

Progressive Disentangled Representation Learning for Fine-Grained Controllable Talking Head Synthesis

Duomin Wang Yu Deng Zixin Yin Heung-Yeung Shum Baoyuan Wang
Xiaobing.AI

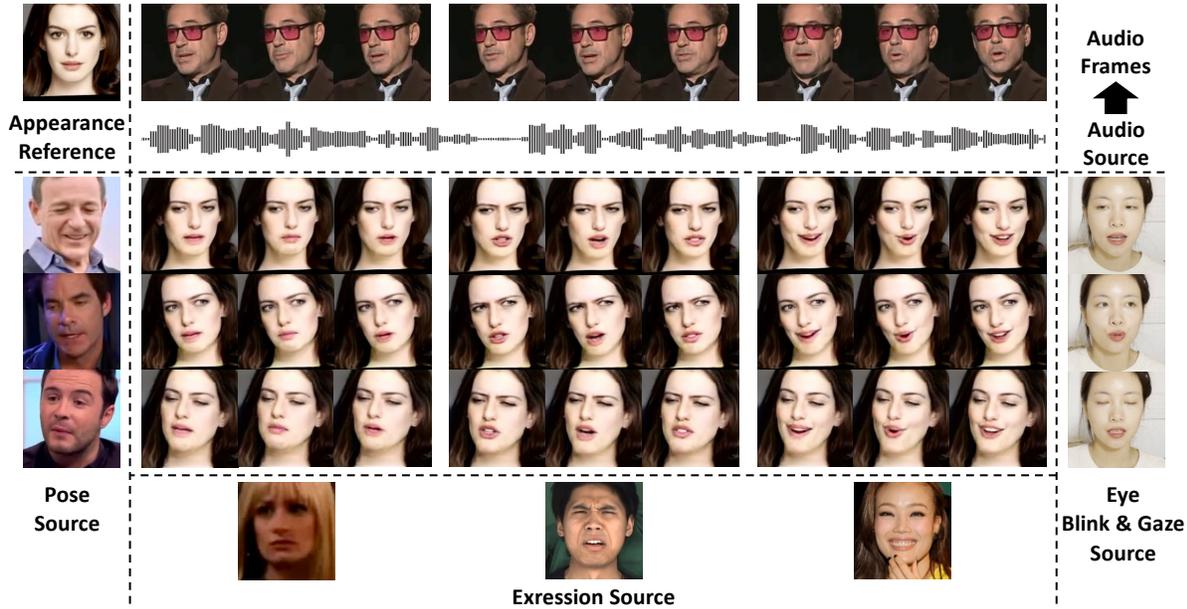


Figure 1. Our method takes an appearance reference as input and generates its talking head with disentangled control over lip motion, head pose, eye gaze&blink, and emotional expression, where the driving signal of lip motion comes from speech audio, and all other motions are controlled by different videos. As shown, it well disentangles all motion factors and achieves precise control over individual motion.

Abstract

We present a novel one-shot talking head synthesis method that achieves disentangled and fine-grained control over lip motion, eye gaze&blink, head pose, and emotional expression. We represent different motions via disentangled latent representations and leverage an image generator to synthesize talking heads from them. To effectively disentangle each motion factor, we propose a progressive disentangled representation learning strategy by separating the factors in a coarse-to-fine manner, where we first extract unified motion feature from the driving signal, and then isolate each fine-grained motion from the unified feature. We introduce motion-specific contrastive learning and regressing for non-emotional motions, and feature-level decorrelation and self-reconstruction for emotional expression, to fully utilize the inherent properties of each motion

factor in unstructured video data to achieve disentanglement. Experiments show that our method provides high quality speech&lip-motion synchronization along with precise and disentangled control over multiple extra facial motions, which can hardly be achieved by previous methods. Project website: <https://dorniwang.github.io/PD-FGC/>

1. Introduction

Talking head synthesis is an indispensable task for creating realistic video avatars and enables multiple applications such as visual dubbing, interactive live streaming, and online meeting. In recent years, researchers have made great progress in one-shot generation of vivid talking heads by leveraging deep learning techniques. Corresponding methods can be mainly divided into audio-driven talking

head synthesis and video-driven face reenactment. Audio-driven methods focus more on accurate lip motion synthesis from audio signals [9, 46, 54, 56]. Video-driven approaches [52, 61] aim to faithfully transfer all facial motions in the source video to target identities and usually treat these motions as a unity without individual control.

We argue that a fine-grained and disentangled control over multiple facial motions is the key to achieving lifelike talking heads, where we can separately control lip motion, head pose, eye motion, and expression, given corresponding respective driving signals. This is not only meaningful from the research aspect which is often known as the disentangled representation learning but also has a great impact on practical applications. Imagining in a real scenario, where we would like to modify the eye gaze of an already synthesized talking head, it could be costly if we cannot solely change it but instead ask an actor to perform a completely new driving motion. Nevertheless, controlling all these factors in a disentangled manner is very challenging. For example, lip motions are highly tangled with emotions by nature, whereas the mouth movement of the same speech can be different under different emotions. There are also insufficient annotated data for large-scale supervised learning to disentangle all these factors. As a result, existing methods either cannot modify certain factors such as eye gaze or expression, or can only change them altogether, or have difficulties providing precise control over individual factors.

In this paper, we propose **Progressive Disentangled Fine-Grained Controllable Talking Head (PD-FGC)** for one-shot talking head generation with disentangled control over lip motion, head pose, eye gaze&blink, and emotional expression¹, where the control signal of lip motion comes from audios, and all other motions can be individually driven by different videos. To this end, our intuition is to learn disentangled latent representation for each motion factor, and leverage an image generator to synthesize talking heads taking these latent representations as input. However, it is very challenging to disentangle all these factors given only in-the-wild video data for training. Therefore, we propose to fully utilize the inherent properties of each motion within the video data with little help of existing prior models. We design a progressive disentangled representation learning strategy to separate each factor control in a coarse-to-fine manner based on their individual properties. It consists of three stages:

1) *Appearance and Motion Disentanglement*. We first learn appearance and motion disentanglement via data augmentation and self-driving [6, 77] to obtain a unified motion feature that records all motions of the driving frame meanwhile excludes appearance information. It serves as a strong starting point for further fine-grained disentanglement.

¹We define the emotional expression as the facial expression that excludes speech-related mouth movement and eye gaze&blink.

2) *Fine-Grained Motion Disentanglement*. Given the unified motion feature, we learn individual motion representation for lip motion, eye gaze&blink, and head pose, via a carefully designed motion-specific contrastive learning scheme as well as the guidance of a 3D pose estimator [15]. Intuitively, speech-only lip motion can be well separated via learning shared information between the unified motion feature and the corresponding audio signal [77]; eye motions can be disentangled by region-level contrastive learning that focuses on eye region only, and head pose can be well defined by 3D rigid transformation.

3) *Expression Disentanglement*. Finally, we turn to the challenging expression separation as the emotional expression is often highly tangled with other motions such as mouth movement. We achieve expression disentanglement via decorrelating it with other motion factors on a feature level, which we find works incredibly well. An image generator is simultaneously learned for self-reconstruction of the driving signals to learn the semantically-meaningful expression representation in a complementary manner.

In summary, our contributions are as follows: **1)** We propose a novel one-shot and fine-grained controllable talking head synthesis method that disentangles appearance, lip motion, head pose, eye blink&gaze, and emotional expression, by leveraging a carefully designed progressive disentangled representation learning strategy. **2)** Motion-specific contrastive learning and feature-level decorrelation are introduced to achieve desired factor disentanglement. **3)** Trained on unstructured video data with limited guidance from prior models, our method can precisely control diverse facial motions given different driving signals, which can hardly be achieved by previous methods.

2. Related work

Audio-driven talking head synthesis. Audio-driven talking head synthesis [3, 4] aims to generate portrait images with synchronized lip motions to the given speech audios. The majority of works [9, 39, 46, 54, 56, 79] focus on controlling only the mouth region and leave other parts unchanged. Some recent works enable control over more facial properties such as eye blink and head pose [8, 57, 70, 74, 75, 78]. More recently, several methods [28, 29, 34, 59] try to introduce emotional expression variations into the synthesis process as it is a crucial property for vivid talking head generation. However, integrating expression control into talking-head synthesis is very challenging due to the lack of expressive data. Some methods [28, 29, 60] build on manually collected emotional talking head dataset [60], yet they cannot well generalize to large-scale scenarios due to the limited data coverage. A recent work GC-AVT [34] leverages in-the-wild data for expressive talking head synthesis. They achieve disentangled control over expression by introducing mouth-region data augmentation to separate lip motion and

other facial expressions. Different from them, we leverage a feature-level decorrelation to disentangle the two factors. Moreover, our method can synthesize arbitrary talking head with a disentangled control over lip motion, head pose, eye gaze&blink, and expressions, while previous methods cannot achieve individual control over all these factors.

Video-driven face reenactment. Video-driven face reenactment targets faithful facial motion transfer between a driving video and a target image. The literature can be mainly divided into warping-based methods [20, 23, 27, 42, 48, 52, 61, 62, 68, 73] and synthesis-based approaches [6, 32, 37, 47, 55, 58, 64, 66]. The warping-based methods predict warping flows between the source and target frames to transform target images or their extracted features to align with source motions. The synthesis-based methods instead learn intermediate representations from input images and directly send them to a generator for image synthesis. The representations can be landmarks [37, 55, 66], 3D face model parameters or meshes [7, 32, 48, 58], or latent features extracted from images [2, 6, 62]. Some recent methods [35, 71] also exploit prior knowledge from a pre-trained 2D GAN [31] for animating face images. Our proposed method also builds on a synthesis-based approach, where we learn disentangled latent representations for multiple facial motions by our designed progressive disentangled representation learning strategy. In addition, different from video-driven approaches, we control the lip motion via audio signals.

