

VisageSynTalk: Unseen Speaker Video-to-Speech Synthesis via Speech-Visage Feature Selection

Joanna Hong[✉], Minsu Kim[✉], and Yong Man Ro[✉]

Image and Video Systems Lab, School of Electrical Engineering, KAIST, South Korea
{joanna2587, ms.k, ymro}@kaist.ac.kr

Abstract. The goal of this work is to reconstruct speech from a silent talking face video. Recent studies have shown impressive performance on synthesizing speech from silent talking face videos. However, they have not explicitly considered on varying identity characteristics of different speakers, which place a challenge in the video-to-speech synthesis, and this becomes more critical in unseen-speaker settings. Our approach is to separate the speech content and the visage-style from a given silent talking face video. By guiding the model to independently focus on modeling the two representations, we can obtain the speech of high intelligibility from the model even when the input video of an unseen subject is given. To this end, we introduce speech-visage selection that separates the speech content and the speaker identity from the visual features of the input video. The disentangled representations are jointly incorporated to synthesize speech through visage-style based synthesizer which generates speech by coating the visage-styles while maintaining the speech content. Thus, the proposed framework brings the advantage of synthesizing the speech containing the right content even with the silent talking face video of an unseen subject. We validate the effectiveness of the proposed framework on the GRID, TCD-TIMIT volunteer, and LRW datasets.

Keywords: Video to Speech Synthesis, Speech-Visage Selection

1 Introduction

Imagine a subway station packed with people, and a middle-aged woman next to you appears to ask you something. It is hard for you to understand her because of the noise of an incoming subway, so you try to follow her by looking at her face and mouth movements and infer what she tries to say. Then, you can finally understand and give an answer to her. These days, people frequently encounter these kinds of situations, not only in real-time but also in silent video conferences, corrupted video messages, and even conversations with a speech-impaired person [3]. In order to help these situations, there has been much research, namely lip-reading, on recognizing speech from silent or audio-corrupted videos.

Video-to-speech synthesis is one of the lip-reading techniques, which reconstructs speech from silent talking face videos. It has the advantage of not requiring extra human annotations (*i.e.*, text), while other conventional text-based

lip-reading techniques need them [2,37]. Nevertheless, video-to-speech synthesis is considered as challenging since it is expected to represent not only the speech content but also the identity characteristics (*e.g.*, voice) of the speaker. Thus, it is difficult to be applied in unseen, even multi-speaker, settings. There has been remarkable progresses in video-to-speech synthesis [8,9,17,22,28,33,40], especially with few speakers. While they have shown impressive performances, they have not explicitly considered the varying identity characteristics of different speakers, thus not investigated well in unseen multi-speaker setting.

To alleviate the challenge, we draw inspiration from human intuition in predicting a silent speech. When a silent talking video – seen or unseen – is given, humans firstly look at the entire appearance that represents the speaker’s character (*e.g.*, gender and age) and then predict the speech sound based on the lip movements [4]. By mimicking the human speech predicting process, we propose to learn to disentangle the lip movements (*i.e.*, speech content) and the visage appearances (*i.e.*, identities) from a silent talking face video and to predict the speech by jointly modeling the two disentangled representations. In doing so, it is promising that the model can reconstruct speech containing correct content from even unseen speaker’s talking face videos.

In this paper, we introduce a novel framework for video-to-speech synthesis. It consists of speech-visage feature selection module that separates speech content and visage-style (*i.e.*, identity) from a given talking face video. The proposed module exploits a deep learning-based feature selection [14,26] with feature transformation and normalization, which is jointly trained with the entire model in an end-to-end manner. The proposed module outputs speech selective masks, each of which contains the distinctive score of the speech content information in the visual feature of a talking face video while leaving out its speaker identity attributes. From the masks, the speech content features and the identity features can be separately driven. With the obtained two distinctive features through the speech-visage feature selection module, we introduce a visage-style based synthesizer, called VS-synthesizer. Inspired by [5,21] [5], the content features are taken into the VS-synthesizer as input, and the encoded content features are sequentially coated with the visage-styles of extracted identity features.

In order to guide the proposed framework, two learning methods are proposed: visual- and audio-identification. In visual-identification learning, we guide the network to produce the same identity features when they are from the same subject and to predict right subject identity from the identity features. Through audio-identification learning, we expect that the network well predicts the correct subject identity from the generated mel-spectrogram, even when the different identity features are coated in the original speech content features.

Through the proposed framework, the model can separately focus on modeling the speech content and generating the speech with target speaker’s appearance. It brings the advantage of synthesizing speech containing the right content even if a silent talking face video of an unseen subject is given. Moreover, the proposed framework can synthesize speech with different visage-styles while maintaining the original content. Our key contributions are as follows: (1)

To the best of our knowledge, it is the first time to directly tackle the challenge induced from varying visage-styles of different speakers, in video-to-speech synthesis by separating the identity and speech content. (2) We design a speech-visage feature selection for masking identity attributes from a talking face video while maintaining speech content, and vice versa. (3) To guarantee the disentanglement of speech content from identity, we propose two learning methods: visual-identification and audio-identification.

2 Related Work

Video to Speech Synthesis. Speech synthesis from silent talking faces is one of the lip-reading techniques that have been consistently studied [28, 43]. The initial approach [9] presented an end-to-end CNN-based model that predicts the speech audio signal from a silent talking face video and significantly improved the performance than the methods using hand-crafted visual features [29]. Another initial work [8] proposed reconstructing the speech representation by using both video frames and dense optical flow fields for capturing the dynamics of lip movements. Lip2Audspec [1] also presented a reconstruction-based video-to-speech synthesis method with autoencoders. 1D GAN-based methods [30, 40] were proposed to directly synthesize a raw waveform from the lip movements video. Lip2Wav [33] introduced a well-known sequence-to-sequence architecture into video-to-speech synthesis to capture the context. Memory [17, 22] proposed to use a multi-modal memory network to associate audio modalities during the inference. Distinct from the previous methods, we try to disentangle the identity characteristics and speech content from a silent talking face video for video-to-speech synthesis.

