

# Deep Variational Generative Models for Audio-visual Speech Separation

Viet-Nhat Nguyen,<sup>1,2</sup> Mostafa Sadeghi,<sup>1</sup> Elisa Ricci,<sup>2</sup> and Xavier Alameda-Pineda,<sup>1</sup> *Senior Member, IEEE*

**Abstract**—In this paper, we are interested in audio-visual speech separation given a single-channel audio recording as well as visual information (lips movements) associated with each speaker. We propose an unsupervised technique based on audio-visual generative modeling of clean speech. More specifically, during training, a latent variable generative model is learned from clean speech spectrograms using a variational auto-encoder (VAE). To better utilize the visual information, the posteriors of the latent variables are inferred from mixed speech (instead of clean speech) as well as the visual data. The visual modality also serves as a prior for latent variables, through a visual network. At test time, the learned generative model (both for speaker-independent and speaker-dependent scenarios) is combined with an unsupervised non-negative matrix factorization (NMF) variance model for background noise. All the latent variables and noise parameters are then estimated by a Monte Carlo expectation-maximization algorithm. Our experiments show that the proposed unsupervised VAE-based method yields better separation performance than NMF-based approaches as well as a supervised deep learning-based technique.

**Index Terms**—Audio-visual speech separation, generative models, conditional variational auto-encoder.

## I. INTRODUCTION

**S**peech (or speaker) separation, sometimes referred to as the “cocktail party problem” [1], is the problem of separating speech sources from a mixture involving several speakers. This is a fundamental task in signal processing with a wide range of applications. In the literature, traditional methods are typically audio-only, for instance based on nonnegative matrix factorization (NMF) [2], [3] or independent component analysis (ICA) [4]. Given that the visual modality can provide very useful information, one important question is how to jointly exploit audio and visual data.

Audio-visual speech separation has been recently addressed in the framework of deep neural network (DNN). For example, speech enhancement, which is a special case of speech separation, has been successfully addressed in [5], where a neural network model is trained to map noisy audio and the associated visual data to clean audio data. Thus, the visual modality plays an important role allowing to separate speech by treating all the undesired speech sources as noise. Some carefully designed neural networks architectures are also proposed in [6] to learn a time-frequency mask, exploiting an additional face landmark detector trained on a separate image dataset. However, the second speech considered as the only source of noise is shared among training set and test set. Perhaps the most similar work close to the problem we are dealing with is [7]. Nevertheless, this method relies on a huge dataset of synchronized audio-visual data including thousands

of hours of video segments from the Web. These DNN-based approaches are versatile to integrate visual information into their architecture. However, they are subject to poor generalization to unseen noise types and require a highly diverse training noise dataset.

Recently, a family of unsupervised methods based on VAEs has been developed to address the speech enhancement problem [8]–[12]. They rely on deep latent variable generative models consisting of two main steps. First, a generative model is trained on clean speech to learn a probabilistic mapping from latent space to the space of (clean) speech spectrogram. This is done with a VAE network equipped with an encoder and a decoder. The encoder approximates the intractable posterior of the latent variables using a Gaussian distribution. The decoder then acts as a generative model by attempting to reconstruct the clean spectrogram. Secondly, after jointly trained with the encoder, the decoder is combined with an unsupervised noise variance model, e.g. NMF, to infer both time-varying loudness (gain values) and the model parameters. To make use of available visual modality, [10] proposed to apply a conditional VAE network [13], where the encoder and decoder are conditioned on the visual information. This approach yields state-of-the-art audio-visual speech enhancement in the unsupervised setting.

In this paper, we inspire from our previous work [10] to address the single-channel audio-visual speech separation problem. We focus on the situation of two individuals speaking simultaneously in the presence of some background noise. The main idea is to learn a generative model for clean speech spectrogram (for both speaker-dependent and speaker-independent cases), and combine it with an NMF model for the background noise at test time. A Monte-Carlo expectation-maximization framework is developed to estimate the individual speech signals. Our experiments confirm the promising performance of the proposed approach when compared with the traditional NMF-based and the recent DNN-based techniques.

The rest of the paper is organized as follows. Section II presents clean speech modeling based on VAE. Then, Section III discusses how to separate speech signals given the learned clean speech model and the observed noisy audio and the associated clean visual data. Finally, the performance of the proposed method is evaluated in Section IV.

