

# Audio- and Gaze-driven Facial Animation of Codec Avatars

Alexander Richard<sup>\*1</sup>, Colin Lea<sup>\*1†</sup>, Shugao Ma<sup>1</sup>, Juergen Gall<sup>2</sup>, Fernando de la Torre<sup>1</sup>, Yaser Sheikh<sup>1</sup>  
<sup>1</sup>Facebook Reality Labs <sup>2</sup>University of Bonn

{richardalex, shugao, ftorre, yasersh}@fb.com, colincsl@gmail.com, gall@iai.uni-bonn.de

## Abstract

Codec Avatars are a recent class of learned, photorealistic face models that accurately represent the geometry and texture of a person in 3D (i.e., for virtual reality), and are almost indistinguishable from video [27]. In this paper we describe the first approach to animate these parametric models in real-time which could be deployed on commodity virtual reality hardware using audio and/or eye tracking. Our goal is to display expressive conversations between individuals that exhibit important social signals such as laughter and excitement solely from latent cues in our lossy input signals. To this end we collected over 5 hours of high frame rate 3D face scans across three participants including traditional neutral speech as well as expressive and conversational speech. We investigate a multimodal fusion approach that dynamically identifies which sensor encoding should animate which parts of the face at any time. See the supplemental video which demonstrates our ability to generate full face motion far beyond the typically neutral lip articulations seen in competing work: <https://research.fb.com/videos/audio-and-gaze-driven-facial-animation-of-codec-avatars/>

## 1. Introduction

Advances in representing photorealistic avatars have greatly improved in recent years [30, 15, 5, 27, 28], however, the ability to animate these avatars in real-time for augmented or virtual reality (AR/VR) applications remains limited [44, 17]. The state of the art in driving these avatars requires a lengthy user-specific setup process [27], custom hardware configurations not amenable to commercial AR/VR, and/or a team of technical artists mapping facial motions from a single user to their own avatar [30, 15]. With ideal, sensor-heavy inputs (i.e. cameras pointed clearly at the face) these approaches can accurately display a user’s facial expressions, but even at best expressive speech tends to be poorly represented [44]. In this work we investigate

<sup>†</sup>author was affiliated with Facebook at time of paper writing  
<sup>\*</sup> indicates equal contribution

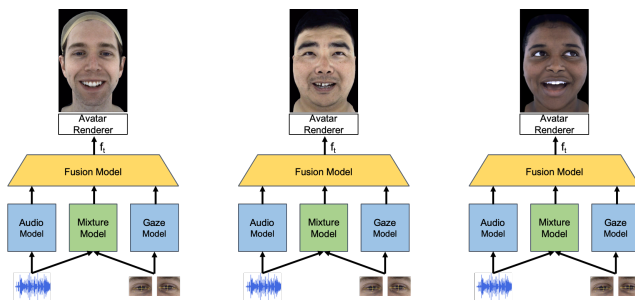


Figure 1. Our Multimodal VAE-based model predicts facial coefficients that animate a photorealistic “Codec” avatar model using only audio and gaze as input. Face images are renders from coefficients generated by our model.

an approach for driving a photorealistic model only using sensors that can be obtained with commodity hardware. Microphones are available on all VR headsets and eye tracking is available on several developer-focused AR/VR devices.<sup>1</sup>

Recent audio-driven facial animation efforts have suffered from a severe lack of data that often consist of only minutes worth of high quality facial capture (i.e., [12, 22]). One goal of this work was to investigate what kind of data is necessary to build an expressive audio-focused animation model. To this end we collected over 5 hours of data consisting of expressive, dyadic conversations across three people which were processed using the capture pipeline in [27] with paired audio, gaze, and facial coefficients per-timestep. This allows us to investigate the impact of training on different data subsets, understand how diverse data actually needs to be, and understand if our input modalities, and the corresponding models, are sufficient for achieving plausibly accurate expressive facial animation.

Multimodal fusion is a challenging problem, especially when using lossy input modalities such as audio and gaze, because there is not enough signal in either modality to accurately predict a facial expression. While there are clear correlations between speech and lip shapes there may not be any indication of when someone smiles, raises their eyebrows, or when they open their mouth to preempt another

<sup>1</sup>e.g., Vive Pro, Magic Leap One, & Qualcomm VR Dev Kit

speaker. Likewise, gaze direction has a clear effect on the eyes but it is unclear how gaze affects the lower face motion. Standard fusion approaches tend to produce more muted facial motions where each input modality plays a more fixed role, *e.g.* gaze only affects the eyes and audio affects the lower face. We describe a simple-yet-effective approach for dynamically updating how the model attends to each sensor to jointly encode correlations such as how motion from eye tracking features affect expressions like smiling.

**Contributions.** Our main contributions are summarized as:

- (a) This is the first paper to provide an extensive study on the variety of expressivity data (*e.g.*, excited conversations, descriptive tasks) required to animate natural facial motion for photorealistic avatars only from audio or from audio and eye tracking. Other approaches rely on tiny amounts of high quality “expressive” data (*e.g.*, 3-5 minutes [22] versus our 5 hours) or moderate amounts of neutral sentence reading [13, 36, 34].
- (b) We are the first to demonstrate a real-time solution for this problem using non-linear, photorealistic full-face models of geometry and texture. Note this is harder than geometry-alone [12, 22] due to non-linearities in texture-based tongue motions and lip articulations.
- (c) We discuss a fundamental issue with deep multimodal models where the network effectively learns to ignore one modality. To overcome, we describe a set of learning techniques, *e.g.* reconstruction of input modalities. We show quantitatively that this improves performance when paired with our dynamic, per-parameter multimodal fusion model, *cf.* Section 5.1.

At run-time our input is synced audio and gaze and output is a vector of facial coefficients for an avatar. We suggest viewing the supplemental video before reading this paper.

## 2. Related Work

Recent audio-driven animation efforts approaches have focused on driving lip articulations [35, 33, 46], full-face geometry [22, 16, 19, 8], and holistic video approaches [10, 45, 32]. Our work focuses on geometry and texture-based full-face animation. We find encoding dynamic changes in texture is critical for realistic lip and tongue motion.

