

000  
001  
002054  
055  
056003 

# Emotion-Aware Talking Head Synthesis with Landmark Translation Networks

057  
058  
059004  
005  
006  
007  
008  
009  
010  
011060  
061  
062  
063  
064  
065

Anonymous CVPR 2021 submission

Paper ID 2925

012 

## Abstract

066

We present a new approach for speech-driven talking head generation. Our goal is to synthesize a talking character that matches the given speech audio and conveys a similar emotion. Our formulation explicitly models the emotion of a speaker and generates plausible talking heads with high visual quality, temporal smoothness, and strong emotional expressiveness. We first extract the speaker’s emotion features from the speech audio clip using an encoder network and generate 3D facial landmarks with lip motions synchronized with the speech audio. Next we modify these 3D landmarks to express corresponding face emotions through edge deformations. Finally, we translate the deformed landmarks sequence into a realistic talking head video by transferring the facial components. We generate the face video by integrating these components and ensuring temporal consistency. We conduct extensive quantitative and qualitative experiments along with user studies, which demonstrate that our method generates talking heads with high visual quality and conveys the perceived emotion.

067  
068036 

## 1. Introduction

069  
070

“Let’s not forget that the little emotions are the great captains of our lives and we obey them without realizing it.”

—Vincent van Gogh

071

As Vincent van Gogh said, emotion is subtle, but it helps people express themselves and understand each other. It serves as an important social signal that conveys people’s physiological states, intentions, and relationships. Humans are sensitive to various kinds of expressions, including voice tone, facial expressions, body gestures, etc. Therefore, animating a realistic emotional talking head has a wide range of applications in both interpersonal and human computer interactions. For example, when having a teleconference, people can drive their own avatar to express both content and emotions through only their voice, without using

072

073

074

075

076

077

078

079

080

081

082

083

084

085

086

087

088

089

090

091

092

093

094

095

096

097

098

099

100

101

102

103

104

105

106

107

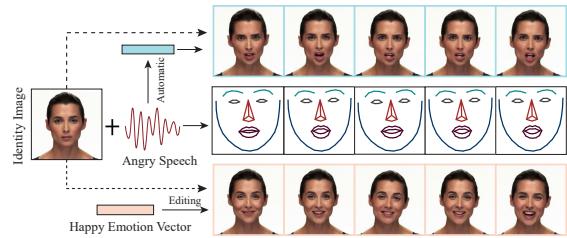


Figure 1. Problem Formulation. Given a neutral face image and an emotional speech clip as input, our model automatically synthesizes a talking head video with expressive emotions of the audio. Our model further enables conditional emotional talking face generation, provided a given emotion feature vector.

any cameras. This is extremely helpful when the participants have limited bandwidth. In addition, emotional talking head generation will enable virtual agents to help people in a more amicable way. For instance, avatars can have calm expressions when they are dealing with humans who are angry or agitated. Similarly, avatars can convey a happy facial expressions while expressing birthday or festival wishes.

Many works have been proposed in the context of speech-driven talking head generation, i.e., modeling the dynamics or movements of the talking character in synchronization with the speech signal. Some works enable facial expression generation by randomization [57] or by analyzing the audio signal [6]. However, their generated facial expressions and head motion may not be expressive enough to convey the emotional state corresponding to the input speech audio. Other methods [40, 26] can generate talking heads that convey a given emotion, but they cannot automatically detect and generate emotions from audio. Some techniques [6, 66] use facial landmarks as an intermediate representation for expression and motion manipulation, as it is easier to control these landmarks as opposed to the pixels. However, such methods may not be able to account for geometric consistency when the facial landmarks are deformed. [8] uses an attention-based mechanism to improve the visual quality. Some methods [48, 60, 65, 8, 57, 6, 66] use discriminators on the integrated face to generate realistic facial images. Overall, there are still many challenges in

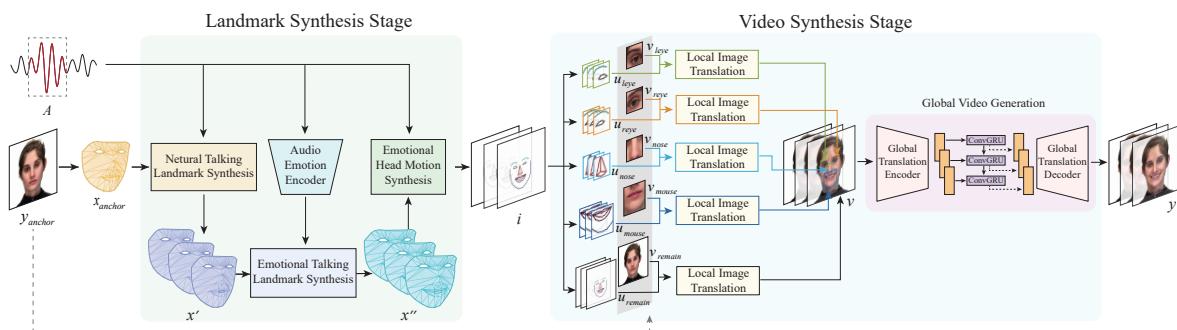
108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118

Figure 2. Illustration of our two-stage pipeline for audio-driven emotional talking head synthesis. The inputs are a speech clip  $A$  and a facial image  $y_{\text{anchor}}$ . The **neutral talking landmark synthesis module** takes extracted facial landmarks  $x_{\text{anchor}}$  and audio feature as input, and generates a talking landmark sequence  $x'$ . The **audio emotion encoder** and **emotion talking landmark synthesis** modules extract emotion features from the speech clip, and generate emotional talking landmark sequence  $x''$ . The **emotional head motion synthesis** module applies transformations and generates landmark sequence  $x$ , which are projected to obtain facial landmark images  $i$ . The **video synthesis stage** decomposes  $i$  and  $y_{\text{anchor}}$  into facial component patches. **Local image translation** module first translates landmark patches to images  $v$ . These are encoded by the **global translation encoder** and sent into a 3-layer **ConvGRU** to be further encoded with temporal dependencies. Finally, they are decoded into emotional talking head frames and used to generate the talking face video  $y$ .

terms of generating talking heads with high visual quality and temporal consistency.

**Main Results:** We present a new approach for speech-driven emotional talking head generation. Given a speech audio clip and a facial image, our method is able to produce a talking head in synchronization with the audio and expresses the emotion corresponding to the audio, e.g., happy, sad, angry, etc. We propose a two-stage pipeline consisting of a *landmark synthesis* stage and a *video synthesis* stage. The landmark synthesis stage computes the facial landmarks of the input character, which convey audio information, including word utterance and encompassed emotions. These landmarks are used to guide video synthesis in the second stage. Our formulation is inspired by the studies in anatomy related to muscle relaxation and contraction for facial movements and expressions [54], we connect landmark points with edges and represent landmarks via edge features, including edge lengths and dihedral angles. The video synthesis stage generates a high-quality talking head video from facial landmarks using local and global computations. Specifically, we first translate component-level features, which helps preserve local facial details. Next, we use a global video translator with a 3-layer *ConvGRU* [2], which improves global face quality and enhances the temporal coherence. The novel components of our work include:

- We propose a general and automatic method that generates talking heads conveying the emotion of the speech audio.
- We use geometry-aware features that accurately control facial landmarks to convey expressive emotions.
- We transform landmark sequences to realistic talking

heads in using local and global methods, that result in higher visual quality and temporal coherence. We also propose a new loss function to measure facial component similarity.

