# Long story short: the summary of (more than) a decade of probabilistic audio-visual learning

Xavier Alameda-Pineda and Radu Horaud
(and a long list of great people)

Inria
*Informatics* *mathematics*

Perception Group

RobotLearn

# Why Probabilistic Audio-Visual Learning?

# Why Probabilistic Audio-Visual Learning?

1. Scientific Challenges:
- ➢ Captured with different sensors
- ➢ Represent different phenomena
- ➢ Have different statistical patterns

# Why Probabilistic Audio-Visual Learning?

1. Scientific Challenges:
- Captured with different sensors
- Represent different phenomena
- Have different statistical patterns

2. From a technological perspective:
- Sensors are mature
- Embedded in numerous devices
- Cheap perception solution

# Why Probabilistic Audio-Visual Learning?

1. Scientific Challenges:
   - Captured with different sensors
   - Represent different phenomena
   - Have different statistical patterns

2. From a technological perspective:
   - Sensors are mature
   - Embedded in numerous devices
   - Cheap perception solution

3. Probabilistic Learning:
   - Usually unsupervised
   - Learn properties of noise/clutter
   - Infer latent variables

# Why Probabilistic Audio-Visual Learning?

1. Scientific Challenges:
- Captured with different sensors
- Represent different phenomena
- Have different statistical patterns

2. From a technological perspective:
- Sensors are mature
- Embedded in numerous devices
- Cheap perception solution

3. Probabilistic Learning:
- Usually unsupervised
- Learn properties of noise/clutter
- Infer latent variables

4. Potential impact:
- Behavior analysis
- Robotic social interaction
- Healthcare, training, security, ...

# Probabilistic Audio-Visual Learning Setting

# Probabilistic Audio-Visual Learning Setting

Picture a social scene with multiple people
chatting and interacting.

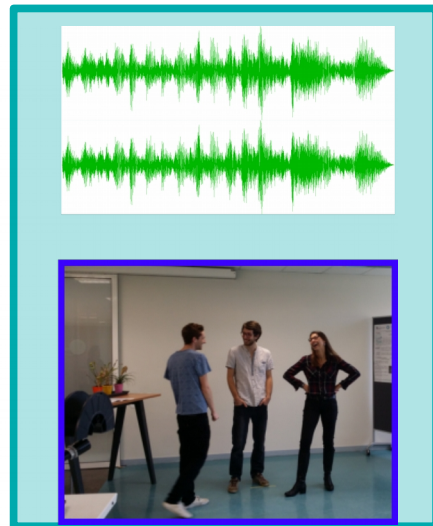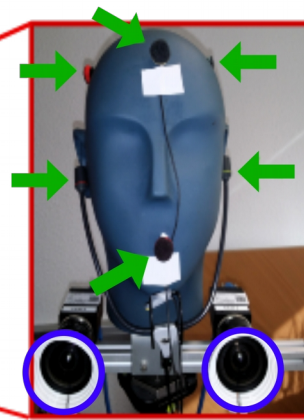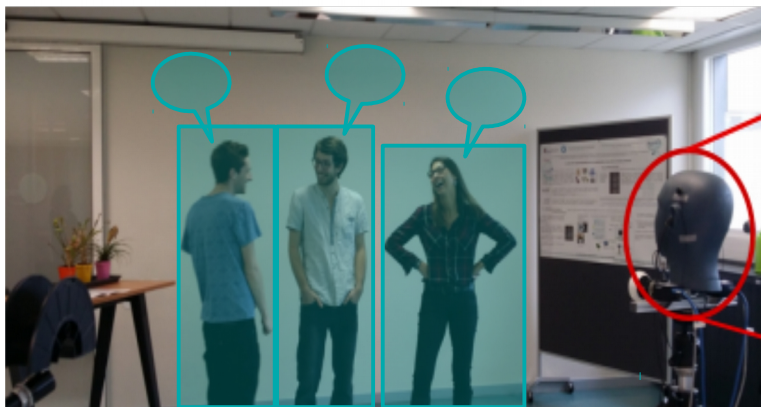# Probabilistic Audio-Visual Learning Setting

Picture a social scene with multiple people chatting and interacting.

A **device** observes de scene with **microphones** and **cameras**.

# Probabilistic Audio-Visual Learning Setting

Picture a social scene with multiple people chatting and interacting.

A **device** observes de scene with **microphones** and **cameras**.

We obtain **auditory (sounds)** and **visual (images)** features.

# Probabilistic Audio-Visual Learning Setting

Picture a social scene with multiple people chatting and interacting.

A **device** observes de scene with **microphones** and **cameras**.

We obtain **auditory (sounds)** and **visual (images)** features.

We would like to infer **latent variables (position, speaking status)**.