**Lower-face Synthesis.** Taylor *et al.* [35] and Suwajanakorn *et al.* [33] generate lower-face animation (offline) by taking low level audio features (Phoneme-based in [35] and MFCCs in [33]) and predicting a set of coefficients corresponding to 2D Active Appearance Models. Results from Taylor *et al.* [35] are reasonable for neutral speech but lack nuanced facial motion or non-neutral expression. Suwajanakorn *et al.* [33] provide compelling videos of former President Barack Obama speaking, however, upper face expression comes from reference video and is not predicted. Zhou *et al.* (VisemeNet) [46] show how data-driven ap-

proaches – using one hour of lower-face landmark data – can be used to drive a set of artist-friendly visemes and jaw and lip (JALI) controls. While they improve lip articulation over their previous JALI model [14], they do not show expression such as smiles, smirks, or non-speech.

**3D Geometry.** Karras *et al.* [22] use 3-5 minutes of tracked 3D geometry, per actor, to generate expressive speech animation using linear predictive coding (LPC) audio features. While generating full-face animation is much harder than lower-only, there is relatively poor lip closure and substantial eyebrow swim. Similarly, Cudeiro *et al.* [12] collected around 3 minutes of speech for each of the 12 participants which they mapped to a learned FLAME [25] geometry model. They achieved good lip closure but because their captures all consisted of neutral speech the results are very monotone. Greenwood *et al.* [19] and Eskimez *et al.* [16] also look at full face animation but only predict sparse landmark-based marker positions which hide a lot of nuance included in high fidelity photorealistic avatars. A key limitation in many of these approaches is that they rely on phoneme- or phoneme-like approaches, which inherently remove stylistic cues important for expressive speech.

**Image-based Animation.** There has been recent interest in animating frontal face images using data in-the-wild (*i.e.*, [10, 41, 32, 42]). Typically these GAN-based papers take a single image and generate a video as if the person is speaking. While impressive, they could not be used for our class of VR use cases which assume parametric models of the face. Brand [7] did some of the earliest work in this area, far pre-dating GANs, by computing trajectories on a manifold of possible facial motions. Chung *et al.* [10] generate full-face animation using cropped frontal images using videos in-the-wild. While their approach is inherently photorealistic, their model seemingly only animates the lower face. Recent work by Vougioukas *et al.* [41], Song *et al.* [32], and Zhou *et al.* [45] have used GANs to add or improve quality of full-face expression for frontal face images.

**Traditional Audio-driven Animation.** Existing audio-driven approaches generate reasonable quality lip animations on low-fidelity stylized avatars [21, 1]. Solutions rely on artist-defined lip shape models (*visemes*) and assume a mapping between phonemes and lip articulations. Extensions to the viseme model look uncanny when applied to photorealistic avatars, in part because of their inability to distinguish between expressions (*i.e.*, talking in an excited versus sad manner) [36, 34, 14, 46]. Photorealistic avatars are frequently built on learned non-linear representations that cannot be combined with artist-created sculptures [27, 35, 22].

Other early work in this area [40, 11] demonstrated audio-driven animation on photorealistic avatars such as with active appearance models. Cao *et al.* [40] shows compelling expressive animation but is based on a motion graph

that requires offline post-processing to time-warp and blend motion snippets. Cosker *et al.* [11] also works offline by synthesizing a coherent texture map and stitching together different regions of the face.

**Multimodal and Time-series Modeling.** We build upon foundational work on variational autoencoders (VAEs) [23, 38, 20] and temporal convolutional networks (TCNs) [37, 24, 3]. Our approach uses the idea of learning a shared latent space among different modalities, which has been explored previously [29, 26]. Unlike [29] our approach does not need explicit regularization of the latent space by an adversarial loss and unlike [26] ours learns a direct mapping from input to target modalities instead of learning a style/domain transfer. See [4] for an overview on multimodal modeling. Our approach attempts to overcome the following phenomenon: with high-capacity multimodal models it is easy to overfit to one modality and thereby ignore another, which has also been investigated in two recent preprints [43, 31]. Wang *et al.* [43] describe an approach that identifies when each modality starts to overfit – using a held-out set – and introduces gradient blending to prevent one modality from dominating the prediction. A preprint by Shi *et al.* [31] describes a Mixture of Experts VAE similar to ours except they assume that each input modality should provide overlapping information: i.e., averaging modality-specific predictions provides a good estimate. In our case audio and gaze are complementary and thus we dynamically, per latent parameter, identify how each modality should be combined based on available signals at the time.

### 3. Multimodal VAE

Our goal is to generate realistic facial motion corresponding to a photorealistic 3D avatar using only audio and gaze as input. We use the deep appearance model of Lombardi *et al.* [27] and denote a sequence of facial coefficients as  $\mathbf{f} = (f_1, \dots, f_T)$  for all  $T$  time steps. Each vector of coefficients  $f_t \in \mathbb{R}^{D_f}$  is decoded into a mesh and texture map using the facial decoder proposed in [27]. We assume corresponding audio and gaze input features  $\mathbf{a} = (a_1, \dots, a_T)$  and  $\mathbf{g} = (g_1, \dots, g_T)$  where  $a_t \in \mathbb{R}^{D_a}$ ,  $g_t \in \mathbb{R}^{D_g}$ . Audio features, gaze, and facial coefficients are sampled at 100 Hz.

A straightforward approach, and one that we use as a baseline, is to train a TCN-based regressor that takes in a sequence of audio and gaze features and simply predicts the facial coefficients for each time-step. This is similar to what is done in [12] but applied to both of our modalities. We show that this approach does not handle nuanced and complementary interactions between each input modality, such as correlations between eye gaze and smiles or speech and blinking. We propose an alternative approach using a specially structured VAE that learns a shared mapping across sensor types and facial configurations, as shown in Figure 2. This model has modality-specific encoders,

