animations with cVAEs. In *EUROGRAPHICS SYMPOSIUM ON COMPUTER ANIMATION (SCA 2022.*

[3] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems* 33 (2020), 12449–12460.

[4] Dan Bigioi, Shubhajit Basak, Hugh Jordan, Rachel McDonnell, and Peter Corcoran. 2023. Speech Driven Video Editing via an Audio-Conditioned Diffusion Model. *arXiv preprint arXiv:2301.04474* (2023).

[5] Yong Cao, Wen C. Tien, Petros Faloutsos, and Frédéric Pighin. 2005. Expressive Speech-Driven Facial Animation. *ACM Trans. Graph.* 24, 4 (oct 2005), 1283–1302. https://doi.org/10.1145/1095878.1095881

[6] Constantinos Charalambous, Zerrin Yumak, and A.F. van der Stappen. 2019. Audio-driven emotional speech animation for interactive virtual characters. *Computer Animation and Virtual Worlds* 30 (2019).

[7] Florinel-Alin Croitoru, Vlad Hondru, Radu Tudor Ionescu, and Mubarak Shah. 2023. Diffusion models in vision: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2023).

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

[9] Radek Danecek, Michael J. Black, and Timo Bolkart. 2022. EMOCA: Emotion Driven Monocular Face Capture and Animation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. 20311–20322.

[10] Radek Daněček, Kiran Chhatre, Shashank Tripathi, Yandong Wen, Michael J. Black, and Timo Bolkart. 2023. Emotional Speech-Driven Animation with Content-Emotion Disentanglement. arXiv:2306.08990 [cs.CV]

[11] DI4D 2023. DI4D. https://di4d.com/.

[12] Chenpng Du, Qi Chen, Tianyu He, Xu Tan, Xie Chen, Kai Yu, Sheng Zhao, and Jiang Bian. 2023. DAE-Talker: High Fidelity Speech-Driven Talking Face Generation with Diffusion Autoencoder. *arXiv preprint arXiv:2303.17550* (2023).

[13] Dynamixyz 2023. Dynamixyz. https://www.dynamixyz.com.

[14] Bernhard Egger, William A. P. Smith, Ayush Tewari, Stefanie Wuhrer, Michael Zollhoefer, Thabo Beeler, Florian Bernard, Timo Bolkart, Adam Kortylewski, Sami Romdhani, Christian Theobalt, Volker Blanz, and Thomas Vetter. 2020. 3D Morphable Face Models—Past, Present, and Future. *ACM Trans. Graph.* 39, 5, Article 157 (jun 2020), 38 pages. https://doi.org/10.1145/3395208

[15] Epic Games 2023. MetaHuman Animator. https://www.unrealengine.com/en-US/blog/delivering-high-quality-facial-animation-in-minutes-metahuman-animator-is-now-available.

[16] Faceware 2023. Faceware. https://facewaretech.com/.

[17] Yingruo Fan, Zhaojiang Lin, Jun Saito, Wenping Wang, and Taku Komura. 2022. FaceFormer: Speech-Driven 3D Facial Animation with Transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 18770–18780.

[18] Yingruo Fan, Zhaojiang Lin, Jun Saito, Wenping Wang, and Taku Komura. 2022. Joint Audio-Text Model for Expressive Speech-Driven 3D Facial Animation. *Proc. ACM Comput. Graph. Interact. Tech.* 5, 1, Article 16 (may 2022), 15 pages. https://doi.org/10.1145/3522615

[19] Gabriele Fanelli, Juergen Gall, Harald Romsdorfer, Thibaut Weise, and Luc Van Gool. 2010. A 3-d audio-visual corpus of affective communication. *IEEE Transactions on Multimedia* 12, 6 (2010), 591–598.

[20] Yao Feng, Haiwen Feng, Michael J. Black, and Timo Bolkart. 2021. Learning an Animatable Detailed 3D Face Model from In-The-Wild Images. *ACM Transactions on Graphics, (Proc. SIGGRAPH)* 40, 8. https://doi.org/10.1145/3450626.3459936

[21] Awni Hannun, Carl Case, Jared Casper, Bryan Catanzaro, Greg Diamos, Erich Elsen, Ryan Prenger, Sanjeev Satheesh, Shubho Sengupta, Adam Coates, et al. 2014. Deep speech: Scaling up end-to-end speech recognition. *arXiv preprint arXiv:1412.5567* (2014).

