

# Unsupervised Probabilistic Learning with Latent Variables

---

**Lecture 2:** The variational EM algorithm, mixtures of VAEs, application to AV speech enhancement.

MLSS Africa 2023, Cape Town

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

Xavier Alameda-Pineda

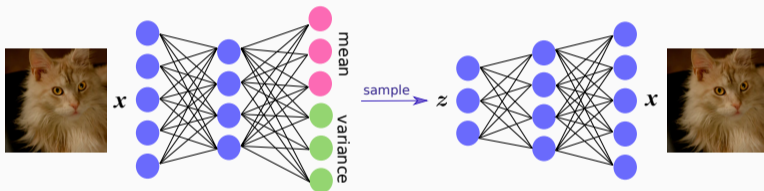
RobotLearn team, Inria at University Grenoble-Alpes,  
Jean-Kuntzman CNRS Laboratory, Multidisciplinary Institute of Artificial Intelligence



## VAE: Summary

**Generative model.** Prior:  $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$  and decoder:  $p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; \mu_{\theta}(\mathbf{z}), \Sigma_{\theta}(\mathbf{z}))$ .

**Inference model (encoder):**  $p_{\theta}(\mathbf{z}|\mathbf{x}) \approx q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \mu_{\phi}(\mathbf{x}), \Sigma_{\phi}(\mathbf{x}))$



**Training criterion** (maximise the evidence lower bound):

$$\mathcal{L}_{\text{ELBO}}(\theta, \phi) = \underbrace{\log p_{\theta}(\mathbf{x}|\hat{\mathbf{z}})}_{\text{Reconstruction}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}))}_{\text{Regularisation}}$$

where  $\hat{\mathbf{z}} \sim q_{\phi}$  is sampled using the reparametrisation trick.

If we recall the formulation for the EM:

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

## Learning - ELBO

If we recall the formulation for the EM:

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

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

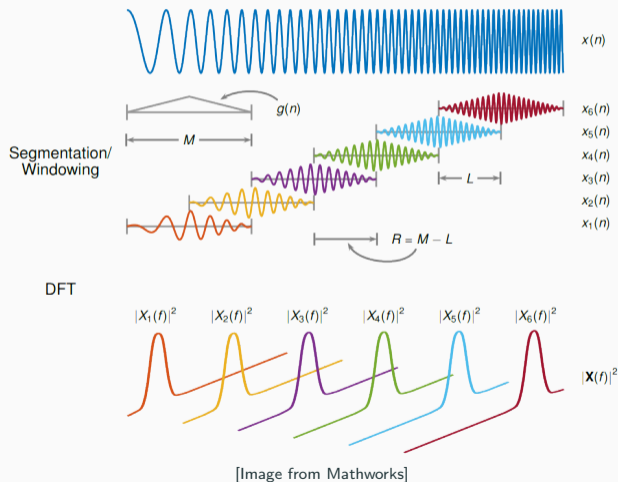
$$\begin{aligned} \log p(\mathbf{x}; \boldsymbol{\theta}, \phi) &\geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \\ \log p(\mathbf{x}; \boldsymbol{\theta}, \phi) &\geq \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}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

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

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

# The short-time Fourier transform (STFT)

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



# Lecture's Outline

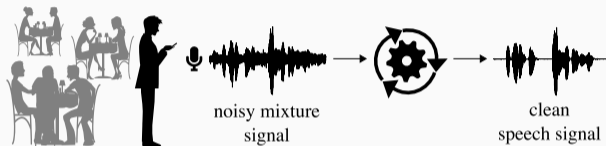
- 1 Unsupervised Audio-visual Speech Enhancement
- 2 Conditional Audio-Visual VAE
- 3 Mixtures of VAEs
- 4 Mixture of Inference Networks VAE

# Unsupervised Audio-visual Speech Enhancement (AV-SE)

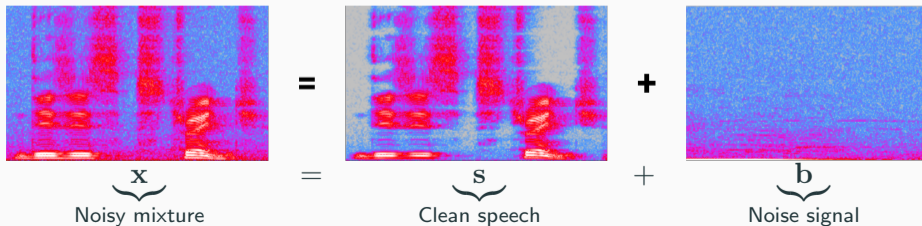
---

# What is Speech Enhancement

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



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





# Audio-visual Speech Enhancement (AV-SE)

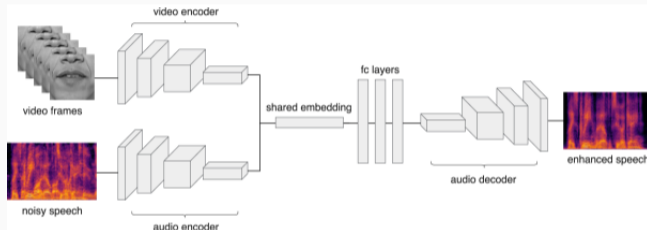
- Visual data, i.e. **lips movements**, provides some complementary information about the unknown speech.
- For **highly noisy audio recordings**, visual information can be very helpful.



*We investigate the use of VAEs to efficiently fuse audio and visual modalities for SE.*

## Supervised vs. Unsupervised AVSE

**Supervised:** Learn a mapping to denoise audio data [Gabbby et al., 2018]:



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

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

## Unsupervised speech enhancement: overview

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



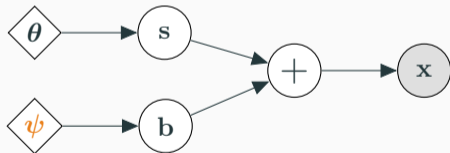
## Unsupervised speech enhancement: overview

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



---

**Test:** learn the **noise parameters** of  $\mathbf{x}$ :

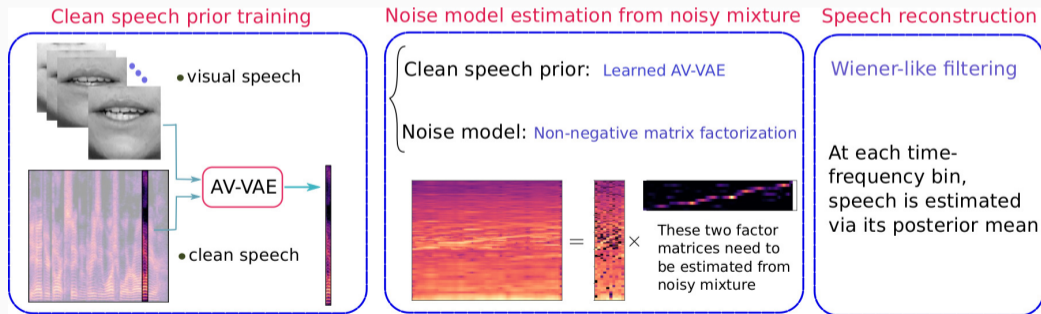


**Note:** at test time,  $\theta$  is frozen, and  $\mathbf{s}$  becomes a latent variable.

# Pre-training, noise estimation & enhancement

Generally speaking our pipeline has three phases:

- 1 Learn a deep generative model of the clean data.
- 2 Estimate the noise parameters from a mixture.
- 3 Enhance the speech signal (Wiener-like filtering).



