

# FONT: FLOW-GUIDED ONE-SHOT TALKING HEAD GENERATION WITH NATURAL HEAD MOTIONS

Jin Liu<sup>1,2</sup>, Xi Wang<sup>1\*</sup>, Xiaomeng Fu<sup>1,2</sup>, Yesheng Chai<sup>1</sup>, Cai Yu<sup>1,2</sup>, Jiao Dai<sup>1\*</sup>, Jizhong Han<sup>1</sup>

<sup>1</sup>Institute of Information Engineering, Chinese Academy of Sciences, Beijing, China

<sup>2</sup>School of Cyber Security, University of Chinese Academy of Sciences, Beijing, China

## ABSTRACT

One-shot talking head generation has received growing attention in recent years, with various creative and practical applications. An ideal natural and vivid generated talking head video should contain natural head pose changes. However, it is challenging to map head pose sequences from driving audio since there exists a natural gap between audio-visual modalities. In this work, we propose a Flow-guided One-shot model that achieves NaTural head motions(FONT) over generated talking heads. Specifically, we design a probabilistic CVAE-based model to predict head pose sequences from driving audio and source face. Then we develop a keypoint predictor that produces unsupervised keypoints describing the facial structure information from the source face, driving audio and pose sequences. Finally, a flow-guided occlusion-aware generator is employed to produce photo-realistic talking head videos from the estimated keypoints and source face. Extensive experimental results prove that FONT generates talking heads with natural head poses and synchronized mouth shapes, outperforming other compared methods.

**Index Terms**— Talking Head Generation, Generative Model, Audio Driven Animation

## 1. INTRODUCTION

Given one source face and driving audio, one-shot talking head generation aims to synthesize a talking head video with reasonable facial animations corresponding to the driving audio [1]. This task receives growing attention since it can be used in a wide range of multimedia applications, e.g. video dubbing, digital avatar animation and short video creation.

Some methods [2, 3] have been proposed to edit the mouth area to achieve lip synchronization. However, they neglect the modeling of head motions, thus generating unnatural talking heads that are far from satisfactory from human observation and practical applications. Therefore, researchers turn to focus on generating talking heads with head pose changes. Recent works [4, 5] choose to introduce an extra auxiliary pose

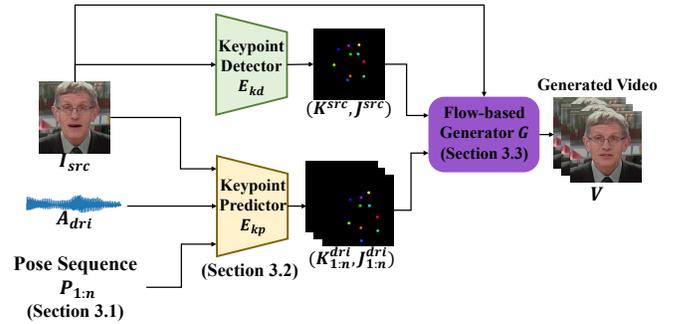


Fig. 1. Overview of the proposed method.

video that guides the head motion changes in the generated talking heads. This formula limits the generalization since it is tedious to find another pose video in one-shot scenario. Hence, some methods try to predict head pose sequence from driving audio.

It is challenging to map driving audio signal into head pose sequence, since there exists natural gap between visual and audio modalities. A great many works [6, 1, 7, 8] are proposed to infer head motions from driving audio and source face. However, they neglect the uncertainty in the head pose prediction task and fail to produce natural head poses. In fact, the mapping from driving audio signal to head pose sequence is inherently a one-to-many problem. In real life, people can behave differently in head poses even speaking the same content. Previous methods adopt deterministic models like LSTM or MLP to perform the task, which fundamentally ignore the uncertainty lying between audio signals and head poses. Furthermore, the lack of facial structure modeling in their generation process also leads to blurry artifacts and poor lip-sync quality.

To solve the above problem, we propose a Flow-guided One-shot talking head generation network with NaTural head motions (FONT). The overall framework is shown in Fig. 1. The driving pose sequences come from the well-designed head prediction module. Detailedly, a probabilistic CVAE-based network is adopted to generate head pose sequences from driving audio and source face, during which the structural similarity loss is imposed instead of MSE loss. The above operations model the uncertainty and the ambiguous

\*Corresponding authors.