[22] Kazi Injamamul Haque and Zerrin Yumak. 2023. FaceXHuBERT: Text-less Speech-driven E(X)pressive 3D Facial Animation Synthesis Using Self-Supervised Speech Representation Learning. In *INTERNATIONAL CONFERENCE ON MULTIMODAL INTERACTION (ICMI '23)* (Paris, France). ACM, New York, NY, USA, 10 pages. https://doi.org/10.1145/3577190.3614157

[23] Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising Diffusion Probabilistic Models. In *Proceedings of the 34th International Conference on Neural Information Processing Systems* (Vancouver, BC, Canada) *(NIPS'20)*. Curran Associates Inc., Red Hook, NY, USA, Article 574, 12 pages.

[24] Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems* 33 (2020), 6840–6851.

[25] Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. 2021. HuBERT: Self-Supervised Speech Representation Learning by Masked Prediction of Hidden Units. https://doi.org/10.48550/ARXIV.2106.07447

[26] JALI 2023. JALI Research. https://jaliresearch.com/.

[27] Xinya Ji, Hang Zhou, Kaisiyuan Wang, Qianyi Wu, Wayne Wu, Feng Xu, and Xun Cao. 2022. EAMM: One-Shot Emotional Talking Face via Audio-Based Emotion-Aware Motion Model. In *ACM SIGGRAPH 2022 Conference Proceedings* (Vancouver, BC, Canada) *(SIGGRAPH '22)*. Association for Computing Machinery, New York, NY, USA, Article 61, 10 pages. https://doi.org/10.1145/3528233.3530745

[28] Patrik Jonell, Taras Kucherenko, Gustav Eje Henter, and Jonas Beskow. 2020. Let's face it: Probabilistic multi-modal interlocutor-aware generation of facial gestures in dyadic settings. In *International Conference on Intelligent Virtual Agents (IVA '20)*. ACM.

[29] Tero Karras, Timo Aila, Samuli Laine, Antti Herva, and Jaakko Lehtinen. 2017. Audio-driven facial animation by joint end-to-end learning of pose and emotion. *ACM Transactions on Graphics (TOG)* 36, 4 (2017), 1–12.

[30] Tianye Li, Timo Bolkart, Michael J Black, Hao Li, and Javier Romero. 2017. Learning a model of facial shape and expression from 4D scans. *ACM Trans. Graph.* 36, 6 (2017), 194–1.

[31] Haiyang Liu, Zihao Zhu, Naoya Iwamoto, Yichen Peng, Zhengqing Li, You Zhou, Elif Bozkurt, and Bo Zheng. 2022. BEAT: A Large-Scale Semantic and Emotional Multi-Modal Dataset for Conversational Gestures Synthesis. In *European conference on computer vision*.

[32] Yuanxun Lu, Jinxiang Chai, and Xun Cao. 2021. Live Speech Portraits: Real-Time Photorealistic Talking-Head Animation. *ACM Transactions on Graphics* 40, 6 (12 2021), 17 pages. https://doi.org/10.1145/3478513.3480484

[33] Araceli Morales, Gemma Piella, and Federico M. Sukno. 2021. Survey on 3D face reconstruction from uncalibrated images. *Computer Science Review* 40 (2021), 100400. https://doi.org/10.1016/j.cosrev.2021.100400

[34] Evonne Ng, Hanbyul Joo, Liwen Hu, Hao Li, , Trevor Darrell, Angjoo Kanazawa, and Shiry Ginosar. 2022. Learning to Listen: Modeling Non-Deterministic Dyadic Facial Motion. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022).

[35] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. 2015. Librispeech: An ASR corpus based on public domain audio books. In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 5206–5210. https://doi.org/10.1109/ICASSP.2015.7178964

[36] Konpat Preechakul, Nattanat Chatthee, Suttisak Wizadwongsa, and Supasorn Suwajanakorn. 2022. Diffusion autoencoders: Toward a meaningful and decodable representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10619–10629.

[37] Prolific 2023. Prolific. https://www.prolific.co.

[38] Qualtrics 2023. Qualtrics. https://www.qualtrics.com.

[39] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog* 1, 8 (2019), 9.

