

Audio-Driven Emotional 3D Talking-Head Generation

Wenqing Wang¹ and Yun Fu^{1,2}

¹ Khoury College of Computer Science, Northeastern University, USA

² Department of Electrical and Computer Engineering, Northeastern University, USA

Abstract—Audio-driven video portrait synthesis is a crucial and useful technology in virtual human interaction and film-making applications. Recent advancements have focused on improving the image fidelity and lip-synchronization. However, generating accurate emotional expressions is an important aspect of realistic talking-head generation, which has remained underexplored in previous works. We present a novel system in this paper for synthesizing high-fidelity, audio-driven video portraits with accurate emotional expressions. Specifically, we utilize a variational autoencoder (VAE)-based audio-to-motion module to generate facial landmarks. These landmarks are concatenated with emotional embeddings to produce emotional landmarks through our motion-to-emotion module. These emotional landmarks are then used to render realistic emotional talking-head video using a Neural Radiance Fields (NeRF)-based emotion-to-video module. Additionally, we propose a pose sampling method that generates natural idle-state (non-speaking) videos in response to silent audio inputs. Extensive experiments demonstrate that our method obtains more accurate emotion generation with higher fidelity.

I. INTRODUCTION

Audio-driven video portraits have become an important technology across a variety of applications, including digital humans, virtual reality, and the entertainment industry. Achieving realistic talking-head synthesis requires not only accurate lip synchronization and high image fidelity but also the accurate generation of emotional expressions. Recently, many works have been proposed to synthesize audio-driven video portraits [21], [6], [25], [8], [37], [33], [28]. However, they often introduce artifacts, produce unrealistic images, or fail to capture the details of the target person. For instance, Wav2Lip [28] demonstrates high performance in lip-synchronization, but it is unable to generate the details or the overall facial movements of the target person, which results in decreased image quality. DaGAN [21] adopts a depth-conditioned generative adversarial network (GAN) for talking-head generation, but its GAN-based renderer struggles with unstable training and challenges in modeling image details. FACIAL [6] employs a GAN integrated with facial implicit attribute learning, but it also suffers from training instability and difficulties in generating intricate details. VideoReTalking [8] directly edits the video frames with the synthesized expressions to create lip-synchronized videos. However, it tends to produce unrealistic results with limited generalization capabilities.

With the adaptation of NeRF [26] in the talking-head generation, NeRF-based methods [7], [10], [30], [27], [9] have demonstrated better performance in producing high-fidelity talking-head videos with accurate rendering of image details. Furthermore, NeRF-based methods AD-NeRF [19],

GeneFace[34], and GeneFace++ [17] manage to improve the lip-synchronization accuracy with realistic image rendering. Despite these advancements, the emotional aspect of generating vivid talking-head videos has been overlooked.

More recent works have begun focusing on generating emotion-aware talking-head videos [18], [22], [13], [15], [32], [11], [31]. EAT [18] employs emotion adaption models to enable controlled emotional video generation. FlowVQTalker [32] synthesizes emotional talking heads using normalizing Flow and Vector-Quantization modeling. EAMM [22] utilizes emotional videos and pose videos to assist the generation of emotional video portraits. FG-EmoTalk [31] generates emotional talking heads using disentangled expression latent code and extracted human facial features. PC-AVS [36] uses pose control to generate talking heads. EVP [15] synthesizes emotional video portraits by disentangling speech content and emotional features.

While these works made improvements in generating emotional expressions, they struggle with several challenges in terms of generating vivid emotions. 1) **Emotion accuracy**. Despite previous methods are able to generate some forms of emotional expressions, their generated emotional expressions are often partially inaccurate or invisible. For instance, some might have corresponding brow expressions but incorrect mouth expressions [18], which can diminish the realness of the synthesized videos. 1) **Identity-preservation**. Previous methods did not take into account preserving the target identity in the process of learning emotional representation. This leads to the loss of identity and facial distortions, which can reduce the realism and fidelity of the generated videos. 3) **Idle state**. Given silent audio, other works generate extra lip movements or unnatural body motions (either have extra movement or are entirely static). This limitation can undermine their ability to produce natural idle-state videos in many applications, such as conversational agents and digital humans.

We propose a system *EmoGene* to handle these three challenges. To address the emotion accuracy challenge, we propose a motion-to-emotion module. This module is composed of a landmark deformation model (LMD), which is a neural network trained on an emotion-labeled video dataset MEAD [12] to accurately generate emotional facial landmarks from the given emotion text label and neutral facial landmarks. This enables accurate and apparent emotion generation. For the identity-preservation challenge, we utilize an emotion-to-video module that consists of NeRF models to preserve the target person’s identity and render high-fidelity videos based on emotional landmarks. To address the idle-

state challenge, we develop a pose sampling method to produce a natural idle-state video from silent audio inputs. This can make the generated talking-head video appear more natural and realistic.

The main contributions of this paper are:

- We developed a three-stage framework for generating audio-driven emotional 3D video portraits with accurate emotional expressions and preserved identity.
- We introduce a landmark deformation model for controlled emotional landmark generation.
- We propose a pose sampling method that produces more realistic idle-state videos.
- Experiments show that our EmoGene outperforms previous works in generating emotional and high-fidelity talking-head videos.

