

# Audio-visual Speech Enhancement based on Deep Generative Models

---

Mostafa Sadeghi

MULTISPEECH team  
Inria Nancy - Grand Est

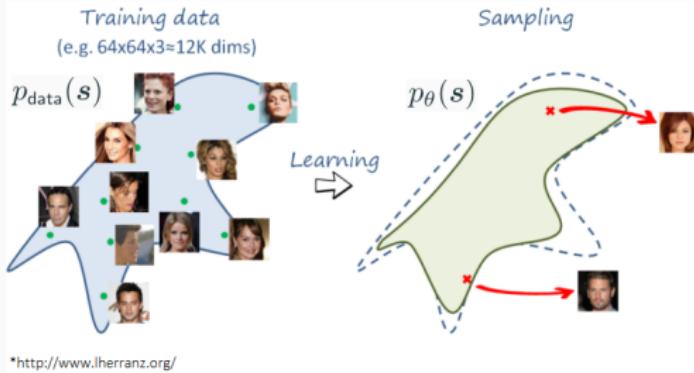
January 2023

- ① Deep Latent Variable Generative Models
- ② Audio-visual Speech Enhancement

# **Deep Latent Variable Generative Models**

---

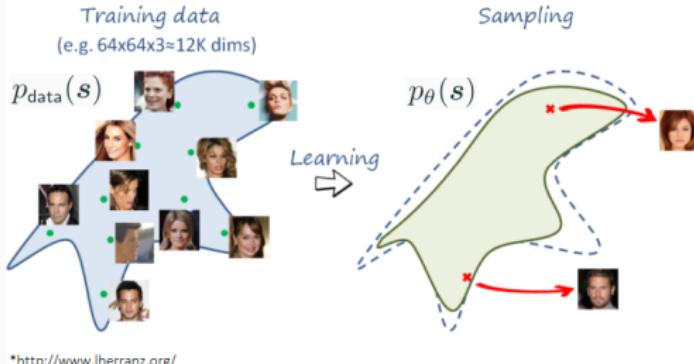
# Generative Models



**Objective:** Learning/simulating a complicated probability distribution of data,  $p_{\text{data}}$ , given some training samples:  $\mathbf{s}_i \sim p_{\text{data}}(\mathbf{s}), \quad i = 1, \dots, N$ .

---

# Generative Models



**Objective:** Learning/simulating a complicated probability distribution of data,  $p_{\text{data}}$ , given some training samples:  $s_i \sim p_{\text{data}}(s), \quad i = 1, \dots, N.$

Learn a parametric distribution  $p_{\theta}(s)$  as close as possible to  $p_{\text{data}}(s)$ :

$$\theta^* = \operatorname{argmin}_{\theta} D_{\text{KL}}\left(p_{\text{data}}(s) \parallel p_{\theta}(s)\right)$$

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

# Latent Variable Generative Models

- $s \in \mathbb{R}^n$ : observed variable
- $z \in \mathbb{R}^\ell$ : latent variable, a concise representation of  $s$  ( $\ell \ll n$ )

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

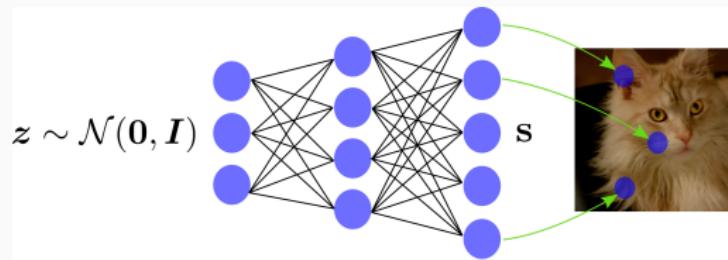
# Latent Variable Generative Models

- $s \in \mathbb{R}^n$ : observed variable
- $z \in \mathbb{R}^\ell$ : latent variable, a concise representation of  $s$  ( $\ell \ll n$ )

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

**Generating new samples:**

Draw  $z_k \sim p_\theta(z)$ , then draw a new sample  $s_k \sim p_\theta(s|z_k)$



# Parameter Estimation: Variational Autoencoder

## Variational Autoencoder (VAE)

Recall the generative model:

$$\begin{cases} p(z) = \mathcal{N}(\mathbf{0}, I) \\ p_{\theta}(s|z) = \mathcal{N}(\mu_{\theta}(z), \Sigma_{\theta}(z)) \end{cases}$$