[40] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. Hierarchical Text-Conditional Image Generation with CLIP Latents. arXiv:2204.06125 [cs.CV]

[41] Ray CGTARIAN 2023. Ray Character Maya Scene by CGTarian. https://www.cgtarian.com/maya-character-rigs/download-free-3d-character-ray.html.

[42] Alexander Richard, Michael Zollhöfer, Yandong Wen, Fernando De la Torre, and Yaser Sheikh. 2021. Meshtalk: 3d face animation from speech using cross-modality disentanglement. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 1173–1182.

[43] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-Resolution Image Synthesis With Latent Diffusion Models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 10684–10695.

[44] Soubhik Sanyal, Timo Bolkart, Haiwen Feng, and Michael Black. 2019. Learning to Regress 3D Face Shape and Expression from an Image without 3D Supervision. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 7763–7772.

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

[46] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. 2021. Score-Based Generative Modeling through Stochastic Differential Equations. In *International Conference on Learning Representations*. https://openreview.net/forum?id=PxTIG12RRHS

[47] Sarah Taylor, Taehwan Kim, Yisong Yue, Moshe Mahler, James Krahe, Anastasio Garcia Rodriguez, Jessica Hodgins, and Iain Matthews. 2017. A deep learning approach for generalized speech animation. *ACM Transactions on Graphics (TOG)* 36, 4 (2017), 1–11.

[48] Sarah L. Taylor, Moshe Mahler, Barry-John Theobald, and Iain Matthews. 2012. Dynamic Units of Visual Speech. In *Proceedings of the ACM SIG-GRAPH/Eurographics Symposium on Computer Animation* (Lausanne, Switzerland) *(SCA '12).* Eurographics Association, Goslar, DEU, 275–284.

[49] Guy Tevet, Sigal Raab, Brian Gordon, Yonatan Shafir, Daniel Cohen-Or, and Amit H Bermano. 2022. Human motion diffusion model. *arXiv preprint arXiv:2209.14916* (2022).

[50] Jonathan Tseng, Rodrigo Castellon, and Karen Liu. 2023. Edge: Editable dance generation from music. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 448–458.

[51] Aaron van den Oord, Oriol Vinyals, and koray kavukcuoglu. 2017. Neural Discrete Representation Learning. In *Advances in Neural Information Processing Systems*, I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2017/file/7a98af17e63a0ac09ce2e96d03992fbc-Paper.pdf

[52] Vicon 2023. Vicon. https://www.vicon.co.

[53] Cheng-hsin Wuu, Ningyuan Zheng, Scott Ardisson, Rohan Bali, Danielle Belko, Eric Brockmeyer, Lucas Evans, Timothy Godisart, Hyowon Ha, Alexander Hypes, et al. 2022. Multiface: A dataset for neural face rendering. *arXiv preprint arXiv:2207.11243* (2022).

[54] Jinbo Xing, Menghan Xia, Yuechen Zhang, Xiaodong Cun, Jue Wang, and Tien-Tsin Wong. 2023. CodeTalker: Speech-Driven 3D Facial Animation with Discrete Motion Prior. *arXiv preprint arXiv:2301.02379* (2023).

[55] Sicheng Yang, Zhiyong Wu, Minglei Li, Zhensong Zhang, Lei Hao, Weihong Bao, Ming Cheng, and Long Xiao. 2023. DiffuseStyleGesture: Stylized Audio-Driven Co-Speech Gesture Generation with Diffusion Models. *arXiv preprint arXiv:2305.04919* (2023).

[56] Hongwei Yi, Hualin Liang, Yifei Liu, Qiong Cao, Yandong Wen, Timo Bolkart, Dacheng Tao, and Michael J Black. 2023. Generating Holistic 3D Human Motion from Speech. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 469–480.

[57] Wenxuan Zhang, Xiaodong Cun, Xuan Wang, Yong Zhang, Xi Shen, Yu Guo, Ying Shan, and Fei Wang. 2023. SadTalker: Learning Realistic 3D Motion Coefficients for Stylized Audio-Driven Single Image Talking Face Animation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8652–8661.

[58] Yang Zhou, Zhan Xu, Chris Landreth, Evangelos Kalogerakis, Subhransu Maji, and Karan Singh. 2018. Visemenet: Audio-driven animator-centric speech animation. *ACM Transactions on Graphics (TOG)* 37, 4 (2018), 1–10.