II. RELATED WORK

Audio-driven emotional video portrait generation leverages a generative model to synthesize the emotional talking-head videos given the corresponding audio. It is related to the previous works on audio-driven video portrait generation and emotion-aware video portrait generation.

Audio-driven video portrait generation. Generating audio-driven video portraits has gained lots of attention in the past few years [28], [30], [21], [6], [25], [37], [33], [8]. Cheng *et al.* [8] introduces a system to edit video expressions based on the input audio. Zhang *et al.* [6] proposes a method to synthesize talking face animation by learning the integrated phonetic, context, and identity information. By utilizing a pre-trained lip synchronization model, [28] manages to generate accurate lip-synchronization. However, it only generates the lip region, which can limit its realness and generalization ability. Zhang *et al.* [35] proposes a talking face generation framework guided by normalizing-flow, but its generation quality depends on the accurate foreground masks. Wav2NeRF [30] leverages wavelet transform and audio-visual cross-modality representations to generate talking heads. Yi *et al.* [19] developed the method to synthesize frames with the re-rendered 3D face animation, which often exhibits unrealistic expressions in its animation results. Following the introduction of NeRF [26], using NeRF to render the talking-head videos attracted more attention in the research community [19], [34], [17], due to their ability to render delicate details. Notably, Ye *et al.* proposed GeneFace++ [17], which is a remarkable audio-driven talking-head generation system that enables realistic rendering and robust lip-synchronization.

Emotion-aware video portrait generation. Building on the concepts above, emotion-aware video portrait generation aims to produce lifelike emotional expressions by incorporating emotional features into the video generation process. To achieve its goal, this task requires the generated video to not only exhibit corresponding audio-lip movements but also display vivid emotions. Previous works explored the different methods to enable controllable emotional expression generation [18], [22], [32], [11], [31], [16], [15]. Gan *et al.* developed EAT [18], an emotional talking-head generation

framework based on an audio-to-expression transformer and emotional adaption networks. EVP [15] proposes a method to synthesize emotional video portraits with disentangled emotion and content features. Ji *et al.* developed EAMM [22] framework to utilize augmented emotional source videos to generate one-shot emotional talking heads. Tan *et al.* [32] produces emotional talking heads by leveraging normalizing Flow and Vector-Quantization modeling. Sung *et al.* [11] generates talking heads capable of expressing laughter using FLAME parameters and vertices. Furthermore, Sun *et al.* [31] developed a framework that employs a disentanglement scheme to isolate emotion latent code and utilizes self-supervised learning to generate emotional talking heads. Additionally, StyleTalk [16] utilizes a speaking style to synthesize video portraits, while Li *et al.* [24] proposes a two-staged system that synthesizes video portraits with facial expressions.

III. METHOD

We present our *EmoGene* framework in this section. EmoGene consists of 3 parts: 1) audio-to-motion: a VAE model that converts the audio features into neutral facial landmarks; 2) motion-to-emotion: a landmark deformation model that transforms the neutral landmarks into emotional landmarks; 3) emotion-to-video: NeRF models that render high-fidelity emotional talking-head video (Figure 1).

A. Audio-to-Motion

We utilize a VAE model [23], [17] to generate audio-driven diverse facial landmarks, as shown in Figure 2. In this module, the VAE is conditioned on the audio features and ground truth (GT) facial landmarks and guided by the lip-synchronization discriminator to learn to generate diverse facial landmarks given an audio input.

Encoder and decoder. Our encoder and decoder are structured as convolutional neural networks. Leveraging a WaveNet-based architecture [5], the convolutional layers are gradually increased to enhance the receptive field. This allows the encoder and decoder to generate sequences with different lengths efficiently.

Training process. In the training process, we utilize Deep 3D Face Reconstruction [4] to extract the ground truth 3DMM facial landmarks from the video, and we selected 68 landmark points to generate detailed facial motion. Then, we extract the HuBERT [14] and pitch features of the given audio. Furthermore, we employ the Monte-Carlo Evidence Lower Bound loss [29] to increase the efficiency of the training process. To guide the VAE training, we evaluate the audio-lip synchronization of the generated landmarks using a SyncNet [3] based pre-trained lip-synchronization discriminator. Hence, the VAE training loss is:

$$\mathcal{L}_{\text{VAE}} = \mathbb{E} [\|l - \hat{l}\|_2^2 + \text{KL}(z | \hat{z}) + \mathcal{L}_{\text{Sync}}(a, \hat{l})]. \quad (1)$$

In this setup, l is the ground truth facial landmark and a represents the input audio. The latent encoding of l is denoted by $\hat{z} = \text{Encoder}(l, a)$, and the latent code z is sampled with