VAE approximates  $p_{\theta}(z|s)$  with a *parametric Gaussian distribution*:

$$q_{\psi}(z|s) = \mathcal{N}(\mu_{\psi}(s), \Sigma_{\psi}(s))$$

# Parameter Estimation: Variational Autoencoder

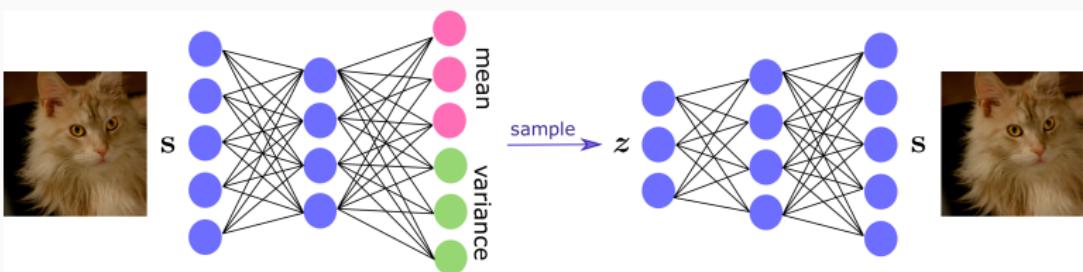
## Variational Autoencoder (VAE)

Recall the generative model:

$$\begin{cases} p(z) = \mathcal{N}(\mathbf{0}, I) \\ p_\theta(s|z) = \mathcal{N}(\mu_\theta(z), \Sigma_\theta(z)) \end{cases}$$

VAE approximates  $p_\theta(z|s)$  with a *parametric Gaussian distribution*:

$$q_\psi(z|s) = \mathcal{N}(\mu_\psi(s), \Sigma_\psi(s))$$



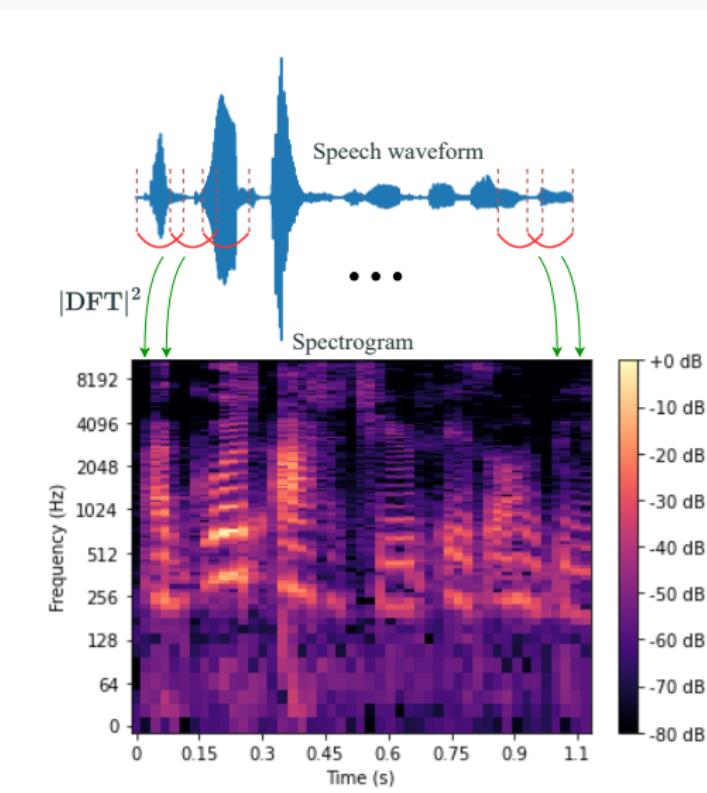
$$\theta^*, \psi^* = \arg \max_{\theta, \psi} \underbrace{\mathbb{E}_{q_\psi(z|s)} [\log p_\theta(s|z)]}_{\text{Reconstruction term}} - \underbrace{D_{\text{KL}}(q_\psi(z|s) \parallel p(z))}_{\text{Regularization term}}$$

## **Audio-visual Speech Enhancement**

---

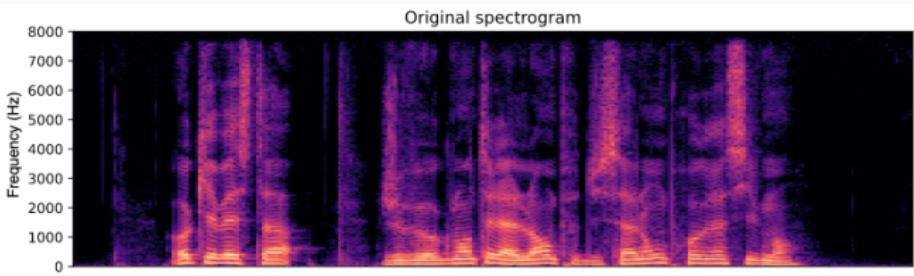
# Speech spectrogram

- A visual representation of the spectrum of frequencies based on Short-time Fourier transform (STFT).
- Apply discrete Fourier transform (DFT) to overlapping segments of speech waveform.
- Arrange the magnitude squared DFT vectors column-wise.