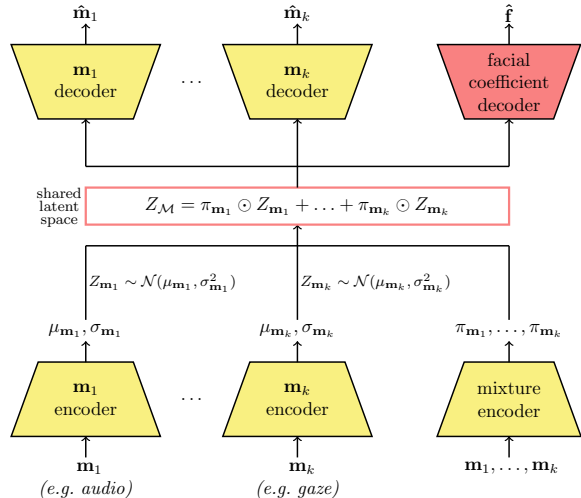


Figure 2. Our Multimodal VAE architecture. Given  $k$  input modalities  $\mathbf{m}_1, \dots, \mathbf{m}_k$  ( $k = 2$  for audio and gaze in our case),  $k$  encoders compute an embedding  $Z_{\mathbf{m}_i}$  for each modality. A mixture encoder outputs a weight for each modality which is then used to compute the shared latent embedding  $Z_{\mathcal{M}}$ . A set of decoders reconstruct all input modalities and facial coefficients  $\mathbf{f}$  from  $Z_{\mathcal{M}}$ .

modality-specific decoders, a facial coefficient decoder, and a mixture encoder that determines how to combine information from each modality. At training time, we force the model to not only predict facial coefficients from audio and gaze input but also to reconstruct both input modalities. This strategy forces the model to focus on eyes and mouth movement and improves the fidelity of the avatar. At test-time we only need the encoders and facial coefficient decoder, however, we find that reconstructing the input modalities at training time improves our model’s ability to generalize to unseen data.

#### 3.1. The Model

We start by setting notation and describing the standard VAE [23] and then extend it to our multimodal VAE.

**Preliminaries.** Consider a VAE that maps input  $\mathbf{x}$  to a latent space  $Z_{\mathbf{x}}$ , i.e. the VAE learns an encoder  $q(Z_{\mathbf{x}}|\mathbf{x})$  and decoder  $p(\mathbf{x}|Z_{\mathbf{x}})$  that maximize the evidence lower bound,

$$\text{ELBO}_{\mathbf{x}} = \mathbb{E}_{Z_{\mathbf{x}}}[\log p(\mathbf{x}|Z_{\mathbf{x}})] - \text{KL}[q(Z_{\mathbf{x}}|\mathbf{x})||p(Z_{\mathbf{x}})] \quad (1)$$

with  $\text{KL}[\cdot]$  denoting the Kullback-Leibler divergence. As in [23], we assume the latent prior  $p(Z_{\mathbf{x}})$  is an isotropic Gaussian with unit variance and the encoder  $q(Z_{\mathbf{x}}|\mathbf{x})$  models a Gaussian distribution with mean  $\mu_{\mathbf{x}}$  and diagonal covariances  $\sigma_{\mathbf{x}}^2$ . Optimizing the decoder  $p$  with the  $\ell_2$ -loss, the maximization of the ELBO is equivalent to minimizing the loss

$$\mathcal{L}_{\mathbf{x}} = \|\mathbf{x} - \hat{\mathbf{x}}\|^2 + \text{KL}[q(Z_{\mathbf{x}}|\mathbf{x})||p(Z_{\mathbf{x}})], \quad (2)$$

where  $\hat{\mathbf{x}}$  is the input reconstructed from the latent embedding  $Z_{\mathbf{x}}$ , *i.e.* the output of the decoder  $p$ . The KL term can be seen as a regularizer on the latent space, pushing it towards an isotropic Gaussian. Note that optimizing the reconstruction error using the  $\ell_2$  loss  $\|\mathbf{x} - \hat{\mathbf{x}}\|^2$  corresponds to maximizing  $p(\mathbf{x}|Z_{\mathbf{x}})$ , assuming the distribution to be an isotropic Gaussian with mean  $\hat{\mathbf{x}}$ .

**Multimodal VAE.** We formulate an alternative VAE architecture, depicted in Figure 2, that encodes multiple input modalities  $\mathbf{m}_i \in \mathcal{M}$  ( $i = 1, \dots, k$ ) into a shared latent space  $Z_{\mathcal{M}}$  using a mixture of per-modality embeddings. The decoder should be able to reconstruct all input modalities from this shared latent space, such that the model is forced to maintain sufficiently detailed information about the input modalities in the shared latent space. The multimodal loss is then  $\mathcal{L} = \mathcal{L}_{\text{rec}} + \mathcal{L}_{\text{KL}}$ , where

$$\mathcal{L}_{\text{rec}} = \sum_{\mathbf{m} \in \mathcal{M}} \|\mathbf{m} - \hat{\mathbf{m}}\|^2, \quad (3)$$

$$\mathcal{L}_{\text{KL}} = \text{KL}[q(Z_{\mathcal{M}}|\mathbf{m}_1, \dots, \mathbf{m}_M) \| p(Z_{\mathcal{M}})]. \quad (4)$$

The latent embedding  $Z_{\mathcal{M}}$  depends on all input modalities and is regularized towards an isotropic Gaussian prior  $p(Z_{\mathcal{M}})$ .

Our input modalities are audio features  $\mathbf{a}$  and gaze features  $\mathbf{g}$ , *i.e.*  $\mathcal{M} = \{\mathbf{a}, \mathbf{g}\}$ . At run-time our goal is to predict a set of facial coefficients, so we add an additional decoder that is optimized to predict the facial coefficients from the joint audio and gaze embedding  $Z_{\mathcal{M}}$ , see Figure 2 for an illustration. Formally, the reconstruction loss from Equation (3) then becomes

$$\mathcal{L}_{\text{rec}} = \|\mathbf{f} - \hat{\mathbf{f}}\|^2 + \sum_{\mathbf{m} \in \mathcal{M}} \|\mathbf{m} - \hat{\mathbf{m}}\|^2, \quad (5)$$

*i.e.* all input modalities and the facial coefficients are reconstructed from the shared latent embedding  $Z_{\mathcal{M}}$ .