We validate our results on two emotional speech video datasets: RAVDESS [35] and CREMA-D [5]. In practice, our method out-performs state-of-the-art methods based on qualitative and quantitative metrics. We also highlight the perceptual benefits based on a user study.

## 2. Related Work

We briefly summarize recent works from three research areas: audio-based talking head generation, facial expression synthesis, and conditional face generation.

### 2.1. Audio-based Talking Head Generation

Previous research has proposed methods for generating talking heads from audio. [18] synthesizes talking heads through motion transfer from a reference motion model and then aligns lip motions with the audio clip. Without referencing a motion model, some methods directly use audio to drive 3D talking head. [26] maps audio features to 3D vertex positions and disentangles an emotion vector from the input audio. At run time, the emotion vector is chosen from the database to drive an emotional talking face. [52] first maps the input audio to a phonemic transcript and then translates the phoneme to corresponding face model parameters. [41, 55] use 3D blender-shapes, which map audio features to blender-shape parameters directly. Similarly, [13] uses the DeepSpeech network [21] as an intermediate representation of audio signal and then regresses parameters of

162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215

216 a FLAME head model [33]. [16, 51] propose methods that  
 217 use audio to control the lower half of the face or lip motion.  
 218 After that, the textured lower half of the face is reattached  
 219 to a talking head.  
 220

221 While 3D model-based methods require post-processing  
 222 for texture mapping, image-based methods can directly generate  
 223 textured talking heads. [34, 12, 7] focus on generating  
 224 talking lips and reattaching them to facial images.  
 225 [4, 48, 65, 8, 57] directly extract visual features of the  
 226 entire face from audio and combine them with the identity  
 227 features. For example, [57] develops a generative adversarial  
 228 network (GAN) based solution. It contains a generator that  
 229 encodes identity, lip motion, and expression features to the  
 230 latent space to generate a plausible talking face. It also has  
 231 three discriminators to enhance video quality, identity  
 232 consistency, and spatial consistency. In addition to generating a  
 233 plausible talking face, [19, 60, 6, 66] propose methods that  
 234 can control the motion of the generated talking head. [6] infer  
 235 the head motion pattern from a clip of input video,  
 236 while [66] learns speaker-aware head motion patterns from  
 237 input audio features. Although current methods can generate  
 238 talking heads with plausible expressions and head motions,  
 239 they cannot express the emotion state.

## 240 2.2. Facial Expression Synthesis

241 Facial expression synthesis has been actively studied in  
 242 computer graphics and computer vision. Traditional data-  
 243 driven methods synthesize facial expressions by reusing  
 244 similar image patches [38] or by retrieving similar expres-  
 245 sion patterns [28, 32] from a large database. Other methods  
 246 synthesize facial expressions by image warping [62, 61, 17,  
 247 63]. For example, [61] adopts an expression flow map to  
 248 perform image warping. [63] proposes a Flow Variational  
 249 Autoencoder that projects image differences to a flow map.  
 250

251 3D deformation-based methods [3, 53, 30, 56] are widely  
 252 used for facial expression synthesis. [53] achieves real-time  
 253 expression transfer by mapping facial deformation of ex-  
 254 pressions to a low-dimensional expression space. [30] syn-  
 255 thetizes high-quality facial expressions by deforming a 3D  
 256 facial model fitted by a 3D Morphable Model (3DMM).  
 257 Their methods not only enable the synthesis of six basic  
 258 expressions, but also support emotion detailed VA (valence  
 259 and arousal) control. [56] constructs SliderGAN which can  
 260 manipulate facial expressions through a continuous de-  
 261 formation space to generate more diverse facial expressions  
 262 than previous methods.

263 Many works use a generative model to synthesize facial  
 264 expressions. [45, 10] disentangle the latent space, enabling  
 265 facial expression synthesis by manipulating latent vectors.  
 266 Other methods adopt conditional generative networks but  
 267 differ in terms of the underlying conditions and features.  
 268 [15, 64, 31] use one-hot labels as conditions to generate cor-  
 269 responding emotions. [50, 42] use action units (AU) as con-

dition labels, which can generate more diverse expressions.  
 270 However, the AU labels are hard to acquire. [47, 44, 49, 59]  
 271 use landmarks as condition labels. Inspired by these works,  
 272 we decide to use facial landmarks as condition labels for  
 273 facial expression synthesis, since landmarks not only con-  
 274 tain information about the facial geometry, but also enable  
 275 the generation of more continuous and diverse facial expres-  
 276 sions.  
 277

## 278 2.3. Conditional Face Generation

279 Conditional GANs [37] generate images that are condi-  
 280 tioned by given attributes. Based on this architecture, many  
 281 methods [24, 58, 27] are able to generate high-quality fa-  
 282 cial images with some conditional attributes. To address the  
 283 lack of supervised data of some instances that have corre-  
 284 sponding attributes, [67, 23] propose self-consistent con-  
 285 straints that enable domain transfer of different attributes.  
 286 Further, [27] introduces StyleGAN, which generates high-  
 287 quality facial images by separating high-level facial at-  
 288 tributes in the semantic space. Some works [20, 9] intro-  
 289 duced component-level face generation. These methods  
 290 first generate facial components with separate conditional  
 291 labels and then fuse the component-level features to gen-  
 292 erate the integrated face. Inspired by these methods, we  
 293 first generate local facial features and then combine them to  
 294 generate an integrated face, enabling local facial component  
 295 control and high-quality facial image generation.  
 296

## 297 3. Our Talking Head Synthesis Algorithm

Symbol	Description
$A$	input speech audio
$y_{anchor}$	input neutral face image
$x_{anchor}$	3D landmarks extracted from $y_{anchor}$
$x'$	neutral talking landmark sequence
$x''$	emotional talking landmark sequence
$x$	emotional talking landmark sequence with head motion
$i$	facial landmark images rasterized from $x$
$u_{tx}$	landmark patch of component at time $t$
$x \in \mathcal{C} := \{leye, reye, mouth, nose, remaining\}$	
$v_{ax}$	original neutral image patch of component $x \in \mathcal{C}$
$v_{tx}$	translated local image patch of component $x \in \mathcal{C}$
$a$	audio feature sequence encoded by AutoVC [43]
$h$	LSTM hidden vector sequence
$\Delta x_t$	predicted landmark relative to $x_{anchor}$ at time $t$
$f$	LogMel Spectrogram feature
$\hat{x}'$	landmark sequence $x'$ in its edge-based representation
$s$	emotion feature vector from <i>audio emotion encoder</i>
$\tilde{f}$	the concatenated feature sequence of $\hat{x}'$ and $s$

320 Table 1. Symbols and notation used in the paper

321 In this section, we present our overall algorithm. Our  
 322 method (Fig. 2) takes a single neutral face image  $y_{anchor}$   
 323 (e.g., 256 × 256 resolution) and a speech audio clip  $A$  as

324 inputs. We generate a talking head sequence  $y$  that is a 2D  
 325 image sequence conveying basic emotions, such as happy,  
 326 sad, angry, disgusted, fear, etc. We ensure that  $y$  conveys an  
 327 emotion corresponding to audio  $A$  and preserves the iden-  
 328 tity of the input image  $y_{anchor}$ . To achieve this goal, we  
 329 propose a two-stage pipeline that utilizes a facial landmark  
 330 sequence as an intermediate representation of a face. These  
 331 landmarks are concise and capable of expressing lip dynam-  
 332 ics, facial expressions and head motions simultaneously.  
 333

334 The *landmark synthesis* stage generates landmark se-  
 335 quences corresponding to lip motion, facial expression, and  
 336 head motion. We first synthesize neutral talking landmark  
 337 sequences  $x'$  by predicting the deformation of landmarks  
 338 from the input landmarks  $x_{anchor}$  driven by the audio  $A$   
 339 (Section 3.1). Next, an emotional talking landmark syn-  
 340 thesis module deforms  $x'$  to a landmark sequence  $x''$  that ex-  
 341 hibits facial expressions conditioned on the emotion fea-  
 342 tures extracted from the input audio  $A$  (Section 3.2). At  
 343 the end of this stage, an emotional head motion synthesis  
 344 module translates and rotates the expressive landmark se-  
 345 quence  $x''$  into a talking landmark sequence  $x$  that moves  
 346 in accordance with the input audio  $A$  (Section 3.3).