Feature Selection. Feature selection has become a focus of many research areas that utilize huge amounts of high-dimensional data. Early works [15, 24] initially surveyed feature selection and extraction techniques for improving learning performance, increasing computational efficiency, decreasing memory storage, and building better generalized models. Among a number of different techniques, deep learning-based feature selection methods are hybrid feature selection methods that ensemble different feature selection algorithms to construct a group of feature subsets [24]. Deep feature selection [25] selected features by imposing a sparse regularization term to select nonzero weights features at the input level. Another work [35] proposed a method to assess which features are more likely to contribute to the classification phase. In recent research for feature selection, attention-based feature selection [14] was proposed to build the correlation that best describes the degree of relevance of the target and features. Most recently, a feature mask module [26] is proposed that considers the relationships between the original features by applying a feature mask normalization. By adopting the feature selection concept, this paper attempts to select identity-relevant and content-relevant features from the visual representations.

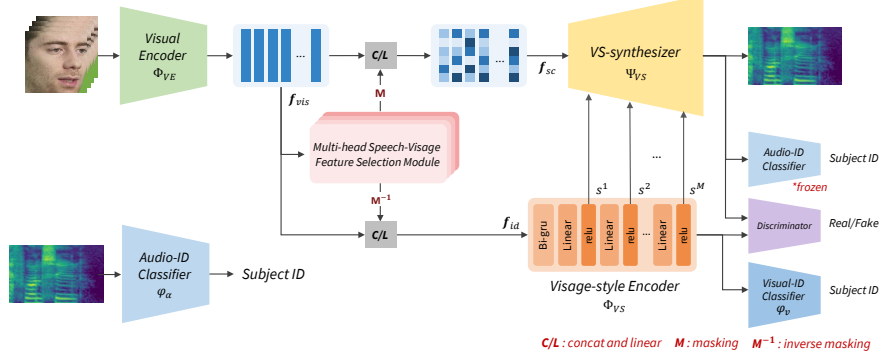


Fig. 1. Overall architecture of the proposed method, containing multi-head speech-visage feature selection and visage-style synthesizer

3 Proposed Method

Suppose we are given a sequence of silent talking face video $\mathbf{x} \in \mathbb{R}^{T \times H \times W \times 3}$ with length T , height H , and width W . The goal of our work is to reconstruct a mel-spectrogram $\mathbf{y} \in \mathbb{R}^{F \times S}$ that matches the input silent talking face frames, where F and S represents the spectral dimension of the mel-spectrogram and the frame length, respectively. The main objective of our learning problem is to disentangle the speech content and the visage-style (*i.e.*, identity) from a silent talking face video, and to synthesize speech by jointly incorporating the two disentangled representations. Hence, it is for enhancing the robustness of the model to unseen speakers and bringing the advantage of generating speech of different visage-styles with fixed speech content. Fig. 1 shows the overview of the proposed framework. It contains two major modules: multi-head speech-visage feature selection and visage-style based synthesizer.

3.1 Speech-visage feature selection

When a silent talking face video is given, humans discriminate the entire appearances of the speaker that represent the speaker’s character (*e.g.*, gender and age) and the lip movements, to associate the speech. Motivated from the human cognitive system [4], speech-visage feature selection module is designed to discriminate between human lip movements and visage-styles.

To this end, a visual encoder Φ_{VE} firstly extracts visual feature \mathbf{f}_{vis} from a silent talking face video \mathbf{x} with the dimension of embedding C ,

$$\mathbf{f}_{vis} = \Phi_{VE}(\mathbf{x}) \in \mathbb{R}^{T \times C}. \quad (1)$$

From \mathbf{f}_{vis} , the proposed speech-visage feature selection module chooses the speech content information while leaving out the identity information, by producing a speech selective mask $\bar{\mathbf{w}}$. Inspired by the modern deep feature selection

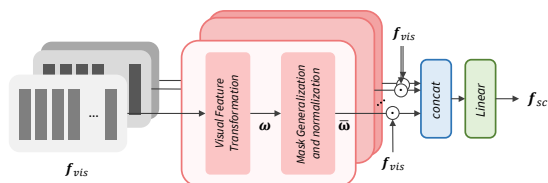


Fig. 2. Multi-head speech-visage feature selection module

method [14], the speech selective mask $\bar{\mathbf{w}}$ is produced with two steps, non-linear transformation and normalization. Firstly, a non-linear transformation, ϕ_{trans} (*i.e.*LSTM), is applied to the visual features \mathbf{f}_{vis} to capture the importance of each feature having on speech content,

$$\mathbf{w} = \phi_{trans}(\mathbf{f}_{vis}) \in \mathbb{R}^{T \times C}. \quad (2)$$