# Conditional Audio-Visual VAE

---

## Conditional Audio-Visual VAE

---

1. Learn a deep generative model

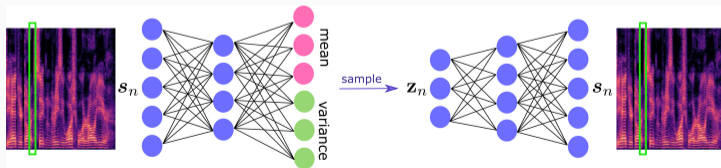
## Generative model (Decoder):

Each clean spectrogram time frame  $s_n$  is assumed to be generated as:

$$s_n | z_n \sim \mathcal{N}_c(\mathbf{0}, \text{diag}(\boldsymbol{\sigma}_s(z_n))), \quad p(z_n) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

## Inference network (Encoder):

$$q(z_n | s_n; \phi) = \mathcal{N}(\boldsymbol{\mu}_z^a(s_n), \text{diag}(\boldsymbol{\sigma}_z^a(s_n)))$$





# Video-only VAE [Sadeghi et al., 2020]

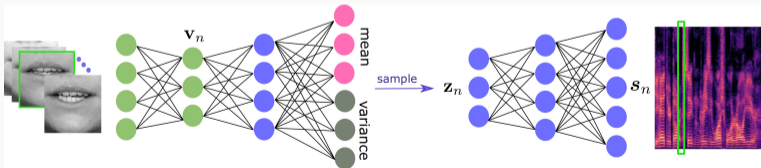
**Generative model (Decoder):**

$$s_n | z_n \sim \mathcal{N}_c(\mathbf{0}, \text{diag}(\boldsymbol{\sigma}_s(z_n))), \quad p(z_n) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

**Inference network (Encoder):**

Infer the posterior using visual data only:

$$q(z_n | \mathbf{v}_n; \phi) = \mathcal{N}(\boldsymbol{\mu}_z^v(\mathbf{v}_n), \text{diag}(\boldsymbol{\sigma}_z^v(\mathbf{v}_n)))$$



▷  $\mathbf{v}_n$  is an embedding for the image of the speaker lips at frame  $n$ .

Inspired by conditional VAE (CVAE), use visual data as some deterministic information

## Generative network (Decoder):

Each clean spectrogram time frame  $s_n$  is assumed to be generated as:

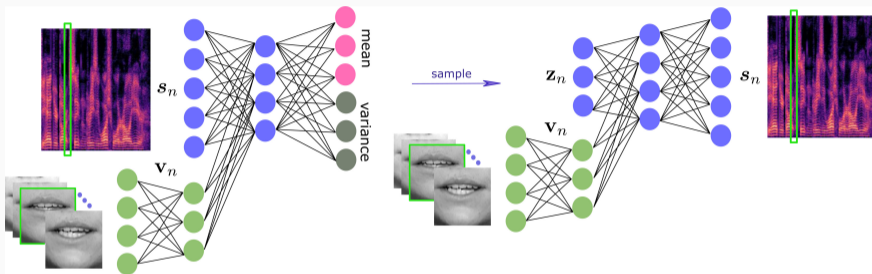
$$\begin{cases} s_n | z_n, v_n & \sim \mathcal{N}_c(\mathbf{0}, \text{diag}(\sigma_s(z_n, v_n))) \\ z_n | v_n & \sim \mathcal{N}(\mu_z(v_n), \text{diag}(\sigma_z(v_n))) \end{cases}$$

## Inference network (Encoder):

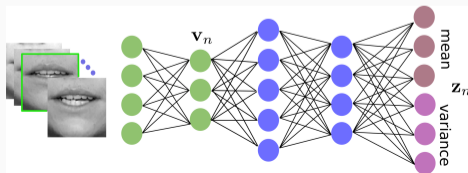
$$q(z_n | s_n, v_n; \phi) = \mathcal{N}(\mu_z^{av}(s_n, v_n), \text{diag}(\sigma_z^{av}(s_n, v_n)))$$

# Audio-visual VAE

## Encoder and decoder networks:



## Prior distribution for latent variables:



## Training samples:

- Clean speech spectrogram time frames:  $\{\mathbf{s}_n\}_{n=0}^{N_{tr}-1}$
- Associated visual data:  $\{\mathbf{v}_n\}_{n=0}^{N_{tr}-1}$

## Optimize the ELBO:

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}) = \frac{1}{N_{tr}} \sum_{n=0}^{N_{tr}-1} \alpha \cdot \mathbb{E}_{p(\mathbf{z}_n|\mathbf{v}_n)} \left[ \ln p(\mathbf{s}_n|\mathbf{z}_n, \mathbf{v}_n; \boldsymbol{\theta}) \right] + \\ (1 - \alpha) \cdot \left( \mathbb{E}_{q(\mathbf{z}_n|\mathbf{s}_n, \mathbf{v}_n; \boldsymbol{\phi})} \left[ \ln p(\mathbf{s}_n|\mathbf{z}_n, \mathbf{v}_n; \boldsymbol{\theta}) \right] - D_{\text{KL}} \left( q(\mathbf{z}_n|\mathbf{s}_n, \mathbf{v}_n; \boldsymbol{\phi}) \parallel p(\mathbf{z}_n|\mathbf{v}_n) \right) \right)$$

- $0 \leq \alpha \leq 1$  gives some reconstruction power to the prior network

## Conditional Audio-Visual VAE

---

2. Estimate the noise parameters from  $x$

# Noise estimation: the model

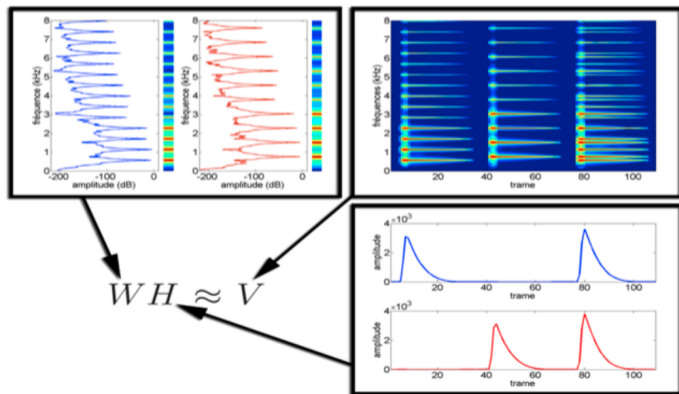
Noisy mixture model:

$$\forall n : \quad \mathbf{x}_n = \mathbf{s}_n + \mathbf{b}_n$$

Noise model:

$$\forall n : \quad \mathbf{b}_n \sim \mathcal{N}_c(\mathbf{0}, \text{diag}(\mathbf{W}_b \mathbf{H}_b[:, n]))$$

Model based on non-negative matrix factorisation:



# Noise estimation: the problem formulation

Noisy mixture model:

$$\forall n : \mathbf{x}_n = \mathbf{s}_n + \mathbf{b}_n$$

Noise model:

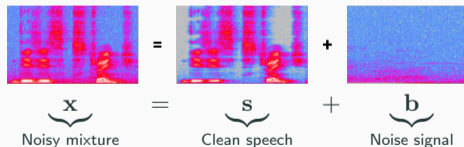
$$\forall n : \mathbf{b}_n \sim \mathcal{N}_c(\mathbf{0}, \text{diag}(\mathbf{W}_b \mathbf{H}_b[:, n]))$$

Clean speech model: Pre-trained generative network:

$$\begin{cases} p(\mathbf{s}_n | \mathbf{z}_n, \mathbf{v}_n) &= \mathcal{N}_c(\mathbf{0}, \text{diag}(\boldsymbol{\sigma}_s(\mathbf{z}_n, \mathbf{v}_n))) \\ p(\mathbf{z}_n | \mathbf{v}_n) &= \mathcal{N}(\boldsymbol{\mu}_z(\mathbf{v}_n), \text{diag}(\boldsymbol{\sigma}_z(\mathbf{v}_n))) \end{cases}$$

Inference:

- ▷ Parameters to be estimated:  $\psi = \{\mathbf{W}_b, \mathbf{H}_b\}$
- ▷ Observed variables:  $\mathbf{x}_n$ 's and  $\mathbf{v}_n$ 's.
- ▷ Latent variables:  $\mathbf{z}_n$ 's and  $\mathbf{s}_n$ 's.



## Noise estimation: the optimisation criterion

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

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



## Noise estimation: the optimisation criterion

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

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

Stop! There is an *easier* way. We can marginalise w.r.t. the  $\mathbf{s}_n$ 's:

$$(p(\mathbf{x}_n, \mathbf{z}_n | \mathbf{v}_n; \psi) = \int_{\mathbf{s}_n} p(\mathbf{x}_n, \mathbf{s}_n, \mathbf{z}_n | \mathbf{v}_n; \psi) d\mathbf{s}_n)$$

$$p(\mathbf{x}_n, \mathbf{z}_n | \mathbf{v}_n; \psi) = \mathcal{N}_c(\mathbf{x}_n; \mathbf{0}, \text{diag}(\mathbf{W}_b \mathbf{H}_b[:, n] + \sigma_s(\mathbf{z}_n, \mathbf{v}_n))) p(\mathbf{z}_n | \mathbf{v}_n)$$

Thus the expected complete-data log-likelihood can be rewritten as:

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

## Noise estimation: the algorithm

One issue (of both  $\mathcal{Q}$ ) is that the posterior is intractable:

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

## Noise estimation: the algorithm

One issue (of both  $Q$ ) is that the posterior is intractable:

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

We propose to approximate the  $Q$  function with Monte-Carlo sampling, leading to a MCEM:

- **E-Step:**  $Q(\boldsymbol{\psi}; \boldsymbol{\psi}^0) \approx \sum_{n=1}^N \frac{1}{R} \sum_{r=1}^R \ln p(\mathbf{x}_n, \mathbf{z}_n^{(r)}, \mathbf{v}_n; \boldsymbol{\psi})$   
The samples  $\{\mathbf{z}_n^{(r)}\}_{r=1, \dots, R}$  are i.i.d. and drawn from  $p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{v}_n; \boldsymbol{\psi}^0)$ .
- **M-Step:**  $\boldsymbol{\psi}^1 \leftarrow \operatorname{argmax}_{\boldsymbol{\psi}} Q(\boldsymbol{\psi}; \boldsymbol{\psi}^0) \Rightarrow$  multiplicative update rules.  
Standard formulae for NMF, not detailed here.

# Conditional Audio-Visual VAE

---

## 3. Speech Enhancement

## Posterior Speech Distribution

Quick reminder: pre-trained decoder  $p(\mathbf{s}_n | \mathbf{z}_n, \mathbf{v}_n)$  and estimated noise parameters  $\psi^* = \{\mathbf{W}_b^*, \mathbf{H}_b^*\}$ . The task is to find the expected value of  $\mathbf{s}_n$  given  $\mathbf{x}_n$  and  $\mathbf{v}_n$ :

$$\hat{\mathbf{s}}_n = \mathbb{E}_{p(\mathbf{s}_n | \mathbf{x}_n, \mathbf{v}_n)}[\mathbf{s}_n] = \int_{\mathbf{s}_n} p(\mathbf{s}_n | \mathbf{x}_n, \mathbf{v}_n) \mathbf{s}_n d\mathbf{s}_n.$$

## Posterior Speech Distribution

Quick reminder: pre-trained decoder  $p(\mathbf{s}_n | \mathbf{z}_n, \mathbf{v}_n)$  and estimated noise parameters  $\boldsymbol{\psi}^* = \{\mathbf{W}_b^*, \mathbf{H}_b^*\}$ . The task is to find the expected value of  $\mathbf{s}_n$  given  $\mathbf{x}_n$  and  $\mathbf{v}_n$ :

$$\hat{\mathbf{s}}_n = \mathbb{E}_{p(\mathbf{s}_n | \mathbf{x}_n, \mathbf{v}_n)}[\mathbf{s}_n] = \int_{\mathbf{s}_n} p(\mathbf{s}_n | \mathbf{x}_n, \mathbf{v}_n) \mathbf{s}_n d\mathbf{s}_n.$$

Problem: there is not close-form for  $p(\mathbf{s}_n | \mathbf{x}_n, \mathbf{v}_n)$  because of  $\mathbf{z}_n$ , thus:

$$\hat{\mathbf{s}}_n = \int_{\mathbf{z}_n} \left( \int_{\mathbf{s}_n} p(\mathbf{s}_n | \mathbf{z}_n, \mathbf{x}_n, \mathbf{v}_n) \mathbf{s}_n d\mathbf{s}_n \right) p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{v}_n) d\mathbf{z}_n.$$

Inner integral: can be computed in closed-form.

Outer integral: approximated with the samples of the MCEM ( $\mathbf{z}_n^{(r)}$ ).

## Posterior Speech Distribution

Quick reminder: pre-trained decoder  $p(\mathbf{s}_n | \mathbf{z}_n, \mathbf{v}_n)$  and estimated noise parameters  $\psi^* = \{\mathbf{W}_b^*, \mathbf{H}_b^*\}$ . The task is to find the expected value of  $\mathbf{s}_n$  given  $\mathbf{x}_n$  and  $\mathbf{v}_n$ :

$$\hat{\mathbf{s}}_n = \mathbb{E}_{p(\mathbf{s}_n | \mathbf{x}_n, \mathbf{v}_n)}[\mathbf{s}_n] = \int_{\mathbf{s}_n} p(\mathbf{s}_n | \mathbf{x}_n, \mathbf{v}_n) \mathbf{s}_n d\mathbf{s}_n.$$

Problem: there is not close-form for  $p(\mathbf{s}_n | \mathbf{x}_n, \mathbf{v}_n)$  because of  $\mathbf{z}_n$ , thus:

$$\hat{\mathbf{s}}_n = \int_{\mathbf{z}_n} \left( \int_{\mathbf{s}_n} p(\mathbf{s}_n | \mathbf{z}_n, \mathbf{x}_n, \mathbf{v}_n) \mathbf{s}_n d\mathbf{s}_n \right) p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{v}_n) d\mathbf{z}_n.$$

Inner integral: can be computed in closed-form.