# Well, OK, but how?

# What is the methodology?

(Apologies if the next slides are a bit dense)

# Unsupervised Probabilistic Learning

# Unsupervised Probabilistic Learning

Observations will be denoted by $\mathbf{x}^a$ and $\mathbf{x}^v$
Latent variables by $\mathbf{z}$

We need to set up a probabilistic
model parametrised by the set $\boldsymbol{\theta}$

$$p_{\boldsymbol{\theta}}(\mathbf{x}^a, \mathbf{x}^v | \mathbf{z}) p_{\boldsymbol{\theta}}(\mathbf{z})$$

# Unsupervised Probabilistic Learning

Observations will be denoted by $\mathbf{x}^a$ and $\mathbf{x}^v$
Latent variables by $\mathbf{z}$

We need to set up a probabilistic
model parametrised by the set $\boldsymbol{\theta}$

$$p_{\boldsymbol{\theta}}(\mathbf{x}^a, \mathbf{x}^v | \mathbf{z}) p_{\boldsymbol{\theta}}(\mathbf{z})$$

**Learning** ↔ maximum likelihood:

$$\boldsymbol{\theta}^* = \arg\max_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x}^a, \mathbf{x}^v)$$

**Inference** ↔ expected value (or mode)

$$\mathbf{z}^* = \mathbb{E}_{p_{\boldsymbol{\theta}^*}(\mathbf{z}|\mathbf{x}^a, \mathbf{x}^v)}\{\mathbf{z}\}$$

# Unsupervised Probabilistic Learning

Observations will be denoted by $\mathbf{x}^a$ and $\mathbf{x}^v$
Latent variables by $\mathbf{z}$

We need to set up a probabilistic
model parametrised by the set $\boldsymbol{\theta}$

$$p_{\boldsymbol{\theta}}(\mathbf{x}^a, \mathbf{x}^v | \mathbf{z}) p_{\boldsymbol{\theta}}(\mathbf{z})$$

**Learning** ↔ maximum likelihood:
$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x}^a, \mathbf{x}^v)$$

**Inference** ↔ expected value (or mode)
$$\mathbf{z}^* = \mathbb{E}_{p_{\boldsymbol{\theta}^*}(\mathbf{z} | \mathbf{x}^a, \mathbf{x}^v)} \{\mathbf{z}\}$$

---

**Examples**: Gaussian mixture models, hidden Markov models, conditional random fields, linear dynamical systems (Kalman filter), probabilistic PCA, variational autoencoders (and dynamical ones), normalising flow, diffusion models, ...

# How to learn and infer? [v1 – Exact EM]

# How to learn and infer? [v1 – Exact EM]

Direct optimisation not analytically solvable:
$$\arg \max_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x})$$

Optimise the expected complete-data log-like:
$$\arg \max_{\boldsymbol{\theta}} \mathbb{E}_{p_{\boldsymbol{\theta}*}(\mathbf{z}|\mathbf{x})} \{\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})\}$$

# How to learn and infer? [v1 – Exact EM]

Direct optimisation not analytically solvable:

$$\arg\max_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x})$$

Optimise the expected complete-data log-like:

$$\arg\max_{\boldsymbol{\theta}} \mathbb{E}_{p_{\boldsymbol{\theta}*}(\mathbf{z}|\mathbf{x})}\{\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})\}$$

**Expectation**: given $\theta^*$, compute:

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \mathbb{E}_{p_{\boldsymbol{\theta}*}(\mathbf{z}|\mathbf{x})}\{\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})\}$$

**Maximisation**: set up the new $\theta^*$ to:

$$\boldsymbol{\theta}^* \leftarrow \arg\max_{\boldsymbol{\theta}} \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$$

# How to learn and infer? [v1 – Exact EM]

Direct optimisation not analytically solvable:

$$\arg\max_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x})$$

Optimise the expected complete-data log-like:

$$\arg\max_{\boldsymbol{\theta}} \mathbb{E}_{p_{\boldsymbol{\theta}*}(\mathbf{z}|\mathbf{x})}\{\log p_{\boldsymbol{\theta}}(\mathbf{x},\mathbf{z})\}$$

**Expectation**: given $\boldsymbol{\theta}^*$, compute:

$$\mathcal{Q}(\boldsymbol{\theta},\boldsymbol{\theta}^*) = \mathbb{E}_{p_{\boldsymbol{\theta}*}(\mathbf{z}|\mathbf{x})}\{\log p_{\boldsymbol{\theta}}(\mathbf{x},\mathbf{z})\}$$

**Maximisation**: set up the new $\boldsymbol{\theta}^*$ to:

$$\boldsymbol{\theta}^* \leftarrow \arg\max_{\boldsymbol{\theta}} \mathcal{Q}(\boldsymbol{\theta},\boldsymbol{\theta}^*)$$

*Finding Audio-Visual Events in Informal Social Gatherings*, ICMI 2011 – **Outstanding Paper Award**