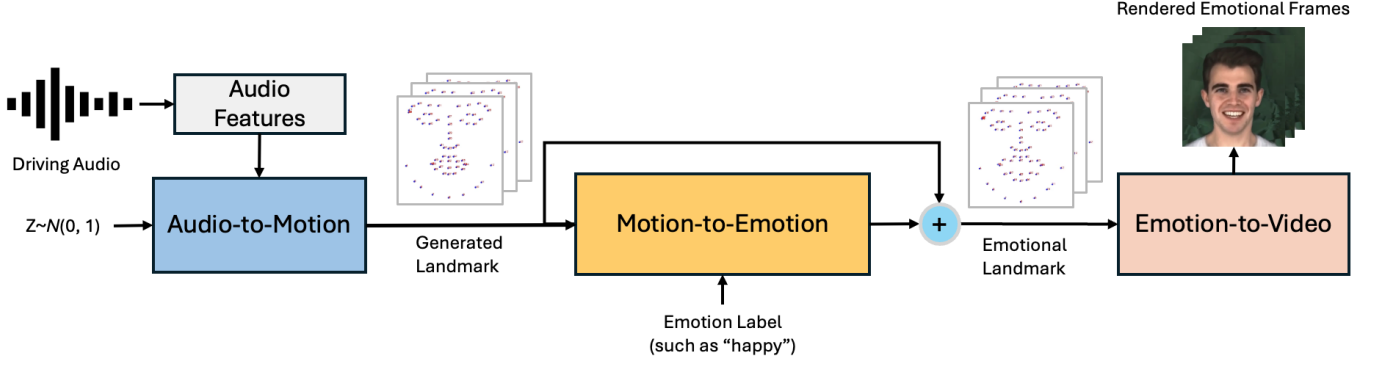


Fig. 1: **The inference pipeline of EmoGene.** 1) In the first stage, the Audio-to-Motion module converts the audio features into neutral facial landmarks. 2) In the second stage, the Motion-to-Emotion module takes the generated neutral landmarks and an emotion label to generate emotional landmarks. 3) In the third stage, the Emotion-to-Video module renders the emotional talking-head video conditioning on the emotional landmarks.

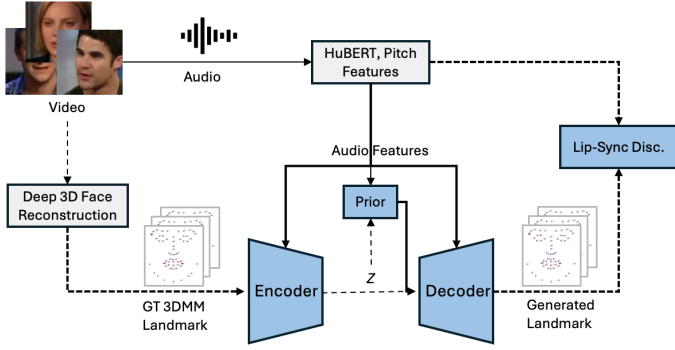


Fig. 2: **The overview of audio-to-motion module.** Dash arrows indicate that the process is only conducted during training.

$z \sim N(0, 1)$. The generated landmark points are represented as $\hat{l} = \text{Decoder}(\hat{z}, a)$. KL is KL divergence, and $\mathcal{L}_{\text{Sync}}$ is the lip-synchronization loss evaluated by the pre-trained lip-synchronization discriminator. During the inference, only the decoder is required to generate audio-driven facial landmarks.

B. Motion-to-Emotion

Without conditioning on emotion, the audio-to-motion module can only generate neutral landmarks. To enable the generation of emotionally expressive facial landmarks, we develop the motion-to-emotion module (Figure 3). This module transforms neutral landmarks into emotional landmarks by employing a landmark deformation model, which generates the emotional landmark deformation differences for computing the emotional facial landmarks.

Landmark deformation model. To deform the neutral landmarks to obtain emotional expressions, we utilize a neural network with 3 fully connected (FC) layers which takes the neutral landmarks and emotional embedding as inputs to generate the emotional landmark deformation differences. The landmark deformation model is defined as:

$$LD(x) = FC(ReLU(FC(ReLU(FC(x))))), \quad (2)$$

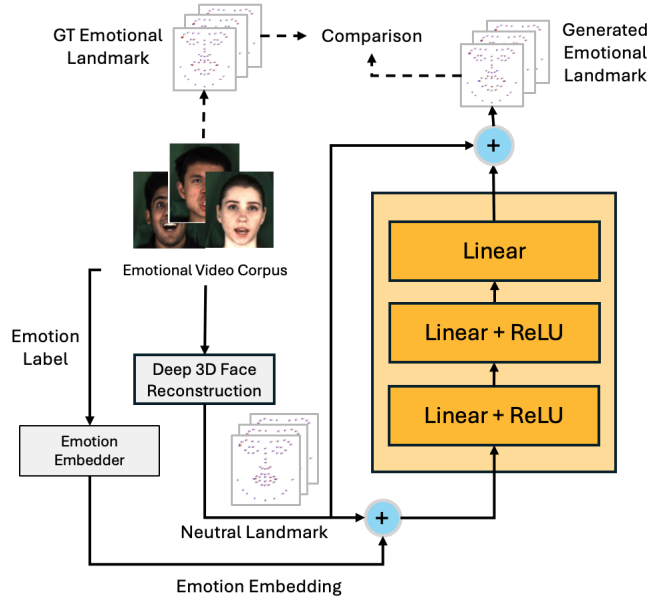


Fig. 3: **The overview of motion-to-emotion module.** Dash arrows indicate that the process is only conducted during training.

where DL represents the landmark deformation model and x denotes the concatenated embedding of the neutral landmarks and emotional embedding. To generate the corresponding emotional landmarks, we concatenate the neutral landmarks with the emotional landmark deformation differences as:

$$\hat{l}_E = \hat{l} \oplus \Delta \hat{l}, \quad (3)$$

where \hat{l}_E is the generated emotional facial landmarks and $\Delta \hat{l}$ is emotional landmark deformation differences from the landmark deformation model.