Next, mask generalization and normalization are performed to prevent the speech selective mask from being biased to the batch-wise input visual features during training [26]. This enables extracting the generalized vector from all samples. For the normalization, the softmax function is utilized to extract the importance score of the speech content of \mathbf{f}_{vis} ,

$$\bar{\mathbf{w}} = \text{Softmax}\left(\frac{1}{B} \sum_{i=1}^B \mathbf{w}_i\right), \quad (3)$$

where \mathbf{w}_i represents the transformed visual feature of i -th sample in the mini-batch size of B . Note that the generalization on mini-batch is performed during training only. Then, the speech selective mask is applied to the embedded visual feature \mathbf{f}_{vis} to select the speech content feature as follows,

$$\mathbf{f}_{sc} = \phi_{sc}(\bar{\mathbf{w}} \odot \mathbf{f}_{vis}) \in \mathbb{R}^{T \times C}, \quad (4)$$

where the ϕ_{sc} is an embedding layer, \odot represents element-wise multiplication, and \mathbf{f}_{sc} represents the selected speech content feature. Since the speech selective mask $\bar{\mathbf{w}}$ only attends to the speech content information, making the mask opposite, $\bar{\mathbf{w}}_c = 1 - \bar{\mathbf{w}}$, can produce the opposite of the speech content, namely the identity feature \mathbf{f}_{id} :

$$\mathbf{f}_{id} = \phi_{id}(\bar{\mathbf{w}}_c \odot \mathbf{f}_{vis}) \in \mathbb{R}^{T \times C}, \quad (5)$$

where ϕ_{id} represents a linear layer that embeds the selected identity feature.

Multi-head speech-visage feature selection. Due to the multiple characteristics, such as gender and age, of a speaker in regard to the identity and speech content, viewing multiple aspects of the visual face features can enable a better

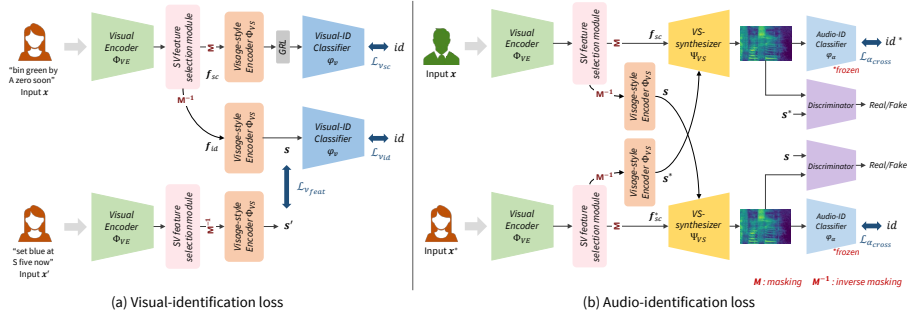


Fig. 3. Visualization of (a) visual-identification loss and (b) audio-identification loss

selection of both the speech content and identity. To enhance the feature selection procedure, the speech-visage feature selection can be employed in a multi-view fashion that produces N different speech selective masks, $\{\bar{w}^1, \dots, \bar{w}^N\}$, as shown in Fig. 2. Similar to the multi-head attention [39], our multi-view design allows the model to jointly consider the information with different aspects (*e.g.*, gender and age). The multi-view speech-visage feature selection procedure can be written as,

$$\mathbf{f}_{sc} = \phi_{sc}([\bar{w}^1 \odot \mathbf{f}_{vis}, \dots, \bar{w}^N \odot \mathbf{f}_{vis}]), \quad (6)$$

where $[,]$ represents concatenation in the channel dimension. Similarly, we also utilize the inverse of multi-view speech selective masks to obtain the identity features,

$$\mathbf{f}_{id} = \phi_{id}([\bar{w}_c^1 \odot \mathbf{f}_{vis}, \dots, \bar{w}_c^N \odot \mathbf{f}_{vis}]). \quad (7)$$

We investigate the effect of using multiple speech selective masks in Section 4.3.

3.2 Visage-style based synthesizer

The speech content features \mathbf{f}_{sc} contain the correct words of speech and the identity features \mathbf{f}_{id} have the visage-style of a certain speaker. During generation, the speech content should be maintained and only the style should be coated. Therefore, our generation objective is similar to that of style transfer [11, 18, 19]. For this purpose, we employ a style-based generator, namely Visage Style-based synthesizer (VS-synthesizer), which reconstructs the mel-spectrogram with respect to the speech content features \mathbf{f}_{sc} clothed in the encoded identity features \mathbf{f}_{id} . For the style encoder, a visage-style encoder Φ_{VS} is introduced to sequentially extract visage-style features $\mathbf{s} = \{s^1, \dots, s^M\}$ from the identity features \mathbf{f}_{id} , where M represents the number of styles which will be embedded into the synthesizer through AdaIN [18, 21]. The speech (*i.e.*, mel-spectrogram) $\mathbf{f}_{mel} \in \mathbb{R}^{F \times S}$ is generated with the following equation,

$$\mathbf{f}_{mel} = \Psi_{VS}(\mathbf{f}_{sc}, \mathbf{s}). \quad (8)$$

To convert the mel-spectrogram into a waveform, we utilize the Griffin-Lim algorithm [13] which is a well known method for converting linear spectrogram into a waveform. Following [42], we use a postnet that learns to convert the mel-spectrogram into a linear spectrogram which is utilized for the Griffin-Lim algorithm. It is trained with the reconstruction loss using ground-truth linear spectrograms.

3.3 Learning to select the speech content

To guide the proposed speech-visage feature selection module to select the speech content feature while leaving out the identity features, we propose two identification learning methods on different modalities, visual and audio.