**Joint Embedding of Audio and Gaze.** Our fusion approach, illustrated in Figure 2, takes an arbitrary set of  $M$  modalities  $\mathbf{m}_i \in \mathcal{M}$  using predicted mean and standard deviation  $\mu_{\mathbf{m}}$  and  $\sigma_{\mathbf{m}}$  for each input modality  $\mathbf{m}$  such that

$$Z_{\mathbf{m}} \sim \mathcal{N}(\mu_{\mathbf{m}}, \sigma_{\mathbf{m}}^2), \quad (6)$$

where  $\mu_{\mathbf{m}}, \sigma_{\mathbf{m}}^2 \in \mathbb{R}^L$  with  $L$  being the dimensionality of the latent space. A separate ‘‘mixture’’ encoder predicts mixing coefficients  $\pi_{\mathbf{m}}$  for each modality.  $\pi_{\mathbf{m}}$  is a vector containing a mixture weight for each component of the latent space and is generated using a softmax layer so each element is positive and sums to one across modalities:  $\sum_{\mathbf{m}} \pi_{\mathbf{m},d} = 1$  for each coefficient  $d$ . The shared embedding  $Z_{\mathcal{M}}$  is defined as a weighted sum of the latent random variables  $Z_{\mathbf{m}}$  of each individual modality,

$$Z_{\mathcal{M}} = \sum_{\mathbf{m} \in \mathcal{M}} \pi_{\mathbf{m}} \odot Z_{\mathbf{m}}, \quad (7)$$

where  $\odot$  is the Hadamard product. Then,

$$Z_{\mathcal{M}} \sim \mathcal{N}\left(\sum_{\mathbf{m} \in \mathcal{M}} \pi_{\mathbf{m}} \odot \mu_{\mathbf{m}}, \sum_{\mathbf{m} \in \mathcal{M}} \pi_{\mathbf{m}}^2 \odot \sigma_{\mathbf{m}}^2\right). \quad (8)$$

During training, we regularize the joint embedding  $Z_{\mathcal{M}}$  to follow an isotropic Gaussian prior using the KL divergence.

This approach was designed with two properties in mind. First, each modality should be disentangled within the latent space such that, if necessary, each modality can simultaneously drive a different part of the face. For example, if a user speaks while darting their eyes around, the audio component should help determine the mouth shapes and the eye gaze should determine the eye shapes and upper face expressions. This suggests that the mixing weights should be defined on a per-coefficient basis such that audio can drive speech-related mouth shapes while eye gaze simultaneously drives eye and upper face shapes. Second, at each point in time the model should be able to dynamically identify which modality is more useful for animating certain facial expressions. When a user is talking then the audio should drive the lip shapes, however, if there is silence then correlations with eye gaze may help indicate other facial expressions, such as smiling. Similarly, if there is a substantial amount of background noise in the audio the model should learn to ignore this signal without explicitly affecting the gaze signal. We achieve this by introducing the a weighting function that outputs the weights  $\pi_{\mathbf{m}}$ , updated at each time-step, to determine the importance of each modality for each latent coefficient.

**Implications of the model formulation.** In contrast to a conventional regression model, this multimodal VAE has several advantages. First, reconstructing the original input modalities along with the facial coefficients forces the model to maintain detailed information about both audio and gaze in the shared latent embedding. We find this results in more accurate lip closure and eye movement of the avatar. Second, the shared and weighted latent space provides some interpretability. More precisely, the model can explicitly decide on the importance of each modality depending on the temporal context. Third, the VAE formulation with the Kullback-Leibler loss on the latent space is more robust against noise and improves the overall quality of the model. An empirical evaluation in Section 5.1 shows the improvements of our proposed multimodal VAE over a conventional regression baseline and ours without the Kullback-Leibler loss.

### 3.2. Network Architecture

Our encoders and decoders are simple temporal convolutional networks (TCNs) [24] with skip connections. Each branch consists of one  $1 \times 1$  convolution to resize the dimensionality of the input (per-frame) to a fixed number of

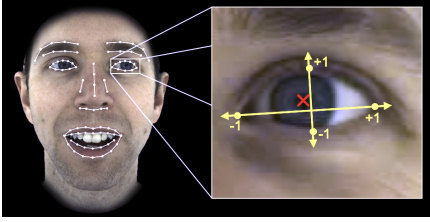


Figure 3. Left: facial landmarks used for evaluation. Right: gaze features are pupil coordinates in a normalized coordinate system defined by the eye corners and iris diameter.

channels (e.g.,  $c = 128$ ), a stack of temporal convolutions,<sup>2</sup> and a set of output  $1 \times 1$  convolutions for generating  $\mu$  and  $\sigma$  (for encoders) or a single output for each decoder. *Leaky ReLU* activations with leakage coefficient 0.2 are used after each convolution and skip connections are used between every set of stacked convolutions. We investigated other architectures, including using mechanisms such as key-value attention [39], but ultimately found that results were only incrementally better than this simpler TCN architecture.

#### 4. Experimental Setup

**Data.** We captured 5 hours of high-density 3D scans of three subjects from different ethnic backgrounds and gender using a multi-camera capture system. Figure 1 shows images of all three subjects. Tracked 3D meshes, a tracked deep active appearance model, and head pose were extracted in a similar manner as described in Lombardi *et al.* [27]. The subjects performed different types of tasks during the capture that comprised a wide variety of facial expressions and social interactions, *e.g.* reading sentences, describing images, summarizing videos, trivia games, and conversations with another person. For perspective, [22] and [12] used 3 - 5 minutes of 3D data per person. A frame-level deep face model was trained following [27] on a 45 minute subset of the data for each subject. Using the resulting appearance encoder, each frame is then mapped to a 256-dimensional vector of facial coefficients. Later, the avatar is rendered using this appearance model which generates texture and geometry from our predicted facial coefficients, *cf.* Figure 1. Two tasks are held out per subject for evaluation which corresponds to roughly 25 minutes of test data each. The remaining data is used for training. The evaluation sequences are (a) a conversational task containing a variety of natural expressions such as laughter and smiles and (b) an image description task with mostly neutral facial expression and a stronger focus on lip synchronization.

**Audio Features.** We extract 80-dimensional mel spectrograms from the raw 16kHz wave signal using the *torchaudio* pyTorch package. This feature extractor computes a spec-

<sup>2</sup>We use 5 layers per stack with kernel length=5 frames and dilation= $2^l$  for layer  $l$ .

rogram using a short-time Fourier transform (STFT) every 10ms over a 50ms Hanning window and warps the resulting 2,048 frequency bins at each timeframe onto an 80-dimensional feature vector using the Mel-scale. Mel spectrograms are widely used in tasks such as speech recognition [9] and audio-visual speech processing [2].