Disentangled representation learning on the face. Disentangled representation learning for faces is a longstanding task and has been widely explored in the literature. Plenty of works [11, 12, 18, 26, 33, 36, 41, 63] focus on unsupervised representation learning, where InfoGAN [12] and β -VAE [26] are two representative works. However, these unsupervised methods cannot guarantee meaningful latent representations well aligned with human perceptions [38]. More recently, several methods [49–51, 65, 76, 80] explore latent space editing of a pre-trained 2D GAN [31] with the help of certain classifiers to achieve disentangled control over desired facial properties. Nevertheless, their controllability is often confined by the linear classifiers and the data distribution of the pre-trained generator. Some methods [16, 19, 22, 47, 48, 67] leverage more powerful prior knowledge such as 3D face model [45] or expression model [47] to guide the representation learning, and develop specific training schemes [19, 67] on structured data to achieve desired factor disentanglement. We achieve disentangled representation learning via a carefully designed progressive training scheme on videos and introduce certain prior models [15] to help with accurate factor control.

3. Method

Given an image of an arbitrary person, our goal is to synthesize a talking-head video of it, where we can separately control different facial motions in each frame, including lip motion, head pose, eye gaze&blink, and emotional expression. We expect the lip motion to be derived from an audio clip and other motions from different respective driving videos. To this end, we propose to represent all controllable facial motions along with their appearance by disentangling latent representations of the input visual and audio signals and learning a corresponding image generator to synthesize desired talking heads from them. We introduce a progressive disentangled representation learning scheme to learn the latent representations in a coarse-to-fine manner, as shown in Fig. 2. We first disentangle the appearance with the facial motions to obtain a unified motion representation that records all motion information (Sec. 3.1). Then, we isolate each fine-grained facial motion, except the expression, from the unified motion feature via motion-specific contrastive learning (Sec. 3.2). Finally, we separate the expression from other motions via feature-level decorrelation, and simultaneously learn the image generator for fine-grained controllable talking head synthesis (Sec. 3.3).

3.1. Appearance and Motion Disentanglement

We argue that to achieve disentanglement over multiple fine-grained motion factors, a primary thing to do is to learn a unified motion representation that records all kinds of motion information meanwhile excludes appearance (*i.e.* identity) information. Such a unified motion feature serves as a strong starting point for further fine-grained factor disentanglement from it. To this end, we follow [6] to disentangle appearance and facial motions.

Specifically, an appearance encoder E_{app} and a motion encoder E_{mot} are introduced to extract corresponding features from an appearance image and a driving frame, respectively. An extra generator G_0 is applied to synthesize a face image with the identity of the appearance image and the facial motion of the driving frame. Self-driving and reconstruction are applied to learn the whole pipeline, where data augmentation is introduced to the motion branch to force the motion encoder to neglect appearance variations and only focus on motion extraction. To further improve the accuracy of extracted motion feature, we introduce a motion reconstruction loss on top of the training losses in [6]:

$$\mathcal{L}_{mot} = \|\phi(I_0) - \phi(I_g)\|_2 + \|\psi(I_0) - \psi(I_g)\|_2, \quad (1)$$

where $\phi(\cdot)$ and $\psi(\cdot)$ are features extracted by the 3D face reconstruction network and the emotion network of [15], I_0 is the image synthesized by G_0 given appearance and motion features, and I_g is the ground truth image. The above training scheme helps us to learn a unified motion feature

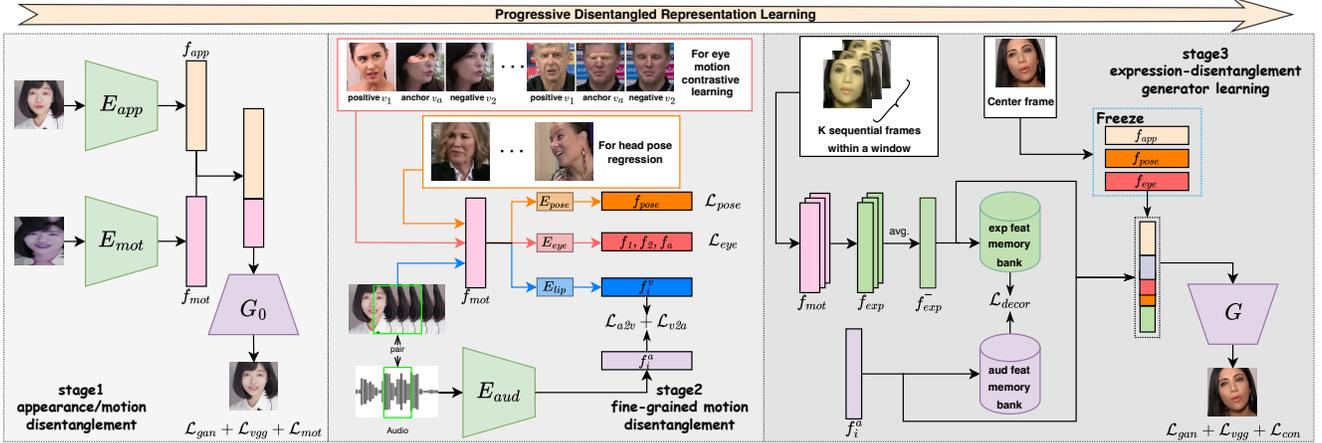


Figure 2. The overview of our method. We achieve factor disentanglement for different facial motions via a progressive disentangled representation learning strategy. We first disentangle appearance with all facial motions to obtain a unified motion feature for further fine-grained disentanglement. Then, we separate each fine-grained motion feature from the unified motion feature via motion-specific contrastive learning and the help of a 3D prior model. Finally, we disentangle expression with other motions by feature-level decorrelation and simultaneously learn an image generator for controllable talking head synthesis.

that faithfully represents all facial movements, which further helps us to achieve fine-grained motion disentanglement to be described in the following sections.

3.2. Fine-Grained Motion Disentanglement

Based on the unified motion feature from the previous stage, we introduce three extra encoders to further extract fine-grained motion features from it, including lip motion feature, eye gaze&blink feature, and head pose feature. The key intuition is to design motion-specific contrastive learning based on the unique property of each individual motion or to leverage the guidance of a prior model if a motion can be well described by it. We do not separate expression in this stage as it can be highly tangled with other factors. We leave its disentanglement in the final stage (Sec. 3.3) where all other factors are effectively separated.

Lip motion contrastive learning. Lip motions can be well separated from other motions by exploring shared information between the unified motion feature and the corresponding speech audio, as shown by previous method [77]. Therefore, we follow [77] to learn the lip motion feature with audio-visual contrastive learning. Given a set of video frames $\{v_i\}$ and their corresponding audio signals $\{a_i\}$, we introduce a lip motion encoder E_{lip} and an audio encoder E_{aud} , and extract lip motion features $\{f_i^v\} = \{E_{lip} \circ E_{mot}(v_i)\}$ and audio features $\{f_i^a\} = \{E_{aud}(a_i)\}$ via the two networks, where E_{mot} is the pre-trained motion encoder from the previous stage. We then construct a positive audio-video pair (f_i^a, f_i^v) and K negative audio-video pairs $(f_i^a, f_k^v), k \neq i$ for each sampled audio feature f_i^a , and vice versa. We enforce the InfoNCE loss [43] following [77] to maximize the similarity between positive pairs

and minimize the similarity between negative pairs:

$$\mathcal{L}_{a2v} = -\log\left[\frac{\exp(\mathcal{S}(f_i^a, f_i^v))}{\exp(\mathcal{S}(f_i^a, f_i^v)) + \sum_{k=1}^K \exp(\mathcal{S}(f_i^a, f_k^v))}\right], \quad (2)$$

$$\mathcal{L}_{v2a} = -\log\left[\frac{\exp(\mathcal{S}(f_i^v, f_i^a))}{\exp(\mathcal{S}(f_i^v, f_i^a)) + \sum_{k=1}^K \exp(\mathcal{S}(f_i^v, f_k^a))}\right], \quad (3)$$

where $\mathcal{S}(\cdot, \cdot)$ is the cosine similarity. This loss ensures that lip motion features predicted by E_{lip} and E_{aud} are close to each other for corresponding video frames and audio. Since the audio signal merely contains lip motion information, it helps with better factor disentanglement. Moreover, we can leverage the audio encoder for audio-driven lip motion synthesis for our controllable talking head.

Eye motion contrastive learning. Eye motions, including eye gaze and blink, are local movements and have limited influence on other facial regions. Therefore, if we substitute the eye region of a person with that of another person to composite a new image, the extracted eye motion feature from it should be identical to that of the latter person. Based on this observation, we design a dedicated contrastive learning scheme to disentangle the eye motions.

Specifically, given two driving frames, namely v_1 and v_2 , we construct an anchor frame v_a by compositing the eye region² of v_1 and other regions of v_2 , as shown in Fig. 2. We introduce an extra encoder E_{eye} to extract eye motion features f_1, f_2 , and f_a from the corresponding unified motion feature of the above frames, and construct a positive pair (f_1, f_a) and a negative pair (f_2, f_a) . We then enforce a

²We use an off-the-shelf method [72] to detect eye landmarks, and warp the eye region of v_1 to align with that of v_2 .

similar InfoNCE loss between these pairs to learn E_{eye} :

$$\mathcal{L}_{eye} = -\log\left[\frac{\exp(\mathcal{S}(f_1, f_a))}{\exp(\mathcal{S}(f_1, f_a)) + \exp(\mathcal{S}(f_2, f_a))}\right]. \quad (4)$$

It helps the eye motion encoder to only focus on eye region motions and neglect the variance of other regions.

