

# Probabilistic generative models for audio-visual processing

---

Xavier Alameda-Pineda

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



# Why unsupervised learning?

Imagine a system for {person tracking, speech denoising, body pose estimation, ...} in:

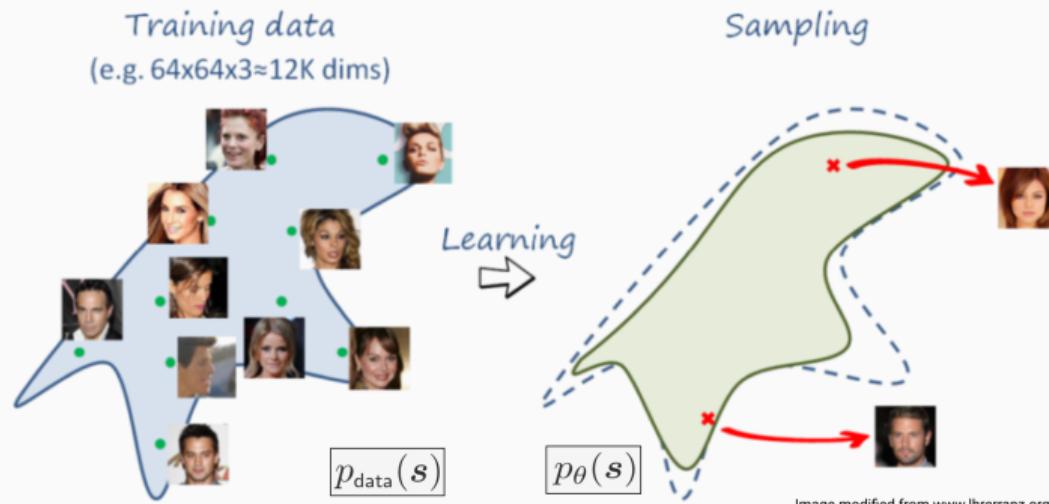


[Images under Creative Commons license]

We would like to avoid re-annotating for every new environment → Interest in unsupervised learning (and/or unsupervised domain adaptation).

## And probabilistic learning?

Probabilistic generative models aim to learn a (parametric) *distribution*  $p_{\theta}(x)$  that approximates the complex data distribution  $p_{\text{data}}(x)$ :



- Once learned, we can (ideally) sample new data.
- We can jointly learn them with other probabilistic models using *maximum likelihood*.

## The Kullback-Leibler divergence and the ML formulation

The Kullback-Leibler (KL) divergence between two distributions writes:

$$D_{\text{KL}}(p(\mathbf{x})\|q(\mathbf{x})) = \underbrace{-\mathbb{E}_{p(\mathbf{x})} \left[ \log \frac{q(\mathbf{x})}{p(\mathbf{x})} \right]}_{\text{Need your help!}} = - \int_{\mathcal{X}} p(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})} d\mathbf{x} \begin{cases} \geq 0 \\ \neq D_{\text{KL}}(q(\mathbf{x})\|p(\mathbf{x})) \end{cases}$$

$$D_{\text{KL}}(p(\mathbf{x})\|q(\mathbf{x})) = 0 \Leftrightarrow p(\mathbf{x}) = q(\mathbf{x}).$$

## The Kullback-Leibler divergence and the ML formulation

The Kullback-Leibler (KL) divergence between two distributions writes:

$$D_{\text{KL}}(p(\mathbf{x}) \| q(\mathbf{x})) = \underbrace{-\mathbb{E}_{p(\mathbf{x})} \left[ \log \frac{q(\mathbf{x})}{p(\mathbf{x})} \right]}_{\text{Need your help!}} = - \int_{\mathcal{X}} p(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})} d\mathbf{x} \begin{cases} \geq 0 \\ \neq D_{\text{KL}}(q(\mathbf{x}) \| p(\mathbf{x})) \end{cases}$$

$$D_{\text{KL}}(p(\mathbf{x}) \| q(\mathbf{x})) = 0 \Leftrightarrow p(\mathbf{x}) = q(\mathbf{x}).$$

Given a training set  $\{\mathbf{x}_i\}_{i=1}^N$ ,  $\mathbf{x}_i \sim p_{\text{data}}(\mathbf{x})$ , ML minimizes the KL divergence:

$$\boldsymbol{\theta}^* = \operatorname{argmin}_{\boldsymbol{\theta}} D_{\text{KL}}(p_{\text{data}}(\mathbf{x}) \| p_{\boldsymbol{\theta}}(\mathbf{x})) = \operatorname{argmin}_{\boldsymbol{\theta}} -\mathbb{E}_{p_{\text{data}}} \left[ \log \frac{p_{\boldsymbol{\theta}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x})} \right]$$

$$= \operatorname{argmax}_{\boldsymbol{\theta}} \mathbb{E}_{p_{\text{data}}} \left[ \log p_{\boldsymbol{\theta}}(\mathbf{x}) \right] \approx \boxed{\operatorname{argmax}_{\boldsymbol{\theta}} \frac{1}{N} \sum_{i=1}^N \log p_{\boldsymbol{\theta}}(\mathbf{x}_i)}$$

## Interest of Latent Variables

- Speech enhancement: noisy speech (observation), clean speech (latent variable)
- Person tracking: detections (observation), person positions (latent variable)
- Representation learning: raw data (observation), representation (latent variable)

## Interest of Latent Variables

- Speech enhancement: noisy speech (observation), clean speech (latent variable)
- Person tracking: detections (observation), person positions (latent variable)
- Representation learning: raw data (observation), representation (latent variable)

Let  $z$  denote the latent variable:

$$\begin{cases} z \sim p_{\theta}(z) \\ x|z \sim p_{\theta}(x|z) \end{cases} \rightarrow p_{\theta}(x) = \int p_{\theta}(x, z) dz = \int p_{\theta}(x|z)p_{\theta}(z) dz$$

## Interest of Latent Variables

- Speech enhancement: noisy speech (observation), clean speech (latent variable)
- Person tracking: detections (observation), person positions (latent variable)
- Representation learning: raw data (observation), representation (latent variable)

Let  $z$  denote the latent variable:

$$\begin{cases} z \sim p_{\theta}(z) \\ x|z \sim p_{\theta}(x|z) \end{cases} \rightarrow p_{\theta}(x) = \int p_{\theta}(x, z) dz = \int p_{\theta}(x|z)p_{\theta}(z) dz$$

**New samples:** Draw  $\hat{z} \sim p_{\theta}(z)$ , then draw a new sample  $\hat{x} \sim p_{\theta}(x|\hat{z})$ .

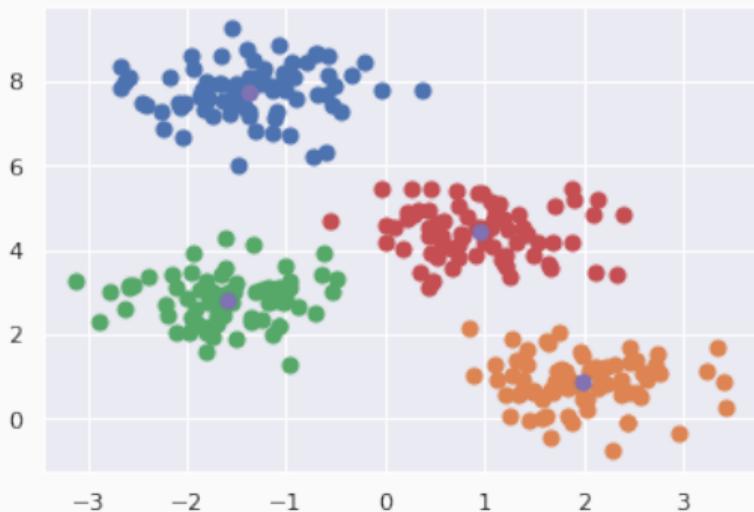
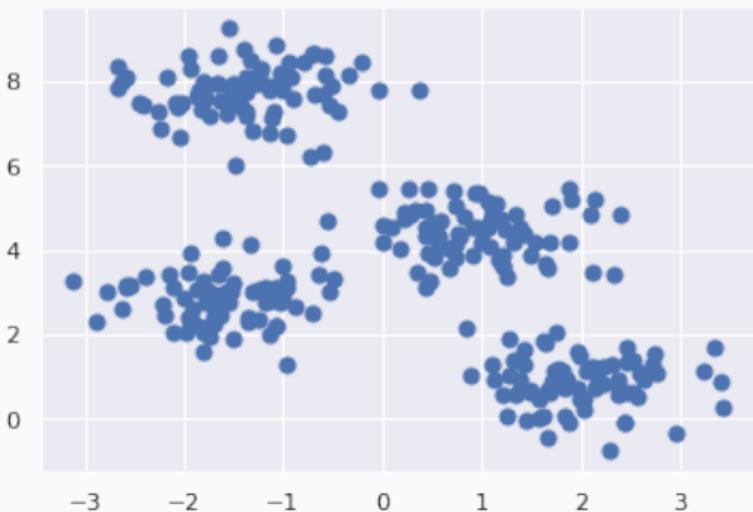
**Learning:** How to estimate  $\theta$  with the ML formulation?

**Inference:** How to get the most likely latent  $p_{\theta}(z|x)$ ?

**From shallow GMMs to deep VAEs**

## Simple example: clustering

Definition: find groups of data points without labels.



## The Gaussian Mixture Model (GMM)

Definition:

- For each  $x_n$  there is a latent variable  $z_n$  taking values from 1 to  $K$ :  $z_n \in \{1, \dots, K\}$ .

## The Gaussian Mixture Model (GMM)

Definition:

- For each  $x_n$  there is a latent variable  $z_n$  taking values from 1 to  $K$ :  $z_n \in \{1, \dots, K\}$ .
- Its prior probability is defined as:  $p(z_n = k) = \pi_k \geq 0$ , with  $\sum_{k=1}^K \pi_k = 1$ .

# The Gaussian Mixture Model (GMM)

Definition:

- For each  $x_n$  there is a latent variable  $z_n$  taking values from 1 to  $K$ :  $z_n \in \{1, \dots, K\}$ .
- Its prior probability is defined as:  $p(z_n = k) = \pi_k \geq 0$ , with  $\sum_{k=1}^K \pi_k = 1$ .
- Given  $z_n$ , the data point is modeled as a multivariate Gaussian:

$$p(x_n | z_n = k) = \mathcal{N}(x_n; \mu_k, \Sigma_k)$$

Advantages:

- ① Having  $\pi_1, \dots, \pi_K$  means that groups can be differently populated.
- ② The shape of the groups is modeled by  $\Sigma_k$ .
- ③ The parameters are:  $\theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$ .

## Maximum likelihood for GMM

Let's compute  $p(\mathbf{x}_n)$

$$p(\mathbf{x}_n) = \sum_{k=1}^K p(\mathbf{x}_n, z_n = k) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k).$$

The log-likelihood:

$$\mathcal{L}(\boldsymbol{\theta} | \mathbf{X}) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k),$$

Computing directly ML by  $\frac{\partial \mathcal{L}}{\partial \pi_k} = 0$ ,  $\frac{\partial \mathcal{L}}{\partial \boldsymbol{\mu}_k} = 0$  or  $\frac{\partial \mathcal{L}}{\partial \boldsymbol{\Sigma}_k} = 0$  is very difficult.

## (Expected complete-data) log-likelihood

We have seen that  $\log p(\mathbf{x})$  does not work well with derivatives. However,  $\log p(\mathbf{x}, \mathbf{z})$  does!

**Problem:**  $\mathbf{z}$  is not observed, thus  $\sum_n \log p(\mathbf{x}_n, \mathbf{z}_n)$  is a random variable.

## (Expected complete-data) log-likelihood

We have seen that  $\log p(\mathbf{x})$  does not work well with derivatives. However,  $\log p(\mathbf{x}, \mathbf{z})$  does!

**Problem:**  $\mathbf{z}$  is not observed, thus  $\sum_n \log p(\mathbf{x}_n, \mathbf{z}_n)$  is a random variable.

Given an initial value of the parameters,  $\boldsymbol{\theta}^0$ , let's take the expectation w.r.t. the posterior distribution  $p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}^0)$  [we will justify this choice later on]:

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^0) = \sum_n \mathbb{E}_{p(\mathbf{z}_n|\mathbf{x}_n; \boldsymbol{\theta}^0)} \log p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta})$$

This function is called: *expected complete-data log-likelihood*, and is the main mathematical object when working with EM algorithms.

## The EM algorithm for GMM

Given  $\boldsymbol{\theta}^0$ , we use the *expected complete-data log-likelihood*  $\mathcal{Q}$ . For iteration  $r = 1, \dots, R$ :

- ① Expectation [E-step]:

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{r-1}) = \mathbb{E}_{p(z|x; \boldsymbol{\theta}^{r-1})} \log p(x, z; \boldsymbol{\theta})$$

- ② Maximisation [M-step]:

$$\boldsymbol{\theta}^r = \arg \max_{\boldsymbol{\theta}} \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{r-1})$$

(observations  $x = \{x_n\}_{n=1}^N$ , latent variables  $z = \{z_n\}_{n=1}^N$ , parameters  $\boldsymbol{\theta} = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$ )

## The EM algorithm for GMM (II)

**E-step** The posterior distribution writes:

$$p(z_n = k | \mathbf{x}_n; \boldsymbol{\theta}^0) = \frac{p(z_n = k; \boldsymbol{\theta}^0)p(\mathbf{x}_n | z_n = k; \boldsymbol{\theta}^0)}{\sum_{\ell} p(z_n = \ell; \boldsymbol{\theta}^0)p(\mathbf{x}_n | z_n = \ell; \boldsymbol{\theta}^0)} = \frac{\pi_k^0 \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_k^0, \boldsymbol{\Sigma}_k^0)}{\sum_{\ell} \pi_{\ell}^0 \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_{\ell}^0, \boldsymbol{\Sigma}_{\ell}^0)} = \eta_{nk}$$

## The EM algorithm for GMM (II)

**E-step** The posterior distribution writes:

$$p(z_n = k | \mathbf{x}_n; \boldsymbol{\theta}^0) = \frac{p(z_n = k; \boldsymbol{\theta}^0)p(\mathbf{x}_n | z_n = k; \boldsymbol{\theta}^0)}{\sum_{\ell} p(z_n = \ell; \boldsymbol{\theta}^0)p(\mathbf{x}_n | z_n = \ell; \boldsymbol{\theta}^0)} = \frac{\pi_k^0 \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_k^0, \boldsymbol{\Sigma}_k^0)}{\sum_{\ell} \pi_{\ell}^0 \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_{\ell}^0, \boldsymbol{\Sigma}_{\ell}^0)} = \eta_{nk}$$

The expected complete-data log-likelihood writes:

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^0) = \sum_{n=1}^N \sum_{k=1}^K \eta_{nk} \log \pi_k \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

## The EM algorithm for GMM (II)

**E-step** The posterior distribution writes:

$$p(\mathbf{z}_n = k | \mathbf{x}_n; \boldsymbol{\theta}^0) = \frac{p(\mathbf{z}_n = k; \boldsymbol{\theta}^0)p(\mathbf{x}_n | \mathbf{z}_n = k; \boldsymbol{\theta}^0)}{\sum_{\ell} p(\mathbf{z}_n = \ell; \boldsymbol{\theta}^0)p(\mathbf{x}_n | \mathbf{z}_n = \ell; \boldsymbol{\theta}^0)} = \frac{\pi_k^0 \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_k^0, \boldsymbol{\Sigma}_k^0)}{\sum_{\ell} \pi_{\ell}^0 \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_{\ell}^0, \boldsymbol{\Sigma}_{\ell}^0)} = \eta_{nk}$$

The expected complete-data log-likelihood writes:

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^0) = \sum_{n=1}^N \sum_{k=1}^K \eta_{nk} \log \pi_k \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

**M-step** The optimal value for  $\boldsymbol{\mu}_k$  writes:

$$\boldsymbol{\mu}_k^* = \frac{1}{\sum_{n=1}^N \eta_{nk}} \sum_{n=1}^N \eta_{nk} \mathbf{x}_n$$

## But why does the EM work?

The main mathematical object in EM is  $\mathcal{Q}$  (The expected complete-data log-likelihood).

What is the relationship with the log-likelihood? Let's take any distribution of  $z$ ,  $q(z)$ :

$$\log p(\mathbf{x}) = \underbrace{\mathbb{E}_{q(\mathbf{z})} \left[ \log \frac{p(\mathbf{x})p(\mathbf{z}|\mathbf{x})}{q(\mathbf{z})} \right]}_{\text{M-step - } \mathcal{Q}} + \underbrace{D_{\text{KL}}(q(\mathbf{z}) \parallel p(\mathbf{z}|\mathbf{x}))}_{\text{E-step - KL}}$$

## But why does the EM work?

The main mathematical object in EM is  $\mathcal{Q}$  (The expected complete-data log-likelihood).

What is the relationship with the log-likelihood? Let's take any distribution of  $z$ ,  $q(z)$ :

$$\log p(\mathbf{x}) = \underbrace{\mathbb{E}_{q(z)} \left[ \log \frac{p(\mathbf{x})p(z|\mathbf{x})}{q(z)} \right]}_{\text{M-step - } \mathcal{Q}} + \underbrace{D_{\text{KL}}(q(z) \parallel p(z|\mathbf{x}))}_{\text{E-step - KL}}$$

Another interpretation. Given  $\theta^0$ :

- ① Set  $q(z) = p(z|\mathbf{x}; \theta^0)$ , minimise KL (thus maximize  $\mathcal{Q}$ ) w.r.t.  $q(z)$ . The E-step reduces the distance between log-likelihood and  $\mathcal{Q}$ .
- ② Maximize  $\mathcal{Q}$  w.r.t.  $\theta$ :

$$\mathcal{Q}(\theta, \theta^0) = \mathbb{E}_{q(z)} \left[ \log \frac{p(\mathbf{x}, z; \theta)}{q(z)} \right].$$

The M-step pushes  $\mathcal{Q}$  and therefore pushes the log-likelihood.

## The Exact EM

$$\log p(\mathbf{x}; \boldsymbol{\theta}) = \underbrace{\mathbb{E}_{q(\mathbf{z})} \left[ \log \frac{p(\mathbf{x})p(\mathbf{z}|\mathbf{x})}{q(\mathbf{z})} \right]}_{\text{M-step - } \mathcal{Q}} + \underbrace{D_{\text{KL}} \left( q(\mathbf{z}) \middle\| p(\mathbf{z}|\mathbf{x}) \right)}_{\text{E-step - KL}}$$

Given  $\boldsymbol{\theta}^0$ :

- ① Set  $q(z) = p(z|\mathbf{x}; \boldsymbol{\theta}^0)$ , and compute  $\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^0)$ .
- ② Maximize  $\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^0)$  w.r.t.  $\boldsymbol{\theta}$ .

