

# LipSync3D: Data-Efficient Learning of Personalized 3D Talking Faces from Video using Pose and Lighting Normalization

Avisek Lahiri<sup>1,2,\*†</sup> Vivek Kwatra<sup>1\*</sup> Christian Frueh<sup>1\*</sup> John Lewis<sup>1</sup> Chris Bregler<sup>1</sup>

<sup>1</sup>Google Research <sup>2</sup>Indian Institute of Technology Kharagpur

{avisek, kwatra, frueh, jplewis, bregler}@google.com

## Abstract

In this paper, we present a video-based learning framework for animating personalized 3D talking faces from audio. We introduce two training-time data normalizations that significantly improve data sample efficiency. First, we isolate and represent faces in a normalized space that decouples 3D geometry, head pose, and texture. This decomposes the prediction problem into regressions over the 3D face shape and the corresponding 2D texture atlas. Second, we leverage facial symmetry and approximate albedo constancy of skin to isolate and remove spatio-temporal lighting variations. Together, these normalizations allow simple networks to generate high fidelity lip-sync videos under novel ambient illumination while training with just a single speaker-specific video. Further, to stabilize temporal dynamics, we introduce an auto-regressive approach that conditions the model on its previous visual state. Human ratings and objective metrics demonstrate that our method outperforms contemporary state-of-the-art audio-driven video reenactment benchmarks in terms of realism, lip-sync and visual quality scores. We illustrate several applications enabled by our framework.

## 1. Introduction

“Talking head” videos, consisting of closeups of a talking person, are widely used in newscasting, video blogs, online courses, etc. Other applications that feature talking faces prominently are face-to-face live chat, 3D avatars and animated characters in games and movies. We present a deep learning approach to synthesize 3D talking faces (both photorealistic and animated) driven by an audio speech signal. We use speaker-specific videos to train our model in a data-efficient manner by employing 3D facial tracking. The resulting system has multiple applications, including video editing, lip-sync for dubbing of videos in a new language, personalized 3D talking avatars in gaming, VR and CGI, as well as compression in multimedia communication.

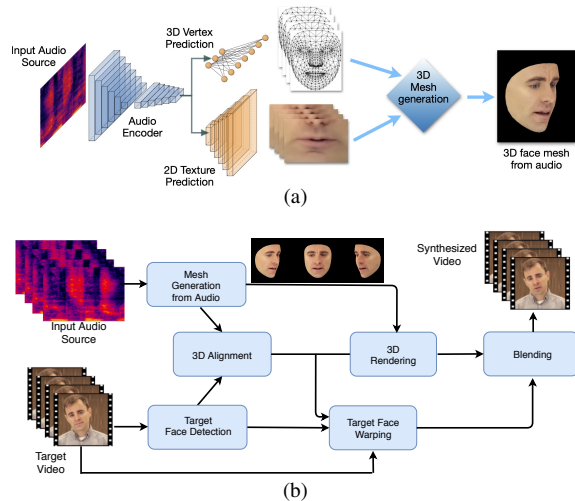


Figure 1: Flow diagram of our approach to (a) generate a dynamically textured 3D face mesh from audio, and (b) insert the generated face mesh into a target video to create a synthesized talking head video from new audio input.

The importance of talking head synthesis has led to a variety of methods in the research literature. Many recent techniques [6, 7, 40, 43, 28, 30] use the approach of regressing facial motion from audio, employing it to deform one or more reference images of the subject. These approaches can inherit the realism of the reference photos, however, the results do not accurately reproduce 3D facial articulation and appearance under general viewpoint and lighting variations. Another body of research predicts 3D facial meshes from audio [38, 13, 19, 11]. These approaches are directly suitable for VR and gaming applications. However, visual realism is often restricted by the quality of texturing. Some recent approaches [32, 33, 14] attempt to bridge the gap by combining 3D prediction with high-quality rendering, but are only able to edit fixed target videos that they train on.

Our work encompasses several of the scenarios mentioned above. We can use 3D information to edit 2D video, including novel videos of the same speaker not seen during training. We can also drive a 3D mesh from audio or text-to-speech (TTS), and synthesize animated characters by predicting face blendshapes.

\*Equal contribution

†Work done while an intern at Google

Next, we highlight some of our key design choices.

**Personalized models:** We train personalized speaker-specific models, instead of building a single universal model to be applied across different people. While universal models like Wav2Lip [30] are easier to reuse for novel speakers, they need large datasets for training and do not adequately capture person-specific idiosyncrasies [5]. Personalized models like ours and NVP [33] produce results with higher visual fidelity, more suitable for editing long speaker-specific videos. Additionally, our model can be trained entirely using a single video of the speaker.

**3D pose normalization:** We use a 3D face detector [20] to obtain the pose and 3D landmarks of the speaker’s face in the video. This information allows us to decompose the face into a normalized 3D mesh and texture atlas, thus decoupling head pose from speech-induced face deformations, *e.g.* lip motion and teeth/tongue appearance.

**Lighting normalization:** We design a *novel* algorithm for removing spatial and temporal lighting variations from the 3D decomposition of the face by exploiting traits such as facial symmetry and albedo constancy of the skin. This lighting normalization removes another confounding factor that can otherwise affect the speech-to-lips mapping.

**Data-efficient learning:** Our model employs an encoder-decoder architecture that computes embeddings from audio spectrograms, and decodes them to predict the decomposed 3D geometry and texture. Pose and lighting normalization allows us to train this model in a data-efficient manner. The model complexity is greatly reduced, since the network is not forced to disentangle unrelated head pose and lighting changes from speech, allowing it to synthesize high quality lip-sync results even from short training videos (2-5 minutes long). Lighting normalization allows training and inference illumination to be different, which obviates the need to train under multiple lighting scenarios. The model predicts *3D talking faces* instead of just a 2D image, even though it learns just from video, broadening its applicability. Finally, pose and lighting normalization can be applied in a backward fashion to align and match the appearance of the synthesized face with novel target videos. See Figure 1 for an overview of our approach.

Our key technical contributions are:

- A method to convert arbitrary talking head video footage into a normalized space that decouples 3D pose, geometry, texture, and lighting, thereby enabling data-efficient learning and versatile high-quality lip-sync synthesis for video and 3D applications.
- A novel algorithm for normalizing facial lighting in video that exploits 3D decomposition and face-specific traits such as symmetry and skin albedo constancy.
- To our best knowledge, this is the first attempt at disentangling pose and lighting from speech via data pre-normalization for personalized models.

- An easy-to-train auto-regressive texture prediction model for temporally smooth video synthesis.
- Human ratings and objective metrics suggest that our method outperforms contemporary audio-driven video reenactment baselines in terms of realism, lip-sync and visual quality scores.

## 2. Related Work

**Audio-driven 3D Mesh Animation:** These methods generate 3D face models driven by input audio or text, but do not necessarily aim for photorealism. In [38], the authors learn a Hidden Markov Model (HMM) to map Mel-frequency Cepstral Coefficients (MFCC) to PCA model parameters. Audio features are mapped to Jali [13] coefficients in [44]. In [19], the authors learn to regress to 3D vertices of a face model conditioned on input audio spectrograms and simultaneously disambiguate variations in facial expressions unexplained by audio. In [18], the authors learn to regress blendshapes of a 3D face using the combined audio-visual embedding from a deep network. VOCA [12] pre-registers subject-specific 3D mesh models using FLAME [26] and then learns (using hours of high quality 4D scans) an offset to that template based on incoming speech, represented with DeepSpeech [15] features.

**Audio-driven Video Synthesis:** These methods aim to generate visually plausible 2D talking head videos, conditioned on novel audio. In [6], an audio-visual correlation loss is used to match lip shapes to speech, while maintaining the identity of the target face. In [7], a two-stage cascaded network is used to first predict 2D facial landmarks from audio, followed by target frame editing conditioned upon these landmarks. In [36], the authors leverage a temporal GAN for synthesizing video conditioned on audio and a reference frame. They improve it further in [37] via a specialized lip-sync discriminator. In contrast to our approach, the above methods fail to produce full-frame outputs; instead they generate normalized cropped faces, whose lips are animated based on input audio and a reference frame.

Among efforts on full-frame synthesis, Video Rewrite [5] was a pioneering work. It represented speech with phonetic labels and used exemplar-based warping for mouth animation. Speech2Vid [8] learns a joint embedding space for representing audio features and the target frame, and uses a shared decoder to transform the embedding into a synthesized frame. X2Face [40] learns to drive a target frame with the head pose and expression of another source video, and it can optionally be also driven by an audio to animate a target frame. A framework to translate an input speech to another language and then modify the original video to match it is presented in [23]. Recently, Wav2Lip [30] reported appreciable lip-sync performance by using a powerful offline lip-sync

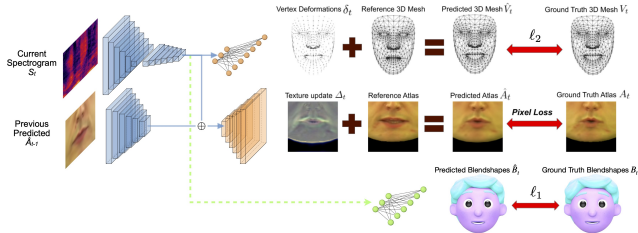


Figure 2: Joint prediction pipeline: geometry and texture models have dedicated decoders but share the audio encoder. The texture model also depends on the previously predicted atlas. Optionally, the audio embedding can drive a 3D CGI character via a blendshape coefficients decoder. Please enlarge to see details.

discriminator [9] as an expert to train their generator. While currently this is one of the best universal models, it lacks the visual fidelity of speaker-specific models.

