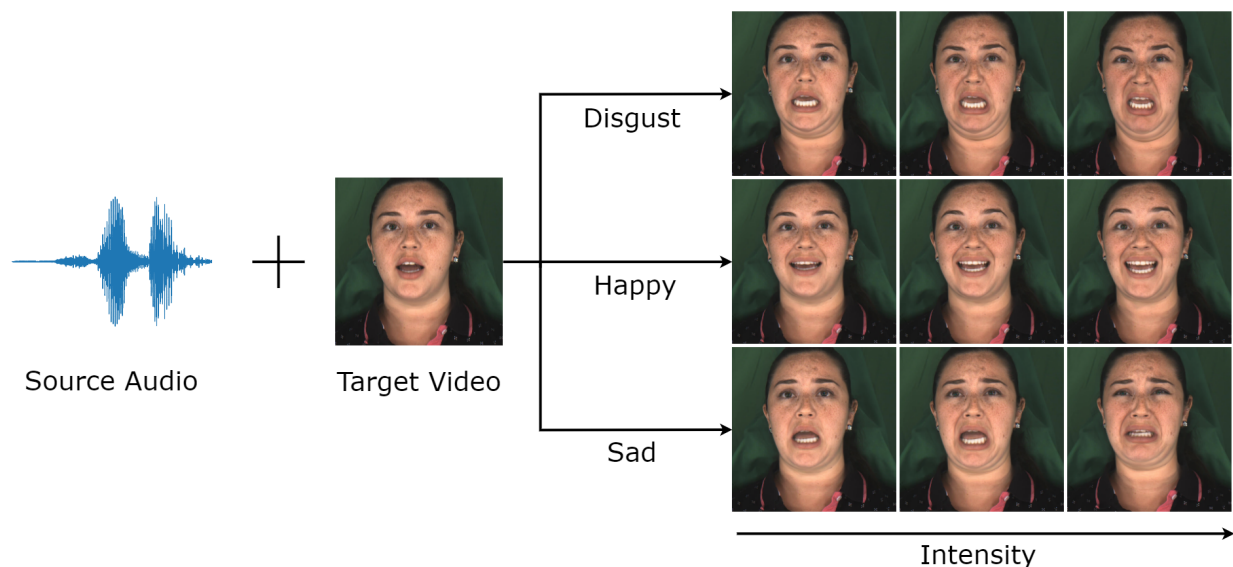


READ Avatars: Realistic Emotion-controllable Audio Driven Avatars

Jack Saunders & Vinay Nambodiri

August 22, 2023



Our method, **READ Avatars** allows for the creation of photo-realistic and lip synchronized video from audio and a reference video, with control over emotion. The same audio can be used to generate videos in multiple emotions. The intensity of each emotion can be directly specified, allowing for fine-grained control over the output.

Abstract

We present READ Avatars, a 3D-based approach for generating 2D avatars that are driven by audio input with direct and granular control over the emotion. Previous methods are unable to achieve realistic animation due to the many-to-many nature of audio to expression mappings. We alleviate this issue by introducing an adversarial loss in the audio-to-expression generation process. This removes the smoothing

effect of regression-based models and helps to improve the realism and expressiveness of the generated avatars. We note furthermore, that audio should be directly utilized when generating mouth interiors and that other 3D-based methods do not attempt this. We address this with audio-conditioned neural textures, which are resolution-independent. To evaluate the performance of our method, we perform quantitative and qualitative experiments, including a user study. We also propose a new metric for comparing

how well an actor’s emotion is reconstructed in the generated avatar. Our results show that our approach outperforms state of the art audio-driven avatar generation methods across several metrics.

1 Introduction

Generating convincing talking head video is a highly desired capability in various applications, such as film and television dubbing, video games and photo-realistic video assistants. While significant progress has been made in this area [3, 8, 12, 15, 16, 21, 25, 28, 30–32, 36, 39, 42], most existing methods produce either low-quality but accurate lip sync using 2D models [3, 25, 30, 36] or high-quality but inconsistent lip sync using 3D models [15, 16, 21, 31, 32, 39]. We hypothesize that two key factors have prevented the development of models that are both high-quality and lip synchronized. The first is that audio to expression is a many-to-many mapping. A given audio can correspond to many lip shapes, and the same lip shapes can produce different audio due to factors such as the larynx. The second factor is that while 3D models improve the visual quality by introducing strong priors, they struggle to represent complex lip shapes and do not model the mouth interior (see Figure 2).

Furthermore, it is desirable to introduce additional signals, such as emotion, to generate more realistic and believable video, and to offer users a level of control over the outputs. A few prior works attempt this [12, 13, 15, 24, 37]. These methods usually either consider emotion as discrete categories [15, 37], which gives semantic control but lacks granularity, or learn latent encodings of emotion [12, 13, 24] which allow for fine-grained control but is not semantic and requires selecting emotions from other sources (video or audio).

In this paper, we introduce READ Avatars, a method for generating talking head video with direct and granular control over emotion, while achieving high levels of lip sync, emotional clarity, and visual quality. We build upon 3D-based approaches, using a morphable model [20] as an intermediate representation of the face and deferred neural rendering [32] to achieve high visual quality. To address the above

issues causing poor lip sync in 3D models, we propose two novel components. First, we add an adversarial loss to the audio-to-expression generator to alleviate the many-to-many mapping issue. Second, we overcome the challenge of representing complex lip shapes and mouth interiors with a morphable model by conditioning a neural texture on audio, encoding audio features on the surface of the mesh using a resolution-independent neural texture based on a SIREN network [29].

In summary our contributions are:

- A novel neural rendering algorithm that leverages neural textures, operates directly on UV coordinates, and can be conditioned on audio, improving the mouth interior.
- The incorporation of a GAN loss into the audio-to-expression network to improve the results by solving the many-to-many issue of audio-to-expression generation.
- A new metric for determining how well an actor’s emotions are captured and reconstructed.