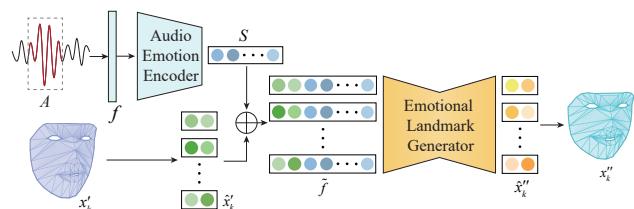
347 For the *video synthesis* stage, we observe that the po-  
 348 sitions of the generated talking head landmarks  $x$  tend to  
 349 change frequently. Thus, facial details will be compro-  
 350 mised if we directly translate the position and facial ex-  
 351 pression of landmarks to facial images. Therefore, to enhance  
 352 the translation quality in terms of facial local details and  
 353 temporal consistency, we adopt a hierarchical network to  
 354 first translate the facial details at a component level, and  
 355 then translate global information including temporal in-  
 356 formation(Section 3.4).

### 3.1. Neutral Talking Landmark Synthesis

357 This module is responsible for driving the static neu-  
 358 tral facial landmark  $x_{anchor} \in \mathbb{R}^{68 \times 3}$ , extracted by Open-  
 359 face [1], to a dynamic landmark sequence  $x' \in \mathbb{R}^{T \times 68 \times 3}$ ,  
 360 which synchronizes with the audio input  $A$ . Following [66],  
 361 we use an LSTM network  $lstm_{neu}$  to model the temporal  
 362 dependency between audio features and landmark move-  
 363 ments. Specifically, we first extract the audio content fea-  
 364 ture sequence  $a$  from  $A$  using the content encoder of Au-  
 365 toVC [43]. Then,  $lstm_{neu}$  takes  $a$  as input and outputs a  
 366 hidden feature sequence  $h$ . Finally, an MLP (Multi-Layer  
 367 Perceptron) transforms  $h$  and  $x_{anchor}$  into landmark dis-  
 368 placements  $\Delta x_t$  during each frame. For each frame, the  
 369 displaced landmarks are calculated by  $x'_t = x_{anchor} +$   
 370  $\Delta x_t, t = 1, \dots, T$ . The output of this module is a landmark  
 371 sequence  $x'$  that utters sentence  $A$ .

### 3.2. Emotional Talking Landmark Synthesis

372 Using a landmark sequence with accurate lip motions  
 373 and a neutral facial expression as input, we want to modify  
 374



375 Figure 3. Our end-to-end network used to generate emotional talk-  
 376 ing landmarks. The network takes speech audio  $A$  and neutral  
 377 talking landmarks  $x_k$  as input. First, the audio emotion encoder  
 378 extracts emotion vector  $s$  from speech audio feature  $f$ , which is  
 379 then concatenated edge-wise with landmark feature  $\hat{x}$ . The  
 380 concatenated feature  $\hat{f}$  is then fed into the *emotional landmark generator*,  
 381 to generate edge feature  $\hat{x}'$  of emotional talking landmarks.  
 382 Finally, the emotional talking landmarks  $x'$  are recovered from  
 383 their edge representation  $\hat{x}'$ .  
 384

385 these landmarks so that they can convey the emotions such  
 386 as happy, sad, angry, fearful, disgusted, etc. from input au-  
 387 dio. To achieve this goal, we propose an end-to-end network  
 388 comprised of an *emotion encoder*  $E$  and an *emotional land-  
 389 mark generator*  $G$  to analyze the emotion of speech audio  
 390 and stylize the neutral landmark, respectively, as illustrated  
 391 in Fig. 3.  
 392

393 Here we assume that the speech audio  $A$  only contains a  
 394 single emotion. To encode the emotion feature from  $A$ , we  
 395 first extract its LogMel Spectrogram feature with Librosa  
 396 [36]. We set the input with windows of size 2048 and a hop  
 397 length of 512. Then the frequency domain is evenly sepa-  
 398 rated into 128 Mel bands. Therefore, for each input audio,  
 399 we obtain its LogMel Spectrogram feature  $f \in \mathbb{R}^{128 \times n_{mel}}$ ,  
 400 where  $n_{mel}$  is the number of LogMel Spectrogram features.  
 401 With this base feature, we use a 6-layer 2D CNN encoder  
 402  $s = E(f)$ , where  $s \in [0, 1]^c$  is the emotion feature vector  
 403 and  $c$  is the number of discrete emotion classes. Our method  
 404 further enables emotional talking face editing. Results will  
 405 be shown in the supplementary material together with the  
 406 network architecture.  
 407

408 The extracted emotion feature vector stylizes the neutral  
 409 facial landmark to express the same emotion as the speech  
 410 audio. Conditioned on it, we are generating a emotional  
 411 landmark sequence with a conditional generator  $G$ .  
 412

413 However, representing the facial landmarks during styl-  
 414 ization remains a challenge. Turning the landmarks into  
 415 an image results in rasterization artifacts, while represent-  
 416 ing the landmarks through point position produces unre-  
 417 sionable, non-smooth results due to lack of constraints. To  
 418 address this, we propose connecting adjacent landmarks  
 419 to form edges, which provide strong spatial constraints.  
 420 Specifically, we use Delaunay triangulation to triangulate  
 421 a template landmark configuration. Then the obtained con-  
 422 nectionivity is also applied to other landmark configurations to  
 423 ensure topological consistency. Landmarks  $x'_k$  can be rep-  
 424

432  
433  
434  
resented as a set of edge lengths and their corresponding  
dihedral angles:  
435

$$\hat{x}'_k = \{L_k, \Theta_k | L_k \in \mathbb{R}_+^e, \Theta_k \in [-\pi, \pi)^e\}, \quad (1)$$

436  
437  
438  
where  $L_k \in \mathbb{R}_+^e$  contains all  $e$  edge lengths, and  $\Theta_k \in [-\pi, \pi)^e$  includes the dihedral angles of all edges.  
439

440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
Inspired by [11], we concatenate the landmark feature  $\hat{x}'_k$   
and emotion feature vector  $s$  edge-wise and feed the combined  
feature  $\tilde{f}_k \in \mathbb{R}^{e \times (2+c)}$  to the generator network  $G$ .  
The generator outputs  $x''_k \in \mathbb{R}^{e \times 2}$ , corresponding to the  
edge length and dihedral angles of the stylized landmarks.  
 $G$  is trained in an adversarial manner, where two discriminators  $D_{cls}$  and  $D_{src}$  discriminate between the classes of  
emotion conveyed by the landmark and whether the generated  
landmarks are real or fake, respectively. Network architectures corresponding to  $G$ ,  $D_{src}$ , and  $D_{cls}$  are detailed  
in the supplementary material.

