

# Variational Audio-Visual Representation Learning

---

Xavier Alameda-Pineda (Xavi)

Keynote Talk @ ACM Multimedia Ottawa, 1st November 2023



# Why Audio-Visual Unsupervised Representation Learning?

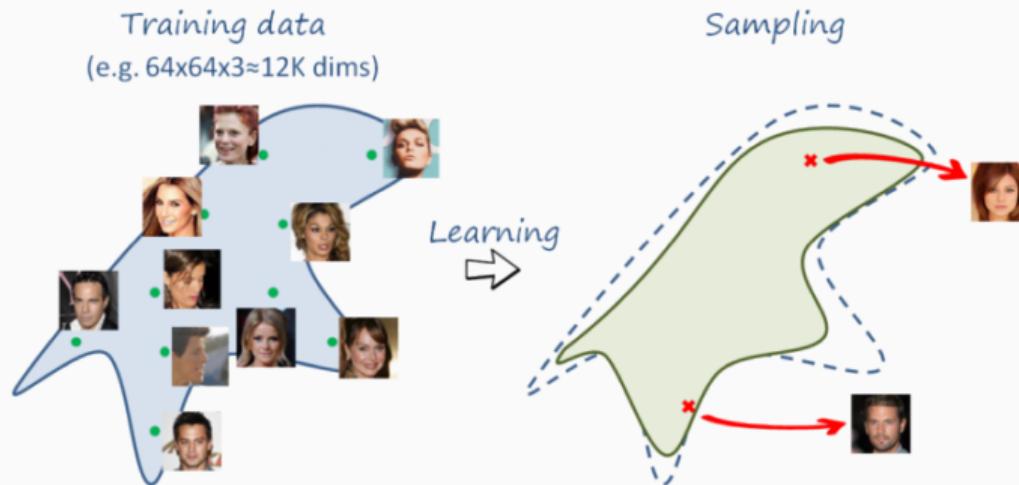


[Images under Creative Commons license]

- Audio: wide field of hearing, only usable when sources are active.
- Video: limited field of view, rich flow of information.
- Unsupervised: avoid the need for human labels (in new environments).

## And probabilistic learning?

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



- We can jointly learn them with other probabilistic models using *maximum likelihood*.

**Disclaimer: Halloween is not over**

www.iconsmind.com



Attendee Discretion is Advised: **Scary Equations Ahead**

## The Kullback-Leibler divergence and the ML formulation

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

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

# The Kullback-Leibler divergence and the ML formulation

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

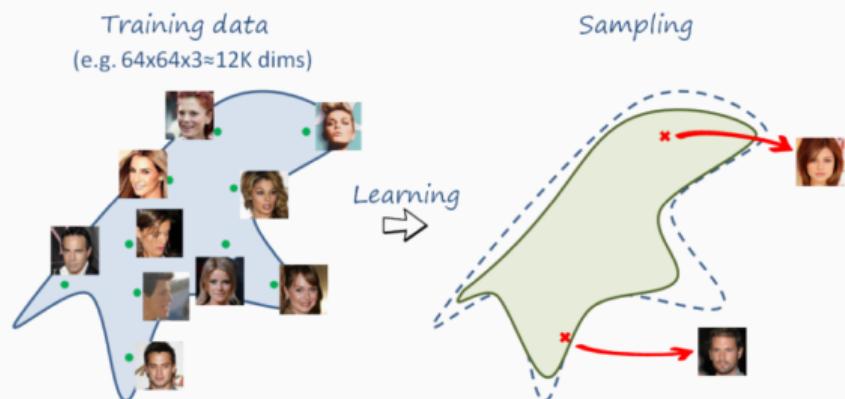
$$D_{\text{KL}}(p(\mathbf{x}) \| q(\mathbf{x})) = -\mathbb{E}_{p(\mathbf{x})} \left[ \log \frac{q(\mathbf{x})}{p(\mathbf{x})} \right] = - \int_{\mathcal{X}} p(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})} d\mathbf{x} \begin{cases} \geq 0 \\ 0 \Leftrightarrow p(\mathbf{x}) = q(\mathbf{x}) \\ \neq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

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

ML minimizes the KL divergence:

$$\boldsymbol{\theta}^* = \operatorname{argmin}_{\boldsymbol{\theta}} D_{\text{KL}}(p_{\text{data}}(\mathbf{x}) \| p_{\boldsymbol{\theta}}(\mathbf{x}))$$

$$\approx \operatorname{argmax}_{\boldsymbol{\theta}} \frac{1}{N} \sum_{i=1}^N \log p_{\boldsymbol{\theta}}(\mathbf{x}_i)$$



## ML formulation with latent variables

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

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

### Exact EM

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

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

Closed-form!

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

---

<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 formulation with latent variables

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

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

### Exact EM

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

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

Closed-form!

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

### Variational EM

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

$q(\mathbf{z}) = q_1(z_1)q_2(z_2)$

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

$D_{\text{KL}} > 0$  but 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 formulation with latent variables

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

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

## Exact EM

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

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

Closed-form!

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

## Variational EM

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

$q(\mathbf{z}) = q_1(z_1)q_2(z_2)$

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

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

## Variational AutoEncoder

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

$q(\mathbf{z}) = q_\phi(\mathbf{z})$

$q_\phi$  optimises ELBO/VLB

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

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

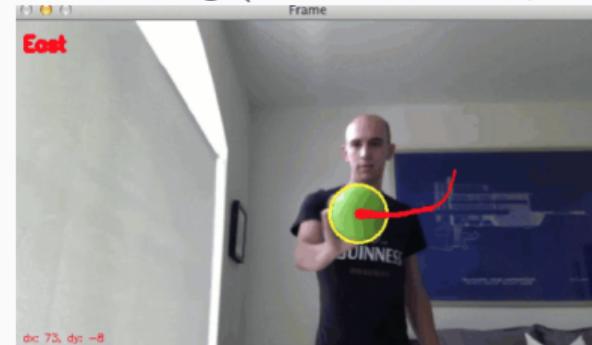
<sup>2</sup>Jordan, M. I., et. al., (1999), Machine Learning.

# Interest of Latent Variables

Speech Enhancement (**noisy** and **clean** signals)



Visual Tracking (**detections**, **true position**)

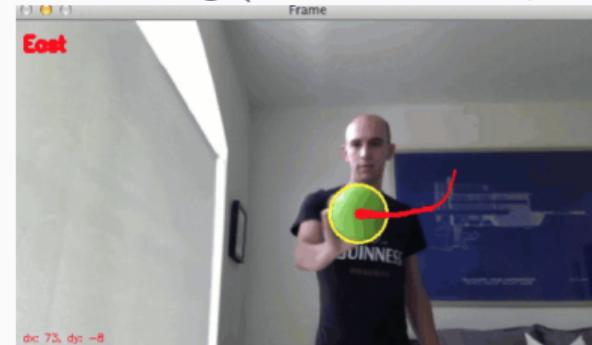


# Interest of Latent Variables

Speech Enhancement (noisy and clean signals)



Visual Tracking (detections, true position)



1. Define the model:

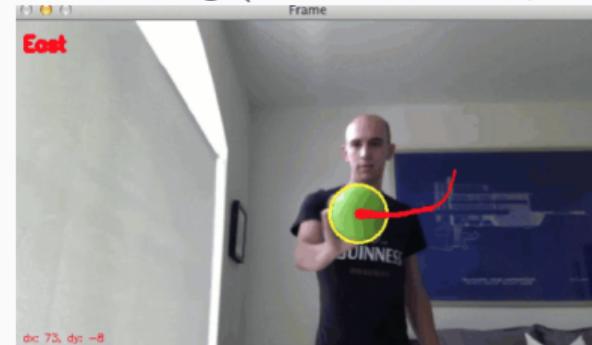
$$\begin{cases} p_{\theta}(s) \& p_{\theta}(x|s) & (\text{prior \& likelihood}) \\ q_{\phi}(s) \approx p_{\theta}(s|x) & (\text{approx. posterior}) \end{cases}$$

# Interest of Latent Variables

Speech Enhancement (noisy and clean signals)



Visual Tracking (detections, true position)



1. Define the model:

$$\begin{cases} p_{\theta}(s) \& p_{\theta}(x|s) & (\text{prior \& likelihood}) \\ q_{\phi}(s) \approx p_{\theta}(s|x) & (\text{approx. posterior}) \end{cases}$$

2. Learning & inference:

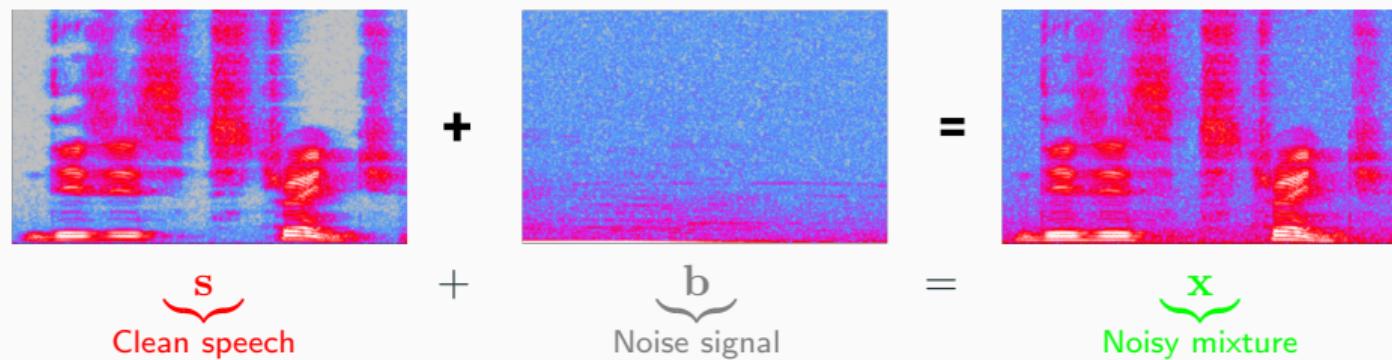
$$\underset{\theta, \phi}{\operatorname{argmax}} \mathcal{L}_{\text{ELBO}}(\theta, \phi)$$

$$\underset{s}{\operatorname{argmax}} q_{\phi^*}(s|x)$$

# Speech Enhancement & Wiener Filter



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



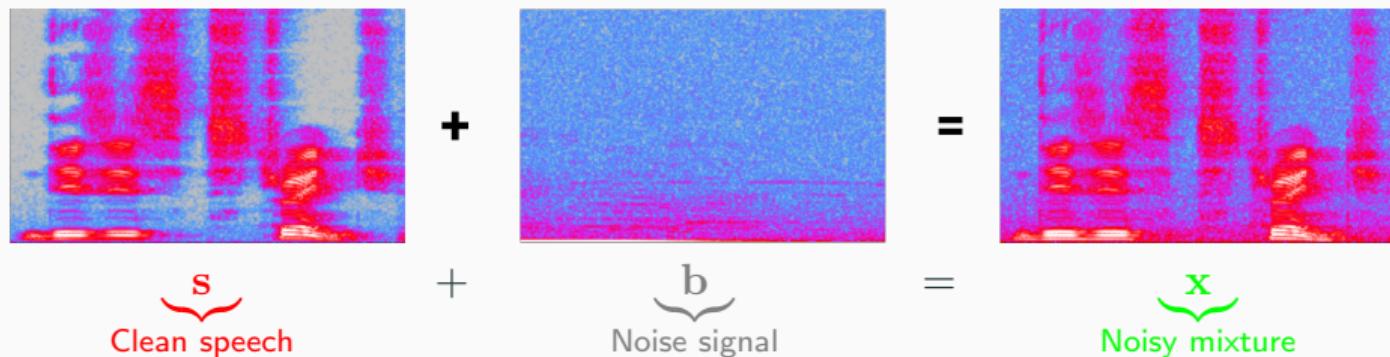
# Speech Enhancement & Wiener Filter



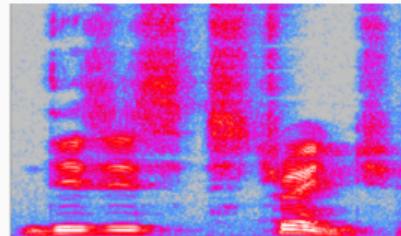
Minimize MSE → Wiener Filter  
(operations are element-wise)

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

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



# Unsupervised Probabilistic SE: paradigm



$\underbrace{s}_{\text{Clean speech}}$

$\downarrow$   
Noise signal

$\downarrow$   
Noisy mixture

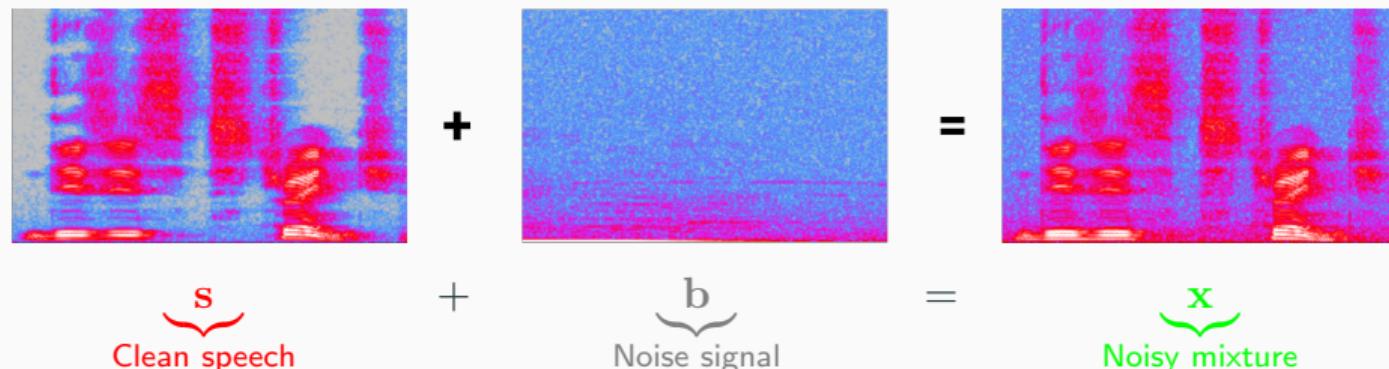
**Train** – Model for clean data:  $\{s_i\}_{i=1}^N$ .



