# Unsupervised Speech Enhancement using Dynamical Variational Auto-Encoders

Xiaoyu Bie,[1] Simon Leglaive,[2] Xavier Alameda-Pineda,[1] *Senior Member, IEEE*, Laurent Girin[3]

*Abstract*—**Dynamical variational auto-encoders (DVAEs) are a class of deep generative models with latent variables, dedicated to time series data modeling. DVAEs can be considered as extensions of the variational autoencoder (VAE) that include the modeling of temporal dependencies between successive observed and/or latent vectors in data sequences. Previous work has shown the interest of DVAEs and their better performance over the VAE for speech signals (spectrogram) modeling. Independently, the VAE has been successfully applied to speech enhancement in noise, in an unsupervised noise-agnostic set-up that does not require the use of a parallel dataset of clean and noisy speech samples for training, but only requires clean speech signals. In this paper, we extend those works to DVAE-based single-channel unsupervised speech enhancement, hence exploiting both speech signals unsupervised representation learning and dynamics modeling. We propose an unsupervised speech enhancement algorithm based on the most general form of DVAEs, that we then adapt to three specific DVAE models to illustrate the versatility of the framework. More precisely, we combine DVAE-based speech priors with a noise model based on nonnegative matrix factorization, and we derive a variational expectation-maximization (VEM) algorithm to perform speech enhancement. Experimental results show that the proposed approach based on DVAEs outperforms its VAE counterpart and a supervised speech enhancement baseline.**

*Index Terms*—**Speech enhancement, dynamical variational autoencoders, nonnegative matrix factorization, variational inference.**

## I. INTRODUCTION

SPEECH enhancement is a classical and fundamental problem in speech processing [1], [2]. It aims at recovering a clean speech signal which has been corrupted by acoustic (background) noise, possibly in a wide variety of scenarios, i.e. different noise types and signal-to-noise power ratios. Classical "signal-processing-based" solutions include spectral subtraction [3], Wiener filtering [4], which use noise (and possibly clean speech) power spectral density estimate(s) obtained from the noisy signal, and the short-time spectral amplitude estimator [5].

Recently, the impressive advances in deep neural networks (DNNs) and deep learning (DL) have opened new possibilities to tackle this task. The most common approach to DL-based speech enhancement, that has been widely studied in the past years, is a *supervised* one, where pairs of noisy and clean speech signals are used during model training. With a large amount of such training data, DNNs can efficiently learn a mapping from noisy speech signals either directly to clean speech signals or to time-frequency denoising masks that are applied to the noisy signal; see a review in [6] and see Section II for other recent examples of works in this line. However, this set of methods can suffer from poor generalization to unseen noise type and acoustic conditions.

In contrast to the supervised methods, *unsupervised* methods, more sparsely considered in the literature, model the clean speech signal, and thus only require a clean speech dataset at training time. These methods estimate the noise properties at test time, while performing speech enhancement.[1] In particular, the variational auto-encoder (VAE) [11], [12] is a deep generative model that has proven successful for unsupervised speech enhancement, both in single-channel [13]–[16] and multi-channel [17]–[19] configuration. Typically, a VAE is trained with clean speech signals to learn a deep latent variable model for clean speech (in the short-term Fourier transform (STFT) domain). At test time, where only noisy speech signals are observed, the pre-trained VAE is combined with a noise model, usually based on nonnegative matrix factorization (NMF) [20], to formulate the mixture model. The latent variables of the clean speech signal, as well as the parameters of the noise model, are then estimated, which enables to design a denoising Wiener-like filter.

Although this approach has shown very interesting results, in particular robustness to different types of noise, almost all previous works considered the original VAE where speech time frames are processed independently. The independence of successive speech frames in a VAE is a too strong and clearly suboptimal assumption, since the model will not exploit the speech signal (spectrogram) dynamics. Yet, in the last years, a number of extensions of the VAE have been proposed to model sequential data, with explicit modeling of the dependencies between successive vectors of latent and observed variables [21]–[29]. In [30], we performed an extensive and comprehensive literature review of these models. We introduced a general class of models called Dynamical Variational Autoencoders (DVAEs) that encompasses and unifies the above-cited temporal VAE extensions, and other possible ones. We also showed that all reviewed DVAEs outperform the VAE in

[1] Inria Grenoble Rhône-Alpes, Univ. Grenoble-Alpes, France
[2] CentraleSupélec, IETR, France
[3] Univ. Grenoble Alpes, Grenoble-INP, GIPSA-lab, France

---

[1]Note that this setting was referred to as "semi-supervised" in previous works on audio source separation [7]–[9]. However, we choose to call it here unsupervised because in the machine learning literature, semi-supervised usually refers to methods that are trained from both labeled and unlabeled datasets (e.g. [10]). In this context, a semi-supervised speech enhancement method would be trained from both a labeled dataset of noisy and clean speech signal pairs, and an unlabeled dataset containing only noisy or clean speech. This is not the case in the present study, where the training dataset only contains clean speech signals.

the speech analysis-resynthesis task, which shows the benefit of considering temporal correlations when modeling speech signal. Although, this opens a large set of possibilities for joint unsupervised speech representation learning and speech dynamics modeling, to the best of our knowledge, only one DVAE model, so-called Recurrent VAE (RVAE), has been used for speech enhancement so far, in our previous work [28].

In the present paper, we present a general framework to use DVAEs in unsupervised speech enhancement: We use the most general formulation of a DVAE to develop the general formulation of DVAE-based unsupervised speech enhancement. Similar to the aforementioned unsupervised methods based on the original VAE, we use DVAEs to model clean speech signals at training time. Pre-training on clean speech can provide a well estimated generative model (also referred to as decoder) and a well estimated inference model (also referred to as encoder). To tackle the speech enhancement problem, we consider the noisy signals as a mixture of clean speech signal (modulated with a frequency-independent gain), and additive noise modeled by NMF. We develop a generic variational expectation-maximization (VEM) [31] algorithm for speech enhancement compatible with all DVAE models. At the E-step, the encoder of a DVAE is fine-tuned to approximate the posterior distribution of the latent variables of clean speech, whereas the NMF parameters and the gain are estimated at M-step. Then the clean speech can be estimated with a probabilistic Wiener filter. To illustrate the versatility of this approach, we apply it to three examples of DVAEs, namely the Deep Kalman Filter (DKF) [23], [24], the Stochastic Recurrent Neural Network (SRNN) [26] and the above-mentioned RVAE.

The rest of the paper is organized as follows. Section II gives further informations on related works. Section III briefly reviews the VAE and DVAE models, rapidly present the above-mentioned three examples of DVAEs, and discusses a few issues of modeling speech signals with DVAEs. Then, in Section IV, we present the proposed DVAE-based speech enhancement algorithm in detail. Section V present a series of experiments conducted with the three example DVAE models. Section VI concludes the paper.

## II. Related Work

Deep learning based speech enhancement models have been widely studied over the past decade. According to the type of input and output data, they can be divided into two categories. In supervised methods, discriminative patterns of speech and noise are learned from training data with noisy-clean speech pairs [6], whereas unsupervised speech enhancement models do not. Since this paper aims at single-channel speech enhancement, the related work on multi-channel are out of the scope of this section.

Deep supervised approaches for speech enhancement were first proposed in [32], [33], where DNNs pretrained with restrictive Boltzmann machine (RBM) were used as binary classifiers to the time-frequency (T-F) spectrogram of noisy speech. Unlike T-F masking methods, [34] proposed a deep auto-encoder (DAE) to learn a mapping from noisy speech to clean speech on the Mel-frequency power spectrum (MFP).

The DAE in [34] is a stack of multiple pretrained auto-encoders (AEs) and will be further fine-tuned after stacking all AEs with pretrained parameters for initialization. Different from previous methods that process the (log-)magnitude spectrum only, [34] proposed to use fully convolutional network (FCN) for speech enhancement on raw waveform. Furthermore, rather than using discriminative approaches, [35] use generative adversarial networks (GAN) to directly learn to generate the clean waveform, where noisy waveform was encoded as additional layer-by-layer conditions for the generator with decoder structure to output a clean speech.

Alternatively, unsupervised methods do not use noisy-clean speech pairs for training, thus avoiding potential robustness issues against unknown acoustic environments. [13] proposed to use VAE to learn a prior distribution on clean speech. At test time, the noise signal is characterized by a Bayesian NMF model whose parameters, as well as the VAE latent variables, are estimated with a Markov chain Monte Carlo (MCMC) algorithm. At each step of the iterative algorithm, the NMF parameters are sampled following a Gibb-sampling-like scheme, whereas the VAE latent variables are sampled with a Metropolis-Hastings algorithm. [14] added a gain parameter for clean speech to provide more robustness with respect to the loudness of the training examples. Different from [13], [14] derived a Monte Carlo expectation-maximization (MCEM) algorithm for inferring the latent variables in the VAE and estimating the mixture model parameters. [15] proposed to use an alpha-stable distribution for the noise model, instead of the Gaussian assumption. [16] developed an iterative variational inference method based on an extended set of latent variables and leveraging the already-trained VAE encoder. This approach is computationally efficient since it does not require sampling or gradient descent at each step of the algorithm. More recently, a guided VAE was proposed in [36], where the VAE-based clean speech signal prior is defined conditionally on the voice activity detection or the ideal binary mask. This guiding information has to be provided by a supervised classifier, separately trained on noisy speech signals.

Most of the VAE-based unsupervised speech enhancement studies published so far focus on exploiting different probability distributions for the latent variables and/or deriving efficient inference and learning algorithms. Very few studies deal with the inherent limitation of the VAE to handle sequential data, that is sequences of samples that are not statistically independent, as is the case of speech data. To the best of our knowledge, only [28] and [37] proposed generative approaches to speech enhancement based on VAE variants that can learn temporal dependencies in the speech model. While [28] proposed a recurrent VAE (RVAE) based on standard recurrent neural networks (RNNs), [37] used stochastic temporal convolutional network (TCNs) [38], [39], allowing the latent variables to have both hierarchical and temporal dependencies.

Besides the above-mentioned previous work on temporal modeling for speech enhancement, other studies have focused on developing extensions of the original VAE for time series. The Deep Kalman Filter (DKF) [23] combines state-space models and VAEs for sequential data modeling. DKF can be

seen as a standard VAE with an additional non-linear first-order Markov model on the latent vectors. It was further developed in [24] by considering an inference model that is consistent with the exact posterior structure, in terms of random variable dependencies. The Stochastic Recurrent Neural Network (SRNN) [26] is another temporal extension of the VAE, with a more complex dynamical generative modeling of the observed and latent data sequences, involving infinite-order temporal dependencies implemented with RNNs. Unlike the original VAE [11], [12] and RVAE [28] which assume independent and identically distributed latent vectors, both DKF and SRNN consider a first-order dependency between successive latent vectors. Actually, RVAE [28], DKF [23], [24], SRNN [26] and several other temporal extensions of the VAE [21], [22], [25], [27], [29] can all be seen as particular instances of a general class of models called Dynamical Variational Auto-encoders (DVAEs), which have been recently reviewed in [30]. These particular models arise from conditional independence assumptions made when factorizing the joint distribution of observed and latent data sequences, using the chain rule. The present paper is the first in-depth study regarding the use of DVAEs, as a general class of models, for unsupervised speech enhancement.

## III. DVAE AND SPEECH MODELING

In this section, we first review the standard VAE [11], [12] and its extensions to temporal models, which are referred to as DVAEs [30]. Then, we briefly introduce speech modeling using DVAE models. In the end of this section, we describe the practical implementation of three typical DVAE models.