## II. CLEAN SPEECH MODELLING

### A. Model

We use the same approach as in [10] to model the clean speech in a multimodal manner with the presence of visual information. Let  $(s, v) = \{s_n \in \mathbb{C}^F, v_n \in \mathbb{R}^M\}_{n=1}^N$  denote our observed training set of  $N$  synchronized frames of audio

<sup>1</sup> Perception team at Inria Grenoble Rhône-Alpes, France.

<sup>2</sup> Università degli Studi di Trento, Italy.

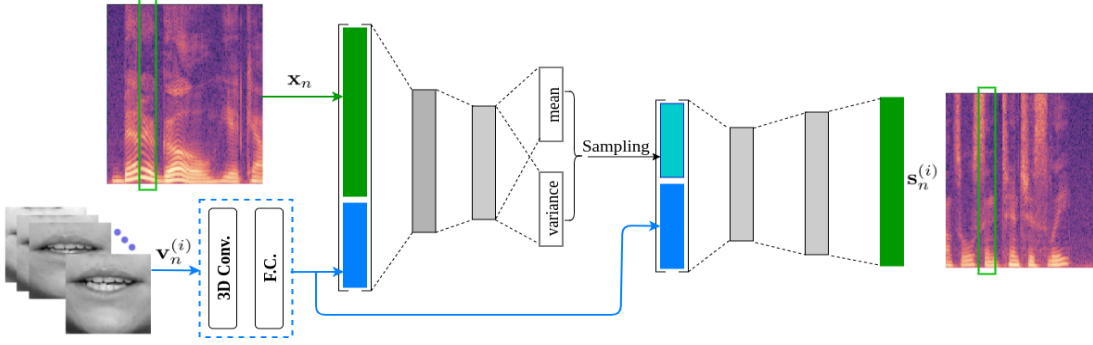


Fig. 1: VAE model architecture, where  $(i)$  indicates speaker number.

and video data where  $s_n$  represents coefficients of short-time Fourier transform (STFT) at time  $n$ , and  $v_n$  is a lip region of interest (ROI) embedding. We assume that  $s$  is generated by a corresponding sequence of latent vectors  $z = \{z_n \in \mathbb{R}^L\}_{n=1}^N$ . Both clean speech  $s$  and latent code  $z$  are conditioned on visual features which serve as a deterministic information. Specifically, we have the following latent variable model:

$$s_n | z_n, v_n; \Theta \sim \mathcal{N}_c(\mathbf{0}, \text{diag}(\sigma_s(z_n, v_n))), \quad (1)$$

$$z_n | v_n; \Gamma \sim \mathcal{N}(\mu_p(v_n), \text{diag}(\sigma_p(v_n))), \quad (2)$$

where  $\mathcal{N}(0, \sigma)$  and  $\mathcal{N}_c(0, \sigma)$  denote, respectively, Gaussian and complex proper Gaussian distributions with zero mean and variance  $\sigma$ . Moreover,  $\sigma_s(\cdot, \cdot) : \mathbb{R}^L \times \mathbb{R}^M \mapsto \mathbb{R}_+^F$  is a neural network with parameters denoted  $\Theta$ . Similarly,  $\mu_p(\cdot) : \mathbb{R}^F \mapsto \mathbb{R}^L$  and  $\sigma_p(\cdot) : \mathbb{R}^F \mapsto \mathbb{R}_+^L$  are neural networks with parameters  $\Gamma$ .

### B. Training

To learn the set of parameters  $\Theta$  and  $\Gamma$ , we follow the principles of VAE. In particular, the posterior of the latent codes is approximated by a parameterized Gaussian distribution:

$$q(z_n | s_n, v_n; \Psi) = \mathcal{N}(\mu_z(s_n, v_n), \text{diag}(\sigma_z(s_n, v_n))), \quad (3)$$

where,  $\mu_z(\cdot, \cdot) : \mathbb{R}_+^F \times \mathbb{R}^M \mapsto \mathbb{R}^L$  and  $\sigma_z(\cdot, \cdot) : \mathbb{R}_+^F \times \mathbb{R}^M \mapsto \mathbb{R}_+^L$  are neural networks, with parameters denoted  $\Psi$ , taking  $\tilde{s}_n \triangleq (|s_{0n}|^2 \dots |s_{F-1n}|^2)^\top$  as input. To learn all the parameters, the following loss function is optimized [10]:

