

Unsupervised Probabilistic Learning with Latent Variables

Lecture 3: Dynamical variational autoencoders, application to speech enhancement.

MLSS Africa 2023, Cape Town

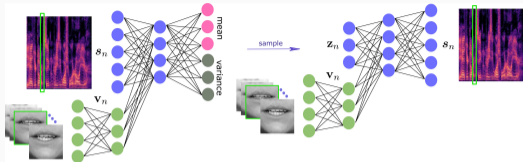
Slides: <https://xavirema.eu/MLSS2023/>

Xavier Alameda-Pineda

RobotLearn team, Inria at University Grenoble-Alpes,
Jean-Kuntzman CNRS Laboratory, Multidisciplinary Institute of Artificial Intelligence

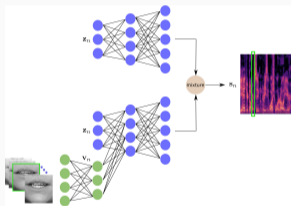


Summary of the three VAE-based models for AV-SE



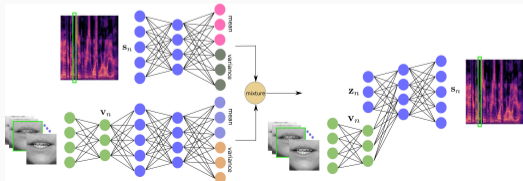
Conditional VAE:

- Training VAE via SGD.
- Systematic AV fusion.
- Learning noise parameters via MCEM.



VAE-MM:

- Training two VAEs via SGD.
- Mixing AV fusion - two speech models.
- Learning noise parameters via VEM.



MIN-VAE:

- Training all 3 networks via VEM+SGD.
- Mixing AV fusion - single speech model.
- Learning noise parameters via VEM.

Motivation for DVAE

All these models process speech data *independently* per each frame, that is regardless of the frame order.

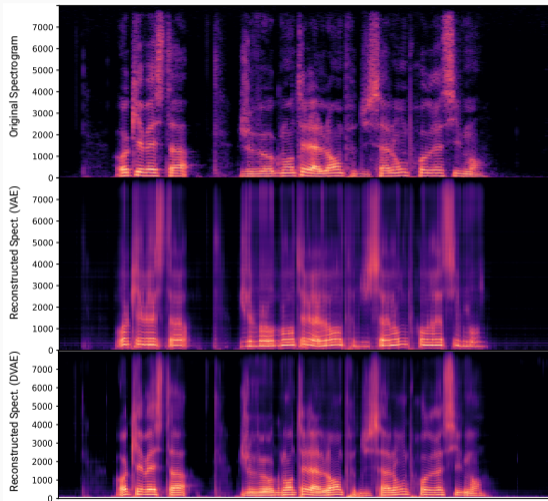
But speech is sequential, and order matters. What can we do?
(we will forget about visual data in this lecture)

Are there VAE-like models able to deal with sequential data?

Yes, we will call them **dynamical variational autoencoders (DVAE)**.

We'll also discuss their use in unsupervised speech enhancement?

Reconstruction?

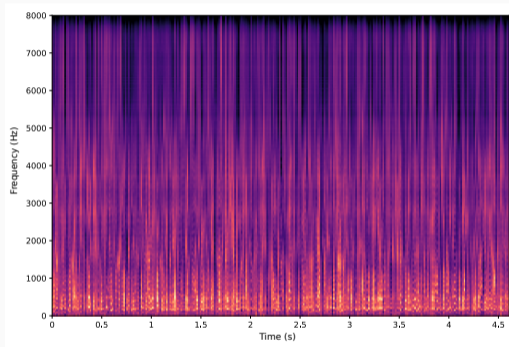


Original (top)

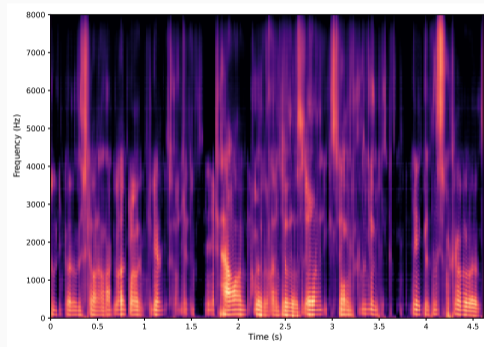
Reconstruction VAE (middle)

Reconstruction DVAE (bottom)

Generation?



Generation VAE (left)



Generation DVAE (right)

Lecture's Outline

- 1 Formalising Dynamical Variational Autoencoders
- 2 Inference network: conditional independence and Markov blankets
- 3 DVAE: general case, implementation and learning
- 4 DVAEs for unsupervised speech enhancement

Formalising Dynamical Variational Autoencoders

Probabilistic Sequential Modeling

We would like to model sequences of observations and latent variables:

Observed sequence: $\mathbf{x}_{1:T} = \{\mathbf{x}_t \in \mathbb{R}^F\}_{t=1}^T$

Latent sequence: $\mathbf{z}_{1:T} = \{\mathbf{z}_t \in \mathbb{R}^D\}_{t=1}^T$

Generative sequential modeling consists in defining the joint distribution with temporal dependencies, rather than the frame-wise joint distribution (as a vanilla VAE does):

$$p_{\theta}^{\text{DVAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) \neq p_{\theta}^{\text{VAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) = \prod_{t=1}^T p_{\theta}(\mathbf{x}_t, \mathbf{z}_t).$$

Chain rule

Using the **chain rule** we can write the joint distribution as a product of conditionals:

$$p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) = p(\mathbf{z}_1)p(\mathbf{x}_1|\mathbf{z}_1) \prod_{t=2}^T p(\mathbf{z}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t-1})p(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t})$$

Generative process:

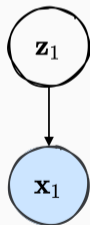


Chain rule

Using the **chain rule** we can write the joint distribution as a product of conditionals:

$$p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) = p(\mathbf{z}_1)p(\mathbf{x}_1|\mathbf{z}_1) \prod_{t=2}^T p(\mathbf{z}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t-1})p(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t})$$

Generative process:

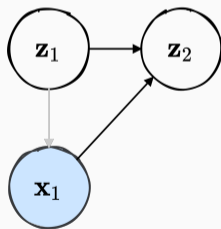


Chain rule

Using the **chain rule** we can write the joint distribution as a product of conditionals:

$$p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) = p(\mathbf{z}_1)p(\mathbf{x}_1|\mathbf{z}_1) \prod_{t=2}^T p(\mathbf{z}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t-1})p(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t})$$

Generative process:

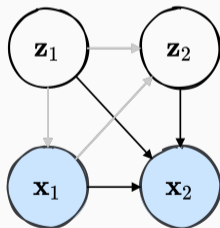


Chain rule

Using the **chain rule** we can write the joint distribution as a product of conditionals:

$$p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) = p(\mathbf{z}_1)p(\mathbf{x}_1|\mathbf{z}_1) \prod_{t=2}^T p(\mathbf{z}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t-1})p(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t})$$

Generative process:

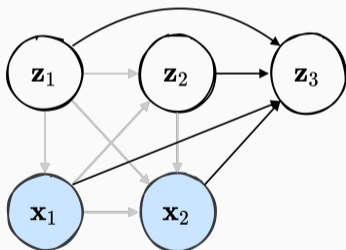


Chain rule

Using the **chain rule** we can write the joint distribution as a product of conditionals:

$$p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) = p(\mathbf{z}_1)p(\mathbf{x}_1|\mathbf{z}_1) \prod_{t=2}^T p(\mathbf{z}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t-1})p(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t})$$

Generative process:

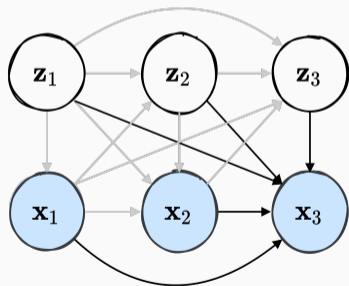


Chain rule

Using the **chain rule** we can write the joint distribution as a product of conditionals:

$$p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) = p(\mathbf{z}_1)p(\mathbf{x}_1|\mathbf{z}_1) \prod_{t=2}^T p(\mathbf{z}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t-1})p(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t})$$

Generative process:

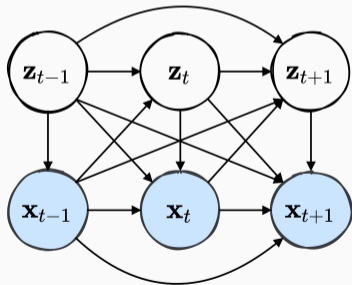


Chain rule

Using the **chain rule** we can write the joint distribution as a product of conditionals:

$$p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) = p(\mathbf{z}_1)p(\mathbf{x}_1|\mathbf{z}_1) \prod_{t=2}^T p(\mathbf{z}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t-1})p(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t})$$

Generative process:



The distribution $p(\mathbf{z}_{t+1}|\mathbf{z}_{1:t}, \mathbf{x}_{1:t})$ is not a *prior* distribution anymore, ...

...it is parametric and it might be auto-regressive (AR).

$p(\mathbf{x}_{t+1}|\mathbf{z}_{1:t+1}, \mathbf{x}_{1:t})$ might be AR as well.

Very well-known particular case: Kalman

A particular case of the above paradigm are linear dynamical systems (a.k.a. Kalman filter). Assuming the following conditional independence hypothesis:

$$p(\mathbf{z}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_{1:t-1}) = p(\mathbf{z}_t | \mathbf{z}_{t-1}) \quad p(\mathbf{x}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_{1:t}) = p(\mathbf{x}_t | \mathbf{z}_t),$$

Very well-known particular case: Kalman

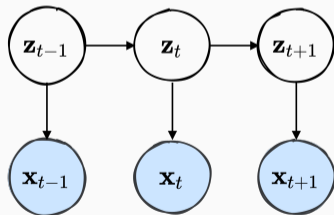
A particular case of the above paradigm are linear dynamical systems (a.k.a. Kalman filter). Assuming the following conditional independence hypothesis:

$$p(\mathbf{z}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_{1:t-1}) = p(\mathbf{z}_t | \mathbf{z}_{t-1}) \quad p(\mathbf{x}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_{1:t}) = p(\mathbf{x}_t | \mathbf{z}_t),$$

the joint distribution simplifies as:

$$p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) = p(\mathbf{z}_1) p(\mathbf{x}_1 | \mathbf{z}_1) \prod_{t=2}^T p(\mathbf{z}_t | \mathbf{z}_{t-1}) p(\mathbf{x}_t | \mathbf{z}_t).$$

This is the family of state-space models (SSMs) introduced by Kalman in 1960 (linear-Gaussian SSMs with continuous state and hidden Markov models with discrete state).



(Deep?) Kalman/Markov model

The standard Kalman (HMM) models, the dependency is linear (left), while in the deep counterparts it is non-linear (right)!

Kalman filter:

$$p_{\theta_z}(z_t|z_{t-1}) = \mathcal{N}(z_t; \mathbf{A}_t z_{t-1}, \mathbf{Q}_t)$$

$$p_{\theta_x}(x_t|z_t) = \mathcal{N}(x_t; \mathbf{B}_t z_t, \mathbf{R}_t),$$

Deep Kalman filter:

$$p_{\theta_z}(z_t|z_{t-1}) = \mathcal{N}(z_t; \mu_{\theta_z}(z_{t-1}), \Sigma_{\theta_z}(z_{t-1}))$$

$$p_{\theta_x}(x_t|z_t) = \mathcal{N}(x_t; \mu_{\theta_x}(z_t), \Sigma_{\theta_x}(z_t)),$$

The non-linearities can be implemented with feed-forward neural networks, as in VAEs.