This research is supported in part by the National Key Research and Development Program of China (No. 2020AAA0140000), and the National Natural Science Foundation of China (No. 61702502).

correspondences between audio and head pose modalities, contributing to natural driving head pose sequences. Then inspired by image animation work FOMM [9], we predict unsupervised keypoints from the source face, driving audio and poses to model the facial structure location. Finally, the occlusion-aware flow-guided generator produces motion flow to indicate the local facial texture variance and generates new talking heads with natural head poses. Moreover, to improve the lip-sync quality, a pre-trained lip-sync discriminator is utilized during the training process.

Our contributions are as follows: (1) We develop a new flow-guided one-shot talking head generation framework that produces natural head motions. (2) A probabilistic CVAE-based network is designed to generate natural head poses from driving audio and source face. (3) We present a flow-guided occlusion-aware generator to produce keypoint-based motion flow indicating facial structure, thus generating natural talking heads. (4) Extensive experimental results prove that our proposed framework achieves the state-of-the-art level compared to other methods.

## 2. RELATED WORK

**One-shot Talking Head Generation.** One-shot talking head generation [10] has long been a significant research topic in the computer vision field. Speech2Vid [11] generates talking faces via an encoder-decoder structure and a refinement module. DAVS [12] and ATVG [2] further improve the quality using disentangled audio-visual representation and external structural information guidance. Wav2Lip [3] applies a pre-trained lip-sync discriminator to improve the generation results. Nevertheless, the above methods merely edit the mouth area and leave other facial regions unchanged, producing unnatural and less realistic talking head videos.

Full-frame talking head generation produces new facial areas but also the neck part of the person, together with the background. MakeitTalk [13] predicts content and speaker-aware displacement on facial landmarks to guide the talking face generation process. To improve the realness, some methods focus on talking heads with natural head poses [4, 7]. However, their lack of face structural modeling and the mouth shape constraint causes identity mismatch and poor lip synchronization performance. However, FONT utilizes motion flow as facial structure information and the lip-syn discriminator to solve the above problem.

**Head Pose Control.** Since there is no explicit head pose information contained in the driving audio signal, it is challenging to achieve head pose control and generate talking heads with natural head motions. Early methods [3, 13] focus on mouth shape accuracy and produce almost still talking heads. Later, PC-AVS [4] first propose to rely on auxiliary pose video to obtain head pose guidance. It limits the generalization of this task since obtaining a long pose video is cumbersome in the one-shot scenario. Several methods turn

to infer pose sequences directly from audio. Audio2Head [7] and AVCT [8] designs a motion-aware LSTM-based network to predict head motions, while HDTF [1] utilizes Multilayer Perceptron to predict head pose coefficients in morphable face model [14]. However, the correspondence between audio and head poses contains uncertainty and predicting head poses from audio is actually an ill-posed problem. In real life, people may act different poses even speaking the same content. Hence, instead of utilizing *deterministic* models like other methods, we choose the *probabilistic* CVAE-based [15] network to model the uncertainty in pose generation.

## 3. METHODOLOGY

The overview of the proposed method is shown in Fig. 1. Driving pose sequence will be predicted first. Then the source face  $I_{src}$ , driving audio and driving pose sequence are fed into the Keypoint Predictor  $E_{kp}$  to predict unsupervised driving keypoints. Then the driving keypoints and source keypoints from  $I_{src}$  are taken as inputs to the Generator and produce final talking head videos.

### 3.1. Head Pose Prediction Module

To generate talking heads with natural head motions, the natural head pose sequences should be predicted first. Different from previous work [6, 7, 8] which adopt deterministic models like LSTM and traditional GAN to generate pose sequences, we design a VAE-based probabilistic model inspired by CVAE [15]. Pose generation is actually an uncertain ill-posed problem since people may behave differently when speaking the same corpus. Hence, the probabilistic model is more suitable for this task.

The head pose prediction module of FONT is shown in Fig. 2. Specifically, a 6-dim vector (i.e., 3 for rotation, 1 for scale and 2 for translation) is adopted as pose information representation for each frame. Specifically, during the training stage, we utilize paired pose clip  $p_{1:t}$ , corresponding audio  $A_p$  and head image  $I_p$  as inputs. They will be fed into the encoder to predict mean and standard deviation values, which will be later used for re-parametrization. Finally, the sampled data,  $I_p$  and  $A_p$  are passed into the decoder to predict pose clip  $\hat{p}_{1:t}$ . The face image and audio are served as the condition information to guide the generation of pose sequence. It is noteworthy that we learn the difference of poses compared to the first frame in  $p_{1:t}$  instead of pose itself. This setting eliminates the influence of the various initial head poses in different pose clips.