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

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

# Unsupervised Probabilistic SE: paradigm

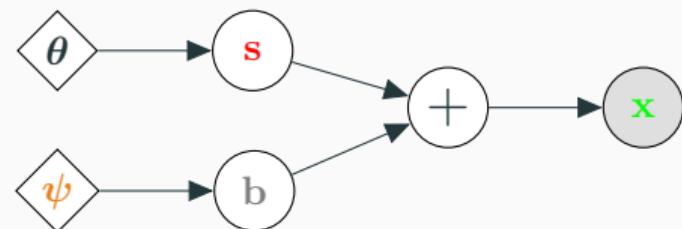


**Train** – Model for clean data:  $\{s_i\}_{i=1}^N$ .



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.*

►Can we jointly learn from AV data for SE?



Mostafa Sadeghi



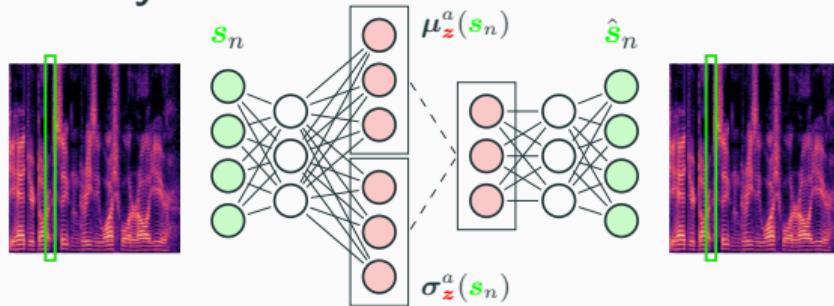
Laurent Girin



Radu Horaud

# Mono-modal VAEs: the baselines

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



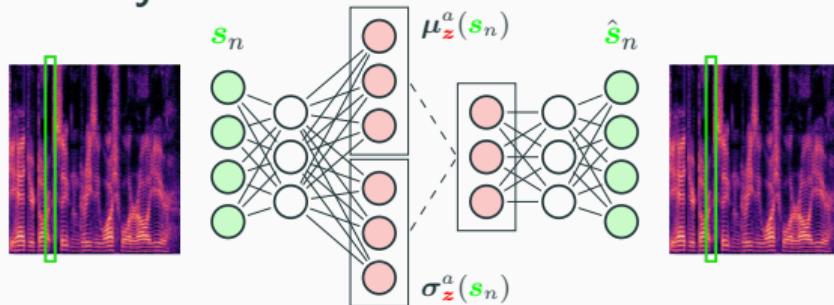
$$p_\theta(s_n|z_n) = \mathcal{N}_c(\mathbf{0}, \text{diag}(\sigma_s(z_n)))$$

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

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

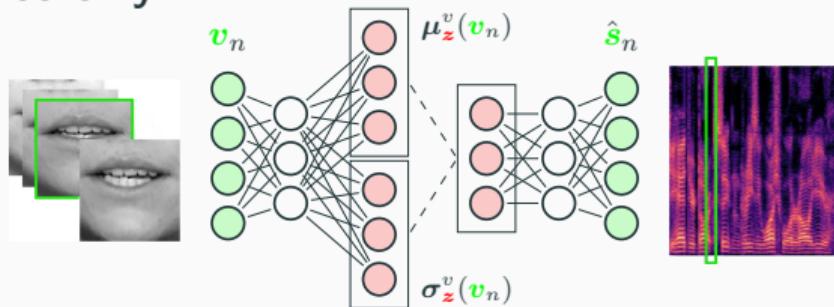
# Mono-modal VAEs: the baselines

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



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

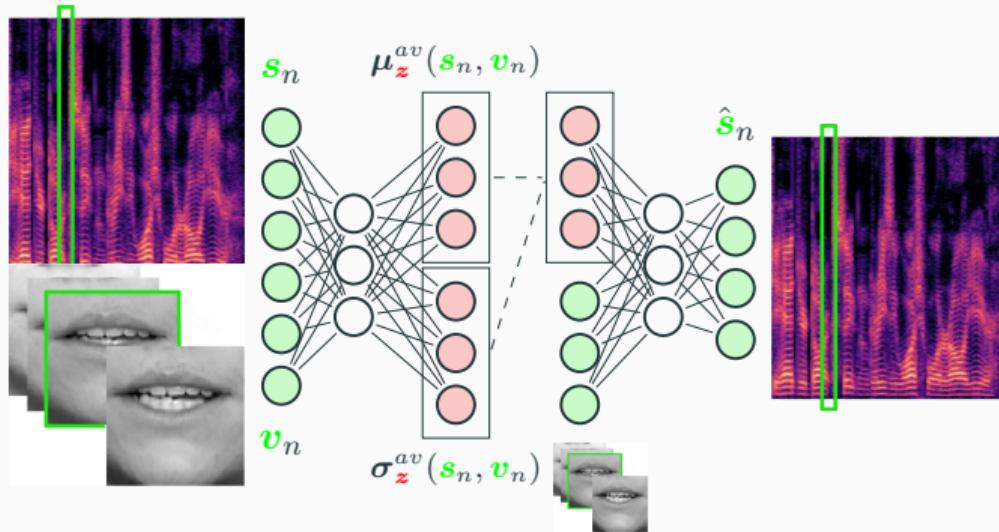
## Video-only:



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

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

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

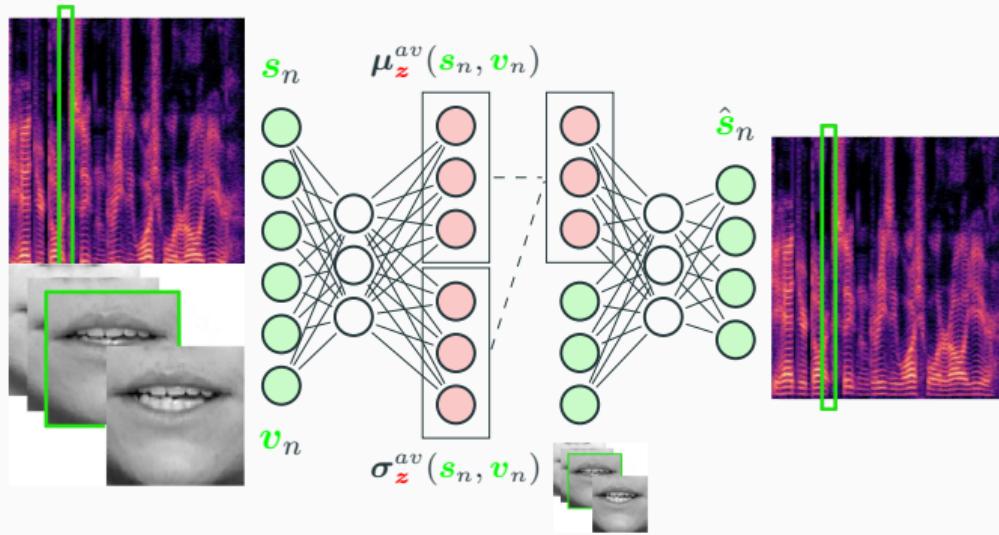


$$p_{\theta}(s_n | z_n, v_n) = \mathcal{N}_c(\mathbf{0}, \text{diag}(\sigma_s(z_n, v_n)))$$

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

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

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



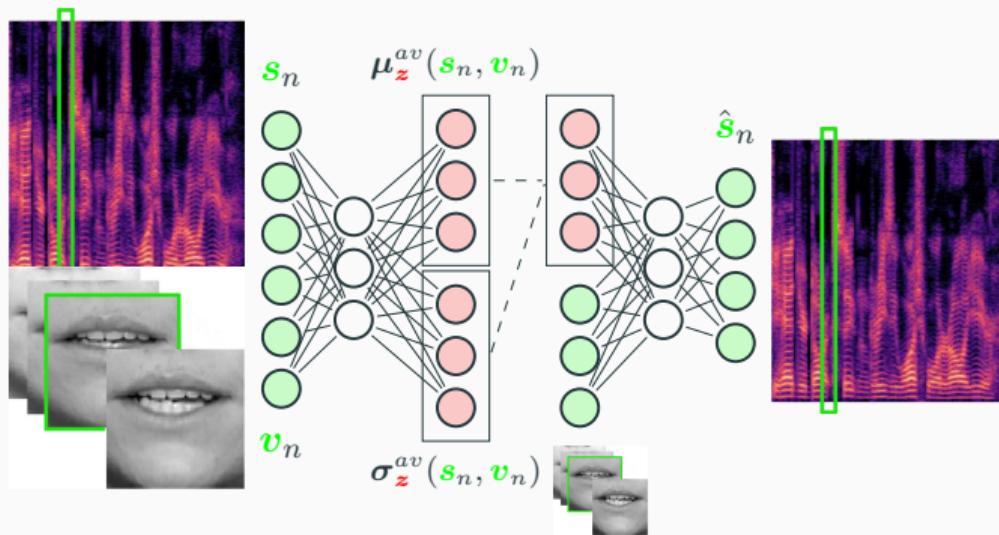
$$p_{\theta}(s_n | z_n, v_n) = \mathcal{N}_c(\mathbf{0}, \text{diag}(\sigma_s(z_n, v_n)))$$

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

What will this learn?

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

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



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

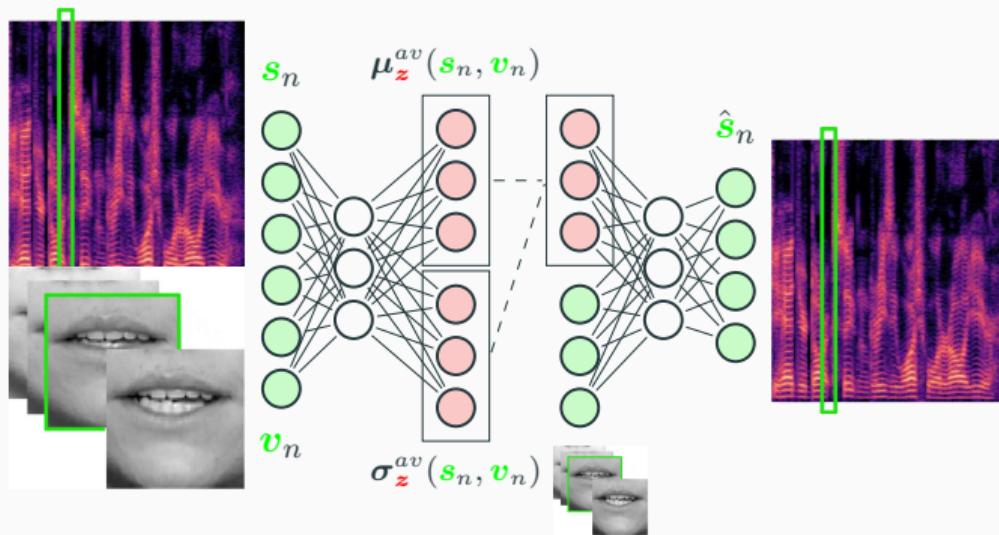
$$q_{\phi}(z_n | v_n, s_n) = \mathcal{N}\left(\mu_z^{av}(v_n, s_n), \text{diag}(\sigma_z^{av}(v_n, s_n))\right)$$

What will this learn?

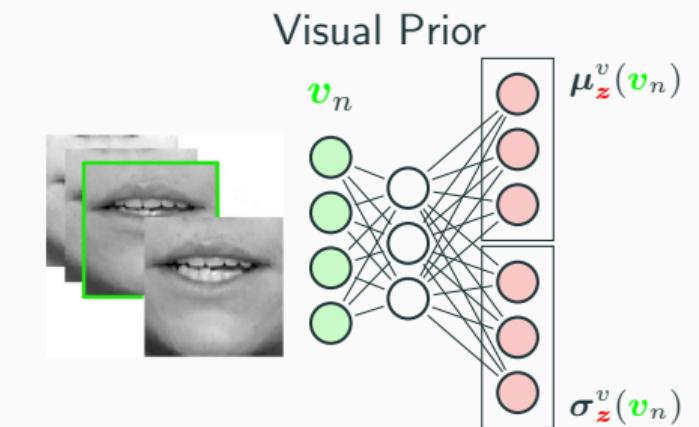
To ignore the video!