## The Exact EM

$$\log p(\mathbf{x}; \boldsymbol{\theta}) = \underbrace{\mathbb{E}_{q(\mathbf{z})} \left[ \log \frac{p(\mathbf{x})p(\mathbf{z}|\mathbf{x})}{q(\mathbf{z})} \right]}_{\text{M-step - } \mathcal{Q}} + \underbrace{D_{\text{KL}} \left( q(\mathbf{z}) \middle\| p(\mathbf{z}|\mathbf{x}) \right)}_{\text{E-step - KL}}$$

Given  $\boldsymbol{\theta}^0$ :

- ① Set  $q(z) = p(z|\mathbf{x}; \boldsymbol{\theta}^0)$ , and compute  $\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^0)$ .
- ② Maximize  $\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^0)$  w.r.t.  $\boldsymbol{\theta}$ .

Several potential issues:

- The posterior distribution does not exist/is not computationally tractable.
- The expectation cannot be taken analytically.
- The maximization w.r.t.  $\boldsymbol{\theta}$  does not have close-form solution.

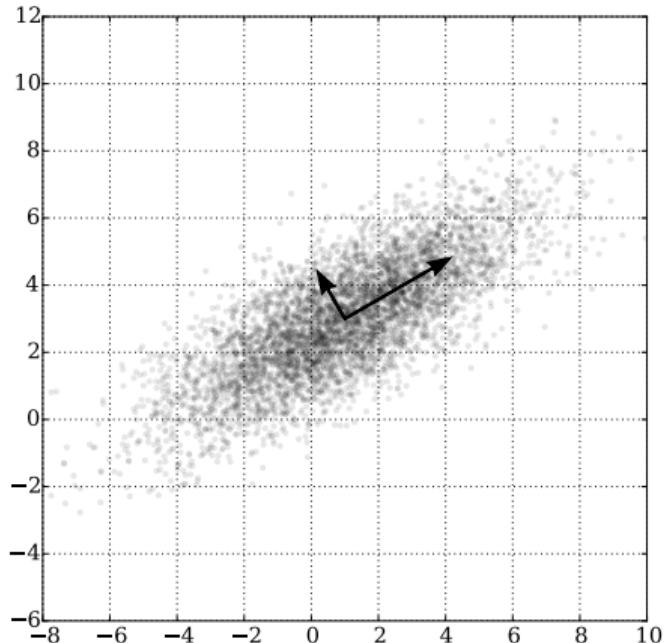
## What about continuous latent variables?

E.g.: extract a representation ( $z$ ) of each ( $x$ ).

Linear model (PPCA) ( $D \ll F$ ):

$$p(z) = \mathcal{N}(z; \mathbf{0}, I), \quad z \in \mathbb{R}^D.$$

$$p(x|z) = \mathcal{N}(x; Az + b, \nu I), \quad x \in \mathbb{R}^F.$$



[Image from Wikimedia Commons]

## What about continuous latent variables?

E.g.: extract a representation ( $z$ ) of each ( $x$ ).

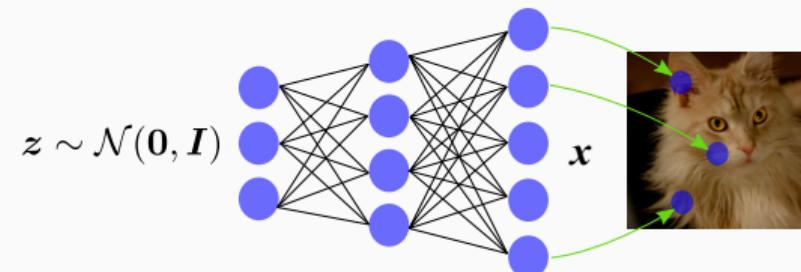
Linear model (PPCA) ( $D \ll F$ ):

$$p(z) = \mathcal{N}(z; \mathbf{0}, I), \quad z \in \mathbb{R}^D.$$

$$p(x|z) = \mathcal{N}(x; Az + b, \nu I), \quad x \in \mathbb{R}^F.$$

Non-linear model (same  $p(z)$ ):

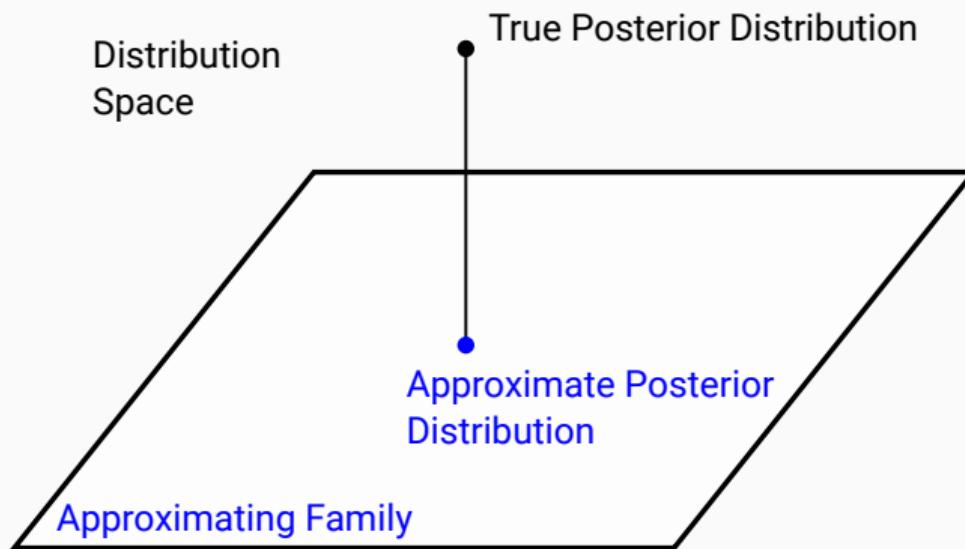
$$p_{\theta}(x|z) = \mathcal{N}\left(x; \mu_{\theta}(z), \Sigma_{\theta}(z)\right)$$



## The posterior distribution

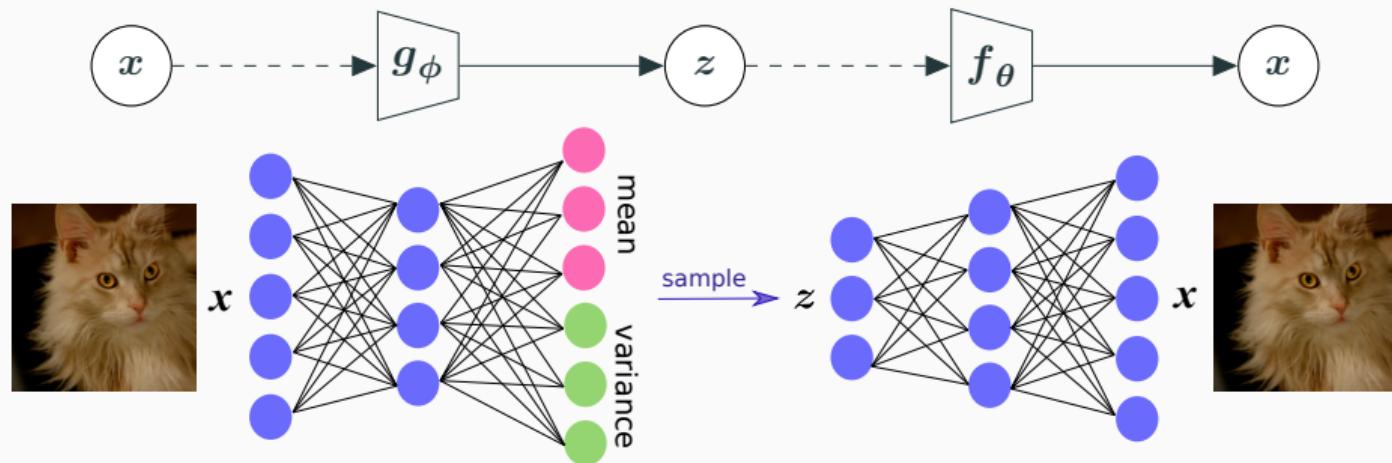
If  $p_{\theta}(x|z)$  is non-linear,  $p_{\theta}(z|x)$  cannot be computed analytically, it needs to be approximated.

The idea is to find the **best candidate within a family of distributions**.



## Overall architecture

If we “chain” the posterior (encoder) and the generative (decoder) model:



$$p(z|x) \approx q(z|x) = \mathcal{N}\left(z; \tilde{\mu}_\phi(x), \tilde{\Sigma}_\phi(x)\right) \quad p_\theta(x|z) = \mathcal{N}\left(x; \mu_\theta(z), \Sigma_\theta(z)\right)$$

But how do we optimise for the parameters  $\theta$  and  $\phi$ ?

## Learning - ELBO

If we recall the formulation for the EM:

$$\log p(\mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z}|\mathbf{x})} \right] + D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x}))$$

## Learning - ELBO

If we recall the formulation for the EM:

$$\log p(\mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z}|\mathbf{x})} \right] + D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x}))$$

Problem: the second term cannot be computed! But it's positive:

$$\begin{aligned} \log p(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\phi}) &\geq \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} \right] \\ \log p(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\phi}) &\geq \underbrace{\mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} \left[ \log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}) \right]}_{\text{Reconstruction}} - \underbrace{D_{\text{KL}}(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}))}_{\text{Regularisation}} \end{aligned}$$

This is known as **Evidence Lower-BOund or ELBO**:  $\mathcal{L}_{\text{ELBO}}(\boldsymbol{\theta}, \boldsymbol{\phi})$ .

Be VERY careful with these expressions: They look alike, but they are NOT the same.

## Learning - Practical comments

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z}|\mathbf{x})} \right] = \mathcal{L}_{\text{ELBO}}(\boldsymbol{\theta}, \boldsymbol{\phi})$$

- One needs to sample from  $q(\mathbf{z}|\mathbf{x})$  to compute the expectation.
- This sampling needs the “reparametrisation trick” to be differentiable.
- Most libraries solve a minimisation problem → we need to provide  $-\mathcal{L}_{\text{ELBO}}(\boldsymbol{\theta}, \boldsymbol{\phi})$ .

## ML with latent variables: EM vs VAE?

ML with latent variable ( $z$ ) leads to EM<sup>1</sup> and VI<sup>2</sup> build from ( $q(z)$  is an arbitrary distribution):

$$\log p(\mathbf{x}) = \underbrace{\mathbb{E}_{q(z)} \left[ \log \frac{p(\mathbf{x})p(z|\mathbf{x})}{q(z)} \right]}_{\text{M-step or VLB}} + \underbrace{D_{\text{KL}}(q(z) \| p(z|\mathbf{x}))}_{\text{E-step}}$$

### Exact EM

Simple  $p(z|\mathbf{x})$

$$q(z) = p(z|\mathbf{x})$$

Closed-form!

$$D_{\text{KL}} = 0$$

### Variational AutoEncoder

$$p(z|\mathbf{x}) ???$$

$$q(z) = q_\phi(z)$$

$q_\phi$  optimises ELBO/VLB

$D_{\text{KL}} > 0$  no closed form!

<sup>1</sup>Dempster, A.P., et. al., (1977), Journal of the Royal Statistical Society.

<sup>2</sup>Jordan, M. I., et. al., (1999), Machine Learning.

## ML with latent variables: EM vs VAE?

ML with latent variable ( $z$ ) leads to EM<sup>1</sup> and VI<sup>2</sup> build from ( $q(z)$  is an arbitrary distribution):

$$\log p(\mathbf{x}) = \underbrace{\mathbb{E}_{q(z)} \left[ \log \frac{p(\mathbf{x})p(z|\mathbf{x})}{q(z)} \right]}_{\text{M-step or VLB}} + \underbrace{D_{\text{KL}}(q(z) \| p(z|\mathbf{x}))}_{\text{E-step}}$$

### Exact EM

Simple  $p(z|\mathbf{x})$

$q(z) = p(z|\mathbf{x})$

Closed-form!

$$D_{\text{KL}} = 0$$

### Variational EM

$p(z|\mathbf{x})$  comp. complex

$$q([\mathbf{z}_1, \mathbf{z}_2]) = q_1(\mathbf{z}_1)q_2(\mathbf{z}_2)$$

$$q_1, q_2 = \operatorname{argmin} D_{\text{KL}}$$

$D_{\text{KL}} > 0$  but closed form!

### Variational AutoEncoder

$p(z|\mathbf{x})$  ???

$$q(z) = q_\phi(z)$$

$q_\phi$  optimises ELBO/VLB

$D_{\text{KL}} > 0$  no closed form!

<sup>1</sup>Dempster, A.P., et. al., (1977), Journal of the Royal Statistical Society.

<sup>2</sup>Jordan, M. I., et. al., (1999), Machine Learning.

# Audio-Visual Unsupervised Speech Enhancement with VAEs



Mostafa  
Sadeghi



Laurent Girin

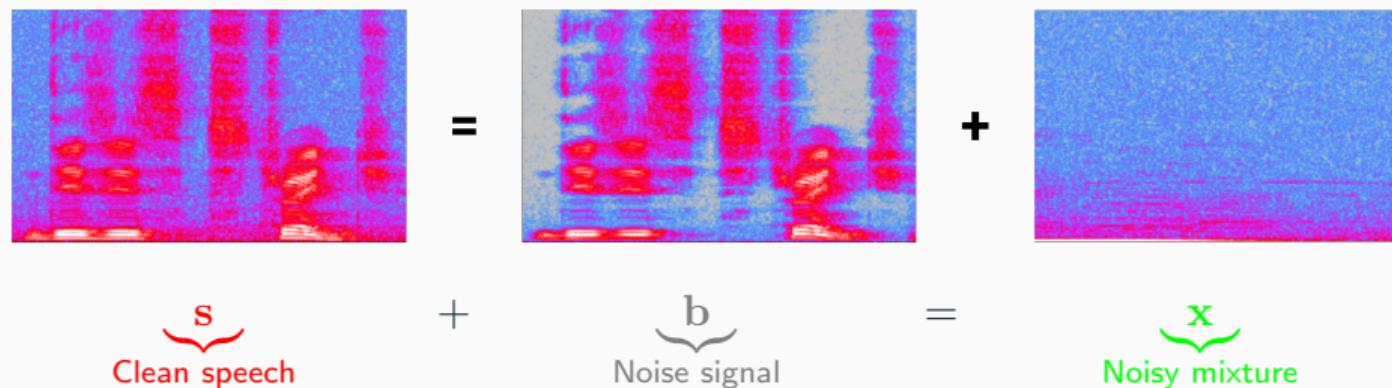


Radu Horaud

# Speech Enhancement & Wiener Filter



Extract the **latent clean speech signal** from the **observed noisy mixture**. (STFT domain)



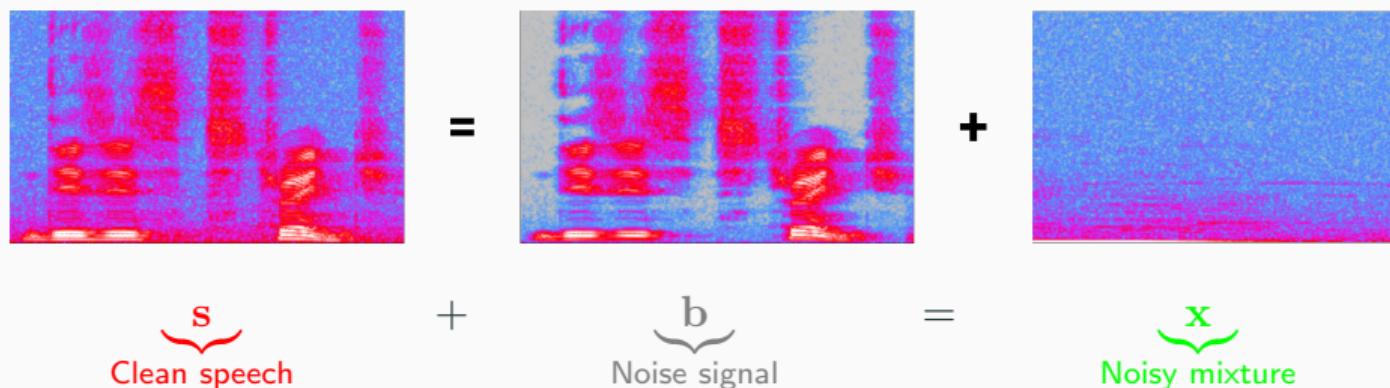
# Speech Enhancement & Wiener Filter



Minimize MSE → Wiener Filter  
(operations are element-wise)

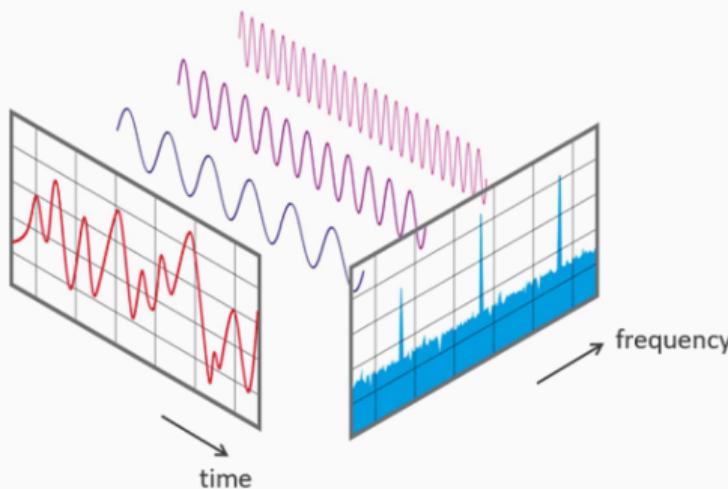
$$\hat{s} = \frac{\sigma_s}{\sigma_s + \sigma_b} x$$

Extract the **latent clean speech signal** from the **observed noisy mixture**. (STFT domain)



## Representing audio: time vs. frequency

Fourier domain decomposes the signal in frequencies:

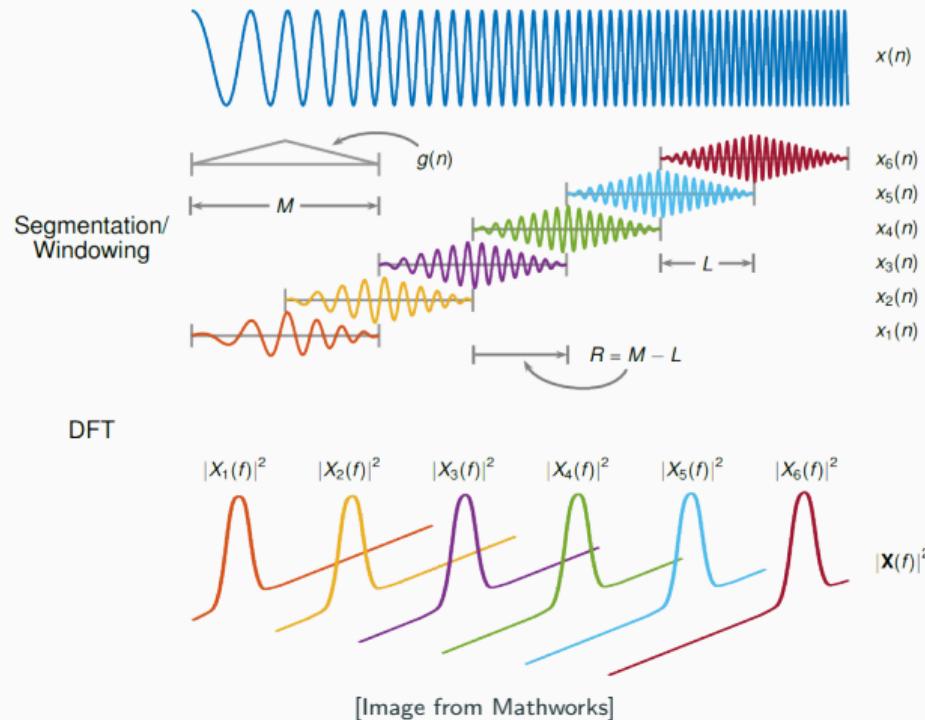


[Image from <https://dev.to/trekhleb/>]

Problem: we have to choose either time **or** frequency.