450  
451  
452  
To train the end-to-end network, we adopt losses as follows.  
First, to have the emotion encoder  $E$  classify the emotion correctly, we use an emotion classification loss:

$$L_{emo} = -\mathbb{E}_f[E(f) \log \hat{s} + (1 - E(f))(1 - \hat{s})], \quad (2)$$

453  
454  
455  
456  
457  
458  
where  $\hat{s}$  is the ground truth emotion label of the input speech  
audio  $A$ . To have the generator network  $G$  generate plausible  
emotional landmarks, we adopt the adversarial loss:

$$L_{adv} = \mathbb{E}_{x'_k} [\log D_{src}(x'_k)] + \mathbb{E}_{x'_k, s} [\log ((1 - D_{src}(G(x'_k, s)))] \quad (3)$$

460  
461  
462  
463  
464  
To ensure generated landmarks convey the target emotion, we introduce a loss term to constrain the style of the generated landmarks:

$$L_{cls} = \mathbb{E}_{x'_k, s} [\|\hat{s} - D_{cls}(G(x'_k, s))\|_2^2]. \quad (4)$$

465  
466  
467  
468  
469  
In addition, to preserve the content, an extra reconstruction constraint is needed. Due to lack of paired stylized data, we adopt cycle consistency loss, as in [67]:

$$L_{rec} = \mathbb{E}_{x'_k, s, s'} [\|x'_k - G(G(x'_k, s), s')\|_2^2], \quad (5)$$

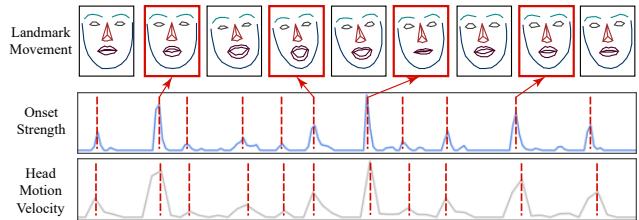
470  
471  
472  
473  
474  
where  $s'$  is the emotion label of landmarks  $x'_k$ . The overall loss is given as:

$$L = L_{emo} + \lambda_{adv} L_{adv} + \lambda_{cls} L_{cls} + \lambda_{rec} L_{rec}. \quad (6)$$

475  
476  
477  
478  
In our implementation,  $\lambda_{adv}$ ,  $\lambda_{cls}$  and  $\lambda_{rec}$  are set to 1, 10, and 10, respectively.

### 479 480 3.3. Emotional Head Motion Synthesis

481  
482  
483  
484  
485  
After generating the talking landmark sequence  $x''$ , we now synthesize head motions to match the landmarks to improve the naturalness of the generated videos. Assuming that each speech audio is a single sentence clip, we first extract the onset strengths  $o$  of the speech audio  $A$  over time



486  
487  
488  
Figure 4. Illustration of the synthesized head motion pattern. The top row is the generated landmark sequence with head motion. Red boxes highlight the talking landmarks with higher velocity. The second row is the onset strength extracted by Librosa [36]. The third row is the head motion velocity computed as the sum of the absolute of row, yaw, pitch angle velocity.  
489  
490  
491  
492  
493

494  
495  
496  
497  
498  
499  
using the Librosa library [36] and only retain those higher than a threshold  $o_{thres}$ . Then we randomize a moving direction of the talking head and set the head motion velocity synchronized with the onset strengths. As shown in Fig. 4, higher onset strengths of speech audio correspond to larger motion velocities. After head motion synthesis, the generated talking landmarks  $x$  are projected to the image plane using the camera matrix, to obtain a talking landmark image sequence  $i$  with head motion.  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509

### 3.4. Hierarchical Landmark Sequence Translation

510  
511  
512  
The second stage of the pipeline (Fig.2) is to translate the talking head landmark sequence. The problem becomes more challenging with the added head motion. We propose a hierarchical translation network with ConvGRU (Gated Recurrent Unit) [2] cells. It translates landmark sequences in a local-to-global manner and ensures the time consistency of the generated image sequences.  
513  
514  
515  
516  
517

518  
519  
520  
521  
As shown in Fig. 2, the translation network comprises two parts. In the first part, we conduct local image translation. We decompose the facial landmark image  $i_t$  and input neutral face image  $y_{anchor}$  into each of five components: left-eye, right-eye, nose, mouth and the remaining parts of the face. For each component, we train an individual *component translation network* to translate landmark patches to images patches. For example, for the left-eye, the procedure of translating landmarks  $u_{t_{leye}}$  to a left eye image patch is formulated as:  
522  
523  
524  
525  
526  
527

$$v_{t_{leye}} = T_{leye}(u_{t_{leye}}, v_{a_{leye}}), \quad (7)$$

528  
529  
530  
531  
where  $T_{leye}$  is a GAN-based translation network of the left-eye component and  $v_{t_{leye}}$  is the translated facial patch image with the same identity feature of  $v_{a_{leye}}$  and the same pose as  $u_{t_{leye}}$ .  
532  
533

534  
535  
536  
537  
538  
539  
In the second part, we conduct global image translation. We place the translated nose, mouth, left-eye, and right-eye patches on the remaining patch of their respective positions of the original landmark. Although the combined facial image may have an inconsistent optical flow and artifacts on the patch borders, the translated image sequence should



Figure 5. Results of our algorithm on RAVDESS dataset corresponding to four different emotions (the bottom four rows), compared with neutral faces (the top row). Each row displays selected frames of the generated video under a specific emotion (labeled on the left). Our results demonstrate vivid expressions and plausible head motions.

Method	RAVDESS				CREMA-D			
	CPBD↑	FID↓	CSIM↑	LMD↓	CPBD↑	FID↓	CSIM↑	LMD↓
Vougioukas et al.[57]	0.243	265.86	0.165	0.901	0.176	135.20	0.214	0.778
Chen et al.[8]	0.042	66.78	0.265	0.989	0.073	103.06	0.458	0.674
Zhou et al.[66]	0.325	77.05	0.457	0.713	0.395	89.025	0.761	0.658
Ours	<b>0.336</b>	<b>14.62</b>	<b>0.516</b>	<b>0.705</b>	<b>0.442</b>	<b>22.25</b>	<b>0.773</b>	<b>0.326</b>

Table 2. Quantitative evaluation results of different methods on RAVDESS and CREMA-D datasets (best results in bold). In all items, we achieve the best results.

be free of artifacts between patches and temporally coherent. To achieve these goals, as shown in Fig. 2, we use a global translation network to generate an integrated face with time consistency. Specifically, a 3-layer encoder network *Enc* extracts image features from the combined facial image, and a 3-layer *ConvGRU* generates sequentially consistent image features. These image features are decoded by a 3-layer decoder network *Dec* to generate a realistic and time-consistent facial image sequence  $y$ . A discriminator  $D_{global}$  is applied to train the global generator in an adversarial manner. The network architecture is detailed in the supplementary material.