Some recent works [32, 33, 31] have focused on 3D model guided video synthesis. In [32] an RNN regresses audio to mouth shape, producing convincing results on President Obama. The approach required very extensive training data however (17 hours). In [33], the DeepSpeech RNN is used to map input speech to audio expression units which then drive a blendshapes-based 3D face model. Finally, a neural renderer [22] is used to render the face model with the audio expressions. Since neural renderer training depends on target illumination, the methods leveraging such rendering [31, 33] suffer from the need for retraining if inference-time lighting conditions change. On the contrary, our method seamlessly adapts to novel lighting.

**Text-based Video Editing:** In [14], the authors present a framework for text based editing of videos (TBE). They first align written transcripts to audio and track each frame to create a face model. During edit operations, a (slow) viseme search is done to find best matching part of training video. This method needs a time-aligned transcript and around one hour of recorded data, and is mostly suitable for small edits. Our method, on the other hand, relies on just the audio signal and can synthesize videos of unrestricted length.

**Actor-driven Video Synthesis:** [34, 22] present techniques for generating and dubbing talking head videos by transferring facial features, such as landmarks or blendshape parameters, from a different actor’s video. These techniques generate impressive results, however they require a video of a surrogate actor to drive synthesis. We emphasize that our approach uses only audio or text-to-speech (TTS) as the driving input, and does not require any actors for dubbing. It is therefore fundamentally different from these methods.

### 3. Method

We now describe the various components of our approach including data extraction and normalization, neural network architecture and training, and finally, inference and synthesis. Figure 2 shows an overview of our model.

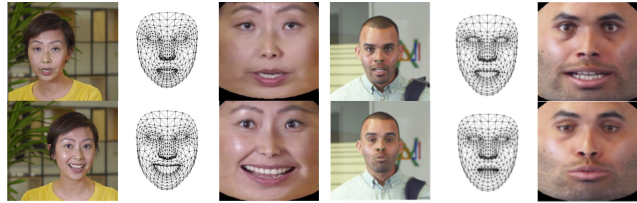


Figure 3: Pose normalization of training data. For each subject– Left: input frames with detected features (see zoomed in); Middle: normalized vertices and triangle mesh; Right: texture atlas which acts as ground truth for texture prediction.

We extract the audio channel from the training video and transform it into frequency-domain spectrograms. These spectrograms are computed using Short-time Fourier transforms (STFT) with a Hann window function [39], over 30ms wide sliding windows that are 10ms apart. We align these STFTs with video frames and stack them across time to create a  $256 \times 24$  complex spectrogram image, spanning 240ms centered around each video frame. Our model predicts the face geometry, texture, and optionally, blendshape coefficients, for each frame based on the audio spectrogram.

The face in the video is tracked using a 3D face landmark detector [20], resulting in 468 facial features, with the depth (z-component) predicted using a deep neural network. We refer to these features as vertices, which are accompanied by a predefined triangulated face mesh with fixed topology.

#### 3.1. Normalizing Training Data

We preprocess the training data to eliminate the effects of head movement and lighting variations, and work with normalized facial geometry and texture. Both training and inference take place in this normalized space.

##### 3.1.1 Pose normalization

For pose normalization, we first select one frame of the input video as a reference frame, and its respective 3D face feature points as reference vertices. The choice of frame is not critical; any frame where the face is sufficiently frontal is suitable. Using the reference vertices, we define a reference cylindrical coordinate system (similar to [4]) with a vertical axis such that most face vertices are equidistant to the axis. We then scale the face size such that the eyes and nose project to fixed locations on this reference cylinder.

Next, for each frame of the training video, we stabilize the rigid head motion (see [3, 24]) to provide a registered 3D mesh suitable for training our geometry model. Specifically, we approximately align the vertices of the upper, more rigid parts of the face with corresponding vertices in the normalized reference using Umeyama’s algorithm [35] and apply the estimated rotation  $\mathbf{R}$ , translation  $\mathbf{t}$  and scale  $c$  to all tracked vertices  $\mathbf{v}$  as  $\hat{\mathbf{r}} = c\mathbf{R}\mathbf{v} + \mathbf{t}$ .

We use these normalized vertices, along with the cylindrical mapping defined above, to create a pose-invariant, *frontalized* projection of the face texture for each video



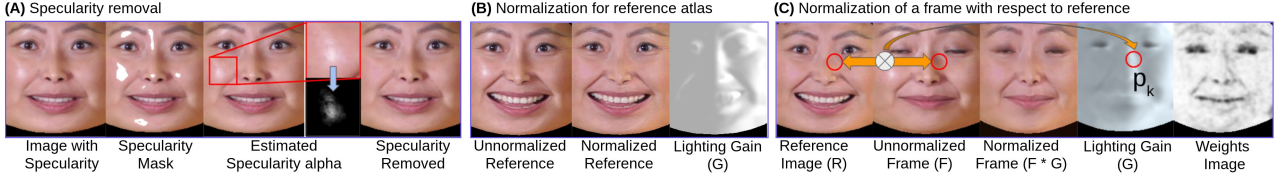


Figure 4: Steps of our proposed lighting normalization during training. (A:) First step is to specularity removal from an input frame. (B:) Second step is self normalization of the reference atlas. (C:) Finally, any given training frame is normalized with respect to the pre-normalized reference atlas of step B.

frame (including the reference frame). Mapping the face vertices to the reference cylinder creates a set of 2D texture coordinates for the face’s surface, which are used to *unroll* its texture. We warp the triangles associated with these coordinates from the source frame onto the texture domain, resulting in a  $256 \times 256$  *texture atlas* that resembles a frontal view of the face, but with the non-rigid features like the lips and mouth moving with the speech. Figure 3 demonstrates the effect of normalization; the head pose is removed, but the moving lip shapes and mouth interior are preserved.

### 3.1.2 Lighting normalization

We normalize the frontalized texture atlas to remove lighting variations, which are mostly caused by head motion or changing illumination. Our lighting normalization algorithm works in two phases. It first exploits facial symmetry to normalize the reference atlas  $R$  *spatially*, removing specularities and lighting variations that run across the face. It then performs a *temporal* normalization across video frames that transforms each frame’s atlas  $F$  to match the illumination of  $R$ . The resulting atlases have a more uniform *albedo-like* appearance, that stays consistent across frames.

We first describe the temporal normalization algorithm, as it is a core component also used during spatial normalization. This algorithm assumes that the two textures  $F$  and  $R$  are pre-aligned geometrically. However, any non-rigid facial movements, *e.g.* from speech, can result in different texture coordinates, and consequently, misalignments between  $R$  and  $F$ . Hence, we first warp  $R$  to align it with  $F$ ’s texture coordinates, employing the same triangle-based warping algorithm used for frontalization.

Given the aligned  $R$  and  $F$ , we estimate a mapping that transforms  $F$  to match the illumination of  $R$ . This mapping is composed of a smooth multiplicative pixel-wise gain  $G$  in the luminance domain, followed by a global channel-wise gain and bias mapping  $\{a, b\}$  in the RGB domain. The resulting normalized texture  $F^n$  is obtained via the following steps: (1)  $(F_y, F_u, F_v) = \text{RGBtoYUV}(F)$ ; (2)  $F_y^l = G * F_y$ ; (3)  $F^l = \text{YUVtoRGB}(F_y^l, F_u, F_v)$ ; (4)  $F^n = aF^l + b$ .

**Gain Estimation:** To estimate the gain  $G$ , we observe that a pair of corresponding pixels at the same location  $k$  in  $F$  and  $R$  should have the same underlying appearance, modulo any change in illumination, since they are in geometric alignment (see Figure 4(C)). This *albedo constancy* assumption, if perfectly satisfied, yields the gain at pixel  $k$

as  $G_k = R_k / F_k$ . However, we note that (a)  $G$  is a smoothly varying illumination map, and (b) albedo constancy may be occasionally violated, *e.g.* in non-skin pixels like the mouth, eyes and nostrils, or where the skin deforms sharply, *e.g.* the nasolabial folds. We account for these factors by, firstly, estimating  $G_k$  over a larger patch  $p_k$  centered around  $k$ , and secondly, employing a robust estimator that weights pixels based on how well they satisfy albedo constancy. We formulate estimating  $G_k$  as minimizing the error:

$$\mathbf{E}_k = \sum_{j \in p_k} W_j \|R_j - G_k * F_j\|^2, \quad (1)$$

where  $W$  is the per-pixel weights image, and solve it using iteratively reweighted least squares (IRLS). In particular, we initialize the weights uniformly, and then update them after each ( $i^{\text{th}}$ ) iteration as:

$$W_k^{i+1} = \exp\left(\frac{-\mathbf{E}_k^i}{T}\right), \quad (2)$$

where  $T$  is a temperature parameter. The weights and gain converge in 5-10 iterations; we use  $T = 0.1$  and a patch size of  $16 \times 16$  pixels for  $256 \times 256$  atlases. Figure 4(C) shows example weights and gain images. Pixels with large error  $\mathbf{E}_k$  get low weights, and implicitly interpolate their gain values from neighboring pixels with higher weights.

To estimate the global color transform  $\{a, b\}$  in closed form, we minimize  $\sum_k W_k \|R_k - aF_k - b\|^2$  over all pixels, with  $W_k$  now fixed to the weights estimated above.