# A SUPPLEMENTARY MATERIAL

## A.1 User studies

In total, 83 survey responses were collected spread among the 3 experiments as follows- 31 for the BIWI survey, 29 for the VOCASET survey, and 23 for the UUDaMM survey. An additional 3 responses were discarded due to the participants failing the attention checks. The surveys were distributed online to different groups that voluntarily took part in the study without being remunerated, managing to collect 38 responses out of the total. The additional 45 responses were collected through the Prolific platform, which facilitates response collection by allowing participants to take part in surveys and get remuneration. We fairly compensated the participants by awarding them the equivalent of $9\$/h$. Most of the participants are adults or young adults, with the following age distribution: 65.85% in the 18-25 age group, 25.61 in the 26-35 age group, and 8.54% in the 36-45 age group. Looking at the gender distribution, 36.59% of the participants identify as female, 54.88% as male, and 8.53% as non-binary/other.

In regards to the user familiarity with the subject, we computed the average reported familiarity of the 3 areas that were questioned and we obtained 2.71 average familiarity with virtual humans, 3.52 average familiarity with 3D animated movies, and 4.06 average familiarity with video games on a 5 point likert scale.

The first 2 surveys present users with 12 pair-wise comparisons. To avoid the impact of random selection by users, we randomly switch the side on which we show our model with our competitors. For each comparison, the user is asked 3 questions. For the first 2 questions we follow previous works [17, 22, 54] and ask the users about lip-sync quality and perceived realism of the animations. Additionally, we add an extra question asking about the animation appropriateness. Survey instruction and an example of the user interface can be seen in Fig. 8a and Fig. 8b respectively.

The questions the users were presented with are as follows:

(1) Comparing the lips of the two animations, which one is more in sync with the audio?
(2) Comparing the full faces of the two animations, which one looks more realistic?
(3) Comparing the full faces of the two animations, which one is more appropriate for the given audio?

## A.2 Datasets

Here, we describe more in detail the datasets that were used in our work. A summary of the datasets can be found in Tab. 7 and Tab. 8.

*A.2.1 BIWI[19].* The dataset is comprised of 14 x 40 x 2 sequences of paired audio and animation. For the creation of these sequences, 14 subjects were asked to read and emote 40 different sentences, each sentence being read 2 times, one time with a neutral expression and one time with a more emotional one. Each sequence is on average approximately 5 seconds long and is captured at $25fps$. The face meshes included are very high-definition comprising 23370 3D vertices, despite only the front of the head being depicted. Based on previous work, we use the same data splits as in [17, 54] and only use the emotional subset of sequences. During training, only 6 subjects (3 female and 3 male) are used, along with 32 spoken sentences per subject. This amounts to a total of 192 sequences and

**Welcome to the facial animation perception study!**

Thank you for taking the time to take this perception study survey on facial animations. During the study, you will be shown **12 pairs of 3D facial animation videos** (each 3-6 seconds in duration) and you will be asked:
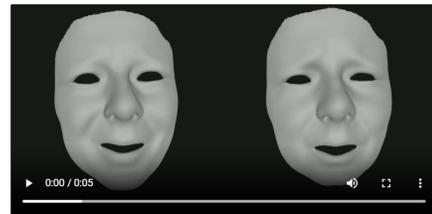
1. which one is more in sync with the audio
2. which one is more realistic based on the accompanying audio
3. which one is more appropriate for the given audio

You will choose either "**Left**" or "**Right**" for each question based on your preference for the two videos.
You may choose different answers for the 3 questions.

Please ensure that your **audio is turned on during the study** and if possible, **watch the videos in full-screen mode for better clarity (watch the videos multiple times if needed)**. The survey will take around **5 minutes** to complete.

**(a) User study instruction.**



| | Left | Right |
|---|---|---|
| Comparing the lips of the two animations, which one is more in sync with the audio? | ○ | ○ |
| Comparing the full faces of the two animations, which one looks more realistic? | ○ | ○ |
| Comparing the full faces of the two animations, which one is more appropriate for the given audio? | ○ | ○ |

**(b) User study UI.**

**Figure 8: Example screenshots of the user studies.**

represents the BIWI-Train dataset. From the remaining 8 sentences of these subjects, 4 are used for validations (24 in total), and 4 for testing (24 in total). We refer to this test set as BIWI-Test-A and will be used to compute objective metrics over our results. With the remaining, 8 subjects and their last 4 sentences, BIWI-Test-B is formed. This dataset represents the primary one we use during the model experiment phase. Furthermore, the hyperparameter tuning along with other experiments and the ablation study are mainly performed on this dataset.