We adopt several loss terms to train the network to generate realistic talking faces. Both local and global translation networks are trained in a supervised manner with ground truth landmark image pairs, so we use mean squared error (MSE) loss to constrain the image reconstruction process. Also, all translation networks are trained in an adversarial manner, so adversarial loss is employed. In addition, we use perceptual loss [25] to constrain the high-level image features. To preserve both local identity consistency between patches in the input neutral face image  $y_{anchor}$  and the gen-

erated facial images  $y$ , we propose a local identity loss inspired by [46]. Specifically, we place the generated component image patch onto its averaged position of this component on a  $224 \times 224$  all-zero tensor and resize the tensor to  $112 \times 112$  with a max-pooling operation. The max-pooled tensor is fed into Arcface [14] to get a feature vector that represents image identity. The identity loss is defined as the MSE between the feature vector of  $y_{anchor}$  and that of  $y$ . These operations ensure that the component is recognized by the Arcface and the loss is capable of back-propagating gradients through the network. With these loss terms, the translation network is able to synthesize high-quality time-consistent identity-preserved talking head image sequences.

## 4. Implementation and Results

**Datasets** We quantitatively evaluate our results on the RAVDESS [35] dataset and CREMA-D [5] dataset. The RAVDESS dataset contains 24 professional actors speaking and singing with eight general emotions: neutral, calm, happy, sad, angry, fearful, surprised, and disgusted. We only use the speaking videos and discard the singing ones. The CREMA-D dataset contains 91 speakers speaking 12



Figure 6. Comparison between our method and other methods on talking head generation on RAVDESS (left) and CREMA-D (right) datasets. Our results are displayed in the bottom two lines where the last line is the crop of the second-last line.

Method	RAVDESS				CREMA-D			
	CPBD↑	FID↓	CSIM↑	LMD↓	CPBD↑	FID↓	CSIM↑	LMD↓
Full Model	<b>0.336</b>	<b>14.621</b>	<b>0.516</b>	<b>0.705</b>	<b>0.442</b>	<b>22.252</b>	<b>0.773</b>	<b>0.326</b>
w/o local translation	0.312	39.059	0.463	0.736	0.437	58.764	0.567	0.387
w/o identity loss	0.319	54.052	0.478	0.775	0.427	60.631	0.496	0.689
w/o convGRU	0.274	52.527	0.491	0.720	0.392	39.083	0.555	0.369

Table 3. Ablation study results. We remove one component from the model in each set of experiment and compare the performance with the rest modules. As is shown in the table, performance is downgraded if any of the modules is removed, thus we prove that each module plays an essential role in our pipeline.

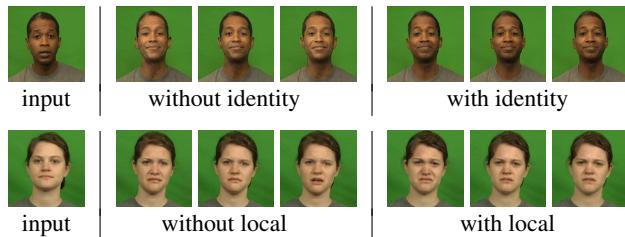


Figure 7. Ablation Studies: We show results with and w/o local identity loss. The identity information is preserved with local identity loss. The last two rows demonstrate the importance of local translation network, which preserves details as eyebrow shape.

sentences in six emotions: angry, disgusted, fearful, happy, neutral, and sad. For data preparation, we sample each video at 30FPS, crop the video by human face center, and resize the cropped image to  $256 \times 256$ . After sampling, the RAVDESS and CREMA-D datasets contain 159,222 and 570,204 video frames, respectively. We split each dataset into training and test sets. The subject IDs of the test set are illustrated in Table 4.

**Implementation Details.** The experiments were carried out on a computer with an Intel i7-6850 CPU, 128GB RAM, and a GTX 1080Ti GPU. The networks in our method are implemented in the PyTorch framework. All networks are

Dataset	Test Subjects
RAVDESS	3, 4, 5, 6
CREMA-D	15, 20, 21, 30, 33, 52, 62, 81, 89

Table 4. The subject IDs of the test sets of the RAVDESS and CREMA-D datasets.

optimized by Adam optimizer [29]. For  $lstm_{neu}$  used in *neutral talking landmark synthesis* module, the hidden size of the LSTM cell is 256. The learning rate of *neutral talking landmark synthesis*, *Audio Emotion Encoder*, *emotional landmark generator*, *local image translation* and *global video generation* module are set to  $1 \times 10^{-4}$ ,  $1 \times 10^{-4}$ ,  $1 \times 10^{-4}$ ,  $2 \times 10^{-4}$ ,  $2 \times 10^{-4}$  respectively. Other network details will be illustrated in supplementary material.

#### 4.1. Results and Comparisons

Figure 5 demonstrates the speech-driven talking head generation results of our method applied to the RAVDESS and CREMA-D datasets. More in-the-wild results are shown in the supplementary material. We demonstrate a character uttering the same sentence with five different emotions: neutral, happy, sad, angry, and disgusted. The talking heads generated by our method clearly display expressive emotions and emotion-aware head motions. The neutral talking head barely has any head motion, happy

756 emotion and sad emotion have slight head motion, and angry emotion and disgusted emotion have obvious head motion.  
 757 The results show that our method is able to generate realistic talking heads conveying the same emotion as the  
 758 input audio. Further, our model is able to realize emotional talking face editing. Results of this functionality will be  
 759 also provided in the supplementary material.  
 760

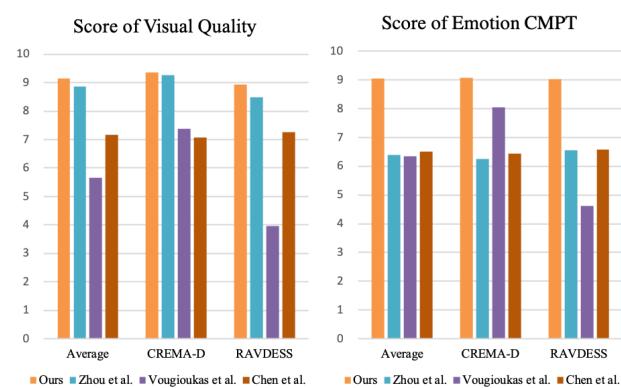
761 Figure 6 compares our method with other state-of-the-art methods [57, 8, 66]. We use an expressive speech audio  
 762 and a neutral face image as inputs. The generated talking heads of all methods are placed in rows. Vougioukas et al.  
 763 [57] generate images of size  $128 \times 96$ , which suffer from poor identity preservation. Chen et al. [8] generate images  
 764 of size  $128 \times 128$ . Their method preserves identity information, but the generated mouth area suffers from poor  
 765 visual quality. Zhou et al. [66] produce talking head images with head motion. All these methods do not generate  
 766 expressive emotions. We compare our methods with other state-of-the-art talking head generation methods quantitatively  
 767 in three aspects: video quality, identity consistency, and word pronouncing accuracy. To evaluate the generated  
 768 video frame quality, we use cumulative probability blur detection (CPBD) [39] to evaluate the frame sharpness and  
 769 Fréchet Inception Distance (FID) [22] to measure the quality and diversity of the generated images. We calculate  
 770 the FID value between the generated image and the ground truth images of the same identity and the same emotion. For  
 771 identity preservation, we follow [6] and use a CSIM criterion, which calculates the cosine similarity between the  
 772 feature vectors extracted from two candidate facial images. To evaluate word pronouncing accuracy, we align the generated  
 773 lip landmarks with ground truth lip landmarks and then compute the landmark distance (LMD) as proposed in [7]. The  
 774 evaluation metrics on the RAVDESS and CREMA-D datasets are shown in Table 2.  
 775