$$\begin{aligned} \mathcal{L}(\Theta, \Gamma, \Psi) = & \alpha \mathbb{E}_{q(z|x, v; \Psi)} [\log p(s|z, v; \Theta)] + \\ & (1 - \alpha) \mathbb{E}_{p(z|v; \Gamma)} [\log p(s|z, v; \Theta)] \\ & + \alpha D_{\text{KL}}(q(z|x, v; \Psi) \| p(z|v; \Gamma)), \end{aligned} \quad (4)$$

in which,  $D_{\text{KL}}(\cdot, \cdot)$  denotes the Kullback-Leibler (KL). Moreover,  $\alpha \in [0, 1]$  controls the contributions of the reconstruction errors computed using the latent code sampled from the encoder and the one sampled from the prior network. This strategy takes into account both the reconstruction from audio-visual and visual-only information, which was originally proposed in [10] and has been proven effective for speech enhancement. The loss function can be straightforwardly optimized via the reparametrization trick and the gradient descent algorithm as in standard VAEs. In order to better utilize the visual modality during training, a mixed speech

(instead of the clean one) is used as the input of the encoder. The architecture of our VAE model is depicted in Figure 1. The encoder-decoder structure allows us to easily derive a speaker-dependent version, in which the generic decoder is substituted by the ones fine-tuned for each individual speaker. Once trained, the decoder(s) are used for speech separation as described in the following.

## III. SPEECH SEPARATION

### A. Unsupervised noise model

We consider an unsupervised NMF-based Gaussian model for the background noise which is simple and proven effective [8]. Specifically, the variance of noise spectrogram could be represented as the product of two non-negative matrices  $\mathbf{W} \in \mathbb{R}_+^{F \times K}$  and  $\mathbf{H} \in \mathbb{R}_+^{K \times N}$ .  $\mathbf{W}$  consists of  $K$  spectral basis vectors and  $\mathbf{H}$  is a matrix of temporal activations, with  $K$  being called the rank of NMF satisfying  $K(F+N) \ll FN$ . This leads to the following distribution for the  $n$ -th noise spectrogram time frame,  $\mathbf{b}_n$ :

$$\mathbf{b}_n \sim \mathcal{N}_c(\mathbf{0}, \text{diag}(\mathbf{W}\mathbf{h}_n)), \quad (5)$$

where,  $\mathbf{h}_n$  denotes the  $n$ -th column of  $\mathbf{H}$ .

### B. Mixture model

We model the observed mixture signal recorded using a single-channel microphone as follows ( $n = 1, \dots, N$ ):

$$\mathbf{x}_n = \sqrt{g_n^{(1)}} \mathbf{s}_n^{(1)} + \sqrt{g_n^{(2)}} \mathbf{s}_n^{(2)} + \mathbf{b}_n. \quad (6)$$

where,  $\mathbf{s}_n^{(1)}$  and  $\mathbf{s}_n^{(2)}$  denote clean speech signals from two speakers. As suggested in [9], the gain parameters  $g_n^{(1)}$  and  $g_n^{(2)}$  are frequency-invariant and provide robustness against the loudness variations of the speech signals along time. Using (1) and (5), and an independence assumption, we have:

$$\begin{aligned} \mathbf{x}_n | z_n^{(1)}, v_n^{(1)}, z_n^{(2)}, v_n^{(2)} \sim & \mathcal{N}_c(\mathbf{0}, g_n^{(1)} \text{diag}(\sigma_s(z_n^{(1)}, v_n^{(1)})) + \\ & g_n^{(2)} \text{diag}(\sigma_s(z_n^{(2)}, v_n^{(2)})) + \text{diag}(\mathbf{W}\mathbf{h}_n)). \end{aligned} \quad (7)$$

### C. Parameters estimation

Let us denote  $\Phi = \{\mathbf{W}, \mathbf{H}, \mathbf{g}^{(1)}, \mathbf{g}^{(2)}\}$  as the collection of parameters we need to estimate from the observed mixture  $\mathbf{x}$  and the visual embeddings  $\mathbf{v}^{(1)}, \mathbf{v}^{(2)}$ . Our ultimate goal is to maximize the log-likelihood of  $\mathbf{x}$ :

$$\max_{\Phi} \int_{\mathbb{Z}} \log p(\mathbf{x} | \mathbf{z}^{(1)}, \mathbf{v}^{(1)}, \mathbf{z}^{(2)}, \mathbf{v}^{(2)}; \Phi) d\mathbf{z}^{(1)} d\mathbf{z}^{(2)}. \quad (8)$$

However, directly optimizing the above expression is not feasible due to the intractability of the integration. As in [8], we develop a Monte Carlo Expectation-Maximization (MCEM) method. The algorithm consists of an expectation (E) step and a maximization (M) step.