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

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



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



$$p_{\theta}(\mathbf{z}_n|\mathbf{v}_n) = \mathcal{N}\left(\boldsymbol{\mu}_{\mathbf{z}}^v(\mathbf{v}_n), \text{diag}(\boldsymbol{\sigma}_{\mathbf{z}}^v(\mathbf{v}_n))\right)$$

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

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

## Learning AV Conditional VAE

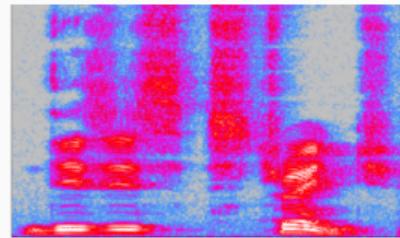
On top of adding a video-only prior, we optimize a modified ELBO:

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}) = & \sum_n (1 - \alpha) \left( \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_n | \mathbf{s}_n, \mathbf{v}_n)} \left[ \ln p_{\boldsymbol{\theta}}(\mathbf{s}_n | \mathbf{z}_n, \mathbf{v}_n) \right] - D_{\text{KL}} \left( q_{\boldsymbol{\phi}}(\mathbf{z}_n | \mathbf{s}_n, \mathbf{v}_n) \| p_{\boldsymbol{\theta}}(\mathbf{z}_n | \mathbf{v}_n) \right) \right) \\ & + \alpha \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{z}_n | \mathbf{v}_n)} \left[ \ln p_{\boldsymbol{\theta}}(\mathbf{s}_n | \mathbf{z}_n, \mathbf{v}_n) \right]\end{aligned}$$

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

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

# AV Conditional VAE for Speech Enhancement

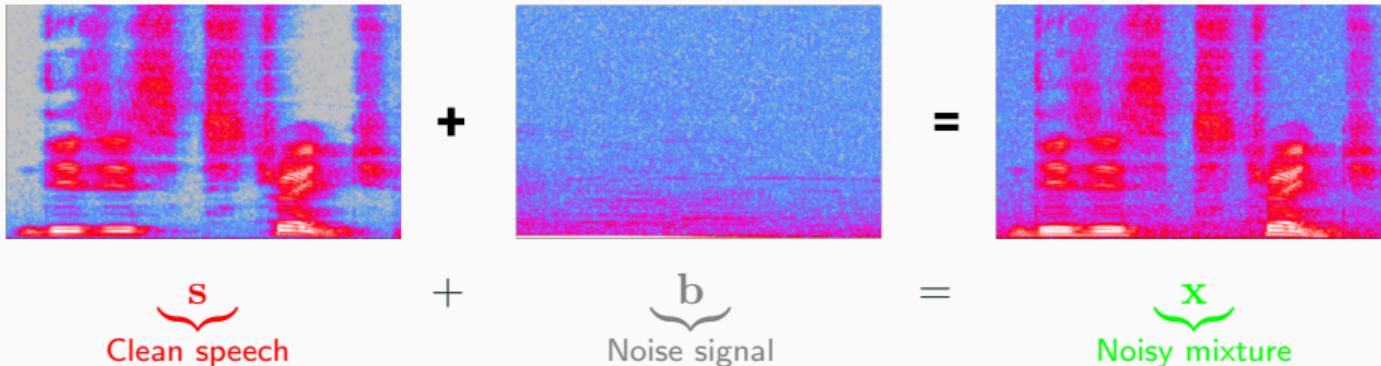


**Train** – Model for **clean data**:  $\{s_i\}_{i=1}^N$ .

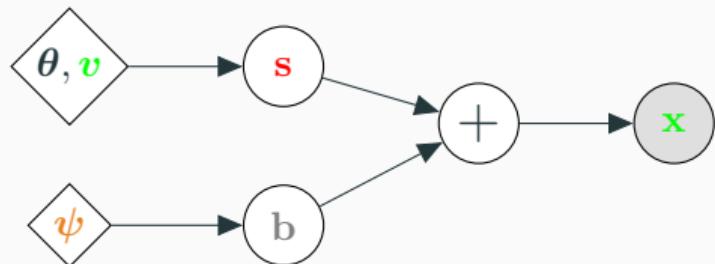


Then freeze  $\theta$  at test/adaptation time.

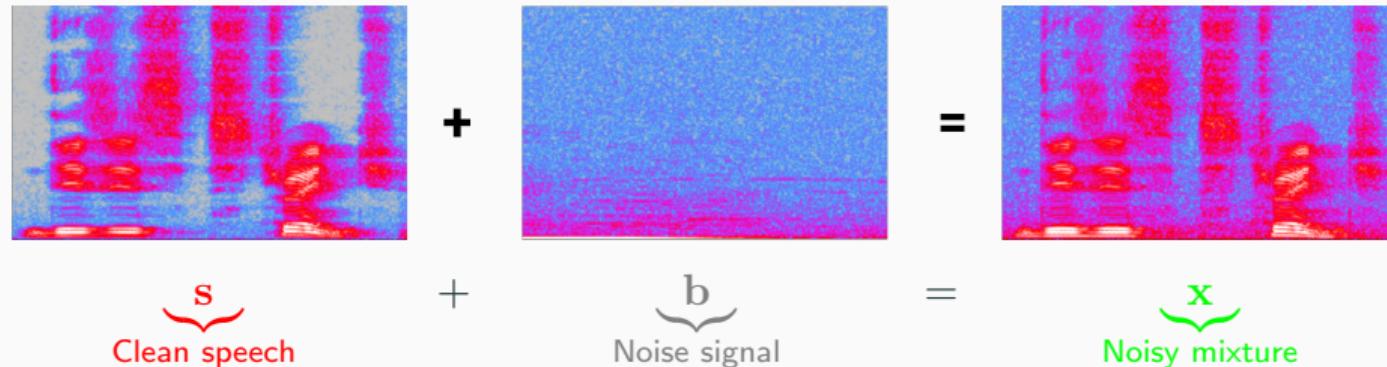
# AV Conditional VAE for Speech Enhancement



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



# AV Conditional VAE for Speech Enhancement



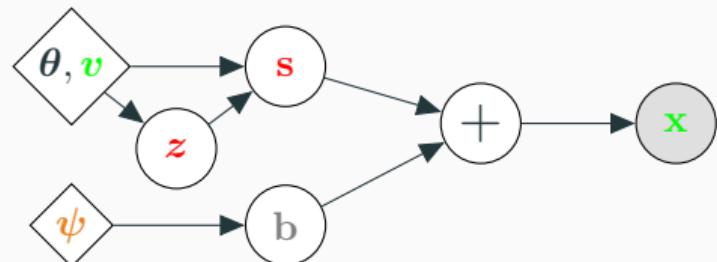
**Speech Enhancement:**

(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!

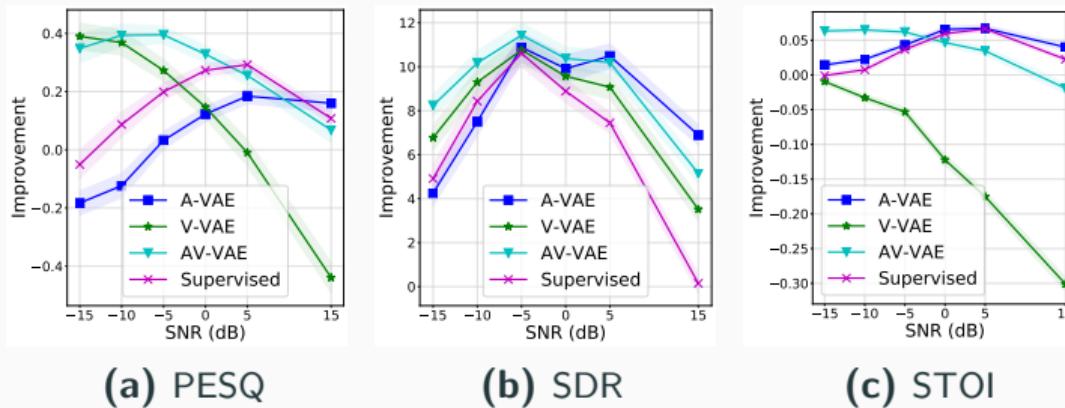
**Test** – learn the **noise parameters** from **noisy samples  $x$**  and estimate the **clean speech  $\hat{s}$** :



## Experiments & Discussion

**NTCD-TIMIT** dataset:<sup>5</sup> AV recordings in controlled conditions, 5h/39 speakers training, 1h/9 speakers test, several noise levels (−15 to 15 dB) and types (car, living room, etc).

**Metrics** improvement w.r.t. the noisy mixture: perceptual evaluation of speech quality (PESQ), signal-to-distortion ratio (SDR), Short-time objective intelligibility (STOI).

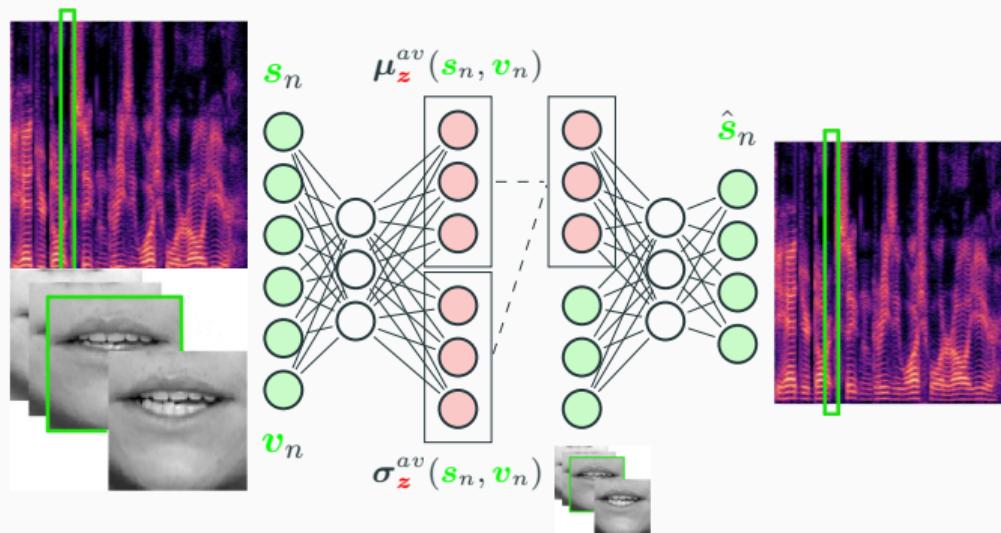


Examples: <https://team.inria.fr/robotlearn/research/av-vae-se/>

<sup>5</sup> Abdelaziz, A.H., (2017), Interspeech.

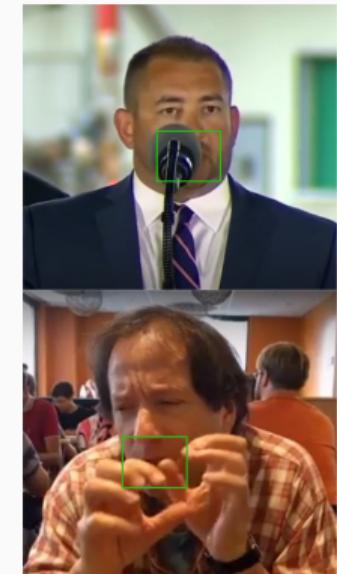
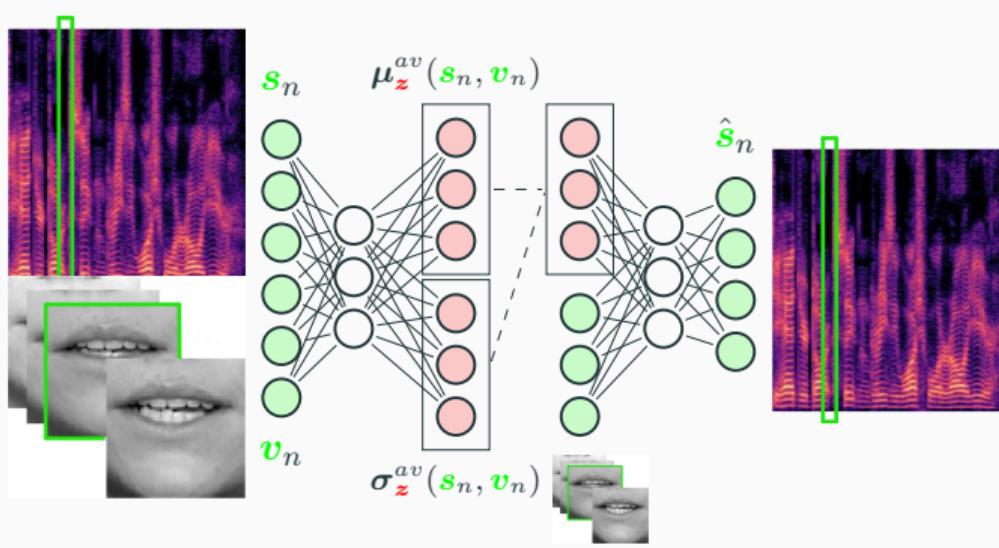
## Limitations: systematic AV fusion

From the model design:



## Limitations: systematic AV fusion

From the model design:



What about clutter?