# The short-time Fourier transform (STFT)

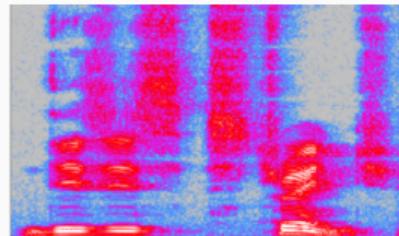
STFT: segment the input signal, and apply DFT to each segment.



## Let's see it!

- Sine frequency sweep (pure sine of increasing frequency).
- Leyenda (piano, harmonics).
- O mio babbino caro (opera, vibrato).
- A few good men (speech, separation).
- ...

# Unsupervised Probabilistic SE: paradigm



   
Clean speech

   
Noise signal

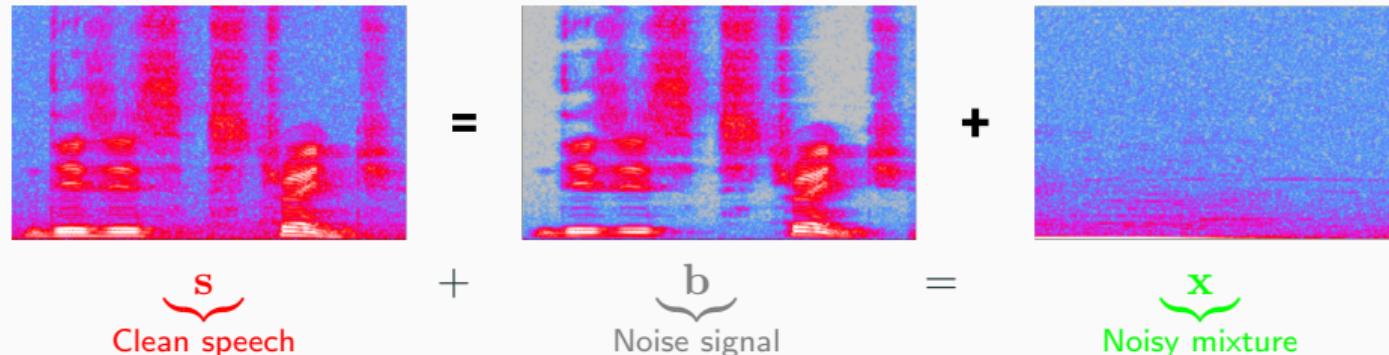
   
Noisy mixture

**Train** – Model for clean data:  $\{\mathbf{s}_i\}_{i=1}^N$ .



Then freeze  $\theta$  at test/adaptation time.

# Unsupervised Probabilistic SE: paradigm

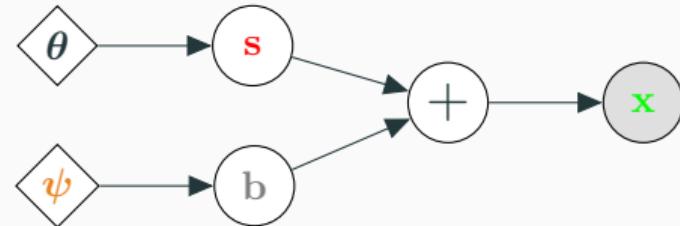


**Train** – Model for clean data:  $\{s_i\}_{i=1}^N$ .



$$p_\theta(s) \approx p_{\text{data}}(s)$$

Then freeze  $\theta$  at test/adaptation time.



**Test** – learn the **noise parameters** from **noisy samples**  $x$  and estimate the **clean speech**  $\hat{s}$ :

# Audio-visual Speech Enhancement (AV-SE)

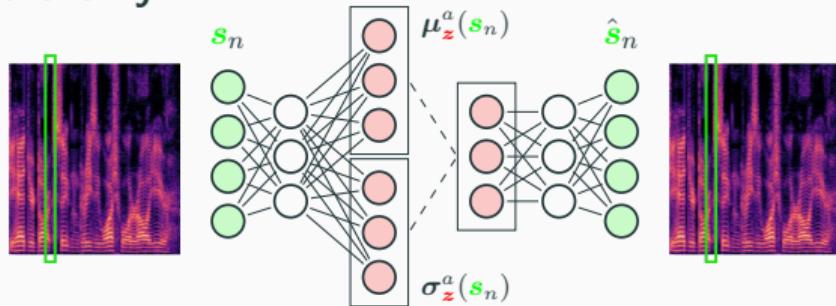
- Visual data (lip motion) provide **complementary information** about the unknown speech.
- For **highly noisy audio recordings**, visual information can be very helpful.



*We investigate the VAE framework to fuse audio and visual data for speech enhancement.*

# Mono-modal VAEs: the baselines

Audio-only:<sup>3</sup>

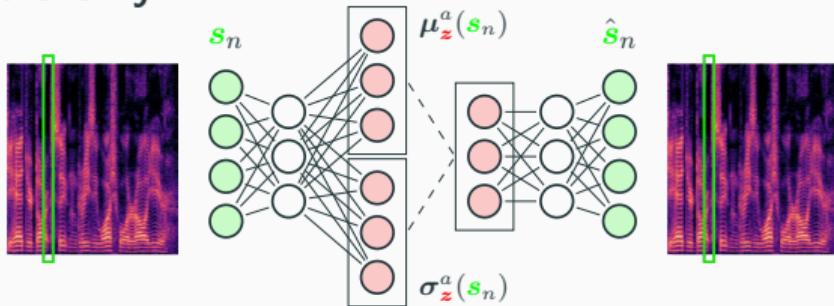


$$p_{\theta}(s_n | z_n) = \mathcal{N}_c(\mathbf{0}, \text{diag}(\sigma_s(z_n)))$$
$$q_{\phi}(z_n | s_n) = \mathcal{N}\left(\mu_z^a(s_n), \text{diag}(\sigma_z^a(s_n))\right)$$

<sup>3</sup>Leglaive, S., et. al., (2018), IEEE MLSP.

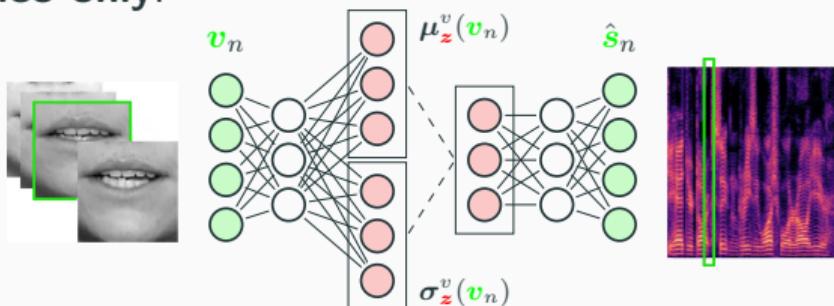
# Mono-modal VAEs: the baselines

## Audio-only:<sup>3</sup>



$$p_{\theta}(s_n | z_n) = \mathcal{N}_c(\mathbf{0}, \text{diag}(\sigma_s(z_n)))$$
$$q_{\phi}(z_n | s_n) = \mathcal{N}\left(\mu_z^a(s_n), \text{diag}(\sigma_z^a(s_n))\right)$$

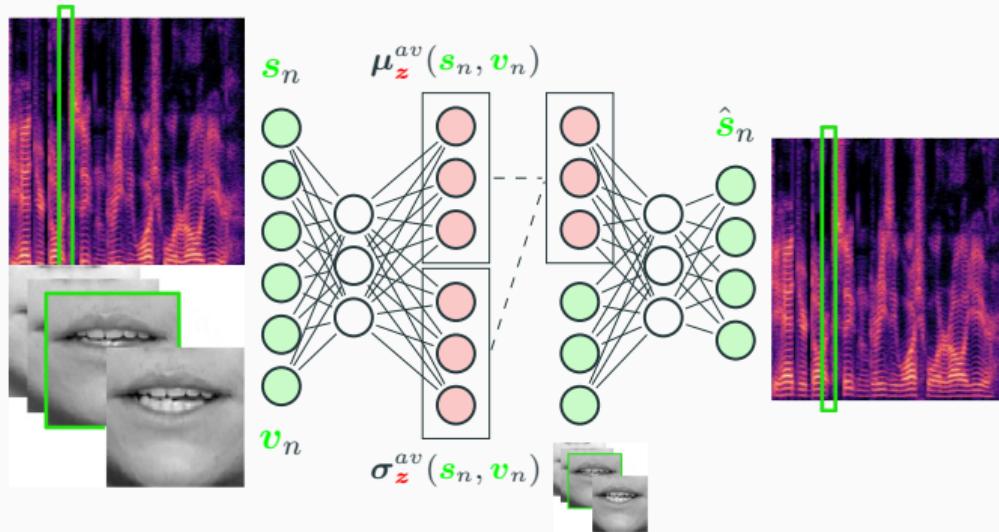
## Video-only:



$$p_{\theta}(s_n | z_n) = \mathcal{N}_c(\mathbf{0}, \text{diag}(\sigma_s(z_n)))$$
$$q_{\phi}(z_n | v_n) = \mathcal{N}\left(\mu_z^v(v_n), \text{diag}(\sigma_z^v(v_n))\right)$$

<sup>3</sup>Leglaive, S., et. al., (2018), IEEE MLSP.

# Audio-visual Conditional VAE<sup>4</sup>

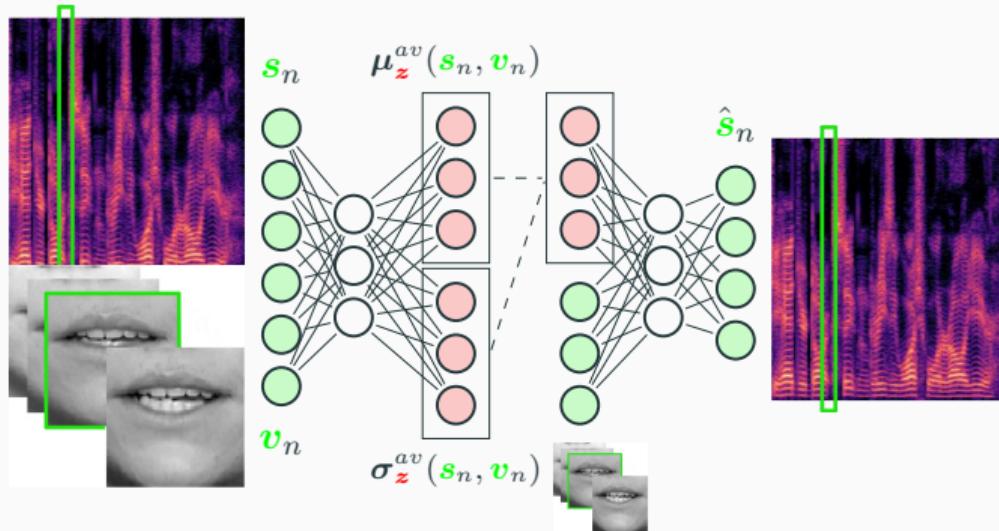


$$p_{\theta}(\mathbf{s}_n | \mathbf{z}_n, \mathbf{v}_n) = \mathcal{N}_c(\mathbf{0}, \text{diag}(\sigma_{\mathbf{s}}(\mathbf{z}_n, \mathbf{v}_n)))$$

$$q_{\phi}(\mathbf{z}_n | \mathbf{v}_n, \mathbf{s}_n) = \mathcal{N}\left(\mu_{\mathbf{z}}^{av}(\mathbf{v}_n, \mathbf{s}_n), \text{diag}(\sigma_{\mathbf{z}}^{av}(\mathbf{v}_n, \mathbf{s}_n))\right)$$

<sup>4</sup>Sadeghi, M., et. al., (2020), IEEE TASLP.

# Audio-visual Conditional VAE<sup>4</sup>



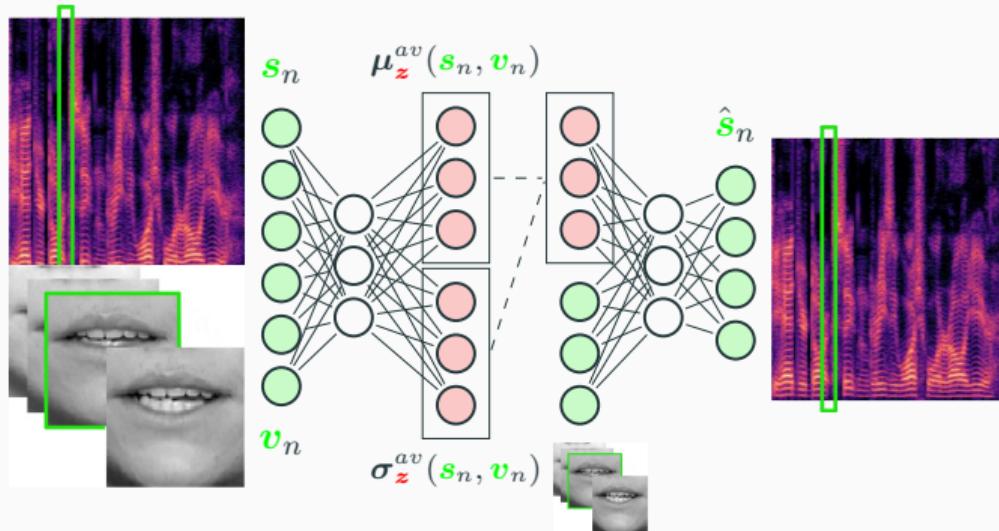
$$p_{\theta}(\mathbf{s}_n | \mathbf{z}_n, \mathbf{v}_n) = \mathcal{N}_c(\mathbf{0}, \text{diag}(\boldsymbol{\sigma}_{\mathbf{s}}(\mathbf{z}_n, \mathbf{v}_n)))$$

$$q_{\phi}(\mathbf{z}_n | \mathbf{v}_n, \mathbf{s}_n) = \mathcal{N}\left(\boldsymbol{\mu}_z^{av}(\mathbf{v}_n, \mathbf{s}_n), \text{diag}(\boldsymbol{\sigma}_z^{av}(\mathbf{v}_n, \mathbf{s}_n))\right)$$

<sup>4</sup>Sadeghi, M., et. al., (2020), IEEE TASLP.

What will this learn?

# Audio-visual Conditional VAE<sup>4</sup>



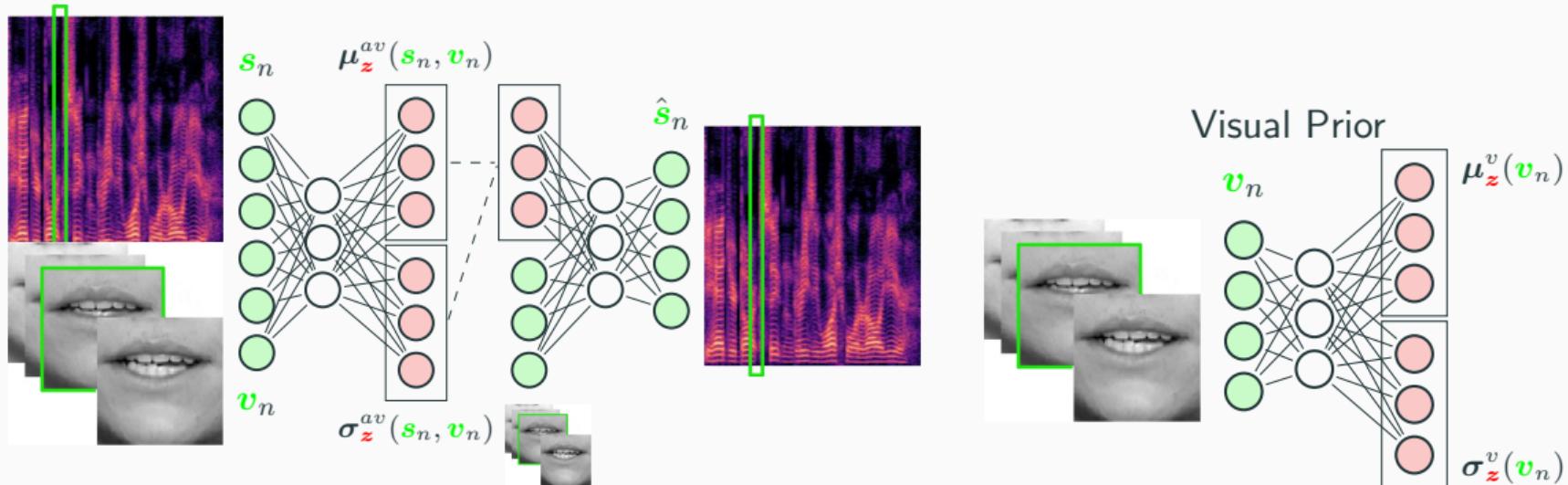
$$p_{\theta}(s_n | z_n, v_n) = \mathcal{N}_c(\mathbf{0}, \text{diag}(\sigma_s(z_n, v_n)))$$

$$q_{\phi}(z_n | v_n, s_n) = \mathcal{N}\left(\mu_z^{av}(v_n, s_n), \text{diag}(\sigma_z^{av}(v_n, s_n))\right)$$

<sup>4</sup>Sadeghi, M., et. al., (2020), IEEE TASLP.