### A. VAEs and DVAEs

To define a VAE, we assume that an observed variable $\mathbf{s}$ of dimension $F$ is generated from an unobserved random variable $\mathbf{z}$ of dimension $L$. This unobserved variable, also called latent variable, has a much lower dimension than the observed variable, i.e. $L \ll F$. Let $p_{\boldsymbol{\theta}}(\mathbf{s}, \mathbf{z}) = p_{\boldsymbol{\theta}_s}(\mathbf{s}|\mathbf{z})p_{\boldsymbol{\theta}_z}(\mathbf{z})$ be the parametric generative model of their joint distribution, where $\boldsymbol{\theta} = \boldsymbol{\theta}_s \cup \boldsymbol{\theta}_z$ denotes the set of parameters. Given a dataset $\mathbf{S} = \{\mathbf{s}_n\}_{n=1}^N$ which consists in $N$ i.i.d. samples of $\mathbf{s}$, we are typically interested in finding $\boldsymbol{\theta}$ that maximizes the log-likelihood of $\mathbf{S}$:

$$\log p_{\boldsymbol{\theta}}(\mathbf{S}) = \sum_{n=1}^N \log p_{\boldsymbol{\theta}}(\mathbf{s}_n) = \sum_{n=1}^N \log \int p_{\boldsymbol{\theta}_s}(\mathbf{s}_n|\mathbf{z}_n)p_{\boldsymbol{\theta}_z}(\mathbf{z}_n)d\mathbf{z}_n, \tag{1}$$

where $\mathbf{z}_n$ is the latent variable corresponding to $\mathbf{s}_n$. In general, a latent vector $\mathbf{z}$ is generated from a very simple prior distribution, typically $p_{\boldsymbol{\theta}_z}(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$,[2] whereas $\mathbf{s}$ is generated by a complex nonlinear function of $\mathbf{z}$, typically a deep neural network (DNN) (and $\boldsymbol{\theta}_s$ is the set of parameters of this DNN). The complexity of the resulting likelihood $p_{\boldsymbol{\theta}_s}(\mathbf{s}|\mathbf{z})$ makes this latent variable model expressive, but it also makes the marginal likelihood (or evidence) in (1) intractable. Therefore,

instead of directly maximizing $\log p_{\boldsymbol{\theta}}(\mathbf{S})$, an inference model $q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{s}) \approx p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{s})$ is introduced, which is also defined by a DNN (of parameters $\boldsymbol{\phi}$), and then the VAE is trained by maximizing the evidence lower bound (ELBO),[3] defined by (for a single pair of observed/latent vectors):

$$\begin{aligned}\mathcal{L}_{\boldsymbol{\theta},\boldsymbol{\phi}}(\mathbf{s}) &= \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{s})}\left[\log p_{\boldsymbol{\theta}}(\mathbf{s}, \mathbf{z}) - \log q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{s})\right] \\ &= \log p_{\boldsymbol{\theta}}(\mathbf{s}) - D_{KL}\left(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{s})||p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{s})\right),\end{aligned} \tag{2}$$

where $D_{KL}(.)$ denotes the Kullback-Leibler (KL) divergence. The KL divergence is always non-negative and thus the ELBO is a lower bound of the intractable log-marginal likelihood of the data. Using the reparameterization trick described in [11], [12], the generative model $p_{\boldsymbol{\theta}}(\mathbf{s}|\mathbf{z})$ and the inference model $q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{s})$ can be trained jointly by maximizing the ELBO over the training dataset.

While standard VAEs assume statistical independence among observations, DVAEs can be seen as a generalization of conventional VAEs for modeling sequential data [30]. A DVAE keeps the global encoder-decoder architecture of the VAE, but considers a sequence of observed random variables $\mathbf{s}_{1:T} = \{\mathbf{s}_t \in \mathbb{R}^F\}_{t=1}^T$ and a corresponding sequence of latent variables $\mathbf{z}_{1:T} = \{\mathbf{z}_t \in \mathbb{R}^L\}_{t=1}^T$. A DVAE is thus defined by the joint probability density function (pdf) of observed and latent variables $p_{\boldsymbol{\theta}}(\mathbf{s}_{1:T}, \mathbf{z}_{1:T})$, which can be factorized using the chain rule:[4]

$$p_{\boldsymbol{\theta}}(\mathbf{s}_{1:T}, \mathbf{z}_{1:T}) = \prod_{t=1}^T p_{\boldsymbol{\theta}}(\mathbf{s}_t, \mathbf{z}_t|\mathbf{s}_{1:t-1}, \mathbf{z}_{1:t-1}) \tag{3}$$

$$= \prod_{t=1}^T p_{\boldsymbol{\theta}_s}(\mathbf{s}_t|\mathbf{s}_{1:t-1}, \mathbf{z}_{1:t})p_{\boldsymbol{\theta}_z}(\mathbf{z}_t|\mathbf{s}_{1:t-1}, \mathbf{z}_{1:t-1}). \tag{4}$$

Similar to VAEs, the exact posterior distribution $p_{\boldsymbol{\theta}}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})$ is not analytically tractable. Consequently, an approximate posterior distribution $q_{\boldsymbol{\phi}}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})$ is introduced and it can be factorized using the chain rule as:

$$q_{\boldsymbol{\phi}}(\mathbf{z}_{1:T}|\mathbf{s}_{1:T}) = \prod_{t=1}^T q_{\boldsymbol{\phi}}(\mathbf{z}_t|\mathbf{z}_{1:t-1}, \mathbf{s}_{1:T}). \tag{5}$$

Chaining the inference and generation, the training of DVAEs still aims to maximize the ELBO of observed data, which is here defined by (for a single observed and latent data sequence):

$$\begin{aligned}\mathcal{L}(\boldsymbol{\phi}, \boldsymbol{\theta}) = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{1:T}|\mathbf{s}_{1:T})}[&\ln p_{\boldsymbol{\theta}}(\mathbf{s}_{1:T}, \mathbf{z}_{1:T}) \\ &- \ln q_{\boldsymbol{\phi}}(\mathbf{z}_{1:T}|\mathbf{s}_{1:T})].\end{aligned} \tag{6}$$

Note that when writing a joint distribution as a product of conditional distributions using the chain rule, a specific ordering of the variables has to be chosen. Among different possibilities, we chose a *causal* ordering to write the factorization in (3) and (4): The generation of $\mathbf{s}_t$ and $\mathbf{z}_t$ uses their past values $\mathbf{s}_{1:t-1}$ and $\mathbf{z}_{1:t-1}$ (plus $\mathbf{z}_t$ for generating $\mathbf{x}_t$). In

---

[2] $\mathcal{N}$ denotes the multivariate Gaussian distribution. We can remark that, here, $\boldsymbol{\theta}_{\mathbf{z}} = \emptyset$.

[3] Also called variational lower bound (VLB) or variational free energy (VFE).

[4] Here and in all the following, we take the convention that $\mathbf{s}_{1:0} = \mathbf{z}_{1:0} = \emptyset$. For $t = 1$ the first term of the product in (3) and (4) is thus $p_{\boldsymbol{\theta}}(\mathbf{s}_1, \mathbf{z}_1)$ and $p_{\boldsymbol{\theta}}(\mathbf{s}_1|\mathbf{z}_1)p_{\boldsymbol{\theta}}(\mathbf{z}_1)$, respectively.

the DVAE literature, almost all models are causal[5] [30]. Each of them is easily expressed as a special case of the general expression (4) where the dependencies in $p_{\boldsymbol{\theta}}(\mathbf{s}_t|\mathbf{s}_{1:t-1},\mathbf{z}_{1:t})$ and $p_{\boldsymbol{\theta}}(\mathbf{z}_t|\mathbf{s}_{1:t-1},\mathbf{z}_{1:t-1})$ are simplified, which may also affect the inference model $q_{\boldsymbol{\phi}}(\mathbf{z}_{1:T}|\mathbf{s}_{1:T})$ in (5). In addition, a given DVAE model can have different implementation with various types of deep neural networks, see [30] for an extensive discussion on this topic.

### B. Speech Modeling Using DVAEs

The VAE and the DVAE class of models have been used to model different kinds of data. In this subsection, we comment on how to use DVAEs for the modeling of speech signals in the short-term Fourier transform (STFT) domain. Fig. 1 illustrates this process. Let $\mathbf{s}_{1:T} = \{\mathbf{s}_t \in \mathbb{C}^F\}_{t=1}^T$ denote a sequence of complex-valued STFT frames, where $t$ is the time-frame index. Each vector $\mathbf{s}_t = \{s_{ft} \in \mathbb{C}\}_{f=1}^F$ represents the speech short-term spectrum at time index $t$, and the index $f$ is the frequency bin. As indicated above, $\mathbf{s}_{1:T}$ is associated with a sequence of latent variables $\mathbf{z}_{1:T} = \{\mathbf{z}_t \in \mathbb{R}^L\}_{t=1}^T$, which has the same sequence length $T$ and a much lower dimension $L \ll F$.

In speech and audio processing, the Fourier coefficients in $\mathbf{s}_t \in \mathbb{C}^F$ are usually assumed to be independent and distributed according to a complex Gaussian circularly symmetric distribution [40] (denoted below by $\mathcal{N}_c$), whose variance vary over time and frequency [5], [41]. The circularly symmetric assumption means that the phase follows a uniform distribution in $[0, 2\pi)$. This local Gaussian model thus exploits the non-stationarity and non-whiteness of the signal in the STFT domain, and it is equivalent to considering that the time-domain audio signal is a locally stationary Gaussian process [42]. Thus, for all time frames $t \in \{1, \ldots, T\}$, the DVAE generative model of speech signal is defined as follows:

$$p_{\boldsymbol{\theta}_{\mathbf{s}}}(\mathbf{s}_t|\mathbf{s}_{1:t-1},\mathbf{z}_{1:t}) = \mathcal{N}_c\left(\mathbf{s}_t; \mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\theta}_{\mathbf{s}},t}\right), \quad (7)$$

$$p_{\boldsymbol{\theta}_{\mathbf{z}}}(\mathbf{z}_t|\mathbf{s}_{1:t-1},\mathbf{z}_{1:t-1}) = \mathcal{N}\left(\mathbf{z}_t; \boldsymbol{\mu}_{\boldsymbol{\theta}_{\mathbf{z}},t}, \boldsymbol{\Sigma}_{\boldsymbol{\theta}_{\mathbf{z}},t}\right), \quad (8)$$

where the diagonal covariance matrix $\boldsymbol{\Sigma}_{\boldsymbol{\theta}_{\mathbf{s}},t} = \text{diag}\{\boldsymbol{v}_{\boldsymbol{\theta}_{\mathbf{s}},t}\}$ is provided by a DNN which takes as input the conditioning variables in (7), namely $(\mathbf{s}_{1:t-1},\mathbf{z}_{1:t-1})$. Similarly, $\boldsymbol{\mu}_{\boldsymbol{\theta}_{\mathbf{z}},t}$ and $\boldsymbol{\Sigma}_{\boldsymbol{\theta}_{\mathbf{z}},t} = \text{diag}\{\boldsymbol{v}_{\boldsymbol{\theta}_{\mathbf{z}},t}\}$ are provided by a DNN which takes as input the conditioning variables in (8), namely $\mathbf{s}_{1:t-1},\mathbf{z}_{1:t-1}$.[6] We collectively denote by $\boldsymbol{\theta} = \{\boldsymbol{\theta}_{\mathbf{s}}, \boldsymbol{\theta}_{\mathbf{z}}\}$ the parameters of the neural networks involved in (7)-(8).

Similarly as for the generative model, the general form of the inference model in a DVAE is given by:

$$q_{\boldsymbol{\phi}}(\mathbf{z}_t|\mathbf{z}_{1:t-1},\mathbf{s}_{1:T}) = \mathcal{N}(\mathbf{z}_t; \boldsymbol{\mu}_{\boldsymbol{\phi},t}, \boldsymbol{\Sigma}_{\boldsymbol{\phi},t}), \quad (9)$$

where $\boldsymbol{\mu}_{\boldsymbol{\phi},t}$ and $\boldsymbol{\Sigma}_{\boldsymbol{\phi},t} = \text{diag}\{\boldsymbol{v}_{\boldsymbol{\phi},t}\}$ are provided by a neural network which takes as input $\mathbf{z}_{1:t-1},\mathbf{s}_{1:T}$ and whose parameters are denoted by $\boldsymbol{\phi}$.

[5]Except for one called Recurrent VAE (RVAE) that has a *non-causal* form, we will specify it later in the speech enhancement section.

[6]It is important to note that the parameters of the distributions involved in a DVAE are always functions of the variables that come after the conditioning bar. In the rest of the paper, we will generally omit to rewrite these variables in the right-hand-side of the probabilistic modeling equations, for concision, but we may punctually make these dependencies explicit when it eases the understanding.

It can be noted that even if the complex-valued vector sequence $\mathbf{s}_{1:t-1}$ or $\mathbf{s}_{1:T}$ is used as a conditioning variable in the distributions in (7)-(9), in practice we use the modulus-squared values of these variables at the encoder and decoder input. In other words, the DVAE encoder and decoder distribution parameters are computed using sequences of vectors with entries equal to $|s_{ft}|^2$, as illustrated in Fig. 1. Note that modulus-squared of data is homogeneous with the decoder output (data variance).