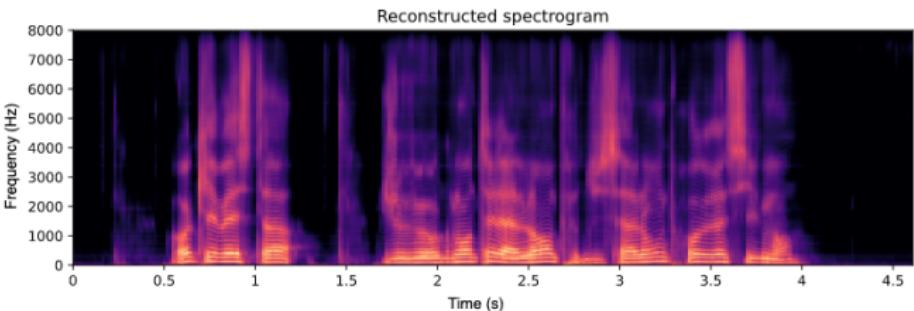
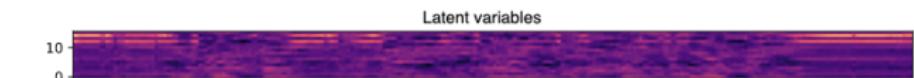


# Speech auto-encoding using VAE

Original signal



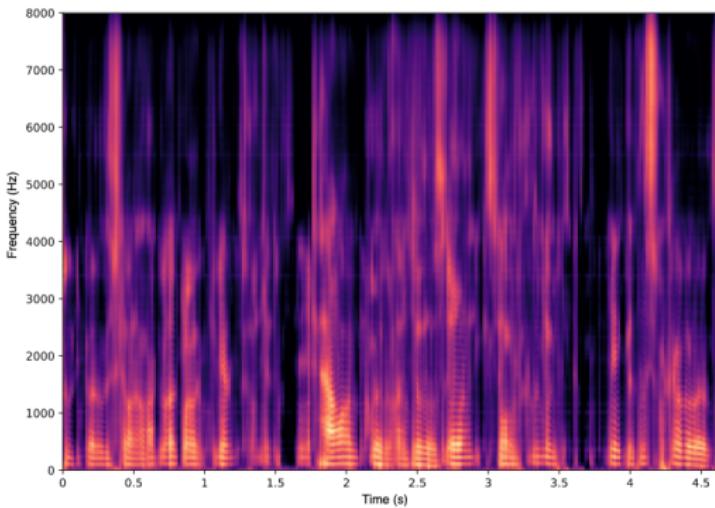
Reconstructed signal



# Speech generation using VAE

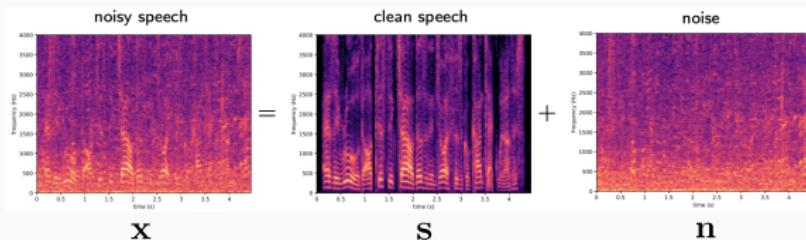
Sample a latent code  $\mathbf{z}$  from the prior, and give to the decoder to get a speech-like signal.

