# A Keypoint Based Enhancement Method for Audio Driven Free View Talking Head Synthesis

Yichen Han, Ya Li, Yingming Gao, Jinlong Xue
*School of Artificial Intelligence*
*Beijing University of Posts and Telecommunications*
Beijing, China
adelacvgaoiro@bupt.edu.cn, yli01@bupt.edu.cn,
yingming.gao@outlook.com, jinlong_xue@bupt.edu.cn

Songpo Wang, Lei Yang
*DeepScience Tech Ltd.*
Beijing, China
wangsongpo@deepscience.cn, yanglei@deepscience.cn

*Abstract*—Audio driven talking head synthesis is a challenging task that attracts increasing attention in recent years. Although existing methods based on 2D landmarks or 3D face models can synthesize accurate lip synchronization and rhythmic head pose for arbitrary identity, they still have limitations, such as the cut feeling in the mouth mapping and the lack of skin highlights. The morphed region is blurry compared to the surrounding face. A Keypoint Based Enhancement (KPBE) method is proposed for audio driven free view talking head synthesis to improve the naturalness of the generated video. Firstly, existing methods were used as the backend to synthesize intermediate results. Then we used keypoint decomposition to extract video synthesis controlling parameters from the backend output and the source image. After that, the controlling parameters were composited to the source keypoints and the driving keypoints. A motion field based method was used to generate the final image from the keypoint representation. With keypoint representation, we overcame the cut feeling in the mouth mapping and the lack of skin highlights. Experiments show that our proposed enhancement method improved the quality of talking-head videos in terms of mean opinion score.

*Index Terms*—talking head generation, speech driven animation

## I. INTRODUCTION

In many applications, such as virtual reality, digital body, video conferencing, and visual dubbing, one-shot audio-driven talking head synthesis is an important component. Early research relied on motion capture by art experts, and could only be used in film and games, which was labor-intensive and time-consuming [5] [23]. In recent years, significant progress has been made in this area, and a number of deep learning methods have been proposed [17] [4] [26] [11] [18] [8] [6] [12] in order to learn the warping from audio to expression. For example, Wav2lip [12] uses a end-to-end framework, and synthesizes lower half of the face. Many methods use 2D facial landmark [2] [16] [4] or 3D head model [1] [18] [15] [3] [21] [13] [8] as a transit medium. Because 2D facial landmark does not contain information about depth perception, most methods using 2D facial landmark [26] [16] do not support viewpoint editing.

In order to achieve more flexible manipulation, neural rendering methods are proposed. They are identity independent, and can change viewpoint in latent space. Neural radiation

fields (NeRF) [6] is proposed to avoid additional intermediate representations. However, it is still a difficult task to control the head pose and expression at the same time. To overcome this limitation, keypoint-based methods are proposed [19] [14]. They can render hair and sunglasses precisely that are not possible for 3D-based methods, and are able to change the viewpoint that are not possible for 2D-based methods.

To overcome the cut feeling in the mouth mapping and the lack of skin highlights, we propose a **Keypoint Based Enhancement (KPBE)** method for audio-driven free view talking head synthesis. Our approach consists of a backend and a frontend. The backend is model-free. Using audio and source images as input, the existing backend methods synthesize intermediate results. The frontend contains five modules: canonical keypoint estimator, appearance feature estimator, head pose and expression estimator, motion field estimator, and generator. Canonical keypoint estimator is for estimating customized keypoints for the different images. Appearance feature estimator is for extracting the appearance features such as skin and color of eyes. Head pose and expression estimator is for extracting head pose and expression from backend output. Specifically, head pose is determined by a rotation matrix and a translation vector. Expression is parameterized by vectors with the same number of the canonical keypoints. Motion field estimator is for compositing the motion vectors in 3D space. Motion vector is obtained via pair of keypoints extracted from source image and driving image. Generator is for generating final video from appearance feature and composited motion field. Using this workflow we can get appearance features that maintain the skin highlights, and the generator avoids the cut feeling in the mouth mapping. In addition, we can change the viewpoint by user-defined head pose matrix.