Given the generative model (7), (8) and the inference model (9), we can develop the objective ELBO function (for one data sequence):

$$\mathcal{L}(\mathbf{s}_{1:T}; \boldsymbol{\theta}, \boldsymbol{\phi})$$

$$= \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{1:T}|\mathbf{s}_{1:T})}\left[\ln p_{\boldsymbol{\theta}}(\mathbf{s}_{1:T},\mathbf{z}_{1:T}) - \ln q_{\boldsymbol{\phi}}(\mathbf{z}_{1:T}|\mathbf{s}_{1:T})\right] \quad (10)$$

$$\overset{c}{=} -\sum_{f,t=1}^{F,T} \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{1:t}|\mathbf{s}_{1:T})}\left[d_{IS}(|s_{ft}|^2, v_{\boldsymbol{\theta}_{\mathbf{s}},ft})\right]$$

$$+ \frac{1}{2}\sum_{l,t=1}^{L,T}\left[\ln\frac{v_{\boldsymbol{\phi},lt}}{v_{\boldsymbol{\theta}_{\mathbf{z}},lt}} - \frac{v_{\boldsymbol{\phi},lt} + (\mu_{\boldsymbol{\phi},lt} - \mu_{\boldsymbol{\theta}_{\mathbf{z}},lt})^2}{v_{\boldsymbol{\theta}_{\mathbf{z}},lt}}\right], \quad (11)$$

where $\overset{c}{=}$ denotes equality up to an additive constant w.r.t. $\boldsymbol{\theta}$ and $\boldsymbol{\phi}$, $d_{IS}(q,p) = q/p - \ln(q/p) - 1$ is the Itakura-Saito (IS) divergence [20], $v_{\boldsymbol{\theta}_{\mathbf{s}},ft} \in \mathbb{R}_+$ is the $f$-th entry of $\boldsymbol{v}_{\boldsymbol{\theta}_{\mathbf{s}},t}$, $\{\mu_{\boldsymbol{\phi},lt} \in \mathbb{R}, \mu_{\boldsymbol{\theta}_{\mathbf{z}},lt} \in \mathbb{R}, v_{\boldsymbol{\phi},lt} \in \mathbb{R}_+, v_{\boldsymbol{\theta}_{\mathbf{z}},lt} \in \mathbb{R}_+\}$ are the $l$-th entry of $\{\boldsymbol{\mu}_{\boldsymbol{\phi},t}, \boldsymbol{\mu}_{\boldsymbol{\theta}_{\mathbf{z}},t}, \boldsymbol{v}_{\boldsymbol{\phi},t}, \boldsymbol{v}_{\boldsymbol{\theta}_{\mathbf{z}},t}\}$, respectively.

### C. DVAEs in practice

In this subsection, we briefly present three representative DVAE models from the literature (see the review in [30]) that can be plugged into the proposed speech enhancement method:

- DKF [23], [24], the simplest DVAE model which is inspired from conventional state-space models;
- RVAE [28], the only model that was presented in the literature with a non-causal form [30];
- SRNN [26], a more complex DVAE model, with an autoregressive generative modeling of the speech signal.

*1) DKF:* DKF generative model follows the structure of a standard state-space model with a first-order Markov model on the latent vectors, i.e. $\mathbf{z}_t$ is generated only from $\mathbf{z}_{t-1}$, and an instantaneous observation model, i.e. $\mathbf{s}_t$ is generated only from $\mathbf{z}_t$. In summary, DKF generative and inference models are defined by:

**Generation:**

$$p_{\boldsymbol{\theta}}(\mathbf{s}_{1:T},\mathbf{z}_{1:T}) = \prod_{t=1}^T p_{\boldsymbol{\theta}_{\mathbf{s}}}(\mathbf{s}_t|\mathbf{z}_t)p_{\boldsymbol{\theta}_{\mathbf{z}}}(\mathbf{z}_t|\mathbf{z}_{t-1}). \quad (12)$$

**Inference:**

$$q_{\boldsymbol{\phi}}(\mathbf{z}_{1:T}|\mathbf{s}_{1:T}) = \prod_{t=1}^T q_{\boldsymbol{\phi}}(\mathbf{z}_t|\mathbf{z}_{t-1},\mathbf{s}_{t:T}). \quad (13)$$

*2) RVAE:* As opposed to DKF, RVAE does not consider any dynamical model for the latent vectors, which are assumed i.i.d., with $p(\mathbf{z}_t) = \mathcal{N}(\mathbf{z}_t; \mathbf{0}, \mathbf{I})$. This prior therefore lacks parameters and does not involve a neural network, as in standard VAEs. However, the observation model in RVAE is
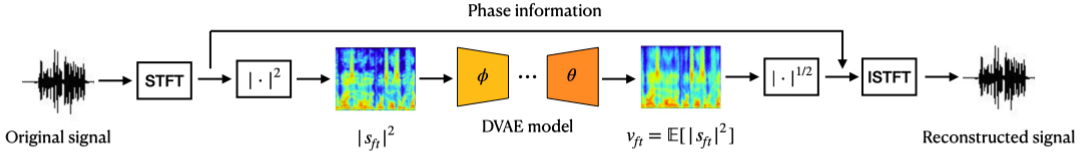
Fig. 1. Speech analysis-resynthesis with DVAEs, in the STFT domain. The speech power spectrogram is used as input to the DVAE, and the output is the variance of a complex Gaussian model in the STFT domain. The audio waveform is reconstructed by inverse STFT using the phase of the original signal.
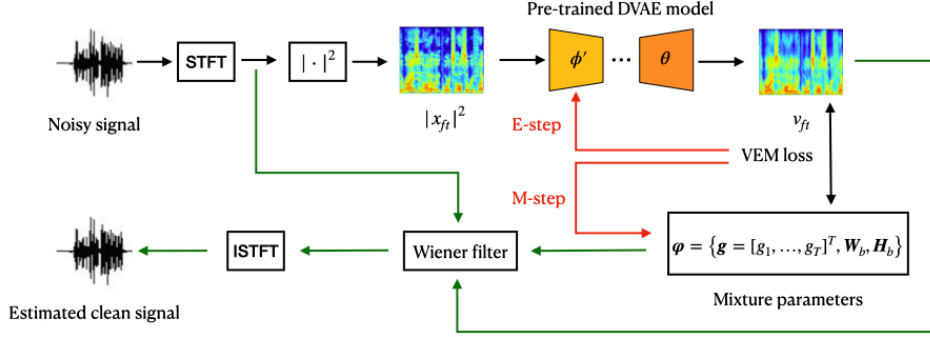


Fig. 2. Overview of the proposed speech enhancement method. The pre-trained DVAE is used within a VEM algorithm for speech enhancement (red arrows). During the E-step the DVAE encoder is fine-tuned (the decoder is fixed), while in the M-step the mixture parameters are estimated. The clean speech signal is estimated by filtering the noisy signal with a Wiener filtering combining the estimated mixture parameters with the DVAE output parameters (green arrows).

more complex than in DKF: in its causal form, $\mathbf{s}_t$ is generated from the current and previous latent variables $\mathbf{z}_{1:t}$, while in its non-causal form, $\mathbf{s}_t$ is generated from the complete sequence of latent vectors $\mathbf{z}_{1:T}$. In short, we have:

**Generation:**

$$p_{\boldsymbol{\theta}}(\mathbf{s}_{1:T}, \mathbf{z}_{1:T}) = \prod_{t=1}^{T} p_{\boldsymbol{\theta}_{\mathbf{s}}}(\mathbf{s}_t | \mathbf{z}_{1:t}) p(\mathbf{z}_t) \quad \text{(causal).} \tag{14}$$

$$p_{\boldsymbol{\theta}}(\mathbf{s}_{1:T}, \mathbf{z}_{1:T}) = \prod_{t=1}^{T} p_{\boldsymbol{\theta}_{\mathbf{s}}}(\mathbf{s}_t | \mathbf{z}_{1:T}) p(\mathbf{z}_t) \quad \text{(non-causal).} \tag{15}$$

**Inference:**

$$q_{\boldsymbol{\phi}}(\mathbf{z}_{1:T} | \mathbf{s}_{1:T}) = \prod_{t=1}^{T} q_{\boldsymbol{\phi}}(\mathbf{z}_t | \mathbf{z}_{1:t-1}, \mathbf{s}_{t:T}) \quad \text{(causal).} \tag{16}$$

$$q_{\boldsymbol{\phi}}(\mathbf{z}_{1:T} | \mathbf{s}_{1:T}) = \prod_{t=1}^{T} q_{\boldsymbol{\phi}}(\mathbf{z}_t | \mathbf{z}_{1:t-1}, \mathbf{s}_{1:T}) \quad \text{(non-causal).} \tag{17}$$

*3) SRNN:* In SRNN, the latent vector $\mathbf{z}_t$ is generated not only from $\mathbf{z}_{t-1}$ (as in DKF), but also from $\mathbf{s}_{1:t-1}$. For the observation model, $\mathbf{s}_t$ is generated not only from $\mathbf{z}_t$ (as in DKF) but also from $\mathbf{s}_{1:t-1}$, so that SRNN actually corresponds to an autoregressive model. The generative and inference models are defined by:

**Generation:**

$$p_{\boldsymbol{\theta}}(\mathbf{s}_{1:T}, \mathbf{z}_{1:T}) = \prod_{t=1}^{T} p_{\boldsymbol{\theta}_{\mathbf{s}}}(\mathbf{s}_t | \mathbf{z}_t, \mathbf{s}_{1:t-1}) p_{\boldsymbol{\theta}_{\mathbf{z}}}(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{s}_{1:t-1}). \tag{18}$$

**Inference:**

$$q_{\boldsymbol{\phi}}(\mathbf{z}_{1:T} | \mathbf{s}_{1:T}) = \prod_{t=1}^{T} q_{\boldsymbol{\phi}}(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{s}_{1:T}). \tag{19}$$

## IV. DVAE FOR SPEECH ENHANCEMENT

This section describes the DVAE-based speech enhancement algorithm, where the clean speech is modeled with DVAE and the noise is modeled with nonnegative matrix factorization (NMF) [20]. It is an extended version of the algorithm proposed in [28] for the RVAE model, with a more general formulation which can be used for speech enhancement with any other DVAE model. The proposed method is illustrated in Fig. 2. In practice, we assume that the DVAE-based clean speech generative and inference models defined in (7)-(8) and (9), respectively, have been learned, i.e. the associated parameters $\boldsymbol{\theta}$ and $\boldsymbol{\phi}$ have been estimated from a dataset of clean speech signals during the training stage (see Section III-B). The objective of speech enhancement is to use this pre-trained model to estimate the clean speech signal when only the noisy mixture is observed. This is done with a variational expectation maximization (VEM) algorithm. We recall that this method is unsupervised, since no clean-noisy speech pairs are used in this method.

In the rest of this section, we first introduce the unsupervised noise and mixture models, then we develop the general strategy to estimate the clean speech signal modeled by a DVAE when only the mixture signal is available, and finally we present the algorithm to estimate the remaining unknown model parameters. Throughout this section, $\mathbf{s}_{1:T} = \{\mathbf{s}_t \in \mathbb{C}^F\}_{t=1}^{T}$, $\mathbf{b}_{1:T} = \{\mathbf{b}_t \in \mathbb{C}^F\}_{t=1}^{T}$, and $\mathbf{x}_{1:T} = \{\mathbf{x}_t \in \mathbb{C}^F\}_{t=1}^{T}$ denote respectively a sequence of STFT clean speech frames, noise frames, and noisy speech frames.

## A. Noise and mixture models

As in [14], [28], we consider a Gaussian noise model with NMF parameterization of the variance [20]. Independently for all time frames $t \in \{1, ..., T\}$, we define:

$$p(\mathbf{b}_t) = \mathcal{N}_c(\mathbf{b}_t; \mathbf{0}, \boldsymbol{\Sigma}_{\mathbf{b},t}), \tag{20}$$

where $\boldsymbol{\Sigma}_{\mathbf{b},t} = \mathrm{diag}\{(\mathbf{W}_b\mathbf{H}_b)_{:,t}\}$ with $\mathbf{W}_b \in \mathbb{R}_+^{F \times K}$ and $\mathbf{H}_b \in \mathbb{R}_+^{K \times T}$. The rank of the factorization $K$ is usually chosen such that $K(F + T) \ll FT$.

We consider that the noisy speech is a mixture of the noise defined in (20) and the clean speech defined in (7), (8):

$$\mathbf{x}_t = \sqrt{g_t}\mathbf{s}_t + \mathbf{b}_t, \tag{21}$$

where $g_t \in \mathbb{R}_+$ is a gain parameter scaling the level of the speech signal at each time frame. It is frame-dependent but frequency-independent. The gain parameter enables to take into account the potentially different loudness between the speech training examples used to learn the (D)VAE parameters and the speech test samples in the noisy sequence that we have to denoise [14].