Generated signal



- Structured as a phoneme sequence, voiced/unvoiced phonemes
- Coarticulation, silences

# Speech Enhancement



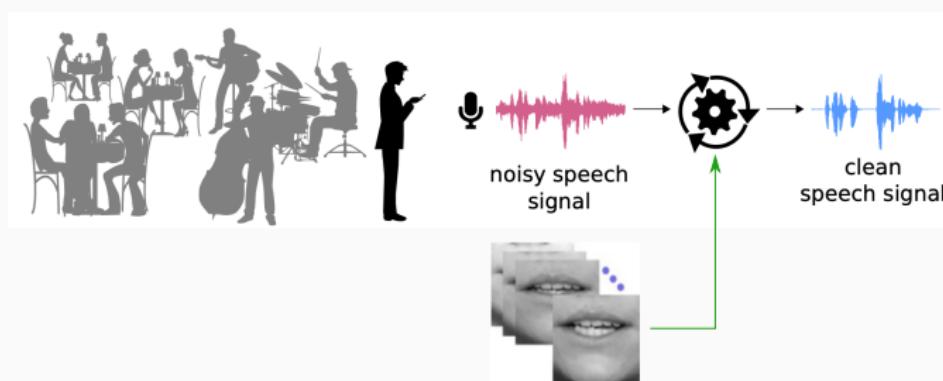
*Improve the quality and intelligibility of the observed noisy speech signal  $\mathbf{x}$ .*

- Close/distant conversations, listening comfort, hearing assistive devices.
- Automatic speech recognition for virtual assistants, social robots.

# Audio-visual Speech Enhancement (AVSE)

Visual modality (**lip movements**):

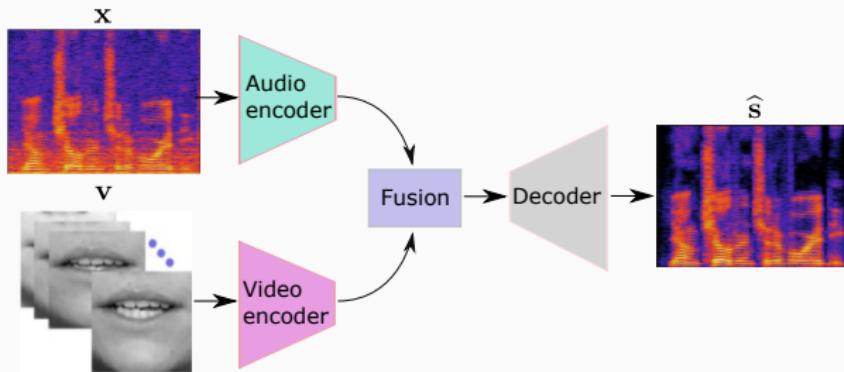
- Correlates well with speech signal (lip reading),
- Very helpful at **highly noisy** environments (unaffected by acoustic noise).



Given **noisy speech** observation  $\mathbf{x} = \mathbf{s} + \mathbf{n}$  & **visual data**  $\mathbf{v}$ , estimate the **clean speech signal**,  $\mathbf{s}$ .

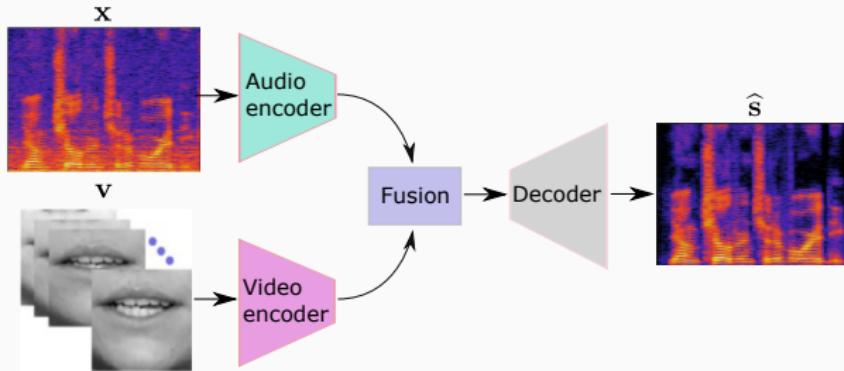
# Supervised (discriminative) AVSE

Model  $p_{\Theta}(s|x, v)$ , and learn  $\Theta$ :



# Supervised (discriminative) AVSE

Model  $p_{\Theta}(s|x, v)$ , and learn  $\Theta$ :



State-of-the-art performance, but ...