**Reference Atlas Normalization using Facial Symmetry:** We first estimate the gain  $G^m$  between the reference  $R$  and its mirror image  $R'$ , using the algorithm described above. This gain represents the illumination change between the left and right half of the face. To obtain a reference with uniform illumination, we compute the symmetrized gain  $G^s = \max(G^m, G^{m'})$ , where  $G^{m'}$  is the mirror image of  $G^m$ , *i.e.* for every symmetric pair of pixels, we make the darker pixel match the brighter one. The normalized reference is then  $R^n = G^s * R$ , as shown in Figure 4(B). Note that our weighting scheme makes the method robust to inherent asymmetries on the face, since any inconsistent pixel pairs will be down-weighted during gain estimation, thereby preserving those asymmetries.

**Specularity Removal:** We remove specularities from the face before normalizing the reference and video frames, since they are not properly modeled as a multiplicative gain, and also lead to duplicate specularities on the reference due



to symmetrization. We model specular image formation as:

$$I = \alpha + (1 - \alpha) * I_c, \quad (3)$$

where  $I$  is the observed image,  $\alpha$  is the specular alpha map and  $I_c$  is the underlying *clean* image without specularities. We first compute a mask, where  $\alpha > 0$ , as pixels whose minimum value across RGB channels in a smoothed  $I$  exceeds the 90<sup>th</sup> percentile intensity across all skin pixels in  $I$ . The face mesh topology is used to identify and restrict computation to skin pixels. We then estimate a *pseudo* clean image  $\tilde{I}_c$  by hole-filling the masked pixels from neighboring pixels, and use it to estimate  $\alpha = (I - \tilde{I}_c)/(1 - \tilde{I}_c)$ .

The final clean image is then  $I_c = (I - \alpha)/(1 - \alpha)$ . Note that our soft alpha computation elegantly handles any erroneous over-estimation of the specular mask (see Figure 4(A)). The above method is specifically tailored for stabilized face textures and is simple and effective, thus we do not require more generalized specular removal techniques [42].

### 3.2. Joint Prediction Model and Training Pipeline

In this section we describe the framework for learning a function  $\mathbf{F}$  to jointly map from domain  $S$  of audio spectrograms to the domains  $V$  of *vertices* and  $A$  of *texture atlases*:  $\mathbf{F} : S \rightarrow V \times A$ , with  $V \in \mathbb{R}^{468 \times 3}$  and  $A \in \mathbb{R}^{128 \times 128 \times 3}$ , where for the purpose of prediction, we crop the texture atlas to a  $128 \times 128$  region around the lips, and only predict these cropped regions. The texture for the upper face is copied over from the reference, or target video frames, depending upon the application. We follow an encoder-decoder architecture for realizing  $\mathbf{F}(\cdot)$ , as shown in Figure 2. It consists of a shared encoder for audio, but separate dedicated decoders for geometry and texture. However, the entire model is trained jointly, end-to-end.

**Audio encoder:** The input at time instant  $t$  is a complex spectrogram,  $S_t \in \mathbb{R}^{256 \times 24 \times 2}$ . Our audio encoder — and face geometry prediction model — is inspired by the one proposed in [19], in which the vertex positions of a fixed-topology face mesh are also modified according to an audio input. However, while [19] used formant preprocessing and autocorrelation layers as input, we directly use complex spectrograms  $S_t$ . Each  $S_t$  tensor is passed through a 12 layer deep encoder network, where the first 6 layers apply 1D convolutions over frequencies (kernel  $3 \times 1$ , stride  $2 \times 1$ ), and the subsequent 6 layers apply 1D convolution over time (kernel  $1 \times 3$ , stride  $1 \times 2$ ), all with leaky ReLU activation, intuitively corresponding to phoneme detection and activation, respectively. This yields a latent code  $L_t^s \in \mathbb{R}^{N_s}$ .

**Geometry decoder:** This decoder maps the latent audio code  $L_t^s$  to vertex *deformations*  $\delta_t$ , which are added to the reference vertices  $V_r$  to obtain the predicted mesh  $\hat{V}_t = V_r + \delta_t$ . It consists of two fully connected layers with 150 and 1404 units, and linear activations, with a dropout layer in the middle. The resulting output is 468 vertices

(1404 = 468  $\times$  3 coordinates). As proposed in [19], we initialize the last layer using PCA over the vertex training data. Further, we impose  $\ell_2$  loss on the vertex positions:  $\mathbf{R}_{\text{geo}} = \|V_t - \hat{V}_t\|_2$ , where  $V_t$  are ground-truth vertices.

**Texture decoder:** This decoder maps the audio code  $L_t^s$  to a texture atlas *update* (difference map)  $\Delta_t$  which is added to the reference atlas  $A_r$  to obtain the predicted atlas,  $\hat{A}_t = A_r + \Delta_t$ . It consists of a fully connected layer to distribute the latent code spatially, followed by progressive up-sampling using convolutional and interpolation layers to generate the  $128 \times 128$  texture update image (see Appendix G.1). We impose an image similarity loss between the predicted and ground-truth atlas  $A_t$ :  $\mathbf{R}_{\text{tex}} = d(A_t, \hat{A}_t)$ , where  $d$  is a visual distance measure. We tried different variants of  $d(\cdot)$  including the  $\ell_1$  loss, Structural Similarity Loss (SSIM), and Gradient Difference Loss (GDL) [27] and found SSIM to perform the best.

**Blendshapes decoder:** To animate CGI characters using audio, we optionally add another decoder to our network that predicts *blendshape* coefficients  $B_t$  in addition to geometry and texture. For training, these blendshapes are derived from vertices  $V_t$  by fitting them to an existing blendshapes basis either via optimization or using a pre-trained model [25]. We use a single fully connected layer to predict coefficients  $\hat{B}_t$  from audio code  $L_t^s$ , and train it using  $\ell_1$  loss  $\mathbf{R}_{\text{bs}} = \|B_t - \hat{B}_t\|_1$  to encourage sparse coefficients.

#### 3.2.1 Auto-regressive (AR) Texture Synthesis:

Ambiguities in facial expressions while speaking (or silent) can result in temporal jitters. We mitigate these by incorporating memory into the network. Rather than using RNNs, we condition the current output of the network ( $A_t$ ) not only on  $S_t$  but also on the previous predicted atlas  $\hat{A}_{t-1}$ , encoding it as a latent code vector  $L_{t-1}^a \in \mathbb{R}^{N_a}$ .  $L_t^s$  and  $L_{t-1}^a$  are combined and passed to the texture decoder to generate the current texture  $\hat{A}_t$  (Figure 2). This appreciably improves the temporal consistency of synthesized results. We can train this AR network satisfactorily via *Teacher Forcing* [41], using previous ground truth atlases. The resulting network  $\mathbf{F}$  is trained end-to-end, minimizing the combined loss  $\mathbf{R} = \mathbf{R}_{\text{tex}} + \alpha_1 \mathbf{R}_{\text{geo}} + \alpha_2 \mathbf{R}_{\text{bs}}$ , where  $\alpha_1 = 3.0$  and  $\alpha_2 = 0.3$  (when enabled). We used hyperparameter search to determine the latent code lengths,  $N_s = 32$  and  $N_a = 2$ .

### 3.3. Inference and Synthesis

**Textured 3D mesh:** During inference, our model predicts geometry and texture from audio input. To convert it to a textured 3D mesh, we project the predicted vertices onto the reference cylinder, and use the resulting 2D locations as texture coordinates. Since our predicted texture atlas is defined on the same cylindrical domain, it is consistent with the computed texture coordinates. The result is a fully textured 3D face mesh, driven by audio input (Figure 1a).

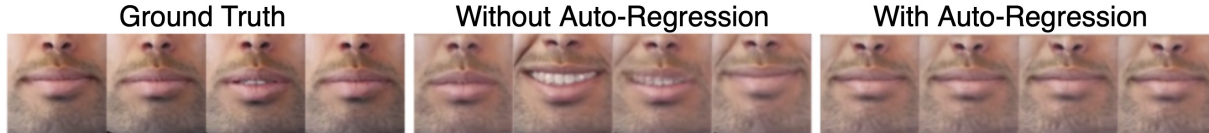


Figure 5: Benefits of proposed auto-regressive (AR) prediction. Left: Four consecutive frames when the subject was silent. Middle: Prediction without AR. Right: Prediction with AR. In absence of AR, the model fluctuates between different visual states, while the AR substantially improves temporal stability.

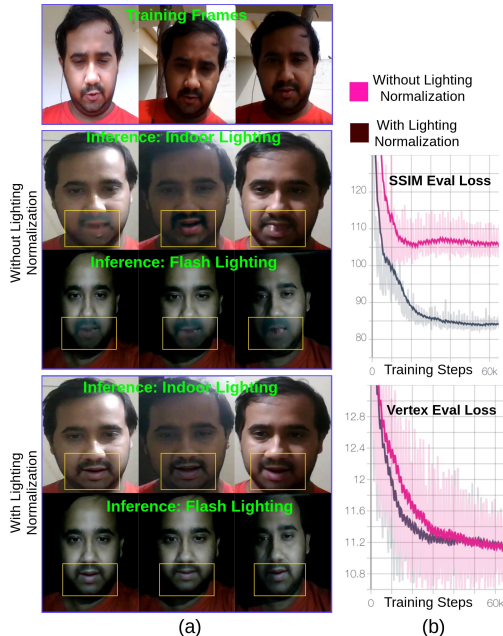


Figure 6: (a) Benefits of the proposed lighting normalization. Top row shows representative training frames in a sunny outdoor setting while we conduct inference under two novel lighting settings which have not been used in training. Note that the proposed lighting normalization enables realistic synthesis under new lighting while absence of lighting normalization yields degraded outputs. (b) Plot of SSIM loss (texture prediction) and vertex loss (geometry prediction) on the evaluation set. Even though both models result in similar lip shapes, the lower SSIM loss of the lighting-normalized model boosts the visual realism and overall lip-sync quality.