Can we learn it?

Recall that, for learning VAE, we needed an auxiliary inference (encoder) network, that approximates the posterior:

This is our next objective!

Inference network: conditional independence and Markov blankets

Sequential Inference Network

As in the VAE, we need to approximate the posterior distribution:

$$p_{\theta}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T}) = q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})$$

Sequential Inference Network

As in the VAE, we need to approximate the posterior distribution:

$$p_{\theta}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T}) = q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})$$

We can always use the Bayes theorem to write:

$$q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T}) = \prod_{t=1}^T q_{\phi}(\mathbf{z}_t|\mathbf{z}_{1:t-1}, \mathbf{x}_{1:T})$$

Can we simplify each of the terms further? It depends on generative model.

Sequential Inference Network

As in the VAE, we need to approximate the posterior distribution:

$$p_{\theta}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T}) = q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})$$

We can always use the Bayes theorem to write:

$$q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T}) = \prod_{t=1}^T q_{\phi}(\mathbf{z}_t|\mathbf{z}_{1:t-1}, \mathbf{x}_{1:T})$$

Can we simplify each of the terms further? It depends on generative model.

For instance, for (deep) kalman filter we have the following result:

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

But how to obtain it?

Conditional Independence

You might very well know what “independence” means:

$$\mathbf{x} \perp\!\!\!\perp \mathbf{y} \Leftrightarrow p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y}) \Leftrightarrow p(\mathbf{x}|\mathbf{y}) = p(\mathbf{x}).$$

Conditional Independence

You might very well know what “independence” means:

$$\mathbf{x} \perp\!\!\!\perp \mathbf{y} \Leftrightarrow p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y}) \Leftrightarrow p(\mathbf{x}|\mathbf{y}) = p(\mathbf{x}).$$

Conditional independence is defined in a similar way, but conditioned to z :

$$\mathbf{x} \perp\!\!\!\perp \mathbf{y} \mid z \Leftrightarrow p(\mathbf{x}, \mathbf{y}|z) = p(\mathbf{x}|z)p(\mathbf{y}|z) \Leftrightarrow p(\mathbf{x}|\mathbf{y}, z) = p(\mathbf{x}|z).$$

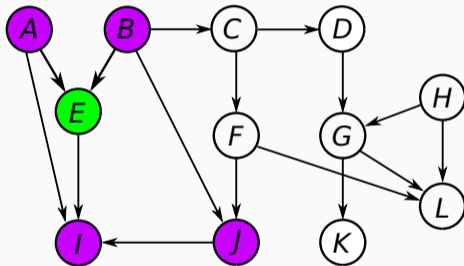
This means that z is *breaking* the link between x and y : if we know z , y does not bring any extra information to describe x .

In models with lots of variables, we cannot test all possible subsets of variables z , we need an *automatic* way to check that.

Markov Blankets

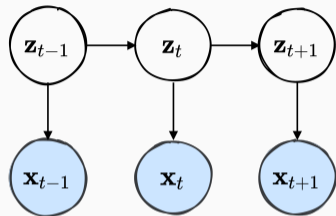
(Disclaimer: we are taking a shortcut. Other key concepts are path blocking and D-separation.)

The **Markov Blanket** of x , $\mathcal{B}(x)$ is the set of nodes isolating x from the rest of the graph.



In the case above: $\mathcal{B}(E) = \{A, B, I, J\}$, that is the set of **parents**, **children** and **co-parents** (other parents of its children).

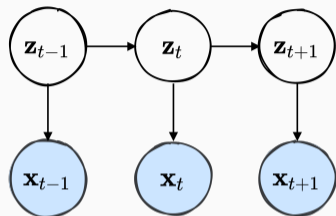
Let's take a look back at DKF



We realise that z_{t-1} is isolating $\{z_t, \mathbf{x}_{t:T}\}$ from the past:
(z_{t-1} is the Markov blanket of $\{z_t, \mathbf{x}_{t:T}\}$)

$$p_{\theta}(z_t | z_{1:t-1}, \mathbf{x}_{1:T}) = p_{\theta}(z_t | z_{t-1}, \mathbf{x}_{t:T})$$

Let's take a look back at DKF



We realise that z_{t-1} is isolating $\{z_t, \mathbf{x}_{t:T}\}$ from the past:
(z_{t-1} is the Markov blanket of $\{z_t, \mathbf{x}_{t:T}\}$)

$$p_{\theta}(z_t | z_{1:t-1}, \mathbf{x}_{1:T}) = p_{\theta}(z_t | z_{t-1}, \mathbf{x}_{t:T})$$

It seems natural to impose the same dependencies in the inference network:

$$p_{\theta}(z_{1:T} | \mathbf{x}_{1:T}) = \prod_{t=1}^T p_{\theta}(z_t | z_{t-1}, \mathbf{x}_{t:T}) \approx \prod_{t=1}^T q_{\phi}(z_t | z_{t-1}, \mathbf{x}_{t:T})$$

Note: It is OK to simplify even further, but you need to know what you are doing.

DVAE: general case, implementation and learning

The general (causal) case

General umbrella for all (causal) DVAEs:

$$p_{\theta}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) = \prod_{t=1}^T p_{\theta_z}(\mathbf{z}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_{1:t-1}) p_{\theta_x}(\mathbf{x}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_{1:t}),$$

where

$$p_{\theta_z}(\mathbf{z}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_{1:t-1}) = \mathcal{N}\left(\mathbf{z}_t; \boldsymbol{\mu}_{\theta_z}(\dots), \text{diag}\{\mathbf{v}_{\theta_z}(\dots)\}\right),$$

$$p_{\theta_x}(\mathbf{x}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_{1:t}) = \mathcal{N}\left(\mathbf{x}_t; \boldsymbol{\mu}_{\theta_x}(\dots), \text{diag}\{\mathbf{v}_{\theta_x}(\dots)\}\right),$$

and $\{\boldsymbol{\mu}_{\theta_z}, \mathbf{v}_{\theta_z}\}$, and $\{\boldsymbol{\mu}_{\theta_x}, \mathbf{v}_{\theta_x}\}$ are **non-linear functions** of the conditioning variables.