Head pose learning. Since head pose can be well defined by a 6D parameter consisting of three Euler angles (*i.e.* pitch, yaw, roll) and 3D translations, we propose to directly regress them via a head pose encoder E_{pose} with the guidance of a 3D face prior model:

$$\mathcal{L}_{pose} = |\mathcal{P}_{pred} - \mathcal{P}_{gt}|_1, \quad (5)$$

where \mathcal{P}_{pred} is the predicted pose parameter by E_{pose} , and \mathcal{P}_{gt} is the ground truth pose parameter obtained by the off-the-shelf 3D face reconstruction model [15].

3.3. Expression Disentanglement

The challenge for expression disentanglement is two-fold. On one hand, the emotional expression can be highly tangled with other motions (*e.g.* mouth movements can be different under different emotions even if the speech contents are identical), which makes it difficult to design motion-specific contrastive learning as done previously. On the other hand, existing expression estimators [15, 17, 40] usually include other motion information in their expression representation, which cannot provide accurate guidance for the expression disentanglement in our scenario. To tackle these challenges, we propose a feature-level decorrelation strategy to disentangle the expression with other motions, along with a self-reconstruction of the driving frame to learn precise expression representation in a complementary manner. The hypothesis behind this is that if an extracted expression feature is independent of the features of other motions, meanwhile its combination with the others can still faithfully reconstruct all facial motions in the driving signal, then it is a precise latent representation of the ground truth expression. We describe the learning strategies in detail below.

In-window decorrelation. We observe that the expression variation in a video sequence is usually less frequent than the changes in other motions. Therefore, if we take the average expression feature within a time window, the other motion information stored in certain dimension of the expression feature should be averaged out, leading to a clean expression feature uncorrelated with other motions. Therefore, given a driving frame, we define a window of size K around it and augment the frames within the window with random rotation, scaling, and color jittering. We then extract the expression features from their corresponding unified motion features via an expression encoder E_{exp} ,

and calculate their average feature as the expression feature for the center driving frame. The average feature will be then sent into a generator G for image synthesis and self-reconstruction described in the following paragraph³.

Lip-motion decorrelation. We further introduce a lip motion decorrelation loss to achieve better expression disentanglement by forcing independence between the expression feature and the lip motion feature (*i.e.* audio feature):

$$\mathcal{L}_{decor} = \frac{1}{D} \sum_{B,D} cor(\bar{F}^e, F^a)^2, \quad (6)$$

where $\bar{F}^e \in \mathbb{R}^{B \times D}$ is a matrix consisting of average expression features within a batch of size B , $F^a \in \mathbb{R}^{B \times D}$ is the corresponding audio feature matrix, D is the feature dimension, and $cor(\cdot, \cdot)$ calculates the feature dimension correlation between the two matrices. In practice, computing the correlation between two variables requires a large batchsize to reach enough accuracy. However, it is difficult to maintain such a large batchsize during training due to memory limitation. To tackle this problem, we maintain two memory banks for the expression feature and the audio feature to compute the correlation, instead of using only the current batch of features. The memory bank always keeps M latest features inside to compute Eq. (6) during training, where M is much larger than the batchsize of each iteration. The gradient will only back-propagate through the current batch of features to update the network weights.

Complementary learning via self-reconstruction. The above two decorrelation strategies ensure feature independence between expression and other motions, yet the extracted expression feature still lacks semantic meaning. Therefore, we leverage an image generator G to take the expression feature along with the features of appearance and other motions as input, and synthesize an image with desired facial motions via self-reconstruction of the driving frame, as shown in Fig. 2. In order to faithfully reconstruct the driving frame, the expression encoder is forced to learn complementary information that is not included in all other motion features, which is exactly the expression information. We enforce multiple losses to learn the expression encoder E_{exp} and the image generator G :

$$\mathcal{L}_{vgg} = \sum_{i=1}^N \|VGG_i(I_f) - VGG_i(I_g)\|_1, \quad (7)$$

where $VGG_i(\cdot)$ is the feature map of the i 's layer in a pre-trained VGG19 [53] network. We also adopt the adversarial loss and the discriminator feature matching loss following [6] to improve the synthesized image quality.

³During inference, we use the expression feature of each frame instead of the average one, as we find the expression encoder learned following our training strategy can well disentangle expression with other motions.

In addition, to ensure that the synthesized image well follows all facial motions of the driving frame, we further introduce a motion-level consistency loss:

$$\mathcal{L}_{con} = \exp(-\mathcal{S}(\mathbf{V}_{lip}(I_f), E_{aud}(a_g))) + \|\mathcal{G}(I_f), \mathcal{G}(I_g)\|_1 + \mathcal{L}_{mot}, \quad (8)$$

where E_{aud} is our audio encoder learned in the previous stage, \mathbf{V}_{lip} is a pre-trained encoder to extract lip motion features from images, $\mathcal{G}(I)$ is a gaze estimator [1], and $\mathcal{S}(\cdot, \cdot)$ is the cosine similarity; I_f , I_g and a_g are our synthesized image, ground truth image, and audio, respectively; \mathcal{L}_{mot} is the motion reconstruction loss defined in Eq. (1).

The above self-reconstruction process, together with the feature-level decorrelation strategy, helps to disentangle the expression feature from the unified motion feature. Moreover, the image generator G learned in this step naturally achieves disentangled and controllable talking head synthesis with all disentangled motion features and the appearance feature as input. We, therefore, take G as our final image generator for talking head synthesis.

4. Experiments

Implementation details. We train our model on VoxCeleb2 [13] dataset and evaluate it on both VoxCeleb2 and Mead [60] dataset. All video frames are aligned following the official annotations [13] and resized to 224×224 . Corresponding audios are extracted from the original videos and converted to Mel-spectrograms. Our appearance encoder E_{app} is implemented as a ResNet50 [24], and the motion encoder E_{mot} takes the same structure as [5]. The audio encoder E_{aud} adopts a ResNetSE34 [30] structure. The encoder for each fine-grained motion factor, including lip motion, head pose, eye gaze&blink, and expression, is implemented as an MLP with ReLU activations. The image generator G_0 in the first stage and the final image generator G are both based on StyleGAN2 [31]. The dimensions of the appearance feature and the unified motion feature are set to 2, 048 and 512, respectively. The dimensions of each fine-grained motion features are 500, 6, 6, and 30 for lip motion (audio), head pose, eye gaze&blink, and expression, respectively. We implement our framework using PyTorch [44], and train it on 8 Tesla V100 GPUs with 32GB memory, using a batchsize of 16 for 50 epochs. See the supplementary materials for more details.

Baselines. We compare our method with existing talking head synthesis methods: **Wav2Lip** [46] that allows only mouth region control; **MakeItTalk** [78] that further introduces random eye blinks and audio-aware head poses; **PC-AVS** [77] with controllable head pose and lip motion; and **EAMM** [28] with disentangled control over lip motion, head pose, and expression. A recent GC-AVT [34] also

achieves expressive talking head synthesis. We do not compare with it since its code is unavailable yet.

4.1. Quantitative Evaluation

We evaluate the image generation quality as well as factor control accuracy of different methods on a self-driving setting, where we use the first frame in a test video to provide appearance and use the following audio and video frames to provide lip motion and other motions, respectively. We use the Fréchet Inception Distances (**FID**) [25] between the synthesized images and the ground truth driving frames to evaluate the image quality. For the accuracy of motion control, we leverage several metrics. We first calculate the facial landmark distance (**LMD**) [10] between the synthesized images and the ground truth to evaluate the overall motion control accuracy. We further calculate the mouth region landmark distance (**LMD_m**) to evaluate the accuracy of lip motion control. We also adopt the Lip Sync Error Confidence (**LSE-C**, also known as Sync_{conf}) [46] to evaluate the lip motion synchronization with the driving audio. Nevertheless, the **LSE-C** value of a method is strongly correlated with its training data, which makes it unfair when comparing methods trained on different data. A recent method [69] also indicates that the **LSE-C** difference between synthesized images of a method and its training data, rather than the absolute value, better reveals lip motion synchronization. Therefore, we propose a normalized confidence score **NLSE-C** for a fair comparison:

$$\text{NLSE-C} = \frac{\text{LSE-C}^{gen} - \text{LSE-C}^{gt}}{\text{LSE-C}^{gt}}, \quad (9)$$

where LSE-C^{gen} is the **LSE-C** score of generated images, and LSE-C^{gt} is the corresponding score of the training data. The new **NLSE-C** measures the relative difference between the generated images and their training data, which better reveals if the synthesized lip motions are reasonable (*i.e.* close to the training data distribution) or not. We find that this new metric better aligns with human perception.

The quantitative results are shown in Tab. 1. Our method yields the best motion control accuracy in terms of **NLSE-C**, **LMD**, and **LMD_m**. We also show competitive image generation quality with other methods. Since Wav2Lip only generates mouth region and copies other regions from input images, we do not evaluate its FID and LMD which can be unfair to other methods.

We further compare our method with the others on expression and head pose control accuracy. For the head pose evaluation, we follow the same self-driving setting as described above, and use a 3D face reconstructor [17] to extract head pose parameters from the synthesized images and calculate their difference (**MSE**) with the ground truth. For the expression evaluation, we find the self-driving setting cannot well evaluate expression controllability as the

Table 1. Quantitative comparison for audio-driven talking head synthesis on VoxCeleb2 [13] and Mead [60]. †: The LSE-C value for the training data of each method is shown in the bracket as a reference.