From (7), (20) and (21), and by assuming the independence of the speech and noise signals, we have for all $t \in \{1, ..., T\}$:

$$p_{\boldsymbol{\theta}_{\mathbf{x}}}(\mathbf{x}_t \mid \mathbf{s}_{1:t-1}, \mathbf{z}_{1:t}) = \mathcal{N}_c(\mathbf{x}_t; \mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\theta}_{\mathbf{x}},t}), \tag{22}$$

where $\boldsymbol{\Sigma}_{\boldsymbol{\theta}_{\mathbf{x}},t} = \mathrm{diag}\{g_t\boldsymbol{v}_{\boldsymbol{\theta}_{\mathbf{s}},t}+(\mathbf{W}_b\mathbf{H}_b)_{:,t}\}$ and $\boldsymbol{\theta}_{\mathbf{x}}$ is the union of the speech generative model parameters $\boldsymbol{\theta}_{\mathbf{s}}$ and the mixture model parameters $\boldsymbol{\varphi} = \{\boldsymbol{g} = [g_1, \ldots, g_T]^T, \mathbf{W}_b, \mathbf{H}_b\}$. As already mentioned in Section III-B, $\boldsymbol{v}_{\boldsymbol{\theta}_{\mathbf{s}},t}$ is actually a function of $(\mathbf{s}_{1:t-1}, \mathbf{z}_{1:t})$. It can be noted that from (20) and (21), it is clear that given the clean speech frame $\mathbf{s}_t$, the noisy speech frame $\mathbf{x}_t$ is characterized by:

$$p_{\boldsymbol{\varphi}}(\mathbf{x}_t \mid \mathbf{s}_t) = \mathcal{N}_c(\mathbf{x}_t; \sqrt{g_t}\mathbf{s}_t, \boldsymbol{\Sigma}_{\mathbf{b},t}). \tag{23}$$

## B. Speech reconstruction

Now we consider the problem of reconstructing the clean speech signal from the observed mixture signal, which consists in computing the following posterior mean:

$$\hat{\mathbf{s}}_t = \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{s}_t|\mathbf{x}_{1:T})}[\mathbf{s}_t]. \tag{24}$$

In practice, we cannot write the posterior $p_{\boldsymbol{\theta}}(\mathbf{s}_t|\mathbf{x}_{1:T})$ analytically, which makes the above expectation intractable. However, leveraging the speech model defined previously, we can approximate it by introducing random variables that are then marginalized.

*1) Introducing the past and current latent variables:* We start from marginalizing with respect to the latent variables $\mathbf{z}_{1:t}$:

$$p_{\boldsymbol{\theta}}(\mathbf{s}_t|\mathbf{x}_{1:T}) = \int p_{\boldsymbol{\theta}}(\mathbf{s}_t|\mathbf{z}_{1:t}, \mathbf{x}_{1:T})p_{\boldsymbol{\theta}}(\mathbf{z}_{1:t}|\mathbf{x}_{1:T})d\mathbf{z}_{1:t}$$

$$= \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{z}_{1:t}|\mathbf{x}_{1:T})}[p_{\boldsymbol{\theta}}(\mathbf{s}_t|\mathbf{z}_{1:t}, \mathbf{x}_{1:T})]. \tag{25}$$

Using (25) to rewrite (24), the estimate of the clean speech signal at time $t$ is given by:

$$\hat{\mathbf{s}}_t = \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{z}_{1:t}|\mathbf{x}_{1:T})}\left[\mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{s}_t|\mathbf{z}_{1:t}, \mathbf{x}_{1:T})}[\mathbf{s}_t]\right]. \tag{26}$$

Let us now focus on the inner expectation, taken with respect to $p_{\boldsymbol{\theta}}(\mathbf{s}_t|\mathbf{z}_{1:t}, \mathbf{x}_{1:T})$. We will come back later on the outer expectation taken with respect to $p_{\boldsymbol{\theta}}(\mathbf{z}_{1:t}|\mathbf{x}_{1:T})$. Using Bayes rule, we have:

$$p_{\boldsymbol{\theta}}(\mathbf{s}_t|\mathbf{z}_{1:t}, \mathbf{x}_{1:T}) = \frac{p_{\boldsymbol{\theta}}(\mathbf{x}_{1:T}|\mathbf{s}_t, \mathbf{z}_{1:t})p_{\boldsymbol{\theta}}(\mathbf{s}_t|\mathbf{z}_{1:t})p_{\boldsymbol{\theta}}(\mathbf{z}_{1:t})}{p_{\boldsymbol{\theta}}(\mathbf{z}_{1:t}, \mathbf{x}_{1:T})} \tag{27}$$

$$\propto p_{\boldsymbol{\theta}}(\mathbf{x}_{1:T}|\mathbf{s}_t, \mathbf{z}_{1:t})p_{\boldsymbol{\theta}}(\mathbf{s}_t|\mathbf{z}_{1:t}) \tag{28}$$

$$\approx p_{\boldsymbol{\theta}}(\mathbf{x}_t|\mathbf{s}_t)p_{\boldsymbol{\theta}}(\mathbf{s}_t|\mathbf{z}_{1:t}). \tag{29}$$

The exact computation of $p_{\boldsymbol{\theta}}(\mathbf{x}_{1:T}|\mathbf{s}_t, \mathbf{z}_{1:t})$ requires the marginalisation of $p_{\boldsymbol{\theta}}(\mathbf{x}_{1:T}, \mathbf{s}_{1:t-1,t+1:T}, \mathbf{z}_{t+1:T}|\mathbf{s}_t, \mathbf{z}_{1:t})$ w.r.t. the undesired variables. This would require not only marginalising from future latent codes, but also from past and future clean speech, which is clearly not feasible. Instead, we consider only the signal mixture model, $p_{\boldsymbol{\theta}}(\mathbf{x}_t|\mathbf{s}_t)$, defined in (23).

*2) Introducing the past speech vectors:* Then, it comes to estimating $p_{\boldsymbol{\theta}}(\mathbf{s}_t|\mathbf{z}_{1:t})$ in (29). To do so, we introduce and then marginalize the past speech vectors $\mathbf{s}_{1:t-1}$:

$$p_{\boldsymbol{\theta}}(\mathbf{s}_t|\mathbf{z}_{1:t}) = \int p_{\boldsymbol{\theta}}(\mathbf{s}_t|\mathbf{s}_{1:t-1}, \mathbf{z}_{1:t})p_{\boldsymbol{\theta}}(\mathbf{s}_{1:t-1}|\mathbf{z}_{1:t})d\mathbf{s}_{1:t-1}$$

$$= \int p_{\boldsymbol{\theta}}(\mathbf{s}_t|\mathbf{s}_{1:t-1}, \mathbf{z}_{1:t})\left[\prod_{\tau=1}^{t-1} p_{\boldsymbol{\theta}}(\mathbf{s}_\tau|\mathbf{s}_{1:\tau-1}, \mathbf{z}_{1:\tau})\right]d\mathbf{s}_{1:t-1}$$

$$= \mathbb{E}_{\prod_{\tau=1}^{t-1} p_{\boldsymbol{\theta}}(\mathbf{s}_\tau|\mathbf{s}_{1:\tau-1}, \mathbf{z}_{1:\tau})}\left[p_{\boldsymbol{\theta}}(\mathbf{s}_t|\mathbf{s}_{1:t-1}, \mathbf{z}_{1:t})\right], \tag{30}$$

where in the second line we used the fact that $s_\tau$ is independent of $z_{\tau+1:t}$ for $\tau < t$.

When computing (30), we are facing two problems: First the expectation is intractable; and second, in a speech enhancement framework, we do not have access to the past ground-truth clean speech vectors $\mathbf{s}_{1:t-1}$.[7] Therefore, we approximate

$$p_{\boldsymbol{\theta}}(\mathbf{s}_t|\mathbf{z}_{1:t}) \approx p_{\boldsymbol{\theta}}(\mathbf{s}_t|\tilde{\mathbf{s}}_{1:t-1}, \mathbf{z}_{1:t})$$
$$= \mathcal{N}_c(\mathbf{s}_t; \mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\theta}_{\mathbf{s}},t}(\tilde{\mathbf{s}}_{1:t-1}, \mathbf{z}_{1:t})), \tag{31}$$

where $\boldsymbol{\Sigma}_{\boldsymbol{\theta}_{\mathbf{s}},t}(\tilde{\mathbf{s}}_{1:t-1}, \mathbf{z}_{1:t}) = \mathrm{diag}\{\boldsymbol{v}_{\boldsymbol{\theta}_{\mathbf{s}},t}(\tilde{\mathbf{s}}_{1:t-1}, \mathbf{z}_{1:t})\}$ and $\tilde{\mathbf{s}}_t$ is computed recursively as $\tilde{\mathbf{s}}_t = \boldsymbol{v}_{\boldsymbol{\theta}_{\mathbf{s}},t}(\tilde{\mathbf{s}}_{1:t-1}, \mathbf{z}_{1:t})$.[8] In practice, it means that the decoder output at time frame $t-1$ is re-injected at the decoder input at the next time frame $t$. This part of the process is only necessary for SRNN, and more generally for any DVAE autoregressive model. For non-autoregressive DVAE models, such as RVAE and DKF, $\boldsymbol{\Sigma}_{\boldsymbol{\theta}_{\mathbf{s}},t}$ is only computed from the sequence of latent vectors.

*3) Computing the conditional posterior:* Substituting (23) and (31) into (29), we have:

$$p_{\boldsymbol{\theta}}(\mathbf{s}_t|\mathbf{z}_{1:t}, \mathbf{x}_{1:T}) \approx \mathcal{N}_c(\mathbf{x}_t; \sqrt{g_t}\mathbf{s}_t, \boldsymbol{\Sigma}_{\mathbf{b},t})\mathcal{N}_c(\mathbf{s}_t; \mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\theta}_{\mathbf{s}},t})$$
$$= \mathcal{N}_c(\mathbf{s}_t; \mathbf{m}_{\mathbf{s},t}, \boldsymbol{\Sigma}_{\mathbf{s},t}), \tag{32}$$

where

$$\mathbf{m}_{\mathbf{s},t} = \sqrt{g_t}\boldsymbol{\Sigma}_{\boldsymbol{\theta}_{\mathbf{s}},t}(g_t\boldsymbol{\Sigma}_{\boldsymbol{\theta}_{\mathbf{s}},t} + \boldsymbol{\Sigma}_{\mathbf{b},t})^{-1}\mathbf{x}_t, \tag{33}$$

$$\boldsymbol{\Sigma}_{\mathbf{s},t} = \boldsymbol{\Sigma}_{\boldsymbol{\theta}_{\mathbf{s}},t}\boldsymbol{\Sigma}_{\mathbf{b},t}(g_t\boldsymbol{\Sigma}_{\boldsymbol{\theta}_{\mathbf{s}},t} + \boldsymbol{\Sigma}_{\mathbf{b},t})^{-1}. \tag{34}$$

---

[7]as opposed to the DVAE training procedure which is done using sequences of clean speech signals.

[8]Here we explicitly write the dependency of the covariance matrix on $\tilde{\mathbf{s}}_{1:t-1}$ and $\mathbf{z}_{1:t}$, to make the use of the DVAE model clear. In the following we will omit it again for concision of presentation.

Finally, from (26), (32) and (33), the estimate of the clean speech signal is given by:

$$\hat{\mathbf{s}}_t \approx \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{z}_{1:t}|\mathbf{x}_{1:T})} \left[ \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{s}_t|\mathbf{z}_{1:t},\mathbf{x}_{1:T})}[s_t] \right]$$
$$\approx \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{z}_{1:t}|\mathbf{x}_{1:T})} \left[ \sqrt{g_t} \boldsymbol{\Sigma}_{\boldsymbol{\theta}_{\mathbf{s}},t} \left( g_t \boldsymbol{\Sigma}_{\boldsymbol{\theta}_{\mathbf{s}},t} + \boldsymbol{\Sigma}_{\mathbf{b},t} \right)^{-1} \right] \mathbf{x}_t, \quad (35)$$

where we recall that $\boldsymbol{\Sigma}_{\boldsymbol{\theta}_{\mathbf{s}},t}$ is actually a function of $(\tilde{\mathbf{s}}_{1:t-1}, \mathbf{z}_{1:t})$. This estimate can be seen as a "probabilistic" Wiener filter, i.e. a Wiener filter averaged over all possible realizations of the latent variables according to their posterior distribution $p_{\boldsymbol{\theta}}(\mathbf{z}_{1:t}|\mathbf{x}_{1:T})$.