The contributions of our work are two folds:

- A keypoint-based model-free enhancement method for audio-driven talking head synthesis is proposed, which can composite the head pose and the lip motion naturally.
- Head pose can be manipulated by user-defined rotation matrix and tranlation vector.

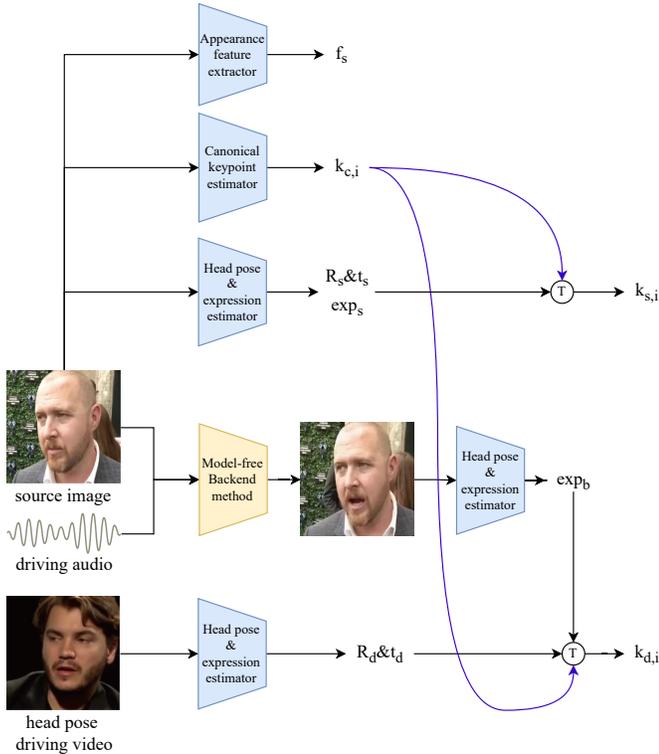The rest of the paper is organized as follows. Section II

Fig. 1. Keypoint decomposition. The appearance features, 3D canonical key points, head pose, and expression are extracted from the source image. Meanwhile, expression is extracted from backend output, and head pose is extracted from the head pose driving video. Applying the corresponding head pose and expression to the canonical keypoint, we get the source keypoint and the driving keypoint.

introduces the related works about talking head synthesis. Section III represents the architecture of our method and details of each part. Experiments and results are shown in Section IV. Section V concludes the paper.

## II. RELATED WORKS

**Audio-driven talking head synthesis**. Driving a talking head with audio is the task of synchronizing image frames of a video with arbitrary audio. There are two lines of work: one is 3D-based and the other is 2D-based. The early 3D-based methods attempt to build the relationship between audio features and lip motions by hand [23] [1], so they need an expert of the field. A well-known person-specific 3D-based work that does not rely on expert is [16]. The authors generated talking-head videos of president Obama, and focused on synthesizing natural head pose and accurate lip motion. However, the method needs a large number of videos of a specific person. 2D-based methods [7] [12] can achieve identity independent expression synthesis by replacing part of the face. Recent 2D-based methods can also synthesize head pose using facial landmarks [24] [22]. And some methods achieve real-time speed synthesis [20]. Latest 3D-based methods can perform identity-independent synthesis [3]. However, realistic hair cannot be generated , because it is difficult for a network to learn the model of these high poly models. All

these methods have only a limited capabilities for manipulating head pose and viewpoint.

**Neural rendering based talking head synthesis**. Sitzmann *et al.* [25] used neural networks to represent the 3D shape or appearance of scenes, and they sampled point set in space to represent the appearance of an object. A method was originally presented by Siarohin *et al.* [14] to get the warping between sparse keypoints and the motion fields from them. Wang *et al.* [19] used 3D-based keypoint warping to overcome the shortcomings of the previous method. We use the similar idea for keypoint decomposition in our method.

Recently, Neural Radience Fields (NeRF) [9] has gained many achievements in neural rendering tasks. They transformed the 3D appearance features to ray sampling results of volume. AD-nerf [6] applied this idea to talking head synthesis. Neural rendering methods achieve free view head pose synthesis, and have the ability to learn the depth information from the 2D image. We use a similar framework for extracting 3D appearance feature.