Method	VoxCeleb2					Mead				
	FID↓	LSE-C†↑	NLSE-C†↓	LMD _m ↓	LMD↓	FID↓	LSE-C†↑	NLSE-C†↓	LMD _m ↓	LMD↓
GT	-	7.35	0	0	0	-	1.76	0	0	0
Wav2Lip [46]	-	9.23	0.183(7.80)	2.54	-	-	9.17	0.176(7.80)	2.54	-
MakeItTalk [78]	19.47	2.03	0.724(7.35)	2.82	6.39	68.35	3.50	0.524(7.35)	2.68	2.56
PC-AVS [77]	14.36	8.21	0.117(7.35)	1.59	2.52	62.85	8.04	0.094(7.35)	2.28	1.92
EAMM Neutral [28]	26.06	4.75	1.699(1.76)	2.29	4.18	50.36	5.34	2.03(1.76)	2.32	2.25
EAMM Emo [28]	27.20	4.55	1.585(1.76)	2.29	4.29	50.49	5.23	1.972(1.76)	2.28	2.27
PD-FGC (Ours)	12.99	7.26	0.012(7.35)	1.15	1.93	73.80	7.24	0.015(7.35)	1.65	1.84

Table 2. Comparison for expression and pose control accuracy.

Method	Expression↓		Pose↓
	VoxCeleb2	Mead	VoxCeleb2
PC-AVS [77]	0.202	0.245	0.0038
EAMM emo [28]	0.196	0.245	0.0196
EAMM neutral [28]	0.192	0.248	0.0203
PD-FGC (Ours)	0.156	0.188	0.0016

Table 3. User study on talking head synthesis.

Method	Lip Sync	Expression	Facial Motion
	Quality↑	Quality↑	Driving Naturalness↑
Wav2Lip [46]	3.50	1.36	1.88
MakeItTalk [78]	1.81	1.89	2.65
PC-AVS [77]	4.46	3.04	3.72
Eamm [28]	1.92	1.77	1.56
PD-FGC (Ours)	4.44	4.38	4.27

appearance reference usually contains similar expressions with the driving frames if they come from the same video clip. Therefore, we conduct a cross-video setting where we use appearance reference and driving frames from different video clips for image synthesis (the driving audio is still from the video clip of the appearance reference). We use a 3D face reconstructor [17] to extract expression parameters in the synthesized images, and compare their difference (MSE) with those of the driving frames to evaluate the control accuracy. As shown in Tab. 2, we achieve the lowest expression and pose control error largely outperforming previous methods. More details are in the supplementary materials.

4.2. Qualitative Evaluation

Fine-grained controllable talking head synthesis. An example of our fine-grained control over synthesized talking head is in Fig. 1. For a given appearance reference, we can control its lip motion, head pose, eye motion, and

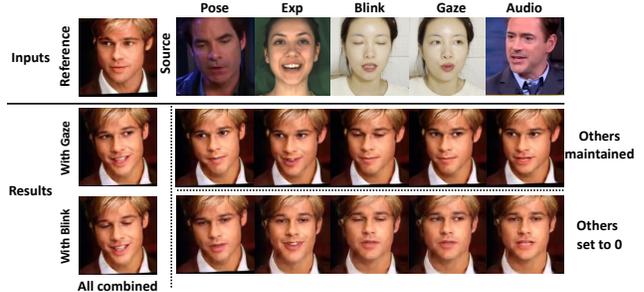


Figure 3. Illustration of our factor disentanglement. **Top:** Controlling one factor while leaving the others unchanged. **bottom:** Controlling one factor and setting the others to zeros.

expression, via different respective audio signals and driving frames, and composite all motions to synthesize a vivid talking head. Our method largely improves the controllability of talking head synthesis upon previous methods, where they cannot achieve separate control over all these factors. More visual results are in the supplementary materials.

Disentangled controllability. We demonstrate factor disentanglement of our method by changing one motion factor at a time given different driving signals in Fig. 3. Our method can independently control the motion of each property to mimic the driving source, and leave all other properties unchanged. Moreover, we can also set all motions to their canonical positions (*i.e.* set features to zero) except the motion to be controlled. These enable our method for diverse downstream applications with different requirements.

Comparison with prior art. We show visual comparisons between our method, PC-AVS [77] and EAMM [28] in Fig. 4. For eye motion and head pose, we adopt the same self-driving setting as in Sec. 4.1. For expression, we use the cross-video setting. We leave the lip motion comparison in the supplementary materials. As shown, our method can well mimic different motions in the driving frames compared to the other methods. PC-AVS [77] can only control head pose besides lip motion, so their synthesized faces



Figure 4. Visual comparison with other methods. The first three columns show the self-driving results. The last column shows the cross-video results. The captions above show the factors that should be focused on in each column.

have different eye motions or expressions compared to the driving frames. It also shows inferior results in head pose control as depicted by the second right column in Fig. 4, due to using an implicit head pose representation instead of the explicit 3D rotation and translation. EAMM [28] cannot well control head pose and eye motions. Moreover, although it can control the expression of a synthesized face, its produced expression is different from the driving source as shown in the last column in Fig. 4.

4.3. User Study

We further conduct user studies for a more comprehensive evaluation. We ask the participants to score from 1 to 5 for the quality of different properties in the synthesized images (5 is the best). The results are in Tab. 3. Our method achieves the best result in expression control quality and facial motion naturalness. And we achieve the second-best result on lip motion synchronization and get very close to the best one (*i.e.* PC-AVS). More details and a user study of factor disentanglement are in the supplementary materials.

4.4. Ablation Study

We conduct an ablation study to validate the efficacy of our proposed feature-level decorrelation in the expression disentanglement stage. We conduct a similar self-driving experiment as in Sec. 4.1, except that we use the first frame

Table 4. Ablation study on expression disentanglement.

Method	Voxceleb2			Voxceleb2	Mead
	LSE-C \uparrow	NLSE-C \downarrow	LMD $_m$ \downarrow	Exp \downarrow	Exp \downarrow
No dis	3.78	0.486	1.81	0.151	0.178
+ In-win	7.60	0.034	1.76	0.159	0.178
+ Decorr	7.02	0.045	1.66	0.157	0.173
All	7.30	0.007	1.27	0.163	0.179

Table 5. Ablation study on window size of the in-window decorrelation strategy.

Size	Voxceleb2			Voxceleb2	Mead
	LSE-C \uparrow	NLSE-C \downarrow	LMD $_m$ \downarrow	Exp \downarrow	Exp \downarrow
7	7.01	0.046	1.81	0.163	0.179
13	7.30	0.007	1.27	0.163	0.179
25	7.23	0.016	1.32	0.164	0.178

in a driving video clip as the expression driving signal instead of using the expression of each frame. The corresponding frames and the audio in the same driving video clip are still used as ground truth to calculate the metrics. In theory, if the expression is well disentangled with other motions, fixing the expression source instead of using the expression in each frame will not influence the lip motion accuracy and should maintain low NLSE-C and LMD $_m$. As shown in Tab. 4, introducing the two decorrelation strategies significantly lowers the quantitative metrics, which indicates better factor disentanglement. And leveraging both of them leads to the best result. We further conduct the cross-video driving experiment similar to Sec. 4.1 to evaluate the expression control accuracy of different alternatives. Our final solution only slightly decreases the expression control accuracy but leads to a large improvement in expression and lip motion disentanglement.

We also study the influence of the window size of the in-window decorrelation in Tab. 5. A window size of 13 yields the best result which is used as our final solution.

5. Conclusion

We presented a fine-grained controllable talking head synthesis method. The core idea is to represent different facial motions via disentangled latent representations. A progressive disentangled representation learning strategy is introduced to separate individual motion factors in a coarse-to-fine manner, by exploring the inherent properties of each factor in unstructured video data. Experiments demonstrated the efficacy of our method on disentangled and fine-grained control of diverse facial motions.

Limitations. Our method mainly focuses on disentangled motion control. The synthesized images may lack fine details and we leave their improvement as future works.