1) *E-step*: In this step, the joint probability needs to be averaged over the posterior probability given the current parameters  $\Phi^*$ , leading to the following lower bound of (8):

$$\mathbb{E}_{p(\mathbf{z}^{(1)}, \mathbf{z}^{(2)} | \mathbf{x}, \mathbf{v}^{(1)}, \mathbf{v}^{(2)}; \Phi^*)} \left[ \log p(\mathbf{x}, \mathbf{z}^{(1)}, \mathbf{v}^{(1)}, \mathbf{z}^{(2)}, \mathbf{v}^{(2)}; \Phi) \right]. \quad (9)$$

Again, the above term cannot be computed analytically. In the MCEM framework, a sequence of  $R$  pairs of latent vectors  $\{(\mathbf{z}_n^{(1,r)}, \mathbf{z}_n^{(2,r)})\}_{r=1}^R$  ( $\forall n$ ) are drawn from the posterior  $p(\mathbf{z}_n^{(1)}, \mathbf{z}_n^{(2)} | \mathbf{x}_n, \mathbf{v}_n^{(1)}, \mathbf{v}_n^{(2)}; \Phi^*)$  to approximate the expected value in (9). Since a direct sampling from this posterior is difficult, we use the MetropolisHastings algorithm [14] which is a Markov Chain Monte Carlo (MCMC) method for obtaining a sequence of random samples from a given distribution. This algorithm starts with an initialization and repeatedly makes a random walk to generate new candidates from a transition kernel. Suppose that the current latent values for the  $n$ -th audio frame are  $\mathbf{z}_n^{(1)}$  and  $\mathbf{z}_n^{(2)}$ . To generate the candidate samples, we choose a symmetric multivariate normal distribution for the transition kernel as it is commonly used:

$$\mathbf{z}_n'^{(1)} | \mathbf{z}_n^{(1)}; \epsilon^2 \sim \mathcal{N}(\mathbf{z}_n^{(1)}, \epsilon^2 \mathbf{I}), \quad (10)$$

and similarly for  $\mathbf{z}_n^{(2)}$ . The new samples are accepted with the following probability:

$$p_a = \min \left( 1, \frac{p(\mathbf{z}_n^{(1)}, \mathbf{z}_n^{(2)} | \mathbf{x}_n, \mathbf{v}_n^{(1)}, \mathbf{v}_n^{(2)}; \Phi^*)}{p(\mathbf{z}_n'^{(1)}, \mathbf{z}_n'^{(2)} | \mathbf{x}_n, \mathbf{v}_n^{(1)}, \mathbf{v}_n^{(2)}; \Phi^*)} \right). \quad (11)$$

Assuming independence of the latent codes of the two speakers, the acceptance ratio in (11) turns out to be:

$$\frac{p(\mathbf{x}_n | \mathbf{z}_n^{(1)}, \mathbf{v}_n^{(1)}, \mathbf{z}_n^{(2)}, \mathbf{v}_n^{(2)}; \Phi^*) p(\mathbf{z}_n^{(1)} | \mathbf{v}_n^{(1)}) p(\mathbf{z}_n^{(2)} | \mathbf{v}_n^{(2)})}{p(\mathbf{x}_n | \mathbf{z}_n'^{(1)}, \mathbf{v}_n^{(1)}, \mathbf{z}_n'^{(2)}, \mathbf{v}_n^{(2)}; \Phi^*) p(\mathbf{z}_n'^{(1)} | \mathbf{v}_n^{(1)}) p(\mathbf{z}_n'^{(2)} | \mathbf{v}_n^{(2)})}.$$

2) *M-step*: To maximize the log likelihood in (9), we adopt the multiplicative update rule as in [10] and generalize it to the case of two clean speech signals:

$$\mathbf{H} \leftarrow \mathbf{H} \odot \left( \frac{\mathbf{W}^\top \left( |\mathbf{X}|^{\odot 2} \odot \sum_{r=1}^R (\mathbf{V}_x^{(r)})^{\odot -2} \right)}{\mathbf{W}^\top \sum_{r=1}^R (\mathbf{V}_x^{(r)})^{\odot -1}} \right)^{\odot \frac{1}{2}}, \quad (12)$$