As for the loss constraints, commonly used reconstruction loss like MSE loss is not suitable for pose generation, since the task is actually an ill-posed one-to-many mapping problem. Therefore, we utilize the Structural Similarity [16] to keep the consistency between the generated and ground truth

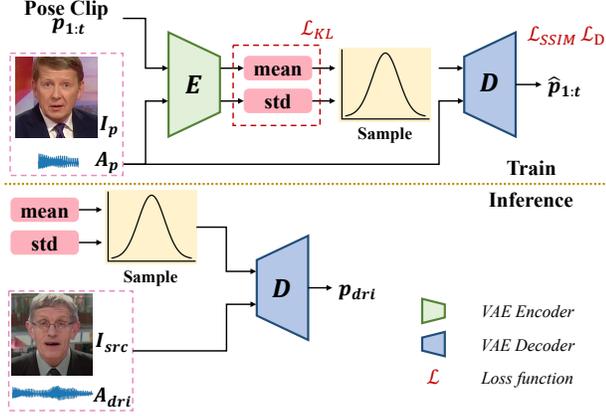


Fig. 2. Overview of head pose prediction module.

pose sequence:

$$\mathcal{L}_{SSIM} = 1 - \frac{(2\mu\hat{\mu} + C_1)(2cov + C_2)}{(\mu^2 + \hat{\mu}^2 + C_1)(\sigma^2 + \hat{\sigma}^2 + C_2)}. \quad (1)$$

$\hat{\mu}$  and  $\hat{\sigma}$  are mean and standard deviation of generated pose sequence while  $\mu$  and  $\sigma$  are that of the ground truth pose.  $cov$  is the covariance between two sequences and  $C$  is the constants to stabilize the division. Meanwhile, to guarantee the similarity between latent space distribution and Gaussian distribution, we define  $\mathcal{L}_{KL}$  as the KL-Divergence between the above two distributions. Furthermore, the discriminator is also adopted to improve the realism of the generated pose.

$$\mathcal{L}_D = \log D(p_{gt}) + \log(1 - D(p)). \quad (2)$$

The overall loss of pose generation is defined by the combination of  $\mathcal{L}_{SSIM}$ ,  $\mathcal{L}_D$  and  $\mathcal{L}_{KL}$ .

During inference, driving audio will be divided into several audio clips. They will be fed into the decoder along with  $I_{src}$  and sampled latent data to produce pose clips. Finally, the pose clips will be stacked together in chronological order and added to the initial head pose to form the driving pose sequence  $P_{1:n}$ .

### 3.2. Keypoint Predictor

Inspired by the widely used image animation work FOMM [9], we choose to use the unsupervised keypoints and their first order dynamics as the structure representation.

As illustrated in Fig. 1, the Keypoint Predictor  $E_{kp}$  takes source face  $I_{src}$ , driving audio  $A_{dri}$  and predicted pose sequence  $p_{1:n}$  as inputs.  $E_{kp}$  first utilized three different encoders to extract the corresponding information. Then the three features are combined and fed into the LSTM-based decoder to recurrently predict the corresponding unsupervised structure representation. At each time step  $t$ , the representation contains learned keypoints  $K_t \in \mathbb{R}^{N \times 2}$  and the corresponding first order dynamics, i.e. jacobians  $J_t \in \mathbb{R}^{N \times 2 \times 2}$  which describes the local affine transformation in the neighborhood area around each keypoint. As for the initial source

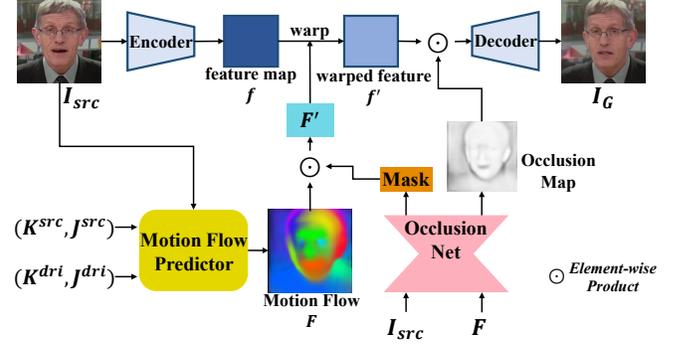


Fig. 3. Overview of the flow-guided generator.