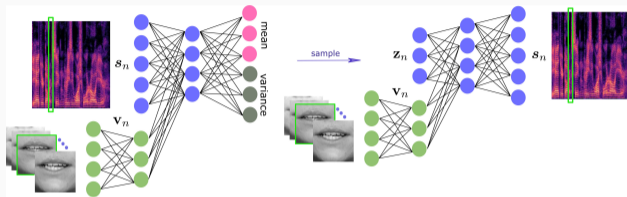
Outer integral: approximated with the samples of the MCEM ( $\mathbf{z}_n^{(r)}$ ).

Clean speech estimated via a **Wiener filter-like** estimator (entry-wise vector operations):

$$\hat{\mathbf{s}}_n = \frac{1}{R} \sum_{r=1}^R \frac{\sigma_s(\mathbf{z}_n^{(r)}, \mathbf{v}_n)}{\sigma_s(\mathbf{z}_n^{(r)}, \mathbf{v}_n) + (\mathbf{W}_b^* \mathbf{H}_b^*)[:, n]} \mathbf{x}_n.$$

# Summary

- 1 Learn a deep generative model of the clean data.



- 2 Estimate the noise parameters from a noisy mixture  $x$  using MCEM:

$$Q(\psi; \psi^0) \approx \sum_{n=1}^N \frac{1}{R} \sum_{r=1}^R \ln p(\mathbf{x}_n, \mathbf{z}_n^{(r)}, \mathbf{v}_n; \psi)$$

- 3 Clean speech estimated via a Wiener filter-like estimator:

$$\hat{\mathbf{s}}_n = \frac{1}{R} \sum_{r=1}^R \frac{\sigma_s(\mathbf{z}_n^{(r)}, \mathbf{v}_n)}{\sigma_s(\mathbf{z}_n^{(r)}, \mathbf{v}_n) + (\mathbf{W}_b^* \mathbf{H}_b^*)[:, n]} \mathbf{x}_n.$$



# Conditional Audio-Visual VAE

---

## Experiments

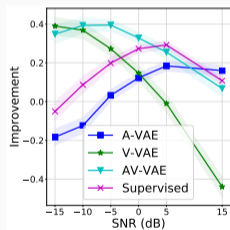
# Settings

- ▷ **NTCD-TIMIT dataset** [Abdelaziz, 2017]
  - Audio-visual recordings in controlled conditions
  - Clean audio as well as noisy versions
  - Frontal video frames with 30 FPS-  $67 \times 67$  lips images
- ▷ **Training set** ( $\sim 5$  hours): 39 speakers  $\times$  98 sentences  $\times$  5 seconds
- ▷ **Test set** ( $\sim 1$  hour): 9 speakers  $\times$  98 sentences  $\times$  5 seconds
- ▷ **Noise levels**:  $-15$  dB,  $-10$  dB,  $-5$  dB, 0 dB, 5 dB and 15 dB
- ▷ **Noise types**: *Living Room (LR), White, Cafe, Car, Babble, and Street*

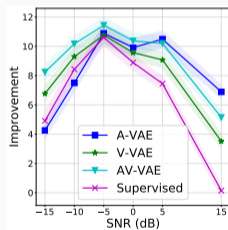
# Results

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

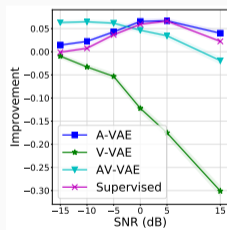
- Perceptual evaluation of speech quality (PESQ).
- Signal-to-distortion ratio (SDR).
- Short-time objective intelligibility (STOI).



(a) PESQ



(b) SDR

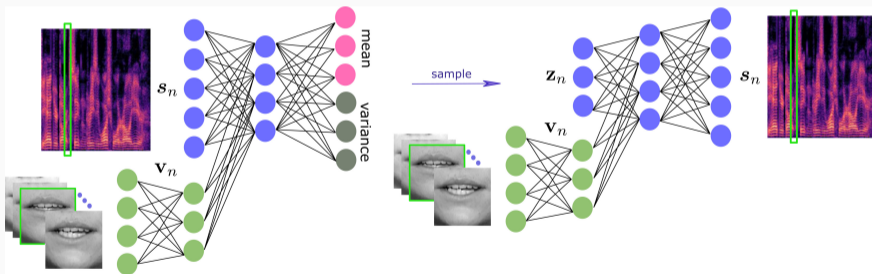


(c) STOI

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

# Systematic AV fusion

The conditional AV-VAE uses ALWAYS visual information.  
This can be a problem when dealing with noisy visual data.



## Mixtures of VAEs

---

# Motivation

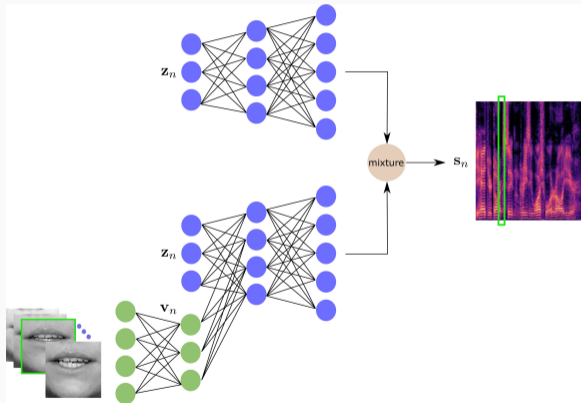
AV-VAE usually yields better results than A-VAE, especially at low SNRs, provided **clean (frontal, non-occluded)** visual data [Sadeghi et al., 2019].



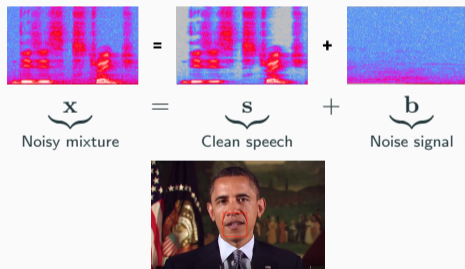
*How to effectively benefit from AV-VAE for in-the-wild video recordings?*

# VAE mixture model (VAE-MM) [Sadeghi et al., 2020b]

A **mixture** of A-VAE plus AV-VAE generative model: use visual information only if clean.



After trained, used for unsupervised AV-SE.