**Gaze Features.** We compute 2D gaze coordinates for each pupil using a normalized coordinate system as shown in Figure 3 (right). The left and right corners of the eye are defined as coordinates  $(-1, 0)$  and  $(1, 0)$ , and the perpendicular axis is scaled using the radius of the iris. This representation is invariant to scaling, translation, and rotation within the image plane. The pupil coordinates were extracted from a frontal view of each user but conceptually this information could also be extracted from a gaze tracker in a VR headset. Accurate upper face expression (*i.e.*, eyebrow motion) likely requires training models directly on raw images from eye-directed cameras in a VR headset. This was investigated by Wei *et al.* [44] using custom multi-camera hardware. The data collection and processing pipelines necessary are highly non-trivial and are beyond our scope.

**Evaluation Metrics.** Subtle facial cues are difficult to quantify but have a huge influence on human perception. While we find the metrics described below to be informative, ultimately we recommend viewing the video in the supplemental material for subjective evaluation.

For quantitative evaluation we render both the ground truth 3D avatar reconstructions and the generated avatars as videos with a resolution of  $960 \times 640$  and measure errors after running a commercial facial landmark tracker (see Figure 3). We report the mean squared error on the landmarks between the original data and the generated avatars for different facial regions, *i.e.* (1) for the 32 landmarks on the mouth, (2) for the 13 landmarks on the nose, (3) for the 20 landmarks around the eyes, (4) for the 18 landmarks on the eyebrows, and (5) averaged over all facial landmarks.

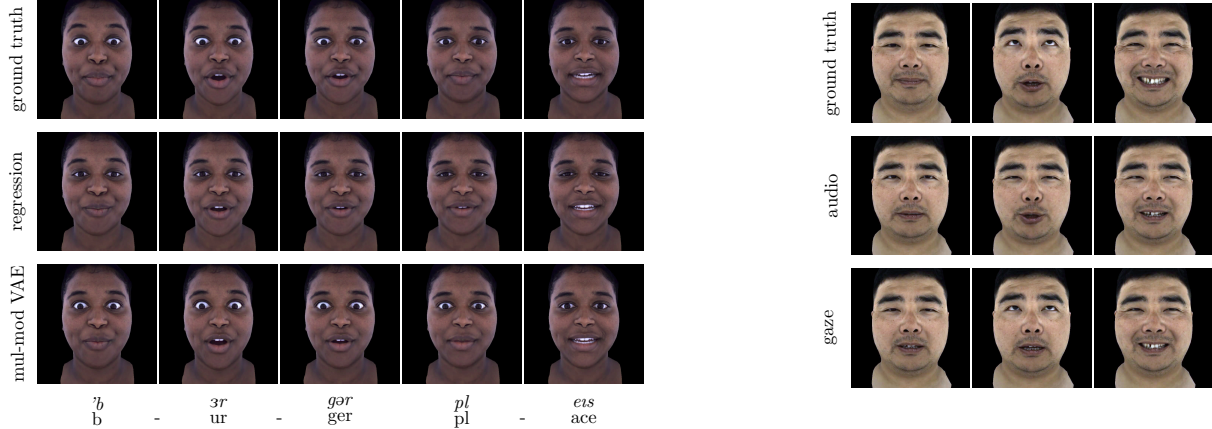
Lip closure is especially important for perceptual quality and is assessed by a second metric. We detect all lip closures in the recorded data by determining when the landmarks on the inner upper and lower lip match (*i.e.* when they do not deviate by more than two pixels each) and compare them to the lip closures detected in the generated avatars. We report the F1-score to emphasize the importance of both high precision and high recall.

#### 5. Experiments & Analysis

We describe results for three ablation studies: a comparison of various models, investigations on training data configurations, and the impact of different input modalities.

Table 1. a-d: effect of the KL loss on different latent embeddings. d vs. e: effect of reconstructing audio and gaze features during training vs. only predicting facial coefficients alone. f: a conventional regression baseline.

		Landmark Error ( $\downarrow$ )				F1-score ( $\uparrow$ )	
		eyebrows	eyes	nose	mouth	all	lip closure
(a)	no KL loss	4.27	6.84	2.27	20.46	15.15	0.557
(b)	KL on $Z_m$ ( $\mathbf{m} \in \{\mathbf{a}, \mathbf{g}\}$ )	4.18	6.08	2.12	19.03	14.12	<b>0.569</b>
(c)	KL on $Z_M$	<b>4.00</b>	5.62	<b>1.99</b>	<b>17.95</b>	<b>13.36</b>	0.546
(d)	KL on $Z_m$ ( $\mathbf{m} \in \{\mathbf{a}, \mathbf{g}\}$ ) and $Z_M$	4.06	<b>5.49</b>	2.06	18.42	13.42	0.521
(e)	no audio/gaze reconstruction	4.18	6.05	2.15	19.42	14.22	0.534
(f)	conventional regression	4.38	7.43	2.34	20.64	15.52	0.462



(a) True, regression, and Multimodal VAE results for utterance “burger place.” The regression baseline fails to close the lips at *pl* and does not capture the wide-opened eyes or raised brows.

(b) Comparison of gaze-only and audio-only models which each encode complementary signals.

Figure 4. Qualitative results of the rendered avatar.

## 5.1. Model Ablations

In this section, we analyze the components of our model and compare the multimodal VAE to a conventional TCN-based regression network. We show (a) that a structured shared latent space leads to better models, (b) that reconstructing the input modalities is beneficial for subtle facial motion such as accurate lip closure, and (c) that our proposed multimodal VAE yields more authentic and expressive facial motion than a conventional regression baseline.

**Structuring the Shared Latent Space.** Conventional VAEs learn a structured latent space by imposing a Gaussian prior on the latent embeddings. For our multimodal VAE, there are several strategies where to apply this prior. The KL loss can either be applied to the per-modality embeddings  $Z_m$ , to the joint embedding  $Z_M$  as proposed in Section 3.1, or to all  $Z_m$  and  $Z_M$ . A comparison of the landmark errors and lip closure scores of the corresponding experiments in Table 1b-d shows that structuring only the per-modality embeddings is beneficial for lip closure but degrades the overall facial geometry and appearance. Imposing the KL loss on the shared embedding  $Z_M$ , on the contrary, leads to consistent improvements in the landmark error in all facial regions at the cost of only a minor degradation in lip closure.

Applying the KL loss on the per-modality embeddings and the shared embedding (Table 1d) is inferior to regularizing  $Z_M$  only. Investigation of Equation (6) and Equation (8) shows that – in case the per-modality and shared embeddings are both regularized towards an isotropic Gaussian – the KL loss is minimized if the mixture weights are either zero or one. Therefore, regularizing both, individual and shared embeddings, strips the model of the ability to interpolate between the embeddings of different modalities.