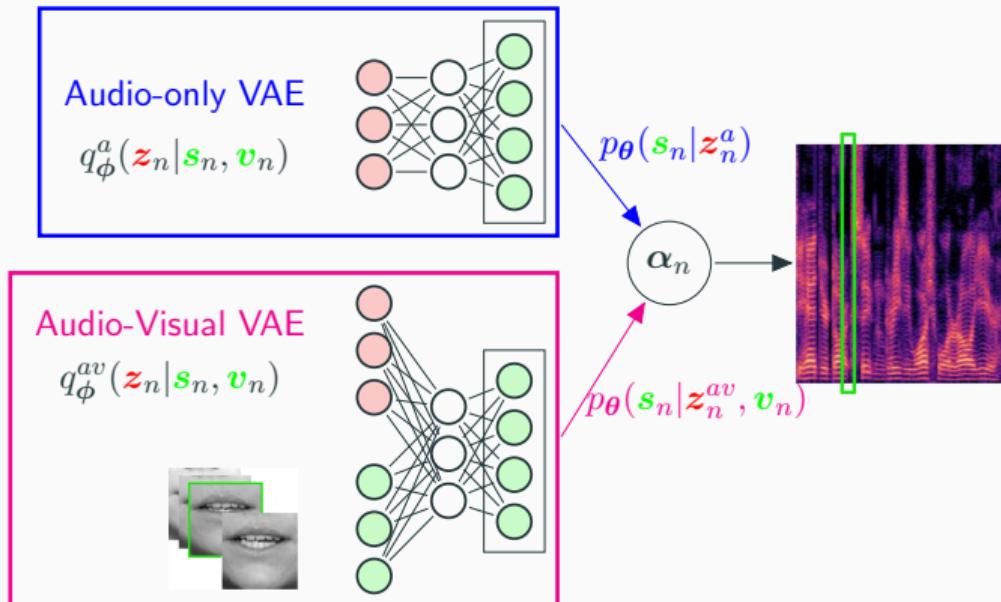
We **ALWAYS** concatenate audio and video information.

- Can we unsupervisedly select what to use?



Mostafa Sadeghi

# VAE Mixture Model<sup>6</sup>



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

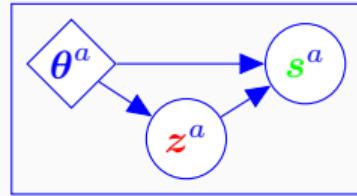
$$\begin{cases} \alpha_n = 1 \leftrightarrow \text{Audio-only} \\ \alpha_n = 0 \leftrightarrow \text{Audio-Visual} \end{cases}$$

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

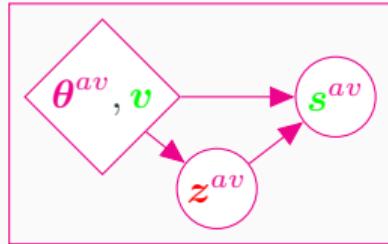
The **Audio-only** and the **Audio-Visual** VAE are **mixed without supervision**.

<sup>6</sup>Sadeghi, M., et. al., (2021), IEEE ICASSP.

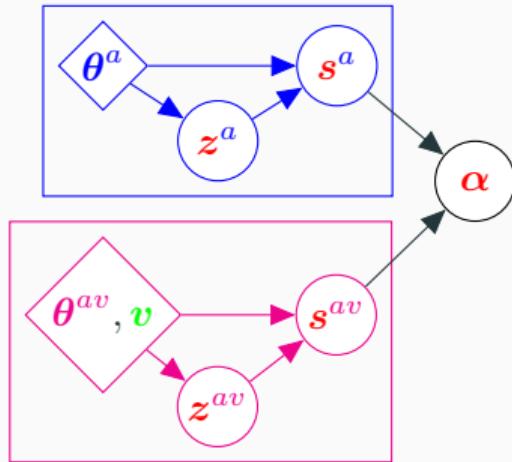
## Train, adaptation & speech enhancement



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

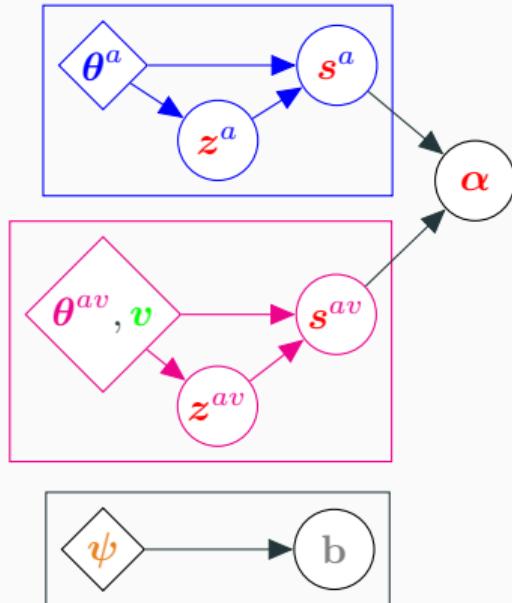


## Train, adaptation & speech enhancement



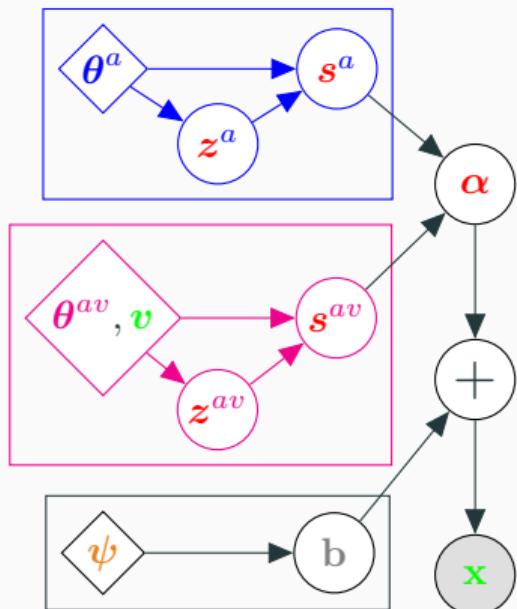
- Train A-VAE (with  $s$ ) and AV-VAE (with  $(s, v)$ ).
- The clean speech is now **latent** ( $s^a, s^{av}$ ), add the **mixing latent** variable ( $\alpha$ ).

## Train, adaptation & speech enhancement



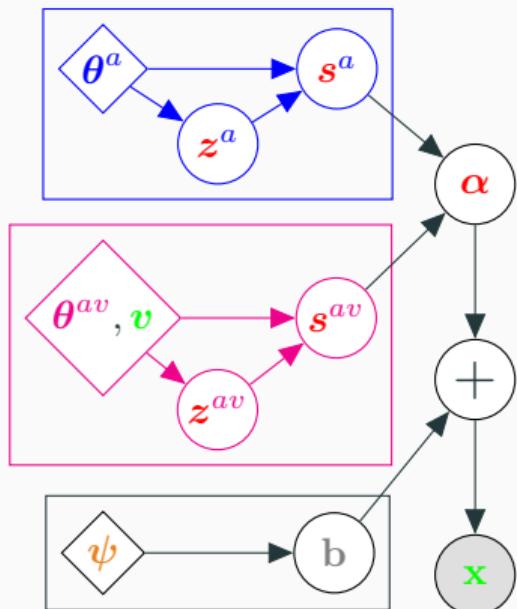
- Train A-VAE (with  $s$ ) and AV-VAE (with  $(s, v)$ ).
- The clean speech is now **latent** ( $s^a, s^{av}$ ), add the **mixing latent** variable ( $\alpha$ ).
- We also consider the noise variable and params.

## 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



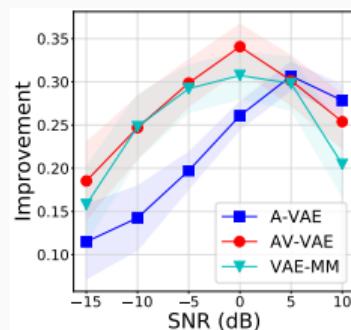
- Train A-VAE (with  $s$ ) and AV-VAE (with  $(s, v)$ ).
- The clean speech is now **latent** ( $s^a, s^{av}$ ), add the **mixing latent** variable ( $\alpha$ ).
- We also consider the noise variable and params.
- Learn  $\psi^*$  and estimate the clean speech. By defining  $\gamma_n(\mathbf{z}_n, \mathbf{v}_n)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

$$\hat{s}_n = \frac{1}{R} \sum_{r=1}^R \frac{\gamma_n(\mathbf{z}_n^{(r)}, \mathbf{v}_n)}{\gamma_n(\mathbf{z}_n^{(r)}, \mathbf{v}_n) + \psi_n^*} \mathbf{x}_n$$

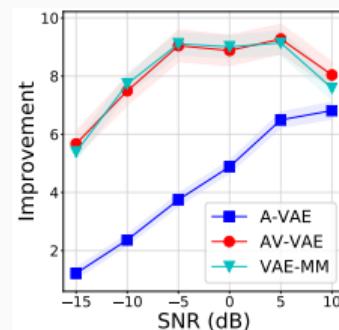
## Settings & Results

- **Data & models:** Same dataset + pre-trained A-VAE and AV-VAE.
- **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)

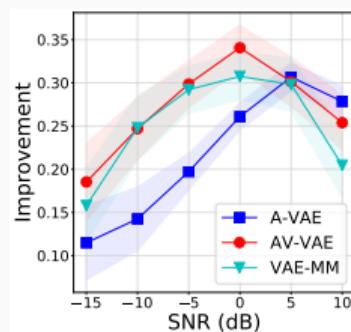


(b) SDR (clean)

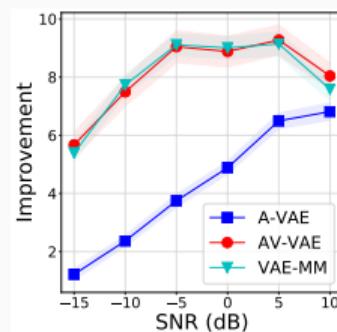
## Settings & Results

- **Data & models:** Same dataset + pre-trained A-VAE and AV-VAE.
- **Setup:** Very similar than in the previous experiments. **Clean and noisy** lips region visual information ( $\sim$  one-third of total video frames per sample).

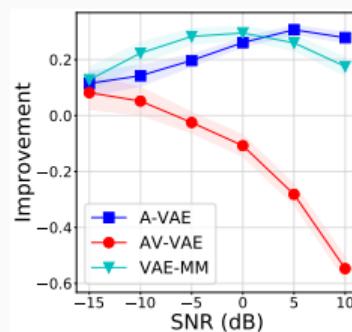
Improvement with respect to the input:



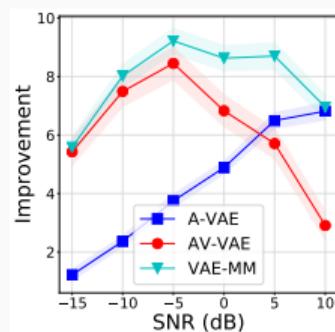
(a) PESQ (clean)



(b) SDR (clean)

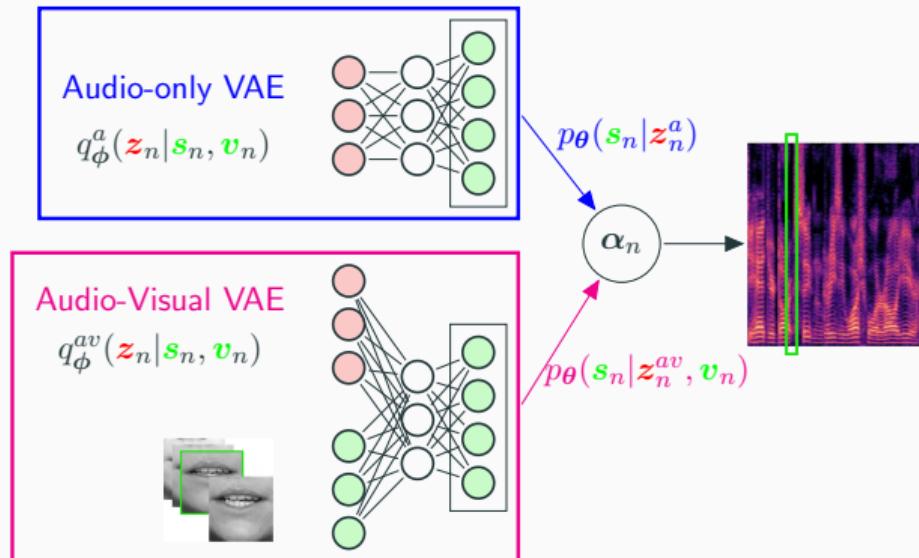


(c) PESQ (noisy)



(d) SDR (noisy)

# VAE Mixture Model: 2 models for the same signal?



For the same input  $s$ , we compute:

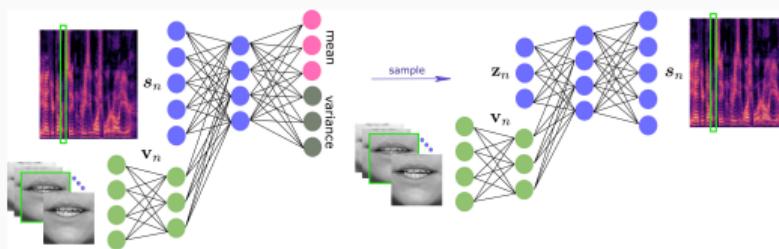
- An **audio-based** reconstruction  $s^a$ .
- An **AV-based** reconstruction  $s^{av}$ .