The expectation in (35) is intractable, but similarly as before we can approximate it by

$$\hat{\mathbf{s}}_t \approx \sqrt{g_t} \boldsymbol{\Sigma}_{\boldsymbol{\theta}_{\mathbf{s}},t} \left( g_t \boldsymbol{\Sigma}_{\boldsymbol{\theta}_{\mathbf{s}},t} + \boldsymbol{\Sigma}_{\mathbf{b},t} \right)^{-1} \mathbf{x}_t, \quad (36)$$

where $\boldsymbol{\Sigma}_{\boldsymbol{\theta}_{\mathbf{s}},t} = \boldsymbol{\Sigma}_{\boldsymbol{\theta}_{\mathbf{s}},t}(\tilde{\mathbf{s}}_{1:t-1}, \mathbf{z}_{1:t})$ and $\tilde{\mathbf{z}}_{1:t}$ is sampled from $p_{\boldsymbol{\theta}}(\mathbf{z}_{1:t}|\mathbf{x}_{1:T})$. In practice, this posterior distribution is also intractable, we thus propose a variational approximation $q_{\boldsymbol{\phi}'}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})$ whose parameters $\phi'$ need to be jointly estimated together with the noisy mixture model parameters $\varphi$, in order to compute the speech estimate in (36). As detailed in the next section, we propose a VEM algorithm to do that. This generalizes the algorithm developed for RVAE in [28] to the whole class of DVAE models.

### C. VEM algorithm for model parameters estimation

Now that we have an expression for the clean speech signal estimate, what remains to be estimated is the set of mixture model parameters $\varphi$ (the NMF noise model parameters and the gains) and the parameters $\phi'$ of the variational distribution $q_{\boldsymbol{\phi}'}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})$. Using a VEM algorithm, we will maximize the following ELBO defined from the noisy speech observations $\mathbf{x}_{1:T}$:

$$\mathcal{L}(\phi', \varphi) = \mathbb{E}_{q_{\boldsymbol{\phi}'}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})} [\ln p_{\varphi}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T})$$
$$- \ln q_{\boldsymbol{\phi}'}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})]. \quad (37)$$

It can be shown [31] that this corresponds to (i) maximizing with respect to $\varphi$ a lower bound of the intractable log-marginal likelihood $\ln p_{\varphi}(\mathbf{x}_{1:T})$, and (ii) minimizing with respect to $\phi'$ the KL divergence between the variational distribution $q_{\boldsymbol{\phi}'}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})$ and the intractable posterior $p_{\boldsymbol{\theta}}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})$. The proposed VEM algorithm thus consists in iterating between the following variational E and M steps.

*1) Variational E-step:* We consider a variational distribution of the same form as the DVAE inference model:

$$q_{\boldsymbol{\phi}'}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T}) = \prod_{t=1}^{T} q_{\boldsymbol{\phi}'}(\mathbf{z}_t|\mathbf{z}_{1:t-1}, \mathbf{x}_{1:T}), \quad (38)$$

where $q_{\boldsymbol{\phi}'}(\mathbf{z}_t|\mathbf{z}_{1:t-1}, \mathbf{x}_{1:T})$ is defined as in (9), except that $\mathbf{s}_t$ is replaced by $\mathbf{x}_t$. The noisy speech frames can be considered as out-of-sample data for the DVAE model trained on clean speech signals [43]. Therefore, similarly as in [28], we can simply fine-tune the pre-trained DVAE inference network on the noisy speech test signal, by maximizing the ELBO in (37) w.r.t. $\phi'$. This objective function can be developed by

marginalizing and sampling, similarly to what was done in the previous sub-section. This leads to the following expression:

$$\mathcal{L}(\phi', \varphi^{\star}) \stackrel{c}{=} - \sum_{f=1}^{F} \sum_{t=1}^{T} \mathbb{E}_{q_{\boldsymbol{\phi}'}} \left[ \ln v_{\boldsymbol{\varphi}^{\star},ft} + \frac{|x_{ft}|^2}{v_{\boldsymbol{\varphi}^{\star},ft}} \right] +$$
$$\frac{1}{2} \sum_{l=1}^{L} \sum_{t=1}^{T} \left[ \ln v_{\boldsymbol{\phi}',lt} - \ln v_{\boldsymbol{\theta}_{\mathbf{z}},lt} - \frac{v_{\boldsymbol{\phi}',lt} + (\mu_{\boldsymbol{\phi}',lt} - \mu_{\boldsymbol{\theta}_{\mathbf{z}},lt})^2}{v_{\boldsymbol{\theta}_{\mathbf{z}},lt}} \right], \quad (39)$$

where $x_{ft}$ denotes the $f$-th entry of $\mathbf{x}_t$ and $\varphi^{\star}$ denotes the current estimate of the mixture model parameters. In practice, the expectation applied on $v_{\boldsymbol{\varphi}^{\star},ft}$ is approximated by using

$$v_{\boldsymbol{\varphi}^{\star},ft} = g_t v_{\boldsymbol{\theta}_{\mathbf{s}},ft}(\tilde{\mathbf{s}}_{1:t-1}, \tilde{\mathbf{z}}_{1:t}) + (\mathbf{W}_b \mathbf{H}_b)_{ft}, \quad (40)$$

where we remind that $\tilde{\mathbf{s}}_{1:t-1}$ is computed recursively from the output of the decoder network as explained after equation (31), and $\tilde{\mathbf{z}}_{1:t}$ is recursively sampled from $q_{\boldsymbol{\phi}'}(\mathbf{z}_{1:t}|\mathbf{x}_{1:T}) = \prod_{\tau=1}^{t} q_{\boldsymbol{\phi}'}(\mathbf{z}_\tau|\mathbf{z}_{1:\tau-1}, \mathbf{x}_{1:T})$ as defined in (38). During the variational E-step, the parameters $\phi'$ are updated with variants of gradient ascent, and we denote by $\phi'^{\star}$ the resulting parameters that will be fixed in the M-step.

We recall that the recursive computation of $\tilde{\mathbf{s}}_t$ is only required for SRNN, and actually for any DVAE model that considers an autoregressive process for the generative speech model. DKF and RVAE, as non-autoregressive models, do not require estimating these quantities, only sampling the latent variables $\tilde{\mathbf{z}}_{1:t}$ is necessary.

*2) M-step:* The M-step consists in maximizing $\mathcal{L}(\phi'^{\star}, \varphi)$ wrt $\varphi$ under a non-negativity constraint. Replacing the intractable expectation in (39) with a Monte Carlo estimate (using actually one single sample), the M-step can be recast as minimizing the following criterion [14]:

$$\mathcal{C}(\varphi) = \sum_{f=1}^{F} \sum_{t=1}^{T} d_{\text{IS}} \left( |x_{ft}|^2, v_{\boldsymbol{\varphi},ft} \right), \quad (41)$$

where $v_{\boldsymbol{\varphi},ft}$ is defined in (40). This optimization problem can be tackled using a majorize-minimize approach [44], which leads to the multiplicative update rules derived in [14] using the methodology proposed in [45]:

$$\mathbf{H}_b \leftarrow \mathbf{H}_b \odot \left[ \frac{\mathbf{W}_b^{\top} \left( |\mathbf{X}|^{\odot 2} \odot (\mathbf{V}_{\mathbf{x}})^{\odot -2} \right)}{\mathbf{W}_b^{\top} (\mathbf{V}_{\mathbf{x}})^{\odot -1}} \right]^{\odot 1/2}; \quad (42)$$

$$\mathbf{W}_b \leftarrow \mathbf{W}_b \odot \left[ \frac{\left( |\mathbf{X}|^{\odot 2} \odot (\mathbf{V}_{\mathbf{x}})^{\odot -2} \right) \mathbf{H}_b^{\top}}{(\mathbf{V}_{\mathbf{x}})^{\odot -1} \mathbf{H}_b^{\top}} \right]^{\odot 1/2}; \quad (43)$$

$$\mathbf{g}^{\top} \leftarrow \mathbf{g}^{\top} \odot \left[ \frac{\mathbf{1}^{\top} \left[ |\mathbf{X}|^{\odot 2} \odot \left( \mathbf{V}_{\mathbf{s}} \odot (\mathbf{V}_{\mathbf{x}})^{\odot -2} \right) \right]}{\mathbf{1}^{\top} \left[ \left( \mathbf{V}_{\mathbf{s}} \odot (\mathbf{V}_{\mathbf{x}})^{\odot -1} \right) \right]} \right]^{\odot 1/2}, \quad (44)$$

where $\odot$ denotes element-wise multiplication and exponentiation, and matrix division is also element-wise, $\mathbf{V}_{\mathbf{s}}, \mathbf{V}_{\mathbf{x}} \in \mathbb{R}_{+}^{F \times T}$ are the matrices of entries $v_{\boldsymbol{\theta}_{\mathbf{s}},ft}$ and $v_{\boldsymbol{\varphi},ft}$ respectively, $\mathbf{X} \in \mathbb{C}^{F \times T}$ is the matrix of entries $x_{ft}$ and $\mathbf{1}$ is an all-ones column vector of dimension $F$. Note that non-negativity is

**Algorithm 1** DVAE-based unsupervised speech enhancement

**Inputs:**
▷ Pre-trained DVAE model: $p_{\boldsymbol{\theta}}(\mathbf{z}_{1:T}, \mathbf{s}_{1:T})$ and $q_{\boldsymbol{\phi}}(\mathbf{z}_{1:T}|\mathbf{s}_{1:T})$
▷ Noisy speech frames $\mathbf{x}_{1:T}$

**Initialization:**
▷ Initialize NMF noise parameters $\mathbf{H}_b$ and $\mathbf{W}_b$ with random nonnegative values
▷ Initialize gain parameters $\mathbf{g} = \mathbf{1}$
▷ Initialize $q_{\boldsymbol{\phi}'}(\mathbf{z}_{1:t}|\mathbf{x}_{1:T})$ with pre-trained inference network $q_{\boldsymbol{\phi}}(\mathbf{z}_{1:t}|\mathbf{s}_{1:T})$

**while** stopping criterion not reached **do**
    **E-step:**
    ▷ Fine-tune $q_{\boldsymbol{\phi}'}(\mathbf{z}_{1:t}|\mathbf{x}_{1:T})$ by maximizing (39) w.r.t. $\phi'$
    ▷ Sample $\tilde{\mathbf{z}}_{1:T}$ from $q_{\boldsymbol{\phi}'}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})$
    ▷ Compute $\boldsymbol{\Sigma}_{\boldsymbol{\theta}_{\mathbf{s}},t}$ for $t = 1$ to $T$ using the DVAE decoder
    **M-step:**
    ▷ Update $\mathbf{H}_b$, $\mathbf{W}_b$ and $\mathbf{g}$ using (42)-(44)
**end while**
**Output:**
▷ Compute the clean speech signal estimate $\hat{\mathbf{s}}_t$ for $t = 1$ to $T$ using (36)

---

ensured provided that these parameters are initialized with non-negative values.

### D. Summary

In summary, the clean speech signal estimation consists in approximating the posterior $p_{\boldsymbol{\theta}}(\mathbf{s}_t|\mathbf{x}_{1:T})$ and taking the mean of the resulting approximate distribution (i.e. the Wiener filter output). The estimation of the involved parameters is made with the VEM algorithm, which consists in iteratively fine-tuning the inference network of the pre-trained DVAEs (E-step) and updating the mixture model parameters $\boldsymbol{\varphi}$ (M-step). The complete proposed speech enhancement method is summarized in Algorithm 1.

## V. Implementation and Experiments

### A. Datasets

*a) Clean speech dataset:* We used the Wall Street Journal dataset (WSJ0) [46], which is a dataset of clean speech recorded from Wall Street Journal news at a sampling rate of 16kHz. We choose the speaker-independent medium vocabulary (5k words) subset. More precisely, the *si_tr_s* subset (around 25 hours), the *si_dt_05* subset (around 2 hours) and the *si_et_05* subset (around 1.5 hours) were used for DVAEs training, validation and testing, respectively.