## Mixtures of VAEs

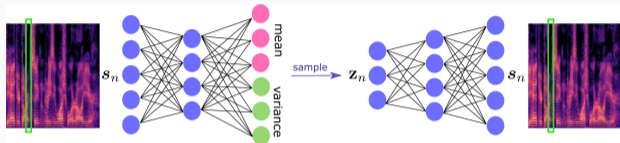
---

1. Learn a deep generative model

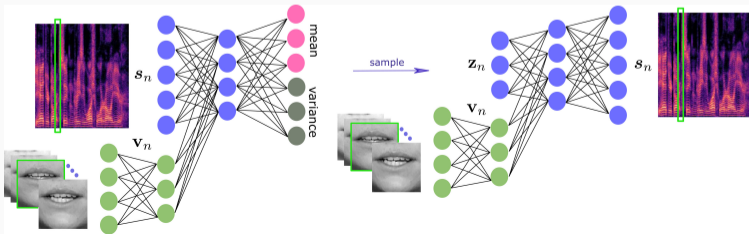


## Pre-training as in the previous model

**Audio-only VAE:**  $s_n | z_n \sim \mathcal{N}(\mathbf{0}, \text{diag}(\sigma_s(z_n)))$  and  $z_n \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$



**Conditional AV VAE:**  $s_n | z_n, v_n \sim \mathcal{N}(\mathbf{0}, \text{diag}(\sigma_s(z_n, v_n)))$  and  $z_n | v_n \sim \mathcal{N}(\mu_z(v_n), \text{diag}(\sigma_z(v_n)))$



**VAE-MM clean speech model:** Combine **A-VAE** with **AV-VAE**:

$$\begin{cases} p(\mathbf{s}_n | \mathbf{z}_n, \mathbf{v}_n, \alpha_n) &= \left[ \mathcal{N}_c(\mathbf{0}, \text{diag}(\boldsymbol{\sigma}_s^a(\mathbf{z}_n))) \right]^{\alpha_n} \times \left[ \mathcal{N}_c(\mathbf{0}, \text{diag}(\boldsymbol{\sigma}_s^{av}(\mathbf{z}_n, \mathbf{v}_n))) \right]^{1-\alpha_n}, \\ p(\mathbf{z}_n | \mathbf{v}_n, \alpha_n) &= \left[ \mathcal{N}(\mathbf{0}, \mathbf{I}) \right]^{\alpha_n} \times \left[ \mathcal{N}(\boldsymbol{\mu}_z^v(\mathbf{v}_n), \text{diag}(\boldsymbol{\sigma}_z^v(\mathbf{v}_n))) \right]^{1-\alpha_n}, \\ p(\alpha_n) &= \pi^{\alpha_n} \times (1 - \pi)^{1-\alpha_n}. \end{cases}$$

$\alpha_n \in \{0, 1\}$  is a latent variable selecting the VAE used in the  $n$ -th frame.

## Mixtures of VAEs

---

2. & 3. Estimate the noise parameters  
from  $x$

# Noise model and optimisation problem

**Noisy mixture model:**

$$\forall n \quad \mathbf{x}_n = \mathbf{s}_n + \mathbf{b}_n$$

**Noise model:**

$$\forall n \quad \mathbf{b}_n \sim \mathcal{N}_c(\mathbf{0}, \text{diag}(\mathbf{W}_b \mathbf{H}_b[:, n]))$$

**Clean speech model:**

Mixture of A-VAE and AV-VAE.

**Inference:**

- ▷ Observed variables:  $\{\mathbf{x}_n, \mathbf{v}_n\}_{n=1}^N$
- ▷ Latent variables:  $\{\mathbf{s}_n, \mathbf{z}_n, \alpha_n\}_{n=1}^N$
- ▷ Parameters to be estimated:  $\psi = \{\mathbf{W}_b, \mathbf{H}_b, \pi\}$

# Noise model and optimisation problem

**Noisy mixture model:**

$$\forall n \quad \mathbf{x}_n = \mathbf{s}_n + \mathbf{b}_n$$

**Noise model:**

$$\forall n \quad \mathbf{b}_n \sim \mathcal{N}_c(\mathbf{0}, \text{diag}(\mathbf{W}_b \mathbf{H}_b[:, n]))$$

**Clean speech model:**

Mixture of A-VAE and AV-VAE.

**Inference:**

- ▷ Observed variables:  $\{\mathbf{x}_n, \mathbf{v}_n\}_{n=1}^N$
- ▷ Latent variables:  $\{\mathbf{s}_n, \mathbf{z}_n, \alpha_n\}_{n=1}^N$
- ▷ Parameters to be estimated:  $\psi = \{\mathbf{W}_b, \mathbf{H}_b, \pi\}$

**Ideally:**

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

## Approximating the posterior distribution

**Challenge:** The posterior is intractable, we propose a variational approximation:  
(in this case, marginalisation w.r.t.  $s_n$  does not simplify things)

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

## Approximating the posterior distribution

**Challenge:** The posterior is intractable, we propose a variational approximation:  
(in this case, marginalisation w.r.t.  $s_n$  does not simplify things)

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

The optimal values for the variational distributions [Bishop, 2006] are obtained with:

**VE  $s_n$ -step:**  $q_{\mathbf{s}}(\mathbf{s}_n) \propto \exp \left( \mathbb{E}_{q_{\mathbf{z}}(\mathbf{z}_n) q_{\alpha}(\alpha_n)} \left[ \log p(\mathbf{x}_n, \mathbf{s}_n, \mathbf{z}_n, \alpha_n | \mathbf{v}_n; \psi^0) \right] \right)$

**VE  $\alpha_n$ -step:**  $q_{\alpha}(\alpha_n) \propto \exp \left( \mathbb{E}_{q_{\mathbf{s}}(\mathbf{s}_n) q_{\mathbf{z}}(\mathbf{z}_n)} \left[ \log p(\mathbf{x}_n, \mathbf{s}_n, \mathbf{z}_n, \alpha_n | \mathbf{v}_n; \psi^0) \right] \right)$

**VE  $z_n$ -step:**  $q_{\mathbf{z}}(\mathbf{z}_n) \propto \exp \left( \mathbb{E}_{q_{\mathbf{s}}(\mathbf{s}_n) q_{\alpha}(\alpha_n)} \left[ \log p(\mathbf{x}_n, \mathbf{s}_n, \mathbf{z}_n, \alpha_n | \mathbf{v}_n; \psi^0) \right] \right)$

Do we have any hope for  $z_n$ ?

## Approximating the posterior distribution

**Challenge:** The posterior is intractable, we propose a variational approximation:  
(in this case, marginalisation w.r.t.  $s_n$  does not simplify things)

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

The optimal values for the variational distributions [Bishop, 2006] are obtained with:

**VE  $s_n$ -step:**  $q_{\mathbf{s}}(\mathbf{s}_n) \propto \exp \left( \mathbb{E}_{q_{\mathbf{z}}(\mathbf{z}_n) q_{\alpha}(\alpha_n)} \left[ \log p(\mathbf{x}_n, \mathbf{s}_n, \mathbf{z}_n, \alpha_n | \mathbf{v}_n; \psi^0) \right] \right)$

**VE  $\alpha_n$ -step:**  $q_{\alpha}(\alpha_n) \propto \exp \left( \mathbb{E}_{q_{\mathbf{s}}(\mathbf{s}_n) q_{\mathbf{z}}(\mathbf{z}_n)} \left[ \log p(\mathbf{x}_n, \mathbf{s}_n, \mathbf{z}_n, \alpha_n | \mathbf{v}_n; \psi^0) \right] \right)$

**VE  $z_n$ -step:**  $q_{\mathbf{z}}(\mathbf{z}_n) \propto \exp \left( \mathbb{E}_{q_{\mathbf{s}}(\mathbf{s}_n) q_{\alpha}(\alpha_n)} \left[ \log p(\mathbf{x}_n, \mathbf{s}_n, \mathbf{z}_n, \alpha_n | \mathbf{v}_n; \psi^0) \right] \right)$

Do we have any hope for  $z_n$ ? Clearly not!  $\rightarrow$  Sampling (from  $q_{\mathbf{z}}(\mathbf{z}_n)$ ):  $\mathbf{z}_n^{(r)}$ .