References

- [1] Ahmed A Abdelrahman, Thorsten Hempel, Aly Khalifa, and Ayoub Al-Hamadi. L2cs-net: Fine-grained gaze estimation in unconstrained environments. *arXiv preprint arXiv:2203.03339*, 2022. 6
- [2] Stella Bounareli, Vasileios Argyriou, and Georgios Tzimiropoulos. Finding directions in gan’s latent space for neural face reenactment. *arXiv preprint arXiv:2202.00046*, 2022. 3
- [3] Matthew Brand. Voice puppetry. In *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*, pages 21–28, 1999. 2
- [4] Christoph Bregler, Michele Covell, and Malcolm Slaney. Video rewrite: Driving visual speech with audio. In *Proceedings of the 24th annual conference on Computer graphics and interactive techniques*, pages 353–360, 1997. 2
- [5] Adrian Bulat and Georgios Tzimiropoulos. How far are we from solving the 2d & 3d face alignment problem?(and a dataset of 230,000 3d facial landmarks). In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1021–1030, 2017. 6
- [6] Egor Burkov, Igor Pasechnik, Artur Grigorev, and Victor Lempitsky. Neural head reenactment with latent pose descriptors. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13786–13795, 2020. 2, 3, 5, 13
- [7] Meng Cao, Haozhi Huang, Hao Wang, Xuan Wang, Li Shen, Sheng Wang, Linchao Bao, Zhifeng Li, and Jiebo Luo. Uni-facegan: A unified framework for temporally consistent facial video editing. *IEEE Transactions on Image Processing*, 30:6107–6116, 2021. 3
- [8] 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*, pages 35–51. Springer, 2020. 2
- [9] Lele Chen, Zhiheng Li, Ross K Maddox, Zhiyao Duan, and Chenliang Xu. Lip movements generation at a glance. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 520–535, 2018. 2
- [10] 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 Conference on Computer Vision and Pattern Recognition*, pages 7832–7841, 2019. 6
- [11] Ricky TQ Chen, Xuechen Li, Roger B Grosse, and David K Duvenaud. Isolating sources of disentanglement in variational autoencoders. *Advances in neural information processing systems*, 31, 2018. 3
- [12] Xi Chen, Yan Duan, Rein Houthoofd, John Schulman, Ilya Sutskever, and Pieter Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. *Advances in neural information processing systems*, 29, 2016. 3
- [13] J. S. Chung, A. Nagrani, and A. Zisserman. Voxceleb2: Deep speaker recognition. In *INTERSPEECH*, 2018. 6, 7, 13
- [14] Joon Son Chung and Andrew Zisserman. Out of time: automated lip sync in the wild. In *Asian conference on computer vision*, pages 251–263. Springer, 2016. 14
- [15] Radek Daněček, Michael J Black, and Timo Bolkart. Emoca: Emotion driven monocular face capture and animation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20311–20322, 2022. 2, 3, 5, 14
- [16] Yu Deng, Jiaolong Yang, Dong Chen, Fang Wen, and Xin Tong. Disentangled and controllable face image generation via 3d imitative-contrastive learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5154–5163, 2020. 3
- [17] 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 *Proceedings of IEEE Computer Vision and Pattern Recognition Workshop on Analysis and Modeling of Faces and Gestures*, 2019. 5, 6, 7, 14
- [18] Zheng Ding, Yifan Xu, Weijian Xu, Gaurav Parmar, Yang Yang, Max Welling, and Zhuowen Tu. Guided variational autoencoder for disentanglement learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7920–7929, 2020. 3
- [19] Chris Donahue, Zachary C Lipton, Akshay Balsubramani, and Julian McAuley. Semantically decomposing the latent spaces of generative adversarial networks. In *ICLR*, 2018. 3
- [20] Nikita Drobyshev, Jenya Chelishev, Taras Khakhulin, Aleksei Ivakhnenko, Victor Lempitsky, and Egor Zakharov. Megaportraits: One-shot megapixel neural head avatars. 2022. 3
- [21] Tobias Fischer, Hyung Jin Chang, and Yiannis Demiris. Rtgene: Real-time eye gaze estimation in natural environments. In *Proceedings of the European conference on computer vision (ECCV)*, pages 334–352, 2018. 14
- [22] Partha Ghosh, Pravir Singh Gupta, Roy Uziel, Anurag Ranjan, Michael J Black, and Timo Bolkart. Gif: Generative interpretable faces. In *2020 International Conference on 3D Vision (3DV)*, pages 868–878. IEEE, 2020. 3
- [23] Kuangxiao Gu, Yuqian Zhou, and Thomas Huang. Flnet: Landmark driven fetching and learning network for faithful talking facial animation synthesis. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 10861–10868, 2020. 3
- [24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 6
- [25] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 6
- [26] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. 2016. 3

- [27] Fa-Ting Hong, Longhao Zhang, Li Shen, and Dan Xu. Depth-aware generative adversarial network for talking head video generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3397–3406, 2022. 3
- [28] Xinya Ji, Hang Zhou, Kaisiyuan Wang, Qianyi Wu, Wayne Wu, Feng Xu, and Xun Cao. Eamm: One-shot emotional talking face via audio-based emotion-aware motion model. In *ACM SIGGRAPH 2022 Conference Proceedings*, SIGGRAPH '22, 2022. 2, 6, 7, 8
- [29] Xinya Ji, Hang Zhou, Kaisiyuan Wang, Wayne Wu, Chen Change Loy, Xun Cao, and Feng Xu. Audio-driven emotional video portraits. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [30] Jee-weon Jung, You Jin Kim, Hee-Soo Heo, Bong-Jin Lee, Youngki Kwon, and Joon Son Chung. Pushing the limits of raw waveform speaker recognition. In *Proc. Interspeech*, 2022. 6
- [31] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020. 3, 6
- [32] Hyeonwoo Kim, Pablo Garrido, Ayush Tewari, Weipeng Xu, Justus Thies, Matthias Niessner, Patrick Pérez, Christian Richardt, Michael Zollhöfer, and Christian Theobalt. Deep video portraits. *ACM Transactions on Graphics (TOG)*, 37(4):1–14, 2018. 3
- [33] Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In *International Conference on Machine Learning*, pages 2649–2658. PMLR, 2018. 3
- [34] Borong Liang, Yan Pan, Zhizhi Guo, Hang Zhou, Zhibin Hong, Xiaoguang Han, Junyu Han, Jingtuo Liu, Errui Ding, and Jingdong Wang. Expressive talking head generation with granular audio-visual control. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3387–3396, June 2022. 2, 6
- [35] Connor Z Lin, David B Lindell, Eric R Chan, and Gordon Wetzstein. 3d gan inversion for controllable portrait image animation. *arXiv preprint arXiv:2203.13441*, 2022. 3
- [36] Zinan Lin, Kiran Koshy Thekumparampil, Giulia C Fanti, and Sewoong Oh. Infogan-cr: Disentangling generative adversarial networks with contrastive regularizers. 2019. 3
- [37] Jin Liu, Peng Chen, Tao Liang, Zhaoxing Li, Cai Yu, Shuqiao Zou, Jiao Dai, and Jizhong Han. Li-net: Large-pose identity-preserving face reenactment network. In *2021 IEEE International Conference on Multimedia and Expo (ICME)*, pages 1–6. IEEE, 2021. 3
- [38] Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Raetsch, Sylvain Gelly, Bernhard Schölkopf, and Olivier Bachem. Challenging common assumptions in the unsupervised learning of disentangled representations. In *international conference on machine learning*, pages 4114–4124. PMLR, 2019. 3
- [39] Yuanxun Lu, Jinxiang Chai, and Xun Cao. Live speech portraits: real-time photorealistic talking-head animation. *ACM Transactions on Graphics (TOG)*, 40(6):1–17, 2021. 2
- [40] Ali Mollahosseini, Behzad Hasani, and Mohammad H Mahoor. Affectnet: A database for facial expression, valence, and arousal computing in the wild. *IEEE Transactions on Affective Computing*, 10(1):18–31, 2017. 5
- [41] Thu Nguyen-Phuoc, Chuan Li, Lucas Theis, Christian Richardt, and Yong-Liang Yang. Hologan: Unsupervised learning of 3d representations from natural images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7588–7597, 2019. 3
- [42] Andrew Zisserman Olivia Wiles, A. Sophia Koepke. X2face: A network for controlling face generation by using images, audio, and pose codes. In *ECCV 2018*, 2018. 3
- [43] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 4, 13
- [44] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019. 6
- [45] Pascal Paysan, Reinhard Knothe, Brian Amberg, Sami Romdhani, and Thomas Vetter. A 3d face model for pose and illumination invariant face recognition. In *2009 sixth IEEE international conference on advanced video and signal based surveillance*, pages 296–301. Ieee, 2009. 3
- [46] KR Prajwal, Rudrabha Mukhopadhyay, Vinay P Namboodiri, 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*, pages 484–492, 2020. 2, 6, 7
- [47] Albert Pumarola, Antonio Agudo, Aleix M Martinez, Alberto Sanfeliu, and Francesc Moreno-Noguer. Ganimation: Anatomically-aware facial animation from a single image. In *Proceedings of the European conference on computer vision (ECCV)*, pages 818–833, 2018. 3
- [48] Yurui Ren, Ge Li, Yuanqi Chen, Thomas H Li, and Shan Liu. Pirenderer: Controllable portrait image generation via semantic neural rendering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13759–13768, 2021. 3
- [49] Yujun Shen, Jinjin Gu, Xiaoou Tang, and Bolei Zhou. Interpreting the latent space of gans for semantic face editing. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9243–9252, 2020. 3
- [50] Yujun Shen, Ceyuan Yang, Xiaoou Tang, and Bolei Zhou. Interfacegan: Interpreting the disentangled face representation learned by gans. *IEEE transactions on pattern analysis and machine intelligence*, 2020. 3
- [51] Yujun Shen and Bolei Zhou. Closed-form factorization of latent semantics in gans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1532–1540, 2021. 3
- [52] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. First order motion model for image animation. *Advances in Neural Information Processing Systems*, 32, 2019. 2, 3