What will this learn?  
To ignore the video!

# Audio-visual Conditional VAE<sup>4</sup>



$$p_{\theta}(\mathbf{s}_n | \mathbf{z}_n, \mathbf{v}_n) = \mathcal{N}_c(\mathbf{0}, \text{diag}(\sigma_s(\mathbf{z}_n, \mathbf{v}_n)))$$

$$p_{\theta}(\mathbf{z}_n | \mathbf{v}_n) = \mathcal{N}\left(\mu_{\mathbf{z}}^v(\mathbf{v}_n), \text{diag}(\sigma_{\mathbf{z}}^v(\mathbf{v}_n))\right)$$

$$q_{\phi}(\mathbf{z}_n | \mathbf{v}_n, \mathbf{s}_n) = \mathcal{N}\left(\mu_{\mathbf{z}}^{av}(\mathbf{v}_n, \mathbf{s}_n), \text{diag}(\sigma_{\mathbf{z}}^{av}(\mathbf{v}_n, \mathbf{s}_n))\right)$$

<sup>4</sup>Sadeghi, M., et. al., (2020), IEEE TASLP.

## Learning AV Conditional VAE

On top of adding a video-only prior, we optimize a modified ELBO:

$$\begin{aligned}\mathcal{L}(\theta, \phi) = & \sum_n (1 - \alpha) \left( \mathbb{E}_{q_\phi(z_n|s_n, v_n)} \left[ \ln p_\theta(s_n|z_n, v_n) \right] - D_{\text{KL}} \left( q_\phi(z_n|s_n, v_n) \parallel p_\theta(z_n|v_n) \right) \right) \\ & + \alpha \mathbb{E}_{p_\theta(z_n|v_n)} \left[ \ln p_\theta(s_n|z_n, v_n) \right]\end{aligned}$$

$\Rightarrow \alpha > 0$  gives some reconstruction power to the visual prior!

This provides the parameters of the clean speech model  $\theta$ . Let's enhance!

# AV Conditional VAE for Speech Enhancement

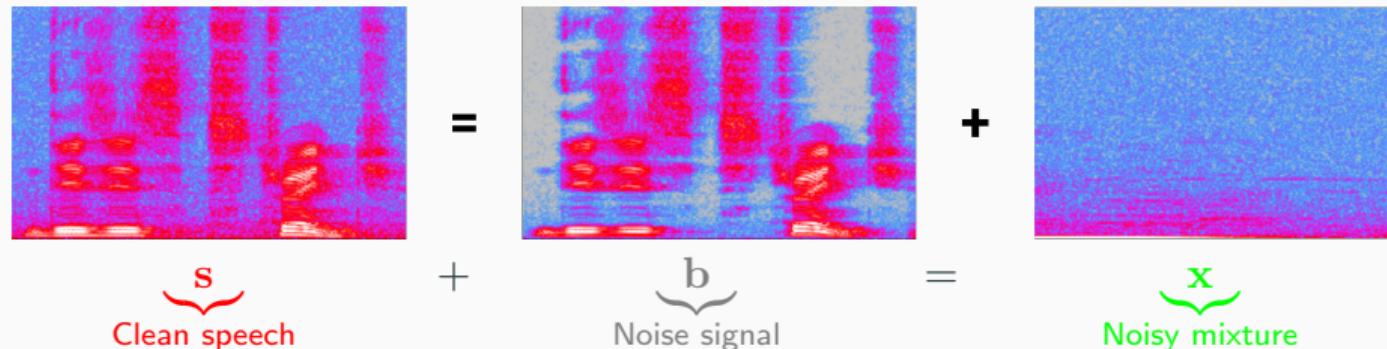


**Train** – Model for **clean data**:  $\{s_i\}_{i=1}^N$ .

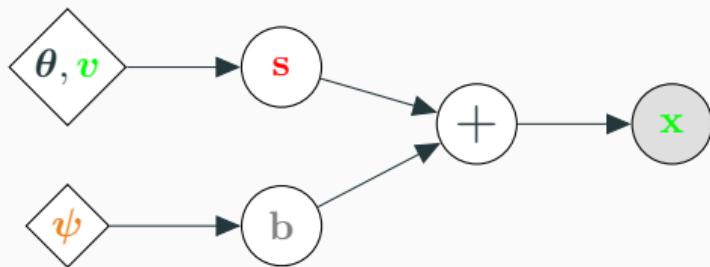


Then freeze  $\theta$  at test/adaptation time.

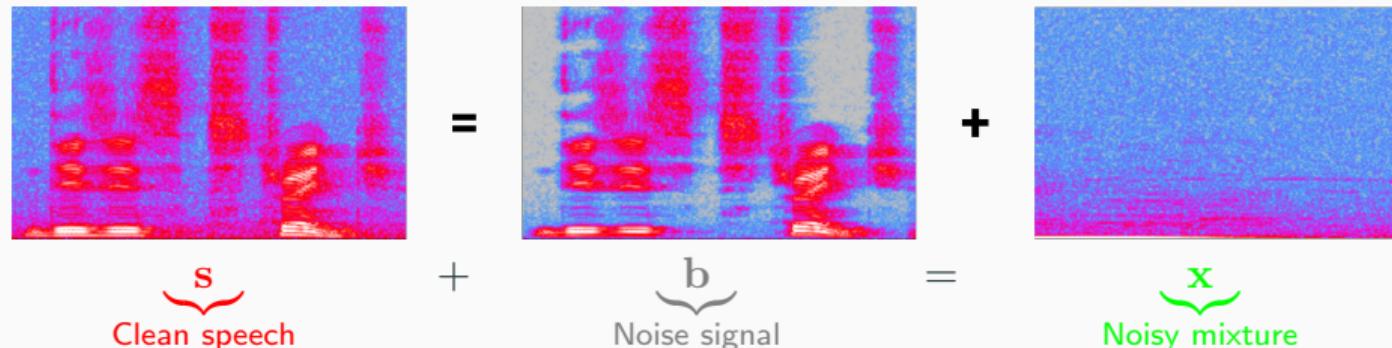
# AV Conditional VAE for Speech Enhancement



**Test** – learn the **noise parameters** from **noisy samples**  $x$  and estimate the **clean speech**  $\hat{s}$ :



# AV Conditional VAE for Speech Enhancement



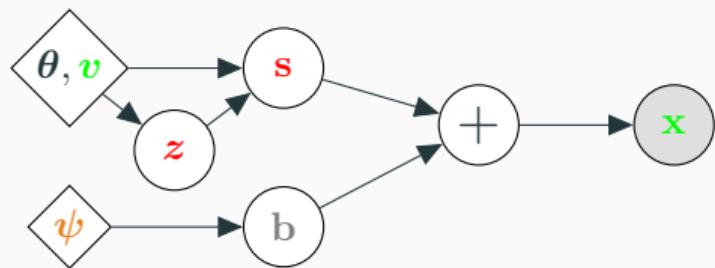
**Speech Enhancement:**

(Sample  $z$  from  $p(z|x_n, v_n)$ )

$$\hat{s}_n = \left( \frac{1}{R} \sum_{r=1}^R \frac{\sigma_s(z^{(r)}, v_n)}{\sigma_s(z^{(r)}, v_n) + \psi_n^*} \right) x_n.$$

Wiener-like filter!

**Test** – learn the **noise parameters** from **noisy samples  $x$**  and estimate the **clean speech  $\hat{s}$** :



## But how do we learn $\psi^*$ ?

We opt for an EM-like solution. Given an initial value of the parameters  $\psi^0$ :

$$Q(\psi, \psi^0) = \sum_n \mathbb{E}_{p(s_n, z_n | \mathbf{x}_n, \mathbf{v}_n; \psi^0, \theta)} [\log p(\mathbf{x}_n, s_n, z_n | \mathbf{v}_n; \psi, \theta)].$$

## But how do we learn $\psi^*$ ?

We opt for an EM-like solution. Given an initial value of the parameters  $\psi^0$ :

$$Q(\psi, \psi^0) = \sum_n \mathbb{E}_{p(s_n, z_n | \mathbf{x}_n, \mathbf{v}_n; \psi^0, \theta)} [\log p(\mathbf{x}_n, s_n, z_n | \mathbf{v}_n; \psi, \theta)].$$

Stop! There is an easier way:  $p(\mathbf{x}_n, z_n | \mathbf{v}_n; \psi) = \int p(\mathbf{x}_n, s_n, z_n | \mathbf{v}_n; \psi) ds_n$

Approximate  $Q$  via Monte-Carlo sampling:

- **E-Step:**  $Q(\psi; \psi^0) \approx \sum_{n=1}^N \frac{1}{R} \sum_{r=1}^R \ln p(\mathbf{x}_n, z_n^{(r)} | \mathbf{v}_n; \psi)$

The samples  $\left\{ z_n^{(r)} \right\}_{r=1, \dots, R}$  are i.i.d. and drawn from  $p(z_n | \mathbf{x}_n, \mathbf{v}_n; \psi^0)$ .

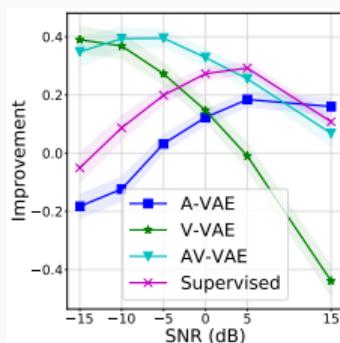
- **M-Step:**  $\psi^1 \leftarrow \operatorname{argmax}_{\psi} Q(\psi; \psi^0) \Rightarrow$  multiplicative update rules.

Standard formulae for NMF, not detailed here.

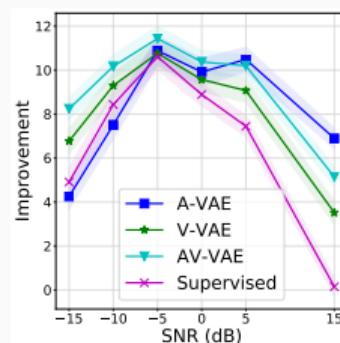
# Results

Performance measures (the higher, the better) – improvement w.r.t. the noisy mixture:

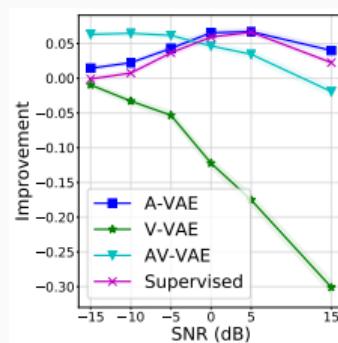
- Perceptual evaluation of speech quality (PESQ).
- Signal-to-distortion ratio (SDR).
- Short-time objective intelligibility (STOI).



(a) PESQ



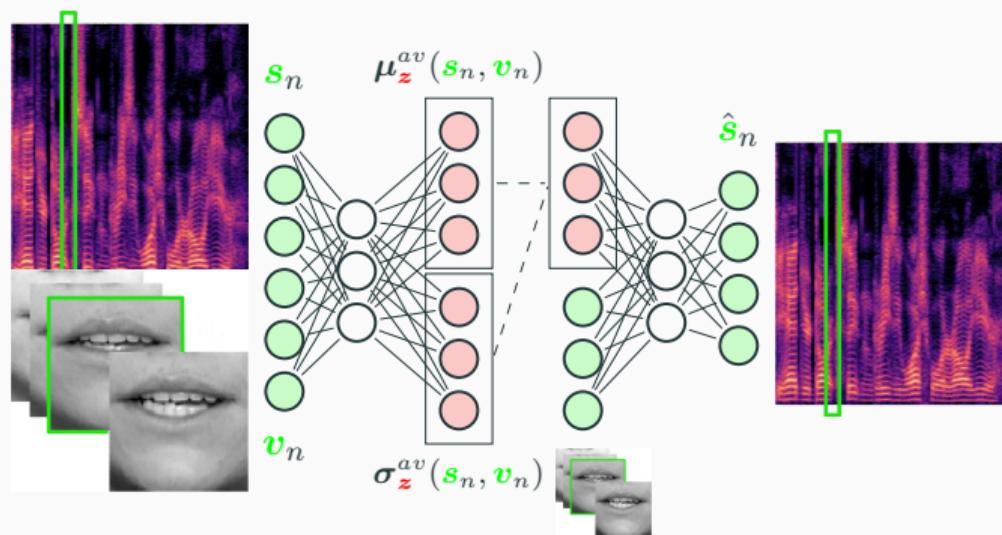
(b) SDR



(c) STOI

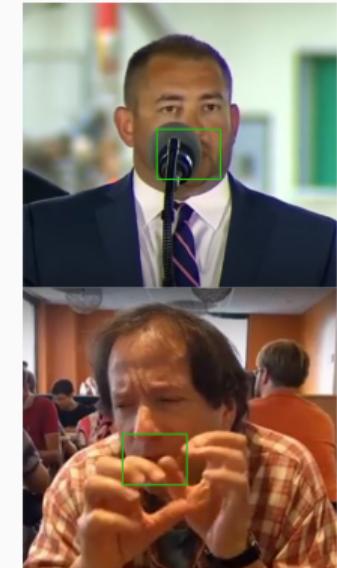
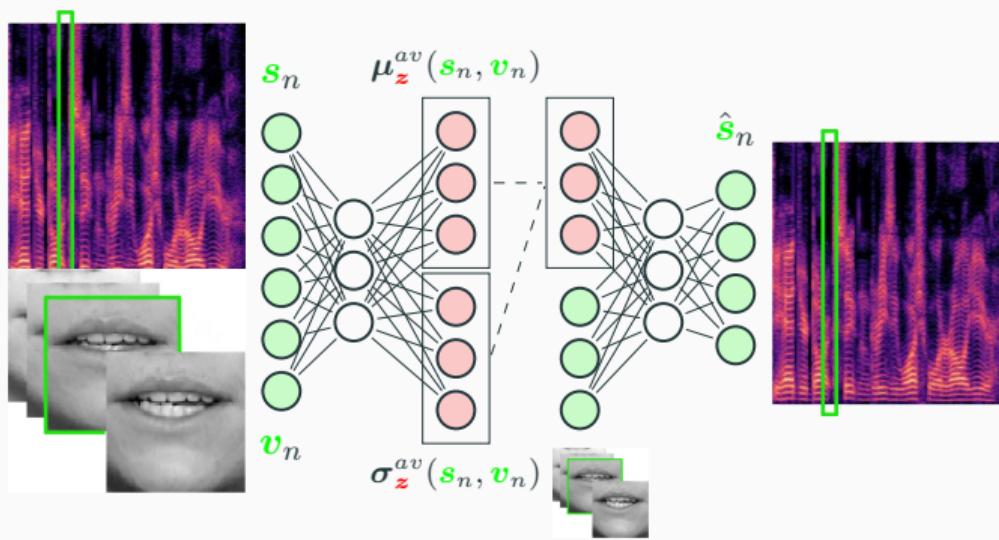
## Limitations: systematic AV fusion

From the model design:



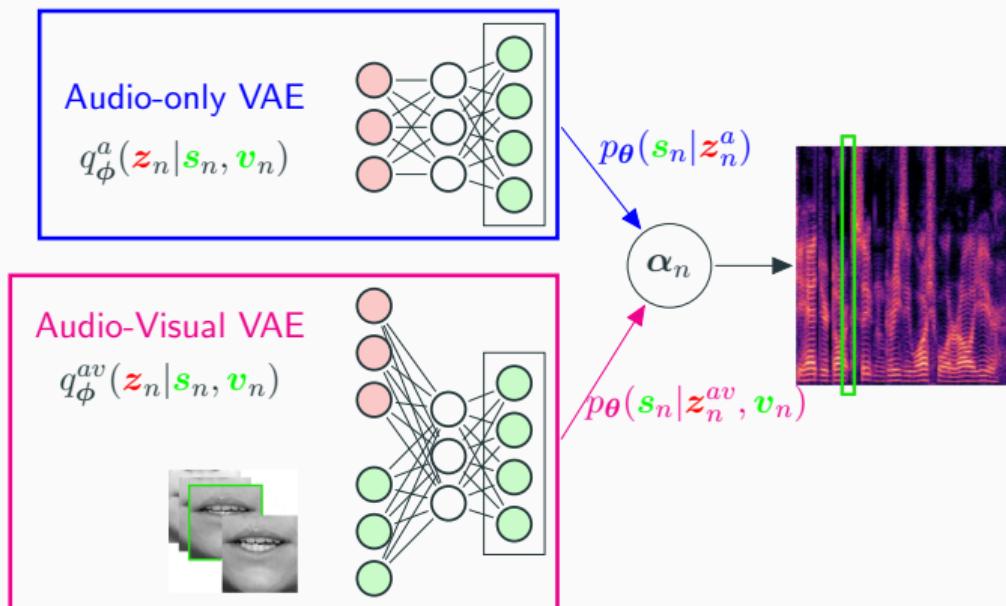
## Limitations: systematic AV fusion

From the model design:



• What about clutter?

# VAE Mixture Model<sup>5</sup>



$\alpha_n$  is the “mixing” latent variable.

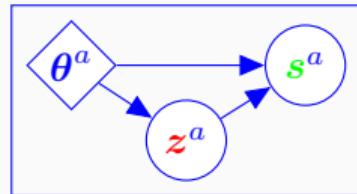
$$\begin{cases} \alpha_n = 1 \leftrightarrow \text{Audio-only} \\ \alpha_n = 0 \leftrightarrow \text{Audio-Visual} \end{cases}$$

Prior:  $p(\alpha_n) = \pi^{\alpha_n} (1 - \pi)^{1 - \alpha_n}$ .