## 776 4.2. Ablation Study

777 We perform the following ablation studies. First, to validate the effectiveness of the local and global structures in  
 778 image translation, we design a comparison group that replaces the component-level image translation network with  
 779 a global translation network, which is denote as ‘w/o local’. Second, to demonstrate the benefits of the identity loss  
 780 in the image translation networks, we design a comparison group that removes this loss in both local translation  
 781 networks, which is denote as ‘w/o identity’. Figure 7 provides some visual examples of images generated when these  
 782 two modules are removed individually. Third, to evaluate the performance of ConvGRU in global video synthesize,  
 783 we design a comparison group that removes the ConvGRU network, which is denote as ‘w/o local convGRU’ (see  
 784 Table 3).  
 785

## 810 4.3. User Study

811 We conduct an online user study for perceptual evaluation. Users were provided with 8 video sets, 4 for  
 812 RAVDESS dataset and 4 for CREMA-D dataset. In each  
 813 video set, users were provided with a single neutral face  
 814 image and four generated talking head videos with audio  
 815 generated by Vougioukas et al. [57], Chen et al. [8], Zhou  
 816 et al. [66] and our method, respectively. Users were asked to  
 817 score each generated video from 0 to 10 in two aspects: vi-  
 818 sual quality and emotional expressiveness. We present the  
 819 results in Table 8 that demonstrate the benefits of our ap-  
 820 proach.  
 821



822 Figure 8. User study results (higher is better) on the visual quality and CMPT, which measures compatibility between audio and visual emotions. Overall, our approach is more promising.  
 823

## 824 5. Conclusion and Limitations

825 We present a new approach to generate geometrically-aware, emotionally-expressive, high-quality talking heads  
 826 driven by speech audio. Our approach computes facial landmarks  
 827 for the character and use them to control the video  
 828 synthesis. We perform local and global computations to im-  
 829 prove the visual quality and maintain temporal coherence.  
 830 We have evaluated the performance on RAVDESS and  
 831 CREMA-D datasets and highlight the benefits over prior  
 832 methods. Our approach has some limitations. The emotion  
 833 expression patterns of different subjects are similar and it  
 834 does not take into account personalized characteristics and  
 835 all arousals. Ultimately, we would like to generate con-  
 836 tinuous pattern of varying emotions based on the valence-  
 837 arousal dimensions. It would be useful to incorporate more  
 838 ideas from anatomy to generate a larger set of face expres-  
 839 sions.  
 840

## 841 References

- [1] Brandon Amos, Bartosz Ludwiczuk, Mahadev Satyanarayanan, et al. Openface: A general-purpose face recognition library with mobile applications. *CMU School of Computer Science*, 6(2), 2016. 4

- 864 [2] Nicolas Ballas, Li Yao, Chris Pal, and Aaron Courville.  
865 Delving deeper into convolutional networks for learning  
866 video representations. *arXiv preprint arXiv:1511.06432*,  
867 2015. 2, 5
- 868 [3] Volker Blanz, Curzio Basso, Tomaso Poggio, and Thomas  
869 Vetter. Reanimating faces in images and video. In *Computer*  
870 *graphics forum*, volume 22, pages 641–650. Wiley Online  
871 Library, 2003. 3
- 872 [4] Christoph Bregler, Michele Covell, and Malcolm Slaney.  
873 Video rewrite: Driving visual speech with audio. In *Pro-*  
874 *ceedings of the 24th annual conference on Computer graph-*  
875 *ics and interactive techniques*, pages 353–360, 1997. 3
- 876 [5] Huawei Cao, David G Cooper, Michael K Keutmann,  
877 Ruben C Gur, Ani Nenkova, and Ragini Verma. Crema-d:  
878 Crowd-sourced emotional multimodal actors dataset. *IEEE*  
879 *transactions on affective computing*, 5(4):377–390, 2014.  
880 2, 6
- 881 [6] Lele Chen, Guofeng Cui, Celong Liu, Zhong Li, Ziyi Kou, Yi  
882 Xu, and Chenliang Xu. Talking-head generation with rhyth-  
883 mic head motion. *arXiv preprint arXiv:2007.08547*, 2020. 1,  
884 3, 8
- 885 [7] Lele Chen, Zhiheng Li, Ross K Maddox, Zhiyao Duan, and  
886 Chenliang Xu. Lip movements generation at a glance. In *Proceedings of the European Conference on Computer Vi-*  
887 *sion (ECCV)*, pages 520–535, 2018. 3, 8
- 888 [8] Lele Chen, Ross K Maddox, Zhiyao Duan, and Chenliang  
889 Xu. Hierarchical cross-modal talking face generation with  
890 dynamic pixel-wise loss. In *Proceedings of the IEEE Con-*  
891 *ference on Computer Vision and Pattern Recognition*, pages  
7832–7841, 2019. 1, 3, 6, 7, 8
- 892 [9] Shu-Yu Chen, Wanchao Su, Lin Gao, Shihong Xia, and  
893 Hongbo Fu. Deepfacedrawing: deep generation of face im-  
894 ages from sketches. *ACM Transactions on Graphics (TOG)*,  
895 39(4):72–1, 2020. 3
- 896 [10] Brian Cheung, Jesse A Livezey, Arjun K Bansal, and  
897 Bruno A Olshausen. Discovering hidden factors of variation  
898 in deep networks. *arXiv preprint arXiv:1412.6583*, 2014. 3
- 899 [11] Yunje Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha,  
900 Sunghun Kim, and Jaegul Choo. Stargan: Unified gener-  
901 ative adversarial networks for multi-domain image-to-image  
902 translation. In *Proceedings of the IEEE conference on*  
903 *computer vision and pattern recognition*, pages 8789–8797,  
2018. 5
- 904 [12] Joon Son Chung, Amir Jamaludin, and Andrew Zisserman.  
905 You said that? *arXiv preprint arXiv:1705.02966*, 2017. 3
- 906 [13] Daniel Cudeiro, Timo Bolkart, Cassidy Laidlaw, Anurag  
907 Ranjan, and Michael J Black. Capture, learning, and syn-  
908 thesis of 3d speaking styles. In *Proceedings of the IEEE Con-*  
909 *ference on Computer Vision and Pattern Recognition*, pages  
10101–10111, 2019. 2
- 910 [14] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos  
911 Zafeiriou. Arcface: Additive angular margin loss for deep  
912 face recognition. In *Proceedings of the IEEE Conference*  
913 *on Computer Vision and Pattern Recognition*, pages 4690–  
914 4699, 2019. 6
- 915 [15] Hui Ding, Kumar Sricharan, and Rama Chellappa. Exprgan:  
916 Facial expression editing with controllable expression inten-  
917 sity. *arXiv preprint arXiv:1709.03842*, 2017. 3
- 918 [16] Bo Fan, Lijuan Wang, Frank K Soong, and Lei Xie. Photo-  
919 real talking head with deep bidirectional lstm. In *2015 IEEE*  
920 *International Conference on Acoustics, Speech and Signal*  
921 *Processing (ICASSP)*, pages 4884–4888. IEEE, 2015. 3
- 922 [17] Pablo Garrido, Levi Valgaerts, Ole Rehmsen, Thorsten  
Thormahlen, Patrick Perez, and Christian Theobalt. Auto-  
923 matic face reenactment. In *Proceedings of the IEEE con-*  
924 *ference on computer vision and pattern recognition*, pages  
4217–4224, 2014. 3
- 925 [18] Pablo Garrido, Levi Valgaerts, Hamid Sarmadi, Ingmar  
Steiner, Kiran Varanasi, Patrick Perez, and Christian  
Theobalt. Vdub: Modifying face video of actors for plau-  
926 sible visual alignment to a dubbed audio track. In *Computer*  
927 *graphics forum*, volume 34, pages 193–204. Wiley Online  
928 Library, 2015. 2
- 929 [19] David Greenwood, Iain Matthews, and Stephen Laycock.  
Joint learning of facial expression and head pose from  
930 speech. *Interspeech*, 2018. 3
- 931 [20] Shuyang Gu, Jianmin Bao, Hao Yang, Dong Chen, Fang  
Wen, and Lu Yuan. Mask-guided portrait editing with condi-  
932 tional gans. In *Proceedings of the IEEE Conference on Com-*  
933 *puter Vision and Pattern Recognition*, pages 3436–3445,  
2019. 3
- 934 [21] Awni Hannun, Carl Case, Jared Casper, Bryan Catanzaro,  
935 Greg Diamos, Erich Elsen, Ryan Prenger, Sanjeev Satheesh,  
Shubho Sengupta, Adam Coates, et al. Deep speech:  
936 Scaling up end-to-end speech recognition. *arXiv preprint*  
937 *arXiv:1412.5567*, 2014. 2
- 938 [22] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner,  
939 Bernhard Nessler, and Sepp Hochreiter. Gans trained by a  
940 two time-scale update rule converge to a local nash equilib-  
941 rium. In *Advances in neural information processing systems*,  
942 pages 6626–6637, 2017. 8
- 943 [23] Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz.  
Multimodal unsupervised image-to-image translation. In  
944 *Proceedings of the European Conference on Computer Vi-*  
945 *sion (ECCV)*, pages 172–189, 2018. 3
- 946 [24] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A  
Efros. Image-to-image translation with conditional adver-  
947 sarial networks. In *Proceedings of the IEEE conference on*  
948 *computer vision and pattern recognition*, pages 1125–1134,  
2017. 3
- 949 [25] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual  
950 losses for real-time style transfer and super-resolution. In  
951 *European conference on computer vision*, pages 694–711.  
Springer, 2016. 6
- 952 [26] Tero Karras, Timo Aila, Samuli Laine, Antti Herva, and  
953 Jaakko Lehtinen. Audio-driven facial animation by joint end-  
954 to-end learning of pose and emotion. *ACM Transactions on*  
955 *Graphics (TOG)*, 36(4):1–12, 2017. 1, 2
- 956 [27] Tero Karras, Samuli Laine, and Timo Aila. A style-based  
957 generator architecture for generative adversarial networks. In  
958 *Proceedings of the IEEE conference on computer vision and*  
959 *pattern recognition*, pages 4401–4410, 2019. 3
- 960 [28] Ira Kemelmacher-Shlizerman, Aditya Sankar, Eli Shecht-  
961 man, and Steven M Seitz. Being john malkovich. In *Euro-*  
962 *pean Conference on Computer Vision*, pages 341–353.  
Springer, 2010. 3