$$\mathbf{W} \leftarrow \mathbf{W} \odot \left( \frac{\left( |\mathbf{X}|^{\odot 2} \odot \sum_{r=1}^R (\mathbf{V}_x^{(r)})^{\odot -2} \right) \mathbf{H}^\top}{\sum_{r=1}^R (\mathbf{V}_x^{(r)})^{\odot -1} \mathbf{H}^\top} \right)^{\odot \frac{1}{2}}, \quad (13)$$

where  $\odot$  denotes entrywise operations.  $\mathbf{V}_x^{(r)}$  is a matrix whose entries are  $g_n^{(1)} \sigma_{s,f}(\mathbf{z}_n^{(1,r)}, \mathbf{v}_n^{(1)}) + g_n^{(2)} \sigma_{s,f}(\mathbf{z}_n^{(2,r)}, \mathbf{v}_n^{(2)}) +$

$(\mathbf{W}\mathbf{H})_{fn}, \forall f, n$ . In a similar fashion, the formula to update the vector of speech gains is as follows:

$$\begin{aligned} (\mathbf{g}^{(i)})^\top &\leftarrow (\mathbf{g}^{(i)})^\top \odot \\ &\left( \frac{\mathbf{1}^\top \left( |\mathbf{X}|^{\odot 2} \odot \sum_{r=1}^R \left( \mathbf{V}_s^{(i,r)} \odot (\mathbf{V}_x^{(r)})^{\odot -2} \right) \right)}{\mathbf{1}^\top \left( \sum_{r=1}^R \left( \mathbf{V}_s^{(i,r)} \odot (\mathbf{V}_x^{(r)})^{\odot -1} \right) \right)} \right)^{\odot \frac{1}{2}}, \end{aligned} \quad (14)$$

in which  $i \in \{1, 2\}$  indicates speaker number and  $\mathbf{1}$  denotes a row vector of size  $F$  with all the elements being 1. Similarly,  $\mathbf{V}_s^{(i,r)}$  consists of entries being  $\sigma_{s,f}(\mathbf{z}_n^{(i,r)}, \mathbf{v}_n^{(i)})$ .

#### D. Speech separation

Having estimated the parameters  $\hat{\Phi} = \{\hat{\mathbf{W}}, \hat{\mathbf{H}}, \hat{\mathbf{g}}^{(1)}, \hat{\mathbf{g}}^{(2)}\}$ , let us denote  $\tilde{s}_{fn}^{(i)} = \sqrt{\hat{g}_n^{(i)}} s_{fn}$  as the scaled version of clean speech spectrograms. Moreover, let  $q_n(\mathbf{z}_n^{(1)}, \mathbf{z}_n^{(2)}) = p(\mathbf{z}_n^{(1)}, \mathbf{z}_n^{(2)} | \mathbf{x}_n, \mathbf{v}_n^{(1)}, \mathbf{v}_n^{(2)}; \hat{\Phi})$  denote the posterior of the latent variables. The clean speech STFT coefficients of speaker 1 are estimated in a probabilistic manner as follows:

$$\begin{aligned} \hat{s}_{fn}^{(1)} &= \mathbb{E}_{p(\tilde{s}_{fn}^{(1)} | \mathbf{x}_n, \mathbf{v}_n^{(1)}, \mathbf{v}_n^{(2)}; \hat{\Phi})} [\tilde{s}_{fn}^{(1)}] \\ &= \mathbb{E}_{q_n(\mathbf{z}_n^{(1)}, \mathbf{z}_n^{(2)})} \left[ \mathbb{E}_{p(\tilde{s}_{fn}^{(1)} | \mathbf{x}_n, \mathbf{z}_n^{(1)}, \mathbf{v}_n^{(1)}, \mathbf{z}_n^{(2)}, \mathbf{v}_n^{(2)}; \hat{\Phi})} [\tilde{s}_{fn}^{(1)}] \right]. \end{aligned} \quad (15)$$