Visual-identification learning. To guide the visage-style features \mathbf{s} obtained from the identity features \mathbf{f}_{id} contain identity-related representation, we apply the identification loss as follows,

$$\mathcal{L}_{v_{id}} = CE(\varphi_v(\Phi_{VS}(\mathbf{f}_{id})), id), \quad (9)$$

where φ_v is a visual-identity classifier, CE represents the cross-entropy loss, and id is the subject identity. Therefore, both the visage-style and identity features can carry the identity-related information. In addition to the identification loss, we sample two input talking face videos with the same subject, \mathbf{x} and \mathbf{x}' . We expect that the two extracted visage-style features from each video, \mathbf{s} and \mathbf{s}' , to be similar, since the visage-style of the same speaker is not varying. Thus, we apply mean squared error objective function as a feature loss,

$$\mathcal{L}_{v_{feat}} = \|\mathbf{s} - \mathbf{s}'\|_2. \quad (10)$$

Finally, to guarantee the disentanglement of speech content and identity representations, the speech content feature \mathbf{f}_{sc} should not contain the identity representations. To achieve this, we adopt an adversarial learning concept that guides the encoder to learn to deceive a classifier. Specifically, Gradient Reversal Layer (GRL) [10] is added before the visual-identity classifier φ_v so that the gradient sign is reversed during back-propagation. The loss function of speech content feature can be written as follows,

$$\mathcal{L}_{v_{sc}} = CE(\varphi_v(grl(\Phi_{VS}(\mathbf{f}_{sc}))), id). \quad (11)$$

Therefore, the visual-identity classifier struggles to find the identity information from \mathbf{f}_{sc} while the speech-visage feature selection module learns to not include the identity information into the speech content features \mathbf{f}_{sc} . Note that we only utilize the last style (*i.e.*, s^M) for the visual-identification learning instead of using all styles to reduce the computational cost. The final visual-identification loss (Fig. 3(a)) is defined as $\mathcal{L}_v = \mathcal{L}_{v_{id}} + \mathcal{L}_{v_{feat}} + \mathcal{L}_{v_{sc}}$.

Audio-identification learning. Although we disentangled the identity features \mathbf{f}_{id} and speech content features \mathbf{f}_{sc} , there is no guidance to properly incorporate the two disentangled representations for generating speech. Therefore, we additionally guide the model with a proposed audio-identification loss at the output side. To this end, a pre-trained audio-identity classifier φ_a is introduced to recognize the subject of the final synthesized mel-spectrogram,

$$\mathcal{L}_{a_{self}} = CE(\varphi_a(\Psi_{VS}(\mathbf{f}_{sc}, \mathbf{s})), id). \quad (12)$$

Moreover, we design a cross speech classification learning (Fig. 3(b)); when two input talking face videos with different subjects \mathbf{x} and \mathbf{x}^* are given, we crossly cloth the visage-style features \mathbf{s} and \mathbf{s}^* into the speech content features of the different subjects, \mathbf{f}_{sc}^* and \mathbf{f}_{sc} , respectively. Therefore, each generated speech should contain the crossly changed visage-style (*i.e.*, identity). This is guided with the following cross-speech classification loss,

$$\begin{aligned} \mathcal{L}_{a_{cross}} &= CE(\varphi_a(\Psi_{VS}(\mathbf{f}_{sc}, \mathbf{s}^*)), id^*) \\ &+ CE(\varphi_a(\Psi_{VS}(\mathbf{f}_{sc}^*, \mathbf{s})), id). \end{aligned} \quad (13)$$

Through the cross-speech classification loss, we can achieve both the disentanglement of speech content and identity and the ability to jointly incorporate the disentangled representations in synthesizing the desired speech. The final audio-identification loss is defined as $\mathcal{L}_a = \mathcal{L}_{a_{self}} + \mathcal{L}_{a_{cross}}$.

3.4 Total loss functions

Adversarial loss. We utilize both unconditional and conditional GAN losses [12, 31], where the former makes the generated mel-spectrogram realistic, and the latter guides the mel-spectrogram to match the final visage-style feature, s^M ,

$$\mathcal{L}_g = \log D(\mathbf{f}_{mel}) + \log D(\mathbf{f}_{mel}, s^M), \quad (14)$$

and the discriminator loss is defined as,

$$\begin{aligned} \mathcal{L}_d &= \log D(\mathbf{y}) + \log(1 - D(\mathbf{f}_{mel})) \\ &+ \log D(\mathbf{y}, s^M) + \log(1 - D(\mathbf{f}_{mel}, s^M)). \end{aligned} \quad (15)$$

Reconstruction loss. Finally, a reconstruction loss is adopted to synthesize the mel-spectrogram containing correct contents. The reconstruction loss is defined as,

$$\mathcal{L}_{recon} = \|\mathbf{y} - \mathbf{f}_{mel}\|_2 + \|\mathbf{y} - \mathbf{f}_{mel}\|_1. \quad (16)$$

Total loss. The total loss function for the generator part is the sum of the pre-defined loss functions with the balancing weights α_1 , α_2 , α_3 , and α_4 ,

$$\mathcal{L}_{tot} = \alpha_1 \mathcal{L}_v + \alpha_2 \mathcal{L}_a + \alpha_3 \mathcal{L}_g + \alpha_4 \mathcal{L}_{recon}. \quad (17)$$

4 Experiments

4.1 Dataset

GRID corpus [7] dataset is the most commonly used dataset for speech reconstruction tasks [1, 8, 9, 28, 30, 33, 40], containing 33 speakers with 6 words taken from a fixed dictionary. Since we focus on the training with a large number of subjects, we conduct experiments on two different settings: 1) multi-speaker independent (unseen) setting where the speakers in the test dataset are unseen and 2) multi-speaker dependent (seen) setting that all 33 speakers are used all training, validation, and evaluation with 90%-5%-5% split, respectively. For multi-speaker independent setting, we follow the same split as [41].

TCD-TIMIT volunteer [16] dataset has 59 speakers with about 100 phonetically rich sentences. Similar to the GRID dataset, we use two experimental settings. We utilize the officially provided data split of the TCD TIMIT dataset. Please note that it is the first time to exploit the TCD-TIMIT volunteer dataset in a video-to-speech task, which was not utilized due to its difficulties.