# How to learn and infer? [v1 – Exact EM]

Direct optimisation not analytically solvable:

$$\arg\max_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x})$$

Optimise the expected complete-data log-like:

$$\arg\max_{\boldsymbol{\theta}} \mathbb{E}_{p_{\boldsymbol{\theta}^*}(\mathbf{z}|\mathbf{x})}\{\log p_{\boldsymbol{\theta}}(\mathbf{x},\mathbf{z})\}$$

**Expectation**: given $\boldsymbol{\theta}^*$, compute:

$$\mathcal{Q}(\boldsymbol{\theta},\boldsymbol{\theta}^*) = \mathbb{E}_{p_{\boldsymbol{\theta}^*}(\mathbf{z}|\mathbf{x})}\{\log p_{\boldsymbol{\theta}}(\mathbf{x},\mathbf{z})\}$$

**Maximisation**: set up the new $\boldsymbol{\theta}^*$ to:

$$\boldsymbol{\theta}^* \leftarrow \arg\max_{\boldsymbol{\theta}} \mathcal{Q}(\boldsymbol{\theta},\boldsymbol{\theta}^*)$$



*Finding Audio-Visual Events in Informal Social Gatherings*, ICMI 2011 – **Outstanding Paper Award**

*Visually-Guided Robot Hearing*, IJRR 2012

# How to learn and infer? [v1 – Exact EM]

Direct optimisation not analytically solvable:

$$\arg \max_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x})$$

**Expectation**: given $\boldsymbol{\theta}^*$, compute:

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \mathbb{E}_{p_{\boldsymbol{\theta}^*}(\mathbf{z}|\mathbf{x})}\{\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})\}$$

Optimise the expected complete-data log-like:

$$\arg \max_{\boldsymbol{\theta}} \mathbb{E}_{p_{\boldsymbol{\theta}^*}(\mathbf{z}|\mathbf{x})}\{\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})\}$$

**Maximisation**: set up the new $\boldsymbol{\theta}^*$ to:

$$\boldsymbol{\theta}^* \leftarrow \arg \max_{\boldsymbol{\theta}} \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$$

*Finding Audio-Visual Events in Informal Social Gatherings*, ICMI 2011 – **Outstanding Paper Award**

*Visually-Guided Robot Hearing*, IJRR 2012

*EM Algorithms for Weighted-Data Clustering for AV Scene Analysis*, TPAMI 2016.

# How to learn and infer? [v1 − Exact EM]

Direct optimisation not analytically solvable:

$$\arg\max_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x})$$

Optimise the expected complete-data log-like:

$$\arg\max_{\boldsymbol{\theta}} \mathbb{E}_{p_{\boldsymbol{\theta}^*}(\mathbf{z}|\mathbf{x})}\{\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})\}$$

**Expectation**: given $\boldsymbol{\theta}^*$, compute:

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \mathbb{E}_{p_{\boldsymbol{\theta}^*}(\mathbf{z}|\mathbf{x})}\{\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})\}$$

**Maximisation**: set up the new $\boldsymbol{\theta}^*$ to:

$$\boldsymbol{\theta}^* \leftarrow \arg\max_{\boldsymbol{\theta}} \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$$



*Finding Audio-Visual Events in Informal Social Gatherings*, ICMI 2011 − **Outstanding Paper Award**

*Visually-Guided Robot Hearing*, IJRR 2012

*EM Algorithms for Weighted-Data Clustering for AV Scene Analysis*, TPAMI 2016.

*Acoustic Space Learning for Sound-source Separation and Localization on Binaural Manifolds*, Neural Systems, 2015 − **Hojjat Adeli Award for Outstanding Contributions in Neural Systems**

# How to learn and infer? [v1 – Exact EM]

Direct optimisation not analytically solvable:

$$\arg \max_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x})$$

**Expectation**: given $\boldsymbol{\theta}^*$, compute:

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \mathbb{E}_{p_{\boldsymbol{\theta}^*}(\mathbf{z}|\mathbf{x})}\{\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})\}$$

**What if we cannot compute the posterior?**

Optimise the expected complete-data log-like:

$$\arg \max_{\boldsymbol{\theta}} \mathbb{E}_{p_{\boldsymbol{\theta}^*}(\mathbf{z}|\mathbf{x})}\{\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})\}$$

**Maximisation**: set up the new $\boldsymbol{\theta}^*$ to:

$$\boldsymbol{\theta}^* \leftarrow \arg \max_{\boldsymbol{\theta}} \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$$

*Finding Audio-Visual Events in Informal Social Gatherings*, ICMI 2011 – **Outstanding Paper Award**

*Visually-Guided Robot Hearing*, IJRR 2012

*EM Algorithms for Weighted-Data Clustering for AV Scene Analysis*, TPAMI 2016.