structure representation, the pre-trained keypoint detector  $E_{kd}$  from FOMM [9] is utilized to provide accurate initial keypoints and first order dynamics. The whole procedure is formulated as:

$$\begin{aligned} (K^{src}, J^{src}) &= E_{kd}(I_{src}), \\ (K_{1:n}^{dri}, J_{1:n}^{dri}) &= E_{kp}(I_{src}, A_{dri}, p_{1:n}). \end{aligned} \quad (3)$$

For the training loss, we regard the  $E_{kd}$  as a teacher network and hope  $E_{kp}$  to learn the knowledge of visual structure representation contained in pre-trained  $E_{kd}$ . We further define the motion representation of the corresponding ground truth video frame extracted by  $E_{kd}$  as supervision, i.e.  $(K^{gt}, J^{gt})$ . Therefore, the loss term of  $E_{kp}$  is as follows:

$$\mathcal{L}_{kp} = \frac{1}{N} \sum_{i=1}^N \left( \|K_i^{kp} - K_i^{gt}\|_1 + \|J_i^{kp} - J_i^{gt}\|_1 \right). \quad (4)$$

In this way, the source and driving structure representation are both successfully obtained.

### 3.3. Flow-guided Generator

As shown in Fig. 3, the flow-guided generator produces talking head  $I_G$  given  $I_{src}$ , source and driving structure representation. It mainly contains the motion flow predictor, the occlusion net, the image encoder and decoder. The motion flow predictor first predicts motion flow  $F$  indicating the variation in each part of the face from source to driving. Then  $I_{src}$  and  $F$  are fed into the Occlusion Net to predict the flow mask and occlusion map. The masked motion flow  $F'$  is utilized to warp the feature map  $f$  of  $I_{src}$  to obtain warped feature  $f'$ . Finally, the occluded feature is sent to the decoder to produce talking head  $I_G$ . In this way, the decoder obtain the source face texture, motion variance and different confidences among the feature map, which all contribute to the accurate generation process. The encoder and decoder consists of several convolutional and up-sampling layers. The occlusion net is based on the hourglass net while the motion flow predictor relies on the numerical calculation between two structure representations. During training, the perceptual loss is utilized between the generated frame  $I_G$  and ground truth frame  $I_{gt}$ :

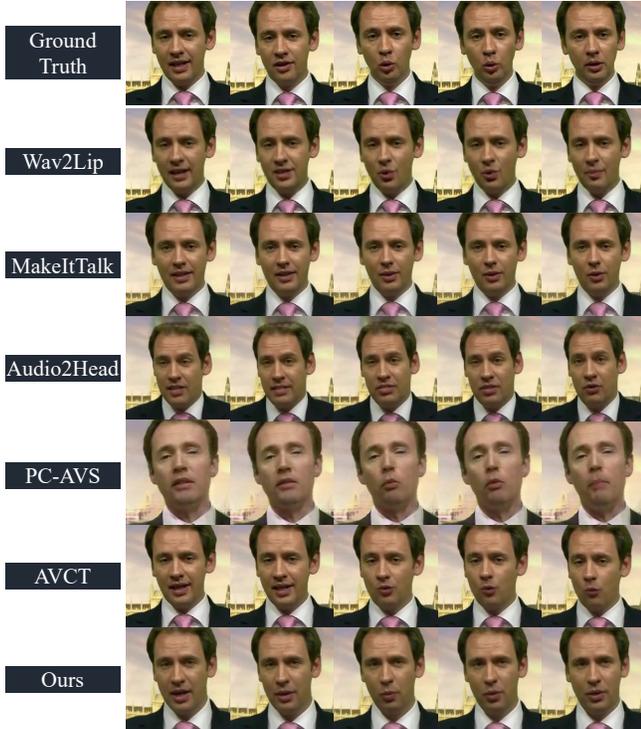


Fig. 4. Qualitative Comparison with other methods.

Table 1. Quantitative comparisons on LRW dataset. The bold and underlined notations represents the Top-2 results.