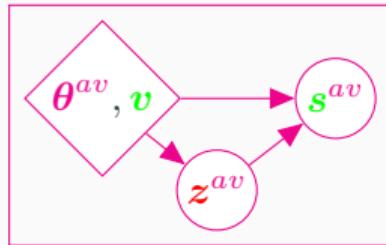
The **Audio-only** and the **Audio-Visual** VAE are **mixed without supervision**.

<sup>5</sup>Sadeghi, M., et. al., (2021), IEEE ICASSP.

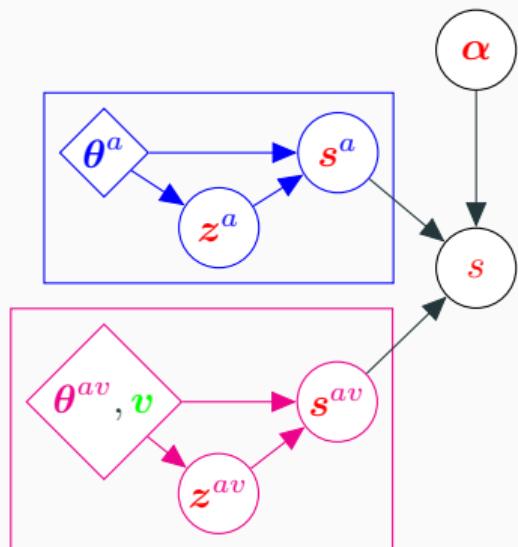
## Train, adaptation & speech enhancement



- Train **A-VAE** (with  $s$ ) and **AV-VAE** (with  $(s, v)$ ).

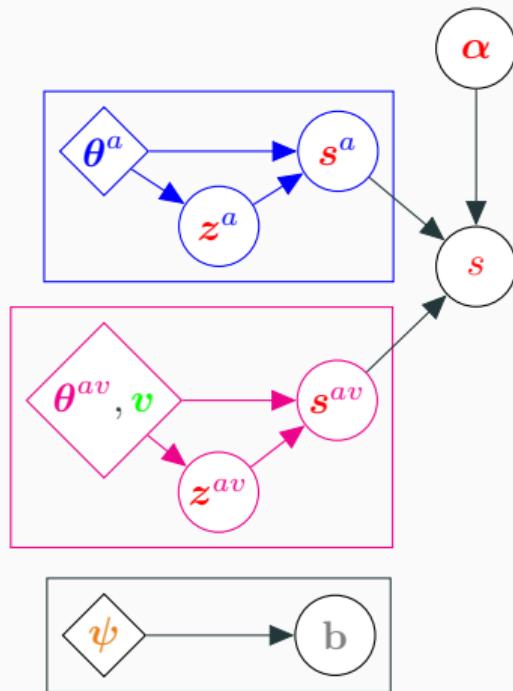


## Train, adaptation & speech enhancement



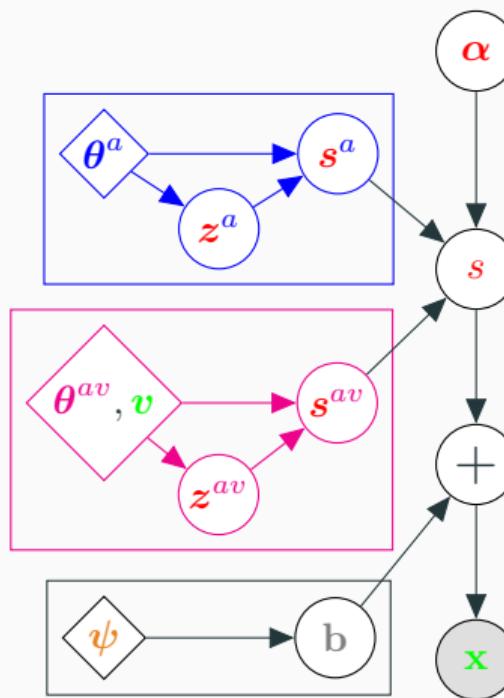
- Train **A-VAE** (with  $s$ ) and **AV-VAE** (with  $(s, v)$ ).
- The clean speech is now **latent** ( $s^a, s^{av}$ ), add the **mixing latent** variable ( $\alpha$ ).

## Train, adaptation & speech enhancement



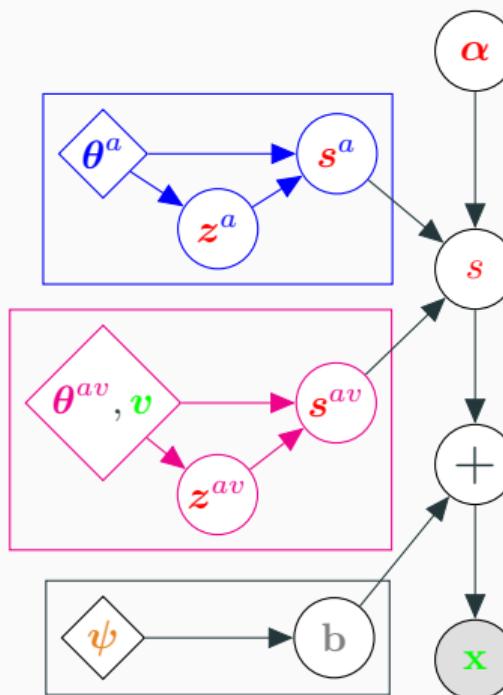
- Train A-VAE (with  $s$ ) and AV-VAE (with  $(s, v)$ ).
- The clean speech is now latent ( $s^a, s^{av}$ ), add the mixing latent variable ( $\alpha$ ).
- We also consider the noise variable and params.

## Train, adaptation & speech enhancement



- Train A-VAE (with  $s$ ) and AV-VAE (with  $(s, v)$ ).
- The clean speech is now **latent** ( $s^a, s^{av}$ ), add the **mixing latent** variable ( $\alpha$ ).
- We also consider the noise variable and params.
- Learn  $\psi^*$  and estimate the clean speech. By defining  $\gamma_n(z_n, v_n) = (\pi_n(\sigma_s^a(z_n))^{-1} + (1 - \pi_n)(\sigma_s^{av}(z_n, v_n))^{-1})^{-1}$

# Train, adaptation & speech enhancement



- Train A-VAE (with  $s$ ) and AV-VAE (with  $(s, v)$ ).
- The clean speech is now **latent** ( $s^a, s^{av}$ ), add the **mixing latent** variable ( $\alpha$ ).
- We also consider the noise variable and params.
- Learn  $\psi^*$  and estimate the clean speech. By defining  $\gamma_n(z_n, v_n) = (\pi_n(\sigma_s^a(z_n))^{-1} + (1 - \pi_n)(\sigma_s^{av}(z_n, v_n))^{-1})^{-1}$ :

$$\hat{s}_n = \frac{1}{R} \sum_{r=1}^R \frac{\gamma_n(z_n^{(r)}, v_n)}{\gamma_n(z_n^{(r)}, v_n) + \psi_n^*} x_n$$

## Learning the noise parameters with VAE-MM

Ideally:

$$Q(\psi, \psi^0) = \sum_{n=1}^N \mathbb{E}_{p(s_n, z_n, \alpha_n | \mathbf{x}_n, \mathbf{v}_n; \psi^0)} [\log p(\mathbf{x}_n, s_n, z_n, \alpha_n | \mathbf{v}_n; \psi)]$$

We need to approximate this posterior:

$$p(s_n, z_n, \alpha_n | \mathbf{x}_n, \mathbf{v}_n; \psi^0) \approx q(s_n, z_n, \alpha_n) = q_s(s_n) q_z(z_n) q_\alpha(\alpha_n).$$

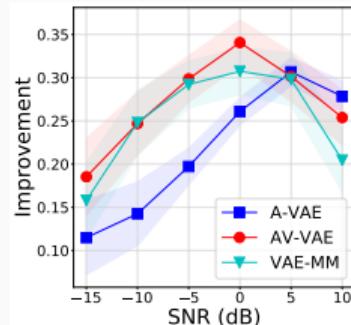
Of course we will need to sample from  $q_z$ , but in that case:

- ①  $q_\alpha$  is a Bernoulli distribution.
- ②  $q_s$  is a complex Gaussian distribution.

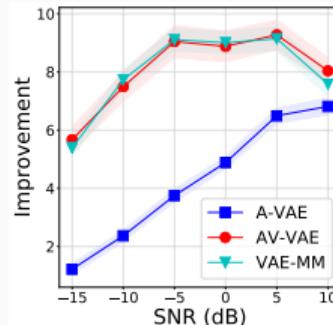
## Settings & Results

- **Noisy+clean speech:** NTCD-TIMIT database [Abdelaziz, 2017]
- **VAE models:** Pre-trained A-VAE and AV-VAE [Sadeghi et al., 2019]
- **Setup:** Very similar than in the previous experiments. **Clean and noisy** lips region visual information ( $\sim$  one-third of total video frames per sample).

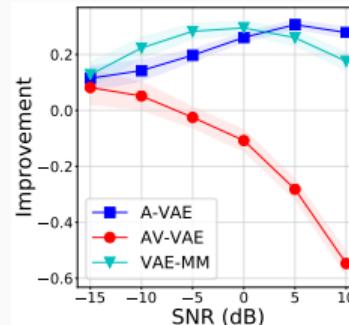
Improvement with respect to the input:



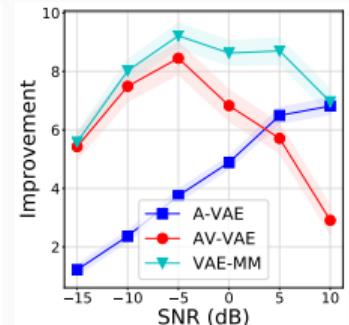
(a) PESQ (clean)



(b) SDR (clean)

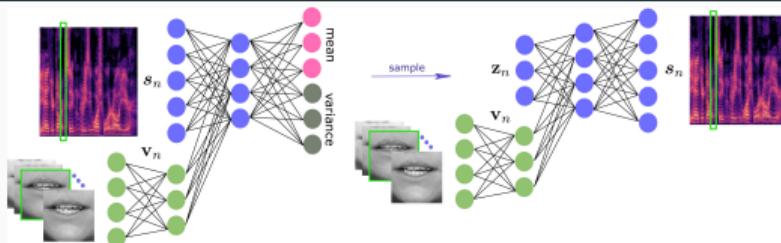


(c) PESQ (noisy)



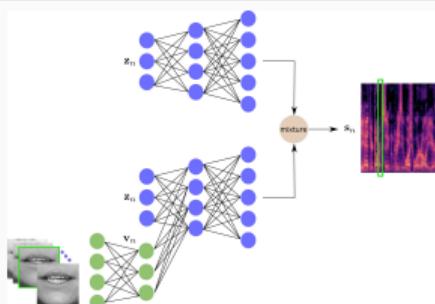
(d) SDR (noisy)

# Conditional VAE, VAE-MM and Mixture of Inference Networks VAE



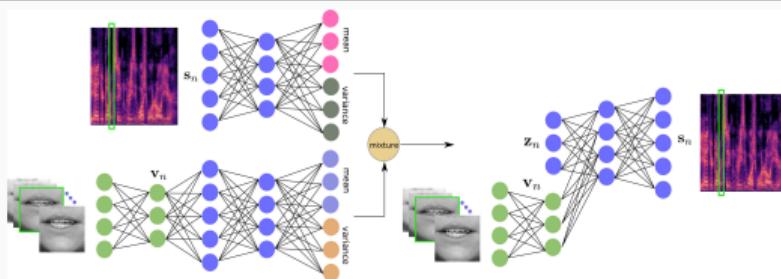
Conditional VAE:

- Training VAE via SGD.
- Systematic AV fusion.
- Enhancement via MECM.



VAE-MM:

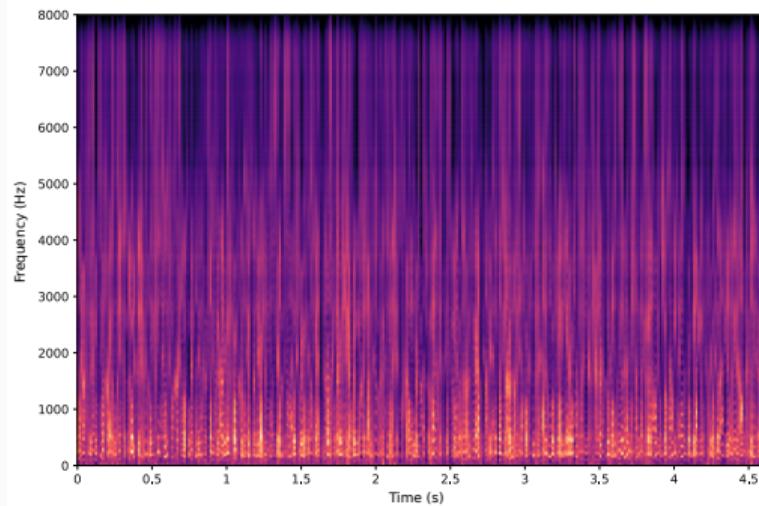
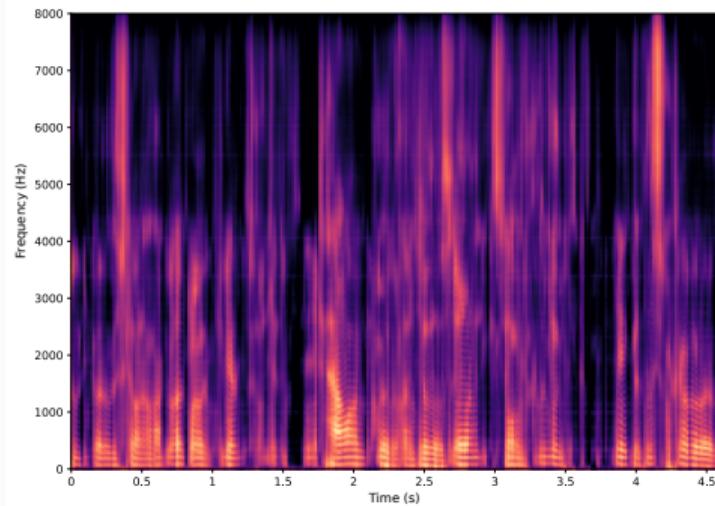
- Training two VAEs via SGD.
- Mixing (AV fusion) via two speech models.
- Enhancement via VEM + sampling.



MIN-VAE:

- Training all 3 networks via VEM+SGD.
- Mixing (AV fusion) via a single speech model.
- Enhancement via VEM + sampling.

## Strong limitation of VAEs: frame modeled independently.



Frames are modeled independently, we need time/sequential modeling!

- **Dynamical VAE (DVAE<sup>6</sup>) – “VAE for sequential modeling”**  
**(family of methods including existing literature)**



Xiaoyu Bie



Laurent Girin



Simon Leglaive



Thomas Hueber



Julien Diard

---

<sup>6</sup>Girin, L., et. al., (2021), FnT Machine Learning – Warning ~ 150 pages!!

## Probabilistic Sequential Modeling (decoder network)

We would like to model sequences of observations ( $\mathbf{x}_{1:T}$ ) and latent variables ( $\mathbf{z}_{1:T}$ ):

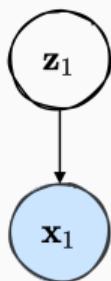
$$p_{\theta}^{\text{DVAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) \neq \prod_{t=1}^T p_{\theta}(\mathbf{x}_t, \mathbf{z}_t) = p_{\theta}^{\text{VAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}).$$



## Probabilistic Sequential Modeling (decoder network)

We would like to model sequences of observations ( $\mathbf{x}_{1:T}$ ) and latent variables ( $\mathbf{z}_{1:T}$ ):

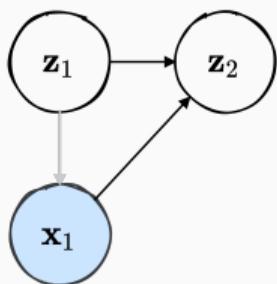
$$p_{\boldsymbol{\theta}}^{\text{DVAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) \neq \prod_{t=1}^T p_{\boldsymbol{\theta}}(\mathbf{x}_t, \mathbf{z}_t) = p_{\boldsymbol{\theta}}^{\text{VAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}).$$



## Probabilistic Sequential Modeling (decoder network)

We would like to model sequences of observations ( $\mathbf{x}_{1:T}$ ) and latent variables ( $\mathbf{z}_{1:T}$ ):

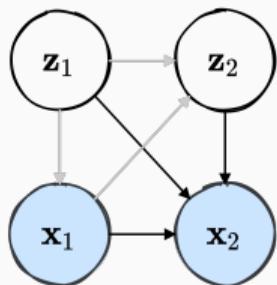
$$p_{\theta}^{\text{DVAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) \neq \prod_{t=1}^T p_{\theta}(\mathbf{x}_t, \mathbf{z}_t) = p_{\theta}^{\text{VAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}).$$



## Probabilistic Sequential Modeling (decoder network)

We would like to model sequences of observations ( $\mathbf{x}_{1:T}$ ) and latent variables ( $\mathbf{z}_{1:T}$ ):

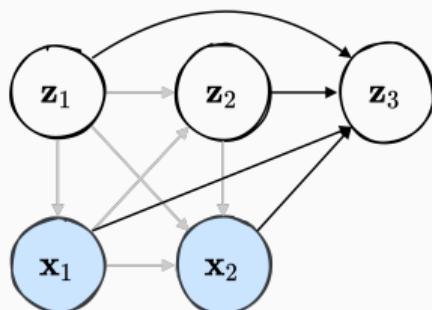
$$p_{\boldsymbol{\theta}}^{\text{DVAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) \neq \prod_{t=1}^T p_{\boldsymbol{\theta}}(\mathbf{x}_t, \mathbf{z}_t) = p_{\boldsymbol{\theta}}^{\text{VAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}).$$



## Probabilistic Sequential Modeling (decoder network)

We would like to model sequences of observations ( $\mathbf{x}_{1:T}$ ) and latent variables ( $\mathbf{z}_{1:T}$ ):

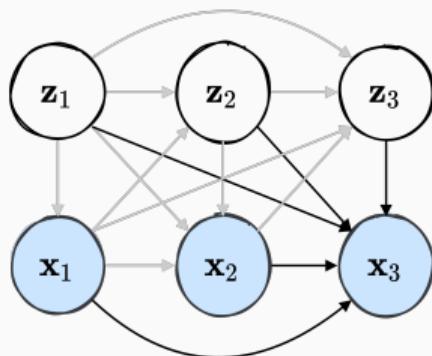
$$p_{\boldsymbol{\theta}}^{\text{DVAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) \neq \prod_{t=1}^T p_{\boldsymbol{\theta}}(\mathbf{x}_t, \mathbf{z}_t) = p_{\boldsymbol{\theta}}^{\text{VAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}).$$



## Probabilistic Sequential Modeling (decoder network)

We would like to model sequences of observations ( $\mathbf{x}_{1:T}$ ) and latent variables ( $\mathbf{z}_{1:T}$ ):

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

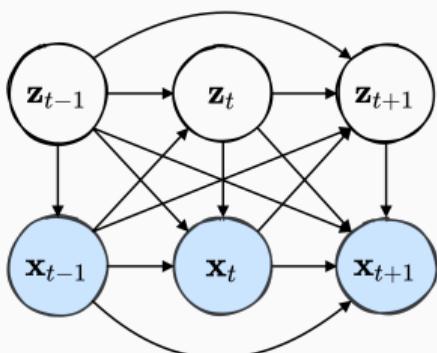


## Probabilistic Sequential Modeling (decoder network)

We would like to model sequences of observations ( $\mathbf{x}_{1:T}$ ) and latent variables ( $\mathbf{z}_{1:T}$ ):

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

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



The *prior* distribution  $p_{\theta}(z_{t+1} | \mathbf{z}_{1:t}, \mathbf{x}_{1:t})$  is now conditional, parametric, and might be auto-regressive (AR).