## III. METHOD

We proposed a **Keypoint Based Enhancement (KPBE)** method for audio-driven free view talking head synthesis. The synthesis framework consists of two steps: keypoint decomposition and generating from keypoint representation, as illustrated in Fig. 1 and Fig. 2, respectively. Specifically, the keypoint decomposition contains two parts. One is the backend which synthesizes the intermediate video from the audio. The other part belongs to the keypoint-based frontend which enhances the result of the backend output. Using keypoint decomposition, it extracts source keypoints and driving keypoints from the inputs. In the step of generating from keypoint representation part, we use the keypoints obtained by above part as input, and use the motion field based generator to synthesize the final image.

### A. Audio driven backend

In this paper, we use PC-AVS and Wav2Lip as our backend methods (indicated by the yellow block in Fig. 1). Driving audio and the source image are used as the input of the backend model. The output of the backend is the intermediate video, and it contains the lip motion information. We extract head pose and expression from the video for enhancement step.

### B. Keypoint-based enhancement frontend

The frontend contains five modules (indicated by the blue blocks in Fig. 1): Canonical keypoint estimator, appearance feature estimator, head pose and expression estimator, motion field estimator and generator.

*1) Canonical keypoint estimator:* Using the canonical keypoint estimator, we can extract canonical keypoint from the source image. Note that canonical keypoint is the keypoint associated only with the identity, and not affected by head pose and expression. These extracted keypoints shall encode a person's head geometry feature in a neutral expression and
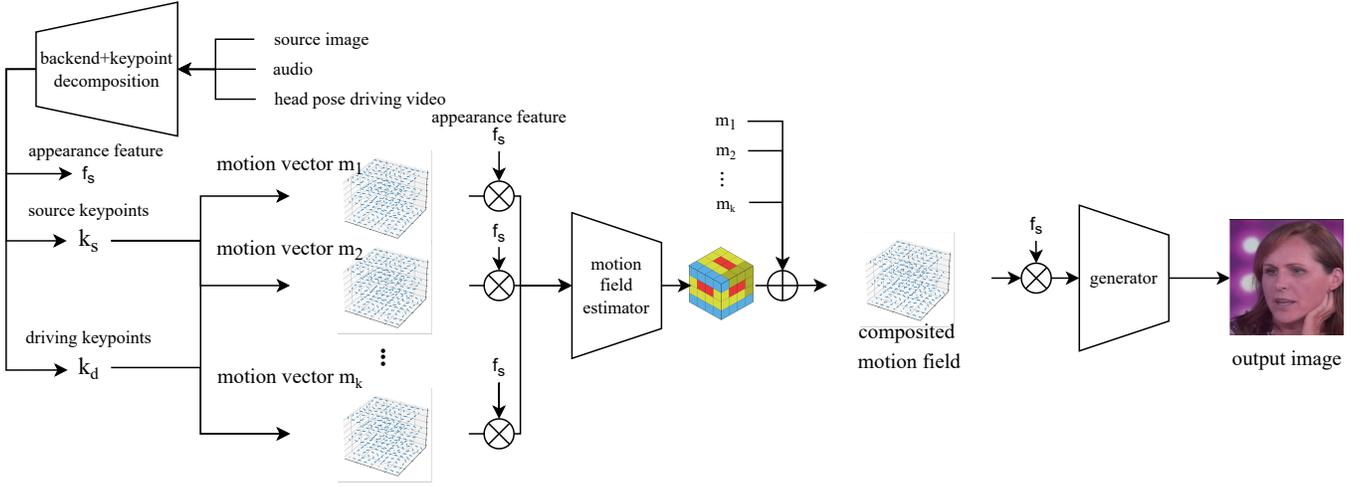
Fig. 2. Generating from keypoint representation. Source keypoints and driving keypoints are derived from keypoint decomposition. We combined each pair keypoints, and used first order approximation to get the motion vectors. The motion field mask is obtained by the motion field estimator. Then the mask is applied on the sum of motion vectors to get the composited motion field. The appearance feature is warped by composited motion field. After that the warped appearance feature is used to generate the final output image.

pose. Specifically, we extracted 15 3D canonical keypoints from the source image in this paper.