The posterior of  $\tilde{s}_{fn}^{(1)}$  can be simplified as follows:

$$\begin{aligned} &p(\tilde{s}_{fn}^{(1)} | \mathbf{x}_n, \mathbf{z}_n^{(1)}, \mathbf{v}_n^{(1)}, \mathbf{z}_n^{(2)}, \mathbf{v}_n^{(2)}; \hat{\Phi}) \\ &\propto p(\mathbf{x}_{fn} | \tilde{s}_{fn}^{(1)}, \mathbf{z}_n^{(2)}, \mathbf{v}_n^{(2)}; \hat{\Phi}) \times p(\tilde{s}_{fn}^{(1)} | \mathbf{z}_n^{(1)}, \mathbf{v}_n^{(1)}) \\ &\sim \mathcal{N}(\mathbf{x}_{fn}; \tilde{s}_{fn}^{(1)}, \hat{g}_n^{(2)} \sigma_{s,f}(\mathbf{z}_n^{(2)}, \mathbf{v}_n^{(2)}) + (\hat{\mathbf{W}}\hat{\mathbf{H}})_{fn}) \times \\ &\quad \mathcal{N}(\tilde{s}_{fn}^{(1)}; 0, \hat{g}_n^{(1)} \sigma_{s,f}(\mathbf{z}_n^{(1)}, \mathbf{v}_n^{(1)})). \end{aligned} \quad (16)$$

By combining (16), (15), we obtain the final estimation of clean speech STFT coefficients:

$$\hat{s}_{fn}^{(1)} = \mathbb{E}_{q_n(\mathbf{z}_n^{(1)}, \mathbf{z}_n^{(2)})} \left[ \frac{\hat{g}_n^{(1)} \sigma_f(\mathbf{z}_n^{(1)}, \mathbf{v}_n^{(1)})}{\hat{V}_f(\mathbf{z}_n^{(1)}, \mathbf{z}_n^{(2)})} \right], \quad (17)$$

and analogously for  $\tilde{s}_{fn}^{(2)}$ . Here,  $\hat{V}_f(\mathbf{z}_n^{(1)}, \mathbf{z}_n^{(2)})$  is the mixture variance at the frequency  $f$  modeled in (6) with learned parameters  $\hat{\Phi}$ . These expectations could be approximated using the same MCMC method as before.

#### IV. EXPERIMENTS

**Dataset.** We evaluate the performance of our method on a mixed-speech version created from the TCD-TIMIT dataset [15] which is widely used in audio-visual (AV) speech processing research. This corpus contains AV recordings of 59 English speakers uttering 98 sentences each. Following [16], the visual data are extracted from raw videos by using a joint face and ROI detector. This results in 30 FPS videos of lips ROIs of size  $67 \times 67$  pixels. The speech signals are sampled at 16 kHz and represented by audio spectral features computed

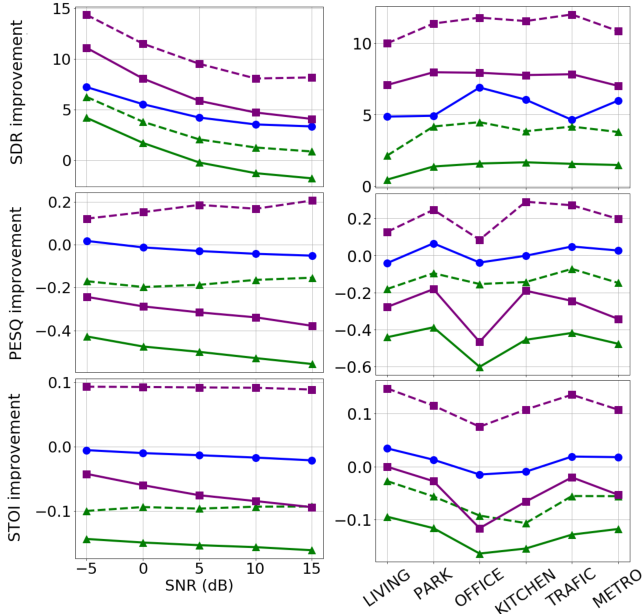


Fig. 2: Average measure improvement on the test set: SDR (top), PESQ (middle) and STOI (bottom) as a function of the SNR level (left) and of the noise type (right). NMF, DNN and CVAE are shown in blue, green and purple respectively. The dashed curves correspond to speaker-dependent models.

using Short-time Fourier transform with a STFT window of 64 ms (hence  $F = 513$ ) and 47.9% overlap.

The dataset is divided into disjoint sets of 42, 8 and 9 speakers for training, validation and testing, respectively. In each set, mixed speech signals are formed by randomly adding two utterances from two different speakers. Six types of noise taken from the DEMAND dataset [17] are added to the mixed speech signals in the test set, with noise levels:  $\{-15, -10, -5, 0, 5\}$  dB, including: DLIVING (living room), NPARK (outside park), OOFFICE (small office), DKITCHEN (inside a kitchen), STRAFFIC (busy traffic intersection) and TMETRO (subway). In the speaker-dependent approach, 80 speech signals from each speaker are dedicated to fine-tune the model and the rest is used for evaluating the performance.