Training process. During the training process, we extract the ground truth facial landmarks from the video and utilize an emotion embedder to construct the emotion embedding from the corresponding emotion label. We concatenate the

neutral landmarks with the emotional embedding, the result of which is then fed into the landmark deformation model to generate the emotional landmark deformation differences. These differences are concatenated with the neutral landmarks to synthesize the emotional landmarks. These generated emotional landmarks are compared with the ground truth emotional landmarks for the corresponding emotion label. The training loss for the landmark deformation model is defined as:

$$\mathcal{L}_{LD} = \mathbb{E} [\|l_E - \hat{l}_E\|_2^2], \quad (4)$$

where l_E is the ground truth emotional landmarks and \hat{l}_E is the generated emotional landmarks.

C. Emotion-to-Video

After we obtain the emotional landmarks from the motion-to-emotion module, we then utilize the emotion-to-video module, which consists of NeRF-based models, to synthesize emotional talking-head frames conditioned on the generated emotional landmarks (Figure 4).

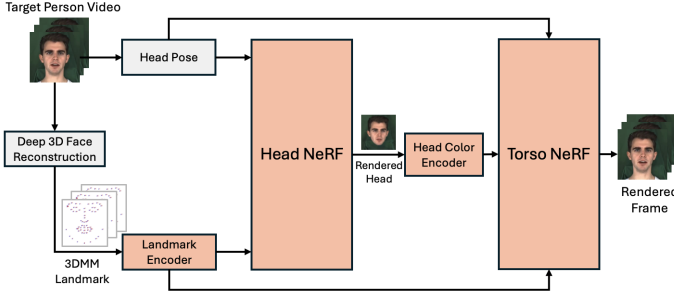


Fig. 4: The overview of emotion-to-video module.

Emotion landmark-conditioned Head-NeRF. To render high-quality talking-head videos, we utilize a NeRF-based model [9], [17] conditioned on the emotional landmarks to dynamically represent a talking head. Different from the original NeRF [26], this model incorporates not only 3D location x and the viewing direction d but also 3D emotional landmarks \hat{l}_E . Based on these parameters, the model function F is defined as:

$$F : (x, d, \hat{l}_E) \rightarrow (c, \sigma), \quad (5)$$

where c represents the predicted RGB color and σ denotes the predicted volume density (for the emitted ray $r(t) = o + t \cdot d$ with the camera origin o) in the 3D neural radiance field.

To render each pixel of the image frame, we aggregate the color c along the ray by following the differentiable volume rendering equation [26]:

$$C(r, \hat{l}_E) = \int_{t_n}^{t_f} \sigma(r(t), \hat{l}_E) \cdot c(r(t), \hat{l}_E, d) \cdot T(t) dt, \quad (6)$$

where C is the pixel color for the ray r . t_n and t_f denote the near and far bounds of the ray. $T(t)$ is the accumulated transmittance along the ray with these bounds:

$$T(t) = \exp \left(- \int_{t_n}^t \sigma(r(s), \hat{l}_E) ds \right), \quad (7)$$

where s is an intermediate point along the path of the ray from the near bound t_n to the current point t .

Torso-NeRF. To enable robust cooperation between the head and torso, we condition the torso-NeRF on the head-NeRF's output color. This enables the torso NeRF to receive more information about the rendered head results, which strengthens the collaborated rendering process of the head and torso and enhances the realness of the generated videos.

Following this, the torsos's function F_{torso} is thus:

$$F_{torso} : (x, C_{head}; d_0, P, \hat{l}_E) \rightarrow (c, \sigma). \quad (8)$$

Here, d_0 denotes the canonical space's view direction, and the head pose $P \in \mathbb{R}^{3 \times 4}$ with a translation vector and rotation matrix.

Training process. To train our NeRF models, we extract the facial landmarks and images from the input video frames. We then utilize the landmark-image pairs to train the NeRF models. To decrease the L_2 reconstruction error between the generated images and the ground truth images, the training loss for the NeRF models is:

$$\mathcal{L}_{NeRF} = \mathbb{E} [\|C(r, \hat{l}_E) - C_{GT}\|_2^2]. \quad (9)$$

D. Pose Sampling Method for Idle State

We propose a pose sampling algorithm to generate natural body and lip movements for the idle-state video. As shown in Figure 6, we first identify the starting idle poses in the original pose tensor, then we replicate these idle poses to construct a collection of idle pose segments with random lengths and fixed gaps (number of non-idle poses) between them (using Algorithm 1). To create a new pose tensor with the idle state, we insert the idle segments into the original tensor (as shown in Figure 5).

Determining the minimum and maximum idle lengths. To generate idle segments, we developed a method to determine the minimum and maximum lengths of contiguous idle segments. Given a pose tensor $X = \{x_1, \dots, x_i, x_{i+1}, \dots, x_{n-1}\}$, this method computes the cosine similarities between each contiguous pair of poses at indices i and $i+1$ as:

$$\text{CosineSimilarity}_i = \frac{\mathbf{x}_i \cdot \mathbf{x}_{i+1}}{\|\mathbf{x}_i\| \|\mathbf{x}_{i+1}\|}. \quad (10)$$

A pair of consecutive poses is considered to belong to an idle segment if their cosine similarity is greater or equal to 1, since a cosine similarity of 1 indicates that two poses are aligned and have maximum similarity. Then, the method calculates the lengths of the idle segments to obtain the minimum and maximum idle segment lengths.