2 Related Works

2.1 Unconditional Audio-Driven Face Models

The task of generating lip-synced video from audio alone, or else from audio and a reference video, known as unconditional audio-driven video generation, has been widely studied and has numerous practical applications, such as dubbing and digital avatars. There are two broad categories of unconditional models: those that use 3D priors and those 2D models that do not.

2D Models: Many approaches to synthesizing talking head videos from audio operate directly in the image or video domain [3, 25, 30]. These methods typically employ an encoder-decoder architecture. ATVG [3] uses audio to control 2D landmarks, which are then used to generate video with attention to highlight the parts that need editing. Wav2Lip [25]

significantly improves lip sync accuracy by minimizing the distance between the audio and generated video according to a pre-trained lip sync detection network. While the lip sync is excellent, the visual quality is poor. Recently, a context-aware transformer [30] was applied to this problem, with an audio-injected refinement network that significantly improves the visual quality. However, all 2D based models to date suffer from limited visual quality. In contrast, our 3D-based method produces much higher quality videos.

3D Guided Face Models: Using explicit 3D supervision, ultra high-quality face models driven by various signals have been created [15, 16, 18, 21, 32, 34, 39]. These methods simplify facial synthesis by modeling the underlying 3D scene with a small set of parameters, such as a 3DMM [2, 5], that can be directly controlled. Despite their high visual quality, these models often lack expressiveness due to the many-to-many mapping problem and the limited lip expressions of the underlying geometry.

Puppetry methods [15, 16], and motion models [12, 28] are able to somewhat solve the many-to-many issue by using a source actor to drive the expressions. This provides a signal, the source actor’s expressions, which is much closer to one-to-one with the target actor. However, it is often undesirable to require a source actor to be filmed, and the resulting video processed every time the model is used. Actor-free methods such as ours are significantly more scalable. Implicit models [7, 8, 23, 43] augment geometry using MLP offsets, allowing for more expressive lip shapes but do not solve the many-to-many problem.

Concurrent work [31] addresses both the many-to-many issue, and the mouth interior using memory networks. However, they rely on reusing explicit pixels from the mouth region, which leads to jitter in the final videos.

2.2 Audio-Driven Face Models with Emotional Control

Only a small number of works have attempted to develop models that allow for explicit control of stylistic attributes, such as emotion, in generated talking head videos.

MEAD [37] introduces an audio-driven model with control over emotion. They trained a network to map audio to landmarks and another to convert input images to the desired emotion, and then used a UNET-based network to combine the upper face with the desired emotion and the generated landmarks to produce the final image. This method can control emotion and intensity, but lacks temporal coherence due to its frame-by-frame nature and has suboptimal lip sync. EVP [13] improves upon MEAD by adding an emotion disentanglement network to separate content and emotion in the audio and using a face-synthesis network based on vid2vid [38] to produce higher quality videos with temporal consistency. However, the lip sync and emotional clarity were not always satisfactory. MEAD and EVP are the most similar to our work. These landmark based models are unable to produce as high quality results as 3D models, they also suffer from unnatural motion without the strong priors of a 3DMM.

Kim et. al. [15] propose a 3D based model with control over style capable of producing highly realistic videos. The model uses a style translation network to convert the animation style of a source actor to that of a target. This, however, requires training a neural renderer for every emotion and a style translation network for every pair of emotions, which quickly becomes intractable for many styles. Similarly EAMM [12] uses a source actor to generate the emotional style, and a motion model to produce the lip motion. This method is unable to model the emotion in the mouth region, and as the emotional style is directly copied from a source actor, the result can appear unnatural on the target actor. NED [24] also manipulates the emotion of a source actor effectively using an emotional manipulation network in the parameter space of a 3DMM. Their method can semantically control the emotion, but also suffers from some unnatural motion owing to a mismatch of emotional style due to the generality of the model.

3 Method

Our method consists of three stages: fitting a 3D Morphable Model (3DMM) to the input videos (Sec-

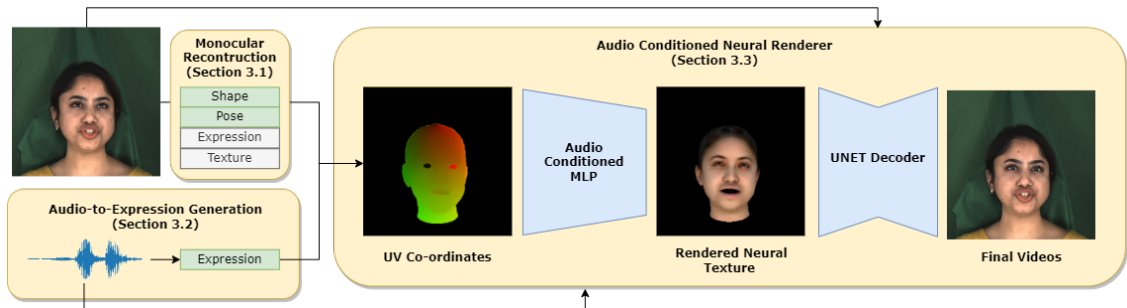


Figure 1: The READ Avatars pipeline. We train 3 separate networks. The first converts audio into expressions (Section 3.2). After rasterization, the second, MLP-based network converts the uv coordinates to neural texture features, conditioning on audio (Section 3.3). The final, UNET-based network takes this rasterized neural texture, and the real video frame, and outputs lip synchronized and emotional video frames.

tion 3.1), generating morphable model parameters from audio using adversarial training (Section 3.2), and training an audio-conditioned deferred neural renderer to produce the final, photo realistic video outputs (Section 3.3). These steps are shown in Figure 1.

In the first stage, we fit a 3D Morphable Model [2, 19, 20, 44] to the input videos using an extension of the Face2Face monocular reconstruction algorithm [34], with the modification of including blink blendshapes for both eyes. We use the implementation provided in Neural Head Avatars [7]. In the second stage, we train a neural network inspired by Pix2Pix [11] to generate morphable model parameters from audio using adversarial training. This allows us to generate realistic animation sequences, even in areas that are not well correlated with audio. At this stage, we introduce a fine-grained emotional label. In the final stage, we train an audio-conditioned deferred neural renderer [33]. Our model uses a SIREN MLP [29] to directly map uv coordinates to texture features, replacing the learned neural texture of previous work [33], and making the task of audio conditioning much easier.