- [53] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 5
- [54] Linsen Song, Wayne Wu, Chen Qian, Ran He, and Chen Change Loy. Everybody’s talkin’: Let me talk as you want. *IEEE Transactions on Information Forensics and Security*, 17:585–598, 2022. 2
- [55] Pu Sun, Yuezun Li, Honggang Qi, and Siwei Lyu. Landmarkgan: Synthesizing faces from landmarks. *Pattern Recognition Letters*, 161:90–98, 2022. 3
- [56] Supasorn Suwajanakorn, Steven M Seitz, and Ira Kemelmacher-Shlizerman. Synthesizing obama: learning lip sync from audio. *ACM Transactions on Graphics (TOG)*, 36(4):1–13, 2017. 2
- [57] Justus Thies, Mohamed Elgharib, Ayush Tewari, Christian Theobalt, and Matthias Nießner. Neural voice puppetry: Audio-driven facial reenactment. *ECCV 2020*, 2020. 2
- [58] Justus Thies, Michael Zollhofer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. Face2face: Real-time face capture and reenactment of rgb videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2387–2395, 2016. 3
- [59] Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic. Realistic speech-driven facial animation with gans. *International Journal of Computer Vision*, 128(5):1398–1413, 2020. 2
- [60] Kaisiyuan Wang, Qianyi Wu, Linsen Song, Zhuoqian Yang, Wayne Wu, Chen Qian, Ran He, Yu Qiao, and Chen Change Loy. Mead: A large-scale audio-visual dataset for emotional talking-face generation. In *European Conference on Computer Vision*, pages 700–717. Springer, 2020. 2, 6, 7, 13
- [61] Ting-Chun Wang, Arun Mallya, and Ming-Yu Liu. One-shot free-view neural talking-head synthesis for video conferencing. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10039–10049, 2021. 2, 3
- [62] Yaohui Wang, Di Yang, Francois Bremond, and Antitza Dantcheva. Latent image animator: Learning to animate images via latent space navigation. In *International Conference on Learning Representations*, 2022. 3
- [63] Yuxiang Wei, Yupeng Shi, Xiao Liu, Zhilong Ji, Yuan Gao, Zhongqin Wu, and Wangmeng Zuo. Orthogonal jacobian regularization for unsupervised disentanglement in image generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6721–6730, October 2021. 3
- [64] Wayne Wu, Yunxuan Zhang, Cheng Li, Chen Qian, and Chen Change Loy. Reenactgan: Learning to reenact faces via boundary transfer. In *Proceedings of the European conference on computer vision (ECCV)*, pages 603–619, 2018. 3
- [65] Zongze Wu, Dani Lischinski, and Eli Shechtman. Stylespace analysis: Disentangled controls for stylegan image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12863–12872, 2021. 3
- [66] Sitao Xiang, Yuming Gu, Pengda Xiang, Mingming He, Koki Nagano, Haiwei Chen, and Hao Li. One-shot identity-preserving portrait reenactment. *arXiv preprint arXiv:2004.12452*, 2020. 3
- [67] Taihong Xiao, Jiapeng Hong, and Jinwen Ma. Elegant: Exchanging latent encodings with gan for transferring multiple face attributes. In *Proceedings of the European conference on computer vision (ECCV)*, pages 168–184, 2018. 3
- [68] Guangming Yao, Yi Yuan, Tianjia Shao, Shuang Li, Shanqi Liu, Yong Liu, Mengmeng Wang, and Kun Zhou. One-shot face reenactment using appearance adaptive normalization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 3172–3180, 2021. 3
- [69] Shunyu Yao, RuiZhe Zhong, Yichao Yan, Guangtao Zhai, and Xiaokang Yang. Dfa-nerf: Personalized talking head generation via disentangled face attributes neural rendering. *arXiv preprint arXiv:2201.00791*, 2022. 6
- [70] Ran Yi, Zipeng Ye, Juyong Zhang, Hujun Bao, and Yong-Jin Liu. Audio-driven talking face video generation with learning-based personalized head pose. *arXiv preprint arXiv:2002.10137*, 2020. 2
- [71] 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. 3
- [72] Baosheng Yu and Dacheng Tao. Heatmap regression via randomized rounding. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. 4
- [73] Xianfang Zeng, Yusu Pan, Mengmeng Wang, Jiangning Zhang, and Yong Liu. Realistic face reenactment via self-supervised disentangling of identity and pose. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 12757–12764, 2020. 3
- [74] Chenxu Zhang, Yifan Zhao, Yifei Huang, Ming Zeng, Saifeng Ni, Madhukar Budagavi, and Xiaohu Guo. Facial: Synthesizing dynamic talking face with implicit attribute learning. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3867–3876, 2021. 2
- [75] Jiangning Zhang, Xianfang Zeng, Chao Xu, Yong Liu, and Hongliang Li. Real-time audio-guided multi-face reenactment. *IEEE Signal Processing Letters*, 2021. 2
- [76] Yutong Zheng, Yu-Kai Huang, Ran Tao, Zhiqiang Shen, and Marios Savvides. Unsupervised disentanglement of linear-encoded facial semantics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3917–3926, 2021. 3
- [77] 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*, pages 4176–4186, 2021. 2, 4, 6, 7, 13
- [78] Yang Zhou, Xintong Han, Eli Shechtman, Jose Echevarria, Evangelos Kalogerakis, and Dingzeyu Li. Makeltalk: speaker-aware talking-head animation. *ACM Transactions on Graphics (TOG)*, 39(6):1–15, 2020. 2, 6, 7
- [79] Hao Zhu, Huaibo Huang, Yi Li, Aihua Zheng, and Ran He. Arbitrary talking face generation via attentional audio-visual

coherence learning. *arXiv preprint arXiv:1812.06589*, 2018.

2

- [80] Peihao Zhu, Rameen Abdal, Yipeng Qin, John Femiani, and Peter Wonka. Improved stylegan embedding: Where are the good latents? *arXiv preprint arXiv:2012.09036*, 2020. 3

Supplementary Material

A. More Implementation Details

A.1. Data Preparation

We train our method on all available videos in the training split of VoxCeleb2 [13] dataset. For evaluation, we use the test split of both VoxCeleb2 and Mead [60] dataset. We randomly sample 500 test video clips from VoxCeleb2, and 460 test clips from the Mead following the official setting.

All video frames are aligned following the official annotations [13], and then resized and center-cropped to 224×224 . Corresponding audios are extracted from the original videos by ffmpeg, and then processed with a sample rate of 16,000 and converted to Mel-spectrograms via FFT. The window size, hop size, and the number of Mel bands are set to 1, 280, 160 and 80, respectively.

A.2. More Training Details

Appearance and motion disentanglement. We follow [6] to learn the appearance encoder E_{app} , motion encoder E_{mot} , and the extra image generator G_0 . Different from [6], for the appearance encoder, we send a single appearance reference as input to obtain the appearance latent feature during training, instead of taking the average latent feature of multiple appearance frames in a video clip. Apart from the original training losses proposed in [6], we further introduce a motion reconstruction loss as described in Sec. 3.1 in the main paper (*i.e.* Eq. (1)). We set the initial learning rates for E_{app} , E_{mot} to $5e^{-5}$. The initial learning rates for G_0 and an extra discriminator for computing the adversarial loss in [6] are set to $5e^{-5}$ and $5e^{-6}$, respectively. The learning rates of all networks are decayed by a rate of 0.5 for every 80,000 iterations. We trained the whole pipeline with a batchsize of 24 for 50 epochs on 8 Tesla V100 GPUs with 32GB memory, which took around 2 weeks.

Lip motion disentanglement. We adopt the audio-visual contrastive learning scheme [77] to learn the lip motion encoder E_{lip} and the audio encoder E_{aud} . The two models are trained on audio-video pairs with the InfoNCE loss [43] as described in the main paper (*i.e.* Eq. (2) and (3)). The original training scheme in [77] utilizes frames from the videos different from those deriving the audio signals to construct the negative pairs, which we found can learn non-lip motion information in the obtained lip motion features. Therefore, we only use the unsynchronized frames and audio from the same video clip as the negative pairs during training. We set the initial learning rates of E_{lip} and E_{aud} to $1e^{-5}$, with a decay rate of 0.93 by every 200,000 iterations. We train

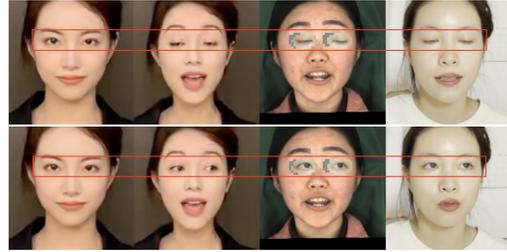


Figure I. Our observation on disentangled eye motion control in the face-reenactment setting in our appearance and motion disentanglement stage. The first column is the appearance reference, the second column is the reenactment result, the third column is the driving source where the eye region comes from the images in the last column. As shown in the figure, the eye motion can be controlled independently without affecting the lip motion in this scenario, which inspires us to design the eye-motion contrastive learning.

the two networks with a batchsize of 32 for 30 epochs. Each item in a batch contains 1 positive pairs and 8 negative pairs. The training took 2 days on 4 Tesla V100 GPUs.

Eye motion disentanglement. The eye motion encoder E_{eye} is learned using our proposed eye-motion contrastive learning described in Sec. 3.2 in the main paper. We describe more details about the motivation behind. Specifically, since our first stage is based on the face reenactment method of [6], we can already synthesize a talking face with the unified motion feature of a driving frame and a given appearance feature via the image generator G_0 . We find that by simply replacing the eye region of the driving frame with a new one bearing different eye blink and gaze, we can achieve a disentangled control of eyes in the synthesized face without affecting other facial motions, as shown in Fig. I. Inspired by this, we formulate the eye-motion contrastive loss in the main paper.

We set the initial learning rate for E_{eye} to $1e^{-5}$, decayed by a rate of 0.5 for every 80,000 iterations. The network is trained with a batchsize of 128 for 30 epochs. The training took 2 days on 4 Tesla V100 GPUs.

Head pose disentanglement. The head pose encoder E_{pose} is learned by regressing the pseudo pose labels as depicted in Sec. 3.2 in the main paper. The learning rate of E_{pose} is also set to $1e^{-5}$ with a decay rate of 0.5 by every 80,000 iterations. The network is trained with a batchsize of 128 for 30 epochs similar to the eye motion encoder. The training took 2 days on 4 Tesla V100 GPUs.

Expression disentanglement. The expression encoder E_{exp} and our final image generator G are learned via our proposed feature-level decorrelation and complementary self-reconstruction in Sec. 3.3 in the main paper. During this stage, all other networks are fixed, including E_{app} , E_{mot} , E_{lip} , E_{aud} , E_{eye} , and E_{pose} . For the in-window decorrelation, we set the window size to 13. For the lip-motion decorrelation, we set the memory bank size to 512 for an accurate estimation of the feature correlation. We set the initial learning rates to $1e^{-5}$ and $2e^{-5}$ for E_{exp} and G , respectively. The learning rate for an extra discriminator to compute the adversarial loss is set to $3.5e^{-6}$. The expression encoder is trained during the first 40,000 iterations and frozen for the following steps. The learning rates for the generator and the discriminator are decayed with a rate of 0.5 by every 80,000 iterations. We use a batchsize of 16 and train all networks for 50 epochs. It took 2 weeks on 8 Tesla V100 GPUs.

A.3. Quantitative Evaluation Details

In Sec. 4.1 in the main paper, we conducted multiple experiments for quantitative metrics calculation (*i.e.* Tab. 1 and 2 in the main paper) under two different settings, namely the *self-driving setting* and the *cross-video setting*.