**Talking head video synthesis:** The pose and lighting normalization transforms (Section 3.1) are invertible, *i.e.* one can render the synthesized face mesh in a different pose under novel lighting, which allows us to procedurally blend it back into a different target video (Figure 1b). Specifically, we warp the textured face mesh to align it with the target face, then apply our lighting normalization algorithm in reverse, *i.e.* on the warped texture, using the target face as reference. One caveat is that the target frame’s area below the chin may not align with the warped synthesized face, due to inconsistent non-rigid deformations of the jaw. Hence, we pre-process each target frame by warping the area below the original chin to match the expected new chin position. To avoid seams at border areas, we gradually blend between the original and new face geometry, and warp the original face in the target frame according to the blended geometry.

**Cartoon rendering:** For stylized visualizations, we can create a cartoon rendering of the textured mesh (or video), by combining bilateral filtering with a line drawing of the facial features. In particular, we identify nose, lips, cheeks

and chin contours in the synthesized face mesh, and draw them prominently over the filtered texture or video frame.

**CGI Characters:** Models trained with the blendshapes decoder also output blendshape coefficients that can drive a CGI character. We combine these predicted blendshapes (that generally affect the lips and mouth) with other blendshapes, such as those controlling head motion and eye gaze, to create lively real-time animations. Please refer to Appendix- J and Fig. 14 for more details.

## 4. Experiments

Our training and inference pipelines were implemented in Tensorflow [1], Python and C++. We trained our models with batch sizes of 128 frames, for 500-1000 epochs, with each epoch spanning the entire training video. Sample training times were between 3-5 hours, depending on video length (usually 2-5min). Average inference times were 3.5ms for vertices, 31ms for texture and 2ms for blendshapes, as measured on a GeForce GTX 1080 GPU. Our research-quality code for blending into target videos takes 50-150ms per frame, depending on the output resolution.

### 4.1. Ablation Studies

**Benefit of Auto-Regressive Prediction:** The auto-regressive texture prediction algorithm stabilizes mouth dynamics considerably. In Figure 5, we show that without auto-regression, the model can produce an unrealistic jittering effect, especially during silent periods.

**Benefit of Lighting Normalization:** We use a short training video (~4 minutes) recorded in an outdoor setting but with varying illumination. However, during inference, we select two novel environments: a) indoor lighting with continuous change of lighting direction, and b) a dark room with a face illuminated by a moving flash light. Some representative frames of models trained with and without lighting normalization are shown in Figure 6(a). Without lighting normalization, the model produces disturbing artifacts around the lip region, exacerbated by the extreme changes in illumination. However, with normalized lighting, the model adapts to widely varying novel illumination conditions. This ability to edit novel videos of the same speaker *on-the-fly* without needing to retrain for new target illumination is a significant benefit. In contrast, neural rendering based approaches [33] require retraining on each new video, because they map 3D face models directly to the facial texture in video without disentangling illumination.



Figure 7: Qualitative comparison on subjects from GRID, CREMA-D and TCD-TIMIT against IJCV'19 and CVPR'19 (latter only available on GRID). Our model is capable of seamlessly blending back into the video instead of animating a normalized cropped frame as in IJCV'19 and CVPR'19..

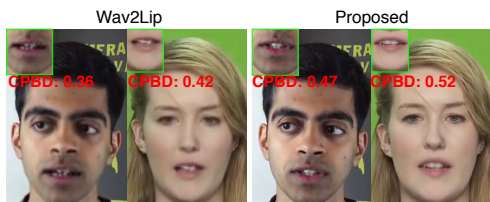


Figure 8: Comparison with Wav2Lip [30]. Our model generates higher resolution outputs (evident by higher CPBD metric [29]) with fewer artifacts compared to Wav2Lip. Examples are provided in accompanying video.

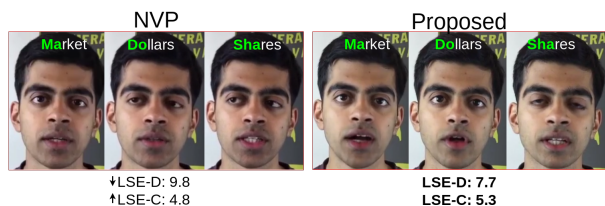


Figure 9: Comparison with NVP [33]. We show that a sequence generated by our method usually has better lip dynamics compared to NVP. The observation is also supported by LSE-D (lower is better) and LSE-C (higher is better) metrics [30] for our model. Examples are provided in accompanying video. Best viewed **zoomed** in.

We also visualize the loss curves on held out evaluation sets in Figure 6(b). With lighting normalization, the SSIM loss (used for texture generation) saturates at a much lower value than without normalization. This supports our hypothesis that lighting normalization results in more data-efficient learning, since it achieves a better loss with the same amount of training data. The vertex loss (responsible for lip dynamics) is similar for both models, because lighting normalization does not directly affect the geometry decoder, but overall lip-sync and visual quality are improved.

## 4.2. Comparison: Self-reenactment

We objectively evaluate our model under the self-reenactment setting (audio same as target video), since it allows us to have access to ground truth facial information. We show experiments with three talking head datasets: GRID [10], TCD-TIMIT [16] and CREMA-D [21].

**Comparing Methods:** We perform quantitative comparisons against state-of-the-art methods whose models/results are publicly available: CVPR'19 [7], IJCV'19 [37], CVPR-W [36]. It is difficult to do an apples-to-apples comparison, since we use personalized models while other techniques use a universal model. However, we minimize this gap by testing on the same 10 subjects from each of the 3 datasets used in IJCV'19 and CVPR'19, and employing the official

evaluation frameworks of these papers. We also compare against other prior methods, but reuse the results already reported by CVPR'19 or IJCV'19. Details of subject IDs are provided in Appendix- H

**Evaluation Metrics:** We follow the trend in recent papers [7, 37, 36], which use **SSIM** (Structural Similarity Index) as a reconstruction metric, **LMD** (Landmark Distance) on mouth features as a shape similarity metric, **CPBD** (Cumulative Probability Blur Detection) [29] as a sharpness metric and **WER** (word error rate) as a content metric to evaluate the correctness of words from reconstructed videos. Following [37], we use a LipNet model [2] pre-trained for lip-reading on GRID dataset [10].

**Observations:** We report the metrics in Figure 10 (left). On LMD and WER, which capture lip-sync, our model is significantly better than any competing method. Also, in terms of reconstruction measures (SSIM, CPBD), our model almost always performs better. CVPR-W and IJCV'19 have a better (though comparable) CPBD on GRID, but it is a low-resolution dataset. On higher resolution TCD-TIMIT and CREMA-D, our CPBD is the best. We also show qualitative comparisons in Figure 7. Note that we synthesize full frame videos, while CVPR'19 and IJCV'19 only generated normalized face crops at a resolution of  $128 \times 128$ , and  $96 \times 128$  respectively. Thus our method is more suitable for practical video applications.

## 4.3. Comparison: Audio-Driven Video Dubbing

In this section we focus on 'audio-driven' video dubbing where the driving audio is different from the target video.

**User Study:** We conducted a user study to quantitatively compare our lip-sync and perceptual quality against the state-of-the-art audio-driven frameworks of Wav2Lip, NVP, IJCV'19 and TBE. In the study, 35 raters were each shown 29 sample clips consisting of synthetic and real videos. For competing methods, we used their released videos or generated results with their pre-trained models. The raters were asked three questions: Q1) Is the video real or fake? Q2) Rate lip-sync quality on a 3-point discrete scale. Q3) Rate visual quality on a 5-point discrete scale. We report the Mean Opinion Scores (MOS) of the questions in Figure 10 (right). As is evident, among the competing methods our method receives the most favorable user ratings.



Methods	GRID				TCD-TIMIT			CREMA-D		
	SSIM(↑)	LMD(↓)	CPBD(↑)	WER(↓)	SSIM	LMD	CPBD	SSIM	LMD	CPBD
CVPR-W'19	0.84	-	<b>0.27</b>	25%	0.69	-	0.25	-	-	-
IJCV'19	0.81	1.32	0.26	23%	0.73	1.81	0.30	0.66	1.70	0.21
CVPR'19	0.81	1.37	0.17	70%	—	—	—	—	—	—
BMVC'17	0.72	—	0.25	58%	0.65	—	0.21	0.70	—	0.21
Chen <i>et al.</i> (ECCV'18)	0.73	1.73	—	—	—	—	—	—	—	—
Wiles <i>et al.</i> (ECCV'18)	0.75	1.60	—	—	—	—	—	—	—	—
Proposed	<b>0.94</b>	<b>0.80</b>	0.25	<b>18%</b>	<b>0.91</b>	<b>1.57</b>	<b>0.40</b>	<b>0.91</b>	<b>1.33</b>	<b>0.26</b>

Method	Is Real ? (% Yes)	Lip-Sync (1-3)	Visual Quality (1-5)
Real	97.6	2.95±0.03	4.55±0.13
IJCV'19	60.7	2.45±0.11	2.49±0.17
Wav2Lip	34.4	1.72±0.12	3.35±0.20
NVP	44.4	1.80±0.13	3.75±0.19
TBE	50.7	2.17±0.17	4.07±0.26
<b>Proposed</b>	<b>71.6</b>	<b>2.46±0.10</b>	<b>4.10±0.15</b>

Figure 10: **Left:** Self-reenactment performance comparison against state-of-the-art benchmarks of CVPR-W'19 [36], IJCV'19 [37], CVPR'19 [7], BMVC'17 [8], Chen *et al.* [6] and Wiles *et al.* [40]. Pre-trained LipNet (for WER) is available only on GRID. Authors of [7] released checkpoint for GRID only. (↑):Higher is better. (↓):Lower is better. Best results are marked in bold. **Right:** Mean Opinion Scores of user study. The statistical significance of these differences in ratings is confirmed by ANOVA with Tukey post-hoc tests. Please see Appendix- A for details.