3.1 Monocular Reconstruction

In the first stage of our method, we aim to find a low-dimensional set of parameters that can model a video

sequence $V = (V_0, \dots, V_n)$. For this purpose, we use the FLAME model [20], which represents explicit 3D geometry using a combination of skinned joints and blendshapes. The FLAME model can be represented as a function \mathcal{V} that maps a set of parameters for shape β , expression ψ , and joint rotations θ onto 5023 3D vertices:

$$\mathcal{V} : \mathbb{R}^{|\beta| \times |\theta| \times |\psi|} \rightarrow \mathbb{R}^{5023 \times 3} \quad (1)$$

A similar function is used to map a set of texture parameters α onto UV-based 2D textures.

We model the rendering process using a full perspective camera that projects a mesh M onto the image plane according to:

$$\hat{I} = \Pi(M, K, \mathbf{R}, \mathbf{t}) \quad (2)$$

Where \mathbf{R} and \mathbf{t} are the rotation and translation of the mesh in world space, and K is the camera intrinsic matrix. We also model the lighting as distant and diffuse, using 3-band spherical harmonics with parameters γ . If we define $\pi = (\alpha, \beta, \theta, \psi, \gamma, \mathbf{R}, \mathbf{t}, K)$ as the set of all parameters, then our objective is to find the optimal $\tilde{\pi}$ that best fits a given image.

To fit the FLAME model to our videos, we adopt the tracking model of Neural Head Avatars [7], which is based on Face2Face [34]. This model uses differentiable rendering in Pytorch3D [26] to minimize the L_1 distance between real and rendered frames, with statistical regularization over the parameters.

We assume that shape, texture, and lighting are fixed for a given actor, as our data is captured under controlled conditions. Therefore, we can first estimate $\pi_{fix} = (\alpha, \gamma, K)$, the parameters that are fixed across all videos for a given subject. These parameters are then fixed for all frames in all videos of the same subject. We can then estimate $\pi_{var} = (\theta, \psi, \mathbf{R}, \mathbf{t})$ on a per-frame basis. These parameters are then the target of the audio-to-parameter generator.

3.2 Audio-to-Parameter Generator

The goal of our method is to animate photo-realistic avatars using audio as the control signal. Previous approaches based on 3D Morphable Models [15, 18, 32, 39] use audio-to-parameter generators to puppeteer a target actor using audio or a source video. These methods rely on neural networks with regression losses to generate a subset of the parameters. Such methods, however, suffer from the many-to-many issue when mapping audio to expressions, as multiple, equally-valid expressions can come from the same audio. Regression based losses mean that weakly correlated parameters such as upper face motion is almost entirely averaged out, while even highly correlated parameters such as the lip and jaw are over-smoothed. To address these issues, we propose an audio-to-parameter generator based on a conditional GAN [6]. Our model is similar to Pix2Pix [11], using a combination of L_1 loss for low frequency parts of the data and an adversarial loss for increased realism.

The input of this network is a section of audio represented as MFCC coefficients, \mathbf{A} , together with an explicit emotion label. The window size of the MFCC is selected to be a multiple of the video frame rate. We use an explicit emotion labeling system to introduce emotion into the generated parameters. For N emotions, we use an $N - 1$ dimensional label, \mathbf{C} , with neutral emotion represented as a zero vector (absence of emotion). Each other emotion is assigned to a dimension and scaled by intensity, with the maximum intensity being 1. This continuous label allows for fine-grained control over the emotion. We distribute the label over the time dimension to obtain



Figure 2: An example of a typical failure in the monocular reconstruction method. (*Left*): The input frame, (*right*): the reconstruction. Note how FLAME model lacks the expressiveness to capture certain mouth shapes, in this case an "O".

$\mathbf{C} = (C_0, \dots, C_n)$. This label is concatenated with the MFCC audio features and serves as the input to the audio to expression generator \mathcal{G}_a , which produces the target parameters for each frame $(\pi_0, \dots \pi_n)$.

$$(\pi_0, \dots \pi_n) = \mathcal{G}_A(\mathbf{A}, \mathbf{C}) \quad (3)$$

The discriminator is conditional, and takes either the real or generated parameters, together with the audio and emotional label and predicts if the given parameters are real or generated.

Both the generator and discriminator networks use an encoder-temporal-decoder model, projecting the audio features into a high-dimensional latent space via a fully connected layer followed by an LSTM [10] and a fully connected decoder to map from this latent space to the parameters. We optimize the objective:

$$\mathcal{L} = \mathcal{L}_1 + \lambda_{\text{GAN}}\mathcal{L}_{\text{GAN}} + \lambda_{\text{vel}}\mathcal{L}_{\text{vel}} \quad (4)$$

where \mathcal{L}_1 is the ℓ_1 distance between the real and predicted parameters, \mathcal{L}_{GAN} is the adversarial loss and \mathcal{L}_{vel} is an ℓ_1 distance between the velocities of the output animation. The velocity loss is known to improve the temporal consistency of speech-driven animation [4]. Each λ is a relative weight, we use $\lambda_{\text{GAN}} = 0.02$ and $\lambda_{\text{vel}} = 100$.

Table 1: **Quantitative comparisons to state-of-the-art.** We compare visual quality using FID, lip sync with LSE-D/C [25] and emotional reconstruction with our metrics A/V-EMD. We compare our results with ATVG [3], MEAD [37] and Audio-driven Emotional Video Portraits [13] (EVP)

Method	LSE-C \uparrow	LSE-D \downarrow	FID \downarrow	A-EMD \downarrow	V-EMD \downarrow
ATVG [3]	5.705	8.731	120.040	0.160	0.239
MEAD [37]	4.080	10.569	38.015	0.974	0.113
EVP [13]	4.061	11.514	43.972	0.119	0.126
Ours	4.431	10.157	13.600	0.0686	0.093

3.3 Audio-Conditioned Neural Renderer

We next consider how to invert the parametric model fitting and produce photo-realistic video. Given a set of parameters (π_i) corresponding to a video \mathbf{V} , we aim to reproduce the video as faithfully as possible. We build upon the idea of neural textures [33], jointly optimizing an image-to-image deferred neural renderer, and a neural texture defined in UV space. However, we find that the FLAME model is not expressive enough to represent complex lip motions (see Figure 2), and that neural textures alone are not sufficient to compensate for this. Furthermore, the FLAME mesh provides no information about the interior of the mouth including the tongue and teeth.