*2) Appearance feature estimator:* With the appearance feature estimator, the appearance features are extracted from the source image. Appearance feature are features related to the appearance of the face, such as skin and hair style. A 3D extraction network is applied to extract 3D appearance feature. It is important to use 3D feature. Only with the 3D appearance feature, can we restore it back to the photo after rotation and translation.

*3) Head pose and expression estimator:* The head pose and expression estimator is for estimating head pose and expression from the image. This module is used in two places in our approach. One is for estimating the expression from the backend output, and the other is for estimating the head pose from the head pose driving video. Specifically, head pose contains rotation matrix and translation vector. Rotation matrix controls the 3D rotate angle of the head, and translation vector controls the center position of the head. Expression is parameterized by vectors with the same number of the canonical keypoint.

With the canonical keypoint, head pose and expression, we give the equation of the 3D keypoint decomposition:

$$k_{s,i} = T(k_{c,i}, R_s, t_s, exp_s) = R_s k_{c,i} + t_s + exp_s \quad (1)$$

$$k_{d,i} = T(k_{c,i}, R_d, t_d, exp_b) = R_d k_{c,i} + t_d + exp_b \quad (2)$$

where $k_{c,i}$ is the *i*-th canonical keypoint extracted from the source image. $R_s$ and $t_s$ are the rotation matrix and the translation vector extracted from the source image. $exp_s$ is the expression extracted from the source image. $k_{s,i}$ is the *i*-th source keypoint. $R_d$ and $t_d$ are the rotation matrix and the translation vector extracted from the head pose driving video. $exp_b$ is the expression extracted from the backend output. $k_{d,i}$ is the *i*-th driving keypoint. As illustrated in Fig. 1, the source

keypoints are obtained by compositing head pose, expression and canonical keypoint extracted from the source image. We get driving keypoint by compositing head pose extracted from the head pose driving video, expression extracted from backend output and canonical keypoint extracted from the source image. Source keypoint and driving keypoint are uniquely determined by identity, expression and head pose. It is of paramount importance for our approach to do the 3D keypoint decomposition. Our approach differs from prior 2D-based audio-driven talking head synthesis methods with regard to the keypoint decomposition. The keypoint decomposition helps learn controllable representations. Note that unlike OSFV [19], our model is audio driven, and only extracts expression feature from the driving video. Our frontend step learns how to composite the feature extracted from different images, and restore the final image. With the keypoint decomposition, we can maximize the retention of appearance feature. This is helpful for reducing the cut feeling in the mouth mapping and the lack of skin highlights.

*4) Motion field estimator:* From the above three modules we can get source keypoints and driving keypoints. As shown in Fig. 2, the source keypoints $k_s$ and the driving keypoints $k_d$ are used as input to get motion vectors. Each pair of keypoints can be used to get a motion vector. A motion vector is a field that defines the movement tendency of every point in the 3D space. It is estimated by the same method of first-order approximation [14]. We apply each motion vector to appearance feature and get the appearance feature warped by a motion vector. All warped appearance features are used as input for motion field estimator. The output is a 3D mask that weights different parts of the head. Then we apply the mask on the sum of the all motion vectors. The output is the composited motion field and a 2D occlusion mask, and we use this field to warp the appearance feature to get the warped appearance

feature.

## C. Generator

Using the warped appearance feature obtained by the motion field estimator as input, the generator transforms the 3D warped appearance feature to the 2D final output image. Note that the warped appearance feature is multiplied by the occlusion mask after a convolution layer. For projecting the 3D feature to the 2D image, the framework of the generator uses a series of 2D upsampling layers.