**Comparison with Wav2Lip [30]:** Unlike other image-based methods, Wav2Lip can paste back the generated face on background video. However, compared to our model, the outputs from Wav2Lip are of low resolution. Also, at high resolution, Wav2Lip produces significant visual artifacts (see Figure 8) and lip-sync starts to degrade.

**Comparison with NVP [33]:** The lip-sync and dynamics of our model are generally better than NVP. The lip movements of NVP are clearly muted compared to our model, as seen in representative frames in Figure 9(a).

## 5. Applications

**Speech/Text-to-Video:** We can create or edit talking head videos for education, advertisement, and entertainment by simply providing new audio transcripts. **“Actor-free” video translation:** while ‘actor-driven’ video translation techniques [22, 34] generally require a professional actor to record the entire translated audio and video, our ‘actor-free’ approach does not need video, and can be driven by either recorded audio, TTS, or voice cloning [17]. **Voice controlled Avatars:** Our model’s blendshapes output can be used to animate CGI characters in real-time, allowing low-bandwidth voice-driven avatars for chat, VR, and games without the need for auxiliary cameras. **Assistive technologies:** Voice-driven 3D faces can support accessibility and educational applications, *e.g.* personified assistants and cartoon animations for visualizing pronunciation.

## 6. Limitations and Conclusion

**Facial expressions:** We do not explicitly handle facial expressions, though our model may implicitly capture correlations between expressions and emotion in the audio track. **Strong movements in the target video:** When synthesized faces are blended back into a target video, emphatic hand or head movement might seem out of place. This has not proved to be a problem in our experiments. **Processing speed:** Our research-quality code, running at highest quality, is slightly slower than real-time. We have presented a data efficient yet robust end-to-end system for synthesizing personalized 3D talking faces, with applications in video creation and editing, 3D gaming and CGI. Our proposed pose and lighting normalization decouples non-essential factors such as head pose and illumination

from speech and enables training our model on a relatively short video of a single person while nevertheless generating high quality lip-sync videos under novel ambient lighting. We envision that our framework is a promising stepping stone towards personalized audio-visual avatars and AI-assisted video content creation.

## 7. Ethical Considerations

Our technology focuses on world-positive use cases and applications. Video translation and dubbing have a variety of beneficial and impactful uses, including making educational lectures, video-blogs, public discourse, and entertainment media accessible to people speaking different languages, and creating personable virtual “assistants” that interact with humans more naturally.

However, we acknowledge the potential for misuse, especially since audiovisual media are often treated as veracious information. We strongly believe that the development of such generative models by *good actors* is crucial for enabling preemptive research on fake content detection and forensics, which would allow them to make early advances and stay ahead of actual malicious attacks. Approaches like ours can also be used to generate counterfactuals for training provenance and digital watermarking techniques. We also emphasize the importance of acting responsibly and taking ownership of synthesized content. To that end, we strive to take special care when sharing videos or other material that have been synthesized or modified using these techniques, by clearly indicating the nature and intent of the edits. Finally, we also believe it is imperative to obtain consent from all performers whose videos are being modified, and be thoughtful and ethical about the content being generated. We follow these guiding principles in our work.

## 8. Acknowledgments

We would like to thank all the performers who graciously allowed us to use their videos for this work, including D. Sculley, Kate Lane, Ed Moreno, Glenn Davis, Martin Aguinis, Caile Collins, Laurence Moroney, Paige Bailey, Ayush Tewari and Yoshua Bengio. We also thank the performers and creators of external, public datasets, as well as all participants of our user study. We would also like to thank our collaborators at Google and Deep

Mind: Paul McCartney, Brian Colonna, Michael Nechyba, Avneesh Sud, Zachary Gleicher, Miaosen Wang, Yi Yang, Yannis Assael, Brendan Shillingford, and Yu Zhang.

## References

- [1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org. [6](#)
- [2] Yannis M Assael, Brendan Shillingford, Shimon Whiteson, and Nando De Freitas. Lipnet: End-to-end sentence-level lipreading. *arXiv preprint arXiv:1611.01599*, 2016. [7](#)
- [3] Thabo Beeler and Derek Bradley. Rigid stabilization of facial expressions. *ACM Trans. Graph.*, 33(4), July 2014. [3](#)
- [4] J. Booth and S. Zafeiriou. Optimal UV spaces for facial morphable model construction. In *IEEE ICIP*, pages 4672–4676, 2014. [3](#)
- [5] Christoph Bregler, Michele Covell, and Malcolm Slaney. Video Rewrite: driving visual speech with audio. In *SIGGRAPH*, volume 97, pages 353–360, 1997. [2](#)
- [6] Lele Chen, Zhiheng Li, Ross K Maddox, Zhiyao Duan, and Chenliang Xu. Lip movements generation at a glance. In *ECCV*, pages 520–535, 2018. [1](#), [2](#), [8](#)
- [7] Lele Chen, Ross K Maddox, Zhiyao Duan, and Chenliang Xu. Hierarchical cross-modal talking face generation with dynamic pixel-wise loss. In *CVPR*, pages 7832–7841, 2019. [1](#), [2](#), [7](#), [8](#)
- [8] Joon Son Chung, Amir Jamaludin, and Andrew Zisserman. You said that? In *BMVC*, 2017. [2](#), [8](#)
- [9] J. S. Chung and A. Zisserman. Out of time: automated lip sync in the wild. In *Workshop on Multi-view Lipreading, ACCV*, 2016. [3](#), [13](#)
- [10] Martin Cooke, Jon Barker, Stuart Cunningham, and Xu Shao. An audio-visual corpus for speech perception and automatic speech recognition. *The Journal of the Acoustical Society of America*, 120(5):2421–2424, November 2006. [7](#), [14](#)
- [11] Daniel Cudeiro, Timo Bolkart, Cassidy Laidlaw, Anurag Ranjan, and Michael J Black. Capture, learning, and synthesis of 3D speaking styles. In *CVPR*, 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 *CVPR*, pages 10101–10111, 2019. [2](#)
- [13] Pif Edwards, Chris Landreth, Eugene Fiume, and Karan Singh. JALI: an animator-centric viseme model for expressive lip synchronization. *ACM TOG*, 35(4):127, 2016. [1](#), [2](#)
- [14] Ohad Fried, Ayush Tewari, Michael Zollhöfer, Adam Finkelstein, Eli Shechtman, Dan B Goldman, Kyle Genova, Zeyu Jin, Christian Theobalt, and Maneesh Agrawala. Text-based editing of talking-head video. *ACM TOG*, 38(4):1–14, 2019. [1](#), [3](#), [11](#)
- [15] Awni Hannun, Carl Case, Jared Casper, Bryan Catanzaro, Greg Diamos, Erich Elsen, Ryan Prenger, Sanjeev Satheesh, Shubho Sengupta, Adam Coates, et al. Deep speech: Scaling up end-to-end speech recognition. *arXiv preprint arXiv:1412.5567*, 2014. [2](#)
- [16] Naomi Harte and Eoin Gillen. TCD-TIMIT: An audio-visual corpus of continuous speech. *IEEE Transactions on Multimedia*, 17(5):603–615, 2015. [7](#), [14](#)
- [17] Wei-Ning Hsu, Yu Zhang, Ron Weiss, Heiga Zen, Yonghui Wu, Yuxuan Wang, Yuan Cao, Ye Jia, Zhifeng Chen, Jonathan Shen, Patrick Nguyen, and Ruoming Pang. Hierarchical generative modeling for controllable speech synthesis. In *International Conference on Learning Representations*, 2019. [8](#), [16](#)
- [18] Ahmed Hussen Abdelaziz, Barry-John Theobald, Paul Dixon, Reinhard Knothe, Nicholas Apostoloff, and Sachin Kajareker. Modality dropout for improved performance-driven talking faces. In *International Conference on Multimodal Interaction*, pages 378–386, 2020. [2](#)
- [19] 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 TOG*, 36(4):94:1–94:12, July 2017. [1](#), [2](#), [5](#)
- [20] Yury Kartynnik, Artsiom Ablavatski, Ivan Grishchenko, and Matthias Grundmann. Real-time facial surface geometry from monocular video on mobile GPUs. In *Third Workshop on Computer Vision for AR/VR, Long Beach, CA*, 2019. [2](#), [3](#)
- [21] Michael K Keutmann, Samantha L Moore, Adam Savitt, and Ruben C Gur. Generating an item pool