*b) Noisy speech dataset:* For evaluating speech enhancement methods, we need a dataset of noisy speech signals along with the ground-truth clean speech, in order to compute objective performance measures. As will be discussed later, we will also compare the proposed unsupervised method with a supervised baseline. The latter requires a labeled dataset of pairs of noisy and clean speech signals not only for testing, but also for training and validation. We therefore created the

QUT-WSJ0 dataset, by mixing the clean speech signals of the WSJ0 dataset with noise signals from the QUT-NOISE dataset [47], including the four noise types {"café", "home", "street", "car"}. We combined the development/enrollment/verification splits of the QUT-NOISE dataset, as detailed in [47], with the above-mentioned training/validation/test splits of the WSJ0 dataset. The amount of data in the resulting QUT-WSJ0 dataset is the same as in the test subset of WSJ0 (around 1.5 hours). The dataset was created as follows. For each utterance of the WSJ0 dataset, we randomly selected a noise file from the QUT-NOISE dataset, in which we randomly selected a segment of noise signal with the same length as the speech utterance, we randomly choose a signal-to-noise ratio (SNR) among $\{-5, 0, 5\}$ dB, and we finally mixed the speech and noise signals at the chosen SNR, using the ITU-R BS.1770-4 protocol [48] to measure the power of the speech and noise signals.[9]

### B. Data preprocessing

The DVAEs were trained with power spectrograms of the speech signals from the WSJ0 dataset, with the following preprocessing. We first removed the silences at the beginning and ending of the files, using a detection threshold of $-30$ dB. The waveform was then normalized by its maximum absolute value. The STFT was computed on the normalized waveform with a 64-ms sine window (1024 samples) and a 25%-overlap (256 samples hop length), resulting in a sequence of 513-dimensional discrete Fourier coefficients (for positive frequencies). We set $T = 50$, meaning that speech utterances of 0.8 s were extracted to train the DVAE models. In summary, each training data sequence is a $50 \times 513$ STFT power spectrogram. This data preprocessing results in a set of $N_{\text{tr}} = 46,578$ training sequences (about 10.3 hours of speech signal) and $N_{\text{val}} = 7,775$ validation sequences (about 1.7 hour). For evaluation, we used the STFT spectrogram of each complete test sequence (still with silence clipping and normalization), which can be of variable length, most often larger than 2 s.

### C. Evaluation metrics

We used three metrics to evaluate the quality of the reconstructed speech signals, either in the analysis-resynthesis or the speech enhancement experiments: The scale-invariant signal-to-distortion ratio (SI-SDR) in dB [49], the perceptual evaluation of speech quality (PESQ) score (in $[-0.5, 4.5]$) [50], and the extended short-time objective intelligibility (ESTOI) score (in $[0, 1]$) [51]. For all measures, a higher value indicates a better result.

### D. Model configurations

Here we present the implementation of the three example DVAEs that we selected to illustrate the proposed DVAE-based speech enhancement algorithm, namely DKF, RVAE and SRNN. As indicated in [30], we can have various implementations for each DVAE model, thus we only present the

---

[9]Note that an SNR computed with this protocol is 2.5 dB lower (in average) than with a simple sum of the squared signal coefficients.
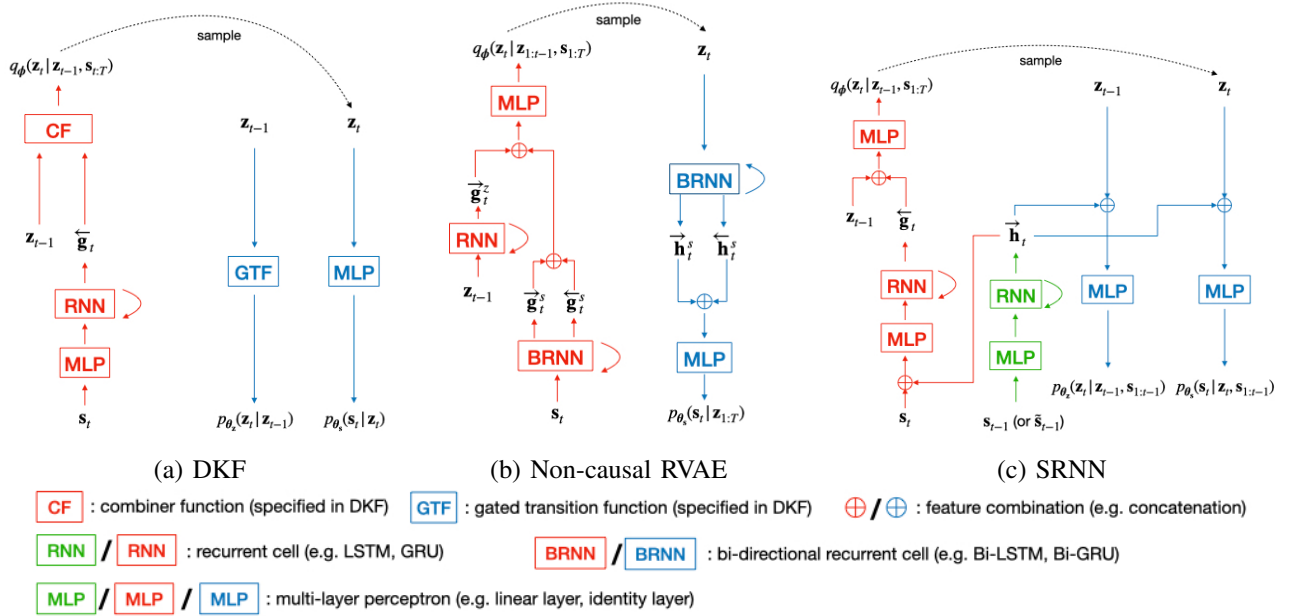
Fig. 3. Model implementation for (a) Deep Kalman Filter (DKF) [23], [24], (b) Recurrent Variational Auto-Encoder [28], and (c) Stochastic Recurrent Neural Network (SRNN) [26]. Each model consists of an inference (encoder) network (in red) and a generation (decoder) network (in blue). SRNN has a shared recurrent module between encoder and decoder (in green), which accumulates the information from past observation vectors.

model configurations that showed the best performance in our experiments (for the latent space dimension selected below).

*1) Dimension of the latent space:* In the present experiments we set $L = 16$. We recall that the data dimension is $F = 513$. We also recall that $\mathbf{z}_t$ is a real-valued vector that is modeled by a Gaussian distribution, so the DNNs modeling $\mathbf{z}_t$ have to output two $L$-dimensional vectors (the mean and variance vectors) for both inference and generation,[10] whereas $\mathbf{s}_t$ is a complex-valued vector modeled by a circular complex Gaussian distribution, which only leaves one $F$-dimensional variance vector to be provided by the decoder DNN. To guarantee the positivity of output variance, we used log-parameterization (the output is the log-variance in $\mathbb{R}$, which is then converted to variance by taking the exponential). Please, note that the last layer to predict mean and log-variance parameters is always a linear layer, where its output dimension corresponds to the one of $\mathbf{z}_t$ (16) or of $\mathbf{s}_t$ (513). We omit this in the following description for simplicity.

*2) DKF:* Fig. 3(a) shows the architecture of DKF in our implementation. We keep the specific *Combiner Function* and *Gated Transition Function* described in [24] for the layers providing the parameters of the inference and generative models of $\mathbf{z}_t$, respectively. Please refer to this paper for the implementation details. For the inference model, we also use a backward long short-term memory (LSTM) [52] layer with an internal state of dimension 128 in order to accumulate the information from $\mathbf{s}_{t:T}$ in (13). Before being fed into the recurrent layer, each vector $\mathbf{s}_t$ passes through an MLP with one hidden layer of dimension 256 and a $\tanh$ activation. The parameters of the generative model of $\mathbf{s}_t$ (i.e. the variance coefficients) are provided by an MLP with 4 hidden layers,

of dimension 32, 64, 128 and 256, with a $\tanh$ activation function.

*3) RVAE:* We implemented the non-causal version of RVAE, and for simplicity, we simply refer to this model as RVAE in all the following. Fig. 3(b) represents its implementation. The inference model includes a bi-directional LSTM (BLSTM) layer with an internal state of dimension 128, which takes as input the complete sequence $\mathbf{s}_{1:T}$. It also includes an RNN which processes the sampled sequence of the past latent vectors $\mathbf{z}_{1:t-1}$. The output of these two layers are then concatenated and mapped into the parameters of the inference model over $\mathbf{z}_t$ by an MLP. The generative part of the model includes a BLSTM layer with an internal state of dimension 128, which takes the sampled $\mathbf{z}_{1:T}$ as input. The output of this bi-directional LSTM layer is finally mapped into the parameters of the generative model over $\mathbf{s}_t$ by the output MLP (in the present RVAE implementation, this MLP is reduced to a single linear layer because this led to the best performance).

*4) SRNN:* SRNN is quite different from the two previous models. As shown in Fig. 3(c), the inference and generative models share an internal RNN state vector $\overrightarrow{\boldsymbol{h}}_t$ (module in green) which is encoding the information from the past observed vectors $\mathbf{s}_{1:t-1}$. This common module is composed of an MLP with one layer of dimension 256 followed by a forward LSTM. The dimension of $\overrightarrow{\boldsymbol{h}}_t$ is 128. For inference, the concatenation of $\overrightarrow{\boldsymbol{h}}_t$ and $\mathbf{s}_t$ is fed into a one-layer MLP of dimension 256 followed by a backward LSTM which provides the vector $\overleftarrow{\mathbf{g}}_t$ (of dimension 128). This vector is then concatenated with the sample of $\mathbf{z}_{t-1}$ and fed into an MLP with two hidden layers of dimension 64 and 32 using a $\tanh$ activation function. For the generative part, we concatenate the shared state $\overrightarrow{\boldsymbol{h}}_t$ along with samples of the latent vector at the previous or current time frame. To generate $\mathbf{z}_t$, we use an MLP

---

[10]except when $\mathbf{z}_t$ is assumed i.i.d. with a standard Gaussian distribution and no DNN is used for its generation.

TABLE I
AVERAGE SPEECH ANALYSIS-RESYNTHESIS RESULTS.

| Models | SI-SDR (dB) | PESQ | ESTOI |
|---|---|---|---|
| VAE | 5.3 | 2.97 | 0.83 |
| DKF | 9.3 | 3.53 | 0.91 |
| RVAE (non-causal) | 8.9 | 3.58 | 0.91 |
| SRNN-TF-True | **11.0** | **3.68** | **0.93** |
| SRNN-TF-Pred | -1.0 | 1.93 | 0.64 |
| SRNN-SS | 7.8 | 3.37 | 0.88 |

TABLE II
SPEECH ENHANCEMENT RESULTS (MEDIAN AND CONFIDENCE INTERVAL).

| Models | SI-SDR (dB) | PESQ | ESTOI |
|---|---|---|---|
| Noisy mixture | $-2.6 \pm 0.5$ | $1.81 \pm 0.03$ | $0.49 \pm 0.01$ |
| VAE | $4.4 \pm 0.4$ | $1.93 \pm 0.05$ | $0.53 \pm 0.01$ |
| DKF | $6.2 \pm 1.0$ | $2.21 \pm 0.05$ | $0.63 \pm 0.01$ |
| RVAE (non-causal) | $\mathbf{6.7 \pm 1.0}$ | $\mathbf{2.38 \pm 0.04}$ | $\mathbf{0.67 \pm 0.01}$ |
| SRNN-TF-Pred | $-2.1 \pm 0.3$ | $1.12 \pm 0.04$ | $0.40 \pm 0.01$ |
| SRNN-SS | $6.6 \pm 0.9$ | $2.24 \pm 0.05$ | $0.64 \pm 0.01$ |
| UMX Original | $4.8 \pm 0.5$ | $2.12 \pm 0.04$ | $0.63 \pm 0.01$ |
| UMX Retrain | $6.2 \pm 0.4$ | $2.20 \pm 0.04$ | $0.65 \pm 0.01$ |

with two hidden layers of dimension 64 and 32, whereas to generate $\mathbf{s}_t$ we use an MLP with one hidden layer of dimension 128. These MLPs also use the `tanh` activation function.