*A.2.2 VOCASET[8].* . The dataset is comprised of 480 sequences of audio-visual pairs which amount to a total length of just 29 minutes. The sequences are recorded at $60fps$ and the facial scans are translated onto the FLAME head topology [30] which is comprised of 5023 3D vertices. Unlike, the BIWI dataset, the mesh includes the whole head and neck of the person, including eyeballs and eyelids. Even though this would technically allow for more expressivity and possibly even eye blinks to be captured, these are scarcely available

| Dataset | BIWI | VOCASET | Multiface |
|---|---|---|---|
| Training Set | • 6 subjects<br>• 32 sequences per subject<br>• Total = 192 sequences | • 8 subjects<br>• 40 sequences per subject<br>• Total = 320 sequences | • 9 subjects<br>• 40 sequences per subject<br>• Total = 360 sequences |
| Validation Set | • 6 seen subjects<br>• 4 sequences per subject<br>• Total = 24 sequences | • 2 unseen subjects<br>• 20 sequences per subject<br>• Total = 40 sequences | • 9 seen subjects<br>• 5 sequences per subject<br>• Total = 45 sequences |
| Test Set A | • 6 seen subjects<br>• 4 sequences per subject<br>• Total = 24 sequences | - | • 9 seen subjects<br>• 5 sequences per subject<br>• Total = 45 sequences |
| Test Set B | • 8 unseen subjects<br>• 4 sequences per subject<br>• 6 conditions per sequence<br>• Total = 192 sequences | • 2 unseen subjects<br>• 20 sequences per subject<br>• 8 conditions per sequence<br>• Total = 320 sequences | • 4 unseen subjects<br>• 5 sequences per subject<br>• 9 conditions per sequence<br>• Total = 180 sequences |

**Table 7: Summary of the dataset split for the 3 vertex based datasets used in our work.**

| Dataset | UUDaMM | BEAT |
|---|---|---|
| Training Set | • 2 subjects<br>• 2029 10-second sequences per subject<br>• Total: 4058 sequences | • 4 subjects<br>• 2175 10-second sequences |
| Validation Set | • 2 subjects<br>• 254 10-second sequences per subject<br>• Total: 508 sequences | • 4 subjects<br>• 274 10-second sequences |
| Test Set | • 2 subjects<br>• 254 10-second sequences per subject<br>• Total: 508 sequences | • 4 subjects<br>• 275 10-second sequences |

**Table 8: Summary of the dataset split for the 2 blendshape based datasets used in our work.**

in the dataset. In terms of data split, we utilise the one proposed by the authors of the dataset and later used by other methods as well [17, 54], 8 of the 12 subjects along with all of their sentences are used for training (320 sequences), 2 subjects with 20 sentences per subject are used for validation (40 sequences), and similarly the last 2 subjects with the last 20 sentences are used for testing (40 sequences). This dataset is used mostly for validating our approach in terms of generalizability and for conducting the user study.

*A.2.3   Multiface[53].* The third vertex based dataset in our research, namely the Multiface dataset, publicly released by *Wuu et al.*[53]. To the best of our knowledge, we are the first ones to use it for the task of facial animation synthesis besides its creators. In comparison to the other datasets, the face meshes included here are more complex in terms of features, containing attributes such as hair, eyelids and facial hair. Moreover, the dataset includes full 3D head scans, including the back of the head as well as the neck.

The publicly released version of the dataset contains a total of 13 subjects out of 250 subjects that was used for training Meshtalk [42]. Even though the full dataset is much larger, comprising a total of 250 subjects, it is not available to the public. For each subject, there are a total of 50 spoken sentences available. Since the sequences are split by subject and sentence, it allows us to have a similar training technique, including the one-hot style embedding.

The authors of the dataset share that the sentences were chosen in such a way that they are phonetically balanced ensuring a good generalisation across the possible phonemes. The animation is sequence of 3D face meshes available at 30 frames per second. Each frame is represented by the full 3D face of the actor, with a total of 6172 3D vertices, including eyelids, neck, as well as different hairstyles for the subjects. Since there is no previous work to follow for this subset of the dataset, we propose our own data split (similar BIWI in terms of number of sequences) and use 9 of the subjects with the first 40 sentences for training (360 sequences) and call this Multiface-Train, the following 5 sentences of the same 9 subjects are used for validation (45 sequences), and the last 5 sentences make up Multiface-Test-A and are used for testing (45 sequences). With the other 4 subjects and their last 5 sentences that were unseen during training, we form Multiface-Test-B.