- Needs a huge audiovisual parallel (noise signal, clean speech) corpus
- Very deep and complex networks

# Unsupervised (generative) AVSE

*Speech enhancement **without** training on noise.*

Model  $p_{\Theta}(\mathbf{s}|\mathbf{x}, \mathbf{v}) \propto \underbrace{p_{\psi}(\mathbf{x}|\mathbf{s}, \mathbf{v})}_{\text{Inference}} \cdot \underbrace{p_{\theta}(\mathbf{s}|\mathbf{v})}_{\text{Training}}$ , and learn  $\Theta = \theta \cup \psi$ :

- **Training** - Learn speech prior distribution  $p_{\theta}(\mathbf{s}|\mathbf{v})$
- **Inference** - Model  $p_{\psi}(\mathbf{x}|\mathbf{s}, \mathbf{v})$ , and infer  $\mathbf{s}$  using  $p_{\theta}(\mathbf{s}|\mathbf{v})$

# Unsupervised (generative) AVSE

*Speech enhancement **without** training on noise.*

Model  $p_{\Theta}(s|x, v) \propto \underbrace{p_{\psi}(x|s, v)}_{\text{Inference}} \cdot \underbrace{p_{\theta}(s|v)}_{\text{Training}}$ , and learn  $\Theta = \theta \cup \psi$ :

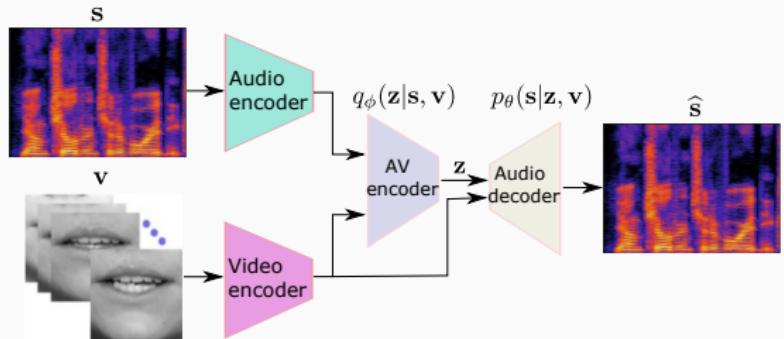
- **Training** - Learn speech prior distribution  $p_{\theta}(s|v)$
- **Inference** - Model  $p_{\psi}(x|s, v)$ , and infer  $s$  using  $p_{\theta}(s|v)$

▷ Advantages over supervised approaches:

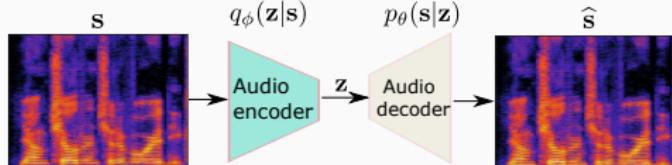
- No need to huge parallel corpora → compact & lightweight models
- Potentially better generalization performance

# VAE architectures

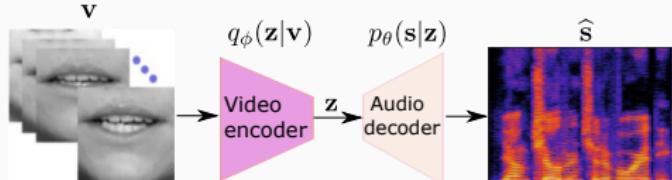
AV-VAE



A-VAE



V-VAE



# Speech Enhancement

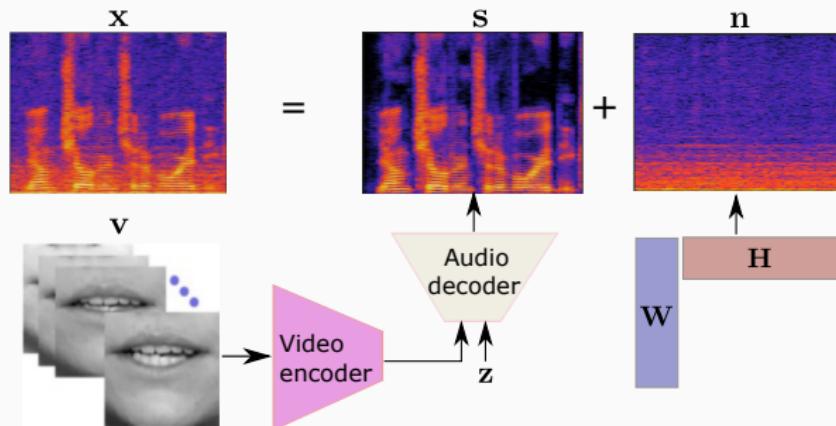
**Observation model:**

$$\forall t : \quad x_t = s_t + n_t$$

**Noise model:**

$$\forall t : \quad n_t \sim \mathcal{N}_c(\mathbf{0}, \text{diag}(\mathbf{WH}[:, t]))$$