### E. SRNN with scheduled sampling

In conventional training of SRNN (and of all autoregressive DVAE models) we use the ground-truth past clean speech vectors $\mathbf{s}_{1:t-1}$ to generate the current one $\mathbf{s}_t$, a strategy sometimes referred to as "teacher forcing" in the literature [53]. We have seen in Section IV-B2 that this is not possible in the enhancement algorithm and thus we have replaced $\mathbf{s}_{1:t-1}$ with its proxy $\tilde{\mathbf{s}}_{1:t-1}$ (recursively computed from the decoder output). To avoid a mismatch between training and speech enhancement conditions, we also trained SRNN using $\tilde{\mathbf{s}}_{1:t-1}$ (instead of $\mathbf{s}_{1:t-1}$) to generate $\mathbf{s}_t$. It is difficult to directly train such a model, so we adopted a "scheduled sampling" approach [54]. We first trained SRNN until convergence using the ground-truth clean speech signal. Then we fine-tuned the model by randomly replacing $\mathbf{s}_{1:t-1}$ with $\tilde{\mathbf{s}}_{1:t-1}$ at the input of the encoder-decoder shared module, when estimating $\mathbf{s}_t$. We replaced 20% of the $\mathbf{s}_{1:t-1}$ vectors for the first 50 epochs, and increased to 40% for the next 50 epochs, and so on, until we completely replaced the ground-truth clean speech signal with generated speech signals. Then we fine-tuned with totally generated speech signals for another 300 epochs. Overall, we fine-tuned the model for 500 epochs. These two training strategies, teacher forcing and scheduled sampling, lead to two different versions of SRNN, referred to as SRNN-TF and SRNN-SS respectively. Moreover, when SRNN is trained with teacher forcing, we test it under two settings: using either the clean speech vectors $\mathbf{s}_{1:t-1}$ or their prediction $\tilde{\mathbf{s}}_{1:t-1}$. We refer to them as SRNN-TF-True and SRNN-TF-Pred respectively. Very importantly, we cannot evaluate SRNN-TF-True for speech enhancement, since SRNN-TF-True requires access to the ground-truth past clean speech vectors. As we will see in Sections V-H and V-I, our experiments demonstrate the interest of training SRNN with scheduled sampling.

### F. Optimizer

For the training of the three models, we use the Adam optimizer [55] with a learning rate of $2e-3$. We trained the models with a batch size of 128 during 300 epochs and kept the one with lowest validation loss. For the fine-tuning to SRNN-SS, we change the learning rate to $1e-3$.

### G. Speech enhancement parameters

For all the methods, the rank of the NMF in the noise model (20) is set to $K = 8$. $\mathbf{W}_b$ and $\mathbf{H}_b$ are randomly initialized

from a uniform distribution in $[0, 1]$ (with a fixed seed to ensure fair comparisons), and $\mathbf{g}$ is initialized with an all-ones vector. For computing (41), we fix the number of samples to $R = 1$. The VEM algorithm is run for 500 iterations. We used Adam [55] for fine-tuning the encoder during the E-step. After pilot experiments conducted on the validation set with different values of the learning rate, this latter was set to $5e-3$ for RVAE and $1e-3$ for DKF and SRNN.

### H. Results of speech analysis-resynthesis

Before we examine the speech enhancement performance, we first rapidly compare the speech modeling capacities of the three selected DVAE models (and the original VAE) in a speech analysis-resynthesis experiment. The overall pipeline is shown in Fig. 1. For the VAE model, we use the architecture proposed in [28]. The mean results averaged over the clean version of the test dataset are presented in Table I (for SI-SDR scores, the noise is the modeling noise, i.e. the difference between original and modeled/reconstructed signal).

We can see from Table I that all DVAE models (except SRNN-TF-Pred) perform largely better than the original VAE, showing the benefits of introducing dynamics into VAE-based speech modeling. We can also see that SRNN-TF-True performs best for speech analysis-resynthesis. This is not surprising since SRNN-TF-True uses the ground-truth past sequence $\mathbf{s}_{1:t-1}$, which is a very strong information, to predict $\mathbf{s}_t$. When at test time $\mathbf{s}_{1:t-1}$ is replaced by the corresponding prediction $\tilde{\mathbf{s}}_{1:t-1}$, i.e. SRNN-TF-Pred, the performance drops drastically, far below VAE. This justifies the interest of the scheduled sampling training strategy for SRNN. Indeed, SRNN-SS performs much better than SRNN-TF-Pred, with scores that are slightly below the ones of RVAE and DKF. Overall, among the models that are applicable to the speech enhancement problem, DKF obtains the best results in terms of SI-SDR score, followed by RVAE, and then SRNN-SS. In terms of PESQ and ESTOI scores, DKF and RVAE are very close. Even if this is not shown in Table I, we can rapidly mention that the non-causal version of RVAE always performs better than the causal version.

### I. Results of speech enhancement

We report in Table II the speech enhancement results on the QUT-WSJ0 test set, averaged over the three tested SNRs ($-5$, 0, 5 dB) and over the 4 noise types, for the different
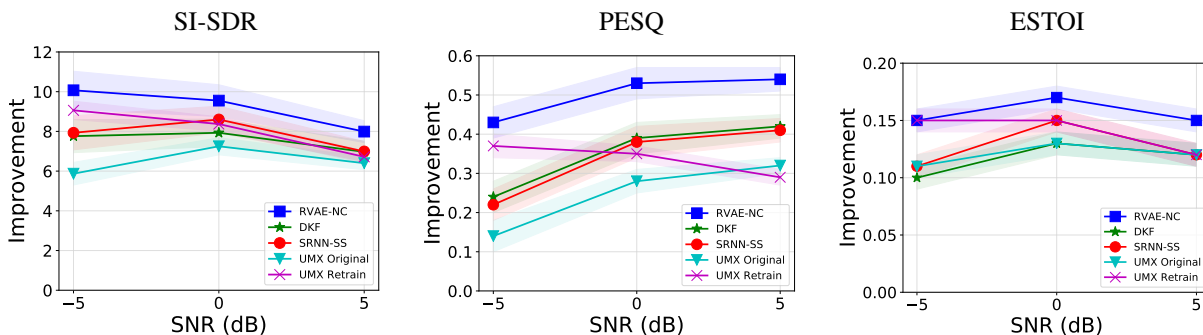
Fig. 4. Performance comparison of different DVAE architecture for speech enhancement. The improvement on SI-SDR (left), PESQ (middle) and ESTOI (right) of the estimated speech w.r.t. the noisy input signal is reported, as a function of the input SNR.
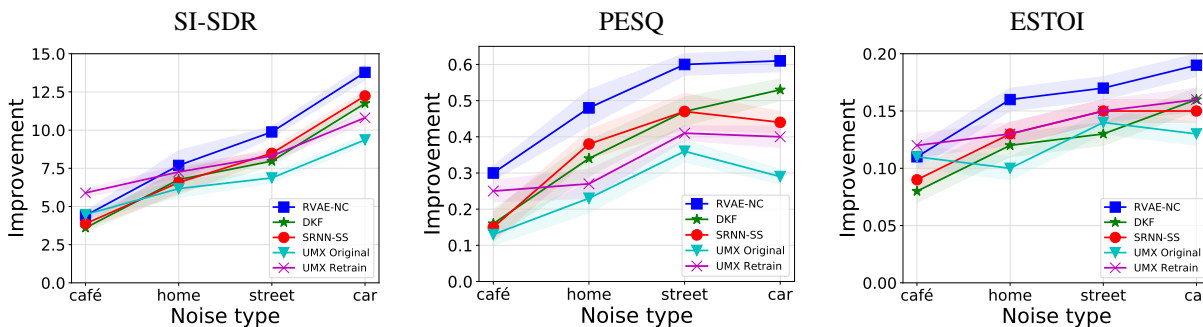


Fig. 5. Performance comparison of different DVAE architecture for speech enhancement. The improvement on SI-SDR (left), PESQ (middle) and ESTOI (right) of the estimated speech w.r.t. the noisy input signal is reported, per noise type.

DVAE models. The results for the original VAE model are the ones reported in [28], which used the same datasets. The confidence interval for the median is defined as 1.57 times the interquartile range divided by the square root of the sample size [56].

We first observe that all DVAE models (except SRNN-TF-Pred) largely outperform the original VAE model, which demonstrates the benefits of introducing dynamics into VAE-based speech modeling for speech enhancement. As could be expected, SRNN-TF-Pred is inefficient for speech enhancement, because of the mismatch between training and speech enhancement conditions. Among the other different tested DVAE models, the (non-causal) RVAE performs best for all evaluation metrics. For example, we have SI-SDR = $6.7 \pm 1.0$ dB for RVAE vs $4.4 \pm 0.4$ dB for VAE, hence an improvement of about $2.3$ dB. SRNN-SS is here just $0.1$ dB below non-causal RVAE, and performs better than DKF ($6.2 \pm 1.0$ dB SI-SDR). PESQ and ESTOI scores follow the same trend. This confirms the relevance and efficiency of the scheduled sampling training strategy, in this speech enhancement context. In terms of PESQ score, the results are quite consistent between analysis-resynthesis and speech enhancement (RVAE has the best PESQ score in both task). We recall that the setting SRNN-TF-True cannot be used for speech enhancement, since the ground-truth clean speech vectors are not available.

We also compare with Open-Unmix (UMX), an open-source supervised method based on a bi-directional LSTM network, which was state-of-the-art for music source separation [57], and was later adapted for speech enhancement [58]. We chose this baseline method because it relies on the same type of recurrent neural networks as the DVAE models considered in this work. "UMX Original" in Table II indicates that we use the model of [58] made available by the authors and trained on the VoiceBank-DEMAND dataset proposed in [59] (and of course we test it on the QUT-WSJO test set). This pre-trained supervised model performs notably worse than the DVAE-based unsupervised methods, e.g. $4.8 \pm 0.5$ dB SI-SDR vs $6.7 \pm 1.0$ dB SI-SDR for RVAE. This is likely to result from a limited generalization capacity, so we also retrained the UMX model on the QUT-WSJ0 training set, so as to have a better match between the training and testing conditions, and a fair comparison with the proposed DVAE-based unsupervised methods. The results are shown in Table II under the label "UMX Retrain". We observe that this retrained model performs better than the original pre-trained model, e.g. $6.2 \pm 0.4$ dB SI-SDR, but the performance is still globally lower than the unsupervised proposed methods (or on par in the case of DKF). Note that the QUT-WSJ0 dataset is particularly challenging for supervised speech enhancement methods, as for instance the "home" category includes "kitchen" noise for training, and "living room" noise for validation and testing, which are quite different sub-categories. The proposed unsupervised speech enhancement algorithms offer a flexibility that is useful when supervised methods have difficulties to generalize to unseen acoustic scenes.

Fig. 4 presents the speech enhancement performance in terms of improvement (or gain) of SI-SDR between input and output, averaged over the 4 noise types but detailed as a function of the input SNR. In all the following, we

do not consider SRNN-TF-Pred anymore, since this model was shown to be inefficient for speech enhancement. The circumflex accent shape of the curves is often encountered when reporting improvement over the noisy input [28], [60]–[62]. Indeed, when the input is very noisy, it is difficult to improve the quality. And when the audio is weakly corrupted, it is also difficult to obtain a significant improvement. This figure shows that RVAE performs best for all SNRs and all metrics. The behavior of the UMX Retrain baseline is quite different from the DVAE-based methods: It performs better than DKF and SRNN-SS (but not better than RVAE) at SNR = −5 dB, but then it decreases below DKF and SRNN-SS as SNR increases. This is particularly visible for the PESQ score.

Finally, we also report the improvement detailed per noise type (but averaged over the 3 SNRs) in Fig. 5. We observe that the algorithms perform quite differently depending on the noise type. As expected, the "café" noise is the most challenging one, whereas the "car" noise, being more stationary, is the easiest to process. Again, RVAE appears to exhibit the best overall performance. UMX seems appropriate for the "café" noise, whereas it is outperformed by DKF and SRNN-SS for the car noise for example (regarding SI-SDR and PESQ scores).

## VI. CONCLUSION

In this paper, we have proposed a general framework for unsupervised speech enhancement based on Dynamical Variational AutoEncoders, or DVAEs. In our framework, DVAEs are used to model the clean speech signal, while the noise is modeled via non-negative matrix factorisation. While DVAEs are pre-trained with clean speech, the noise parameters are estimated at test time, together with the clean speech, from the noisy speech observations. To achieve that, we have derived a VEM algorithm for the most general formulation of a DVAE model, which can then be easily adapted to particular instances of DVAEs. We have illustrated this principle with DKF, RVAE and SRNN, and this can be extended to other DVAE models, e.g. STORN [22] or VRNN [25].

We have compared the performance obtained with those three example DVAEs, both for speech analysis-resynthesis and speech enhancement. For speech enhancement, the proposed approach was shown to outperform a supervised baseline method using the same kind of recurrent networks. The (non-causal) RVAE provided the best performance among the tested DVAEs. The SRNN architecture also holds great potential, provided that it is trained with scheduled sampling in order to reduce the gap between the training and speech enhancement conditions. If this gap could be further decreased, we believe that it could have even better performance than the non-causal RVAE model. This aspect should be further investigated, possibly including other autoregressive models in the DVAE family (e.g. VRNN [25]).