In the self-driving setting, we use all test clips described in Appendix A.1 for evaluation. We set the first frame in each video as the appearance reference, and drive it using the video frames and the corresponding audio from the same video clip. The audio signals are used to drive the lip motion and the video frames for other motions. Since the source and the target are from the same video, we can easily use the driving frames as the ground truth to evaluate the performance of each method.

In the cross-video setting, we use the first frame from 100 randomly sampled test video clips as an appearance reference and use the first frame from another 100 random test video clip as the driving frame to control all non-lip motions. We still use the audio signals from the video clip of the corresponding appearance frame to control the lip motion. The cross-video setting is designed to evaluate the expression control performance, where we extract the expression parameters of the synthesized videos and their corresponding driving frames using a 3D face reconstructor [15], and compare the expression parameter difference. This helps us to evaluate if a method can precisely transfer the expression from a source to a target. By contrast, in the self-driving setting, since the source and the target are from the same video clip, their expressions are usually the same. Under this circumstance, if a method well mimics the expression motion of the appearance reference, it is difficult to judge whether it successfully transfers the source expression to the target or merely copies the expression from the appearance reference.

A.4. User Study Details

We conduct two user studies to evaluate the motion control performance. In the first experiment, we ask participants to evaluate the accuracy of lip motion synchronization and expression control, as well as the naturalness of all facial motions. We generate 120 videos using 12 random appearance references and 10 random driving clips and randomly select 35 synthesized videos out of them for evaluation. Fifteen participants are asked to score from 1 to 5 for the quality of different properties in the synthesized videos (5 is the best). The corresponding results are in Tab. 3 in the main paper.

In the second experiment, we ask the same group of participants to evaluate the disentanglement controllability of our method. We generate 5 videos using an appearance reference and 3 randomly selected driving videos for the head pose, expression, and eye motion, respectively. In each synthesized video, only one motion factor is controlled by the driving source and all other factors remain unchanged. The participants are asked to score from 1 to 5 for the variation level of each motion in the synthesized videos (5 indicates the largest variation, and 1 means nearly unchanged). The corresponding results are in Tab. VII and discussed in Appendix B.2.

B. More Results

B.1. Fine-Grained Controllable Talking Heads

Figure IV and V show more talking head synthesis results by our method. Our method well mimics the motions from different driving sources and combines them to generate vivid talking heads. **Animations can be found in the accompanying video.**

B.2. Disentangled Controllability

We quantitatively evaluate the disentangled controllability of our method. To this end, we generate talking head images by only varying one motion factor and setting other factors to zeros (*i.e.* canonical positions). We then extract corresponding motion features from the synthesized results and compute the variance of each motion factor in a video clip. Ideally, if different motions are perfectly disentangled, the computed variances will be close to zero for all motions except the one being controlled.

In practice, we use off-the-shelf models to extract each motion feature from our synthesized images. For eye gaze and blink, we use the model of [21]. For expression and pose, we use a 3D face reconstruction model [17]. For lip motion, we use the model of [14]. The variance of each motion factor \mathbf{f} is computed using the following equation:

$$var(\mathbf{f}) = \frac{1}{N} \sum_{i=1}^N \frac{1}{M_i} \sum_{j=1}^{M_i} \|\mathbf{f}_{ij} - \bar{\mathbf{f}}_i\|_2, \quad (I)$$

Table VI. Quantitative evaluation on factor disentanglement of our method. In each row, we compute the variance of a motion feature extracted from the synthesized videos when controlling different individual motion factors.

Variance	Control property				
	lip	pose	blink	gaze	exp
Speech lip motion	11.24	5.16	0.83	0.74	3.76
Head pose	0.0091	0.1597	0.0041	0.0045	0.0088
Eye blink	0.00038	0.00389	0.06657	0.00089	0.00225
Eye gaze	0.089	0.100	0.095	0.105	0.088
Expression	3.07	3.07	2.98	2.93	3.59

Table VII. User study on factor disentanglement of our method.

Variance	Control property				
	lip	pose	blink	gaze	exp
lip	4.7	1.1	1	1	1.1
pose	1.1	4.6	1	1.2	1.3
blink	1.1	1.1	4.1	1.5	1
gaze	1	1.1	1.4	4.4	1
exp	1.3	1.1	1	1	3.7

where f_{ij} is the corresponding extracted motion feature of the j -th frame in the i -th video clip, \bar{f}_i is mean of f_{ij} , N is the number of test videos, and M_i is the length of each video clip.

Table VI shows the computed variance of each motion factor. Each row shows the variance of a single motion factor under different motion control. As shown, the variance of a factor reaches the maximum when the controlling factor is the same with it, and largely decreases when controlled under a different motion factor. This indicates that our method can disentangle different motion controls so that they have a minor influence on each other.

However, the computed variance can still be large in some cases (*e.g.* the left four columns in the last row in Tab. VI). This is due to that the off-the-shelf motion feature extractors are not perfect and can be influenced by variations of other motions when extracting a certain motion feature. Therefore, we refer the readers to the accompanying video to examine the disentanglement ability of our method. We also conduct a user study to better evaluate the factor disentanglement. The results are in Tab. VII (see Appendix A.4 for detailed description). As shown, the variance score is close to 5 when the factor for variance calculation and the factor to be controlled are the same, and close to 1 when they are different, which reveals the disentangled controllability of our method.

B.3. Expression Interpolation

We further investigate the expressive ability of our learned expression feature. We show expression interpolation results by linearly interpolating two expression features

from different expression sources. As shown in Fig. II, our method can smoothly transfer between two different expressions. The synthesized images at interpolated points also have natural expressions. This indicates that our method learns a reasonable expression latent space that supports continuous expression control.

B.4. Comparison with the prior methods

We show the lip motion synthesis comparison in Fig. III. The images are synthesized under the self-driving setting. As depicted, the lip motion generated by our method is natural and closer to the ground truth compare to the alternatives. **See the accompanying video for animations.**

B.5. Ablation Study

Motion reconstruction loss. We further conduct an ablation study to validate the efficacy of our motion reconstruction loss proposed in Sec. 3.1 in the main paper. As shown in Fig. VI, with the motion reconstruction loss, facial motions in the synthesized images contain more details and are closer to the driving sources. By contrast, removing the motion reconstruction loss leads to poor reenactment results for driving sources with rich expressions. As a result, the motion reconstruction loss is important for obtaining an informative unified motion feature to achieve accurate motion control.

C. Ethics Consideration

Our method enables precise and disentangled control over multiple facial motions for vivid talking head gener-

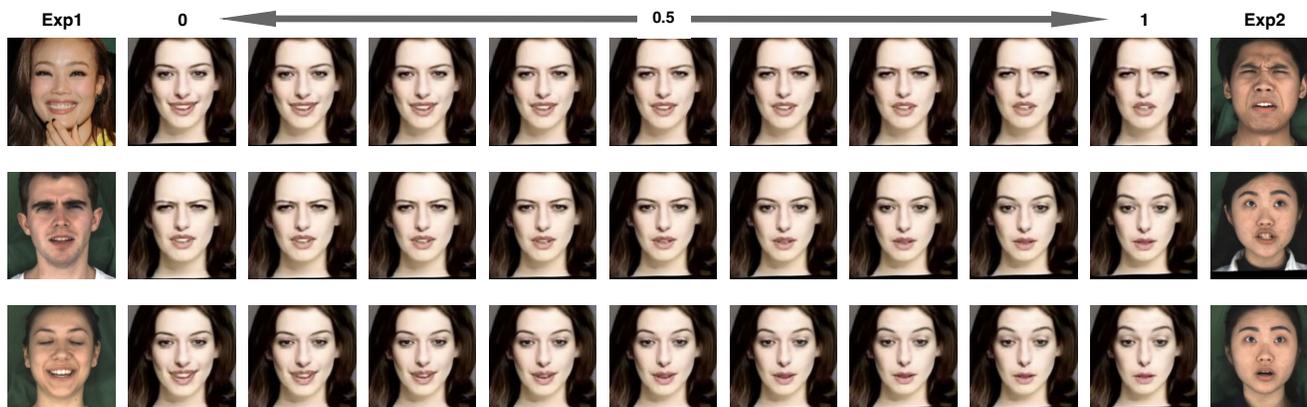


Figure II. Expression interpolation by our method.

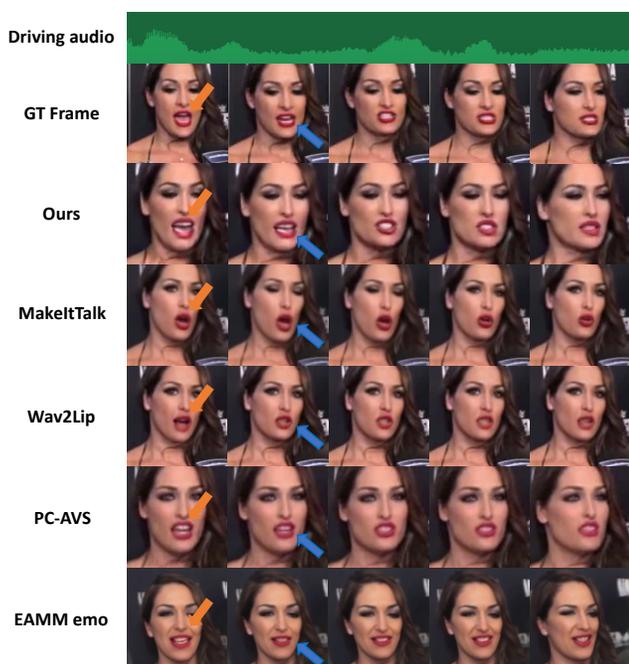


Figure III. Comparison on lip motion control. Images are synthesized under the self-driving setting where the lip motion is driven by the audio signal. Our method yields the best result.

ation. While the major goal of it is to synthesize virtual avatar for applications like live streaming, it can be misused to create deceptive and harmful content of real people. Especially, one may use it to synthesize fake videos of celebrities. We do not condone using our method for generating misleading information that could harm people's reputations. We also suggest investigating advanced forgery detection methods to identify the synthesized fake images and videos to prevent illegal usage.