To address this issue, we propose audio-conditioned neural textures. The aim is to encode audio information on the surface of the mesh to allow for more complex lip shapes and mouth interiors to be learned. We replace the static, learned neural texture with a SIREN MLP [29] texture network \mathcal{T} , which maps a uv coordinate of the rasterized meshes directly to a feature vector. This bypasses the need for texture lookup which is slow and limits the resolution of the neural textures. The use of a network also allows us to easily condition on audio by simply concatenating audio features to the uv coordinates.

We do not want the audio in the network to be biased by emotion or identity, as this will prevent us easily changing the emotion of the generated videos. To remove such information, we use the output of a Wav2Vec2 network [1, 40] pretrained to predict phoneme probabilities at a 50 fps. The 50 most com-

mon phonemes comprise over 99% of the audio data, therefore we restrict the features to these only. Next, we resample these probabilities to 60 fps, twice the frame rate of the video. We take a window of W frames centered at the target video frame and use a small neural network $\mathcal{A} : \mathbb{R}^{2W \times 50}$ to encode the audio over this window into a single vector $\mathbf{a}_{enc} \in \mathbb{R}^{N_a}$, where N_a is the dimension of the encoded audio and is a hyperparameter. The encoded audio vector is then used to condition the neural texture.

The audio encoder consists of several fully connected layers, followed by temporal convolutions, a reshaping layer that removes the time dimension, and a final fully connected decoder. The now encoded audio vector \mathbf{a}_{enc} is concatenated with the UV coordinates obtained during rasterization, which serves as the input to our texture network.

The output of the texture network is a multi-channel image, which appears as a rasterization of the mesh with a neural texture. This texture encodes audio information on the surface. Similar to [33], we use a 16-channel neural texture, enabling representation of higher-order lighting effects. This rasterized image is then passed through a UNET [27] based deferred neural renderer \mathcal{R} , which produces a photorealistic final frame, leveraging the audio features encoded on the mesh.

To address the issue of jitter in the final video, we include additional renderings for the frames in a window of length W_R centered on the target frame. These additional renderings are based on the rasterized UV coordinates of the parameters from the frames on either side of the target frame. This results in an input with $16 + 2W_R$ channels for the UNET

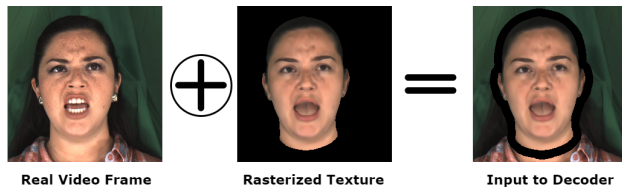


Figure 3: The input to the decoder network consists of a rendered neural texture and a real frame with a border of black pixels. Note that the texture cannot represent the mouth interior, but it is generated by the decoder.

decoder, which is able to smooth out the jitter over a window of frames. We have found this approach to be effective in reducing jitter in the final video.

In order to produce realistic and temporally consistent background in our videos, we blend the output of our texture network \hat{V}_t with the original video frames V_t . We do this by using the alpha channel from the rendered mesh as a mask to separate the foreground and background. The foreground mask, α is expanded by a fixed number of pixels to obtain α_{exp} . This expanded mask is zeroed out of the real frame and the foreground mask α is used to fill in these pixels with the rendered mesh. This process gives an input consisting of the rendered mesh in the real frame, with a border that the decoder can inpaint. This is best shown in Figure 3.

$$V_t^* = \alpha \hat{V}_t + (1 - \alpha_{\text{exp}}) V_t \quad (5)$$

To train the audio encoder, texture network, and deferred neural renderer, we optimize the following objective function end-to-end:

$$\mathcal{L}(V^*, V) = \lambda_1 \mathcal{L}_1 + \lambda_{\text{VGG}} \mathcal{L}_{\text{VGG}} + \lambda_{\text{GAN}} \mathcal{L}_{\text{GAN}} \quad (6)$$

Here, \mathcal{L}_1 is the ℓ_1 distance between the real and generated frames, \mathcal{L}_{VGG} is a VGG-based style loss [14], and \mathcal{L}_{GAN} is an adversarial loss. The hyperparameters λ_1 , λ_{VGG} , and λ_{GAN} are used to weight the importance of each loss. We use $\lambda_1 = \lambda_{\text{VGG}} = 1$ and $\lambda_{\text{GAN}} = 0.01$.

3.4 Implementation Details

To prepare the data for our method, we first crop every frame to a square shape of 256 pixels. We do this by estimating a bounding box for each frame with padding, then finding the smallest square that covers the union of these boxes. We reshape this square to the desired resolution. We implement our pipeline in Pytorch and Pytorch3D. Our models are trained on a single NVIDIA RTX3080 graphics card. All networks are optimized using Adam [17] with a learning rate of 0.0001. The renderer is trained for 5 epochs, taking about 15 hours, while the audio-to-expression generator is trained for 10 epochs taking around 5 hours. We use the LSGAN formulation for all adversarial training [22].

4 Results

4.1 Dataset

We use the MEAD dataset [37] for our experiments, which includes 60 actors (both male and female) speaking 30 sentences in 8 different emotions at 3 levels of intensity, recorded from multiple angles. For this work, we only use the front-facing camera footage. We follow the train-test split outlined in MEAD and train models for 4 of these subjects. Figures ?? and 7 show a selection of results coming from our method, across multiple emotions and intensities.

4.2 Quantitative metrics

In order to evaluate the performance of our method, we consider three qualities: visual quality, lip sync, and emotional clarity.