Generally speaking, the inference with DVAEs is non-causal, meaning that past and future noisy speech observations are required to enhance a given speech frame. Causal DVAE-based speech enhancement can also be investigated, but is out of the scope of this paper. Other future works include

introducing other powerful encoder-decoder networks, e.g. the temporal convolutional networks (TCN) [39] and the Transformer [63], in the present unsupervised speech enhancement framework. The DVAE models may be further boosted with more expressive latent variables, e.g. introducing hierarchical multi-scale structure and normalizing flows [64]. We also plan to extend the proposed DVAE-based speech enhancement framework to a multi-modal framework, using the speaker's lips motion and visual appearance, in the continuation of VAE-based audio-visual speech enhancement [62].

## REFERENCES

[1] J. Benesty, S. Makino, and J. Chen, *Speech enhancement*. Springer Science & Business Media, 2006.

[2] P. C. Loizou, *Speech enhancement: Theory and practice*. CRC press, 2013.

[3] S. Boll, "Suppression of acoustic noise in speech using spectral subtraction," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 27, no. 2, pp. 113–120, 1979.

[4] J. S. Lim and A. V. Oppenheim, "Enhancement and bandwidth compression of noisy speech," *Proceedings of the IEEE*, vol. 67, no. 12, pp. 1586–1604, 1979.

[5] Y. Ephraim and D. Malah, "Speech enhancement using a minimum-mean square error short-time spectral amplitude estimator," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 32, no. 6, pp. 1109–1121, 1984.

[6] D. Wang and J. Chen, "Supervised speech separation based on deep learning: An overview," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 26, no. 10, pp. 1702–1726, 2018.

[7] P. Smaragdis, B. Raj, and M. Shashanka, "Supervised and semi-supervised separation of sounds from single-channel mixtures," in *International Conference on Independent Component Analysis and Signal Separation*, Charleston, USA, 2007.

[8] G. J. Mysore and P. Smaragdis, "A non-negative approach to semi-supervised separation of speech from noise with the use of temporal dynamics," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Prague, Czech Republic, 2011.

[9] N. Mohammadiha, P. Smaragdis, and A. Leijon, "Supervised and unsupervised speech enhancement using nonnegative matrix factorization," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 21, no. 10, pp. 2140–2151, 2013.

[10] D. P. Kingma, S. Mohamed, D. Jimenez Rezende, and M. Welling, "Semi-supervised learning with deep generative models," in *Advances in Neural Information Processing Systems (NeurIPS)*, Montreal, Canada, 2014.

[11] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," in *International Conference on Learning Representations (ICLR)*, Banff, Canada, 2014.

[12] D. J. Rezende, S. Mohamed, and D. Wierstra, "Stochastic backpropagation and approximate inference in deep generative models," in *International Conference on Machine Learning (ICML)*, Beijing, China, 2014.

[13] 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 *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Calgary, Canada, 2018.

[14] S. Leglaive, L. Girin, and R. Horaud, "A variance modeling framework based on variational autoencoders for speech enhancement," in *IEEE International Workshop on Machine Learning for Signal Processing (MLSP)*, Aalborg, Denmark, 2018.

[15] S. Leglaive, U. Şimşekli, A. Liutkus, L. Girin, and R. Horaud, "Speech enhancement with variational autoencoders and alpha-stable distributions," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Brighton, UK, 2019.

[16] M. Pariente, A. Deleforge, and E. Vincent, "A statistically principled and computationally efficient approach to speech enhancement using variational autoencoders," *arXiv preprint arXiv:1905.01209*, 2019.

[17] K. Sekiguchi, Y. Bando, K. Yoshii, and T. Kawahara, "Bayesian multi-channel speech enhancement with a deep speech prior," in *Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA)*, Honolulu, USA, 2018.

[18] S. Leglaive, L. Girin, and R. Horaud, "Semi-supervised multichannel speech enhancement with variational autoencoders and non-negative matrix factorization," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Brighton, UK, 2019.

[19] M. Fontaine, A. A. Nugraha, R. Badeau, K. Yoshii, and A. Liutkus, "Cauchy multichannel speech enhancement with a deep speech prior," in *European Signal Processing Conference (EUSIPCO)*, A Coruna, Spain, 2019.

[20] C. Févotte, N. Bertin, and J.-L. Durrieu, "Nonnegative matrix factorization with the Itakura-Saito divergence: With application to music analysis," *Neural Computation*, vol. 21, no. 3, pp. 793–830, 2009.

[21] O. Fabius and J. R. van Amersfoort, "Variational recurrent autoencoders," *arXiv preprint arXiv:1412.6581*, 2014.

[22] J. Bayer and C. Osendorfer, "Learning stochastic recurrent networks," *arXiv preprint arXiv:1411.7610*, 2014.

[23] R. G. Krishnan, U. Shalit, and D. Sontag, "Deep Kalman filters," *arXiv preprint arXiv:1511.05121*, 2015.

[24] R. Krishnan, U. Shalit, and D. Sontag, "Structured inference networks for nonlinear state space models," in *AAAI Conference on Artificial Intelligence (AAAI)*, San Francisco, USA, 2017.

[25] J. Chung, K. Kastner, L. Dinh, K. Goel, A. Courville, and Y. Bengio, "A recurrent latent variable model for sequential data," in *Advances in Neural Information Processing Systems (NeurIPS)*, Montreal, Canada, 2015.

[26] M. Fraccaro, S. K. Sønderby, U. Paquet, and O. Winther, "Sequential neural models with stochastic layers," in *Advances in Neural Information Processing Systems (NeurIPS)*, Barcelona, Spain, 2016.

[27] Y. Li and S. Mandt, "Disentangled sequential autoencoder," in *International Conference on Machine Learning (ICML)*, Stockholm, Sweden, 2018.

[28] S. Leglaive, X. Alameda-Pineda, L. Girin, and R. Horaud, "A recurrent variational autoencoder for speech enhancement," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Barcelona, Spain, 2020.

[29] M. Fraccaro, S. Kamronn, U. Paquet, and O. Winther, "A disentangled recognition and nonlinear dynamics model for unsupervised learning," in *Advances in Neural Information Processing Systems (NeurIPS)*, Long Beach, USA, 2017.

[30] L. Girin, S. Leglaive, X. Bie, J. Diard, T. Hueber, and X. Alameda-Pineda, "Dynamical variational autoencoders: A comprehensive review," *arXiv preprint arXiv:2008.12595*, 2020.

[31] R. M. Neal and G. E. Hinton, "A view of the EM algorithm that justifies incremental, sparse, and other variants," in *Learning in graphical models*. Springer, 1998, pp. 355–368.

[32] Y. Wang and D. Wang, "Boosting classification based speech separation using temporal dynamics," in *Conference of the International Speech Communication Association (INTERSPEECH)*, Portland, USA, 2012.

[33] ——, "Cocktail party processing via structured prediction," *Advances in Neural Information Processing Systems (NeurIPS)*, 2012.

[34] X. Lu, Y. Tsao, S. Matsuda, and C. Hori, "Speech enhancement based on deep denoising autoencoder." in *Conference of the International Speech Communication Association (INTERSPEECH)*, Lyon, France, 2013.

[35] S. Pascual, A. Bonafonte, and J. Serra, "Segan: Speech enhancement generative adversarial network," *arXiv preprint arXiv:1703.09452*, 2017.

[36] G. Carbajal, J. Richter, and T. Gerkmann, "Guided variational autoencoder for speech enhancement with a supervised classifier," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Online Event / Toronto, Canada, 2021.

[37] J. Richter, G. Carbajal, and T. Gerkmann, "Speech enhancement with stochastic temporal convolutional networks," *Conference of the International Speech Communication Association (INTERSPEECH)*, 2020.

[38] E. Aksan and O. Hilliges, "STCN: Stochastic temporal convolutional networks," in *International Conference on Learning Representations (ICLR)*, New Orleans, USA, 2018.

[39] C. Lea, R. Vidal, A. Reiter, and G. D. Hager, "Temporal convolutional networks: A unified approach to action segmentation," in *European Conference on Computer Vision (ECCV)*, Amsterdam, The Netherlands, 2016.

[40] F. D. Neeser and J. L. Massey, "Proper complex random processes with applications to information theory," *IEEE Transactions on Information Theory*, vol. 39, no. 4, pp. 1293–1302, 1993.

[41] E. Vincent, M. G. Jafari, S. A. Abdallah, M. D. Plumbley, and M. E. Davies, "Probabilistic modeling paradigms for audio source separation," in *Machine Audition: Principles, Algorithms and Systems*. IGI global, 2011, pp. 162–185.

[42] A. Liutkus, R. Badeau, and G. Richard, "Gaussian processes for under-determined source separation," *IEEE Transactions on Signal Processing*, vol. 59, no. 7, pp. 3155–3167, 2011.

[43] P.-A. Mattei and J. Frellsen, "Refit your encoder when new data comes by," in *NeurIPS Workshop on Bayesian Deep Learning*, Montreal, Canada, 2018.

[44] D. R. Hunter and K. Lange, "A tutorial on MM algorithms," *The American Statistician*, vol. 58, no. 1, pp. 30–37, 2004.

[45] C. Févotte and J. Idier, "Algorithms for nonnegative matrix factorization with the $\beta$-divergence," *Neural Computation*, vol. 23, no. 9, pp. 2421–2456, 2011.

[46] J. Garofolo, D. Graff, D. Paul, and D. Pallett, "CSR-I (WSJ0) Sennheiser LDC93S6B. https://catalog.ldc.upenn.edu/ldc93s6b," *Philadelphia: Linguistic Data Consortium*, 1993.

[47] D. Dean, A. Kanagasundaram, H. Ghaemmaghami, M. H. Rahman, and S. Sridharan, "The QUT-NOISE-SRE protocol for the evaluation of noisy speaker recognition," in *Conference of the International Speech Communication Association (INTERSPEECH)*, Dresden, Germany, 2015.

[48] ITU-R, "Recommendation BS.1770-4: Algorithms to measure audio programme loudness and true-peak audio level," *BS Series*, 2011.

[49] J. Le Roux, S. Wisdom, H. Erdogan, and J. R. Hershey, "SDR: Half-baked or well done?" in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Brighton, UK, 2019.

[50] A. Rix, J. Beerends, M. Hollier, and A. Hekstra, "Perceptual evaluation of speech quality (PESQ): A new method for speech quality assessment of telephone networks and codecs," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Salt Lake City, USA, 2001.

[51] C. H. Taal, R. C. Hendriks, R. Heusdens, and J. Jensen, "An algorithm for intelligibility prediction of time–frequency weighted noisy speech," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 19, no. 7, pp. 2125–2136, 2011.

[52] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

[53] R. J. Williams and D. Zipser, "A learning algorithm for continually running fully recurrent neural networks," *Neural Computation*, vol. 1, no. 2, pp. 270–280, 1989.

[54] S. Bengio, O. Vinyals, N. Jaitly, and N. Shazeer, "Scheduled sampling for sequence prediction with recurrent neural networks," *arXiv preprint arXiv:1506.03099*, 2015.

[55] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *International Conference on Learning Representations (ICLR)*, San Diego, USA, 2015.

[56] R. McGill, J. W. Tukey, and W. A. Larsen, "Variations of box plots," *The American Statistician*, vol. 32, no. 1, pp. 12–16, 1978.

[57] F.-R. Stöter, S. Uhlich, A. Liutkus, and Y. Mitsufuji, "Open-Unmix: A reference implementation for music source separation," *Journal of Open Source Software*, vol. 4, no. 41, p. 1667, 2019.

[58] S. Uhlich and Y. Mitsufuji, "Open-Unmix for speech enhancement (UMX SE)," May 2020. [Online]. Available: https://doi.org/10.5281/zenodo.3786908

[59] C. Valentini-Botinhao, X. Wang, S. Takaki, and J. Yamagishi, "Speech enhancement for a noise-robust text-to-speech synthesis system using deep recurrent neural networks," in *Conference of the International Speech Communication Association (INTERSPEECH)*, San Francisco, USA, 2016.

[60] M. Sadeghi and X. Alameda-Pineda, "Switching variational autoencoders for noise-agnostic audio-visual speech enhancement," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Online Event / Toronto, Canada, 2021.

[61] ——, "Mixture of inference networks for VAE-based audio-visual speech enhancement," *IEEE Transactions on Signal Processing*, vol. 69, pp. 1899–1909, 2021.

[62] 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.

[63] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.

[64] A. Vahdat and J. Kautz, "NVAE: A deep hierarchical variational autoencoder," in *Advances in Neural Information Processing Systems (NeurIPS)*, Online Event / Vancouver, Canada, 2020.