## D. Training

Because we hope that the expression from the backend can be used in frontend, the loss of this step consists of perceptual loss and expression loss. Loss functions are shown as follows:

$$\mathcal{L}_p = \sum_{i \in Pyr} \|VGG(Pyr_i G) - VGG(Pyr_i s)\|_1 \quad (3)$$

$$\mathcal{L}_{exp} = \|E_{d,i}\|_1 \quad (4)$$

$\mathcal{L}_p$ is the perceptual loss that helps to produce sharp-looking outputs. VGG is a face recognition model [27]. A pre-trained VGG framework is used to get the features of the final output image and the source image, and the L1 loss is computed between the features. $G$ is the generated final image. $s$ is the source image. We use a pyramid structure, $Pyr_i G$ and $Pyr_i s$ means the $i_t h$ scaled images. First features are extracted using the VGG network. After that, two images are fed into a downsample layer, then we use the VGG network to extract features from the sownsampled images and compute the L1 loss once again. We repeat this process five times, and compute L1 loss in different resolutions. Losses of different resolutions are added to compute the pyramid perceptual loss.

Expression loss $\mathcal{L}_{exp}$ is for penalizing the large keypoint deformation because the expression is a relatively small change compared to the head pose. Using the expression loss, we can clamp the magnitude of the keypoint deformation.

Our model is trained in two steps. First, we used the loss functions as the loss functions used in OSFV [19] to train the frontend. Second, we added the backend to fine-tune the frontend. Our model was trained on Voxceleb1 which consists of talking head videos, and each video contains single person. The input of the second step is a source image sampled from video and a driving image sampled from backend output.

## IV. Experiments

### A. Dataset

Our experiment was based on Voxceleb1 [10]. We used the same preprocessing method in FOMM [14], which has 18672 videos for train and 525 videos for test. Voxceleb1 had different resolution videos and all videos were converted to 25 fps in our training. Since there were many extreme head pose cases, different light conditions, and different ethnicity in the dataset, our model was robust.

### B. Implementation detail

Canonical keypoint estimator, appearance feature estimator, head pose and expression estimator, motion field estimator and generator are with the same structure of the OSFV [19]. We fine tuned the pretrained backend and frontend model with our training loss.

We used the ADAM optimizer with $\beta_1 = 0.5$ and $\beta_2 = 0.999$ and trained on 4 NVIDIA 24GB GTX3090 GPUs. The model was trained for 200 epochs, and each epoch needs 2 hours.

### C. Metrics

We used PSNR and SSIM [28] as metrics to quantify the faithfulness of the synthesized videos. SSIM measures the structural similarity between ground truth and synthesized videos. Compared with PSNR, it is more robust to global illumination changes. Given a reference image $s$ and a test image $g$, both of which have a size of M×N, the PSNR between s and g is defined by:

$$\mathrm{PSNR}(s,g) = 10\log_{10}\left(255^2/MSE(s,g)\right)$$

where

$$\mathrm{MSE}(s,g) = \frac{1}{MN}\sum_{i=1}^{M}\sum_{j=1}^{N}(s_{ij} - g_{ij})^2$$

PSNR is used to measure image reconstruction quality. It computes the mean squared error between the input and the output images. SSIM is defined as:

$$\mathrm{SSIM}(s,g) = l(s,g)c(s,g)s(s,g)$$

where

$$\begin{cases} l(s,g) = \frac{2\mu_s\mu_g + C_1}{\mu_s^2 + \mu_g^2 + C_1} \\ c(s,g) = \frac{2\sigma_s\sigma_g + C_2}{\sigma_s^2 + \sigma_g^2 + C_2} \\ s(s,g) = \frac{\sigma_{sg} + C_3}{\sigma_s\sigma_g + C_3} \end{cases} \quad (5)$$

$\mu_s$ and $\mu_g$ are the mean luminance of the images. $\sigma_s$ and $\sigma_g$ are the variance of the luminance. $\sigma_{sg}$ is the covariance between the two images s and g. The values of SSIM are from 0 to 1. Value 0 means no relation between images, and 1 means two images are the same image.

The test videos for our experiment were generated by random choosing video from the dataset as the ground truth. We randomly chose one frame from it as the source image, and use the same video as the head pose driving video. We calculate these two metrics on the test set, and evaluate the mean score.