Training the multimodal VAE without a KL loss leads to an unstructured latent space that is not particularly suitable for interpolation between the observed training embeddings and therefore does not generalize well enough on unseen data. We empirically observe this issue in Table 1a, where the landmark error is consistently worse than for any of the KL regularized variants in lines b-d.

**Reconstructing the Input Modalities.** A comparison of line c and e in Table 1 reveals that reconstructing the input modalities improves performance. If the shared latent space does not contain sufficient information about gaze and audio, the model tends to overfit the training data and fails to generate accurate lip closure, mouth shape, or eye movement. We find the reconstruction of the input modalities particularly important to generate subtle and fine-grained

Table 2. Comparison of audio and gaze models on tasks that are neutral (image description) and expressive (conversation). Note the impact of gaze on estimating the mouth shape and the impact of audio on accurate lip closure.

		Landmark Error ( $\downarrow$ )					F1-score ( $\uparrow$ )
		eyebrows	eyes	nose	mouth	all	<i>lip closure</i>
<i>conversation</i>	only audio	5.03	13.70	4.19	32.67	25.70	0.437
	only gaze	4.20	<b>5.92</b>	3.11	32.94	20.03	0.062
	audio + gaze	<b>3.66</b>	6.05	<b>2.49</b>	<b>22.95</b>	<b>15.71</b>	<b>0.450</b>
<i>descriptive</i>	only audio	5.05	9.88	1.73	13.08	14.05	0.650
	only gaze	4.78	5.27	1.89	22.95	14.88	0.075
	audio + gaze	<b>4.49</b>	<b>5.01</b>	<b>1.30</b>	<b>10.79</b>	<b>9.99</b>	<b>0.677</b>

facial motion and expressions.

**Multimodal VAE vs. Regression Baseline.** We trained a baseline regression-based TCN that receives the concatenated audio and gaze features as input and predicts the facial coefficients directly. The size of the TCN is comparable to our multimodal VAEs. Conceptually this is similar to what Cudeiro *et al.* [12] used for audio-alone except in our case we use audio and gaze. We find the baseline performs significantly worse than the proposed multimodal VAE in landmark errors and lip closure (Table 1c and f). The regression model tends towards more neutral expressions and fails to capture more articulated expressions such as excitement with widely opened eyes and raised eyebrows or subtle lip motion and accurate lip closure for plosives such as “p” (Figure 4a).

One strength of the multimodal VAE is its capability to dynamically decide on the weight each modality gets. Figure 5 shows the mixture weights of the audio modality for each component of the latent space over a one second long snippet. Interestingly, the model learns to use some latent components exclusively for a single modality. Specifically, the blue horizontal stripes represent components for which the audio weight is always zero, *i.e.* components that exclusively encode gaze information. The yellow horizontal stripes are purely dedicated to the audio modality. Many latent components use a mixture of the input modalities that varies over time depending on the temporal context.

## 5.2. Modality Ablations

### Subject Comparisons.

We compare our full system with single modality versions to analyze the impact each signal has on encoding different types of facial expressions. Experiments were evaluated on “image description” and “conversational” tasks with a single-input version of the architecture outlined in Figure 2, *i.e.* without the weights encoder and decoders for the held-out modalities. See results in Table 2. As expected, the landmark error around the eyes is large when only using audio and low when using the gaze input. It was also of no surprise that the gaze modality fails to predict lip closure. Performance for the mouth shape estimation was more surprising. For the image description task, the audio modality

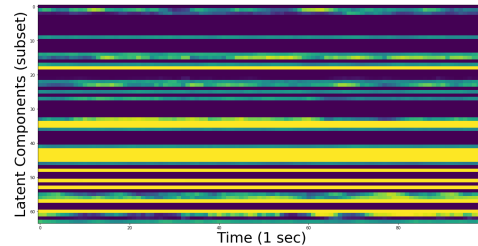


Figure 5. Mixture weights (per parameter) from a 1 second test clip. Some latent components (yellow rows) always use the audio embedding, others (blue rows) always use the gaze embedding. The rest (varying colors) change dynamically based on the input.

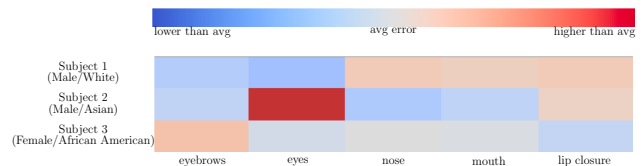


Figure 6. Relative errors per subject compared to the errors averaged over all subjects. Red indicates errors that are proportionally larger than the averaged errors, blue is lower than average.

accurately models the mouth shape estimation and the gaze modality fails to achieve comparably good results. For the conversational task, surprisingly both modalities have huge errors for the mouth shape. In this task the user frequently smiles, laughs, and gives the other person quizzical expressions. While the audio may pick up on speech related lip motion and loud laughter, we find that non-verbal, voiceless expressions like smiling are more reliably predicted from gaze, see Figure 4b for an example. Within the given data captures, smiles and laughter have strong co-occurring gaze patterns where the user squints their eyes and looking downwards, and can therefore be picked up without audio input. While each modality is indeed complementary, combining both results substantially improved performance as shown in Table 2.

### 5.3. Data Ablations

It should be no surprise that the kind of data used for training is critical for the results. In this section, we investigate the impact of different subjects, monotone versus expressive tasks, and different audio features.

Table 3. Impact of training data on different test tasks for landmark error and lip closure (Subject 1).

train data	test data	
	descriptive #2	conversation #2
	↓ landmarks / lips ↑	↓ landmarks / lips ↑
descriptive #1	8.38 / 0.715	35.66 / 0.234
conversation #1	13.22 / 0.429	14.56 / 0.216
all	8.97 / 0.712	16.06 / 0.275

Figure 6 shows a per-subject breakdown of the errors relative to the averaged errors over all subjects. Red means the error is proportionally larger than for the average over subjects, blue means it is lower. The largest outlier is the eye error for Subject 2 who has smaller eyes than the other subjects and a dark iris which makes pupil tracking less reliable. It is also interesting to observe results from Subject 3, a trained actress, in the supplementary video. She heavily moves her eyebrows during expressive speech, which is hard to synthesize correctly from only audio and gaze and leads to increased errors in that facial region. We attribute Subject 1’s slightly increased mouth landmark errors and lip closure errors to his frequent open-mouthed smiles and his tendency during conversations to open his mouth while silent to show an intent to speak. These voiceless mouth movements particularly impact lip closure and mouth shape prediction.