**Baseline methods.** To assess the performance of our proposed method, we compare it with two other methods: an NMF-based [3] approach which is a dictionary learning and speaker-dependent method, and a DNN-based [5] approach trained to directly reconstruct spectrograms of constituent utterances. While the former is audio-only and semi-supervised, the latter is audio-visual and fully supervised.

**Architecture and parameter settings.** The architecture of our AV-CVAE is similar to that of [10], which consists of a visual feature extractor and a standard encoder-decoder. In contrast to [10] which uses a fully-connected network for the visual feature extractor, we propose to use 3 layers of 3D-convolution, followed by 2 fully-connected layers. The encoder has 2 fully-connected layers outputting a 128-dimensional latent space. The decoder is symmetric.

As suggested in [9], all the dictionaries for clean speech signals have rank  $K_{speech} = 64$ , yielding the best speech separation performance. Moreover, the rank of the noise

TABLE I: Average improvement for different configurations.

Gender Meas.	Same (Male)			Different			Same (Female)		
	SNR	PESQ	STOI	SNR	PESQ	STOI	SNR	PESQ	STOI
NMF	3.47	-0.21	-0.04	6.27	0.12	0.02	2.40	-0.17	-0.07
DNN-i	-0.95	-0.49	-0.16	1.62	-0.45	-0.15	-0.47	-0.64	-0.16
DNN-d	0.78	-0.06	-0.04	1.92	-0.07	-0.04	0.77	-0.17	-0.06
CVAE-i	5.27	-0.55	-0.10	7.89	-0.19	-0.05	5.72	-0.32	-0.09
CVAE-d	9.43	0.01	0.10	11.27	0.29	0.10	9.47	0.09	0.06

dictionary is arbitrarily fixed to  $K_{noise} = 10$ . To assess the speech separation performance of the method proposed in [5], we adapt its available implementation for speech enhancement to speech separation. To this end, a second speech signal is added to the main speech instead of noise, to form a mixed audio. The rest is kept the same.

**Training protocol.** The proposed VAE architecture is trained with Adam optimizer [18] with a constant learning rate of  $10^{-4}$  for over 45 epochs as determined by the separation performance in the validation set. The trade-off weight  $\alpha$  in loss function (4) is set to 0.9. In the speaker-dependent scenario, the VAE is fine-tuned only for 1 epoch to avoid overfitting. The batch size is set to 256.

**Results.** We evaluate the speech separation performance in terms of signal-to-distortion ratio (SDR) [19], the short-time objective intelligibility (STOI) [20], and the perceptual evaluation of speech quality (PESQ) [21] as standard metrics. We report all the measurements in terms of improvement compared with the obtained scores evaluated on the unprocessed mixture signal. We report our results in Figure 2 as a function of the noise power and the noise type. In terms of SDR, our CVAE-based methods exhibit superior performance compared to the DNN-based counterparts and to the NMF. In terms of PESQ, the proposed unsupervised approach outperforms the DNN-based counterparts, but only the CVAE-based speaker-dependent outperforms the NMF, which is speaker-dependent too. A similar behavior is observed for STOI, with a small difference. In some cases, the speaker independent CVAE outperforms the speaker-dependent DNN, which is quite a remarkable achievement. Similar conclusions are obtained when the methods are compared in different gender configurations, see Table I. One clear phenomenon observed is that the two speakers of different gender are easily separated (by all methods) than two of the same gender. Some audio examples are available on <https://team.inria.fr/perception/research/avss/>.

## V. CONCLUSIONS

In this paper, we presented an unsupervised variational generative modeling framework, inspired by [10], to address the single-channel AV speech separation problem. The proposed method consists of a training phase to learn an audio-visual generative model based on VAE for clean speech, and a test phase to estimate clean speech signals from observed single-channel audio as well as the visual information. To model the background noise, a variance model based on NMF is used. An MCEM method is developed for estimating noise parameters and clean speech signals. Our experiments showed that the proposed CVAE method exhibits better performance than a DNN-based baseline on the TCD-TIMIT dataset, both for speaker-independent and speaker-dependent scenarios.