LRW [6] dataset contains up to 1000 utterances of 500 different words, spoken by manifold speakers. Since the original dataset does not provide identity information, we clustered and labeled the speaker information of LRW. Total 17,580 speakers are labeled; train, validation, and evaluation splits are newly generated so that the subjects are completely separated among three splits (20 for test and validation, respectively, and the rest for train). It is also the first time to utilize the identity information with the multi-speaker independent (unseen) splits. The details and splits are available in supplementary materials.

4.2 Implementation details

For both GRID and TCD-TIMIT volunteer datasets, we center-crop [44] and resize the video frames to 96×96 , and 128×128 for LRW dataset. All of the audio in the dataset are resampled to 16kHz. We convert the mel-spectrogram so that the length of the mel-spectrogram is 4 times longer than that of the video frames. The architectural details of each module can be found in the supplementary materials. We use the Adam optimizer [23] with 0.0001 learning rate, discretely decaying half at step 20000, 40000, and 60000. We choose the number N of multi-head masks to 6 and 9 for multi-speaker independent setting and multi-speaker dependent setting, respectively. The number of styles is set to 3 (*i.e.*, $M = 3$). The hyperparameters α_1 , α_2 , α_3 , and α_4 are 1.0, 1.0, 1.0, and 50.0, respectively. For computing, we use a single Titan-RTX GPU.

For the evaluation, we use three standard speech quality metrics: Short Time Objective Intelligibility (STOI) [38], Extended Short Time Objective Intelligibility (ESTOI) [20] for estimating the intelligibility and Perceptual Evaluation of Speech Quality (PESQ) [34]. To verify our generated speech, we conduct a human subjective study through mean opinion scores of naturalness, content accuracy, and voice matching.

Table 1. Performance comparison in multi-speaker independent setting on GRID

Method	STOI	ESTOI	PESQ
GAN-based [40]	0.445	0.188	1.240
Vocoder-based [28]	0.537	0.227	1.230
Lip2Wav [33]	0.522	0.251	1.284
VV-Memory [17]	0.550	0.275	1.346
End-to-end GAN [30]	0.553	0.269	1.372
Proposed model	0.567	0.308	1.373

Table 2. Performance comparison in multi-speaker independent setting on TCD-TIMIT volunteer dataset

Method	STOI	ESTOI	PESQ
Lip2Wav [33]	0.456	0.210	1.375
VV-Memory [17]	0.450	0.212	1.382
Proposed model	0.478	0.217	1.410

4.3 Experimental results

Results in multi-speaker independent setting. To verify the robustness of the proposed framework to unseen speakers, we conduct the experiments on a multi-speaker independent setting of the GRID and TCD-TIMIT volunteer datasets, where unseen subjects are utilized for testing. Table 1 elaborates the performance comparisons on the GRID dataset. We can clearly see that the proposed method outperforms the state-of-the-art performances. For the TCD-TIMIT volunteer dataset, shown in the upper part of Table 2, our proposed method achieved 0.478, 0.217, and 1.410, in STOI, ESTOI, and PESQ, respectively, outperforming the previous works [17, 33].

We additionally conduct a human subjective study through mean opinion scores (MOS) for naturalness, intelligibility, and voice matching. Naturalness evaluates how natural the synthetic speech is compared to the actual human voice, and intelligibility evaluates how clear words in the synthetic speech sound compared to the actual transcription. For the above two measures, naturalness and intelligibility, we follow the exactly same protocol of the previous works [17, 33]. We additionally measure voice matching part that determines how well the results of the proposed model matches the voice of the target speaker. We use 20 samples obtained from the multi-speaker independent setting of the GRID dataset and ask 16 participants to evaluate 6 different approaches and the ground truth in a 5-point scale. The mean scores with 95% confidence intervals are shown in Table 3. Our method achieves the score of 2.96, 3.35, and 3.34 for naturalness, intelligibility, and voice matching, respectively, which are the best among the state-of-the-art methods. Especially, the highest intelligibility means the proposed framework can generate speech containing the right content by disentangling the speech content from the identity representations. Moreover, from the voice matching, we verify that the model can synthesize the proper voices that follow the visages of the subjects even if the subjects are not seen before.

Table 3. MOS results comparison of the previous methods [22, 28, 30, 33, 40], the proposed method, and the ground truth

Method	Naturalness	Intelligibility	Voice Matching
GAN-based [40]	1.94±0.22	1.74±0.21	1.37±0.17
Vocoder-based [28]	1.98±0.16	1.68±0.25	1.15±0.11
Lip2Wav [33]	2.71±0.25	2.64±0.24	2.71±0.23
VV-Memory [17]	2.91±0.19	2.80±0.23	2.85±0.26
End-to-end GAN [30]	2.68±0.22	2.76±0.26	2.18±0.19
Proposed model	2.96±0.28	3.35±0.34	3.34±0.27
Actual Voice	4.28±0.40	4.73±0.41	-

Table 4. Performance comparison in multi-speaker dependent setting on GRID corpus

Method	STOI	ESTOI	PESQ
End-to-end GAN [30]	0.647	0.436	1.777
Proposed model	0.667	0.502	1.868

Table 5. Performance comparison in multi-speaker dependent setting on TCD-TIMIT volunteer dataset

Method	STOI	ESTOI	PESQ
Lip2Wav [33]	0.524	0.303	1.545
VV-Memory [17]	0.555	0.356	1.584
Proposed model	0.557	0.352	1.587

Results in multi-speaker dependent setting. To verify that the effectiveness of the proposed method in a multi-speaker dependent setting, we conduct experiments on the full data of the GRID dataset and the TCD-TIMIT volunteer dataset. Table 4 shows the comparison results on the GRID dataset with the previous state-of-the-art method [30]. The results on the TCD-TIMIT volunteer dataset are shown in Table 5. The proposed method achieves the best performances except for ESTOI, but it shows comparable performance with [17]. The results in the multi-speaker dependent setting show that the proposed method is effective not only for an unseen speaker but also for multi-speaker.