*Acoustic Space Learning for Sound-source Separation and Localization on Binaural Manifolds*, Neural Systems, 2015 – **Hojjat Adeli Award for Outstanding Contributions in Neural Systems**

# How to learn and infer? [v2 – Variational EM]

Close look to the log-likelihood, for any distribution          :

# How to learn and infer? [v2 – Variational EM]

Close look to the log-likelihood, for any distribution $q_\phi(\mathbf{z}|\mathbf{x})$:

$$\log p_{\boldsymbol\theta}(\mathbf{x}) = \mathbb{E}_{q_\phi}(\mathbf{z}|\mathbf{x})\{\log p_{\boldsymbol\theta}(\mathbf{x}, \mathbf{z})\} + \mathcal{H}_{q_\phi} + \mathcal{D}_{\mathsf{KL}}(q_\phi \| p_{\boldsymbol\theta}(\mathbf{z}|\mathbf{x}))$$

# How to learn and infer? [v2 – Variational EM]

Close look to the log-likelihood, for any distribution $q_\phi(\mathbf{z}|\mathbf{x})$:

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) = \mathbb{E}_{q_\phi}(\mathbf{z}|\mathbf{x})\{\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})\} + \mathcal{H}_{q_\phi} + \mathcal{D}_{\mathsf{KL}}(q_\phi \| p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}))$$

In the exact EM, $q_\phi(\mathbf{z}|\mathbf{x}) = p_{\boldsymbol{\theta}^*}(\mathbf{z}|\mathbf{x})$ and $\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ is a tight lower bound. In some cases (not too complex), $p_{\boldsymbol{\theta}^*}(\mathbf{z}|\mathbf{x})$ is computationally intractable. Break the hidden in two (or +).

# How to learn and infer? [v2 – Variational EM]

Close look to the log-likelihood, for any distribution $q_\phi(\mathbf{z}|\mathbf{x})$:

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) = \mathbb{E}_{q_\phi}(\mathbf{z}|\mathbf{x})\{\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})\} + \mathcal{H}_{q_\phi} + \mathcal{D}_{\mathsf{KL}}(q_\phi\|p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}))$$

In the exact EM, $q_\phi(\mathbf{z}|\mathbf{x}) = p_{\boldsymbol{\theta}^*}(\mathbf{z}|\mathbf{x})$ and $\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ is a tight lower bound. In some cases (not too complex), $p_{\boldsymbol{\theta}^*}(\mathbf{z}|\mathbf{x})$ is computationally intractable. Break the hidden in two (or +).

$$\mathbf{z} = (\mathbf{z}_1, \mathbf{z}_2) \quad \text{and} \quad q_\phi(\mathbf{z}) = q_{\phi_1}(\mathbf{z}_1)q_{\phi_2}(\mathbf{z}_2) \qquad q^*_{\phi_1}, q^*_{\phi_2} = \arg\min \mathcal{D}_{\mathsf{KL}}(q_{\phi_1}q_{\phi_2}\|p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}))$$

# How to learn and infer? [v2 – Variational EM]

Close look to the log-likelihood, for any distribution $q_\phi(\mathbf{z}|\mathbf{x})$:

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) = \mathbb{E}_{q_\phi}(\mathbf{z}|\mathbf{x})\{\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})\} + \mathcal{H}_{q_\phi} + \mathcal{D}_{\mathsf{KL}}(q_\phi \| p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}))$$

In the exact EM, $q_\phi(\mathbf{z}|\mathbf{x}) = p_{\boldsymbol{\theta}^*}(\mathbf{z}|\mathbf{x})$ and $\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ is a tight lower bound. In some cases (not too complex), $p_{\boldsymbol{\theta}^*}(\mathbf{z}|\mathbf{x})$ is computationally intractable. Break the hidden in two (or +).

$$\mathbf{z} = (\mathbf{z}_1, \mathbf{z}_2) \quad \text{and} \quad q_\phi(\mathbf{z}) = q_{\phi_1}(\mathbf{z}_1)q_{\phi_2}(\mathbf{z}_2) \qquad q^*_{\phi_1}, q^*_{\phi_2} = \arg\min \mathcal{D}_{\mathsf{KL}}(q_{\phi_1}q_{\phi_2} \| p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}))$$



*A variational EM algorithm for the separation of moving sound sources*, IEEE WASPAA 2015 –
**Best Student Paper Award**

# How to learn and infer? [v2 – Variational EM]

Close look to the log-likelihood, for any distribution $q_\phi(\mathbf{z}|\mathbf{x})$:

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) = \mathbb{E}_{q_\phi}(\mathbf{z}|\mathbf{x})\{\log p_{\boldsymbol{\theta}}(\mathbf{x},\mathbf{z})\} + \mathcal{H}_{q_\phi} + \mathcal{D}_{\mathsf{KL}}(q_\phi \| p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}))$$