- for translational social cognition research: methodology and initial validation. *Behavior research methods*, 47(1):228–234, 2015. 7, 14
- [22] H. Kim, P. Garrido, A. Tewari, W. Xu, J. Thies, N. Nießner, P. Pérez, C. Richardt, M. Zollhöfer, and C. Theobalt. Deep Video Portraits. *ACM TOG*, 2018. 3, 8, 16
- [23] Prajwal KR, Rudrabha Mukhopadhyay, Jerin Philip, Abhishek Jha, Vinay Namboodiri, and CV Jawahar. Towards automatic face-to-face translation. In *ACM Multimedia*, pages 1428–1436, 2019. 2, 13
- [24] Mathieu Lamarre, J.P. Lewis, and Etienne Danvoye. Face stabilization by mode pursuit for avatar construction. In *Image and Vision Computing*, pages 1–6. IEEE, 2018. 3
- [25] J.P. Lewis, Ken Anjyo, Taehyun Rhee, Mengjie Zhang, Frédéric H. Pighin, and Zhigang Deng. Practice and theory of blendshape facial models. In Sylvain Lefebvre and Michela Spagnuolo, editors, *Eurographics - State of the Art Reports*. Eurographics Association, 2014. 5
- [26] Tianye Li, Timo Bolkart, Michael J Black, Hao Li, and Javier Romero. Learning a model of facial shape and expression from 4D scans. *ACM TOG*, 36(6):194, 2017. 2
- [27] Michael Mathieu, Camille Couprie, and Yann Lecun. Deep multi-scale video prediction beyond mean square error. *ICLR*, 2016. 5
- [28] Gaurav Mittal and Baoyuan Wang. Animating face using disentangled audio representations. In *WACV*, 2019. 1
- [29] Niranjan D Narvekar and Lina J Karam. A no-reference perceptual image sharpness metric based on a cumulative probability of blur detection. In *International Workshop on Quality of Multimedia Experience*, pages 87–91. IEEE, 2009. 7
- [30] KR Prajwal, Rudrabha Mukhopadhyay, Vinay P Namboodiri, and CV Jawahar. A lip sync expert is all you need for speech to lip generation in the wild. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 484–492, 2020. 1, 2, 7, 8, 11
- [31] Linsen Song, Wayne Wu, Chen Qian, Chen Qian, and Chen Change Loy. Everybody’s talkin’: Let me talk as you want. *arXiv preprint*, arXiv:, 2020. 3
- [32] Supasorn Suwajanakorn, Steven M Seitz, and Ira Kemelmacher-Shlizerman. Synthesizing Obama: learning lip sync from audio. *ACM (TOG)*, 36(4):95, 2017. 1, 3, 13
- [33] Justus Thies, Mohamed Elgharib, Ayush Tewari, Christian Theobalt, and Matthias Nießner. Neural voice puppetry: Audio-driven facial reenactment. In *ECCV*, pages 716–731. Springer, 2020. 1, 2, 3, 6, 7, 8, 11
- [34] Justus Thies, Michael Zollhofer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. Face2face: Real-time face capture and reenactment of rgb videos. In *CVPR*, pages 2387–2395, 2016. 3, 8, 16
- [35] Shinji Umeyama. Least-squares estimation of transformation parameters between two point patterns. *IEEE TPAMI*, 13(4):376–380, 1991. 3
- [36] Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic. End-to-end speech-driven realistic facial animation with temporal gans. In *CVPR Workshops*, pages 37–40, 2019. 2, 7, 8
- [37] Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic. Realistic speech-driven facial animation with GANs. *IJCV*, pages 1–16, 2019. 2, 7, 8, 11, 14
- [38] Lijuan Wang, Wei Han, Frank K Soong, and Qiang Huo. Text driven 3D photo-realistic talking head. In *Twelfth Annual Conference of the International Speech Communication Association*, 2011. 1, 2
- [39] Wikipedia. Short-time Fourier transform. [https://en.wikipedia.org/wiki/Short-time\\_Fourier\\_transform](https://en.wikipedia.org/wiki/Short-time_Fourier_transform), 2019. 3
- [40] Olivia Wiles, A Sophia Koepke, and Andrew Zisserman. X2face: A network for controlling face generation using images, audio, and pose codes. In *ECCV*, pages 670–686, 2018. 1, 2, 8
- [41] R. J. Williams and D. Zipser. A learning algorithm for continually running fully recurrent neural networks. *Neural Computation*, 1(2):270–280, 1989. 5, 14
- [42] Qingxiong Yang, Jinhui Tang, and Narendra Ahuja. Efficient and robust specular highlight removal. *IEEE TPAMI*, 37(6):1304–1311, 2014. 5
- [43] Hang Zhou, Yu Liu, Ziwei Liu, Ping Luo, and Xiang Wang. Talking face generation by adversarially disentangled audio-visual representation. In *AAAI*, volume 33, pages 9299–9306, 2019. 1
- [44] Yang Zhou, Zhan Xu, Chris Landreth, Evangelos Kalogerakis, Subhransu Maji, and Karan Singh. Visemenet: Audio-driven animator-centric speech animation. *ACM TOG*, 37(4):161, 2018. 2



## Appendix

### A. User Study Analysis

We conducted a user study to quantitatively compare our lip-sync and perceptual quality against the state-of-the-art audio-driven frameworks of Wav2Lip [30], NVP [33], IJCV'19 [37] and TBE [14]. In the study, N=35 raters were each shown a total 29 sample clips consisting of synthetic and real videos. For competing methods, we used their released videos (NVP, TBE) or generated results with their pre-trained models (IJCV'19, Wav2Lip). The raters were primarily drawn from a pool of subjects without research expertise, supplemented with a minority (N=14) who were researchers. The subgroup of researchers included some having familiarity with computer vision topics but none were expert on speech-driven animation. The raters were asked three questions: **Q1: Is the video real or fake?** **Q2: Rate the quality of the lip-sync, i.e. how well does the motion of the lips match the audio, on a 3-point (discrete) scale (poor, acceptable, great).** **Q3: Rate the picture quality, e.g. naturalness, resolution, and consistency of the video, on a discrete 5-point scale from 1-5 (poor-great).**

Figure 12 shows, for each question, the percentage of raters who selected each rating. We report the Mean Opinion Scores (MOS) of the questions in Table 1. As is evident, among the competing methods, our method receives the most favorable user ratings.

We performed a statistical analysis to confirm the significance of these ratings. For Q1 (only) we excluded the text-to-speech results from consideration, since it was straightforward to judge the videos as “fake” due to the computer-generated speech. However, questions Q2 and Q3 are still relevant in the text-to-speech case, since it is possible to rate the quality of lip-sync and overall image naturalness even when the voice is clearly synthetic.

The statistical analysis confirms that the differences in real-fake ratings on Q1 are significant (Kruskal-Wallis test,  $\chi^2 = 158, p < 1e-04$ ), and our method outperforms the other methods after adjusting for multiple comparisons (Tukey’s Honest Significant Differences (HSD) IJCV p adj.= .003, Wav2Lip p adj.<1e-04, NVP p adj.<1e-04). For Q2 the differences in ratings are significant (Kruskal-Wallis  $\chi^2 = 279, p < 1e-04$ ), and our method outperforms most competing methods with statistical significance after adjusting for the multiple tests (Tukey’s HSD TBE p adj.= .035, Wav2Lip p adj.<1e-04, NVP p adj.<1e-04), however the difference versus IJCV'19 is not significant. For Q3 (Kruskal-Wallis  $\chi^2 = 248, p < 1e-04$ ) our method outperforms most competing methods with statistical significance after adjustment for multiple comparison (HSD IJCV p adj.<1e-04, Wav2lip p adj.<1e-04, NVP p adj.= 0.04 however the comparison with TBE is not significant.

These significance results for Q2 and Q3 (and in particu-

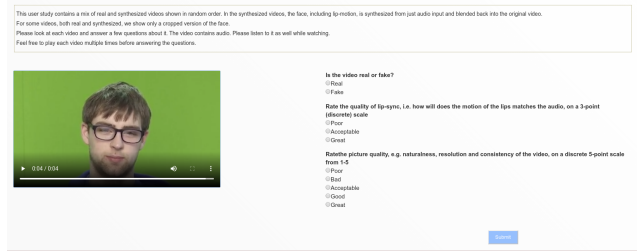


Figure 11: Screen shot from our user study.

Method	Is Real ? (% Yes)	[no TTS] Is Real ? (% Yes)	Lip-Sync (1-3)	Visual Quality (1-5)
Real	97.6	97.6	2.95±0.03	4.55±0.13
IJCV'19	60.7	60.7	2.45±0.11	2.49±0.17
Wav2Lip	34.4	37.6	1.72±0.12	3.35±0.20
NVP	44.4	50.0	1.80±0.13	3.75±0.19
TBE	50.7	n/a	2.17±0.17	4.07±0.26
<b>Proposed</b>	<b>71.6</b>	<b>77.25</b>	<b>2.46±0.10</b>	<b>4.10±0.15</b>

Table 1: User study analysis. Column 1: percentage of “real” ratings by category. Column 2: percentage of “real” ratings with text-to-speech driven results removed. Column 3: mean opinion score of lip-sync quality. Column 4: mean opinion score of picture quality.

lar the lack of significance for IJCV and TBE respectively) are plausible given cursory examination of the videos. The results of IJCV'19 show good quality lip-sync but the overall image quality is limited, thus explaining its good performance on Q2 but poor performance on Q3. TBE operates in part by re-mixing input video frames so it results in high picture quality by definition (Q3), but its lip-sync quality is poorer than our method and that of IJCV'19.

### B. Comparison Notes on Text Based Editing [14]

According to the user study, among the 3D model based methods, Text-based-Editing (TBE) is the second-best method (following our method). However, our framework has some distinct training and inference time advantages over TBE:

- TBE was trained on a training corpus of more than 1 hour of video recording. Our model was trained on ~7 minutes of data in this case.
- TBE assumes an accurate text transcript and uses phoneme based alignment tools to align the text with audio. In contrast, our model only requires a speech signal as input.
- The average training time of TBE is 42 hours. Our typical training time is somewhere in between 3-5 hours.