Which kind of networks can be use to construct these functions?

The general (causal) case

General umbrella for all (causal) DVAEs:

$$p_{\theta}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) = \prod_{t=1}^T p_{\theta_z}(\mathbf{z}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_{1:t-1}) p_{\theta_x}(\mathbf{x}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_{1:t}),$$

where

$$p_{\theta_z}(\mathbf{z}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_{1:t-1}) = \mathcal{N}\left(\mathbf{z}_t; \boldsymbol{\mu}_{\theta_z}(\dots), \text{diag}\{\mathbf{v}_{\theta_z}(\dots)\}\right),$$

$$p_{\theta_x}(\mathbf{x}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_{1:t}) = \mathcal{N}\left(\mathbf{x}_t; \boldsymbol{\mu}_{\theta_x}(\dots), \text{diag}\{\mathbf{v}_{\theta_x}(\dots)\}\right),$$

and $\{\boldsymbol{\mu}_{\theta_z}, \mathbf{v}_{\theta_z}\}$, and $\{\boldsymbol{\mu}_{\theta_x}, \mathbf{v}_{\theta_x}\}$ are **non-linear functions** of the conditioning variables.

Which kind of networks can be use to construct these functions? RNN, Transformers, ...

Example with vanilla RNN

Recall the DVAE generative model:

$$p_{\theta}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) = \prod_{t=1}^T p_{\theta_z}(\mathbf{z}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_{1:t-1}) p_{\theta_x}(\mathbf{x}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_{1:t}).$$

The conditional distributions are parametrized, for instance:

- $\mathbf{h}_t = \sigma(\mathbf{W}_{xh}\mathbf{x}_{t-1} + \mathbf{W}_{zh}\mathbf{z}_{t-1} + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h)$,
- $p_{\theta_z}(\mathbf{z}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_{1:t-1}) = \mathcal{N}(\mathbf{z}_t; \mu_{\theta_z}(\mathbf{h}_t), \text{diag}\{\mathbf{v}_{\theta_z}(\mathbf{h}_t)\})$,
- $p_{\theta_x}(\mathbf{x}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_{1:t}) = \mathcal{N}(\mathbf{x}_t; \mu_{\theta_x}(\mathbf{z}_t, \mathbf{h}_t), \text{diag}\{\mathbf{v}_{\theta_x}(\mathbf{z}_t, \mathbf{h}_t)\})$.

In this example, the same RNN is used for \mathbf{x} and for \mathbf{z} . This is an arbitrary choice.

Inference for the general case

In the general case, we cannot simplify the dependencies of the inference model:

$$q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T}) = \prod_{t=1}^T q_{\phi}(\mathbf{z}_t|\mathbf{z}_{1:t-1}, \mathbf{x}_{1:T}),$$

where typically we have ($\tilde{\boldsymbol{\mu}}_{\phi}$, $\tilde{\mathbf{v}}_{\phi}$ are non-linear functions of $\mathbf{z}_{1:t-1}$, $\mathbf{x}_{1:T}$):

$$q_{\phi}(\mathbf{z}_t|\mathbf{z}_{1:t-1}, \mathbf{x}_{1:T}) = \mathcal{N}\left(\mathbf{z}_t; \tilde{\boldsymbol{\mu}}_{\phi}(\dots), \text{diag}\{\tilde{\mathbf{v}}_{\phi}(\dots)\}\right).$$

Inference for the general case

In the general case, we cannot simplify the dependencies of the inference model:

$$q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T}) = \prod_{t=1}^T q_{\phi}(\mathbf{z}_t|\mathbf{z}_{1:t-1}, \mathbf{x}_{1:T}),$$

where typically we have $(\tilde{\boldsymbol{\mu}}_{\phi}, \tilde{\mathbf{v}}_{\phi})$ are non-linear functions of $\mathbf{z}_{1:t-1}, \mathbf{x}_{1:T}$:

$$q_{\phi}(\mathbf{z}_t|\mathbf{z}_{1:t-1}, \mathbf{x}_{1:T}) = \mathcal{N}\left(\mathbf{z}_t; \tilde{\boldsymbol{\mu}}_{\phi}(\dots), \text{diag}\{\tilde{\mathbf{v}}_{\phi}(\dots)\}\right).$$

One possible parametrization of the conditional posterior of \mathbf{z}_t is given as follows:

$$q_{\phi}(\mathbf{z}_t|\mathbf{z}_{1:t-1}, \mathbf{x}_{1:T}) = \mathcal{N}\left(\mathbf{z}_t; \tilde{\boldsymbol{\mu}}_{\phi}\left(\overset{\rightarrow}{\mathbf{h}}_t, \overset{\leftarrow}{\mathbf{h}}_t\right), \text{diag}\left\{\tilde{\mathbf{v}}_{\phi}\left(\overset{\rightarrow}{\mathbf{h}}_t, \overset{\leftarrow}{\mathbf{h}}_t\right)\right\}\right),$$

where:

- $\overset{\rightarrow}{\mathbf{h}}_t = \sigma(\overset{\rightarrow}{\mathbf{W}}_{xh}\mathbf{x}_{t-1} + \overset{\rightarrow}{\mathbf{W}}_{zh}\mathbf{z}_{t-1} + \overset{\rightarrow}{\mathbf{W}}_{hh}\overset{\rightarrow}{\mathbf{h}}_{t-1} + \overset{\rightarrow}{\mathbf{b}}_h)$ – encodes causal dependencies.
- $\overset{\leftarrow}{\mathbf{h}}_t = \sigma(\overset{\leftarrow}{\mathbf{W}}_{xh}\mathbf{x}_t + \overset{\leftarrow}{\mathbf{W}}_{hh}\overset{\leftarrow}{\mathbf{h}}_{t+1} + \overset{\leftarrow}{\mathbf{b}}_h)$ – encodes non-causal dependencies.