*A.2.4   Utrecht University Dyadic Multimodal Motion Capture Dataset (UUDaMM).* Our in-house UUDaMM dataset consists of synchronously captured dyadic conversations between 2 actors in terms of four modalities- gesture, face, audio and text. The dataset contains 9 hours and 41 minutes of recorded conversations between two speakers, of which 6 hours and 53 minutes represent active speech sequences, in which at least one of the two actors is speaking. For our work in hand, we only discuss on the facial data in this document,

| Diffusion Steps | MVE ↓ x$10^{-3}$mm | LVE ↓ x$10^{-4}$mm | FDD ↓ x$10^{-5}$mm | Diversity ↑ |
|---|---|---|---|---|
| 100 | 7.2391 | 4.7703 | 5.0705 | 0.8236 |
| 250 | 6.9618 | 4.5515 | 4.2407 | 0.8331 |
| 500 | **6.8507** | **4.4364** | 4.3212 | 0.8446 |
| 750 | 7.1290 | 4.6428 | **4.1268** | **0.8725** |
| 1000 | 7.0387 | 4.6897 | 4.5889 | 0.8617 |

**Table 9: Evaluation metrics computed over the test results of different numbers for the diffusion timesteps**

but we consider that the way the dataset was collected, is also helpful for the task of facial animation via multi-modal learning approaches as well as for modelling dyadic interaction.

The facial performance capture, done with Dynamixyz[13] system at 120 FPS, is solved and retargeted to a virtual 3D character in a Maya scene. We use a publicly available character[41] to ensure ease of use among the research community. The character comes with artist generated blendshapes as well as intuitive facial controls that drive the blendshapes. After the facial performance is solved and retargeted, we use a python script to extract the temporal facial control values to form the training dataset. A similar python script can be used to import the inferred data into the Maya scene as well.

*A.2.5 BEAT[31].* The second blendshape based dataset we use is the BEAT dataset [31]. Just like UUDaMM, the dataset also contains body motion data along with facial capture, the difference being that the facial animations are encoded as sequences of blendshape weights instead of facial controls.

iPhone 12 Pro is used to capture the facial data of the actors, encoded as 52 blendshape weights defined in Apple's ARKit. The frame rate of the facial capture is 60 FPS. The full dataset includes 30 participants, of which half are female and half male. Each participant was asked to read 118 predefined texts, each resulting in a one-minute recording, in various emotional ways in order to cover multiple variations of emotional speech. The 8 emotions captured during the collection are neutral, happiness, anger, sadness, contempt, surprise, fear, and disgust. An additional 12 recordings of 10 minutes each were captured for each participant in which they perform free-form conversations with an off-screen director. Furthermore, the dataset contains sequences spoken in 4 different languages: English, Chinese, Japanese, and Spanish, also including native and non-native English speakers. The total size of the recordings amounts to about 76 hours. For our training, we use a 16 hour subset of the dataset comprising native English speaking sequences of 4 subjects. All these features make the dataset to be the most diverse out of the ones we consider, the only downfall being that "true" dyadic conversations are not recorded. In terms of pre-processing, we do not apply any transformations to the values of the features themselves, and merely split the sequences so that the data follows a similar format to that present in the other datasets. After splitting the provided sequences into 10-second segments, we obtain a total of 11398 that are used for training. In terms of data split, we employ an 80-10-10, training-validation-test split.
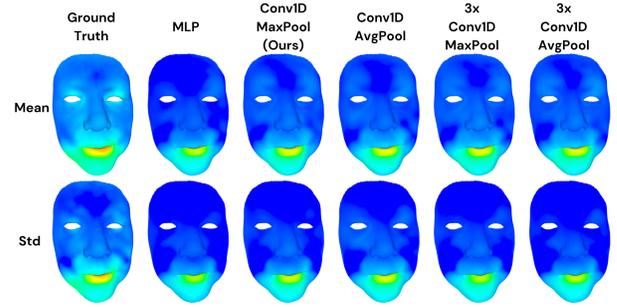