In addition, we also conducted a subjective test to evaluate the performance of our proposed method where subjects were recruited to rate the mean opinion scores (MOS) of lip synchronization, head pose naturalness, and video realness. The value ranges from 1 to 5 (higher is better).

## D. Objective evaluation

We chose Wav2lip [12] and PC-AVS [22] as the backend method, and compared our method with the results before enhancement for quantitative evaluation. As shown in Table I, both methods have improvement on SSIM and PSNR. As illustrated in Fig. 3, we can see that our method has a better color balance than the PC-AVS. The PC-AVS results had a pale skin, and our method results had a normal skin. Wav2lip has a block of blurring arround the mouth, and our method result does not have the blurring.

## E. Subjective evaluation

We conducted a subjective test to compare different methods. Each method synthesized 13 videos which were then used as stimuli. The subjective test was performed online where ten subjects were recruited to rate MOS of lip synchronization, head pose naturalness, and video realness. The subjects were asked to wear headphones and complete the experiment in a quiet environment. Table II shows that with KPBE enhancement, the lip synchronization, head pose naturalness, and video realness of the enhanced videos are improved. Fig. 3 shows synthesized video frames by PC-AVS, Wav2lip, and their KPBE enhancements. First line and second line are the results of the PC-AVS and our method. PC-AVS could synthesize pretty good head pose and lip motion, but they sometimes synthesized people without highlights on skin. When the head had a relatively large rotation, the face would be deficient. In contrast our model had a better face highlights and had no facial deficient.

The significant enhancement is on the Wav2lip method, and their method often synthesizes accurate lip motion but with low visual quality on the mouth part. Our method makes it overcome this shortcoming. From the third line and forth line in Fig. 3, we can see that the result of Wav2lip has some blurring at the mouth part. There is a dividing line between the neck and the face. Our method result has a clearer mouth, and there is no dividing line.
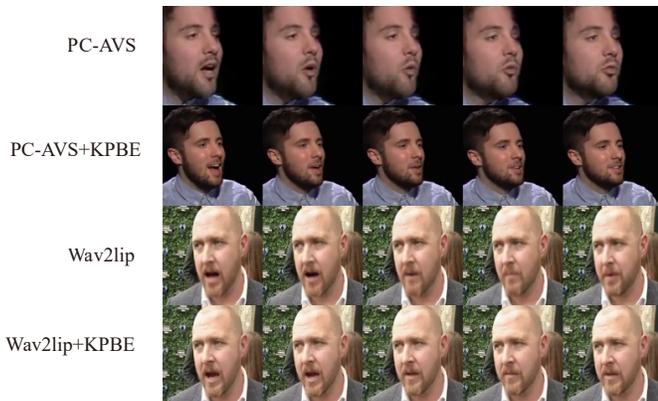


Fig. 3. Synthesized video frames by PC-AVS, Wav2lip, and our method.



Fig. 4. Novel viewpoint video synthesis. With the user-defined head pose matrix, we can reenact the audio-driven talking head video from a novel viewpoint.

TABLE I
OBJECTIVE RESULTS. WE CALCULATED THE METRICS ON THE TEST SET, AND EVALUATE THE MEAN SCORE.

| Method | SSIM↑ | PSNR↑ |
|---|---|---|
| Wav2lip | 0.85 | 28.67 |
| Wav2lip+KPBE(Ours) | **0.90** | **30.37** |
| PC-AVS | 0.80 | 22.19 |
| PC-AVS+KPBE(Ours) | **0.82** | **23.64** |

TABLE II
MEAN OPINION SCORES. VALUE FROM 1 TO 5, HIGHER IS BETTER.

| Method | Lip Sync | Head Pose Naturalness | Video Realness |
|---|---|---|---|
| Wav2lip [12] | 4.11 | 4.03 | 4.09 |
| Wav2lip+KPBE(Ours) | **4.14** | **4.14** | **4.27** |
| PC-AVS [22] | 3.44 | 3.12 | 3.25 |
| PC-AVS+KPBE(Ours) | **3.70** | **3.31** | **3.41** |
| Ground_truth | 4.59 | 4.50 | 4.73 |

## F. Novel viewpoint video synthesis

To test the free-view video generation capabilities of our model, we synthesized the video from a novel viewpoint with the user-defined rotation matrix and translation matrix. The results are shown in Fig. 4. Using the head pose matrix, we can freely control the rotation of the head pose. Using the translation matrix, we can control the displacement of the head.