Not necessarily equal.

This is strange!!! We have proposed a model with a single decoder.<sup>7</sup>

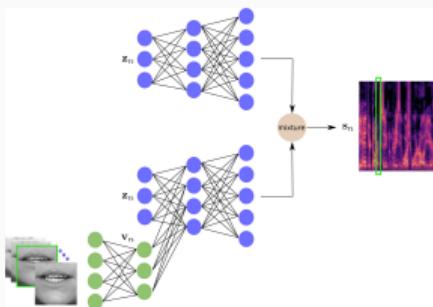
<sup>7</sup>Sadeghi, M., et. al., (2021), IEEE TSP.

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



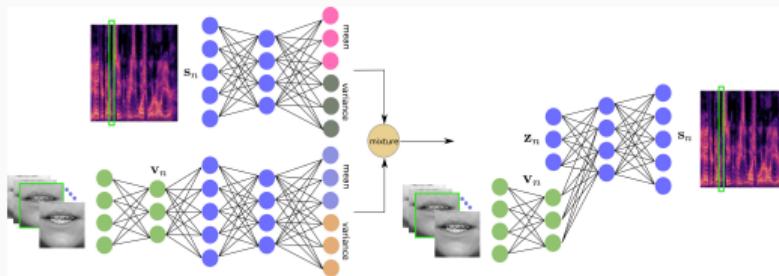
Conditional VAE:

- Training VAE via SGD.
- Systematic AV fusion.
- Appeared at IEEE TASLP in 2020.



VAE-MM:

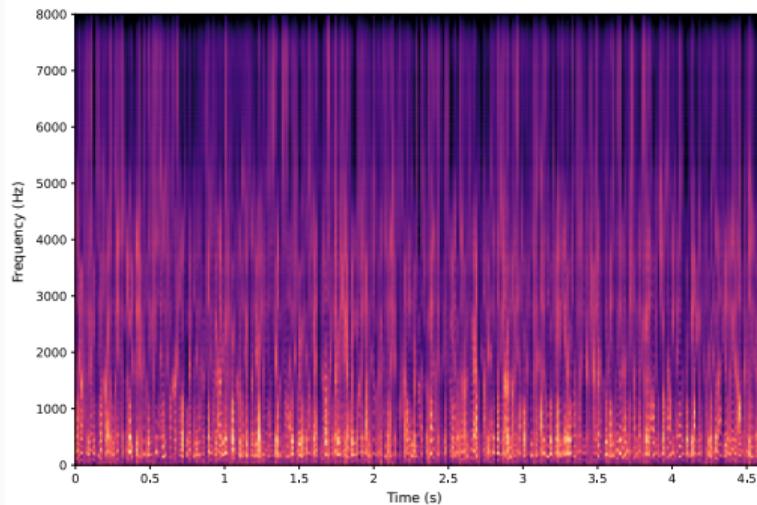
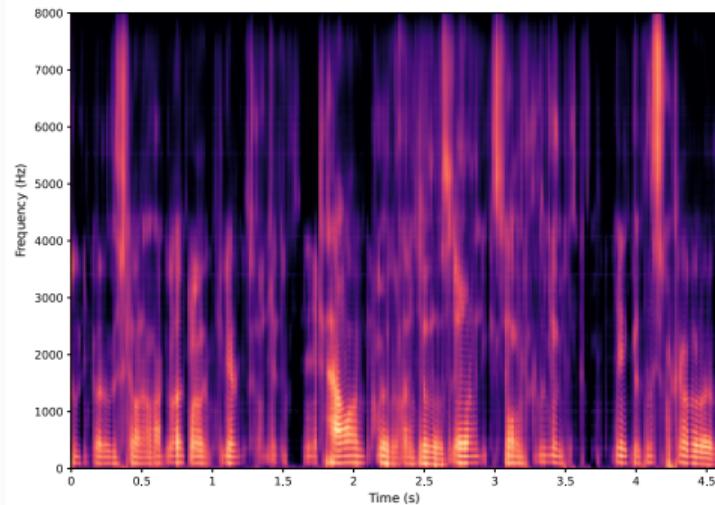
- Training two VAEs via SGD.
- Mixing (AV fusion) via two speech models.
- Appeared at IEEE ICASSP in 2021.



MIN-VAE:

- Training all 3 networks via VEM+SGD.
- Mixing (AV fusion) via a single speech model.
- Appeared at IEEE TSP in 2021.

## Strong limitation of VAEs: frame modeled independently.



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

- **Dynamical VAE (DVAE<sup>8</sup>) – “VAE for sequential modeling”**  
**(family of methods including existing literature)**



**Xiaoyu Bie**

**Laurent Girin**

**Simon Leglaive**

**Thomas Hueber**

**Julien Diard**

<sup>8</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 ( $\mathbf{z}_{1:T}$ ):

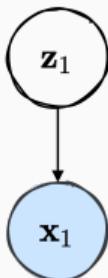
$$p_{\theta}^{\text{DVAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) \neq \prod_{t=1}^T p_{\theta}(\mathbf{x}_t, \mathbf{z}_t) = p_{\theta}^{\text{VAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}).$$



## Probabilistic Sequential Modeling (decoder network)

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

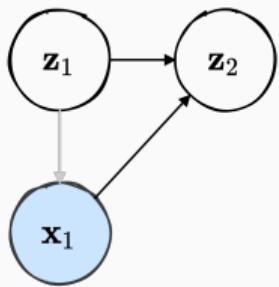
$$p_{\theta}^{\text{DVAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) \neq \prod_{t=1}^T p_{\theta}(\mathbf{x}_t, \mathbf{z}_t) = p_{\theta}^{\text{VAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}).$$



## Probabilistic Sequential Modeling (decoder network)

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

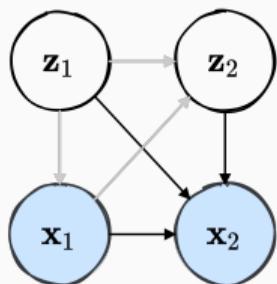
$$p_{\theta}^{\text{DVAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) \neq \prod_{t=1}^T p_{\theta}(\mathbf{x}_t, \mathbf{z}_t) = p_{\theta}^{\text{VAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}).$$



## Probabilistic Sequential Modeling (decoder network)

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

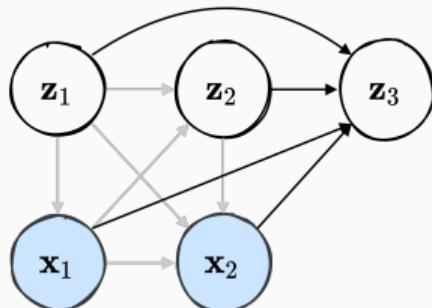
$$p_{\theta}^{\text{DVAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) \neq \prod_{t=1}^T p_{\theta}(\mathbf{x}_t, \mathbf{z}_t) = p_{\theta}^{\text{VAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}).$$



## Probabilistic Sequential Modeling (decoder network)

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

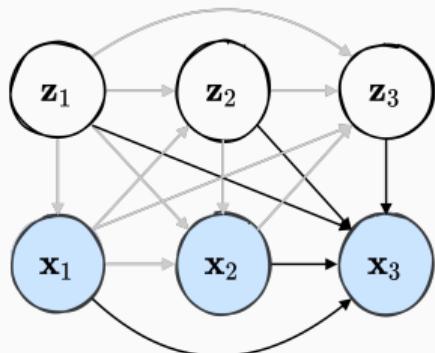
$$p_{\theta}^{\text{DVAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) \neq \prod_{t=1}^T p_{\theta}(\mathbf{x}_t, \mathbf{z}_t) = p_{\theta}^{\text{VAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}).$$



## Probabilistic Sequential Modeling (decoder network)

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

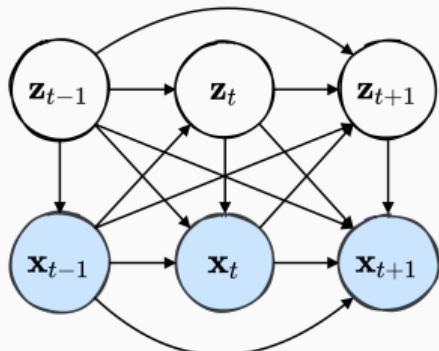
$$p_{\theta}^{\text{DVAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) \neq \prod_{t=1}^T p_{\theta}(\mathbf{x}_t, \mathbf{z}_t) = p_{\theta}^{\text{VAE}}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}).$$



## 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:

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

## Probabilistic Sequential Inference (encoder network)

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

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

We can always use the Bayes theorem to write:

$$q_{\phi}(z_{1:T} | x_{1:T}) = \prod_{t=1}^T q_{\phi}(z_t | z_{1:t-1}, x_{1:T})$$

Can we simplify each of the terms further? It depends on generative model. Use **D-separation**<sup>9</sup> to find out the true dependencies, then choose whether to keep them (e.g. to ensure causality).

---

<sup>9</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>10</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:

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

There are many possible implementations<sup>10</sup> for the same probabilistic dependencies!!!

---

<sup>10</sup>Several DVAEs @ <https://github.com/XiaoyuBIE1994/DVAE-speech>

## Learning: ELBO

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

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

- 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:

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

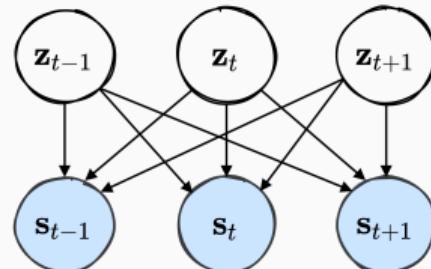
- The reconstruction and regularisation terms are evaluated at every frame  $t$ .
- Because of the model, the KL term depends on previous latent variables  $\Rightarrow$  sampling.
- The sampling occurs sequentially and cannot be parallelized!

## Application to Unsupervised Probabilistic SE (revisit)

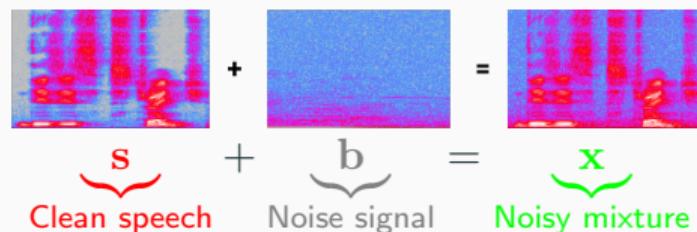


**Train** – Model  $\theta$  for clean data:  $\{s_{1:T}\}$ .

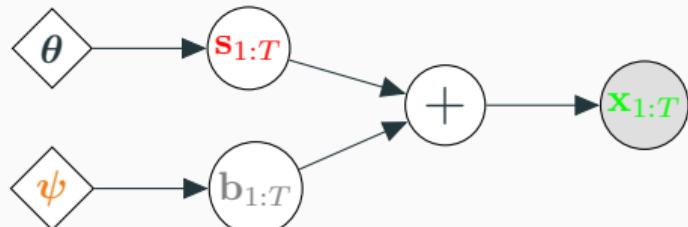
$$\theta \rightarrow s_{1:T} \quad p_\theta(s_{1:T}) \approx p_{\text{data}}(s_{1:T})$$



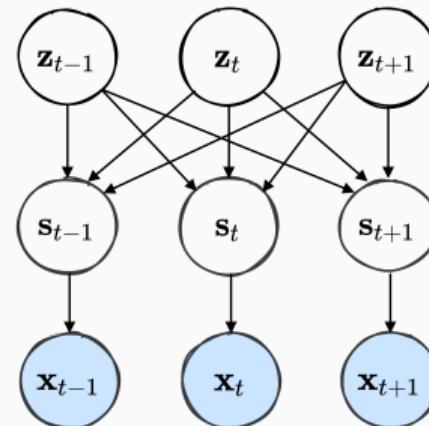
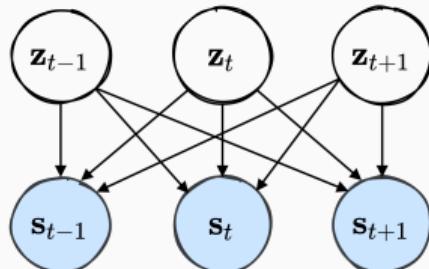
# Application to Unsupervised Probabilistic SE (revisit)



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



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



# Results on the Wall Street Journal (WSJ) and VoiceBank (VB)<sup>11</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	5.7	3.6		
UMX* 2020	Full			N/A				
MetricGAN+* (Interspeech'21)	Full							
Noisy mixture	-	WSJ+QUT	Same	N/A	5.7	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>11</sup>Bie, X., et. al., (2022), IEEE TASLP.