For visual quality, we use the Fréchet Inception Distance (FID) metric to measure the similarity of the generated frames to the ground truth. We crop the frames tightly around the face region in order to avoid biasing the results towards our method, which uses the ground truth background. To measure lip synchronization, we use the Lip Sync Error (LSE) metrics introduced in wav2lip [25]. These metrics are calculated using a pre-trained syncnet and include

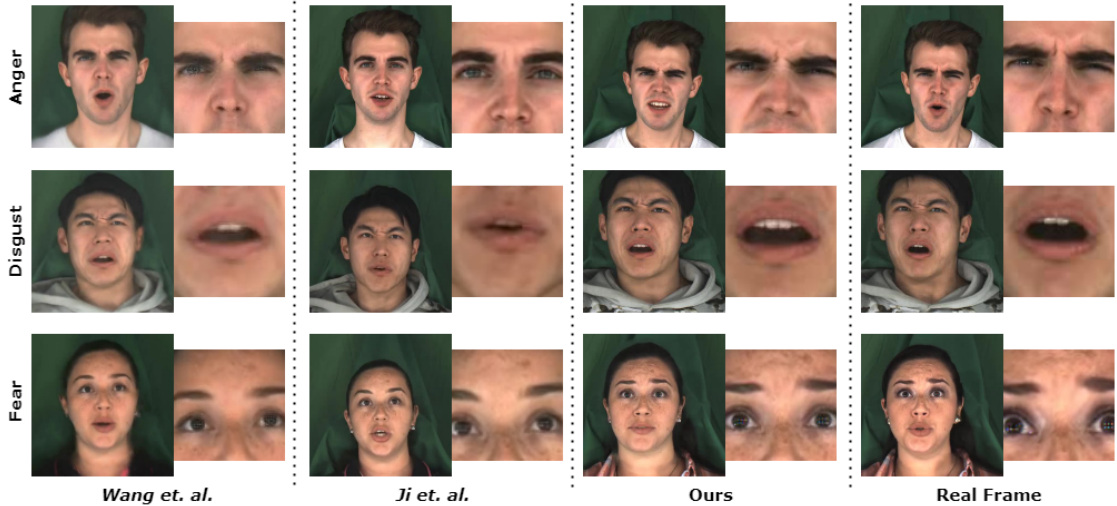


Figure 4: We compare our results to those of MEAD [37] and EVP [13]. Our results are of much higher visual quality than those of MEAD, the zoomed-in regions demonstrate that our method produces more convincing and accurate emotions compared to EVP.

LSE-D, which measures the minimum distance between audio and video features, and LSE-C, which measures the confidence that the audio and video are synchronized. To measure emotional clarity, it is not enough to measure differences at the frame level, as the intensity of emotion naturally varies over a video. We therefore introduce a new metric for emotional clarity that measures the differences between distributions of emotion. This metric is based on a pre-trained EmoNet model [35]. We predict the valence and arousal for each frame, and compare the distributions of these values between the generated and ground truth videos. We approximate the distance between these distributions using the Earth Movers Distance, and compute this distance for each subject and emotion separately, taking the average to obtain the valence Emotional Mean Distance (V-EMD) and arousal Emotional Mean Distance (A-EMD) metrics.

4.3 Comparisons to State-of-the-Art

We compare our method to several state-of-the-art audio-driven avatar models. Our main comparisons are with Audio-Driven Expressive Video Generation

(EVP) [13] and Multimodal Emotion-aware Dataset (MEAD) [37]. MEAD uses an audio-to-landmark LSTM, an emotion transformer to alter the audio-driven landmarks to any given emotion, and a final UNET-based model to produce output frames from the emotional landmarks. AudioDrivenEVP improves on this approach by designing a disentanglement model to separate audio into emotion and content, which is then used with a landmark alignment method to control for pose, and a video-to-video network that produces high-quality and temporally stable video from landmarks. We also compare to ATVG [3], a 2D-based method that excels in lip synchronization but has poor visual quality, and is unable to edit emotion.

Quantitative: The results of these comparisons are shown in Table 1. Our results outperform all competitors on visual quality (FID) and emotional reconstruction (A/V-EMD). While our method is not able to reach the lip-sync quality of the unconditional ATVG [3], it has far better visual quality and emotional clarity. We outperform both methods capable of controlling emotion: MEAD [37] and EVP [13] on lip sync.

Table 2: Results of the user study. We ask users to select their preference between our video and each of the competitors for four criteria. Where ours is preferred strongly (weakly) we denote the result ++ (+), where there is no preference, 0 and where the other method is preferred strongly (weakly), -- (-). The data in all but the rightmost column is in percentages. Note the data is rounded and may not add to 100%.

Statement	--	-	0	+	++	mean
Ours ζ MEAD (lip-sync)	1	12	14	28	44	+1.02
Ours ζ MEAD (visual quality)	0	3	9	23	65	+1.49
Ours ζ MEAD (naturalness)	1	1	15	23	65	+1.39
Ours ζ MEAD (emotion)	1	3	21	36	38	+1.06
Ours ζ EVP (lip-sync)	7	16	22	42	11	+0.50
Ours ζ EVP (visual quality)	2	22	23	42	11	+0.39
Ours ζ EVP (naturalness)	7	17	18	26	38	+0.49
Ours ζ EVP (emotion)	4	8	22	26	38	+0.87
Ours ζ Real (lip-sync)	69	28	3	1	0	-1.64
Ours ζ Real (visual quality)	33	41	23	2	0	-1.05
Ours ζ Real (naturalness)	53	38	8	1	0	-1.43
Ours ζ Real (emotion)	40	38	19	3	1	-1.13

Table 3: **Ablation study.** We compare visual quality using FID, lip sync with LSE-D/C [25] and emotional reconstruction with our metrics A/V-EMD. We compare our full model to the same model both without the GAN loss in the audio-to-expression generator and without the audio conditioned neural texture.