## Variational E-steps

**VE  $\alpha_n$ -step:** It follows that  $q_\alpha$  is a Bernoulli distribution with parameter:

( $g(x) = 1/(1 + \exp(-x))$ ) denotes the sigmoid function)

$$\pi_n = g\left(\frac{1}{R} \sum_{r=1}^R \mathbb{E}_{q_s(\mathbf{s}_n)} \left[ \log \frac{p(\mathbf{s}_n, \mathbf{z}_n^{(r)} | \alpha_n = 1)}{p(\mathbf{s}_n, \mathbf{z}_n^{(r)} | \mathbf{v}_n, \alpha_n = 0)} \right] + \log \frac{\pi}{1 - \pi}\right)$$

Audio (numerator) vs. audio-visual (denominator).

## Variational E-steps

**VE  $\alpha_n$ -step:** It follows that  $q_\alpha$  is a Bernoulli distribution with parameter:

( $g(x) = 1/(1 + \exp(-x))$ ) denotes the sigmoid function)

$$\pi_n = g\left(\frac{1}{R} \sum_{r=1}^R \mathbb{E}_{q_s(\mathbf{s}_n)} \left[ \log \frac{p(\mathbf{s}_n, \mathbf{z}_n^{(r)} | \alpha_n = 1)}{p(\mathbf{s}_n, \mathbf{z}_n^{(r)} | \mathbf{v}_n, \alpha_n = 0)} \right] + \log \frac{\pi}{1 - \pi}\right)$$

Audio (numerator) vs. audio-visual (denominator).

**VE  $s_n$ -step** (speech estimation),  $q_s$  is a complex Gaussian with:

$$\hat{\mathbf{s}}_n = \frac{1}{R} \sum_{r=1}^R \frac{\gamma_n(\mathbf{z}_n^{(r)}, \mathbf{v}_n)}{\gamma_n(\mathbf{z}_n^{(r)}, \mathbf{v}_n) + \mathbf{W}_b^* \mathbf{H}_b^*[:, n]} \mathbf{x}_n \quad \left( \gamma_n(\mathbf{z}_n^{(r)}, \mathbf{v}_n) \leftrightarrow \sigma_s(\mathbf{z}_n^{(r)}, \mathbf{v}_n) \right)$$

## Variational E-steps

**VE  $\alpha_n$ -step:** It follows that  $q_\alpha$  is a Bernoulli distribution with parameter:

( $g(x) = 1/(1 + \exp(-x))$ ) denotes the sigmoid function)

$$\pi_n = g\left(\frac{1}{R} \sum_{r=1}^R \mathbb{E}_{q_s(\mathbf{s}_n)} \left[ \log \frac{p(\mathbf{s}_n, \mathbf{z}_n^{(r)} | \alpha_n = 1)}{p(\mathbf{s}_n, \mathbf{z}_n^{(r)} | \mathbf{v}_n, \alpha_n = 0)} \right] + \log \frac{\pi}{1 - \pi}\right)$$

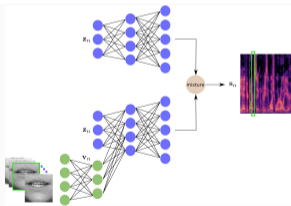
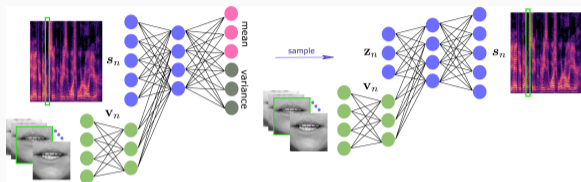
Audio (numerator) vs. audio-visual (denominator).

**VE  $s_n$ -step** (speech estimation),  $q_s$  is a complex Gaussian with:

$$\hat{\mathbf{s}}_n = \frac{1}{R} \sum_{r=1}^R \frac{\gamma_n(\mathbf{z}_n^{(r)}, \mathbf{v}_n)}{\gamma_n(\mathbf{z}_n^{(r)}, \mathbf{v}_n) + \mathbf{W}_b^* \mathbf{H}_b^*[:, n]} \mathbf{x}_n \quad \left( \gamma_n(\mathbf{z}_n^{(r)}, \mathbf{v}_n) \leftrightarrow \sigma_s(\mathbf{z}_n^{(r)}, \mathbf{v}_n) \right)$$

with  $\gamma_n(\mathbf{z}_n, \mathbf{v}_n) = \left( \pi_n \frac{1}{\sigma_s^a(\mathbf{z}_n)} + (1 - \pi_n) \frac{1}{\sigma_s^{av}(\mathbf{z}_n, \mathbf{v}_n)} \right)^{-1}$ .

## Summary so far



### Conditional VAE:

- Training VAE via SGD.
- Learning noise parameters via MCEM.

### VAE-MM:

- Training two VAEs via SGD.
- Learning noise parameters via VEM.

# Mixtures of VAEs

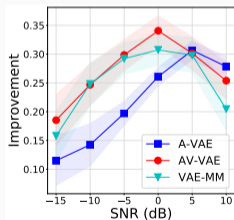
---

## Experiments

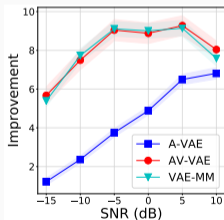
# Settings & Results

- **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).

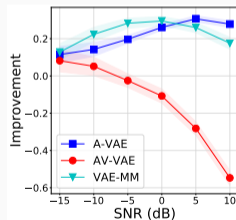
Improvement with respect to the input:



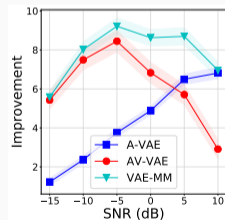
(a) PESQ (clean)



(b) SDR (clean)



(c) PESQ (noisy)

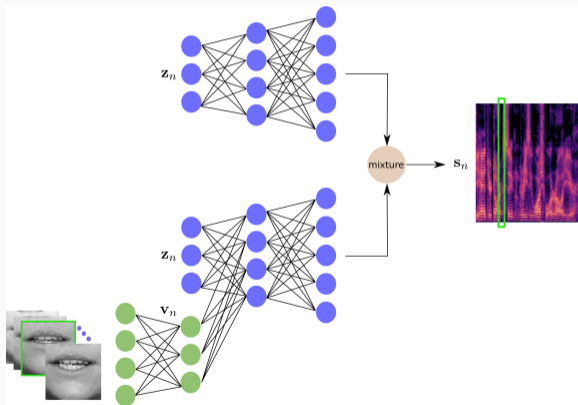


(d) SDR (noisy)

## Two different models for the clean speech!

VAE-MM can effectively choose when to use visual information.

However, it has two different probabilistic models for the SAME speech signal!