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

## Probabilistic Sequential Inference (encoder network)

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

$$p_{\theta}(z_{1:T} | x_{1:T}) = q_{\phi}(z_{1:T} | x_{1:T})$$

---

<sup>7</sup>Bishop, C. M., (2006).

## Probabilistic Sequential Inference (encoder 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. Use **D-separation**<sup>7</sup> to find out the true dependencies, then choose whether to keep them (e.g. to ensure causality).

---

<sup>7</sup>Bishop, C. M., (2006).

## Implementation (both decoder and encoder)

Let's take the general DVAE decoder:

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

---

<sup>8</sup>Several DVAEs @ <https://github.com/XiaoyuBIE1994/DVAE-speech>

## Implementation (both decoder and encoder)

Let's take the general DVAE decoder:

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

The generative distributions can be implemented with an RNN:

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

There are **many possible** implementations<sup>8</sup> for the same probabilistic dependencies!!!

---

<sup>8</sup>Several DVAEs @ <https://github.com/XiaoyuBIE1994/DVAE-speech>

## Learning: ELBO

The objective is build in the same way as VAEs, but looks different:

$$\begin{aligned}\mathcal{L}(\mathbf{x}_{1:T}; \phi, \theta) &\stackrel{?}{=} \sum_{t=1}^T \mathbb{E}_{q_\phi(z_t | \mathbf{x}_{1:T})} \underbrace{\ln p_{\theta_x}(\mathbf{x}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_{1:t})}_{\text{Reconstruction}} \\ &\quad - \sum_{t=1}^T \underbrace{D_{\text{KL}}\left(q_\phi(z_t | \mathbf{z}_{1:t-1}, \mathbf{x}_{1:T}) \parallel p_{\theta_z}(z_t | \mathbf{x}_{1:t-1}, \mathbf{z}_{1:t-1})\right)}_{\text{Regularization}}\end{aligned}$$

- The reconstruction and regularisation terms are evaluated at every frame  $t$ .

## Learning: ELBO

The objective is build in the same way as VAEs, but looks different:

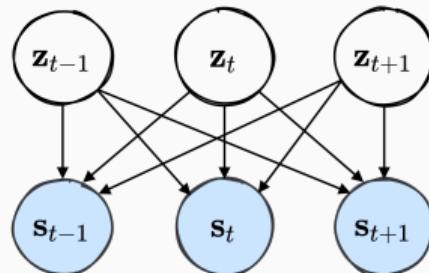
$$\begin{aligned}\mathcal{L}(\mathbf{x}_{1:T}; \boldsymbol{\phi}, \boldsymbol{\theta}) = & \sum_{t=1}^T \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{1:t} | \mathbf{x}_{1:T})} \underbrace{[\ln p_{\boldsymbol{\theta}_x}(\mathbf{x}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_{1:t})]}_{\text{Reconstruction}} \\ & - \sum_{t=1}^T \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{1:t-1} | \mathbf{x}_{1:T})} \left[ \underbrace{D_{\text{KL}}\left(q_{\boldsymbol{\phi}}(\mathbf{z}_t | \mathbf{z}_{1:t-1}, \mathbf{x}_{1:T}) \parallel p_{\boldsymbol{\theta}_z}(\mathbf{z}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_{1:t-1})\right)}_{\text{Regularization}} \right]\end{aligned}$$

- 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 parallelized!

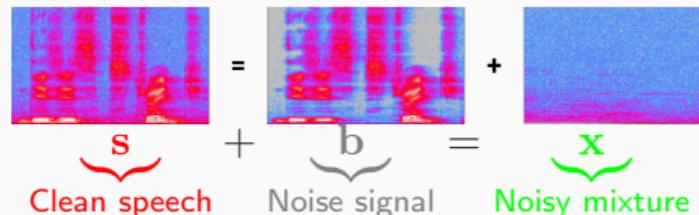
## Application to Unsupervised Probabilistic SE (revisit)



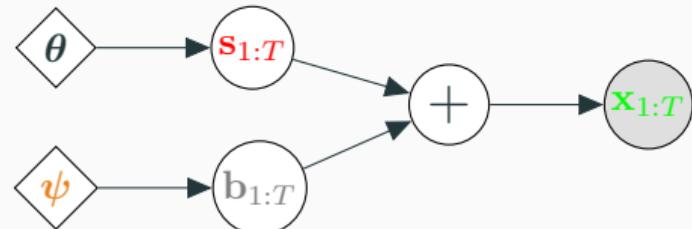
Train – Model  $\theta$  for clean data:  $\{s_{1:T}\}$ .



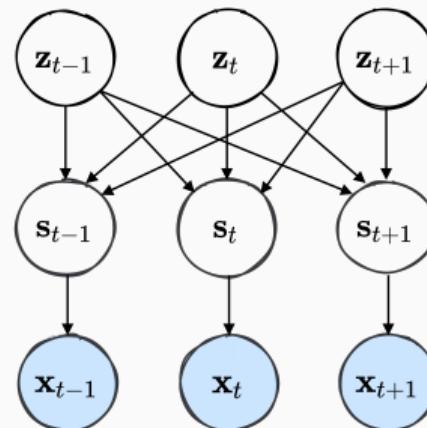
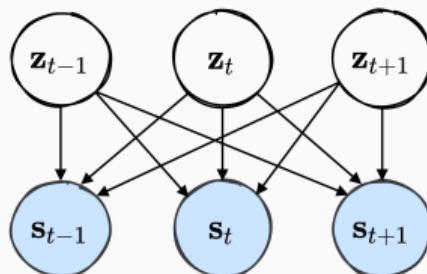
# Application to Unsupervised Probabilistic SE (revisit)



Test – learn the noise parameters from noisy samples  $\mathbf{x}$  and estimate the clean speech  $\hat{\mathbf{s}}_{1:T}$ :



Train – Model  $\theta$  for clean data:  $\{s_{1:T}\}$ .



## Results on the Wall Street Journal (WSJ) and VoiceBank (VB)<sup>9</sup>

Method	Superv.	Test	Train	SI-SDR (dB)	Train	SI-SDR (dB)
Noisy mixture	-	-	-	-2.6	-	-2.6
VAE-VEM (ICASSP'20)	None	WSJ+QUT	Same	5.0	N/A	<b>5.8</b>
RVAE-VEM (Proposed)	None			<b>5.8</b>		
MetricGAN-U (ICASSP'22)	Partial	WSJ+QUT	Same	N/A	<u>5.7</u>	3.6
UMX* 2020	Full			<u>5.7</u>		
MetricGAN+* (Interspeech'21)	Full			3.6		
Noisy mixture	-	WSJ+QUT	Same	N/A	<b>5.8</b>	3.6
VAE-VEM (ICASSP'20)	None					
RVAE-VEM (Proposed)	None					
NyTT EUSIPCO'21	Partial/Xtra					
MetricGAN-U (half) ICASSP'22	Partial					
UMX 2020	Full					
MetricGAN+ Interspeech'21	Full					

Results in dB: **best** and second best per section.

<sup>9</sup>Bie, X., et. al., (2022), IEEE TASLP.

# Results on the Wall Street Journal (WSJ) and VoiceBank (VB)<sup>9</sup>

Method	Superv.	Test	Train	SI-SDR (dB)	Train	SI-SDR (dB)
Noisy mixture	-		-	-2.6	-	-2.6
VAE-VEM (ICASSP'20)	None	WSJ+QUT		5.0		
RVAE-VEM (Proposed)	None			<b>5.8</b>		
MetricGAN-U (ICASSP'22)	Partial	Same		N/A		
UMX* 2020	Full			<u>5.7</u>		
MetricGAN+* (Interspeech'21)	Full			3.6		
Noisy mixture	-		-	8.4	-	8.4
VAE-VEM (ICASSP'20)	None	VB-DMD		16.4		
RVAE-VEM (Proposed)	None			<u>17.1</u>		
NyTT EUSIPCO'21	Partial/Xtra	Same		<b>17.7</b>		
MetricGAN-U (half) ICASSP'22	Partial			8.2		
UMX 2020	Full			14.0		
MetricGAN+ Interspeech'21	Full			8.5		

Results in dB: **best** and second best per section.

<sup>9</sup>Bie, X., et. al., (2022), IEEE TASLP.

# Results on the Wall Street Journal (WSJ) and VoiceBank (VB)<sup>9</sup>

Method	Superv.	Test	Train	SI-SDR (dB)	Train	SI-SDR (dB)
Noisy mixture	-		-	-2.6	-	-2.6
VAE-VEM (ICASSP'20)	None	WSJ+QUT		5.0		3.8
RVAE-VEM (Proposed)	None			<b>5.8</b>		<b>4.3</b>
MetricGAN-U (ICASSP'22)	Partial	Same		N/A		-1.6
UMX* 2020	Full			<u>5.7</u>	Different	<u>4.1</u>
MetricGAN+* (Interspeech'21)	Full			3.6		1.8
Noisy mixture	-		-	8.4	-	8.4
VAE-VEM (ICASSP'20)	None	VB-DMD		16.4		
RVAE-VEM (Proposed)	None			<u>17.1</u>		
NyTT EUSIPCO'21	Partial/Xtra			<b>17.7</b>		
MetricGAN-U (half) ICASSP'22	Partial	Same		8.2		
UMX 2020	Full			14.0		
MetricGAN+ Interspeech'21	Full			8.5		

Results in dB: **best** and second best per section.

<sup>9</sup>Bie, X., et. al., (2022), IEEE TASLP.

## Results on the Wall Street Journal (WSJ) and VoiceBank (VB)<sup>9</sup>

Method	Superv.	Test	Train	SI-SDR (dB)	Train	SI-SDR (dB)
Noisy mixture	-		-	-2.6	-	-2.6
VAE-VEM (ICASSP'20)	None	WSJ+QUT	Same	5.0	Different	3.8
RVAE-VEM (Proposed)	None			<b>5.8</b>		<b>4.3</b>
MetricGAN-U (ICASSP'22)	Partial	VB-DMD	Same	N/A	Different	-1.6
UMX* 2020	Full			<u>5.7</u>		<u>4.1</u>
MetricGAN+* (Interspeech'21)	Full			3.6		1.8
Noisy mixture	-		-	8.4	-	8.4
VAE-VEM (ICASSP'20)	None	VB-DMD	Same	16.4	Different	<u>15.0</u>
RVAE-VEM (Proposed)	None			<u>17.1</u>		<b>17.3</b>
NyTT EUSIPCO'21	Partial/Xtra	VB-DMD	Same	<b>17.7</b>	Different	N/A
MetricGAN-U (half) ICASSP'22	Partial			8.2		N/A
UMX 2020	Full			14.0		10.4
MetricGAN+ Interspeech'21	Full			8.5		3.9

Results in dB: **best** and second best per section.

<sup>9</sup>Bie, X., et. al., (2022), IEEE TASLP.

DVAEs are good (great!) at temporal modeling!

All implementations use RNN (or variants).

► What about transformers?



Xiaoyu Bie



Wen Guo



Simon Leglaive



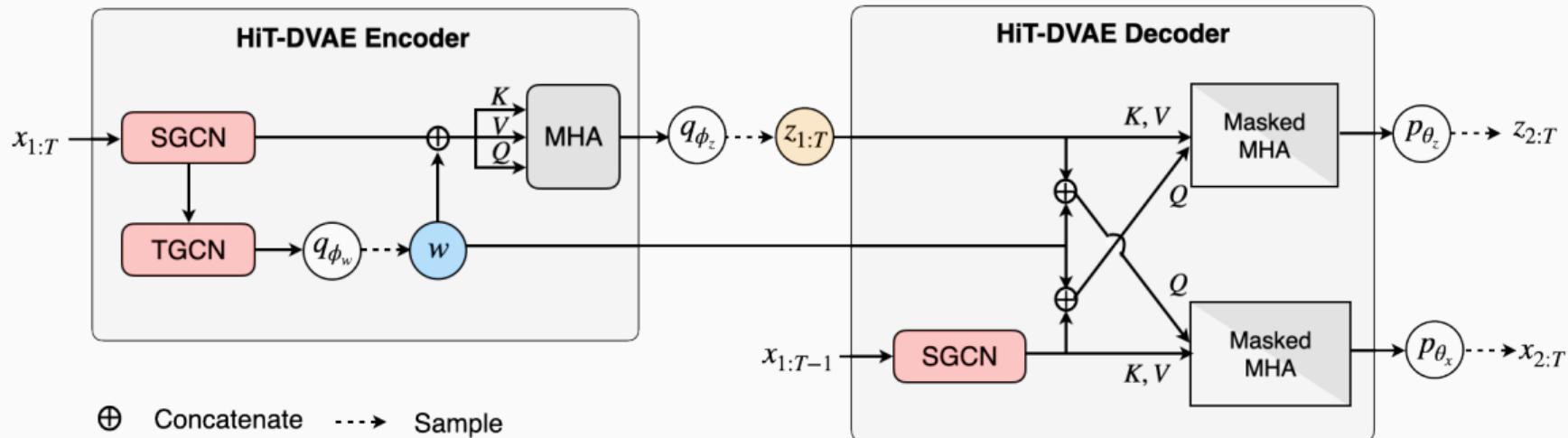
Francesc Moreno



Laurent Girin

# Introducing HiT-DVAE<sup>10</sup>

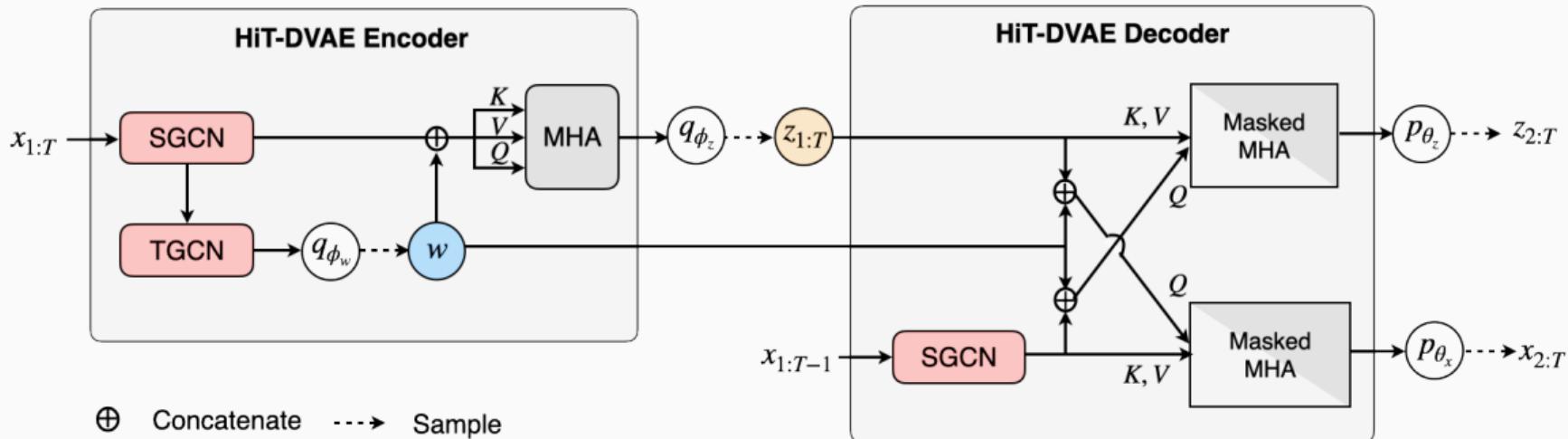
Hierarchical Transformer DVAE – two levels of latent variables: static  $w$  and dynamic  $z_{1:T}$ .



<sup>10</sup>Bie, X., et. al., (2023), Pre-print.

# Introducing HiT-DVAE<sup>10</sup>

Hierarchical Transformer DVAE – two levels of latent variables: static  $w$  and dynamic  $z_{1:T}$ .

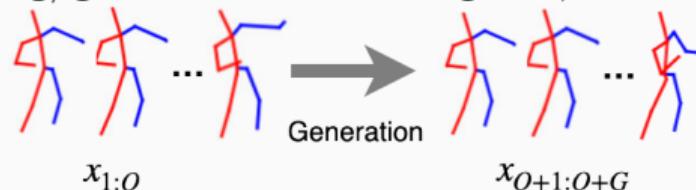


There is something **odd** in the way Q, K, V are used...

<sup>10</sup>Bie, X., et. al., (2023), Pre-print.