Method	LSE-C \uparrow	LSE-D \downarrow	FID \downarrow	A-EMD \downarrow	V-EMD \downarrow
Ours	4.431	10.157	13.600	0.069	0.093
Ours w/o GAN loss	4.047	10.446	12.587	0.079	0.090
Ours w/o audio texture	4.175	10.398	15.96	0.069	0.96

Qualitative: Figure 4 shows our results in comparison to MEAD [37] and EVP [13]. Our results show clearly better visual quality than MEAD. Compared with EVP our method is capable of producing emotion that is much clearer and more closely matches the real videos, the expanded regions highlight this. In particular, it can be seen that the eyebrows convey the target emotion far better in our method. Additional results can be found in the supplementary video that further demonstrate the advantages of our method.

User Study: To gauge the subjective quality of our generated avatars, we conducted a user study. We selected four subjects and generated five videos in each of the eight emotions for a total of 40 videos. The background and pose parameters for these videos

were taken from the longest video of the target emotion in the training set. We conducted our user study to compare our work to that of MEAD [37], EVP [13] and the real videos. We perform a two alternative forced choice study, pairing each of our videos with its counterpart from the alternatives. The users are shown both of the videos in a random order, together with the target emotion. We ask users to select which of the two videos is better in four categories: lip sync, visual quality, naturalness and emotional clarity. Users are able to specify if they prefer our video strongly ++, weakly +, the other video strongly -- or weakly - or if they find them equal, 0. A total of 10 users completed the study. The results are shown in Table 2. Our method strongly outperforms MEAD across all categories. Compared with EVP,

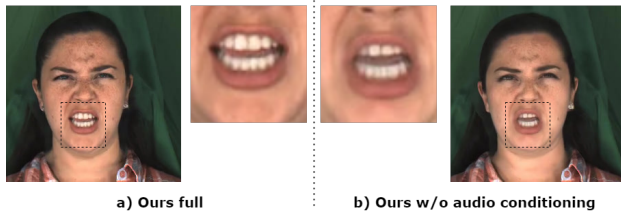


Figure 5: The addition of audio conditioning in the neural texture improves the quality of the resulting frame, particularly in the mouth interior.

our method is also preferred across all categories. However, this preference is weaker for lip sync and natural quality, but much stronger for emotional clarity. The user study shows that our work still does not reach the quality of real video, suggesting there is still room for future work.

4.4 Ablation Study

We also perform an ablation study for both the adversarial loss and the audio conditioned neural texture. The results of this comparison are shown in Table 3. The inclusion the adversarial loss improves the lip-sync at the cost of a small loss in visual quality. We note that the adversarial loss has little effect on the valence metric, but a stronger effect on the arousal. We hypothesize that this difference is due to the fluctuation in the intensity of emotion being smoothed with a pure regression loss. For the audio-conditioned neural texture, we compare our work to a static, neural texture [33]. Our method improves both the visual quality and lip-sync, with small improvements in the emotional reconstruction. As expected, the improvements of our audio conditioning are most notable in the mouth interior. This is because the audio allows the decoder to disambiguate the multiple mouth interiors that could be represented by the same underlying morphable model geometry. Figure 5 shows this improvement.

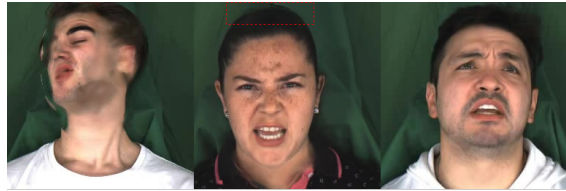


Figure 6: Failure cases of our method. **Left & middle:** When the pose of the target video is significantly different from the training data, artifacts occur. **Right:** When the tracking is inaccurate, our model produces blurry results.

5 Conclusion

We present a new method for producing audio driven avatars with control over emotion. We have used a 3D-based pipeline with the addition of an adversarial loss in the audio-to-expression generator and an audio-conditioned, resolution independent neural texture. Our method alleviates the many-to-many problem in conditioned, audio-driven video generation and surpasses state-of-the-art for lip sync and visual quality, as well as emotional reconstruction, as highlighted by our novel metric. Our comprehensive solution can be used for diverse applications.

Limitations: Our model sometimes suffers when using extreme poses (Figure 6). Furthermore, as we use reference videos to control for pose and background, length of the generated videos is limited. Future work will look to address arbitrary length video generation, potentially by considering pose generation. It is also worth investigating other models that address many-to-many generation, such as diffusion models [9, 41].

Ethical Implications: The ability to create synthetic digital humans comes with a serious potential for misuse. In particular, works such as ours could be altered in order to produce convincing misinformation. For this reason, we do not make our pipeline available to the general public. However, we are willing to share code with other researchers.

References

- [1] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in Neural Information Processing Systems*, 33:12449–12460, 2020.
- [2] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces. In *Proceedings of the 26th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '99, page 187–194, USA, 1999. ACM Press/Addison-Wesley Publishing Co.
- [3] Lele Chen, Ross K. Maddox, Zhiyao Duan, and Chenliang Xu. Hierarchical cross-modal talking face generation with dynamic pixel-wise loss. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [4] Daniel Cudeiro, Timo Bolkart, Cassidy Laidlaw, Anurag Ranjan, and Michael Black. Capture, learning, and synthesis of 3D speaking styles. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 10101–10111, 2019.
- [5] Bernhard Egger, William A. P. Smith, Ayush Tewari, Stefanie Wuhrer, Michael Zollhofer, Thabo Beeler, Florian Bernard, Timo Bolkart, Adam Kortylewski, Sami Romdhani, Christian Theobalt, Volker Blanz, and Thomas Vetter. 3d morphable face models—past, present, and future. *ACM Trans. Graph.*, 2020.
- [6] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks, 2014.
- [7] Philip-William Grassal, Malte Prinzler, Titus Leistner, Carsten Rother, Matthias Nießner, and Justus Thies. Neural head avatars from monocular rgb videos. *arXiv preprint arXiv:2112.01554*, 2021.
- [8] Yudong Guo, Keyu Chen, Sen Liang, Yongjin Liu, Hujun Bao, and Juyong Zhang. Ad-nerf: Audio driven neural radiance fields for talking head synthesis. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.
- [9] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *arXiv preprint arXiv:2006.11239*, 2020.
- [10] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [11] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros. Image-to-image translation with conditional adversarial networks, 2016.
- [12] Xinya Ji, Hang Zhou, Kaisiyuan Wang, Qianyi Wu, Wayne Wu, Feng Xu, and Xun Cao. Eamm: One-shot emotional talking face via audio-based emotion-aware motion model. In *ACM SIGGRAPH 2022 Conference Proceedings*, SIGGRAPH '22, New York, NY, USA, 2022. Association for Computing Machinery.
- [13] Xinya Ji, Hang Zhou, Kaisiyuan Wang, Wayne Wu, Chen Change Loy, Xun Cao, and Feng Xu. Audio-driven emotional video portraits. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [14] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pages 694–711. Springer, 2016.
- [15] Hyeonwoo Kim, Mohamed Elgharib, Hans-Peter Zollöfer, Michael Seidel, Thabo Beeler, Christian Richardt, and Christian Theobalt. Neural style-preserving visual dubbing. *ACM Transactions on Graphics (TOG)*, 38(6):178:1–13, 2019.
- [16] Hyeonwoo Kim, Pablo Garrido, Ayush Tewari, Weipeng Xu, Justus Thies, Matthias Nießner, Patrick Pérez, Christian Richardt, Michael Zollöfer, and Christian Theobalt. Deep video