In the exact EM, $q_\phi(\mathbf{z}|\mathbf{x}) = p_{\boldsymbol{\theta}^*}(\mathbf{z}|\mathbf{x})$ and $\mathcal{Q}(\boldsymbol{\theta},\boldsymbol{\theta}^*)$ is a tight lower bound. In some cases (not too complex), $p_{\boldsymbol{\theta}^*}(\mathbf{z}|\mathbf{x})$ is computationally intractable. Break the hidden in two (or +).

$$\mathbf{z} = (\mathbf{z}_1, \mathbf{z}_2) \quad \text{and} \quad q_\phi(\mathbf{z}) = q_{\phi_1}(\mathbf{z}_1)q_{\phi_2}(\mathbf{z}_2) \qquad q_{\phi_1}^*, q_{\phi_2}^* = \arg\min \mathcal{D}_{\mathsf{KL}}(q_{\phi_1}q_{\phi_2} \| p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}))$$

---



*A variational EM algorithm for the separation of moving sound sources*, IEEE WASPAA 2015 –
**Best Student Paper Award**

*A Variational EM Algorithm for the Separation of Time-Varying Convolutive Audio Mixtures*,
IEEE/ACM TASLP 2016

# How to learn and infer? [v2 – Variational EM]

Close look to the log-likelihood, for any distribution $q_\phi(\mathbf{z}|\mathbf{x})$:

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) = \mathbb{E}_{q_\phi}(\mathbf{z}|\mathbf{x})\{\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})\} + \mathcal{H}_{q_\phi} + \mathcal{D}_{\mathsf{KL}}(q_\phi \| p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}))$$

In the exact EM, $q_\phi(\mathbf{z}|\mathbf{x}) = p_{\boldsymbol{\theta}^*}(\mathbf{z}|\mathbf{x})$ and $\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ is a tight lower bound. In some cases (not too complex), $p_{\boldsymbol{\theta}^*}(\mathbf{z}|\mathbf{x})$ is computationally intractable. Break the hidden in two (or +).

$$\mathbf{z} = (\mathbf{z}_1, \mathbf{z}_2) \quad \text{and} \quad q_\phi(\mathbf{z}) = q_{\phi_1}(\mathbf{z}_1) q_{\phi_2}(\mathbf{z}_2) \qquad q_{\phi_1}^*, q_{\phi_2}^* = \arg\min \mathcal{D}_{\mathsf{KL}}(q_{\phi_1} q_{\phi_2} \| p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}))$$

*A variational EM algorithm for the separation of moving sound sources*, IEEE WASPAA 2015 – **Best Student Paper Award**

*A Variational EM Algorithm for the Separation of Time-Varying Convolutive Audio Mixtures*, IEEE/ACM TASLP 2016

*Variational Bayesian Inference for Audio-Visual Tracking of Multiple Speakers*, IEEE TPAMI 2019.

# How to learn and infer? [v2 – Variational EM]

Close look to the log-likelihood, for any distribution $q_\phi(\mathbf{z}|\mathbf{x})$:

$$\log p_{\boldsymbol\theta}(\mathbf{x}) = \mathbb{E}_{q_\phi}(\mathbf{z}|\mathbf{x})\{\log p_{\boldsymbol\theta}(\mathbf{x},\mathbf{z})\} + \mathcal{H}_{q_\phi} + \mathcal{D}_{\mathsf{KL}}(q_\phi\|p_{\boldsymbol\theta}(\mathbf{z}|\mathbf{x}))$$

In the exact EM, $q_\phi(\mathbf{z}|\mathbf{x}) = p_{\boldsymbol\theta^*}(\mathbf{z}|\mathbf{x})$ and $\mathcal{Q}(\boldsymbol\theta,\boldsymbol\theta^*)$ is a tight lower bound. In some cases (not too complex), $p_{\boldsymbol\theta^*}(\mathbf{z}|\mathbf{x})$ is computationally intractable. Break the hidden in two (or +).

$$\mathbf{z} = (\mathbf{z}_1, \mathbf{z}_2) \quad \text{and} \quad q_\phi(\mathbf{z}) = q_{\phi_1}(\mathbf{z}_1)q_{\phi_2}(\mathbf{z}_2) \qquad q_{\phi_1}^*, q_{\phi_2}^* = \arg\min \mathcal{D}_{\mathsf{KL}}(q_{\phi_1}q_{\phi_2}\|p_{\boldsymbol\theta}(\mathbf{z}|\mathbf{x}))$$



*A variational EM algorithm for the separation of moving sound sources*, IEEE WASPAA 2015 – **Best Student Paper Award**

*A Variational EM Algorithm for the Separation of Time-Varying Convolutive Audio Mixtures*, IEEE/ACM TASLP 2016

*Variational Bayesian Inference for Audio-Visual Tracking of Multiple Speakers*, IEEE TPAMI 2019.