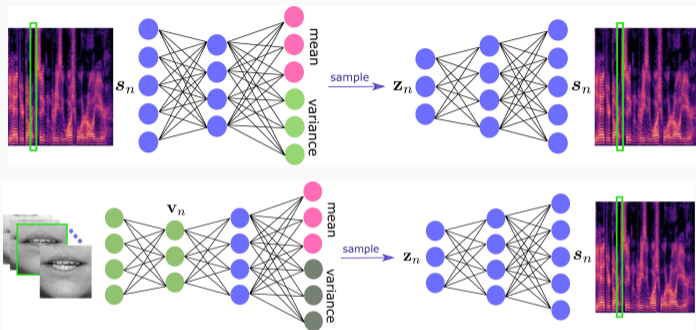


# Mixture of Inference Networks VAE

---



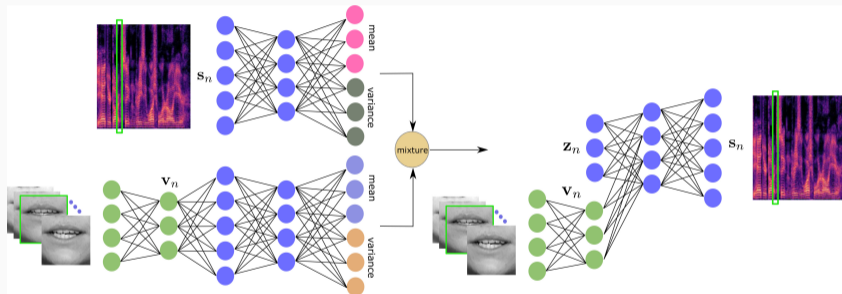
# Motivation



How to exploit the mixture model mechanism while having a single generative model for the clean speech signal?

# Mixture of Inference Networks VAE (MIN-VAE) [Sadeghi et al., 2021]

Train a mixture of audio and visual inference networks:



- The shared generative model (decoder) is trained using both audio and visual latent codes
- Once trained, the latent codes at test phase can be initialized using the visual encoder

# Mixture of Inference Networks: Generative Model

The generative model for each  $s_n$ :

$$s_n | z_n, v_n \sim \mathcal{N}_c(\mathbf{0}, \text{diag}(\sigma_s(z_n, v_n)))$$
$$z_n | \alpha_n \sim \left[ \mathcal{N}(\boldsymbol{\mu}_a, \sigma_a \mathbf{I}) \right]^{\alpha_n} \cdot \left[ \mathcal{N}(\boldsymbol{\mu}_v, \sigma_v \mathbf{I}) \right]^{1-\alpha_n}$$
$$\alpha_n \sim \pi^{\alpha_n} \times (1 - \pi)^{1-\alpha_n}$$

- Each latent code,  $z_n$ , is generated either from an **audio** or from a **video** prior,
- The audio and video priors are parametrized by  $(\boldsymbol{\mu}_a, \sigma_a)$  and  $(\boldsymbol{\mu}_v, \sigma_v)$ ,
- A **mixing latent variable**  $\alpha_n \in \{0, 1\}$  describes whether  $z_n$  corresponds to the **audio** or to the **video** prior.

## Mixture of Inference Networks VAE

---

1. Learn a deep generative model

# Training MIN-VAE

- ▷ Observed variables:  $\mathbf{s} = \{\mathbf{s}_n\}_{n=1}^{N_{tr}}$  and  $\mathbf{v} = \{\mathbf{v}_n\}_{n=1}^{N_{tr}}$
  - ▷ Latent variables:  $\mathbf{z} = \{\mathbf{z}_n\}_{n=1}^{N_{tr}}$  and  $\boldsymbol{\alpha} = \{\alpha_n\}_{n=1}^{N_{tr}}$
  - ▷ Parameters to be estimated:  $\underbrace{\boldsymbol{\theta}}_{\text{decoder}}, \underbrace{\phi_a, \phi_v}_{\text{encoders}}, \underbrace{\mu_a, \sigma_a}_{\text{audio prior}}, \underbrace{\mu_v, \sigma_v}_{\text{visual prior}}, \underbrace{\pi}_{\text{mixing prior}}$
- 

We target a lower bound on the data log-likelihood:

$$\sum_{n=1}^N \mathbb{E}_{p(\mathbf{z}_n, \alpha_n | \mathbf{s}_n, \mathbf{v}_n)} [\log p(\mathbf{s}_n, \mathbf{z}_n, \alpha_n | \mathbf{v}_n)]$$

What is going to happen to the posterior  $p(\mathbf{z}_n, \alpha_n | \mathbf{s}_n, \mathbf{v}_n)$ ?

## Training MIN-VAE: approximating the posterior

Posterior distribution of the latent variables:

$$p(\mathbf{z}_n, \alpha_n | \mathbf{s}_n, \mathbf{v}_n) = \underbrace{p(\mathbf{z}_n | \mathbf{s}_n, \mathbf{v}_n, \alpha_n)}_{\text{Code posterior}} \cdot \underbrace{p(\alpha_n | \mathbf{s}_n, \mathbf{v}_n)}_{\text{Mixing posterior}}.$$

The two posteriors in the RHS are intractable.

# Training MIN-VAE: approximating the posterior

Posterior distribution of the latent variables:

$$p(\mathbf{z}_n, \alpha_n | \mathbf{s}_n, \mathbf{v}_n) = \underbrace{p(\mathbf{z}_n | \mathbf{s}_n, \mathbf{v}_n, \alpha_n)}_{\text{Code posterior}} \cdot \underbrace{p(\alpha_n | \mathbf{s}_n, \mathbf{v}_n)}_{\text{Mixing posterior}}.$$

The two posteriors in the RHS are intractable.

We assume the following variational approximation for the codes:

$$q(\mathbf{z}_n | \mathbf{s}_n, \mathbf{v}_n, \alpha_n; \phi) = \begin{cases} q(\mathbf{z}_n | \mathbf{s}_n; \phi_a) & \alpha_n = 1 \quad (\text{audio encoder}) \\ q(\mathbf{z}_n | \mathbf{v}_n; \phi_v) & \alpha_n = 0 \quad (\text{video encoder}) \end{cases}$$

A variational distribution, denoted  $q(\alpha_n)$ , is considered for  $p(\alpha_n | \mathbf{s}_n, \mathbf{v}_n)$ .

Final approximation:

$$p(\mathbf{z}_n, \alpha_n | \mathbf{s}_n, \mathbf{v}_n) \approx q(\mathbf{z}_n | \mathbf{s}_n, \mathbf{v}_n, \alpha_n; \phi) \cdot q(\alpha_n)$$

## Training MIN-VAE: the algorithm

We can obtain a variational EM (details omitted):

- E-step:  $q(\alpha_n)$  is a Bernoulli distribution with parameter:

$$\pi_n = \pi \cdot g\left(J_n(\alpha_n = 1) - J_n(\alpha_n = 0)\right),$$

where  $J_n(\alpha_n)$  is the VAE cost of representing the  $n$ -th vector with audio ( $\alpha_n = 1$ ) or video ( $\alpha_n = 0$ ) encoding ( $g$  is the sigmoid).