- portraits. *ACM Transactions on Graphics (TOG)*, 37(4):163, 2018.
- [17] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- [18] Avisek Lahiri, Vivek Kwatra, Christian Frueh, John Lewis, and Chris Bregler. Lipsync3d: Data-efficient learning of personalized 3d talking faces from video using pose and lighting normalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [19] J. P. Lewis, Ken Anjyo, Taehyun Rhee, Mengjie Zhang, Fred Pighin, and Zhigang Deng. Practice and Theory of Blendshape Facial Models. In Sylvain Lefebvre and Michela Spagnuolo, editors, *Eurographics 2014 - State of the Art Reports*. The Eurographics Association, 2014.
- [20] Tianye Li, Timo Bolkart, Michael J. Black, Hao Li, and Javier Romero. Learning a model of facial shape and expression from 4D scans. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)*, 36(6):194:1–194:17, 2017.
- [21] Yuanxun Lu, Jinxiang Chai, and Xun Cao. Live speech portraits: Real-time photorealistic talking-head animation. *ACM Trans. Graph.*, 40(6), dec 2021.
- [22] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley. Least squares generative adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2794–2802, 2017.
- [23] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020.
- [24] Foivos Paraperas Papantoniou, Panagiotis P. Filntisis, Petros Maragos, and Anastasios Roussos. Neural emotion director: Speech-preserving semantic control of facial expressions in "in-the-wild" videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [25] K R Prajwal, Rudrabha Mukhopadhyay, Vinay P. Namboodiri, and C.V. 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, MM '20*, page 484–492, New York, NY, USA, 2020. Association for Computing Machinery.
- [26] Nikhila Ravi, Jeremy Reizenstein, David Novotny, Taylor Gordon, Wan-Yen Lo, Justin Johnson, and Georgia Gkioxari. Accelerating 3d deep learning with pytorch3d. *arXiv:2007.08501*, 2020.
- [27] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In Nassir Navab, Joachim Hornegger, William M. Wells, and Alejandro F. Frangi, editors, *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, pages 234–241, Cham, 2015. Springer International Publishing.
- [28] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. First order motion model for image animation. In *Conference on Neural Information Processing Systems (NeurIPS)*, December 2019.
- [29] Vincent Sitzmann, Julien N.P. Martel, Alexander W. Bergman, David B. Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. In *Proc. NeurIPS*, 2020.
- [30] Yasheng Sun, Hang Zhou, Kaisiyuan Wang, Qianyi Wu, Zhibin Hong, Jingtuo Liu, Errui Ding, Jingdong Wang, Ziwei Liu, and Koike

- Hideki. Masked lip-sync prediction by audio-visual contextual exploitation in transformers. In *SIGGRAPH Asia 2022 Conference Papers*, SA '22, New York, NY, USA, 2022. Association for Computing Machinery.
- [31] Anni Tang, Tianyu He, Xu Tan, Jun Ling, Runnan Li, Sheng Zhao, Li Song, and Jiang Bian. Memories are one-to-many mapping alleviators in talking face generation, 2022.
- [32] Justus Thies, Mohamed Elgharib, Ayush Tewari, Christian Theobalt, and Matthias Nießner. Neural voice puppetry: Audio-driven facial reenactment. *ECCV 2020*, 2020.
- [33] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: Image synthesis using neural textures. *ACM Transactions on Graphics (TOG)*, 38(4):1–12, 2019.
- [34] Justus Thies, Michael Zollhöfer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. Face2face: Real-time face capture and reenactment of rgb videos. *Commun. ACM*, 62(1):96–104, dec 2018.
- [35] Antoine Toisoul, Jean Kossaifi, Adrian Bulat, Georgios Tzimiropoulos, and Maja Pantic. Estimation of continuous valence and arousal levels from faces in naturalistic conditions. *Nature Machine Intelligence*, 2021.
- [36] Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic. End-to-end speech-driven facial animation with temporal gans. *arXiv preprint arXiv:1805.09313*, 2018.
- [37] Kaisiyuan Wang, Qianyi Wu, Linsen Song, Zhuoqian Yang, Wayne Wu, Chen Qian, Ran He, Yu Qiao, and Chen Change Loy. Mead: A large-scale audio-visual dataset for emotional talking-face generation. In *ECCV*, 2020.
- [38] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. Video-to-video synthesis. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- [39] Xin Wen, Miao Wang, Christian Richardt, Ze-Yin Chen, and Shi-Min Hu. Photorealistic audio-driven video portraits. *IEEE Transactions on Visualization and Computer Graphics*, 26(12):3457–3466, 2020.
- [40] Qiantong Xu, Alexei Baevski, and Michael Auli. Simple and effective zero-shot cross-lingual phoneme recognition. *arXiv preprint arXiv:2109.11680*, 2021.
- [41] Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Yingxia Shao, Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and applications. *arXiv preprint arXiv:2209.00796*, 2022.
- [42] Xinwei Yao, Ohad Fried, Kayvon Fatahalian, and Maneesh Agrawala. Iterative text-based editing of talking-heads using neural retargeting. *ACM Trans. Graph.*, 40(3), aug 2021.
- [43] Yufeng Zheng, Victoria Fernández Abrevaya, Marcel C. Bühler, Xu Chen, Michael J. Black, and Otmar Hilliges. I M Avatar: Implicit morphable head avatars from videos. In *Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [44] M. Zollhöfer, J. Thies, P. Garrido, D. Bradley, Thabo Beeler, P. Pérez, M. Stamminger, M. Nießner, and C. Theobalt. State of the art on monocular 3d face reconstruction, tracking, and applications. *Computer Graphics Forum*, 37:523–550, 05 2018.