Learning: ELBO

The objective function is slightly different compared with standard VAEs:

$$\mathcal{L}(\mathbf{x}_{1:T}; \phi, \theta) = \sum_{t=1}^T \mathbb{E}_{q_{\phi}(\mathbf{z}_{1:t}|\mathbf{x}_{1:T})} [\ln p_{\theta_{\mathbf{x}}}(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t})] \\ - \sum_{t=1}^T \mathbb{E}_{q_{\phi}(\mathbf{z}_{1:t-1}|\mathbf{x}_{1:T})} \left[D_{\text{KL}} \left(q_{\phi}(\mathbf{z}_t|\mathbf{z}_{1:t-1}, \mathbf{x}_{1:T}) \parallel p_{\theta_{\mathbf{z}}}(\mathbf{z}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t-1}) \right) \right]$$

Learning: ELBO

The objective function is slightly different compared with standard VAEs:

$$\mathcal{L}(\mathbf{x}_{1:T}; \phi, \theta) = \sum_{t=1}^T \mathbb{E}_{q_{\phi}(\mathbf{z}_{1:t}|\mathbf{x}_{1:T})} [\ln p_{\theta_{\mathbf{x}}}(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t})] \\ - \sum_{t=1}^T \mathbb{E}_{q_{\phi}(\mathbf{z}_{1:t-1}|\mathbf{x}_{1:T})} \left[D_{\text{KL}} \left(q_{\phi}(\mathbf{z}_t|\mathbf{z}_{1:t-1}, \mathbf{x}_{1:T}) \parallel p_{\theta_{\mathbf{z}}}(\mathbf{z}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t-1}) \right) \right]$$

- The reconstruction and regularisation terms are evaluated at every frame t .
- Because of the model, the KL term depends on previous latent variables \Rightarrow sampling.
- The sampling occurs sequentially and cannot be paralelized!

DVAE Review & Code

Six (of the many) different models in the literature belong to the DVAE family:

- STORN J. Bayer and C. Osendorfer, Learning Stochastic Recurrent Networks, arXiv, 2014
- VRNN J. Chung et al., A recurrent latent variable model for sequential data, NeurIPS, 2015
- SRNN* M. Fraccaro et al., Sequential neural models with stochastic layers, NeurIPS, 2016
- DMM* R. Krishnan et al., Structured Inference Networks for Nonlinear State Space Models, AAAI, 2017
- DSAE Y. Li and S Mandt, Disentangled sequential autoencoder, ICML, 2018.
- RVAE* S. Leglaive et al., A recurrent variational autoencoder for speech enhancement, ICASSP 2020

Different conditional independence assumptions.

* the encoder is compliant with the true posterior.

DVAE Review & Code

Six (of the many) different models in the literature belong to the DVAE family:

STORN	J. Bayer and C. Osendorfer, Learning Stochastic Recurrent Networks, arXiv, 2014
VRNN	J. Chung et al., A recurrent latent variable model for sequential data, NeurIPS, 2015
SRNN*	M. Fraccaro et al., Sequential neural models with stochastic layers, NeurIPS, 2016
DMM*	R. Krishnan et al., Structured Inference Networks for Nonlinear State Space Models, AAAI, 2017
DSAE	Y. Li and S Mandt, Disentangled sequential autoencoder, ICML, 2018.
RVAE*	S. Leglaive et al., A recurrent variational autoencoder for speech enhancement, ICASSP 2020

Different conditional independence assumptions.

* the encoder is compliant with the true posterior.

Check “Dynamical Variational Autoencoders: A Comprehensive Review”, FnT ML:

- review and discuss these models (and more) with unified notations,
- compare their performance for analysis/resynthesis of speech signals,
- provide an implementation: <https://github.com/XiaoyuBIE1994/DVAE-speech>.

Results on speech analysis-resynthesis (reconstruction)

RESULTS OF THE SPEECH ANALYSIS-RESYNTHESIS EXPERIMENT,
AVERAGED OVER THE TEST SUBSET OF WSJ0 AND VOICEBANK.

Models	Dataset	SI-SDR (dB)	PESQ MOS	PESQ WB	PESQ NB	ESTOI
VAE	WSJ0	8.0	3.33	2.95	3.31	0.89
DKF	WSJ0	9.0	3.55	3.39	3.61	0.91
RVAE	WSJ0	9.8	3.65	3.57	3.75	0.92
SRNN	WSJ0	8.2	3.48	3.24	3.52	0.90
VAE	VB	8.6	3.22	2.79	3.15	0.88
DKF	VB	9.4	3.35	2.96	3.34	0.90
RVAE	VB	9.6	3.41	3.00	3.42	0.90
SRNN	VB	9.1	3.39	2.99	3.39	0.89

Various metrics, the SI-SDR is the scale-invariant signal-to-noise ratio.

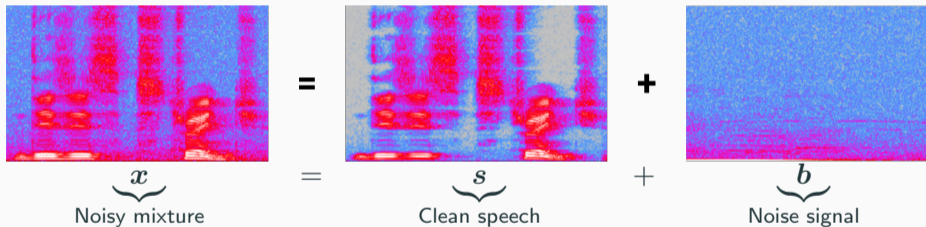
DKF/**RVAE**/SRNN outperform the VAE, which does not model cross-frame dependencies.