Method	SSIM $\uparrow$	CPBD $\uparrow$	LMD $\downarrow$	LSE-C $\uparrow$
Wav2Lip	<u>0.812</u>	0.172	5.73	<b>7.237</b>
MakeItTalk	0.796	0.161	7.13	3.141
Audio2Head	0.743	0.168	7.34	2.135
PC-AVS	0.778	<u>0.185</u>	3.93	6.420
AVCT	0.805	0.181	<u>3.56</u>	6.567
Ground Truth	1.000	0.189	0.00	6.876
Ours	<b>0.825</b>	<b>0.187</b>	<b>3.48</b>	<u>6.572</u>

$$\mathcal{L}_{per} = \sum_{i=1}^l \|VGG_i(I_G) - VGG_i(I_{gt})\|_1, \quad (5)$$

where  $VGG(\cdot)$  denotes the  $i_{th}$  channel feature of the pre-trained VGG network. Furthermore, to improve the lip-sync quality, we adopt a pre-trained discriminator to predict the embedding of corresponding audio and video. The discriminator [3] is trained to judge the synchronization between randomly sampled audio-visual pairs. We adopt the cosine-similarity between audio and video embedding  $a$  and  $v$  extracted by the discriminator as the lip-sync loss to indicate the probability of whether the pair is in-sync.

$$\mathcal{L}_{sync} = \frac{v \cdot a}{\max(\|v\|_2 \cdot \|a\|_2, \epsilon)} \quad (6)$$



Fig. 5. Large-pose qualitative comparison results.

Table 2. Quantitative comparisons on HDTF dataset.

Method	SSIM $\uparrow$	CPBD $\uparrow$	LMD $\downarrow$	LSE-C $\uparrow$
Wav2Lip	<u>0.786</u>	<b>0.176</b>	2.89	6.97
MakeItTalk	0.751	0.132	5.46	4.87
Audio2Head	0.735	0.145	4.83	3.90
PC-AVS	0.762	0.164	3.18	<u>7.18</u>
AVCT	0.769	0.167	<u>2.71</u>	7.09
Ground Truth	1.000	0.181	0.00	8.58
Ours	<b>0.789</b>	<u>0.169</u>	<b>2.69</b>	<b>7.22</b>

## 4. EXPERIMENTS

### 4.1. Experimental Settings

**Datasets.** We evaluate our method on LRW [17] and HDTF [1] datasets. The LRW dataset contains over 1000 short utterances of each 500 different words and all the videos are extracted from BBC television in the wild. The HDTF dataset is a large in the wild audio-visual dataset that consists of long utterances of over 300 subjects.

**Implementation Details.** The face video frames are cropped to  $256 \times 256$  size at 25 FPS and the audio is pre-processed into 16kHz. We compute 28-dim MFCC feature with a window size of 10ms to produce a  $28 \times 12$  feature for each video frame.

As for training, our method is trained in stages. The generator is trained with keypoint predictor after the latter gets stable results. The ADAM optimizer is adopted with an initial learning rate as  $2 \times 10^{-4}$ , which linearly decreases to  $2 \times 10^{-5}$ . We train our model on 1 Tesla V100 GPU and each part requires 0.5, 2 and 3 days for training respectively.

### 4.2. Experimental Results

**Evaluation Metrics.** The performance is evaluated on image quality and lip-sync quality. The SSIM [16] and Cumulative Probability of Blur Detection (CPBD) [18] scores are utilized to judge the quality of talking head frames. For lip-sync quality, the Landmark Distance(LMD) and Lip-Sync Error-Confidence(LSE-C) are applied. LMD means the average Euclidean distance between corresponding facial land-



**Fig. 6.** Qualitative results driven by the same audio and different source faces on HDTF.

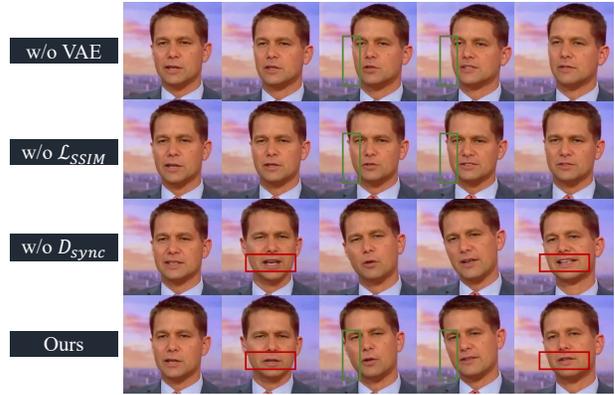
marks. LSE-C is the confidence score of the correspondence between audio and video features extracted from pre-trained SyncNet [19].