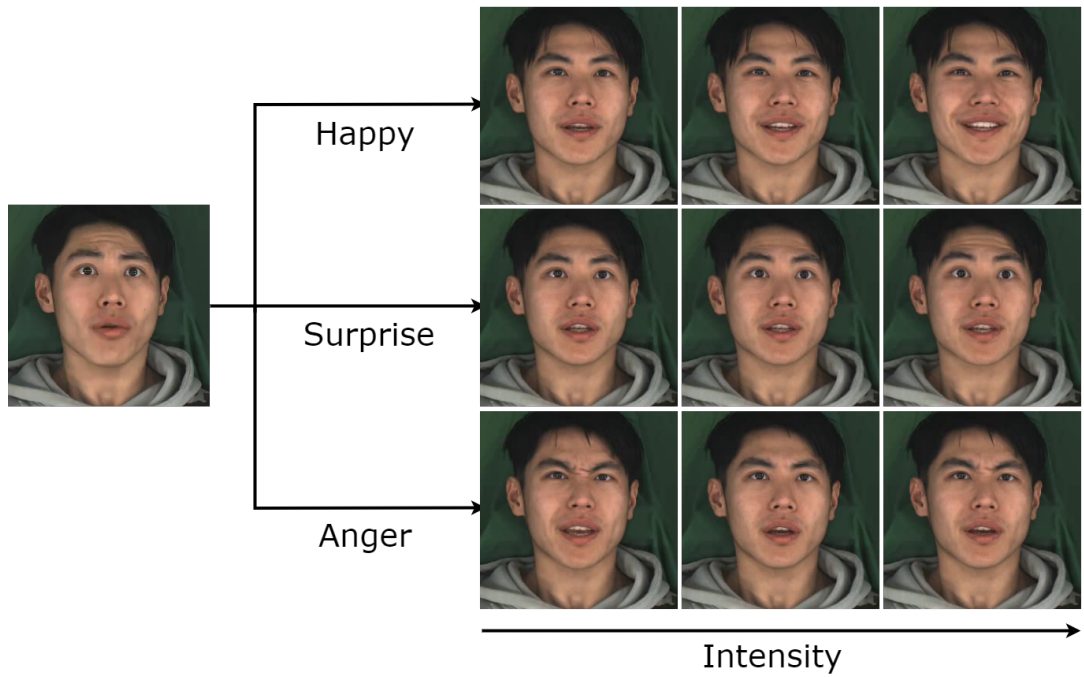
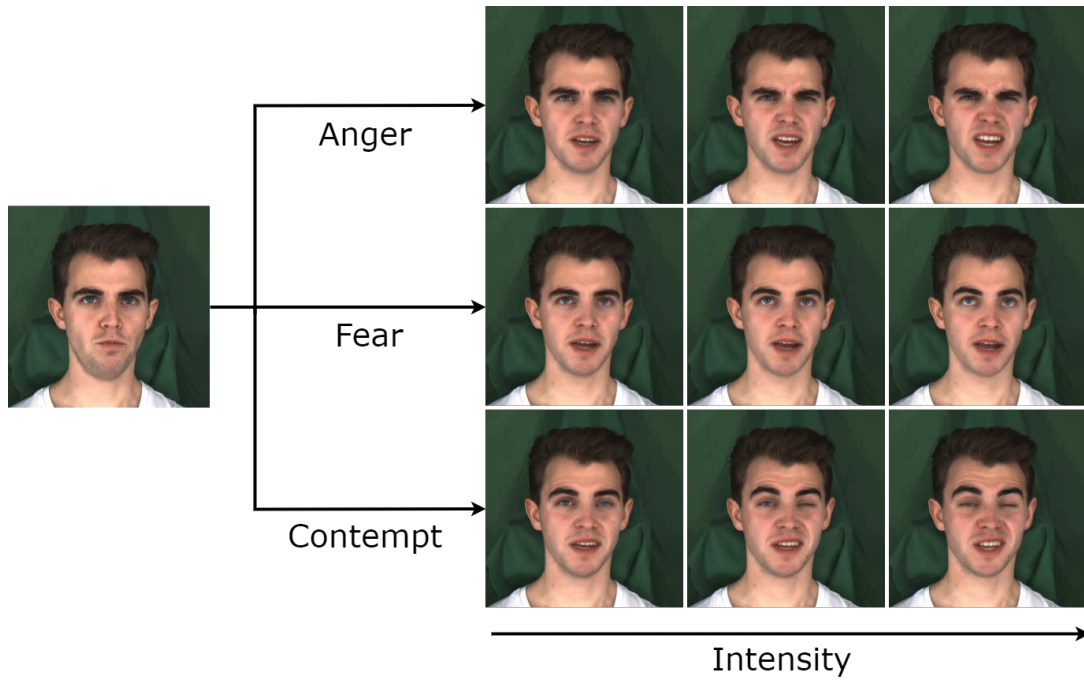


Figure 7: Our method allows for fine-grained control on multiple subjects. Here we show two subjects in with three emotions and three levels of intensity.