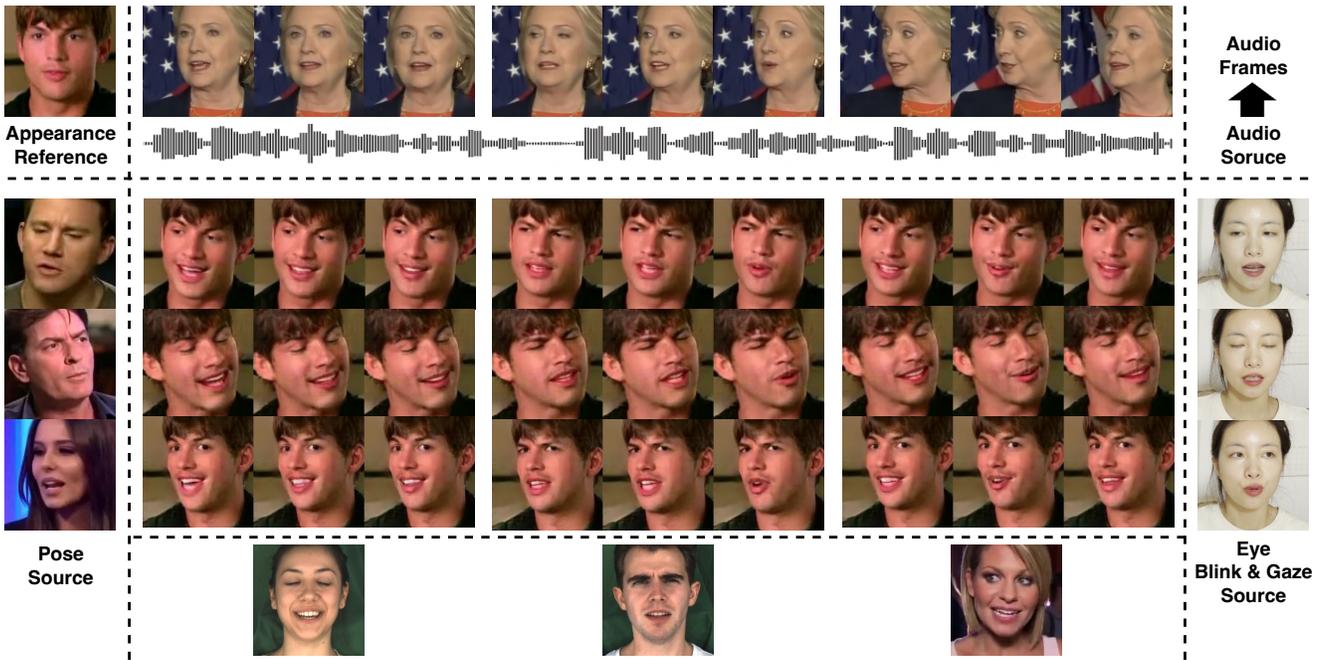


Figure IV. Fine-grained controllable talking heads synthesized by our method.

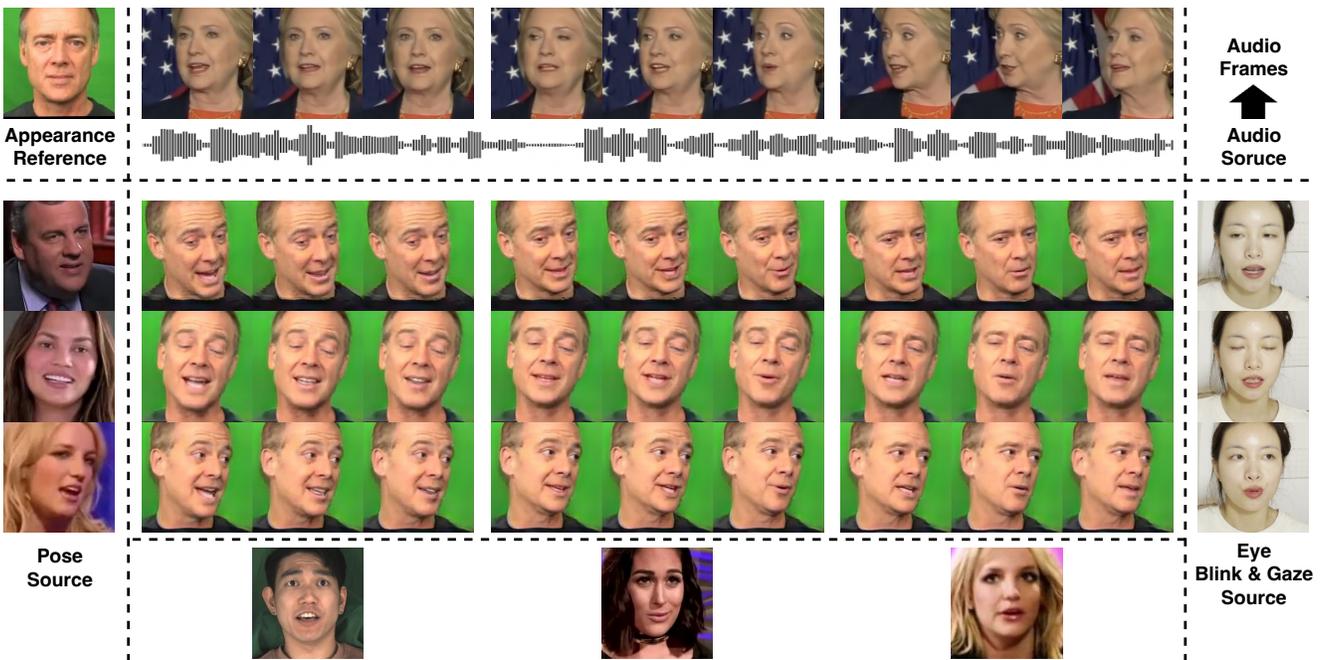


Figure V. Fine-grained controllable talking heads synthesized by our method.

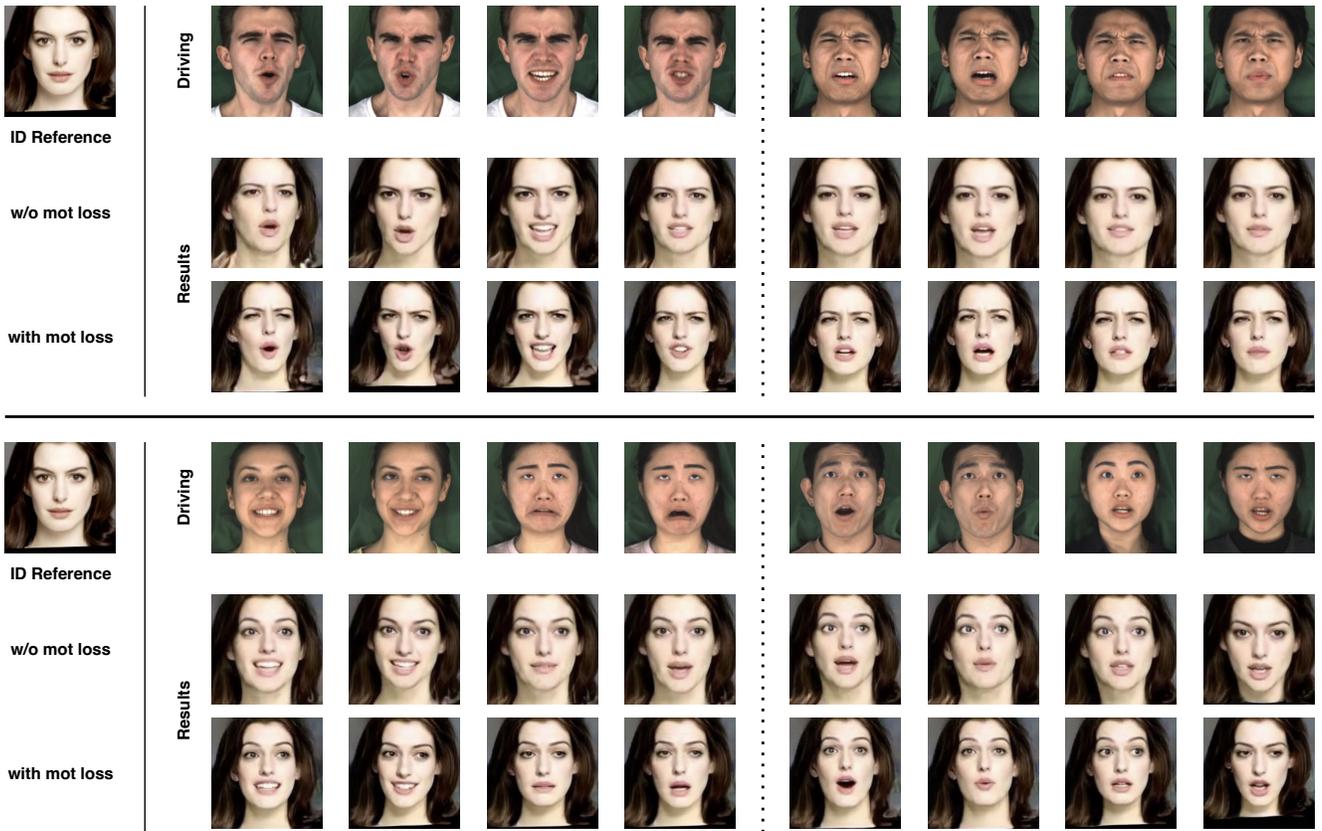


Figure VI. Ablation study on the motion reconstruction loss in the appearance&motion disentanglement learning.