**Figure 9: Ablation on the noise encoder.**

| Noise Encoder | MVE ↓ x $10^{-3}$mm | LVE ↓ x $10^{-4}$mm | FDD ↓ x $10^{-5}$mm |
|---|---|---|---|
| MLP | 7.1728 | 4.9453 | 3.7748 |
| Conv1D + MaxPool (Ours) | **6.8088** | **4.2985** | 3.9100 |
| Conv1D + AvgPool | 6.8735 | 4.5766 | 4.2383 |
| 3 x (Conv1D + MaxPool) | 6.9217 | 4.5241 | **3.6020** |
| 3 x (Conv1D + AvgPool) | 6.8415 | 4.4130 | 3.9621 |

**Table 10: Results of different types of encoder for the noise**

## A.3 Additional Ablation on Diffusion Steps

We experimented with different diffusion step numbers as can be seen in Tab. 9. Analysing the visual results we observe that the number of diffusion steps does not influence the model's capacity of producing acceptable animations, all of the configurations producing correct lip-sync animations. The differences are mostly noticed when it comes to the general expressivity of the animations. The results obtained with just 100 diffusion steps are generally the worst both objectively and subjectively, while we do not see a lot of visual differences between the results obtained with the other tested values. Both 500 and 750 diffusion time steps yield acceptable results, however increasing the value beyond that worsens the results.

## A.4 Additional Ablation on Noise Encoder

For V-FaceDiffuser, due to the large number of 3D vertices in the data, we introduced a noise encoder in order to reduce the high dimensionality of the noise input to a low dimensional latent representation. We experimented with different configurations for the noise encoder as follows-

- **MLP.** A series of 2 fully connected layers.
- **Conv1D + MaxPool (Ours)** One fully connected layer followed by a single-dimensional convolution and a max pooling layer.
- **Conv1D + AvgPool** One fully connected layer followed by a single-dimensional convolution and an average pooling layer.
- **3 x Conv1D + MaxPool** Same as the second item but three times.

(a) An example of how the generated animation can be edited by an animator after automatically generated by our model.



(b) By using the BEAT dataset, FaceDiffuser is able to generate ARKit blendshape animations that can be used to animate a large variety of 3D characters that are ARKIT Blendshape enabled such as Ready Player Me avatars and Epic Games' MetaHumans.

Figure 10: Use cases of our approach in existing animation production workflows.

- **3 x Conv1D + AvgPool** Same as the third item but three times.

By analysing the results in both Tab. 10 and the heatmaps in Fig. 9, we can see a trend in the expressivity of the animations, viewed as a higher activation of the face, can be seen as the noise encoder becomes more complex. Some more complex configurations (i.e.

row 4 and 5 in Tab. 10) were also tried, showing better results in terms of the mean motion of the face. However, these come at the detriment of the error metrics compared to ground truth data. Out of the experimented configurations, our proposed choice performed the best. We can also notice that the Conv1D approaches are performing better than the MLP, both by looking at the results in Tab. 10 and the visual representation in Fig. 9. Furthermore, the method using max pooling performs slightly better than the average pooling one.

## A.5 Additional Use Cases

*A.5.1 Animation Editing.* By using a generalised animation encoding such as the Apple ARKit blendshape expression space, we are able to easily edit the resulting animations to better emulate our desired results. This can be done by simply updating the animation curves as can be seen in Fig. 10a and in the accompanying supplementary video. A simple change like moving one of the curves upwards would change the entire sequence expression and could be used to generate various expressions. For example, an animator might choose to generate more obvious mouth movements by applying a filter over the *mouthOpen* blendshape, which would make mouth opens and closures more extreme. Considering that accurate lip movements are automatically generated by a data-driven model like ours, the animators then have freedom of customising those animations to better fit their desires.

*A.5.2 Animation Transfer.* Another identifiable use case of generating ARKit blendhshape animations (with the model trained on BEAT dataset) is the transferability of such animations between different characters having the same blendshapes or semantically the same/similar blendshapes with different names. No additional retargeting steps are required if the blendshape names are the same as ARKit ones while a trivial script that maps the blendshape names between source data and target data would solve the retargeting and therefore the transfer is easy even for novice animators. This opens a wide-range of opportunities as there is a vast array of different faces that can be animated by using our model. However, due to the limited capability in terms of expressiveness in ARKit blendshapes, applying them to high fidelity photorealistic human characters may cause the effect known as the uncanny valley effect.