## V. CONCLUSIONS

Audio-driven talking head synthesis attracts increasing attention in recent years. We proposed a novel keypoint-based enhancement method for audio-driven talking head synthesis to generate enhanced video with better lighting balance and more natural expressions and head poses. Using keypoint decomposition, our method could disentangle the image features into appearance features, canonical keypoints, and head pose matrix. We extracted expression from the backend output and head pose from the head pose driving video. After that, the motion field based generator was used to generate the final image. Experiments results verified that our method can reduce the cut feeling in the mouth mapping and the lack of skin highlights. Limited by the low resolution and language propensity of the dataset, our model cannot synthesize high-resolution videos and accurate lip movements in some languages. In the future,

we will focus on improving the resolution and cross-language accuracy in audio-driven talking head synthesis.

## REFERENCES

[1] Anderson, R., Stenger, B., Wan, V. & Cipolla, R. An Expressive Text-Driven 3D Talking Head. *ACM SIGGRAPH 2013 Posters*. pp. 1-1 (2013)

[2] Chen, L., Maddox, R., Duan, Z. & Xu, C. Hierarchical Cross-Modal Talking Face Generation with Dynamic Pixel-Wise Loss. *Proceedings Of The IEEE/CVF Conference On Computer Vision And Pattern Recognition*. pp. 7832-7841 (2019)

[3] Chen, L., Cui, G., Liu, C., Li, Z., Kou, Z., Xu, Y. & Xu, C. Talking-Head Generation with Rhythmic Head Motion. *European Conference On Computer Vision*. pp. 35-51 (2020)

[4] Das, D., Biswas, S., Sinha, S. & Bhowmick, B. Speech-Driven Facial Animation Using Cascaded Gans for Learning of Motion and Texture. *European Conference On Computer Vision*. pp. 408-424 (2020)

[5] Edwards, P., Landreth, C., Fiume, E. & Singh, K. JALI: An Animator-Centric Viseme Model for Expressive Lip Synchronization. *ACM Transactions On Graphics (TOG)*. **35**, 127:1-127:11 (2016,7)

[6] Guo, Y., Chen, K., Liang, S., Liu, Y., Bao, H. & Zhang, J. Ad-Nerf: Audio Driven Neural Radiance Fields for Talking Head Synthesis. *Proceedings Of The IEEE/CVF International Conference On Computer Vision*. pp. 5784-5794 (2021)

[7] Jamaludin, A., Chung, J. & Zisserman, A. You Said That?: Synthesising Talking Faces from Audio. *International Journal Of Computer Vision*. **127**, 1767-1779 (2019)

[8] Ji, X., Zhou, H., Wang, K., Wu, W., Loy, C., Cao, X. & Xu, F. Audio-Driven Emotional Video Portraits. *Proceedings Of The IEEE/CVF Conference On Computer Vision And Pattern Recognition*. pp. 14080-14089 (2021)

[9] Mildenhall, B., Srinivasan, P., Tancik, M., Barron, J., Ramamoorthi, R. & Ng, R. Nerf: Representing Scenes as Neural Radiance Fields for View Synthesis. *European Conference On Computer Vision*. pp. 405-421 (2020)

[10] Nagrani, A., Chung, J. & Zisserman, A. VoxCeleb: A Large-Scale Speaker Identification Dataset. *Interspeech 2017*. pp. 2616-2620 (2017,8)

[11] Pham, H., Cheung, S. & Pavlovic, V. Speech-Driven 3D Facial Animation with Implicit Emotional Awareness: A Deep Learning Approach. *Proceedings Of The IEEE Conference On Computer Vision And Pattern Recognition Workshops*. pp. 80-88 (2017)