*Variational Inference and Learning of Piecewise-linear Dynamical Systems*, IEEE TNNLS, 2021.

# How to learn and infer? [v3 – VAE]

# How to learn and infer? [v3 – VAE]

What if the posterior $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$ does not have analytic form?

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) = \mathbb{E}_{q_{\phi}}(\mathbf{z}|\mathbf{x})\{\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})\} + \mathcal{H}_{q_{\phi}} + \mathcal{D}_{\mathsf{KL}}(q_{\phi}\|p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}))$$

# How to learn and infer? [v3 – VAE]

What if the posterior $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$ does not have analytic form?

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) = \mathbb{E}_{q_{\phi}}(\mathbf{z}|\mathbf{x})\{\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})\} + \mathcal{H}_{q_{\phi}} + \mathcal{D}_{\mathsf{KL}}(q_{\phi}\|p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}))$$

Seminal case, the generative model $p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})$ is implemented with a deep neural network.

# How to learn and infer? [v3 – VAE]

What if the posterior $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$ does not have analytic form?

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) = \mathbb{E}_{q_\phi}(\mathbf{z}|\mathbf{x})\{\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})\} + \mathcal{H}_{q_\phi} + \mathcal{D}_{\text{KL}}(q_\phi \| p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}))$$

Seminal case, the generative model $p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})$ is implemented with a deep neural network.
The posterior is approximated with a different neural network: $q_\phi(\mathbf{z}|\mathbf{x}) \approx p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) \geq \mathbb{E}_{q_\phi}(\mathbf{z}|\mathbf{x})\{\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})\} - \mathcal{D}_{\text{KL}}(q_\phi \| p_{\boldsymbol{\theta}}(\mathbf{z})) = \mathcal{L}_{\text{ELBO}}(\boldsymbol{\theta}, \phi)$$

# How to learn and infer? [v3 – VAE]

What if the posterior $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$ does not have analytic form?

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) = \mathbb{E}_{q_\phi}(\mathbf{z}|\mathbf{x})\{\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})\} + \mathcal{H}_{q_\phi} + \mathcal{D}_{\mathsf{KL}}(q_\phi \| p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}))$$

Seminal case, the generative model $p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})$ is implemented with a deep neural network.
The posterior is approximated with a different neural network: $q_\phi(\mathbf{z}|\mathbf{x}) \approx p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) \geq \mathbb{E}_{q_\phi}(\mathbf{z}|\mathbf{x})\{\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})\} - \mathcal{D}_{\mathsf{KL}}(q_\phi \| p_{\boldsymbol{\theta}}(\mathbf{z})) = \mathcal{L}_{\mathsf{ELBO}}(\boldsymbol{\theta}, \phi)$$

# How to learn and infer? [v3 – VAE]

What if the posterior $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$ does not have analytic form?

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) = \mathbb{E}_{q_{\phi}}(\mathbf{z}|\mathbf{x})\{\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})\} + \mathcal{H}_{q_{\phi}} + \mathcal{D}_{\mathsf{KL}}(q_{\phi}\|p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}))$$

Seminal case, the generative model $p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})$ is implemented with a deep neural network.
The posterior is approximated with a different neural network: $q_{\phi}(\mathbf{z}|\mathbf{x}) \approx p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) \geq \mathbb{E}_{q_{\phi}}(\mathbf{z}|\mathbf{x})\{\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})\} - \mathcal{D}_{\mathsf{KL}}(q_{\phi}\|p_{\boldsymbol{\theta}}(\mathbf{z})) = \mathcal{L}_{\mathsf{ELBO}}(\boldsymbol{\theta}, \phi)$$

*Audio-visual speech enhancement using conditional variational auto-encoders*, IEEE/ACM TASLP 2020

*Mixture of inference networks for VAE-based audio-visual speech enhancement*, IEEE TSP 2021

# How to learn and infer? [v3 – VAE]

What if the posterior $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$ does not have analytic form?

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) = \mathbb{E}_{q_{\phi}}(\mathbf{z}|\mathbf{x})\{\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})\} + \mathcal{H}_{q_{\phi}} + \mathcal{D}_{\mathsf{KL}}(q_{\phi}\|p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}))$$

Seminal case, the generative model $p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})$ is implemented with a deep neural network.
The posterior is approximated with a different neural network: $q_{\phi}(\mathbf{z}|\mathbf{x}) \approx p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) \geq \mathbb{E}_{q_{\phi}}(\mathbf{z}|\mathbf{x})\{\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})\} - \mathcal{D}_{\mathsf{KL}}(q_{\phi}\|p_{\boldsymbol{\theta}}(\mathbf{z})) = \mathcal{L}_{\mathsf{ELBO}}(\boldsymbol{\theta}, \phi)$$

*Audio-visual speech enhancement using conditional variational auto-encoders*, IEEE/ACM TASLP 2020

*Mixture of inference networks for VAE-based audio-visual speech enhancement*, IEEE TSP 2021

*Deep variational generative models for audio-visual speech separation*, IEEE MLSP 2021

# How to learn and infer? [v3 – VAE]

What if the posterior $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$ does not have analytic form?