Results on dataset with a large number of subjects. We additionally conduct an experiment on LRW dataset which contains 17,580 subjects to verify the generalization of the proposed model to new large unseen speakers. Table 6 shows the performance in multi-speaker independent (unseen) setting on LRW. This even indicates the comparable performance to the results reported in Lip2Wav [33] (0.543 STOI, 0.344 ESTOI, and 1.197 PESQ) which has performed the experiments on LRW dataset with the original seen setting that contain overlapped subjects in all train, validation, and test splits. This proves that our model works well on dataset with a very large number of subjects with diverse vocabulary, thus generalizing our model’s performance. The audio samples of the generated speech of LRW are available in supplementary materials.

Qualitative results. We visualize the generated mel-spectrogram with the ground truth ones and those from the previous works [17, 33]. Fig. 5(a) indicates the generated mel-spectrogram from the multi-speaker independent setting of the GRID and TCD-TIMIT datasets, respectively. Additionally, Fig. 4

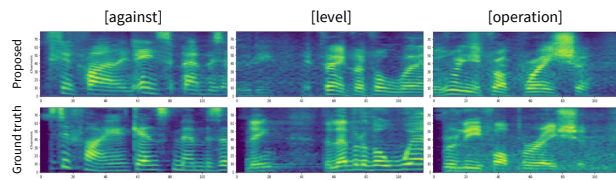


Fig. 4. Qualitative results of generated mel-spectrogram of ground truth and the proposed method on LRW

Table 6. Performance in multi-speaker independent setting on LRW

Proposed model	
STOI	0.555
ESTOI	0.305
PESQ	1.264

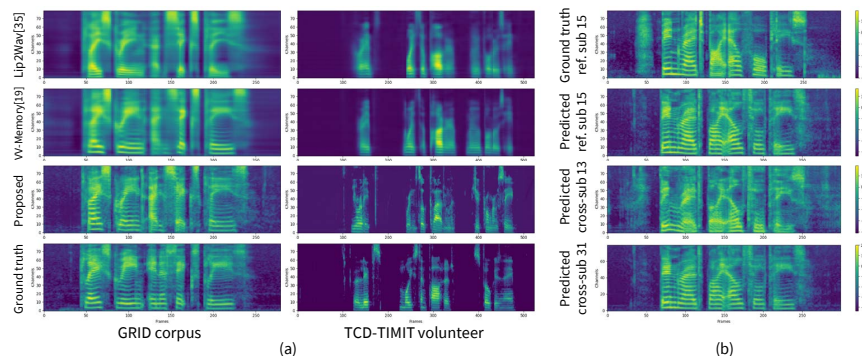


Fig. 5. Qualitative results of (a) generated mel-spectrogram of ground truth, the proposed method, [17], and [33] in multi-speaker independent setting of GRID corpus and TCD-TIMIT datasets and (b) the ground truth and the generated mel-spectrogram by changing the reference speaking-style features of subject id 15 (female) with that of subject id 13 (male), and that of subject id 31 (female)

shows the generated mel-spectrogram of words *against*, *level*, and *operation* in LRW dataset with the ground truth ones. It is clearly shown that the generated mel-spectrograms from the proposed method are visually well-matched with the ground truth mel-spectrograms.

One of our contribution is that we can synthesize speech with different visage-styles by altering the identity features f_{id} with others. Fig. 5(b) shows the results of the generated mel-spectrogram with different visage-style features, subject id 13 and 31, which are originally from the subject id 15 of the GRID corpus dataset. When we generate with male speaker’s visage-style (*i.e.*, subject id 13) we can observe that the overall frequency of generated mel-spectrogram becomes lower, which means the proposed method can reflect the changed identity features. The audio samples are provided in the supplementary materials.

Ablation study. We analyze the effectiveness of the proposed architecture through ablation studies. We firstly verify two proposed learning methods, visual- and audio-identification, that help to guide the speech-visage feature selection module. Then, we examine that the multi-head speech-visage feature selection

Table 7. Ablation study in multi-speaker independent setting on GRID dataset

Baseline	Proposed Method			STOI	ESTOI	PESQ
	\mathcal{L}_v	\mathcal{L}_a	Multi-head			
✓	✗	✗	✗	0.521	0.247	1.288
✓	✓	✗	✗	0.532	0.289	1.299
✓	✓	✓	✗	0.556	0.291	1.360
✓	✓	✓	✓	0.567	0.308	1.373

Table 8. Analysis on different number of speech selective masks in multi-speaker dependent setting on GRID dataset

Metric	N=1	N=3	N=6	N=9
STOI	0.651	0.648	0.653	0.667
ESTOI	0.489	0.480	0.486	0.502
PESQ	1.706	1.738	1.767	1.868

technique is more beneficial than the single speech-visage feature selection. Table 7 shows the ablation results in the multi-speaker independent setting using the GRID dataset. The baseline is the model that does not apply the speech-visage feature selection, so \mathbf{f}_{vis} are taken in to both VS-synthesizer and visage-style encoder. After applying the speech-visage feature selection, the performances increases when both visual- and audio- identification learning methods are adopted. The highest performances are obtained when multiple selections are adopted with 6 heads in the feature selection. The result shows that the multiple masks help the module to discover various attributes of the input visual features, thus yielding better separation of the speech content and identity, which are finally beneficial to reconstruct the speech of diverse speakers.