# Results on the Wall Street Journal (WSJ) and VoiceBank (VB)<sup>11</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	<u>5.8</u>
RVAE-VEM (Proposed)	None			<u>5.8</u>		
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	-	-	-	8.4	-	8.4
VAE-VEM (ICASSP'20)	None	VB-DMD	Same	16.4	<u>17.1</u>	<u>17.7</u>
RVAE-VEM (Proposed)	None			<u>17.1</u>		
NyTT EUSIPCO'21	Partial/Xtra	VB-DMD	Same	<u>17.7</u>	8.2	14.0
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>11</sup>Bie, X., et. al., (2022), IEEE TASLP.

# Results on the Wall Street Journal (WSJ) and VoiceBank (VB)<sup>11</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		Same	<b>5.8</b>		<b>4.3</b>
MetricGAN-U (ICASSP'22)	Partial			N/A		-1.6
UMX* 2020	Full	Same		<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	Same		<b>17.7</b>		
MetricGAN-U (half) ICASSP'22	Partial			8.2		
UMX 2020	Full	Same		14.0		
MetricGAN+ Interspeech'21	Full			8.5		

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

<sup>11</sup>Bie, X., et. al., (2022), IEEE TASLP.

# Results on the Wall Street Journal (WSJ) and VoiceBank (VB)<sup>11</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		Same	<b>5.8</b>		<b>4.3</b>
MetricGAN-U (ICASSP'22)	Partial			N/A		-1.6
UMX* 2020	Full	VB-DMD		<u>5.7</u>	Different	<u>4.1</u>
MetricGAN+* (Interspeech'21)	Full		Same	3.6		1.8
Noisy mixture	-		-	8.4	-	8.4
VAE-VEM (ICASSP'20)	None	VB-DMD		16.4		<u>15.0</u>
RVAE-VEM (Proposed)	None		Same	<u>17.1</u>		<b>17.3</b>
NyTT EUSIPCO'21	Partial/Xtra			<b>17.7</b>		N/A
MetricGAN-U (half) ICASSP'22	Partial	VB-DMD		8.2	Different	N/A
UMX 2020	Full		Same	14.0		10.4
MetricGAN+ Interspeech'21	Full			8.5		3.9

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

<sup>11</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



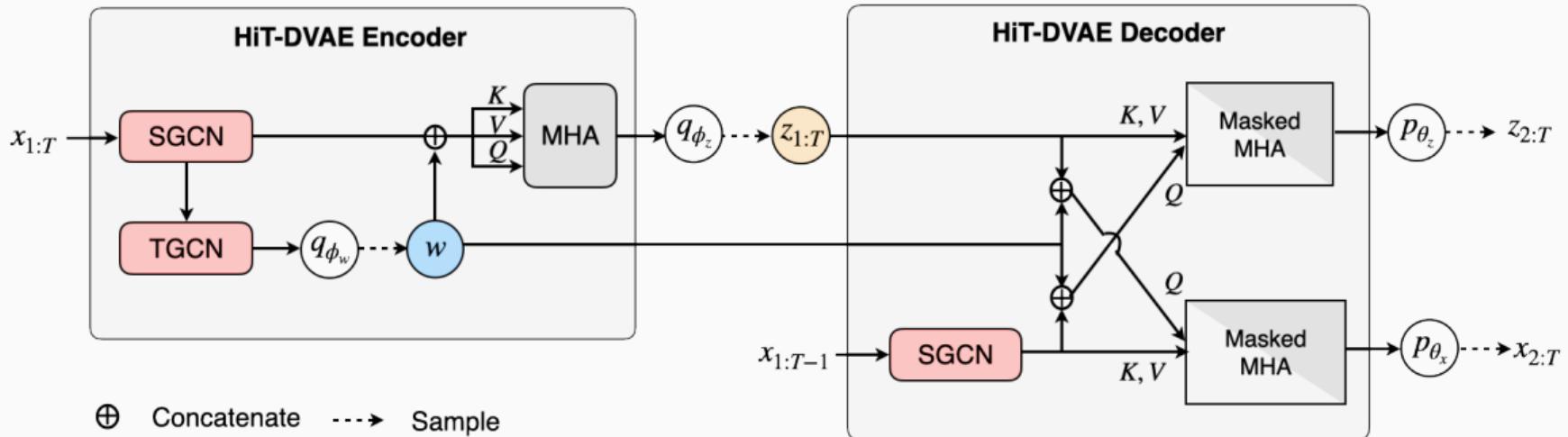
Francesc Moreno



Laurent Girin

# Introducing HiT-DVAE<sup>12</sup>

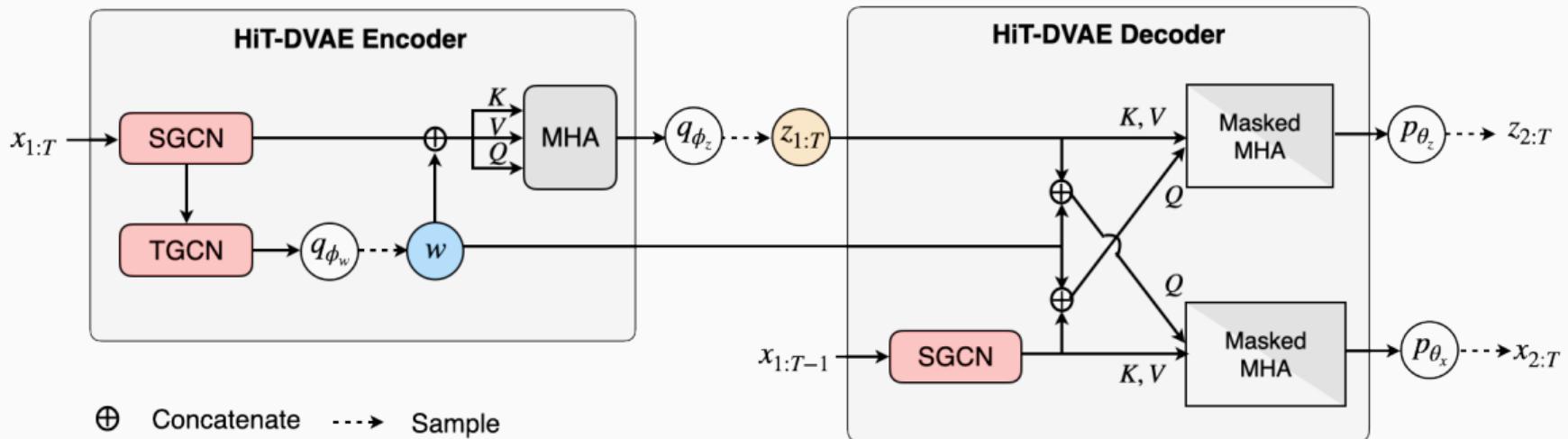
Hierarchical Transformer DVAE – two levels of latent variables: **static  $w$**  and **dynamic  $z_{1:T}$** .



<sup>12</sup>Bie, X., et. al., (2023), Pre-print.

# Introducing HiT-DVAE<sup>12</sup>

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

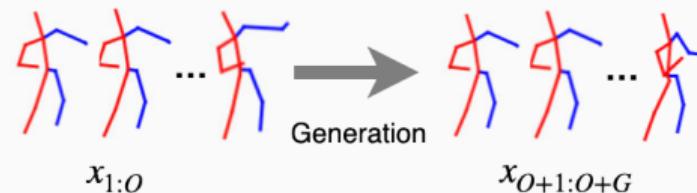


There is something **odd** in the way  $Q$ ,  $K$ ,  $V$  are used...

<sup>12</sup>Bie, X., et. al., (2023), Pre-print.

## HIT-DVAE: Modified transformer architecture

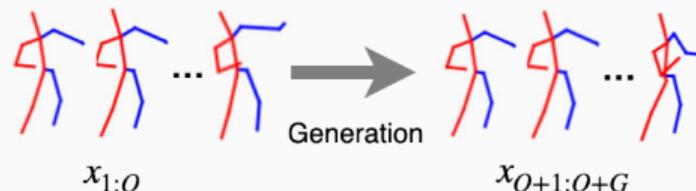
Task: human motion modeling/generation. Predicting  $x_{O+1}$  from  $x_{1:O}$  is very easy.<sup>13</sup>



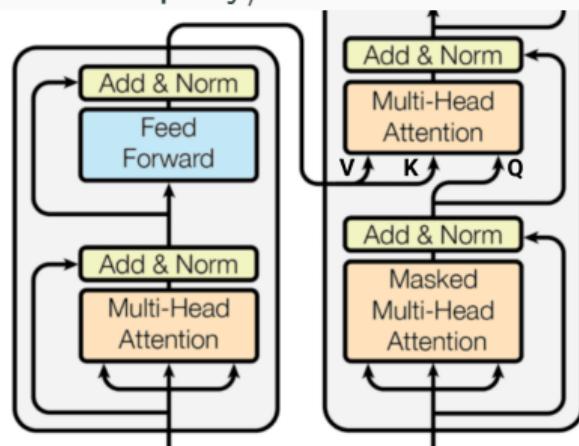
<sup>13</sup>Guo., W., et. al., (2021), IEEE WACV.

# HIT-DVAE: Modified transformer architecture

Task: human motion modeling/generation. Predicting  $x_{O+1}$  from  $x_{1:O}$  is very easy.<sup>13</sup>



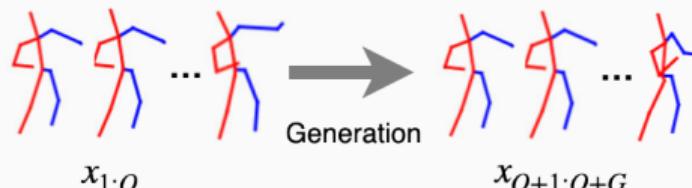
Standard query/residual connection:



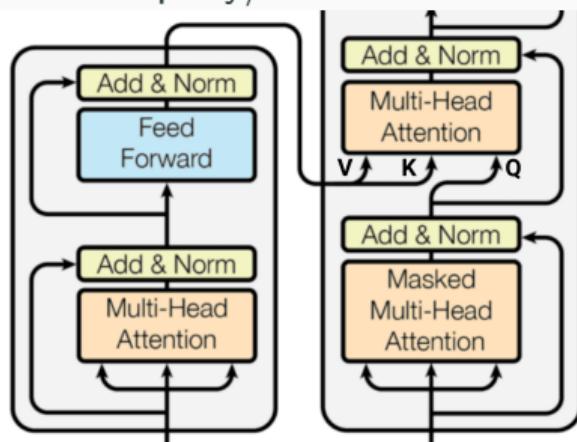
<sup>13</sup>Guo., W., et. al., (2021), IEEE WACV.

# HIT-DVAE: Modified transformer architecture

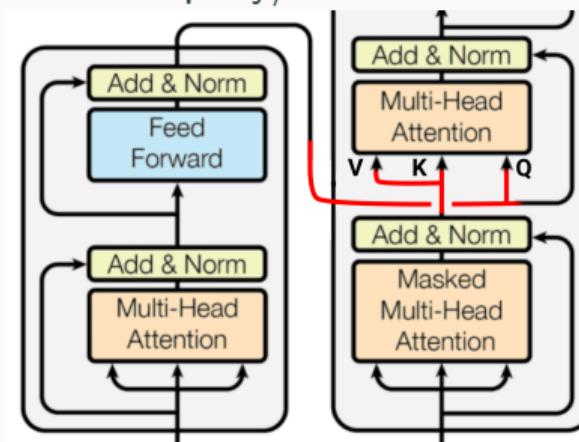
Task: human motion modeling/generation. Predicting  $x_{O+1}$  from  $x_{1:O}$  is very easy.<sup>13</sup>



Standard query/residual connection:



HIT-DVAE query/residual connection:



<sup>13</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>14</sup>

<sup>14</sup>Lin, X., et. al., (2023), IEEE ICASSP.

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

► What about multiple modalities?