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) = \mathbb{E}_{q_\phi}(\mathbf{z}|\mathbf{x})\{\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})\} + \mathcal{H}_{q_\phi} + \mathcal{D}_{\mathsf{KL}}(q_\phi \| p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}))$$

Seminal case, the generative model $p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})$ is implemented with a deep neural network.
The posterior is approximated with a different neural network: $q_\phi(\mathbf{z}|\mathbf{x}) \approx p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) \geq \mathbb{E}_{q_\phi}(\mathbf{z}|\mathbf{x})\{\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})\} - \mathcal{D}_{\mathsf{KL}}(q_\phi \| p_{\boldsymbol{\theta}}(\mathbf{z})) = \mathcal{L}_{\mathsf{ELBO}}(\boldsymbol{\theta}, \phi)$$



*Audio-visual speech enhancement using conditional variational auto-encoders*, IEEE/ACM TASLP 2020

*Mixture of inference networks for VAE-based audio-visual speech enhancement*, IEEE TSP 2021

*Deep variational generative models for audio-visual speech separation*, IEEE MLSP 2021

*Switching variational auto-encoders for noise-agnostic audio-visual speech enhancement*, IEEE ICASSP, 2021.

# How to learn and infer? [v4 – Dynamical VAE]

# How to learn and infer? [v4 – Dynamical VAE]

How to model temporal dependencies with deep generative models?

# How to learn and infer? [v4 – Dynamical VAE]

How to model temporal dependencies with deep generative models?

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}_{1:T}) \geq \mathbb{E}_{q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})}\{\log p_{\boldsymbol{\theta}}(\mathbf{x}_{1:T}|\mathbf{z}_{1:T})\} - \mathcal{D}_{\mathsf{KL}}(q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})\|p_{\boldsymbol{\theta}}(\mathbf{z}_{1:T}))$$

Both the generative model and approximate posterior need to model temporal dependencies implemented via, e.g., recurrent or transformer networks.

# How to learn and infer? [v4 – Dynamical VAE]

How to model temporal dependencies with deep generative models?

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}_{1:T}) \geq \mathbb{E}_{q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})}\{\log p_{\boldsymbol{\theta}}(\mathbf{x}_{1:T}|\mathbf{z}_{1:T})\} - \mathcal{D}_{\mathsf{KL}}(q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})\|p_{\boldsymbol{\theta}}(\mathbf{z}_{1:T}))$$

Both the generative model and approximate posterior need to model temporal dependencies implemented via, e.g., recurrent or transformer networks.

The dependencies of $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$ can be implemented in $q_{\phi}(\mathbf{z}|\mathbf{x})$ (still an approximation).

# How to learn and infer? [v4 – Dynamical VAE]

How to model temporal dependencies with deep generative models?

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}_{1:T}) \geq \mathbb{E}_{q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})}\{\log p_{\boldsymbol{\theta}}(\mathbf{x}_{1:T}|\mathbf{z}_{1:T})\} - \mathcal{D}_{\mathsf{KL}}(q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})\|p_{\boldsymbol{\theta}}(\mathbf{z}_{1:T}))$$

Both the generative model and approximate posterior need to model temporal dependencies implemented via, e.g., recurrent or transformer networks.

The dependencies of $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$ can be implemented in $q_{\phi}(\mathbf{z}|\mathbf{x})$ (still an approximation).



*A recurrent variational autoencoder for speech enhancement*, IEEE ICASSP 2020

# How to learn and infer? [v4 – Dynamical VAE]

How to model temporal dependencies with deep generative models?

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}_{1:T}) \geq \mathbb{E}_{q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})}\{\log p_{\boldsymbol{\theta}}(\mathbf{x}_{1:T}|\mathbf{z}_{1:T})\} - \mathcal{D}_{\mathsf{KL}}(q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})\|p_{\boldsymbol{\theta}}(\mathbf{z}_{1:T}))$$

Both the generative model and approximate posterior need to model temporal dependencies implemented via, e.g., recurrent or transformer networks.

The dependencies of $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$ can be implemented in $q_{\phi}(\mathbf{z}|\mathbf{x})$ (still an approximation).

*A recurrent variational autoencoder for speech enhancement*, IEEE ICASSP 2020

*Dynamical variational autoencoders: A comprehensive review*, FnT Machine Learning 2021

# How to learn and infer? [v4 – Dynamical VAE]

How to model temporal dependencies with deep generative models?