Generating idle segments with random lengths and fixed gaps. As shown in Algorithm 1, this algorithm generates the $(start, end)$ pairs of the idle segments, where $start$ and end are the starting and ending indices of the idle segments in the pose tensor. To synthesize a natural idle-state

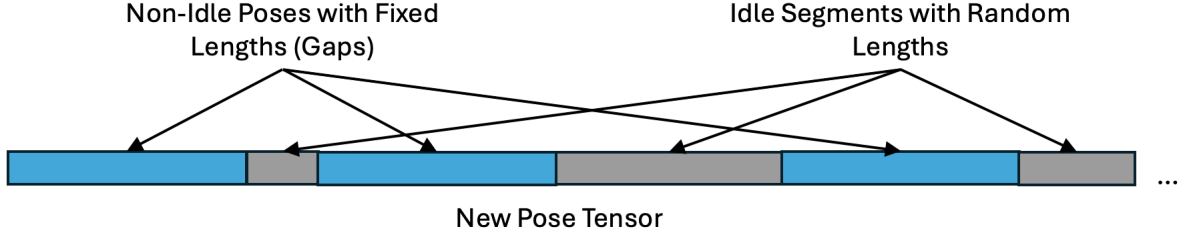


Fig. 5: **New pose tensor reconstruction.** To construct the new pose tensor, we insert the idle segments after their corresponding non-idle pose tensors.

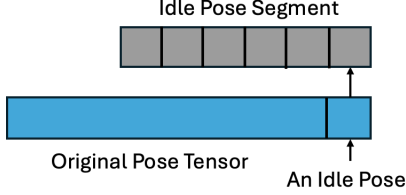


Fig. 6: **Idle pose segment reconstruction.** We construct each idle segment by identifying and replicating an idle pose of the original pose tensor.

motion, this algorithm iterates through the pose tensor and randomly determines the length of each idle segment within the constraints of the minimum and maximum idle segment lengths. This random sampling method creates variety in the idle motion, which helps to generate natural idle-state videos. After determining each idle segment, the algorithm keeps a fixed gap before the start of the next idle segment. This algorithm generates idle motion periods that mimic natural waiting or resting states, which can increase the realness of the idle-state videos.

Inserting idle segments into pose tensor. To enable the idle state, we create a new pose tensor by inserting the idle segments after their corresponding non-idle pose tensors (Figure 5). This produces natural idle-state motions and smooth motion-stillness transitions for the given silent audio, which makes the rendered videos more realistic.

Algorithm 1 Generate idle segments with random lengths and fixed gaps.

```

1: function GENERATESEGMENTS( $n, \text{min\_len}, \text{max\_len}, \text{fixed\_gap}$ )
2:   segments  $\leftarrow$  []
3:   current_position  $\leftarrow$  0
4:   while current_position  $< n$  do
5:      $l \leftarrow$  random between min_len and min(max_len,  $n - \text{current\_position}$ )
6:     start  $\leftarrow$  current_position, end  $\leftarrow$  start +  $l - 1$ 
7:     if end + fixed_gap  $\geq n$  then break
8:     end if
9:     segments.append((start, end))
10:    current_position  $\leftarrow$  end + 1 + fixed_gap
11:   end while
12:   return segments
13: end function

```

IV. EXPERIMENTS

A. Experimental Settings

Datasets. We trained the audio-to-motion module on VoxCeleb2 [2], which contains 6,112 identities with over 1 million utterances, to learn a generalized audio-to-motion mapping. For learning the landmark deformation model, we utilized the MEAD dataset [12], which consists of labeled emotional talking videos of 60 identities with 8 emotions (surprise, sad, neutral, happy, hear, contempt, disgust, angry) for each utterance. Based on identity, we split the dataset into training and testing sets. To train the NeRF models, we adopted the MEAD videos and collected videos, and each training video is around 3 to 6 minutes with 512x512 resolution and 25 FPS.

Comparisons. The 3 notable works we compare with EmoGene are: 1) EAT [18], which leverages a pre-trained transformer and emotion adaptation models to generate emotional expressions; 2) Wav2Lip [28], which adopts a pre-trained sync-expert to synchronize lip movements; 3) GeneFace++ [17], which uses auxiliary audio features and an efficient NeRF to render talking-head videos.

Implementation details. EmoGene is trained on 1 NVIDIA RTX A6000. The landmark deformation model and VAE take around 40K to converge (about 14 hours). The NeRF models are trained for 400K iterations in total (about 12 hours).

B. Quantitative Evaluation

Evaluation metrics. We adopt the SSIM, PSNR, and FID [20] to measure the image quality and fidelity of the generated videos. We employ the landmark distance (LMD) [1] to measure audio-lip synchronization and emotional expression accuracy.

Evaluation results. From the results (Table I), we observe that EmoGene outperforms other methods in achieving high image quality with the best SSIM and PSNR scores. Furthermore, our method archives the second-best LMD score for generating emotional expressions. However, Wav2Lip only generates the lip region but not the entire face. This narrower generation scope of Wav2Lip might enable it to have most of its generated image landmarks overlapped with the ground truth image landmarks, which can lead to an inflated LMD score.