## HIT-DVAE: Modified transformer architecture

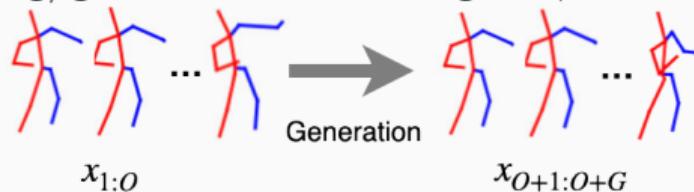
Task: human motion modeling/generation. Predicting  $x_{O+1}$  from  $x_{1:O}$  is very easy.<sup>11</sup>



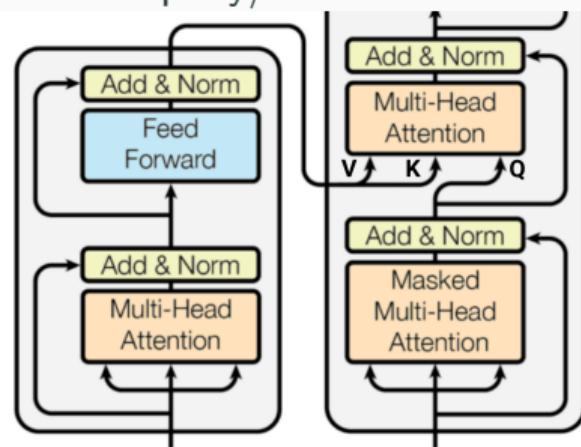
<sup>11</sup>Guo., W., et. al., (2021), IEEE WACV.

# HIT-DVAE: Modified transformer architecture

Task: human motion modeling/generation. Predicting  $x_{O+1}$  from  $x_{1:O}$  is very easy.<sup>11</sup>



Standard query/residual connection:



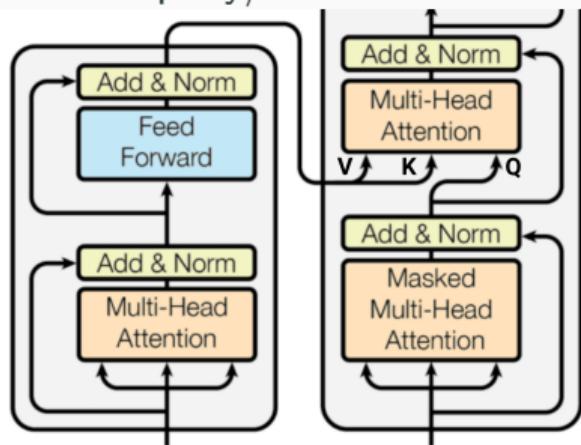
<sup>11</sup>Guo., W., et. al., (2021), IEEE WACV.

# HIT-DVAE: Modified transformer architecture

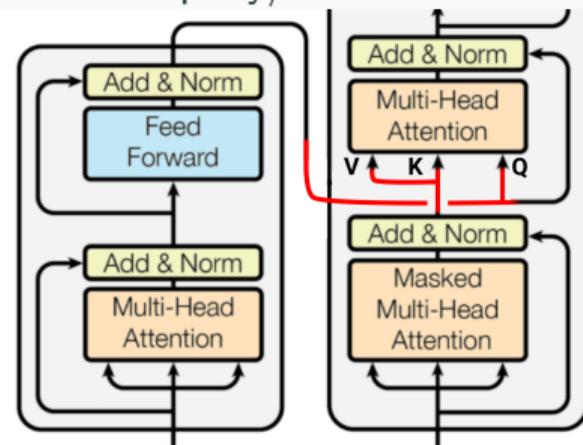
Task: human motion modeling/generation. Predicting  $x_{O+1}$  from  $x_{1:O}$  is very easy.<sup>11</sup>



Standard query/residual connection:



HIT-DVAE query/residual connection:



<sup>11</sup>Guo., W., et. al., (2021), IEEE WACV.

# HIT-DVAE: Results

	Start	GT	End pose of 10 sample		Start	GT	End pose of 10 sample
Ours							
gsps							
DLow							
HumanEva-I, Box						Human3.6m, Discussion	
Ours							
gsps							
DLow							
HumanEva-I, Walk						Human3.6m, Greeting	

Good trade-off between diversity, quality and accuracy. Also useful for audio modeling.<sup>12</sup>

<sup>12</sup>Lin, X., et. al., (2023), IEEE ICASSP.

**DVAEs can model mono-modal sequences.**

- ▶ What about multiple modalities?



Samir Sadok



Simon Leglaive

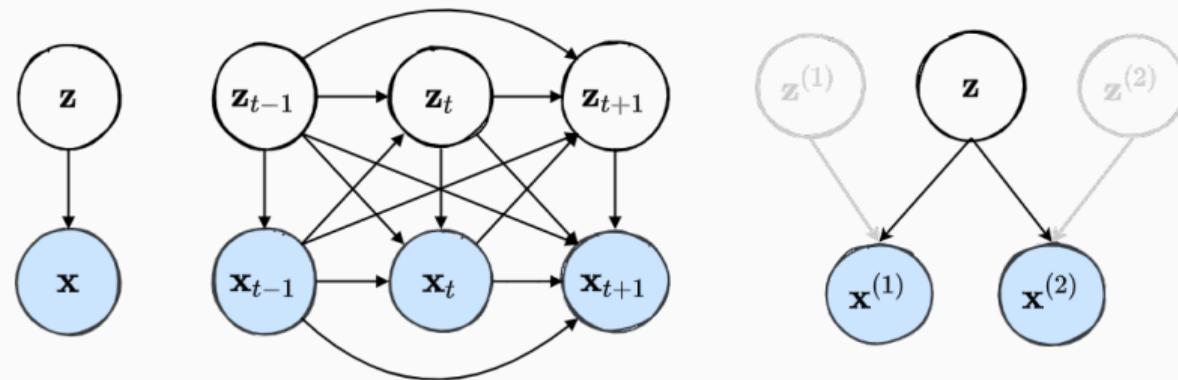


Renaud Séguier



Laurent Girin

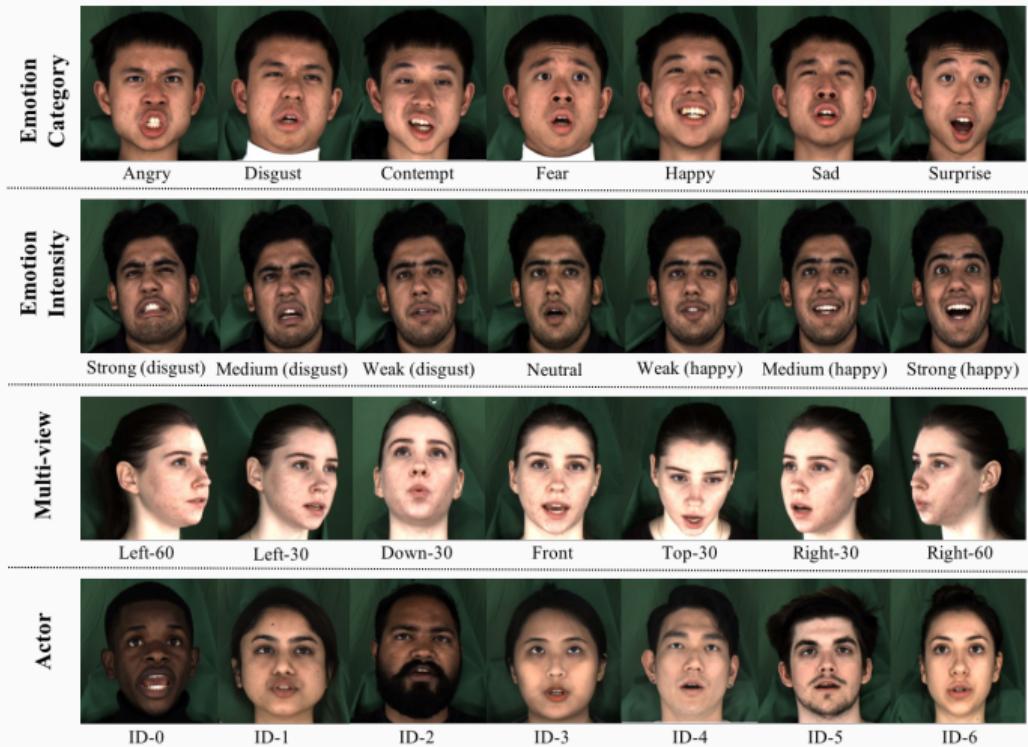
In other words:



VAEs (left) can be multi-modal<sup>13</sup> (right). Can DVAEs (middle) be multi-modal too?

<sup>13</sup>Sutter, T. M., et. al., (2021), ICLR.

# Task: emotional audio-visual speech modeling.<sup>14</sup>

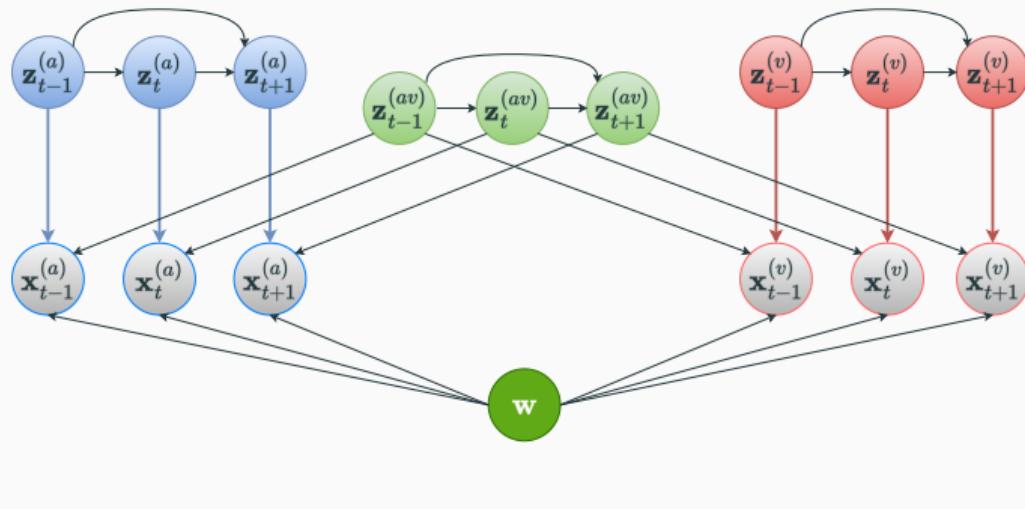


What should we model?

- Static AV  
(ID, emotion)
- Dynamic AV  
(lip-audio corr.)
- Dynamic A  
(other audio features)
- Dynamic V  
(eye AUs).

<sup>14</sup>Wang, K., et. al., (2020), ECCV.

# Introducing MDVAEs<sup>15</sup>

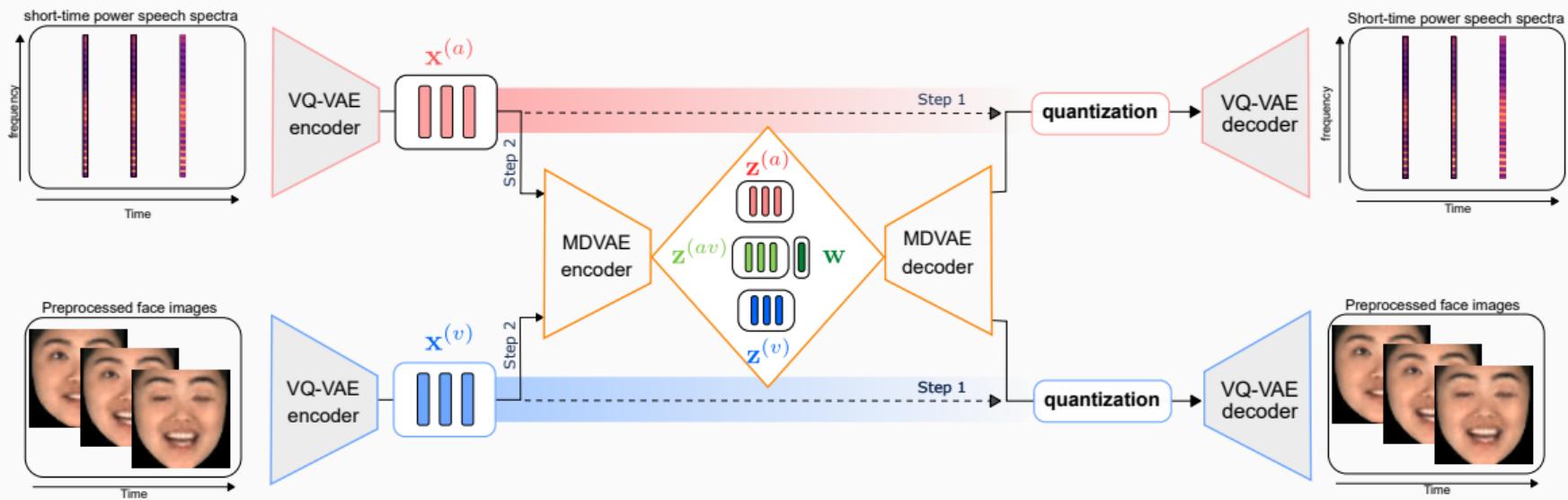


What should we model?

- Static AV –  $w$   
(ID, emotion)
- Dynamic AV –  $z^{av}$   
(lip-audio corr.)
- Dynamic A –  $z^a$   
(other audio features)
- Dynamic V –  $z^v$   
(eye AUs).

<sup>15</sup>Sadok, S., et. al., (2024), Neural Networks.

# The VQ-MDVAE Architecture



(i) Quantize auditory and visual features. (ii) Use MDVAE to model the quantized features.

## MDVAE: Some fun things

Let's see a couple of videos on:

- latent variable transfer,
- latent variable interpolation.

Earlier, we combined (static) VAEs with mixing latents.

- ▶ Is it possible/useful with DVAEs?

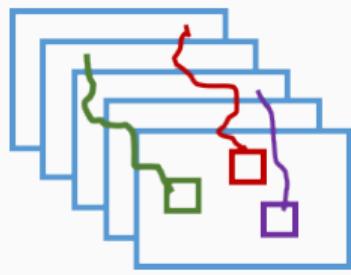


Xiaoyu Lin

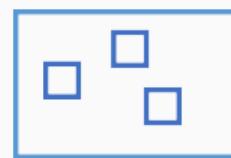


Laurent Girin

## Motivating application: unsupervised multiple object tracking



Tracklets @  $t - 1$



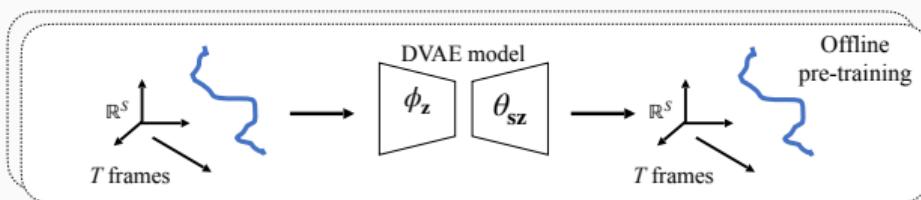
Detections @  $t$



Desired result

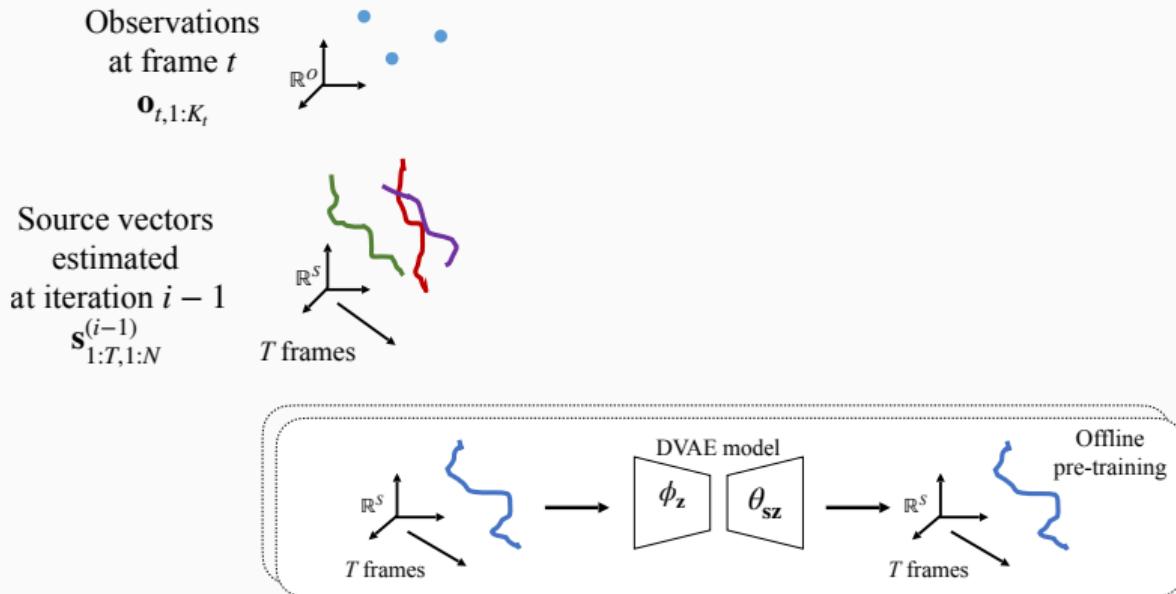
**Question:** Can we model the non-linear dynamics with a DVAE, and have an assignment mechanism within the same ML formulation?

# Introducing MixDVAE<sup>16</sup>



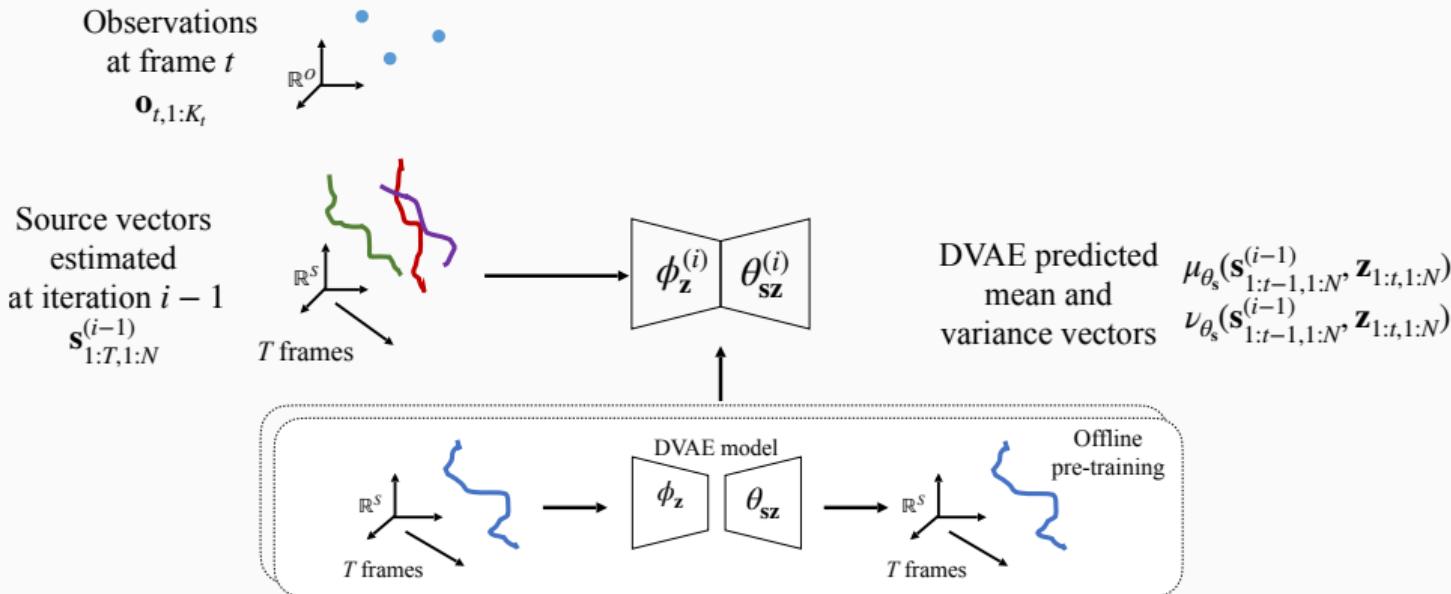
<sup>16</sup>Lin, X., et. al., (2023), Under review.