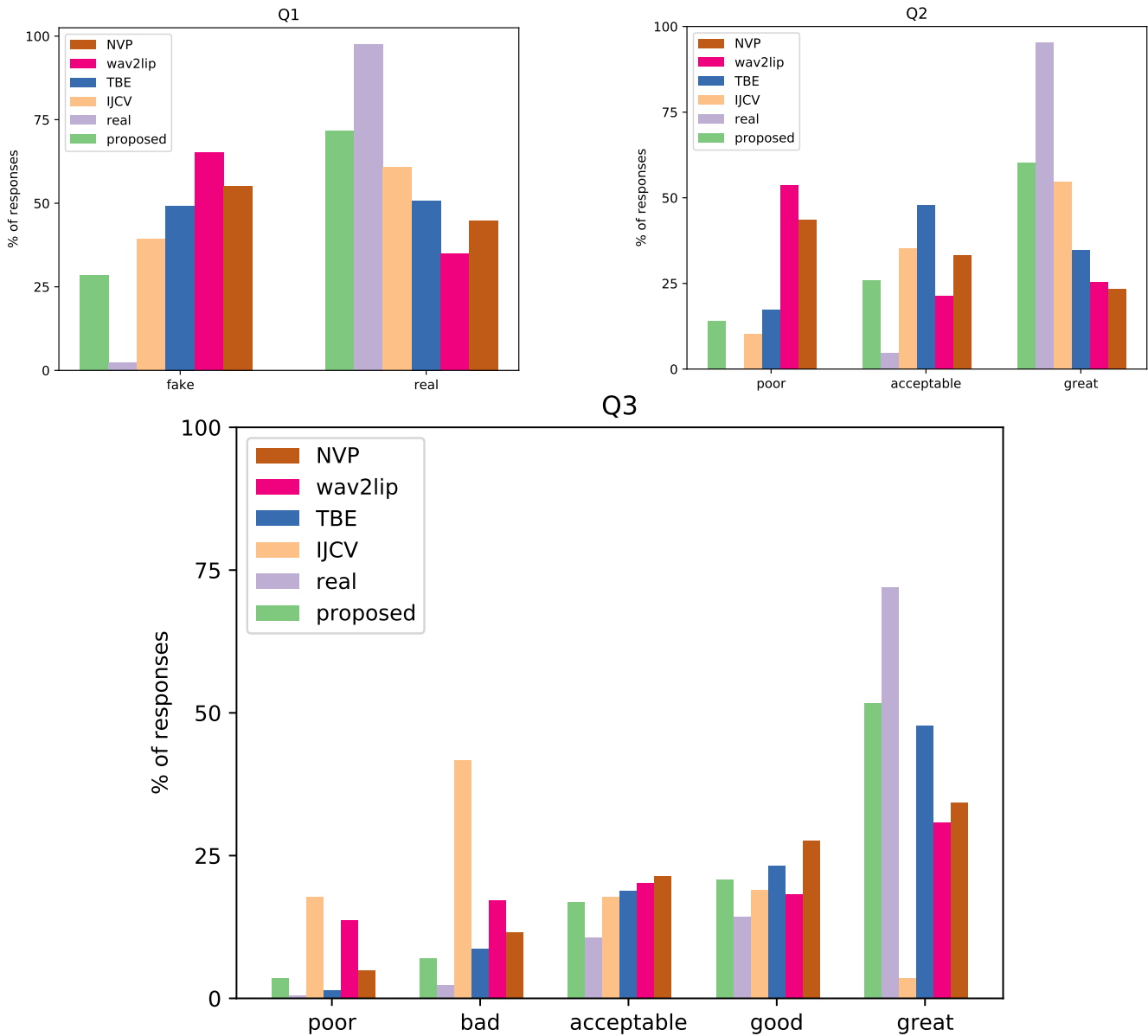


Figure 12: Raw user study results. **Q1**: Percentage of real/fake ratings for each of the six video categories. Viewers believe our synthetic videos (proposed) are real roughly two-thirds of the time. **Q2**: Ratings of lip-sync quality for the six video categories, expressed as percentages. Our synthetic videos (proposed) are perceived as at least comparable to those of IJCV while being clearly superior to other competing methods in lip-sync quality. **Q3**: Ratings of picture quality for the six video categories, expressed as percentages. Here our method greatly outperforms IJCV (and NVP and wav2lip) while being at least comparable to TBE. It can be seen that our method is the best performer across the three questions.

- The inference speed of TBE is significantly slower. This is mainly attributed to the costly viseme search (~5 minutes for 3 words). Our method executes within a few tens of milliseconds.
- TBE also relies on neural rendering for learning to generate facial texture based on the illumination in the training sequence. It is thus not apt for operating under

new ambient lighting without further retraining. Our framework is capable of seamlessly adapting to novel lighting conditions during inference.

### C. Limited comparison with Suwajanakorn *et al.* [32]

The work by Suwajanakorn *et al.* also involves training a personalized talking face model. However, [32] only synthesized results for a single person (former U.S. President Barack Obama), using hours of training video. While our model is perfectly capable of similar synthesis, we consciously refrain from training on living political personalities due to ethical considerations. Hence, we are unable to directly compare our work with that of Suwajanakorn *et al.*. Nevertheless, our framework offers the following advantages:

- Suwajanakorn *et al.* trained on 14 hours of weekly President addresses recorded between 2009-2016. In contrast, our framework just requires  $\sim 5$  minutes of training video.
- Suwajanakorn *et al.* demonstrate their outputs only under the specific studio lighting setup of the President’s office. Their texture generation network is not designed with the goal of handling diverse ambient lighting. In contrast, our network disentangles and normalizes the effects of illumination, thereby enabling inference under diverse lighting conditions.
- Typical training data pre-processing time of Suwajanakorn *et al.* is around 2 weeks on 10 cluster nodes of Intel Xeon E5530. In contrast, our combined pre-processing and training takes only about 3-5 hours on a single system equipped with a NVIDIA P1000 GPU.

### D. Discussion on LSE metrics

We have used the official code release<sup>1</sup> by the authors of Wav2Lip for evaluating the automated LSE metrics, LSE-D (lower is better) and LSE-C (higher is better). We faced two issues while using this metric:

**(a):** Even though user study suggest that the lip-sync quality of Wav2Lip is usually inferior to our model, the LSE metrics are always better for Wav2Lip. We observed this pattern over a range of different videos. LSE metrics are computed from paired audio-visual representation coming from a SyncNet [9] network. Wav2Lip also leverages a SyncNet architecture and audio-visual representations as a lip-sync loss during training. We hypothesize that since a similar architecture is used both as a training loss and for scoring, Wav2Lip may be biased to do particularly well on this metric. Thus, we refrained from reporting LSE metrics for Wav2Lip. However, for other methods, the metric yields numbers consistent with human evaluations of lip-sync.

**(b):** The code sometimes fails to detect faces even under normal illumination and thereby does not give LSE metrics. So, we could not report LSE metrics on all videos.

<sup>1</sup> <https://github.com/Rudrabha/Wav2Lip>

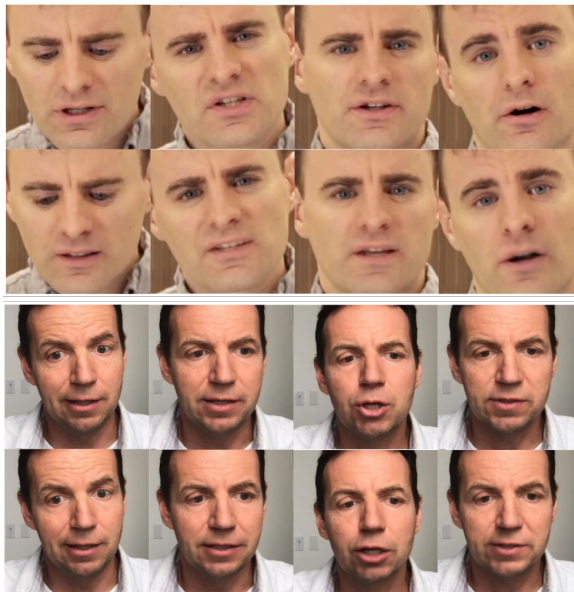


Figure 13: Comparing our result against ground-truth. For each subject, top row is the original sequence of frames, while the bottom row is the resynthesized sequence.

### E. Leaving out Wav2Lip for Self-reenactment Comparisons

While comparing methods for self-reenactment tasks (Figure 10), we do not include Wav2Lip among the competing methods. Along with the current audio, Wav2Lip also feeds in the sequence of target frames with the lip region unmasked. In a self-reenactment setting, the input target frames are same as the final expected output from the network. Hence, Wav2Lip would have an unfair advantage over our method, since our framework is entirely audio driven and only utilizes masked target frames (around the lip region) for pasting back the synthesized output. Therefore, we do not compare Wav2Lip for self-reenactment results. A similar advantage is also available to LipGAN [23] which is a precursor to the framework of Wav2Lip.

### F. Selecting Code Length for Audio and Previous Atlas

In this section we report studies to determine the code length for encoding the current time step’s spectrogram and previous time step’s predicted atlas. Since we wish to automatically determine acceptable settings of these parameters, the study was conducted for a self-reenactment task in which we have access to ground truth facial information.