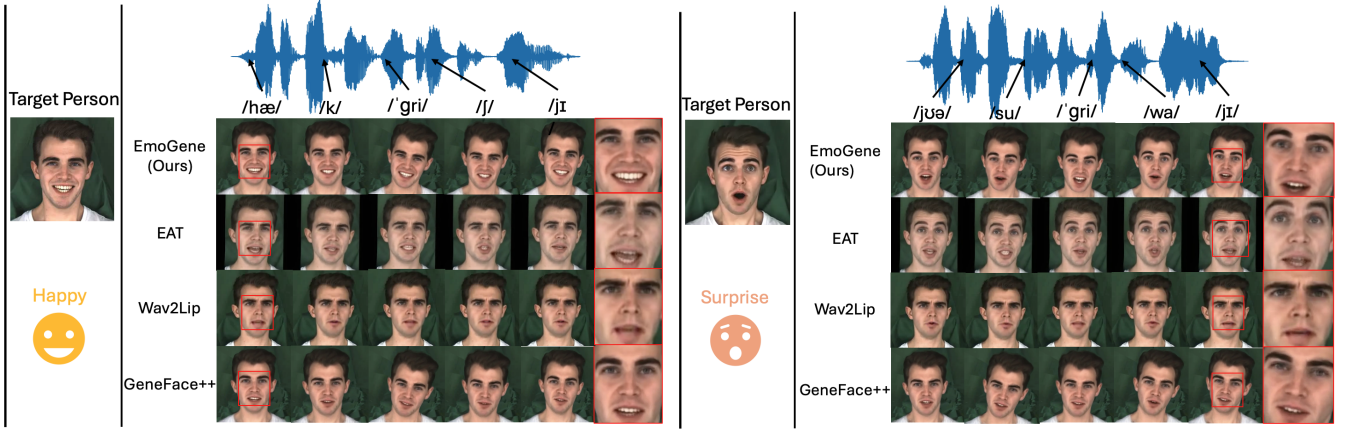


Fig. 7: **Qualitative comparison of generated keyframe results.** Results of *happy* (left) and *surprise* (right). The corresponding phonetic symbols are shown on the top row.

TABLE I: Quantitative evaluation results. Best results are in **red**. Second best results are in **blue**.

| Method/Score | SSIM \uparrow | PSNR \uparrow | LMD \downarrow | FID \downarrow |
|-----------------|-----------------|-----------------|------------------|------------------|
| EAT [18] | 0.479 | 16.940 | 12.371 | 142.655 |
| GeneFace++ [17] | 0.729 | 19.687 | 11.815 | 128.659 |
| Wav2Lip [28] | 0.623 | 19.428 | 10.429 | 96.935 |
| EmoGene (Ours) | 0.730 | 19.698 | 11.687 | 128.743 |
| Ground Truth | 1.000 | ∞ | 0.000 | 0.000 |

C. Qualitative Evaluation

For qualitative evaluation, we show the keyframes of two emotion-specific clips in Figure 7. We observe that although Wav2Lip and GeneFace++ show good lip-synchronization, they are not able to generate emotional expression. While EAT is able to synthesize regional emotional expressions, it struggles to generate accurate emotional expressions across all facial regions, including brows, eyes, and mouth, which can lead to unnatural expressions. Furthermore, the generated images of EAT contain distortions of facial features, which results in identity loss and reduced image fidelity. For instance, during its generation of emotional expressions, its generated facial structures exhibit visible differences in overall facial structures compared to the target person. In comparison, EmoGene shows accurate facial emotional expressions and lip-synchronization with high identity-preservation.

In addition, we show the effectiveness of our pose sampling method for generating natural idle-state videos. We present the keyframes of the video clip generated by EmoGene (with pose sampling) and Wav2Lip (without pose sampling) in Figure 8. We observe that given the silent audio, the pose sampling algorithm allows EmoGene to have more stable and natural body and lip movements. In contrast, Wav2Lip exhibits larger head movements and extra lip motions, despite having silent audio as input.

User study. We have 20 evaluators to evaluate the generated videos. For each of the 4 methods, we produce 3 video clips across 8 emotions, resulting in a total of 96 videos. We employ the Mean Opinion Score (MOS) for the evalua-

tions, with ratings spanning from 1 (Bad) to 5 (Excellent). Evaluators are instructed to rate each video according to 4 criteria: 1) emotional accuracy; 2) lip synchronization; 3) video realness; and 4) video quality.

We observe from the results (Table II) that: 1) EmoGene outperforms the other methods in emotional accuracy, video realness, and video quality. 2) EmoGene has lower perceived lip synchronization accuracy compared to other methods. This could be due to the emotional landmark deformation process of EmoGene, which might potentially affect its lip-synchronization accuracy when generating emotional landmarks.

TABLE II: User study results. Best results are in **bold**.

| Criteria/Method | EmoGene (Ours) | EAT | GeneFace++ | Wav2Lip |
|-----------------|----------------|-------|------------|--------------|
| Emotion Acc. | 2.513 | 2.470 | 2.435 | 2.214 |
| Lip Sync. | 2.638 | 2.663 | 2.863 | 3.081 |
| Video Realness | 2.440 | 2.404 | 2.423 | 2.018 |
| Video Quality | 2.438 | 2.366 | 2.436 | 2.131 |

D. Ablation Study

We conduct an ablation study in this section to demonstrate the essential role of the landmark deformation model in EmoGene.