DVAEs for unsupervised speech enhancement

Unsupervised Speech Enhancement

Speech Enhancement: Remove the background noise from the observed mixture speech.

Short-time Fourier transform (STFT) is a time-frequency (matrix) representation.



Unsupervised (our work): Learn a **generative audio-visual model for clean speech** and combine it with an **unsupervised noise model** at test time

Unsupervised SE is **more flexible** since it adapts to various noises.

Unsupervised speech enhancement: overview

Train a **generative speech model** with clean data: $\{\mathbf{s}_i\}_{i=1}^N$.

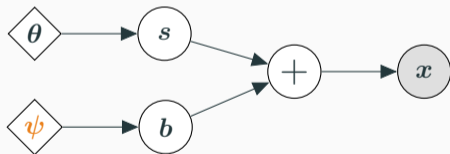


Unsupervised speech enhancement: overview

Train a **generative speech model** with clean data: $\{s_i\}_{i=1}^N$.



Test: learn the **noise parameters** of x :



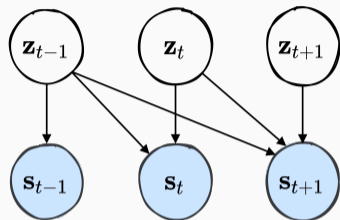
Note: at test time, θ is frozen, and s becomes a latent variable. Let's detail a bit more RVAE!

Recurrent Variational Autoencoder (RVAE)

Causal version (C-RVAE):

$$p_{\theta_s}(\mathbf{s}_{1:T}, \mathbf{z}_{1:T}) = \prod_{t=1}^T p(\mathbf{z}_t) p_{\theta_s}(\mathbf{s}_t | \mathbf{z}_{1:t}),$$

where $p(\mathbf{z}_t) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $p_{\theta_s}(\mathbf{s}_t | \mathbf{z}_{1:t}) = \mathcal{N}_c(\mathbf{s}_t; \mathbf{0}, \text{diag}\{v_{\theta_s}(\mathbf{z}_{1:t})\})$

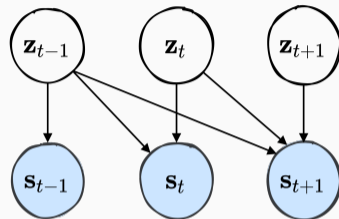


Recurrent Variational Autoencoder (RVAE)

Causal version (C-RVAE):

$$p_{\theta_s}(\mathbf{s}_{1:T}, \mathbf{z}_{1:T}) = \prod_{t=1}^T p(\mathbf{z}_t) p_{\theta_s}(\mathbf{s}_t | \mathbf{z}_{1:t}),$$

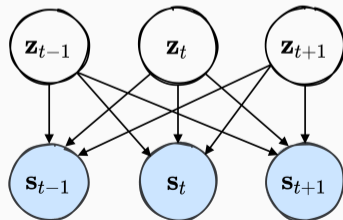
where $p(\mathbf{z}_t) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $p_{\theta_s}(\mathbf{s}_t | \mathbf{z}_{1:t}) = \mathcal{N}_c(\mathbf{s}_t; \mathbf{0}, \text{diag}\{\mathbf{v}_{\theta_s}(\mathbf{z}_{1:t})\})$



Non-causal version (NC-RVAE)

$$p_{\theta_s}(\mathbf{s}_{1:T}, \mathbf{z}_{1:T}) = \prod_{t=1}^T p(\mathbf{z}_t) p_{\theta_s}(\mathbf{s}_t | \mathbf{z}_{1:T}),$$

where $p(\mathbf{z}_t) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $p_{\theta_s}(\mathbf{s}_t | \mathbf{z}_{1:T}) = \mathcal{N}_c(\mathbf{s}_t; \mathbf{0}, \text{diag}\{\mathbf{v}_{\theta_s}(\mathbf{z}_{1:T})\})$



Noisy speech model

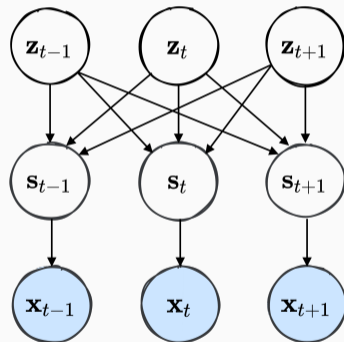
At **test time**, the clean speech signal is **latent**, and we consider the following generative model:

$$p_{\theta}(\mathbf{x}_{1:T}, \mathbf{s}_{1:T}, \mathbf{z}_{1:T}) = \prod_{t=1}^T p(\mathbf{z}_t) p_{\theta_s}(\mathbf{s}_t | \mathbf{z}_{1:T}) p_{\theta_x}(\mathbf{x}_t | \mathbf{s}_t),$$

where \mathbf{x}_t is a **noisy speech** frame, whose likelihood is defined by:

$$p_{\theta_x}(\mathbf{x}_t | \mathbf{s}_t) = \mathcal{N}_c(\mathbf{x}_t; \mathbf{s}_t, \text{diag}\{(\mathbf{W}\mathbf{H})_{:,t}\}),$$

with $\theta_x = \{\mathbf{W} \in \mathbb{R}_+^{F \times K}, \mathbf{H} \in \mathbb{R}_+^{K \times T}\}$.