**Clean speech model:** Trained generative (decoder) network.



# Speech Enhancement

## Inference:

- ▷ Parameters to be estimated:  $\psi = \{\mathbf{W}, \mathbf{H}\}$
- ▷ Observed variables:  $\{(\mathbf{x}_t, \mathbf{v}_t)\}_{t=1}^T$
- ▷ Latent variables:  $\mathbf{z} = \{\mathbf{z}_t\}_{t=1}^T$
- ▷ Likelihood:

$$p_\psi(\mathbf{x}_t | \mathbf{z}_t, \mathbf{v}_t) = \mathcal{N}_c \left( \mathbf{0}, \text{diag}(\boldsymbol{\sigma}_\theta^{av}(\mathbf{z}_t, \mathbf{v}_t)) + \text{diag}(\mathbf{W}\mathbf{H}[:, t]) \right)$$

## Parameter estimation:

$$\psi^* = \underset{\psi}{\operatorname{argmax}} \log p_\psi(\mathbf{x} | \mathbf{v}) = \underset{\psi}{\operatorname{argmax}} \int \log p_\psi(\mathbf{x}, \mathbf{z} | \mathbf{v}) d\mathbf{z}$$

# Parameter Estimation

## Expectation Maximization (EM)

From an initialization  $\psi^{(0)}$  of the parameters, iterate:

- **E-Step:**  $Q(\psi|\psi^{(k)}) = \mathbb{E}_{p_{\psi^{(k)}}(\mathbf{z}|\mathbf{x}, \mathbf{v})}[\log p_{\psi}(\mathbf{x}, \mathbf{z}, \mathbf{v})].$

Intractable expectation  $\rightarrow$  Markov chain Monte Carlo method.

$$Q(\psi|\psi^{(k)}) \approx \frac{1}{R} \sum_{r=1}^R \log p_{\psi}(\mathbf{x}, \mathbf{z}^{(r)}, \mathbf{v})$$

$$\{\mathbf{z}^{(r)}\}_{r=1}^R \sim p(\mathbf{z}|\mathbf{x}, \mathbf{v}; \boldsymbol{\theta}_u^*)$$

- **M-Step:**  $\psi^{(k+1)} \leftarrow \operatorname{argmax}_{\psi} Q(\psi|\psi^{(k)}).$

# Speech Estimation

Once the parameters are estimated, the speech STFT frames are estimated as follows ( $\forall f, t$ ):

$$\begin{aligned}\hat{s}_{ft} &= \mathbb{E}_{p_{\psi^*}(s_{ft}|x_{ft}, \mathbf{v}_t)}[s_{ft}] \\ &= \mathbb{E}_{p_{\psi^*}(\mathbf{z}_t|\mathbf{x}_t, \mathbf{v}_t)} \left[ \mathbb{E}_{p_{\psi^*}(s_{ft}|\mathbf{z}_t, \mathbf{v}_t, \mathbf{x}_t)}[s_{ft}] \right] \\ &= \mathbb{E}_{p_{\psi^*}(\mathbf{z}_t|\mathbf{x}_t, \mathbf{v}_t)} \left[ \frac{\sigma_{\theta,f}^{av}(\mathbf{z}_t, \mathbf{v}_t)}{\sigma_{\theta,f}^{av}(\mathbf{z}_t, \mathbf{v}_t) + (\mathbf{W}^* \mathbf{H}^*)_{f,t}} \right] \cdot x_{ft}.\end{aligned}$$

where,  $\psi^*$  denotes the set of estimated parameters by the EM method.

# Examples

Noisy

A-VAE

V-VAE

AV-VAE

## References

- [1] Kingma, D.P. and Welling, M., 2019. An introduction to variational autoencoders. *Foundations and Trends in Machine Learning*, 12(4), pp.307-392.
- [2] Girin, L., Leglaive, S., Bie, X., Diard, J., Hueber, T. and Alameda-Pineda, X., 2021. Dynamical Variational Autoencoders: A Comprehensive Review. *Foundations and Trends in Machine Learning*, 15(1-2), pp.1-175.
- [3] Sadeghi, M., Leglaive, S., Alameda-Pineda, X., Girin, L. and Horaud, R., 2020. Audio-visual speech enhancement using conditional variational auto-encoders. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 28, pp.1788-1800.
- [4] Kang, Z., Sadeghi, M., Horaud, R., Donley, J., Kumar, A. and Alameda-Pineda, X., 2022. Expression-preserving face frontalization improves visually assisted speech processing. *International Journal of Computer Vision (IJCV)*.

Thank you for your attention