- 972 [29] Diederik Kingma and Jimmy Ba. Adam: A method for  
973 stochastic optimization. *Computerence*, 2014. 7 1026  
974 [30] Dimitrios Kollias, Shiyang Cheng, Evangelos Ververas,  
975 Irene Kotsia, and Stefanos Zafeiriou. Deep neural network  
976 augmentation: Generating faces for affect analysis. *International  
977 Journal of Computer Vision*, pages 1–30, 2020. 3 1027  
978 [31] Ying-Hsiu Lai and Shang-Hong Lai. Emotion-preserving  
979 representation learning via generative adversarial network  
980 for multi-view facial expression recognition. In *2018 13th  
981 IEEE International Conference on Automatic Face & Gesture  
982 Recognition (FG 2018)*, pages 263–270. IEEE, 2018. 3 1028  
983 [32] Kai Li, Feng Xu, Jue Wang, Qionghai Dai, and Yebin Liu. A  
984 data-driven approach for facial expression synthesis in video.  
985 In *2012 IEEE Conference on Computer Vision and Pattern  
986 Recognition*, pages 57–64. IEEE, 2012. 3 1029  
987 [33] Tianye Li, Timo Bolkart, Michael J Black, Hao Li, and Javier  
988 Romero. Learning a model of facial shape and expression  
989 from 4d scans. *ACM Trans. Graph.*, 36(6):194–1, 2017. 3 1030  
990 [34] Kang Liu and Joern Ostermann. Realistic facial expression  
991 synthesis for an image-based talking head. In *2011 IEEE  
992 International Conference on Multimedia and Expo*, pages 1–  
993 6. IEEE, 2011. 3 1031  
994 [35] Steven R. Livingstone and Frank A. Russo. The Ryerson  
995 Audio-Visual Database of Emotional Speech and Song  
996 (RAVDESS), Apr. 2018. Funding Information Natural Sciences  
997 and Engineering Research Council of Canada: 2012-  
998 341583 Hear the world research chair in music and emotional  
999 speech from Phonak. 2, 6 1032  
1000 [36] Brian McFee, Colin Raffel, Dawen Liang, Daniel PW Ellis,  
1001 Matt McVicar, Eric Battenberg, and Oriol Nieto. librosa:  
1002 Audio and music signal analysis in python. In *Proceedings  
1003 of the 14th python in science conference*, volume 8, 2015. 4, 1033  
1004 5 1034  
1005 [37] Mehdi Mirza and Simon Osindero. Conditional generative  
1006 adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014. 3 1035  
1007 [38] Umar Mohammed, Simon JD Prince, and Jan Kautz. Visio-  
1008 lization: generating novel facial images. *ACM Transactions  
1009 on Graphics (TOG)*, 28(3):1–8, 2009. 3 1036  
1010 [39] Niranjan D Narvekar and Lina J Karam. A no-reference  
1011 perceptual image sharpness metric based on a cumulative  
1012 probability of blur detection. In *2009 International Workshop  
1013 on Quality of Multimedia Experience*, pages 87–91. IEEE,  
1014 2009. 8 1037  
1015 [40] Hai X Pham, Samuel Cheung, and Vladimir Pavlovic.  
1016 Speech-driven 3d facial animation with implicit emotional  
1017 awareness: A deep learning approach. In *Proceedings of the  
1018 IEEE Conference on Computer Vision and Pattern Recog-  
1019 nition Workshops*, pages 80–88, 2017. 1 1038  
1020 [41] Hai Xuan Pham, Yuting Wang, and Vladimir Pavlovic.  
1021 End-to-end learning for 3d facial animation from speech. In  
1022 *Proceedings of the 20th ACM International Conference on Mul-  
1023 timodal Interaction*, pages 361–365, 2018. 2 1039  
1024 [42] Albert Pumarola, Antonio Agudo, Aleix M Martinez,  
1025 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 1040  
1026 [43] Kaizhi Qian, Yang Zhang, Shiyu Chang, Xuesong Yang, and  
Mark Hasegawa-Johnson. AutoVC: Zero-shot voice style  
transfer with only autoencoder loss. volume 97 of *Proceed-  
1027 ings of Machine Learning Research*, pages 5210–5219, Long  
Beach, California, USA, 09–15 Jun 2019. PMLR. 3, 4 1041  
1028 [44] Fengchun Qiao, Naiming Yao, Zirui Jiao, Zihao Li, Hui  
Chen, and Hongan Wang. Geometry-contrastive gan for  
facial expression transfer. *arXiv preprint arXiv:1802.01822*,  
2018. 3 1042  
1029 [45] Scott Reed, Kihyuk Sohn, Yuting Zhang, and Honglak Lee.  
1030 Learning to disentangle factors of variation with manifold in-  
teraction. In *International Conference on Machine Learning*,  
1031 pages 1431–1439, 2014. 3 1043  
1032 [46] Elad Richardson, Yuval Alaluf, Or Patashnik, Yotam Nitzan,  
Yaniv Azar, Stav Shapiro, and Daniel Cohen-Or. Encoding  
in style: a stylegan encoder for image-to-image translation.  
*arXiv preprint arXiv:2008.00951*, 2020. 6 1044  
1033 [47] Lingxiao Song, Zhihe Lu, Ran He, Zhenan Sun, and Tieniu  
Tan. Geometry guided adversarial facial expression syn-  
thesis. In *Proceedings of the 26th ACM international conference  
1034 on Multimedia*, pages 627–635, 2018. 3 1045  
1035 [48] Yang Song, Jingwen Zhu, Dawei Li, Xiaolong Wang, and  
Hairong Qi. Talking face generation by conditional recur-  
rent adversarial network. *arXiv preprint arXiv:1804.04786*,  
2018. 1, 3 1046  
1036 [49] Kritaphat Songsri-in and Stefanos Zafeiriou. Face video  
generation from a single image and landmarks. *arXiv preprint  
arXiv:1904.11521*, 2019. 3 1047  
1037 [50] Joshua M Susskind, Geoffrey E Hinton, Javier R Movellan,  
and Adam K Anderson. Generating facial expressions with  
deep belief nets. *Affective Computing, Emotion Modelling,  
Synthesis and Recognition*, pages 421–440, 2008. 3 1048  
1038 [51] Supasorn Suwajanakorn, Steven M Seitz, and Ira  
Kemelmacher-Shlizerman. Synthesizing obama: learn-  
ing lip sync from audio. *ACM Transactions on Graphics  
(TOG)*, 36(4):1–13, 2017. 3 1049  
1039 [52] Sarah Taylor, Taehwan Kim, Yisong Yue, Moshe Mahler,  
James Krahe, Anastasio Garcia Rodriguez, Jessica Hodgins,  
and Iain Matthews. A deep learning approach for generalized  
speech animation. *ACM Transactions on Graphics (TOG)*,  
36(4):1–11, 2017. 2 1050  
1040 [53] Justus Thies, Michael Zollhofer, Marc Stamminger, Chris-  
tian Theobalt, and Matthias Nießner. Face2face: Real-time  
face capture and reenactment of rgb videos. In *Proceed-  
1041 ings of the IEEE conference on computer vision and pattern  
recognition*, pages 2387–2395, 2016. 3 1052  
1042 [54] Ying-Li Tian, Takeo Kanade, and Jeffrey F Cohn. Facial  
expression analysis. In *Handbook of face recognition*, pages  
247–275. Springer, 2005. 2 1053  
1043 [55] Panagiotis Tzirakis, Athanasios Papaioannou, Alexander  
Lattas, Michail Tarasiou, Björn Schuller, and Stefanos  
Zafeiriou. Synthesising 3d facial motion from “in-the-wild”  
speech. *arXiv preprint arXiv:1904.07002*, 2019. 2 1054  
1044 [56] Evangelos Ververas and Stefanos Zafeiriou. Slidergan:  
Synthesizing expressive face images by sliding 3d blendshape  
parameters. *International Journal of Computer Vision*, pages  
1–22, 2020. 3 1055  
1045