**Training Data Characteristics.** While deep neural networks excel at interpolating within the training distribution, they typically struggle to extrapolate beyond it. This is what has limited most prior work to generating monotone-looking speech; typical datasets simply include neutral read sentences as their training data. We show that a diverse dataset with both expressive and descriptive speech is key to not only accurate but also authentic facial animation.

Table 3 shows the performance of our approach when training the same model on three different subsets of the data. The *descriptive* tasks tend to be largely monotone, the *conversational* tasks tend to be more lively, and *all* includes a mixture of many types of expressive and descriptive speech. We evaluate on held out descriptive speech and expressive conversation to illustrate how the nature of training data affects different test scenarios. As expected, when training on the *descriptive* tasks and testing on *conversational* tasks the landmark error is poor, because the model does not generate as many expressions like smiling and laughter. Likewise when trained on *conversational* tasks the *descriptive* results deteriorate since all generated facial motion is more extreme than it should be. Overall, we find that descriptive data is responsible for accurate lip closure whereas conversational data is crucial to capture a wide range of expressions. As one would hope, when training on both, descriptive *and* expressive speech, the model generates both muted and expressive results when needed.

Table 4. Impact of different audio features (Subject 1).

audio features	Mouth Error (↓) (landmark error)	lip closure (↑) (F1 score)
phonemes	20.74	0.251
mel spectrograms	20.53	0.465

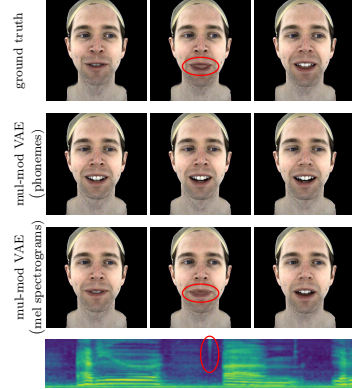


Figure 7. Our model encodes subtle audio cues such as wetting your lips with your tongue, which are only represented by a short high-frequency segment in the spectrogram.

**Phonemes vs. Mel Spectrograms.** Phonemes are a popular mid-level representation in speech processing [6, 18] and also find application in face animation [35, 14, 46]. We therefore compare results when using phonemes versus mel spectrograms. Phoneme models are trained to robustly recognize which out of 43 phonemes a user is saying and inherently abstract away the subtleties of speech. Figure 7 shows that we are capable of picking up very nuanced sounds such as licking your lips which is encoded by a tiny blip in the high frequency information. This is not encoded in the phoneme-based model. While this example may not generalize to low quality microphones or noisy environments, it highlights one of many signals that are lost from phoneme-based encodings. We find that the kind of training data, as described above, has a larger impact than the choice of audio features, however, Table 4 indicates a substantial improvement in lip closure using mel spectrograms and a modest improvement in expressivity. This is shown in more detail in the supplemental video.

## 6. Conclusion

In this work we showed that with a sufficient amount of expressive animation data we are able to map from raw audio to expressive facial animation. We find that in general our approach to multi-modal fusion is able to overcome limitations with models overfitting to individual sensors and improves animation performance. In the supplemental video we show that the quality of lip articulations from our audio-only solution can even surpass the quality from video-based solutions such as [44] which uses mouth and eye cameras. Future work may look at combining our



audio-driven model with this video-based solution.

**Ethics Remarks.** This work, and work on photorealistic avatars more broadly, have strong implications on privacy and should be approached with caution if considering real-world use cases. Restricting avatar access, for example with biometrics, is critical for preventing misrepresentation. Furthermore users should be made aware of how their avatar is being portrayed, potentially in real-time, to prevent issues where the predicted facial expression does not match the user’s intent.