Effectiveness of multi-heads. To analyze the effect of different number of speech selective masks from the multi-head speech-visage feature selection module, we check the performances by differing the number of heads in multi-speaker dependent setting on the GRID dataset, shown in Table 8. While the proposed method with the single speech selective mask achieves the reasonable performance compared to [30] in Table 4, the 9 speech selective masks helps the proposed model attaining the highest performances. This means that the sufficient number of the speech selective masks enables our model to separate the speech content and identity.

We additionally visualize the representations of speech content features \mathbf{f}_{sc} and identity features \mathbf{f}_{id} in multi-speaker independent setting on the GRID dataset. Fig. 6(a) shows t-SNE [27] visualization of two features from the single speech-visage feature selection procedure, $N=1$, and Fig. 6(b) shows the two features from $N=6$. Each color represents a different subject identity. We can observe that the identity feature \mathbf{f}_{id} tends to be clustered with the same identity while the speech content feature \mathbf{f}_{sc} does not, confirming the proposed framework is effective for disentangling the two factors. Moreover, when we increase

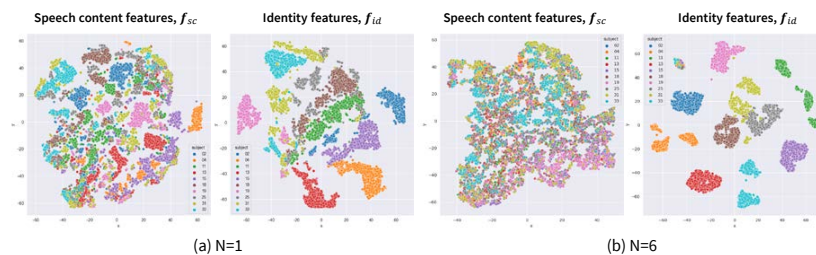


Fig. 6. t-SNE [27] visualization of speech content features f_{sc} and identity features f_{id} of (a) single speech visage feature selection procedure ($N=1$) and (b) multi-head speech visage feature selection procedure ($N=6$) in regard to the subject ids

Table 9. The Equal Error Rate (EER) for evaluating the content-voice disentanglement quality

EER (%)	f_{id}	f_{sc}
$N=1$	29.44	33.84
$N=6$	16.90	46.48

the number of heads for the speech-visage selection module, the disentanglement is further strengthened as seen in the better-clustered identity features f_{id} .

Speaker verification on disentangled features. Finally, we perform the speaker verification on the disentangled identity features f_{id} and the speech content features f_{sc} in multi-speaker independent setting on the GRID. We quantitatively evaluate the content-voice disentanglement quality using the Equal Error Rate (EER) (The lower the EER value, the higher the accuracy) which is commonly used for identity verification. Following [32], we find the EER of f_{id} to be 29.44% and that of f_{sc} to be 33.84% for $N=1$, and 16.90% and that of f_{sc} to be 46.48% for $N=6$, shown in Table 9. The results show that the proposed method can well disentangle the identity and speech content representations. With the greater N , the model can disentangle the two features more clearly.

5 Conclusion

We propose a novel video-to-speech synthesis framework with the speech-visage feature selection, visage-style based synthesizer, and two learning methods. The speech-visage feature selection separates the speech content and speaker identity, and the visage-style based synthesizer utilizes them to adequately reconstruct speech from silent talking face videos. The experimental results on benchmark databases show that the proposed method effectively synthesizes the speech from silent talking face videos of unseen speakers. [36]

Acknowledgement This work was supported by the National Research Foundation of Korea (NRF) grant funded by the Korea government (MSIT) (No. NRF-2022R1A2C2005529).

References

1. Akbari, H., Arora, H., Cao, L., Mesgarani, N.: Lip2audspec: Speech reconstruction from silent lip movements video. In: 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). pp. 2516–2520. IEEE (2018)
2. Assael, Y.M., Shillingford, B., Whiteson, S., De Freitas, N.: Lipnet: End-to-end sentence-level lipreading. arXiv preprint arXiv:1611.01599 (2016)
3. Burnham, D., Campbell, R., Away, G., Dodd, B.: Hearing eye II: the psychology of speechreading and auditory-visual speech. Psychology Press (2013)
4. Chen, T.: Audiovisual speech processing. IEEE signal processing magazine **18**(1), 9–21 (2001)
5. Chen, Y.H., Wu, D.Y., Wu, T.H., Lee, H.y.: Again-vc: A one-shot voice conversion using activation guidance and adaptive instance normalization. In: ICASSP 2021–2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). pp. 5954–5958. IEEE (2021)
6. Chung, J.S., Zisserman, A.: Lip reading in the wild. In: Asian conference on computer vision. pp. 87–103. Springer (2016)
7. Cooke, M., Barker, J., Cunningham, S., Shao, X.: An audio-visual corpus for speech perception and automatic speech recognition. The Journal of the Acoustical Society of America **120**(5), 2421–2424 (2006)
8. Ephrat, A., Halperin, T., Peleg, S.: Improved speech reconstruction from silent video. In: Proceedings of the IEEE International Conference on Computer Vision Workshops. pp. 455–462 (2017)
9. Ephrat, A., Peleg, S.: Vid2speech: speech reconstruction from silent video. In: 2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). pp. 5095–5099. IEEE (2017)
10. Ganin, Y., Lempitsky, V.: Unsupervised domain adaptation by backpropagation. In: International conference on machine learning. pp. 1180–1189. PMLR (2015)
11. Gatys, L.A., Ecker, A.S., Bethge, M.: Image style transfer using convolutional neural networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2414–2423 (2016)
12. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. Advances in neural information processing systems **27** (2014)
13. Griffin, D., Lim, J.: Signal estimation from modified short-time fourier transform. IEEE Transactions on acoustics, speech, and signal processing **32**(2), 236–243 (1984)
14. Gui, N., Ge, D., Hu, Z.: Afs: An attention-based mechanism for supervised feature selection. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 33, pp. 3705–3713 (2019)
15. Guyon, I., Elisseeff, A.: An introduction to variable and feature selection. Journal of machine learning research **3**(Mar), 1157–1182 (2003)
16. Harte, N., Gillen, E.: Tcd-timit: An audio-visual corpus of continuous speech. IEEE Transactions on Multimedia **17**(5), 603–615 (2015)
17. Hong, J., Kim, M., Park, S.J., Ro, Y.M.: Speech reconstruction with reminiscent sound via visual voice memory. IEEE/ACM Transactions on Audio, Speech, and Language Processing **29**, 3654–3667 (2021)
18. Huang, X., Belongie, S.: Arbitrary style transfer in real-time with adaptive instance normalization. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1501–1510 (2017)