## Training MIN-VAE: the algorithm

We can obtain a variational EM (details omitted):

- E-step:  $q(\alpha_n)$  is a Bernoulli distribution with parameter:

$$\pi_n = \pi \cdot g\left(J_n(\alpha_n = 1) - J_n(\alpha_n = 0)\right),$$

where  $J_n(\alpha_n)$  is the VAE cost of representing the  $n$ -th vector with audio ( $\alpha_n = 1$ ) or video ( $\alpha_n = 0$ ) encoding ( $g$  is the sigmoid).

- M-step: the parameters are updated by maximizing:

$$\sum_{n=1}^{N_{tr}} \mathbb{E}_{q(\alpha_n)} \left[ J_n(\alpha_n) \right] - D_{\text{KL}}\left(q(\alpha_n) \parallel p(\alpha_n)\right)$$

$$\pi = \frac{1}{N_{tr}} \sum_{n=1}^{N_{tr}} \pi_n \text{ and SGD for } \theta, \phi_a, \phi_v, \mu_a, \mu_v, \sigma_a, \sigma_v.$$

## Mixture of Inference Networks VAE

---

2. & 3. Estimate the noise parameters  
from  $x$

# Parameter and Speech Estimation

Noisy speech model:  $\forall n : \mathbf{x}_n = \mathbf{s}_n + \mathbf{b}_n$

Noise model:  $\forall n : \mathbf{b}_n \sim \mathcal{N}_c(\mathbf{0}, \text{diag}(\mathbf{W}_b \mathbf{H}_b[:, n]))$

Clean speech model: Trained MIN-VAE

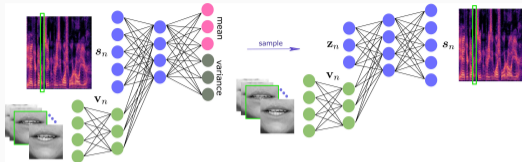
$$\begin{aligned} \mathbf{s}_n | \mathbf{z}_n, \mathbf{v}_n &\sim \mathcal{N}_c(\mathbf{0}, \text{diag}(\boldsymbol{\sigma}_s(\mathbf{z}_n, \mathbf{v}_n))) \\ \mathbf{z}_n | \alpha_n &\sim \left[ \mathcal{N}(\boldsymbol{\mu}_a, \sigma_a \mathbf{I}) \right]^{\alpha_n} \cdot \left[ \mathcal{N}(\boldsymbol{\mu}_v, \sigma_v \mathbf{I}) \right]^{1-\alpha_n} \\ \alpha_n &\sim \pi^{\alpha_n} \times (1 - \pi)^{1-\alpha_n} \end{aligned}$$

---

**Inference:** (we skip the VEM as it follows a similar principle than the previous one)

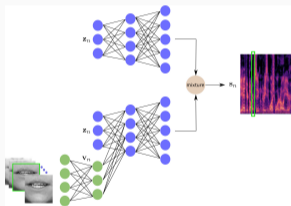
- ▷ Observed variables:  $\{\mathbf{x}_n, \mathbf{v}_n\}_{n=1}^N$
- ▷ Latent variables:  $\{\mathbf{s}_n, \mathbf{z}_n, \alpha_n\}_{n=1}^N$
- ▷ Parameters to be estimated:  $\boldsymbol{\psi} = \{\mathbf{W}_b, \mathbf{H}_b, \pi\}$

# Summary of the three models



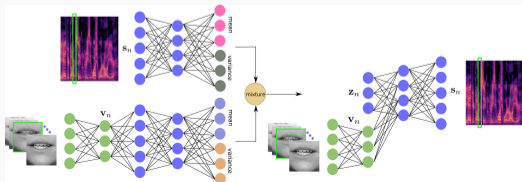
## Conditional VAE:

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



## VAE-MM:

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



## MIN-VAE:

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

# Mixture of Inference Networks VAE

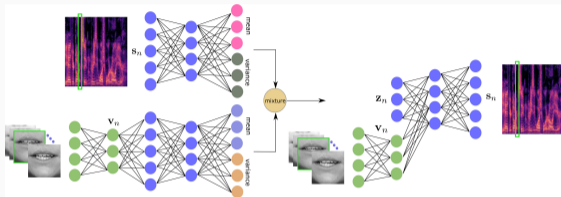
---

## Experiments

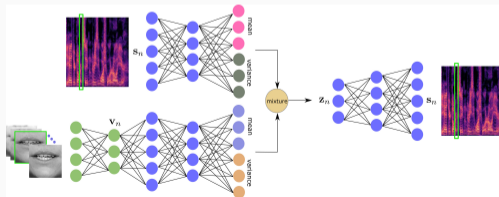
# Settings

Settings similar as before with the NTCD-TIMIT dataset. Two possible models.

MIN-VAE-v1:

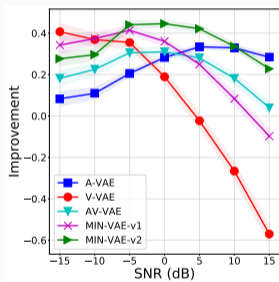


MIN-VAE-v2:

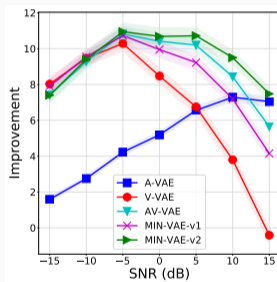


# Results

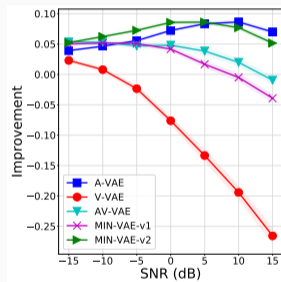
Adding noise to one-third of spectrograms input to audio-encoder:



(a) PESQ



(b) SDR



(c) STOI

## Conclusion

- Variational autoencoders are an excellent framework for unsupervised audio-visual SE.
- However, (V)EM-based inference may be slow. Re-using the trained encoders is an alternative.
- Extending the mixture VAE to more than two encoders, to include other sources of information.
- Huge limitation so far: STFT speech frames are processed independently (not sequentially).

---

**Lecture 3:** Dynamical variational autoencoders (VAE's that can process sequences).



## References

- 1 D. P. Kingma and M. Welling, "Auto-encoding variational Bayes," ICLR, 2014.
- 2 Y. Bando et al., "Statistical speech enhancement based on probabilistic integration of variational autoencoder and non-negative matrix factorization," in Proc. ICASSP, 2018.
- 3 S. Leglaive et al., "A variance modeling framework based on variational autoencoders for speech enhancement," in Proc. MLSP, 2018.
- 4 M. Sadeghi et al., "Audio-visual Speech Enhancement Using Conditional Variational Auto-Encoder," IEEE TASLP, 2020.
- 5 M. Sadeghi and X. Alameda-Pineda, "Robust Unsupervised Audio-Visual Speech Enhancement using a Mixture of VAEs," IEEE ICASSP 2020.
- 6 M. Sadeghi and X. Alameda-Pineda, "Mixture of Inference Networks for VAE-based Audio-visual Speech Enhancement", IEEE TST, 2021.
- 7 C. Bishop, Pattern Recognition and Machine Learning, Springer-Verlag Berlin, Heidelberg, 2006.
- 8 A. H. Abdelaziz, "NTCD-TIMIT: A new database and baseline for noise-robust audio-visual speech recognition," in Proc. INTERSPEECH, 2017.