## References

- [1] Amazon sumerian, 2018. <https://aws.amazon.com/sumerian/>. Last accessed 2 Aug 2020.
- [2] Triantafyllos Afouras, Joon Son Chung, Andrew Senior, Oriol Vinyals, and Andrew Zisserman. Deep audio-visual speech recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.
- [3] Shaojie Bai, J. Zico Kolter, and Vladlen Koltun. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *arXiv:1803.01271*, 2018.
- [4] Tadas Baltrusaitis, Chaitanya Ahuja, and Louis-Philippe Morency. Multimodal machine learning: A survey and taxonomy. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(2):423–443, 2019.
- [5] James Bancroft, Nafees Bin Zafar, Sean Comer, Takashi Kuribayashi, Jonathan Litt, and Thomas Miller. Mica: a photoreal character for spatial computing. In *ACM SIGGRAPH 2019 Talks*, 2019.
- [6] Maximilian Bisani and Hermann Ney. Joint-sequence models for grapheme-to-phoneme conversion. *Speech communication*, 50(5):434–451, 2008.
- [7] Matthew Brand. Voice puppetry. In *ACM Transaction on Graphics*, 1999.
- [8] Xin Chen, Chen Cao, Zehao Xue, and Wei Chu. Joint audio-video driven facial animation. In *IEEE Int. Conf. on Acoustics, Speech and Signal Processing*, pages 3046–3050, 2018.
- [9] Chung-Cheng Chiu, Tara N. Sainath, Yonghui Wu, Rohit Prabhavalkar, Patrick Nguyen, Zhifeng Chen, Anjali Kannan, Ron J. Weiss, Kanishka Rao, Ekaterina Gonina, et al. State-of-the-art speech recognition with sequence-to-sequence models. In *IEEE Int. Conf. on Acoustics, Speech and Signal Processing*, pages 4774–4778, 2018.
- [10] Joon Son Chung, Amir Jamaludin, and Andrew Zisserman. You said that? In *British Machine Vision Conference*, 2017.
- [11] Darren Cosker, David Marshall, Paul L. Rosin, and Yulia Hicks. Video realistic talking heads using hierarchical non-linear speech-appearance model. In *MIRAGE*, 2003.
- [12] Daniel Cudeiro, Timo Bolkart, Cassidy Laidlaw, Anurag Ranjan, and Michael J. Black. Capture, learning, and synthesis of 3D speaking styles. *IEEE Conf. on Computer Vision and Pattern Recognition*, 2019.
- [13] Daniel Cudeiro, Timo Bolkart, Cassidy Laidlaw, Anurag Ranjan, and Michael J. Black. Capture, learning, and synthesis of 3d speaking styles. In *IEEE Conf. on Computer Vision and Pattern Recognition*, 2019.
- [14] Pif Edwards, Chris Landreth, Eugene Fiume, and Karan Singh. JALI: An animator-centric viseme model for expressive lip synchronization. *ACM Transaction on Graphics*, 35(4), 2016.
- [15] Epic Games, 3Lateral, Tencent, and Vicon. Siren: World’s first digital human. In *Game Developers Conference*, 2018.
- [16] Sefik Emre Eskimez, Ross K. Maddox, Chenliang Xu, and Zhiyao Duan. Generating talking face landmarks from speech. In *Int. Conf. on Latent Variable Analysis and Signal Separation*, pages 372–381, 2018.
- [17] Christian Frueh, Avneesh Sud, and Vivek Kwatra. Headset removal for virtual and mixed reality. In *ACM SIGGRAPH 2017 Talks*, SIGGRAPH ’17, pages 80:1–80:2, New York, NY, USA, 2017. ACM.
- [18] Alex Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. Speech recognition with deep recurrent neural networks. In *IEEE Int. Conf. on Acoustics, Speech and Signal Processing*, pages 6645–6649, 2013.
- [19] David Greenwood, Iain Matthews, and Stephen D. Lacey. Joint learning of facial expression and head pose from speech. In *Interspeech*, pages 2484–2488, 2018.
- [20] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. In *Int. Conf. on Learning Representations*, 2017.
- [21] Sam Johnson, Colin Lea, and Ronit Kassis. Tech note: Enhancing oculus lipsync with deep learning. 2018. [developer.oculus.com/blog/tech-note-enhancing-oculus-lipsync-with-deep-learning/](https://developer.oculus.com/blog/tech-note-enhancing-oculus-lipsync-with-deep-learning/). Last accessed 2 Aug 2020.
- [22] 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 Transaction on Graphics*, 36(4), 2017.
- [23] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In *Int. Conf. on Learning Representations*, 2014.
- [24] Colin Lea, Michael D. Flynn, Rene Vidal, Austin Reiter, and Gregory D. Hager. Temporal convolutional networks for action segmentation and detection. In *IEEE Conf. on Computer Vision and Pattern Recognition*, 2017.
- [25] Tianye Li, Timo Bolkart, Michael J. Black, Hao Li, and Javier Romero. Learning a model of facial shape and expression from 4D scans. *ACM Transaction on Graphics*, 36(6), 2017.
- [26] Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. In *Advances in Neural Information Processing Systems*, pages 700–708, 2017.
- [27] Stephen Lombardi, Jason Saragih, Tomas Simon, and Yaser Sheikh. Deep appearance models for face rendering. *ACM Transaction on Graphics*, 37(4), 2018.
- [28] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *ACM Transaction on Graphics*, 38(4), 2019.
- [29] Noam Mor, Lior Wolf, Adam Polyak, and Yaniv Taigman. A universal music translation network. *Int. Conf. on Learning Representations*, 2019.

- [30] Mike Seymour, Chris Evans, and Kim Libreri. Meet mike: Epic avatars. *ACM SIGGRAPH 2017 VR Village*, 2017.
- [31] Yuge Shi, Narayanaswamy Siddharth, Brooks Paige, and Philip Torr. Variational mixture-of-experts autoencoders for multi-modal deep generative models. In *Advances in Neural Information Processing Systems*, 2019.
- [32] Yang Song, Jingwen Zhu, Xiaolong Wang, and Hairong Qi. Talking face generation by conditional recurrent adversarial network. In *Int. Joint Conf. on Artificial Intelligence*, 2019.
- [33] Supasorn Suwajanakorn, Steven M. Seitz, and Ira Kemelmacher-Shlizerman. Synthesizing obama: Learning lip sync from audio. *ACM Transaction on Graphics*, 36(4), 2017.
- [34] Sarah Taylor, Akihiro Kato, Iain A. Matthews, and Ben P. Milner. Audio-to-visual speech conversion using deep neural networks. In *Interspeech*, pages 1482–1486, 2016.
- [35] Sarah Taylor, Taehwan Kim, Yisong Yue, Moshe Mahler, James Krahe, Anastasio Garcia Rodriguez, Jessica Hodgins, and Iain Matthews. A deep learning approach for generalized speech animation. *ACM Transaction on Graphics*, 36(4), 2017.
- [36] Sarah L. Taylor, Moshe Mahler, Barry-John Theobald, and Iain Matthews. Dynamic units of visual speech. In *Proc. of the ACM SIGGRAPH/Eurographics Symposium on Computer Animation*, pages 275–284, 2012.
- [37] Aäron Van Den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew W. Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. In *ISCA Speech Synthesis Workshop*, page 125, 2016.
- [38] Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. In *Advances in Neural Information Processing Systems*, pages 6306–6315, 2017.
- [39] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008, 2017.
- [40] Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic. Expressive speech-driven facial animation. In *ACM Transaction on Graphics*, 2005.
- [41] Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic. End-to-end speech-driven facial animation with temporal gans. In *British Machine Vision Conference*, 2018.
- [42] Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic. Realistic speech-driven facial animation with gans. In *International Journal on Computer Vision*, 2019.
- [43] Weiyao Wang, Du Tran, and Matt Feiszli. What makes training multi-modal networks hard? *arXiv:1905.12681*, 2019.
- [44] Shih-En Wei, Jason Saragih, Tomas Simon, Adam W. Harley, Stephen Lombardi, Michal Purdoch, Alexander Hypes, Dawei Wang, Hernan Badino, and Yaser Sheikh. VR facial animation via multiview image translation. *ACM Transaction on Graphics*, 38(4), 2019.
- [45] Hang Zhou, Yu Liu, Liu Liu, Ping Luo, and Xiaogang Wang. Talking face generation by adversarially disentangled audio-visual representation. In *AAAI Conf. on Artificial Intelligence*, 2019.
- [46] Yang Zhou, Zhan Xu, Chris Landreth, Evangelos Kalogerakis, Subhansu Maji, and Karan Singh. Visemenet: Audio-driven animator-centric speech animation. *ACM Transaction on Graphics*, 37(4), 2018.