# Introducing MixDVAE<sup>16</sup>



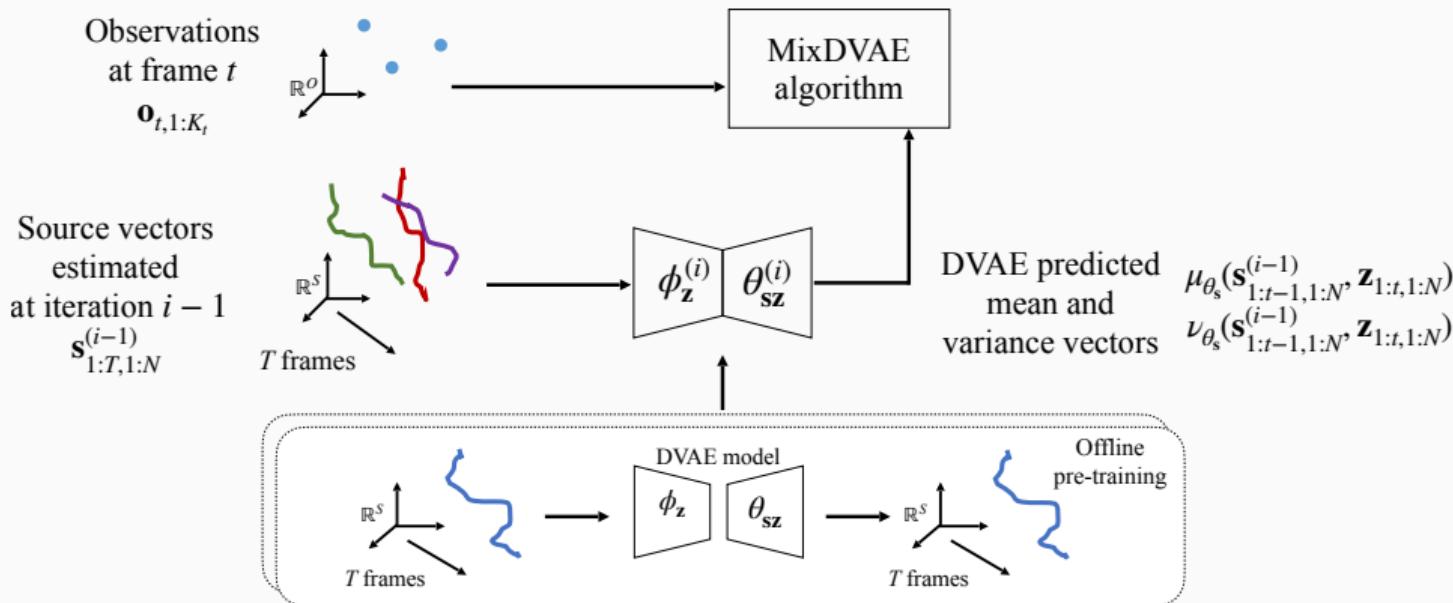
<sup>16</sup>Lin, X., et. al., (2023), Under review.

# Introducing MixDVAE<sup>16</sup>



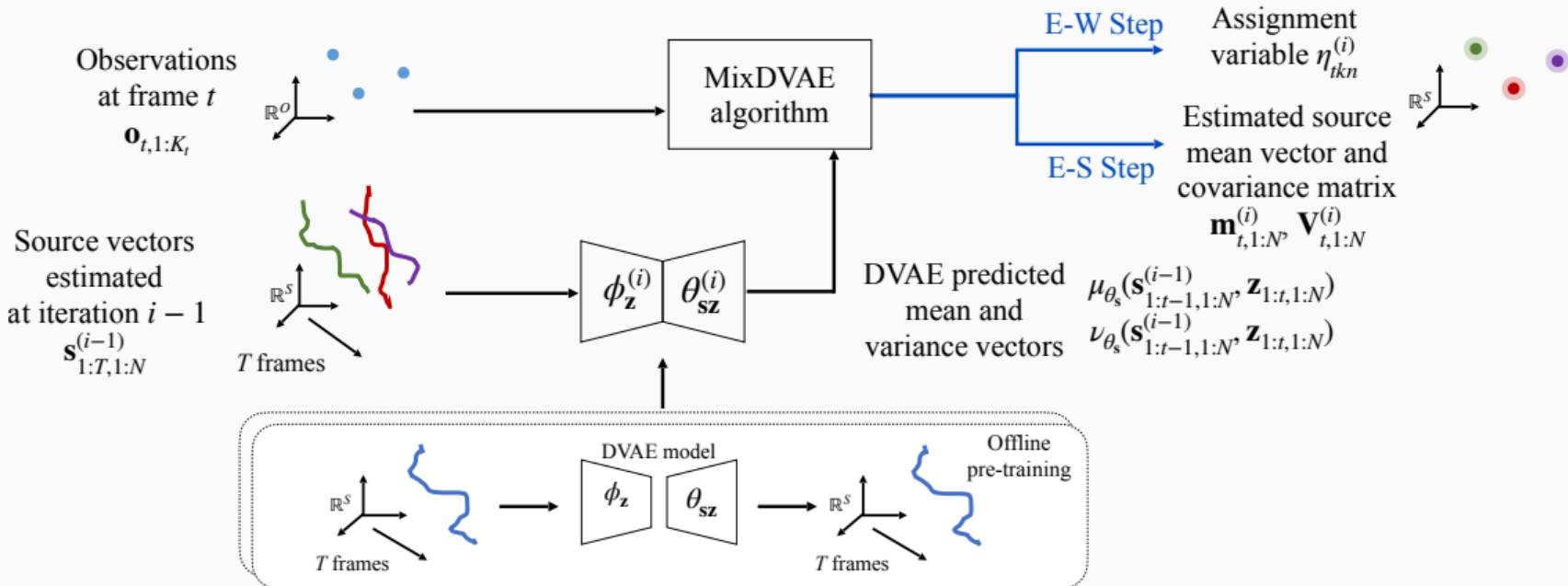
<sup>16</sup>Lin, X., et. al., (2023), Under review.

# Introducing MixDVAE<sup>16</sup>



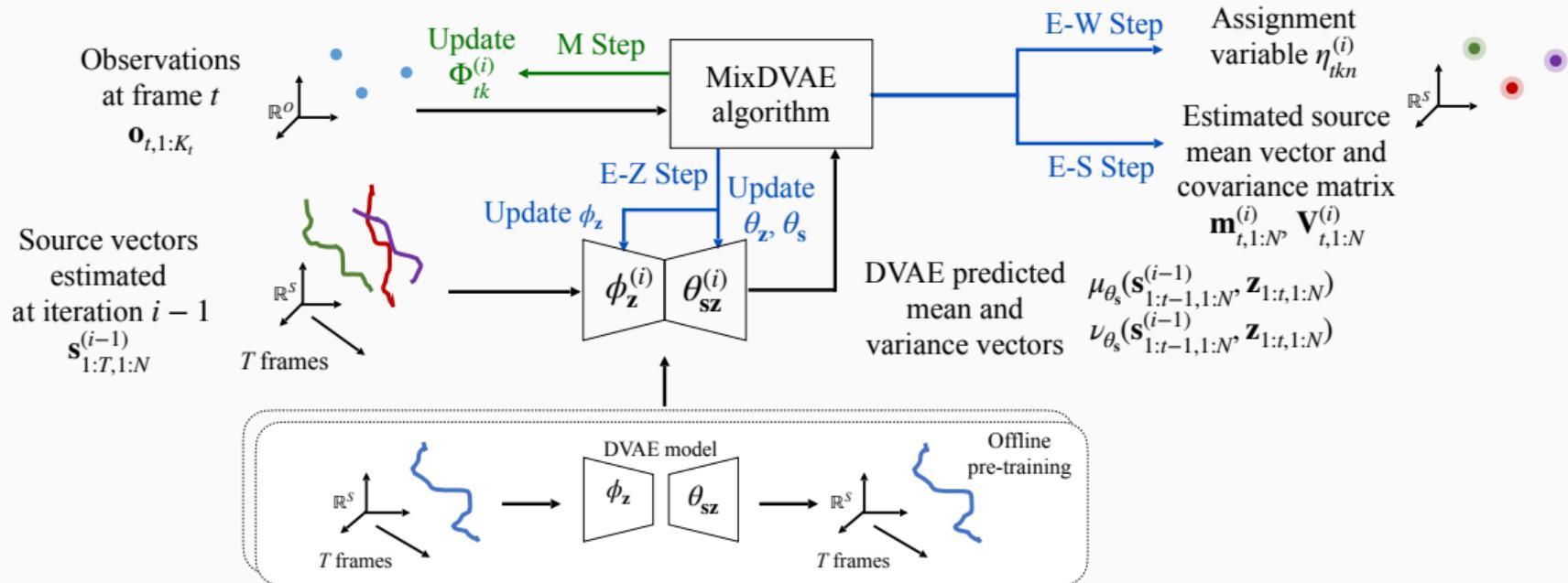
<sup>16</sup>Lin, X., et. al., (2023), Under review.

# Introducing MixDVAE<sup>16</sup>



<sup>16</sup>Lin, X., et. al., (2023), Under review.

# Introducing MixDVAE<sup>16</sup>



<sup>16</sup>Lin, X., et. al., (2023), Under review.

## Quick discussion & results

GMM update ( $K$  observations):

$$\boldsymbol{\mu}_n = \sum_k \underbrace{\eta_{kn}}_{\text{assign.}} \mathbf{o}_k$$

## Quick discussion & results

GMM update ( $K$  observations):

$$\boldsymbol{\mu}_n = \sum_k \underbrace{\eta_{kn}}_{\text{assign.}} \mathbf{o}_k$$

Kalman update (sequence of  $T$  obs):

$$\mathbf{s}_t = \underbrace{\mathbf{P}_t \mathbf{o}_t}_{\text{update}} + \underbrace{\mathbf{T}_t \mathbf{s}_{t-1}}_{\text{prediction}}$$

## Quick discussion & results

GMM update ( $K$  observations):

$$\boldsymbol{\mu}_n = \sum_k \underbrace{\eta_{kn}}_{\text{assign.}} \mathbf{o}_k$$

Kalman update (sequence of  $T$  obs):

$$\mathbf{s}_t = \underbrace{\mathbf{P}_t \mathbf{o}_t}_{\text{update}} + \underbrace{\mathbf{T}_t \mathbf{s}_{t-1}}_{\text{prediction}}$$

MixDVAE update is a combination:

$$\mathbf{s}_{tn} = \underbrace{\mathbf{P}_t \sum_k \eta_{kn} \mathbf{o}_{tk}}_{\text{assig. \& update}} + \underbrace{\mathbf{T}_t(\mathbf{s}_{1:t-1})}_{\text{non-lin. prediction}}$$

## Quick discussion & results

GMM update ( $K$  observations):

$$\boldsymbol{\mu}_n = \sum_k \underbrace{\eta_{kn}}_{\text{assign.}} \mathbf{o}_k$$

Kalman update (sequence of  $T$  obs):

$$\mathbf{s}_t = \underbrace{\mathbf{P}_t \mathbf{o}_t}_{\text{update}} + \underbrace{\mathbf{T}_t \mathbf{s}_{t-1}}_{\text{prediction}}$$

MixDVAE update is a combination:  $\mathbf{s}_{tn} = \underbrace{\mathbf{P}_t \sum_k \eta_{kn} \mathbf{o}_{tk}}_{\text{assig. \& update}} + \underbrace{\mathbf{T}_t(\mathbf{s}_{1:t-1})}_{\text{non-lin. prediction}}$

- Results in unsupervised MOT and in semi-blind source separation.
- Work in progress: fine-tuning, complexity, learning from noise, ...

# VAE/DVAE Summary

	Mono-modal	Multi-modal	Mixtures
Static	VAE [Kingma'14] VQ-VAE [van den Oord'17]		
Dynamic			

# VAE/DVAE Summary

	Mono-modal	Multi-modal	Mixtures
Static	VAE [Kingma'14]	CVAE [Sadeghi'20]	
	VQ-VAE [van den Oord'17]	MVAE [Sutter'21]	
Dynamic			

# VAE/DVAE Summary

	Mono-modal	Multi-modal	Mixtures
Static	VAE [Kingma'14]	CVAE [Sadeghi'20]	VAE-MM [Sadeghi'20]
	VQ-VAE [van den Oord'17]	MVAE [Sutter'21]	MIN-VAE [Sadeghi'21]
Dynamic			

# VAE/DVAE Summary

	Mono-modal	Multi-modal	Mixtures
Static	VAE [Kingma'14]	CVAE [Sadeghi'20]	VAE-MM [Sadeghi'20]
	VQ-VAE [van den Oord'17]	MVAE [Sutter'21]	MIN-VAE [Sadeghi'21]
Dynamic	DVAE [Girin'21]		
	Sw-VAE [Sadeghi'21]		

# VAE/DVAE Summary

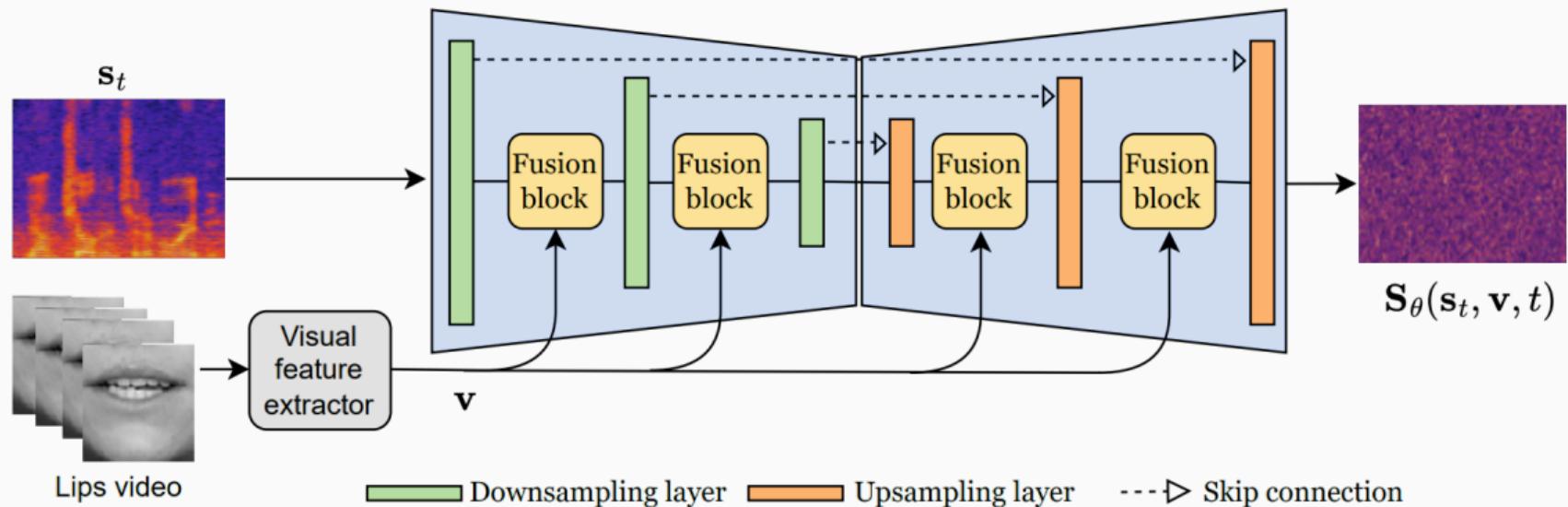
	Mono-modal	Multi-modal	Mixtures
Static	VAE [Kingma'14] VQ-VAE [van den Oord'17]	CVAE [Sadeghi'20] MVAE [Sutter'21]	VAE-MM [Sadeghi'20] MIN-VAE [Sadeghi'21]
Dynamic	DVAE [Girin'21] Sw-VAE [Sadeghi'21]	VQ-MDVAE [Sadok'23]	

# VAE/DVAE Summary

	Mono-modal	Multi-modal	Mixtures
Static	VAE [Kingma'14] VQ-VAE [van den Oord'17]	CVAE [Sadeghi'20] MVAE [Sutter'21]	VAE-MM [Sadeghi'20] MIN-VAE [Sadeghi'21]
Dynamic	DVAE [Girin'21] Sw-VAE [Sadeghi'21]	VQ-MDVAE [Sadok'23]	MixDVAE [Lin'23]

# Limited to VAE? NO!

We started investigated unsupervised AVSE with diffusion models:



Still need an EM algorithm at test time...

## What else is there to be done?

Plenty!!!

- Systematic AV for diffusion works.
- Can we design mixtures of diffusion models?
- What about latent diffusion?
- Can we have mixtures combined with latent diffusion?
- What about temporal diffusion models? And their mixtures?
- And beyond diffusion? e.g. flow matching?

## **Summary & conclusions**

## Take-home messages

- Latent variable modeling has endless possibilities.
- Auto-grad cannot save you.
- Analytic results can bring interpretability.
- Many possibilities are yet to be explored.
- Uncommon skills/knowledge (e.g. VI) can make your profile unique.

Thanks...



for bearing with me!

to my colleagues & collaborators!

for your challenging questions  
& interesting discussion.

We are also interested in meta-learning, reinforcement learning, domain adaptation, etc.

I am around until Thursday!