Samir Sadok



Simon Leglaise

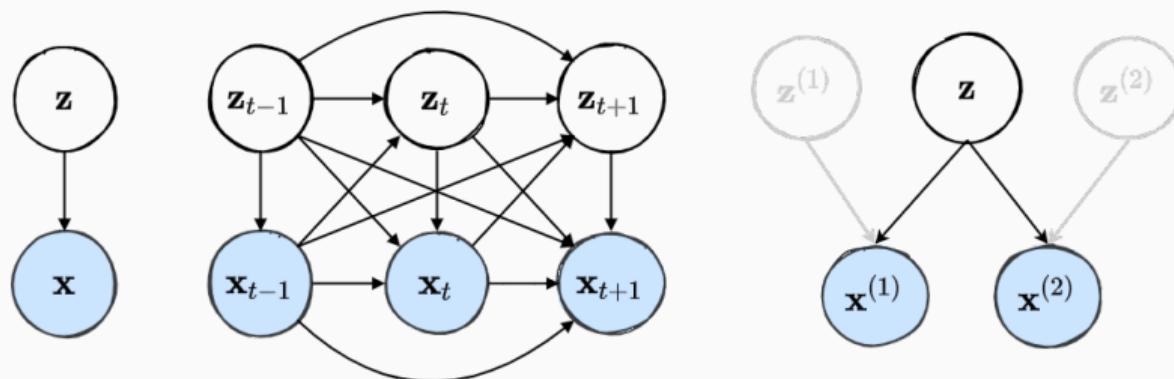


Renaud Séguier



Laurent Girin

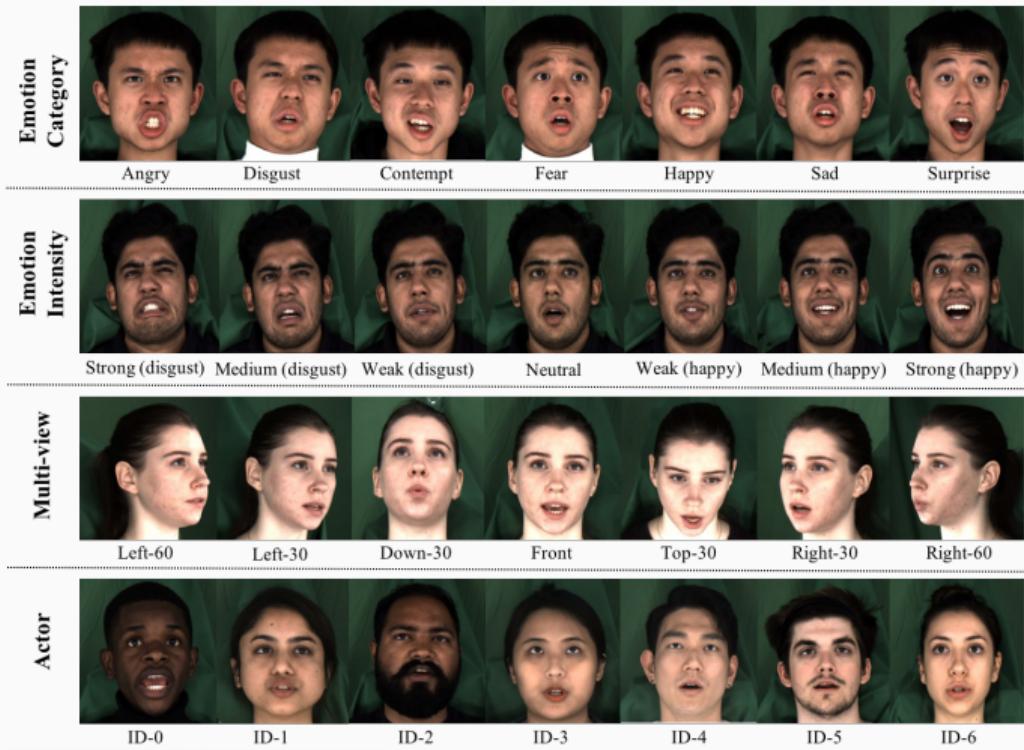
In other words:



VAEs (left) can be multi-modal<sup>15</sup> (right). Can DVAEs (middle) be multi-modal too?

<sup>15</sup>Sutter, T. M., et. al., (2021), ICLR.

# Task: emotional audio-visual speech modeling.<sup>16</sup>

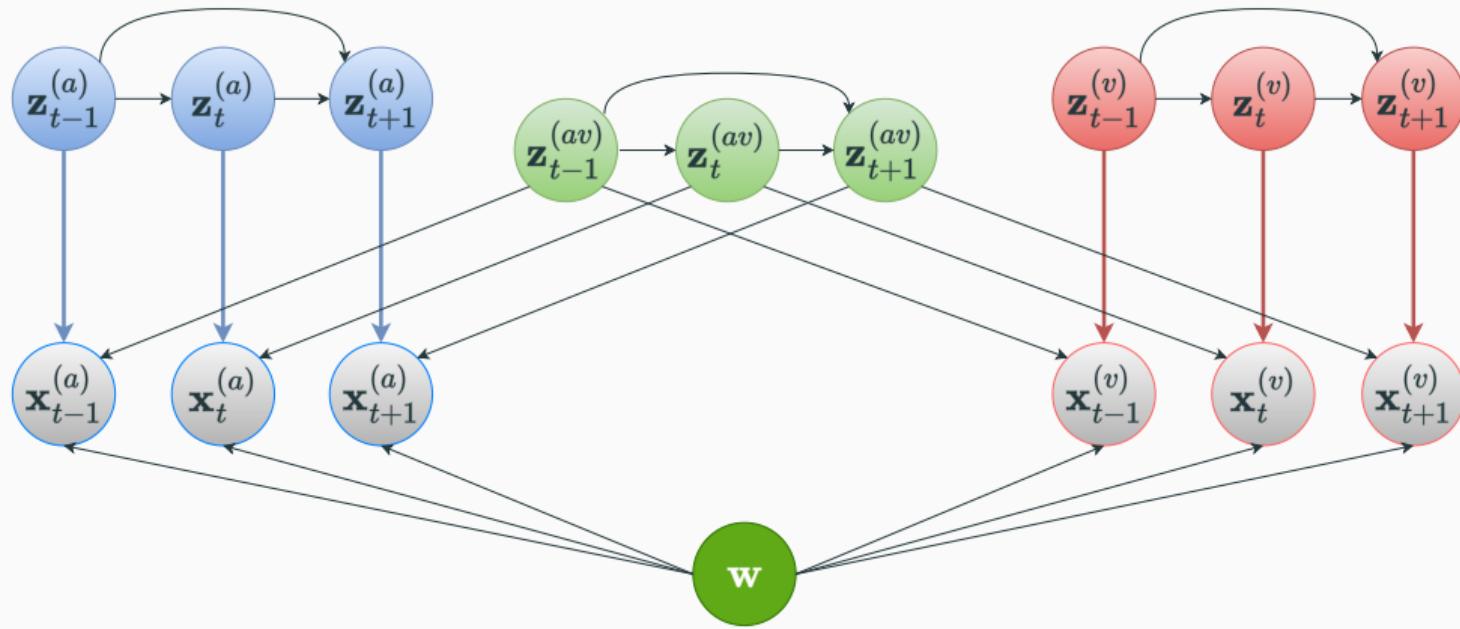


What should we model?

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

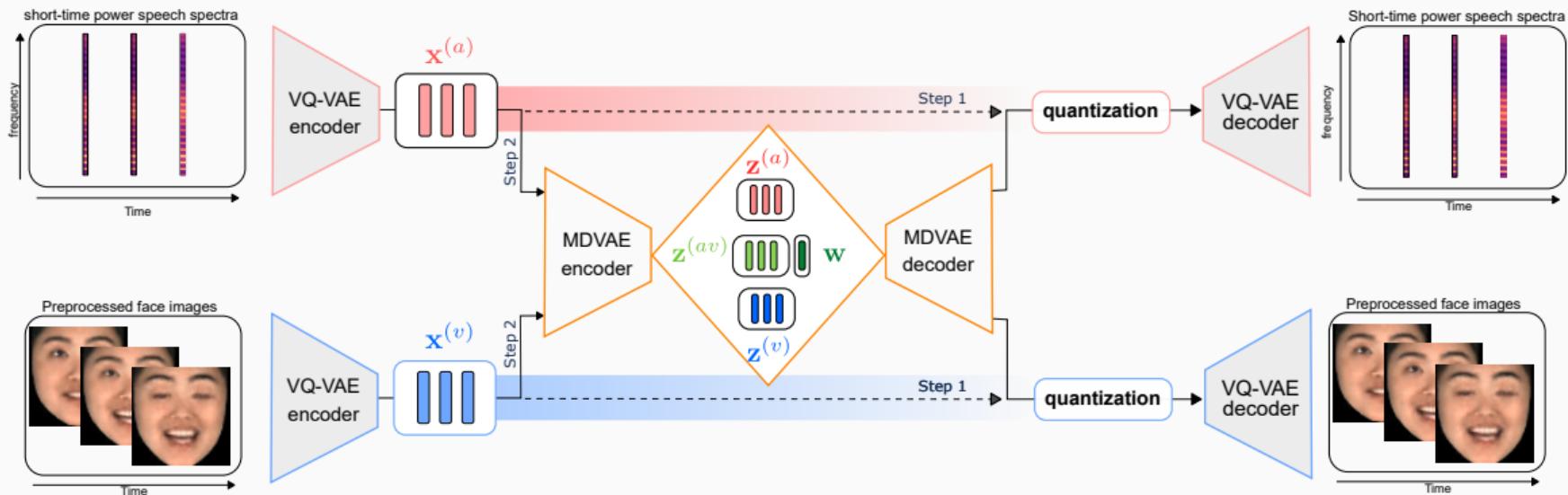
<sup>16</sup>Wang, K., et. al., (2020), ECCV.

# Introducing MDVAEs<sup>17</sup>



<sup>17</sup>Sadok, S., et. al., (2023), Under review 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?

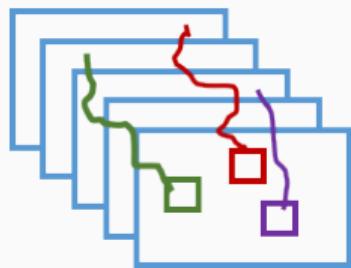


Xiaoyu Lin



Laurent Girin

## Motivating application: unsupervised multiple object tracking



Tracklets @  $t - 1$



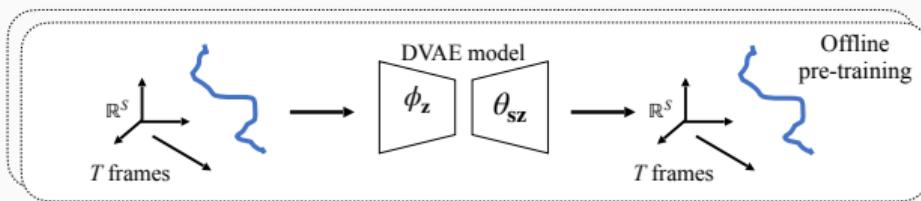
Detections @  $t$



Desired result

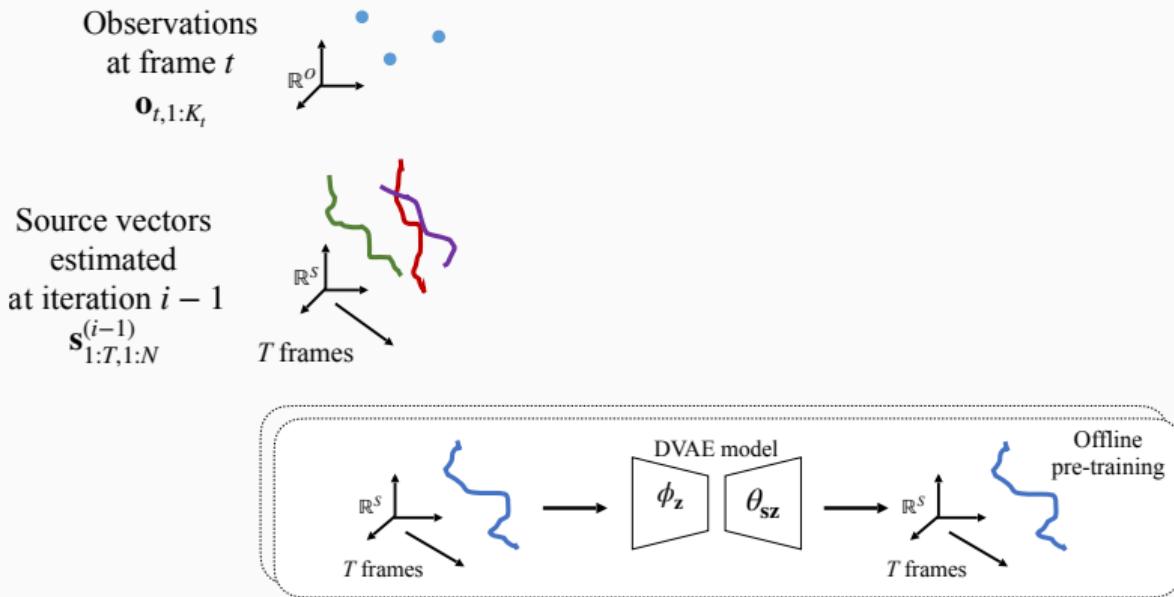
**Question:** Can we model the non-linear dynamics with a DVAE, and have an assignment mechanism within the same ML formulation?