19. Huang, X., Liu, M.Y., Belongie, S., Kautz, J.: Multimodal unsupervised image-to-image translation. In: Proceedings of the European conference on computer vision (ECCV). pp. 172–189 (2018)
20. Jensen, J., Taal, C.H.: An algorithm for predicting the intelligibility of speech masked by modulated noise maskers. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* **24**(11), 2009–2022 (2016)
21. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4401–4410 (2019)
22. Kim, M., Hong, J., Park, S.J., Ro, Y.M.: Multi-modality associative bridging through memory: Speech sound recollected from face video. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 296–306 (2021)
23. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
24. Li, J., Cheng, K., Wang, S., Morstatter, F., Trevino, R.P., Tang, J., Liu, H.: Feature selection: A data perspective. *ACM Computing Surveys (CSUR)* **50**(6), 1–45 (2017)
25. Li, Y., Chen, C.Y., Wasserman, W.W.: Deep feature selection: theory and application to identify enhancers and promoters. *Journal of Computational Biology* **23**(5), 322–336 (2016)
26. Liao, Y., Latty, R., Yang, B.: Feature selection using batch-wise attenuation and feature mask normalization. In: 2021 International Joint Conference on Neural Networks (IJCNN). pp. 1–9. IEEE (2021)
27. Van der Maaten, L., Hinton, G.: Visualizing data using t-sne. *Journal of machine learning research* **9**(11) (2008)
28. Michelsanti, D., Slizovskaia, O., Haro, G., Gómez, E., Tan, Z.H., Jensen, J.: Vocoder-based speech synthesis from silent videos. In: Interspeech 2020. pp. 3530–3534 (2020)
29. Milner, B., Le Cornu, T.: Reconstructing intelligible audio speech from visual speech features. *Interspeech 2015* (2015)
30. Mira, R., Vougioukas, K., Ma, P., Petridis, S., Schuller, B.W., Pantic, M.: End-to-end video-to-speech synthesis using generative adversarial networks. arXiv preprint arXiv:2104.13332 (2021)
31. Mirza, M., Osindero, S.: Conditional generative adversarial nets. arXiv preprint arXiv:1411.1784 (2014)
32. Nagrani, A., Chung, J.S., Zisserman, A.: Voxceleb: a large-scale speaker identification dataset. arXiv preprint arXiv:1706.08612 (2017)
33. Prajwal, K., Mukhopadhyay, R., Namboodiri, V.P., Jawahar, C.: Learning individual speaking styles for accurate lip to speech synthesis. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 13796–13805 (2020)
34. Rix, A., Beerends, J., Hollier, M., Hekstra, A.: Perceptual evaluation of speech quality (pesq)-a new method for speech quality assessment of telephone networks and codecs. In: 2001 IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings (Cat. No.01CH37221). vol. 2, pp. 749–752 vol.2 (2001). <https://doi.org/10.1109/ICASSP.2001.941023>
35. Roy, D., Murty, K.S.R., Mohan, C.K.: Feature selection using deep neural networks. In: 2015 International Joint Conference on Neural Networks (IJCNN). pp. 1–6. IEEE (2015)

36. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Grad-cam: Visual explanations from deep networks via gradient-based localization. In: Proceedings of the IEEE international conference on computer vision. pp. 618–626 (2017)
37. Stafylakis, T., Tzimiropoulos, G.: Combining residual networks with lstms for lipreading. arXiv preprint arXiv:1703.04105 (2017)
38. Taal, C.H., Hendriks, R.C., Heusdens, R., Jensen, J.: A short-time objective intelligibility measure for time-frequency weighted noisy speech. In: 2010 IEEE international conference on acoustics, speech and signal processing. pp. 4214–4217. IEEE (2010)
39. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: Advances in neural information processing systems. pp. 5998–6008 (2017)
40. Vougioukas, K., Ma, P., Petridis, S., Pantic, M.: Video-driven speech reconstruction using generative adversarial networks. arXiv preprint arXiv:1906.06301 (2019)
41. Vougioukas, K., Petridis, S., Pantic, M.: End-to-end speech-driven facial animation with temporal gans. arXiv preprint arXiv:1805.09313 (2018)
42. Wang, Y., Skerry-Ryan, R., Stanton, D., Wu, Y., Weiss, R.J., Jaitly, N., Yang, Z., Xiao, Y., Chen, Z., Bengio, S., et al.: Tacotron: Towards end-to-end speech synthesis. arXiv preprint arXiv:1703.10135 (2017)
43. Yadav, R., Sardana, A., Namboodiri, V.P., Hegde, R.M.: Speech prediction in silent videos using variational autoencoders. In: ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). pp. 7048–7052. IEEE (2021)
44. Zhang, S., Zhu, X., Lei, Z., Shi, H., Wang, X., Li, S.Z.: S3fd: Single shot scale-invariant face detector. In: Proceedings of the IEEE international conference on computer vision. pp. 192–201 (2017)