Landmark deformation model. We test the setting of w/o landmark deformation model. In this setting, we remove the emotional landmark deformation from the neutral landmarks to evaluate how the absence of emotional deformation affects the performance of our framework. The results of this setting are shown in Table III. Due to the removal of emotional deformations, the NeRF is solely conditioned on the neutral landmarks, which do not exhibit any emotional cues. Thus, the conditioning on neutral landmarks restricts the NeRF from generating emotional expressions. Consequently, we observe a notable decrease in PSNR, LMD, and FID scores, which shows the critical role of the landmark deformation model in generating emotional expressions in the framework.

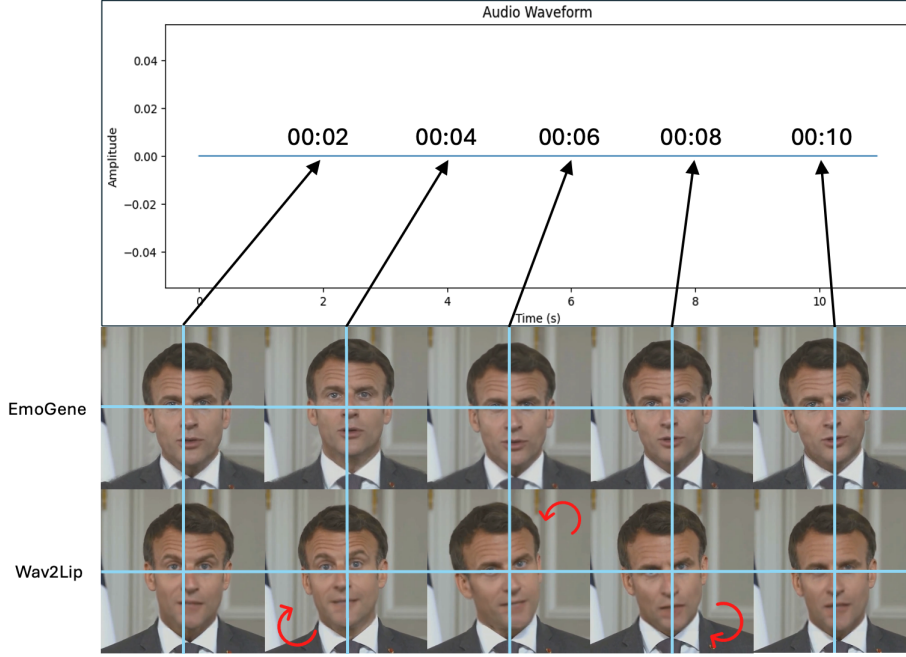


Fig. 8: **Qualitative comparison of pose sampling results.** Given silent driving audio, the results of EmoGene (with pose sampling) and Wav2Lip (without pose sampling) are on the top and bottom rows. The red rotation arrows indicate the head movements.

TABLE III: Ablation study results. Best results are in **bold**.

| Setting/Score | SSIM \uparrow | PSNR \uparrow | LMD \downarrow | FID \downarrow |
|---------------|-----------------|-----------------|------------------|------------------|
| EmoGene | 0.730 | 19.698 | 11.687 | 128.743 |
| w/o LDM | 0.730 | 19.691 | 11.894 | 129.403 |

ACKNOWLEDGEMENT

This project was supported by Bithuman, Inc.. The work was partially done when Wenqing Wang was an intern at Bithuman, Inc..

V. CONCLUSION

We present EmoGene to synthesize accurate emotional video portraits. A landmark deformation model is proposed to synthesize the robust emotional landmarks to condition the NeRF model to render vivid emotional expressions. Furthermore, a pose sampling method is introduced to generate natural idle-state videos given silent audio. Extensive experiments show that our method produces more accurate emotional expressions with preserved identity and natural movements.

Limitations. 1) The performance of our method is limited by the emotional training data, in terms of diversity and size. 2) The landmark deformation data can potentially impact the lip-synchronization accuracy of the generated videos.

Future directions. For future works, one direction is to expand the emotional training dataset with a broader range of emotional expressions. To improve lip-synchronization while generating a vivid emotion generation, future works can also consider refining the emotion transformation process by integrating other methods, such as expressive GAN-based models, and conducting extensive user studies to obtain more insights into the user-perceived emotional generation quality and image fidelity of the emotional video portraits.