# Introducing MixDVAE<sup>18</sup>



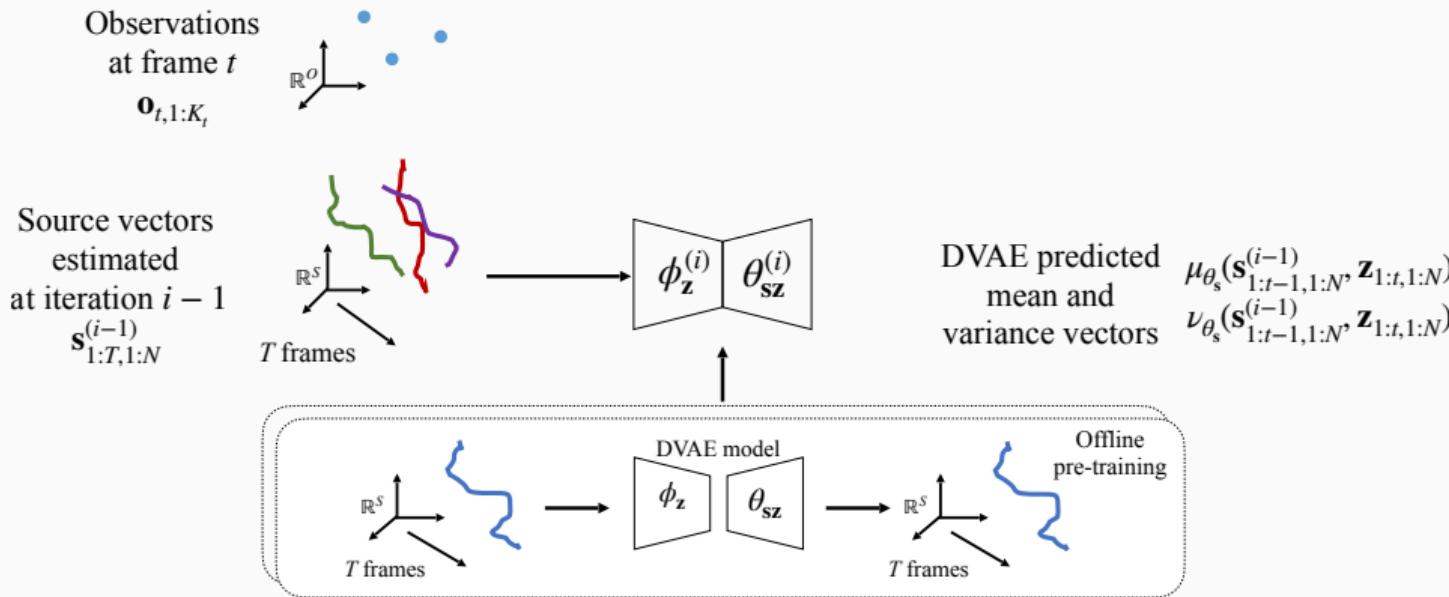
<sup>18</sup>Lin, X., et. al., (2023), Under review.

# Introducing MixDVAE<sup>18</sup>



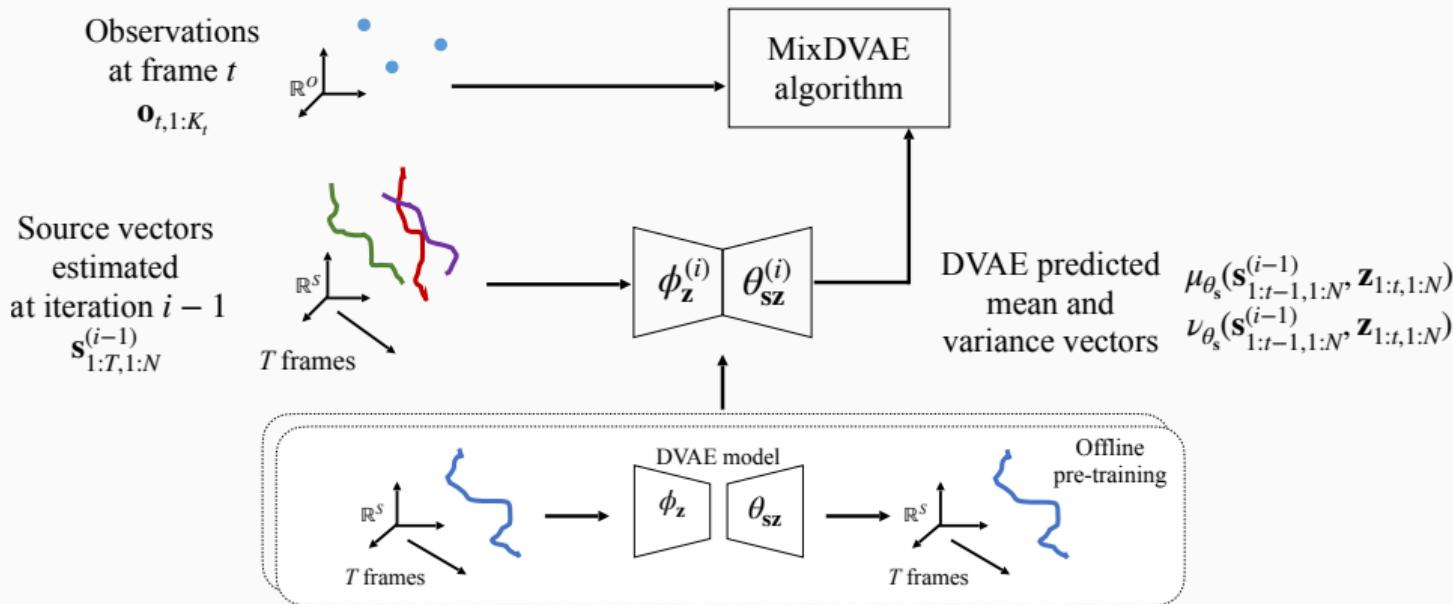
<sup>18</sup>Lin, X., et. al., (2023), Under review.

# Introducing MixDVAE<sup>18</sup>



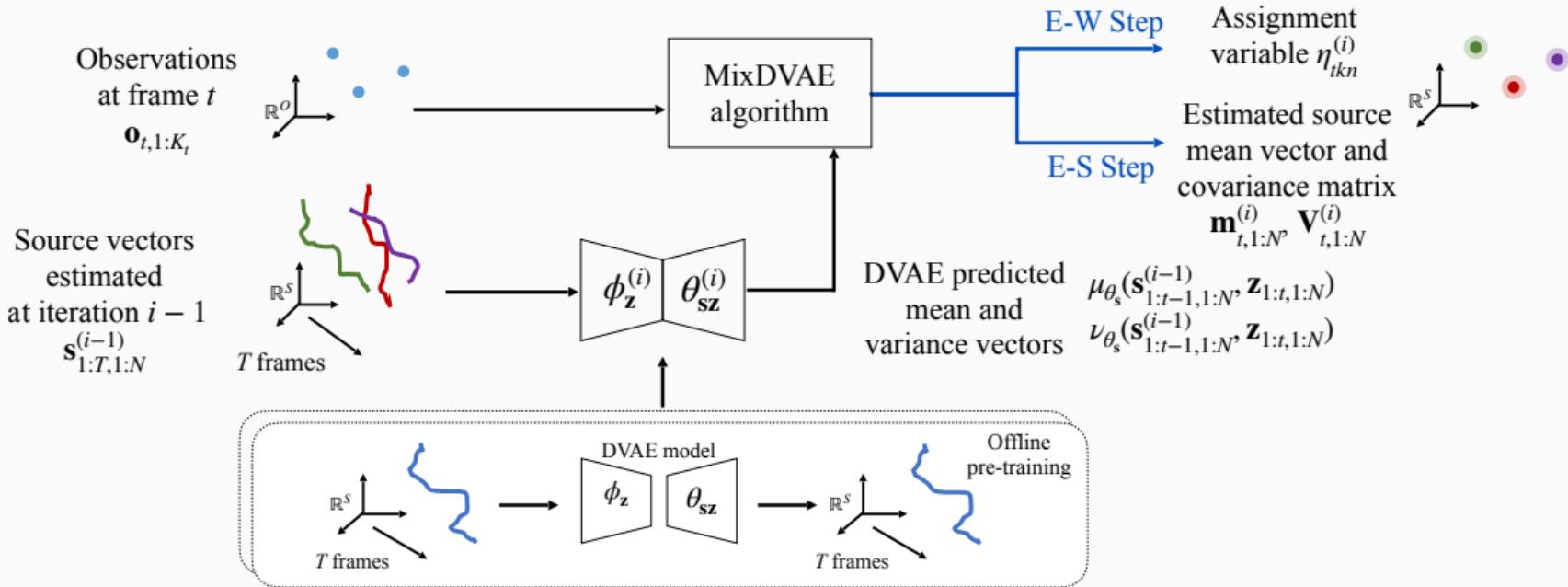
<sup>18</sup>Lin, X., et. al., (2023), Under review.

# Introducing MixDVAE<sup>18</sup>



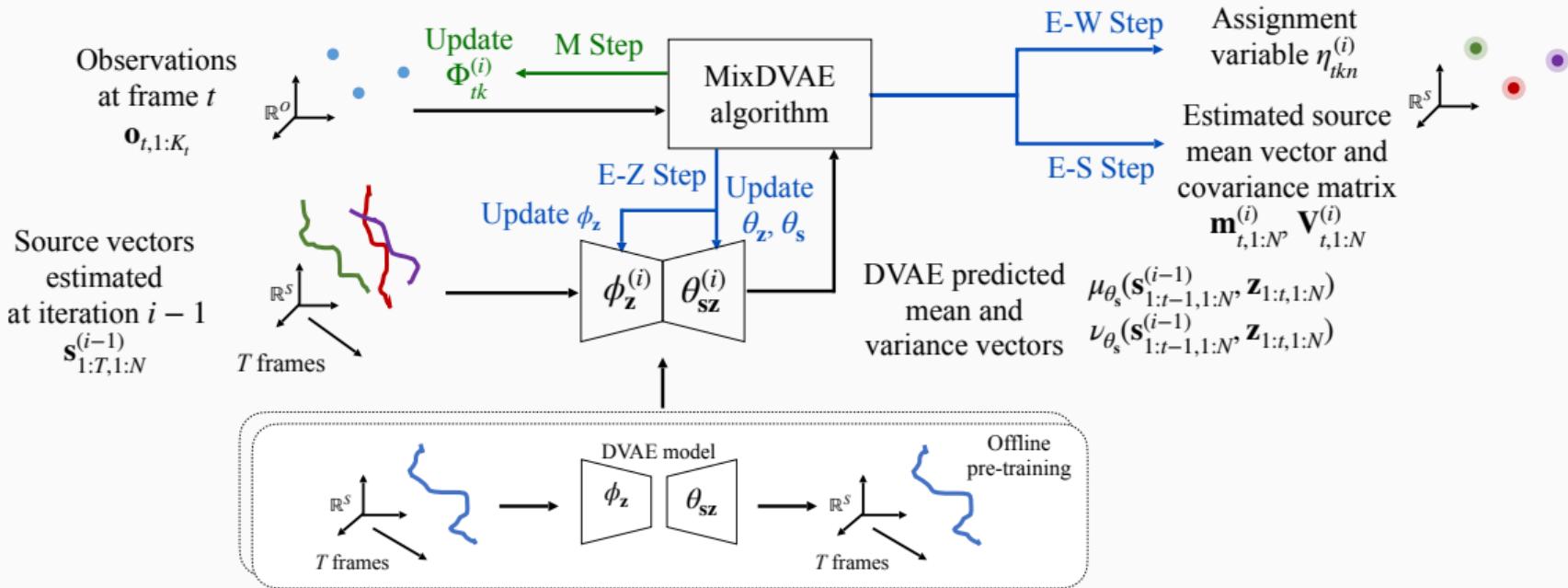
<sup>18</sup>Lin, X., et. al., (2023), Under review.

# Introducing MixDVAE<sup>18</sup>



<sup>18</sup>Lin, X., et. al., (2023), Under review.

# Introducing MixDVAE<sup>18</sup>



<sup>18</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, ...

# Summary

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

# Summary

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

## 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			

---

## Summary

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

## Summary

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

## Summary

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

## Conclusions & Future Work

- Variational inference provides a general framework for multi-modal unsupervised learning.
- It is specially suitable for low-data regimes.
- Recent models (VAE, DVAE, VQ-VAE) can be combined with more classical ones (mixture models, HMM) within the same probabilistic paradigm.
- The probabilistic/maximum likelihood paradigm provides a principled way for disentangling representations and interpreting them.

Open questions:

- Computational complexity is an issue when involving EM/VEM algorithms.
- These methods MUST be extended/rethought for other modalities.
- Understand and develop opportunities with other families (Diffusion, Flows, ...).

## Social Robotics, Artificial Intelligence and Multimedia

Grenoble (FR) 19th-23rd, February, 2024. **No fee!**

Prof. Hatice Gunes, University of Cambridge

Prof. Gabriel Skantze, KTH Stockholm

Prof. Raja Chatila, Sorbonne Université

Prof. Marc Hanheide, University of Lincoln

Prof. Xuesu Xiao, George Mason University

Prof. Antonios Gasteratos, Democritus U. of Thrace

Prof. Wenwu Wang, University of Surrey

Dr. Vasiliki Charisi, JRC European Comission



# 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.