EM Algorithm for Parameter Learning

Expectation-maximization (EM)-like algorithm:

- Approximate $p_{\theta}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})$
- Estimate $\theta_{\mathbf{x}} = \{\mathbf{W} \in \mathbb{R}_+^{F \times K}, \mathbf{H} \in \mathbb{R}_+^{K \times T}\}$

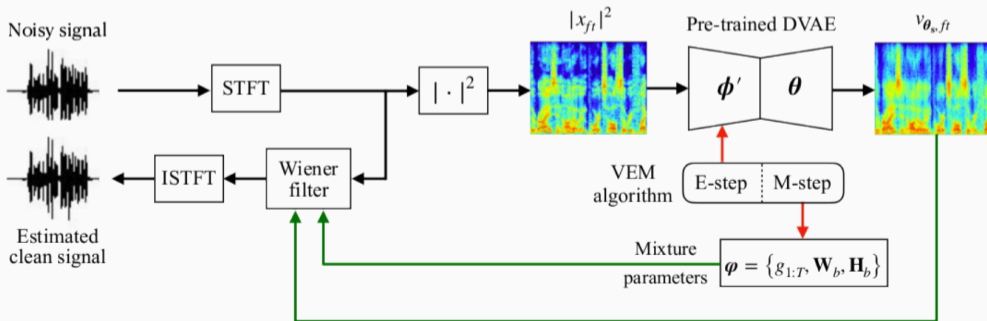
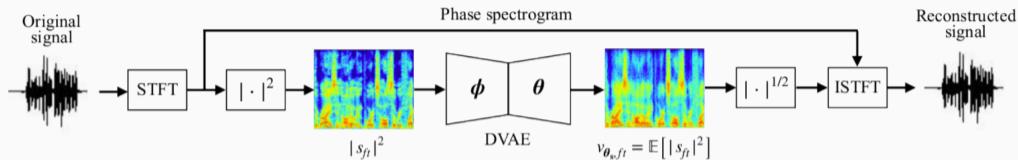
Three strategies for the intractable E-step:

- Markov Chain Monte Carlo
- Variational inference]: $p_{\theta}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T}) \approx q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})$
- $p_{\theta}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T}) \approx \delta(\mathbf{z}_{1:T} - \mathbf{z}_{1:T}^{\text{MAP}})$

Posterior mean estimate of the speech signal with Wiener-like filtering (element-wise ops):

$$\hat{\mathbf{s}}_t = \mathbb{E}_{p_{\theta}(\mathbf{s}_t|\mathbf{x}_{1:T})}[\mathbf{s}_t] = \mathbb{E}_{p_{\theta}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})} \left[\frac{\mathbf{v}_{\theta_s}(\mathbf{z}_{1:T})}{\mathbf{v}_{\theta_s}(\mathbf{z}_{1:T}) + (\mathbf{W}\mathbf{H})_{:,t}} \right] \mathbf{x}_t.$$

Pre-training and Parameter Learning



Results on the Wall Street Journal (WSJ) and VoiceBank (VB)

Method	Superv.	Test subset	Train subset	SI-SDR (dB)	Train subset	SI-SDR (dB)
Noisy mixture	-	WSJ0-QUT	-	-2.6	-	-2.6
VAE-VEM ICASSP'20	UA	WSJ0-QUT	WSJ0	5.0		
Proposed DKF-VEM	UA	WSJ0-QUT	WSJ0	5.1		
Proposed RVAE-VEM	UA	WSJ0-QUT	WSJ0	5.8		
Proposed SRNN-VEM	UA	WSJ0-QUT	WSJ0	5.2		
MetricGAN-U (full) ICASSP'22	UD	WSJ0-QUT	WSJ0-QUT	N/A		
MetricGAN-U (half) ICASSP'22	UD	WSJ0-QUT	WSJ0-QUT	N/A		
UMX* 2020	S	WSJ0-QUT	WSJ0-QUT	5.7		
MetricGAN+* Interspeech'21	S	WSJ0-QUT	WSJ0-QUT	3.6		
Noisy mixture	-					
VAE-VEM ICASSP'20	UA					
Proposed DKF-VEM	UA					
Proposed RVAE-VEM	UA					
Proposed SRNN-VEM	UA					
NyTT EUSIPCO'21	UD					
NyTT EUSIPCO'21	UD					
MetricGAN-U (full) ICASSP'22	UD					
MetricGAN-U (half) ICASSP'22	UD					
UMX 2020	S					
MetricGAN+ Interspeech'21	S					

Results on the Wall Street Journal (WSJ) and VoiceBank (VB)

Method	Superv.	Test subset	Train subset	SI-SDR (dB)	Train subset	SI-SDR (dB)
Noisy mixture	-	WSJ0-QUT	-	-2.6	-	-2.6
VAE-VEM ICASSP'20	UA	WSJ0-QUT	WSJ0	5.0		
Proposed DKF-VEM	UA	WSJ0-QUT	WSJ0	5.1		
Proposed RVAE-VEM	UA	WSJ0-QUT	WSJ0	5.8		
Proposed SRNN-VEM	UA	WSJ0-QUT	WSJ0	5.2		
MetricGAN-U (full) ICASSP'22	UD	WSJ0-QUT	WSJ0-QUT	N/A		
MetricGAN-U (half) ICASSP'22	UD	WSJ0-QUT	WSJ0-QUT	N/A		
UMX* 2020	S	WSJ0-QUT	WSJ0-QUT	5.7		
MetricGAN+* Interspeech'21	S	WSJ0-QUT	WSJ0-QUT	3.6		
Noisy mixture	-	VB-DMD	-	8.4	-	8.4
VAE-VEM ICASSP'20	UA	VB-DMD	VB	16.4		
Proposed DKF-VEM	UA	VB-DMD	VB	16.9		
Proposed RVAE-VEM	UA	VB-DMD	VB	17.1		
Proposed SRNN-VEM	UA	VB-DMD	VB	14.2		
NyTT EUSIPCO'21	UD	VB-DMD	VB-DMD (Xtra)	17.7		
NyTT EUSIPCO'21	UD	VB-DMD	VB-DMD	12.1		
MetricGAN-U (full) ICASSP'22	UD	VB-DMD	VB-DMD	6.5		
MetricGAN-U (half) ICASSP'22	UD	VB-DMD	VB-DMD	8.2		
UMX 2020	S	VB-DMD	VB-DMD	14.0		
MetricGAN+ Interspeech'21	S	VB-DMD	VB-DMD	8.5		