[12] Prajwal, K., Mukhopadhyay, R., Namboodiri, V. & Jawahar, C. 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*. pp. 484-492 (2020)

[13] Richard, A., Lea, C., Ma, S., Gall, J., De la Torre, F. & Sheikh, Y. Audio- and Gaze-Driven Facial Animation of Codec Avatars. *Proceedings Of The IEEE/CVF Winter Conference On Applications Of Computer Vision*. pp. 41-50 (2021)

[14] Siarohin, A., Lathuilière, S., Tulyakov, S., Ricci, E. & Sebe, N. First Order Motion Model for Image Animation. *Advances In Neural Information Processing Systems*. **32** (2019)

[15] Song, L., Wu, W., Qian, C., He, R. & Loy, C. Everybody's Talkin': Let Me Talk as You Want. *IEEE Transactions On Information Forensics And Security*. (2022)

[16] Suwajanakorn, S., Seitz, S. & Kemelmacher-Shlizerman, I. Synthesizing Obama: Learning Lip Sync from Audio. *ACM Transactions On Graphics (ToG)*. **36**, 1-13 (2017)

[17] Taylor, S., Kim, T., Yue, Y., Mahler, M., Krahe, J., Rodriguez, A., Hodgins, J. & Matthews, I. A Deep Learning Approach for Generalized Speech Animation. *ACM Transactions On Graphics (TOG)*. **36**, 1-11 (2017)

[18] Thies, J., Elgharib, M., Tewari, A., Theobalt, C. & Nießner, M. Neural Voice Puppetry: Audio-driven Facial Reenactment. *European Conference On Computer Vision*. pp. 716-731 (2020)

[19] Wang, T., Mallya, A. & Liu, M. One-Shot Free-View Neural Talking-Head Synthesis for Video Conferencing. *Proceedings Of The IEEE/CVF Conference On Computer Vision And Pattern Recognition*. pp. 10039-10049 (2021)

[20] Zakharov, E., Ivakhnenko, A., Shysheya, A. & Lempitsky, V. Fast Bi-Layer Neural Synthesis of One-Shot Realistic Head Avatars. *European Conference On Computer Vision*. pp. 524-540 (2020)

[21] Zhou, Y., Han, X., Shechtman, E., Echevarria, J., Kalogerakis, E. & Li, D. MakeltTalk: Speaker-Aware Talking-Head Animation. *ACM Transactions On Graphics (TOG)*. **39**, 1-15 (2020,12)

[22] Zhou, H., Sun, Y., Wu, W., Loy, C., Wang, X. & Liu, Z. Pose-Controllable Talking Face Generation by Implicitly Modularized Audio-Visual Representation. *Proceedings Of The IEEE/CVF Conference On Computer Vision And Pattern Recognition*. pp. 4176-4186 (2021)

[23] Zhou, Y., Xu, Z., Landreth, C., Kalogerakis, E., Maji, S. & Singh, K. Visemenet: Audio-Driven Animator-Centric Speech Animation. *ACM Trans. Graph.*. **37**, 161:1-161:10 (2018,7)

[24] Zhu, H., Huang, H., Li, Y., Zheng, A. & He, R. Arbitrary Talking Face Generation via Attentional Audio-Visual Coherence Learning. *Proceedings Of The Twenty-Ninth International Conference On International Joint Conferences On Artificial Intelligence*. pp. 2362-2368 (2021)

[25] Sitzmann, V., Zollhöfer, M. & Wetzstein, G. Scene representation networks: Continuous 3d-structure-aware neural scene representations. *Advances In Neural Information Processing Systems*. **32** (2019)

[26] Lu, Y., Chai, J. & Cao, X. Live Speech Portraits: Real-Time Photorealistic Talking-Head Animation. *ACM Trans. Graph.*. **40**, 220:1-220:17 (2021,12)

[27] Omkar, M., Vedaldi, A., Zisserman, A. & Others Deep face recognition. *Bmvc*. **1**, 6 (2015)

[28] Wang, Z., Bovik, A., Sheikh, H. & Simoncelli, E. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions On Image Processing*. **13**, 600-612 (2004)