## REFERENCES

- [1] Simon Haykin and Zhe Chen, “The cocktail party problem,” *Neural Computation*, vol. 17, no. 9, pp. 1875–1902, 2005. 1
- [2] Daniel D. Lee and H. Sebastian Seung, “Algorithms for non-negative matrix factorization,” in *Proceedings of the 13th International Conference on Neural Information Processing Systems*, Cambridge, MA, USA, 2000, NIPS00, p. 535541, MIT Press. 1
- [3] Cdric Fvotte, Emmanuel Vincent, and Alexey Ozerov, *Single-channel audio source separation with NMF: divergences, constraints and algorithms*, 01 2017. 1, 4
- [4] Jen-Tzung Chien and Bo-Cheng Chen, “A new independent component analysis for speech recognition and separation,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 14, no. 4, pp. 1245–1254, 2006. 1
- [5] Aviv Gabbay, Asaph Shamir, and Shmuel Peleg, “Visual speech enhancement,” in *Interspeech*. 2018, pp. 1170–1174, ISCA. 1, 4
- [6] G. Morrone, S. Bergamaschi, L. Pasa, L. Fadiga, V. Tikhonoff, and L. Badino, “Face landmark-based speaker-independent audio-visual speech enhancement in multi-talker environments,” in *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 6900–6904. 1
- [7] Ariel Ephrat, Inbar Mosseri, Oran Lang, Tali Dekel, Kevin Wilson, Avinatan Hassidim, William T. Freeman, and Michael Rubinstein, “Looking to listen at the cocktail party: A speaker-independent audio-visual model for speech separation,” *ACM Trans. Graph.*, vol. 37, no. 4, July 2018. 1
- [8] Y. Bando, M. Mimura, K. Itoyama, K. Yoshii, and T. Kawahara, “Statistical speech enhancement based on probabilistic integration of variational autoencoder and non-negative matrix factorization,” in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2018, pp. 716–720. 1, 2, 3
- [9] S. Leglaive, L. Girin, and R. Horaud, “A variance modeling framework based on variational autoencoders for speech enhancement,” in *2018 IEEE 28th International Workshop on Machine Learning for Signal Processing (MLSP)*, 2018, pp. 1–6. 1, 2, 4
- [10] M. Sadeghi, S. Leglaive, X. Alameda-Pineda, L. Girin, and R. Horaud, “Audio-visual speech enhancement using conditional variational autoencoders,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 28, pp. 1788–1800, 2020. 1, 2, 3, 4
- [11] M. Sadeghi and X. Alameda-Pineda, “Robust unsupervised audio-visual speech enhancement using a mixture of variational autoencoders,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020. 1
- [12] M. Sadeghi and X. Alameda-Pineda, “Mixture of inference networks for vae-based audio-visual speech enhancement,” *arXiv preprint <https://arxiv.org/abs/1912.10647>*, 2020. 1
- [13] Kihyuk Sohn, Honglak Lee, and Xinchen Yan, “Learning structured output representation using deep conditional generative models,” in *Advances in Neural Information Processing Systems* 28, C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, Eds., pp. 3483–3491. Curran Associates, Inc., 2015. 1
- [14] W. K. Hastings, “Monte Carlo sampling methods using Markov chains and their applications,” *Biometrika*, vol. 57, no. 1, pp. 97–109, 04 1970. 3
- [15] N. Harte and E. Gillen, “Tcd-timit: An audio-visual corpus of continuous speech,” *Multimedia, IEEE Transactions*, vol. 17, no. 5, pp. 603–615, 2015. 3
- [16] Ahmed Hussen Abdelaziz, “Ntcd-timit,” May 2017. 3
- [17] Joachim Thiemann, Nobutaka Ito, and Emmanuel Vincent, “DEMAND: a collection of multi-channel recordings of acoustic noise in diverse environments,” June 2013, Supported by Inria under the Associate Team Program VERSAMUS. 4
- [18] Diederik P. Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun, Eds., 2015. 4
- [19] E. Vincent, R. Gribonval, and C. Fevotte, “Performance measurement in blind audio source separation,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 14, no. 4, pp. 1462–1469, 2006. 4
- [20] C. H. Taal, R. C. Hendriks, R. Heusdens, and J. Jensen, “A short-time objective intelligibility measure for time-frequency weighted noisy speech,” in *2010 IEEE International Conference on Acoustics, Speech and Signal Processing*, 2010, pp. 4214–4217. 4
- [21] A. W. Rix, J. G. Beerends, M. P. Hollier, and A. P. Hekstra, “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)*, 2001, vol. 2, pp. 749–752 vol.2. 4