Results on the Wall Street Journal (WSJ) and VoiceBank (VB)

Method	Superv.	Test subset	Train subset	SI-SDR (dB)	Train subset	SI-SDR (dB)
Noisy mixture	-	WSJ0-QUT	-	-2.6	-	-2.6
VAE-VEM ICASSP'20	UA	WSJ0-QUT	WSJ0	5.0	VB	3.8
Proposed DKF-VEM	UA	WSJ0-QUT	WSJ0	5.1	VB	3.5
Proposed RVAE-VEM	UA	WSJ0-QUT	WSJ0	5.8	VB	4.3
Proposed SRNN-VEM	UA	WSJ0-QUT	WSJ0	5.2	VB	4.6
MetricGAN-U (full) ICASSP'22	UD	WSJ0-QUT	WSJ0-QUT	N/A	VB-DMD	-2.3
MetricGAN-U (half) ICASSP'22	UD	WSJ0-QUT	WSJ0-QUT	N/A	VB-DMD	-1.6
UMX* 2020	S	WSJ0-QUT	WSJ0-QUT	5.7	VB-DMD	4.1
MetricGAN+* Interspeech'21	S	WSJ0-QUT	WSJ0-QUT	3.6	VB-DMD	1.8
Noisy mixture	-	VB-DMD	-	8.4	-	8.4
VAE-VEM ICASSP'20	UA	VB-DMD	VB	16.4		
Proposed DKF-VEM	UA	VB-DMD	VB	16.9		
Proposed RVAE-VEM	UA	VB-DMD	VB	17.1		
Proposed SRNN-VEM	UA	VB-DMD	VB	14.2		
NyTT EUSIPCO'21	UD	VB-DMD	VB-DMD (Xtra)	17.7		
NyTT EUSIPCO'21	UD	VB-DMD	VB-DMD	12.1		
MetricGAN-U (full) ICASSP'22	UD	VB-DMD	VB-DMD	6.5		
MetricGAN-U (half) ICASSP'22	UD	VB-DMD	VB-DMD	8.2		
UMX 2020	S	VB-DMD	VB-DMD	14.0		
MetricGAN+ Interspeech'21	S	VB-DMD	VB-DMD	8.5		

Results on the Wall Street Journal (WSJ) and VoiceBank (VB)

Method	Superv.	Test subset	Train subset	SI-SDR (dB)	Train subset	SI-SDR (dB)
Noisy mixture	-	WSJ0-QUT	-	-2.6	-	-2.6
VAE-VEM ICASSP'20	UA	WSJ0-QUT	WSJ0	5.0	VB	3.8
Proposed DKF-VEM	UA	WSJ0-QUT	WSJ0	5.1	VB	3.5
Proposed RVAE-VEM	UA	WSJ0-QUT	WSJ0	5.8	VB	4.3
Proposed SRNN-VEM	UA	WSJ0-QUT	WSJ0	5.2	VB	4.6
MetricGAN-U (full) ICASSP'22	UD	WSJ0-QUT	WSJ0-QUT	N/A	VB-DMD	-2.3
MetricGAN-U (half) ICASSP'22	UD	WSJ0-QUT	WSJ0-QUT	N/A	VB-DMD	-1.6
UMX* 2020	S	WSJ0-QUT	WSJ0-QUT	5.7	VB-DMD	4.1
MetricGAN+* Interspeech'21	S	WSJ0-QUT	WSJ0-QUT	3.6	VB-DMD	1.8
Noisy mixture	-	VB-DMD	-	8.4	-	8.4
VAE-VEM ICASSP'20	UA	VB-DMD	VB	16.4	WSJ0	15.0
Proposed DKF-VEM	UA	VB-DMD	VB	16.9	WSJ0	16.8
Proposed RVAE-VEM	UA	VB-DMD	VB	17.1	WSJ0	textbf17.3
Proposed SRNN-VEM	UA	VB-DMD	VB	14.2	WSJ0	16.8
NyTT EUSIPCO'21	UD	VB-DMD	VB-DMD (Xtra)	17.7	WSJ0-QUT	N/A
NyTT EUSIPCO'21	UD	VB-DMD	VB-DMD	12.1	WSJ0-QUT	N/A
MetricGAN-U (full) ICASSP'22	UD	VB-DMD	VB-DMD	6.5	WSJ0-QUT	N/A
MetricGAN-U (half) ICASSP'22	UD	VB-DMD	VB-DMD	8.2	WSJ0-QUT	N/A
UMX 2020	S	VB-DMD	VB-DMD	14.0	WSJ0+QUT	10.4
MetricGAN+ Interspeech'21	S	VB-DMD	VB-DMD	8.5	WSJ0+QUT	3.9

Conclusions & Resources

- Unsupervised speech enhancement methods can be more “flexible”.
- Useful for unsupervised sequential deep modeling.
- Can be combined with other probabilistic models (at training or test time).
- Sequential sampling is not very efficient.

Samples: <https://team.inria.fr/robotlearn/unsup-se-dvae/>

Code: https://github.com/XiaoyuBIE1994/DVAE_SE

Paper: <https://arxiv.org/abs/2106.12271>

DVAE's can be used for other data/tasks: check
[https://team.inria.fr/robotlearn/dvae/!](https://team.inria.fr/robotlearn/dvae/)