**Quantitative Results.** We choose several state-of-the-art methods as comparison, i.e. Wav2Lip [3], MakeItTalk [13], Audio2Head [7], PC-AVS [4] and AVCT [8]. The frames of each method are generated using their official code. The head poses of Wav2Lip and PC-AVS are fixed since they can not obtain head poses from audio. The Ground Truth results are also added for better comparison.

Detailed results on LRW and HDTF can be found in Table 1 and Table 2. FONT achieves the best performance under most of the evaluation metrics on both datasets. As Wav2Lip merely edits the mouth area, it achieves better CPBD score on HDTF. Furthermore, as mentioned in by PC-AVS [4], the leading LSE-C only means that Wav2Lip is comparable to the ground truth, not better. The LMD score also proves high-level lip-synchronization of our method. Overall, the above results prove that FONT generates high-quality talking heads.

**Qualitative Results.** The qualitative comparison results are shown in Fig. 4. All the frames are generated using the same source face and driving audio. It indicates that FONT generates talking heads with natural head motions, accurate mouth shape and identity information. Specifically, Wav2Lip generates fixed faces and blurry mouth areas. Though MakeItTalk and Audio2Head produce head pose changes, they fail to preserve the lip synchronization corresponding to driving audio. PC-AVS can not preserve the identity information of the source face compared to ground truth. AVCT produces obvious visual artifacts in the background area and sometimes fails to produce an accurate mouth shape.

We also show comparison results in large-pose faces, as shown in Fig. 5. Other methods displays wired facial shape change and obvious identity mismatch problem, while FONT generates natural head motions while obtaining high-level image quality. Furthermore, Fig. 6 displays the qualitative results of FONT on the HDTF dataset. It displays the synced video that provides driving audio and generated talking head



**Fig. 7.** Qualitative ablation study results. The red and green rectangles indicate the difference of mouth shape and head motion, respectively.

**Table 3.** Numerical ablation study results.

Method	SSIM $\uparrow$	CPBD $\uparrow$	LMD $\downarrow$	LSE-C $\uparrow$
w/o VAE	0.746	0.160	5.48	7.18
w/o $\mathcal{L}_{SSIM}$	0.738	0.158	6.72	6.79
w/o $D_{sync}$	0.752	0.166	3.46	4.28
Ours	<b>0.789</b>	<b>0.169</b>	<b>2.69</b>	<b>7.22</b>

videos under different source faces. The results indicate that HDTF produces natural head motions while maintaining high-level lip-syn quality. Please see dynamic demos in the supplementary materials for better comparison.

**Ablation Results.** To evaluate the performance of each component in FONT, we conduct the ablation study on several variants: (1) replace the probabilistic VAE-based model with the deterministic LSTM-based model in head pose generation (**w/o VAE**), (2) replace the SSIM loss into the traditional L1 loss (**w/o  $\mathcal{L}_{SSIM}$** ) and (3) remove the lip-sync discriminator from the generator (**w/o  $D_{sync}$** ). The results are shown in Table 3. Given that the SSIM relates to image pixel accuracy, pose accuracy and image quality become worse when removing the above module. As all the variants share basically the same flow-guided generation pattern, most of them achieve similar CPBD scores. The model **w/o  $D_{sync}$**  show a poor LSE-C score indicating bad lip synchronization. The model **w/o VAE** and **w/o  $\mathcal{L}_{SSIM}$**  fail to produce natural head pose, leading to bad LMD score.

Moreover, we show qualitative ablation results in Fig. 7. The red and green rectangles mark the difference between each generated frame. The model **w/o VAE** and **w/o  $\mathcal{L}_{SSIM}$**  fail to produce dynamic natural head motions and tend to produce average still talking heads. Without  $D_{sync}$ , the mouth shape accuracy also decreases, as the red rectangle shows. Overall, we see the contribution of each component in FONT.

## 5. CONCLUSION

In this paper, we present FONT, a flow-guided one-shot model that generates talking heads with natural head motions. The head pose sequence is first predicted by a well-designed probabilistic VAE-based model. After getting the driving pose sequence, we utilize self-supervised keypoints to predict motion flow as face structure representation from the source face and driving audio. Finally, the occlusion-aware flow-guided generator produces talking heads. Both quantitative and qualitative experiments demonstrate that we obtain talking heads with natural poses and high-level lip-sync quality compared with other methods.