REFERENCES

- [1] L. Chen, Z. Li, R. Maddox, Z. Duan, and C. Xu. Lip movements generation at a glance. *Proceedings of the European Conference on Computer Vision*, pages 520–535, 2018.
- [2] J. Chuang, A. Nagrani, and A. Zisserman. Voxceleb2: Deep speaker recognition. *INTERSPEECH*, 2018.
- [3] J. Chung and A. Zisserman. Out of time: automated lip sync in the wild. in workshop on multi-view lip-reading. *Asian Conference on Computer Vision*, pages 251–263, 2017.
- [4] Y. Deng, J. Yang, S. Xu, D. Chen, Y. Jia, and X. Tong. Accurate 3d face reconstruction with weakly-supervised learning: From single image to image set. *IEEE Computer Vision and Pattern Recognition Workshops*, 2019.
- [5] A. O. et al. Wavenet: A generative model for raw audio. *arXiv:1609.03499*, 2016.
- [6] C. Z. et al. Facial: Synthesizing dynamic talking face with implicit attribute learning. *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3867–3876, 2021.
- [7] D. L. et al. Ae-nerf: Audio enhanced neural radiance field for few shot talking head synthesis. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(4):3037–3045, 2024.
- [8] K. C. et al. Videoretalking: Audio-based lip synchronization for talking head video editing in the wild. *SIGGRAPH-ASIA*, (30):1–9, 2022.
- [9] K. P. et al. Hypernerf: a higher-dimensional representation for topologically varying neural radiance fields. *ACM Transactions on Graphics*, 40(238):1–12, 2021.
- [10] K. P. et al. Nerfies: Deformable neural radiance fields. *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5865–5874, 2021.
- [11] K. S. et al. Laughtalk: Expressive 3d talking head generation with laughter. *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 6404–6413, 2024.
- [12] K. W. et al. Mead: A large-scale audio-visual dataset for emotional talking-face generation. *European Conference on Computer Vision*, 12366:700–717, 2020.
- [13] N. D. et al. Emoportraits: Emotion-enhanced multimodal one-shot head avatars. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8498–8507, 2024.
- [14] W. H. et al. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. *arXiv:1312.6114*, 29:3451–3460, 2021.
- [15] X. J. et al. Audio-driven emotional video portraits. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14080–14089, 2021.
- [16] Y. M. et al. Speech driven talking face generation from a single image and an emotion condition. *IEEE Transactions on Multimedia*, pages 3480–3490, 2021.
- [17] Z. Y. et al. Geneface++: Generalized and stable real-time audio-driven 3d talking face generation. *arXiv preprint arXiv:2305.00787*, 2023.
- [18] Y. Gan, Z. Yang, X. Yue, L. Sun, and Y. Yang. Efficient emotional adaptation for audio-driven talking-head generation. *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22634–22645, 2023.
- [19] Y. Guo, K. Chen, S. Liang, Y. Liu, H. Bao, and J. Zhang. Ad-nerf: Audio driven neural radiance fields for talking head synthesis. *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5784–5794, 2020.
- [20] M. Heusel, H. Ramsauer, T. Unterthiner, and B. Nessler. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Neural Information Processing Systems*, 2017.
- [21] F. Hong, L. Zhang, L. Shen, and D. Xu. Depth-aware generative adversarial network for talking head video generation. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3397–3406, 2022.
- [22] X. Ji, H. Zhou, K. Wang, Q. Wu, W. Wu, F. Xu, and X. Cao. Eamm: One-shot emotional talking face via audio-based emotion-aware motion model. *ACM SIGGRAPH 2022 Conference Proceedings*, (61):1–10, 2022.
- [23] D. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv:1312.6114*, 2013.
- [24] L. Li, S. Wang, Z. Zhang, Y. Ding, Y. Zheng, X. Yu, and C. Fan. Write-a-speaker: Text-based emotional and rhythmic talking-head generation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(1):1911–1920, 2021.
- [25] Y. Lu, J. Chai, and X. Cao. Live speech portraits: Real-time photorealistic talking-head animation. *ACM Transactions on Graphics*, 40(220):1–17, 2021.
- [26] B. Mildenhall, P. Srinivasan, M. Tancik, J. Barron, R. Ramamoorthi, and R. Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 405–421, 2020.
- [27] Z. Peng. SyncTalk: The devil is in the synchronization for talking head synthesis. *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 666–676, 2024.
- [28] K. Prajwal, R. Mukhopadhyay, V. Nambodiri, and C. Jawahar. A lip sync expert is all you need for speech to lip generation in the wild. *Proceedings of the 28th ACM International Conference on Multimedia*, pages 484–492, 2020.
- [29] X. Ran, M. Xu, L. Mei, Q. Xu, and Q. Liu. Detecting out-of-distribution samples via variational auto-encoder with reliable uncertainty estimation. *Neural Networks*, 145:199–208, 2022.
- [30] A. Shin, J. Lee, J. Hwang, Y. Kim, and G. Park. Wav2nerf: Audio-driven realistic talking head generation via wavelet-based nerf. *Image and Vision Computing*, 148, 2024.
- [31] Z. Sun, Y. Xuan, F. Liu, and Y. Xiang. Fg-emotalk: Talking head video generation with fine-grained controllable facial expressions. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(5):5043–5051, 2024.
- [32] S. Tan, B. Ji, and Y. Pan. Flowvqtalker: High-quality emotional talking face generation through normalizing flow and quantization. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26317–26327, 2024.
- [33] S. Wang, L. Li, Y. Ding, and X. Yu. One-shot talking face generation from single-speaker audio-visual correlation learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(3):2531–2539, 2020.
- [34] Z. Ye, Z. Jiang, Y. Ren, J. Liu, J. He, and Z. Zhao. Geneface: Generalized and high-fidelity audio-driven 3d talking face synthesis. *arXiv preprint arXiv:2301.13430*, 2023.
- [35] Z. Zhang, L. Li, Y. Ding, and C. Fan. Flow-guided one-shot talking face generation with a high-resolution audio-visual dataset. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3661–3670, 2021.
- [36] H. Zhou, Y. Sun, W. Wu, C. Loy, X. Wang, and Z. Liu. Pose-controllable talking face generation by implicitly modularized audio-visual representation. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4174–4184, 2021.
- [37] Y. Zhou, X. Han, E. Shechtman, J. Echevarria, E. Kalogerakis, and D. Li. Makeittalk: Speaker-aware talking-head animation. *ACM Transactions on Graphics*, (221):1–15, 2020.