For this study, we curated a custom dataset of subjects selected from YouTube instructional videos, webcam recordings, and studio conversations. The custom dataset had around 10,000 audio-synchronized frames. Sample



frames and corresponding reconstructions from two such subjects are shown in Figure 13.

### F.1. Selection of Previous Atlas Code Length, $N_a$

We conducted an ablation study to determine the length  $N_a$  of the latent code  $L_{t-1}^A$ . The code length governs the contribution of the previous visual state to the current frame. With increasing code length, the model starts to incorrectly neglect the current audio input, instead basing its output mostly on the previous state. In Table 2 we report the average metrics over different  $N_a$ . Note that  $N_a = 0$  signifies a model trained without auto-regression. Both SSIM and LMD improve when using auto-regression initially but deteriorate as we increase  $N_a$ , with  $N_a = 2$  giving the best results.

Metric	← AR Code Length: $N_a$ →			
	0	2	8	16
SSIM	0.92	<b>0.93</b>	0.901	0.889
LMD	2.00	<b>1.88</b>	2.91	2.16

Table 2: Parameter sweep for selecting latent code lengths,  $N_a$ , for encoding previous time step atlas based on LMD (lower is better) and SSIM (higher is better) metrics. Best results are marked in bold.

### F.2. Selection of Audio Code Length, $N_S$

As mentioned earlier, we used a 32-dimensional vector for encoding the audio spectrogram. This number was chosen by performing a parameter sweep over  $N_S \in \{8, 16, 32, 64, 128\}$  on our custom dataset. In Table 3 we compare the SSIM and LMD metrics, averaged over the subjects in the dataset. We observe that  $N_S = 32$  yields the most favorable performance. Hence, we use it as our default choice when encoding the spectrogram.

Metric	← Audio Code Length: $N_S$ →				
	8	16	32	64	128
LMD	1.83	1.87	<b>1.77</b>	1.81	1.82
SSIM	0.907	0.914	<b>0.917</b>	0.917	0.910

Table 3: Parameter sweep for selecting audio latent code length based on LMD and SSIM metrics. Best results are marked in bold.

## G. Network Architectures

In the main paper, we describe the architecture of the audio encoder, which computes the latent code from audio

spectrograms, and the geometry decoder, which computes the 3D vertices from the audio latent code. Here, we describe the additional network components of our model.

### G.1. Texture Decoder Architecture

In Table 4 we present details of the texture decoder. The input to the decoder is either a 32D vector (conditioned on only audio latent code), or 34D (conditioned on audio latent code + previous time step atlas latent code). This is followed by a series of convolution and upsampling layers.

### G.2. Auto-regressive Encoder Architecture

In our auto-regressive architecture, the previous atlas at the last time step is encoded as an additional latent vector (along with the audio encoded vector). The output is a latent vector of length  $N_a = 2$ . The encoder architecture for the auto-regressive model is shown in Table 5.

**Training by “Teacher Forcing”:** As stated in the main paper, during training we do not provide the actual previous predicted atlas as input to the auto-regressive model, since that would entail a recursion in the network. Instead we follow the *Teacher Forcing* [41] paradigm of training the network with the ground truth previous atlas.

In our initial experiments, we implemented a recursive network and fed in the actual predicted previous atlas to the model during training. The reconstruction quality of that approach was worse than using ground truth atlases during training (*i.e.* Teacher Forcing), however. Note that during inference, the *predicted* previous atlas is fed to the model, because the ground truth atlas is not known at that time. Also, for predicting the first frame, we provide an ‘*all-zeros*’ image as a proxy for previous frame because there is no previous frame to start with. To handle this case, we train the model by feeding it with ‘*all-zeros*’ for the previous atlas with a probability of 20%. This trains the model to reconstruct the atlas both with and without the knowledge of previous time step’s visual state.

## H. Subject Details: GRID, CREMA-D and TCD-TIMIT

For self-reenactment studies we performed experiments on GRID [10], TCD TIMIT [16] and CREMA-D [21] datasets. Following the exact setting in [37], we select the same set of 10 subjects (see Table 6) from each of the datasets).

## I. Sharpness of Synthesized Lip Region

In main paper, we mentioned that our method is capable of generating high quality lip-sync, and we objectively established this with the commonly used CPBD metric. The metric was evaluated on the entire face. While it is true

Input	Type	Kernel	Stride	Channels	Outputs
Input (Latent Vector): = 32D (only audio) = 34D (Audio + AutoRegressive)					
Latent Vector	FC	-	-	-	16384
Reshape: $4 \times 4 \times 1024$					
2× Bilinear Upsample					
$8 \times 8 \times 1024$	Conv	$3 \times 3$	$1 \times 1$	512	$8 \times 8 \times 512$
2× Bilinear Upsample					
$16 \times 16 \times 512$	Conv	$3 \times 3$	$1 \times 1$	256	$16 \times 16 \times 256$
2× Bilinear Upsample					
$32 \times 32 \times 256$	Conv	$3 \times 3$	$1 \times 1$	128	$32 \times 32 \times 128$
2× Bilinear Upsample					
$64 \times 64 \times 128$	Conv	$3 \times 3$	$1 \times 1$	64	$64 \times 64 \times 64$
2× Bilinear Upsample					
$128 \times 128 \times 64$	Conv	$5 \times 5$	$1 \times 1$	3	$128 \times 128 \times 3$

Table 4: Architecture of the texture decoder. Each fully connected (FC) and convolution layer is followed by a ReLU non-linearity, while only the last convolution layer is followed by a tanh non-linearity. The length of the input latent vector depends on the mode of the experiment.

Input	Type	Kernel	Stride	Channels	Outputs
Input (RGB): = $128 \times 128 \times 3$					
Input	Conv	$5 \times 5$	$2 \times 2$	128	$64 \times 64 \times 128$
$64 \times 64 \times 128$	Conv	$5 \times 5$	$2 \times 2$	256	$32 \times 32 \times 256$
$32 \times 32 \times 256$	Conv	$5 \times 5$	$2 \times 2$	512	$16 \times 16 \times 512$
$16 \times 16 \times 512$	Conv	$5 \times 5$	$2 \times 2$	1024	$8 \times 8 \times 1024$
$8 \times 8 \times 1024$	Conv	$5 \times 5$	$2 \times 2$	2048	$4 \times 4 \times 2048$
$4 \times 4 \times 2048$	FC	-	-	-	2

Table 5: Architecture of the encoder for the previous atlas in auto-regressive mode. The encoder input is an RGB image and output is a latent vector. Each convolution layer is followed by a ReLU non-linearity. The last fully-connected (FC) layer is followed by a tanh non-linearity.

that the sharpness of the final composite full face is a primary factor of photo-realism, we also acknowledge that our method benefits from copying the texture from upper part of face from target frames.

Here we focus on determining the sharpness of only the lower half of the face (below the nostrils). In Table 7, we compare against Wav2Lip, LipGAN, NVP and TBE on the user study videos. Even on the lower mouth region, our method attains better CPBD scores across all competing methods.

## J. Applications

Our approach of generating textured 3D geometry enables us to address a broader variety of applications than purely image-based or 3D-only techniques, as discussed here. Sample screenshots from some of these applications

are shown in Figure 14.

**3D Talking Avatars:** 3D avatars can make multiplayer online games and Virtual Reality (VR) environments more social and engaging. They may also be employed for audio/video chat applications and virtual visual assistants. While such avatars can be driven by a video feed from a web-cam or head-mounted camera, the ability to generate a 3D talking face from just audio obviates the need for any auxiliary camera device, and also helps preserve privacy, while reducing bandwidth requirements at the same time. Our technique supports generating both 3D textured faces as well as CGI avatars for these applications.

**Video creation and editing:** Our approach can be used for editing videos, *e.g.* to insert new content in an online course, or to correct an error without the cumbersome

Dataset	Test Subject ID
GRID	2, 4, 11, 13, 15, 18, 19, 25, 31, 33
TCD-TMIT	8, 9, 15, 18, 25, 28, 33, 41, 55, 56
CREMA-D	15, 20, 21, 30, 33, 52, 62, 81, 82, 89

Table 6: IDs of subjects used for self-reenactment experiment.

LipGAN	Proposed	Wav2Lip	Proposed	TBE	Proposed	NVP	Proposed
0.07	0.14	0.06	0.13	0.12	0.18	0.10	0.18

Table 7: Comparing CPBD metrics (on lower half of faces) of competing methods against our proposed method. For each method, we select common pairs of videos for the competing and our method from the pool of user study videos.



Figure 14: Sample applications enabled by our 3D talking face generation pipeline.

and sometimes impossible procedure of re-shooting the whole video under original conditions. Instead, a new audio transcript may be recorded for the edited portion, followed by applying our synthesis technique to modify the corresponding video segment. Such a speech-to-video or text-to-video system may be useful in multiple domains such as education, advertising and entertainment. We can also generate cartoon renderings for these videos, which may be preferred in some applications, *e.g.* sketch videos, stylized animations, or assistive technologies such as pronunciation visualization.

**Video translation and dubbing:** Even though we train

our models on videos in a single language, they are surprisingly robust to both different languages as well as text-to-speech (TTS) audio at inference time. Using available transcripts or a speech recognition system to obtain captions, and subsequently a text-to-speech system to generate audio, we can automatically translate and lip-sync existing videos into different languages. In conjunction with appropriate video re-timing and voice-cloning [17], the resulting videos look fairly convincing. We have employed our approach for translating and dubbing videos from English to Spanish or Mandarin, and vice-versa. Notably, in contrast to narrator-driven techniques [22, 34], our approach for video dubbing does not require a human actor in the loop, and is therefore more scalable across languages.