For ethical considerations, FONT is intended for the video editing industry and focuses on world-positive use cases and applications. We believe the proper usage of this technique will enhance the development of artificial intelligence research and relevant multimedia applications. To ensure proper use, we will release our codes and contribute to deep-fake detection research.

## 6. REFERENCES

- [1] Zhimeng Zhang, Lincheng Li, Yu Ding, and Changjie Fan, “Flow-guided one-shot talking face generation with a high-resolution audio-visual dataset,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 3661–3670.
- [2] 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*, 2019, pp. 7832–7841.
- [3] KR Prajwal, Rudrabha Mukhopadhyay, Vinay P Nambodiri, and CV Jawahar, “A lip sync expert is all you need for speech to lip generation in the wild,” in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 484–492.
- [4] Hang Zhou, Yasheng Sun, Wayne Wu, Chen Change Loy, Xiaogang Wang, and Ziwei Liu, “Pose-controllable talking face generation by implicitly modularized audio-visual representation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 4176–4186.
- [5] Fei Yin, Yong Zhang, Xiaodong Cun, Mingdeng Cao, Yanbo Fan, Xuan Wang, Qingyan Bai, Baoyuan Wu, Jue Wang, and Yujiu Yang, “Styleheat: One-shot high-resolution editable talking face generation via pretrained stylegan,” *arXiv preprint arXiv:2203.04036*, 2022.
- [6] Lele Chen, Guofeng Cui, Celong Liu, Zhong Li, Ziyi Kou, Yi Xu, and Chenliang Xu, “Talking-head generation with rhythmic head motion,” in *European Conference on Computer Vision*. Springer, 2020, pp. 35–51.
- [7] Suzhen Wang, Lincheng Li, Yu Ding, Changjie Fan, and Xin Yu, “Audio2head: Audio-driven one-shot talking-head generation with natural head motion,” *arXiv preprint arXiv:2107.09293*, 2021.
- [8] Suzhen Wang, Lincheng Li, Yu Ding, and Xin Yu, “One-shot talking face generation from single-speaker audio-visual correlation learning,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2022, vol. 36, pp. 2531–2539.
- [9] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe, “First order motion model for image animation,” *Advances in Neural Information Processing Systems*, vol. 32, pp. 7137–7147, 2019.
- [10] Lele Chen, Guofeng Cui, Ziyi Kou, Haitian Zheng, and Chenliang Xu, “What comprises a good talking-head video generation?: A survey and benchmark,” *arXiv preprint arXiv:2005.03201*, 2020.
- [11] Joon Son Chung, Amir Jamaludin, and Andrew Zisserman, “You said that?,” *arXiv preprint arXiv:1705.02966*, 2017.
- [12] Hang Zhou, Yu Liu, Ziwei Liu, Ping Luo, and Xiaogang Wang, “Talking face generation by adversarially disentangled audio-visual representation,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019, vol. 33, pp. 9299–9306.
- [13] Yang Zhou, Xintong Han, Eli Shechtman, Jose Echevarria, Evangelos Kalogerakis, and Dingzeyu Li, “Makelttalk: speaker-aware talking-head animation,” *ACM Transactions on Graphics (TOG)*, vol. 39, no. 6, pp. 1–15, 2020.
- [14] Yu Deng, Jiaolong Yang, Sicheng Xu, Dong Chen, Yunde Jia, and Xin Tong, “Accurate 3d face reconstruction with weakly-supervised learning: From single image to image set,” in *IEEE Computer Vision and Pattern Recognition Workshops*, 2019.
- [15] Carl Doersch, “Tutorial on variational autoencoders,” *arXiv preprint arXiv:1606.05908*, 2016.
- [16] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli, “Image quality assessment: from error visibility to structural similarity,” *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004.
- [17] Joon Son Chung and Andrew Zisserman, “Lip reading in the wild,” in *Asian conference on computer vision*. Springer, 2016, pp. 87–103.
- [18] Niranjan D Narvekar and Lina J Karam, “A no-reference perceptual image sharpness metric based on a cumulative probability of blur detection,” in *2009 International Workshop on Quality of Multimedia Experience*. IEEE, 2009, pp. 87–91.
- [19] Joon Son Chung and Andrew Zisserman, “Out of time: automated lip sync in the wild,” in *Asian conference on computer vision*. Springer, 2016, pp. 251–263.