- 1080 [57] Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic. 1134  
1081 Realistic speech-driven facial animation with gans. *International Journal of Computer Vision*, pages 1–16, 2019. 1, 3, 1135  
1082 6, 7, 8 1136  
1083 1137  
1084 [58] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, 1138  
1085 Jan Kautz, and Bryan Catanzaro. High-resolution image syn- 1139  
1086 thesis and semantic manipulation with conditional gans. In 1140  
1087 *Proceedings of the IEEE conference on computer vision and 1141  
1088 pattern recognition*, pages 8798–8807, 2018. 3 1142  
1089 [59] Wei Wang, Xavier Alameda-Pineda, Dan Xu, Pascal Fua, 1143  
1090 Elisa Ricci, and Nicu Sebe. Every smile is unique: 1144  
1091 Landmark-guided diverse smile generation. In *Proceedings 1145  
1092 of the IEEE Conference on Computer Vision and Pattern 1146  
1093 Recognition*, pages 7083–7092, 2018. 3 1147  
1094 [60] Olivia Wiles, A Sophia Koepke, and Andrew Zisserman. 1148  
1095 X2face: A network for controlling face generation using 1149  
1096 images, audio, and pose codes. In *Proceedings of the Eu- 1150  
1097 ropean conference on computer vision (ECCV)*, pages 670– 1151  
1098 686, 2018. 1, 3 1152  
1099 [61] Fei Yang, Lubomir Bourdev, Eli Shechtman, Jue Wang, and 1153  
1100 Dimitris Metaxas. Facial expression editing in video using a 1154  
1101 temporally-smooth factorization. In *2012 IEEE Conference 1155  
1102 on Computer Vision and Pattern Recognition*, pages 861– 1156  
1103 868. IEEE, 2012. 3 1157  
1104 [62] Fei Yang, Jue Wang, Eli Shechtman, Lubomir Bourdev, and 1158  
1105 Dimitri Metaxas. Expression flow for 3d-aware face com- 1159  
1106 ponent transfer. In *ACM SIGGRAPH 2011 papers*, pages 1–10. 1160  
1107 2011. 3 1161  
1108 [63] Raymond Yeh, Ziwei Liu, Dan B Goldman, and Aseem 1162  
1109 Agarwala. Semantic facial expression editing using autoen- 1163  
1110 coded flow. *arXiv preprint arXiv:1611.09961*, 2016. 3 1164  
1111 [64] Feifei Zhang, Tianzhu Zhang, Qirong Mao, and Changsheng 1165  
1112 Xu. Joint pose and expression modeling for facial expres- 1166  
1113 sion recognition. In *Proceedings of the IEEE conference on 1167  
1114 computer vision and pattern recognition*, pages 3359–3368, 1168  
1115 2018. 3 1169  
1116 [65] Hang Zhou, Yu Liu, Ziwei Liu, Ping Luo, and Xiaogang 1170  
1117 Wang. Talking face generation by adversarially disentan- 1171  
1118 gled audio-visual representation. In *Proceedings of the 1172  
1119 AAAI Conference on Artificial Intelligence*, volume 33, pages 1173  
1120 9299–9306, 2019. 1, 3 1174  
1121 [66] Yang Zhou, Dingzeyu Li, Xintong Han, Evangelos Kaloger- 1175  
1122 akis, Eli Shechtman, and Jose Echevarria. Makeittalk: 1176  
1123 Speaker-aware talking head animation. *arXiv preprint 1177  
1124 arXiv:2004.12992*, 2020. 1, 3, 4, 6, 7, 8 1178  
1125 [67] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A 1179  
1126 Efros. Unpaired image-to-image translation using cycle- 1180  
1127 consistent adversarial networks. In *Proceedings of the IEEE 1181  
1128 international conference on computer vision*, pages 2223– 1182  
1129 2232, 2017. 3, 5 1183  
1130 1184  
1131 1185  
1132 1186  
1133 1187