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}_{1:T}) \geq \mathbb{E}_{q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})}\{\log p_{\boldsymbol{\theta}}(\mathbf{x}_{1:T}|\mathbf{z}_{1:T})\} - \mathcal{D}_{\mathsf{KL}}(q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})\|p_{\boldsymbol{\theta}}(\mathbf{z}_{1:T}))$$

Both the generative model and approximate posterior need to model temporal dependencies implemented via, e.g., recurrent or transformer networks.

The dependencies of $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$ can be implemented in $q_{\phi}(\mathbf{z}|\mathbf{x})$ (still an approximation).



*A recurrent variational autoencoder for speech enhancement*, IEEE ICASSP 2020

*Dynamical variational autoencoders: A comprehensive review*, FnT Machine Learning 2021

*HiT-DVAE: Human Motion Generation via Hierarchical Transformer Dynamical VAE*

# How to learn and infer? [v4 – Dynamical VAE]

How to model temporal dependencies with deep generative models?

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}_{1:T}) \geq \mathbb{E}_{q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})} \{\log p_{\boldsymbol{\theta}}(\mathbf{x}_{1:T}|\mathbf{z}_{1:T})\} - \mathcal{D}_{\mathsf{KL}}(q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T}) \| p_{\boldsymbol{\theta}}(\mathbf{z}_{1:T}))$$

Both the generative model and approximate posterior need to model temporal dependencies implemented via, e.g., recurrent or transformer networks.

The dependencies of $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$ can be implemented in $q_{\phi}(\mathbf{z}|\mathbf{x})$ (still an approximation).



*A recurrent variational autoencoder for speech enhancement*, IEEE ICASSP 2020

*Dynamical variational autoencoders: A comprehensive review*, FnT Machine Learning 2021

*HiT-DVAE: Human Motion Generation via Hierarchical Transformer Dynamical VAE*

*Unsupervised speech enhancement using dynamical variational autoencoders*, IEEE TASLP, 2022

# How to learn and infer? [v4 – Dynamical VAE]

How to model temporal dependencies with deep generative models?

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}_{1:T}) \geq \mathbb{E}_{q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})}\{\log p_{\boldsymbol{\theta}}(\mathbf{x}_{1:T}|\mathbf{z}_{1:T})\} - \mathcal{D}_{\mathsf{KL}}(q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})\|p_{\boldsymbol{\theta}}(\mathbf{z}_{1:T}))$$

Both the generative model and approximate posterior need to model temporal dependencies implemented via, e.g., recurrent or transformer networks.

The dependencies of $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$ can be implemented in $q_{\phi}(\mathbf{z}|\mathbf{x})$ (still an approximation).

## How is all this possible?

*A recurrent variational autoencoder for speech enhancement*, IEEE ICASSP 2020

*Dynamical variational autoencoders: A comprehensive review*, FnT Machine Learning 2021

*HiT-DVAE: Human Motion Generation via Hierarchical Transformer Dynamical VAE*

*Unsupervised speech enhancement using dynamical variational autoencoders*, IEEE TASLP, 2022

# Scientific life lessons

# (Leçons de vie scientifique)

# Back to 2009, Barcelona...

# Back to 2009, Barcelona...





*La paciència és la mare de la ciència.* (catalan proverb)

# Masters Thesis



Safe environment!

# Masters Thesis



Safe environment!

*Take pride in what you do, but*
*don't let your pride guide you.*

# PhD Thesis – 1ˢᵗ Submission



Want to do things on your own, but perhaps not ready ;)

# PhD Thesis – 1st Submission



Want to do things on your own, but perhaps not ready ;)

*Perseverance is key.*

# PhD Thesis – Paper Writing



"Writing" journal paper by concatenating two short papers (obtaining a long complex paper)

# PhD Thesis – Paper Writing



"Writing" journal paper by concatenating two short papers (obtaining a long complex paper)

*Writing is 1/3 of your research time.*

# PhD Thesis – Focus

We could do that, and that, and that, ...

# PhD Thesis – Focus

We could do that, and that, and that, ...



*...yeah, sure, but focus.*

# PhD Thesis – Defence

*Tell a story...*

# PhD Thesis – Defence

with highlights,

*Tell a story...*

# PhD Thesis – Defence

with highlights,

unexpected features,

*Tell a story...*

# PhD Thesis – Defence

*Tell a story...*

with highlights,

unexpected features,

limitations and how to overcome them,

# PhD Thesis – Defence

with highlights,

unexpected features,

*Tell a story...*
*...and be proud of your work!*

limitations and how
to overcome them,

ICMI 2011

EUSIPCO 2012

Humanoids 2012

Humanoids 2012

JMUI 2013

WASPAA 2013

Humanoids 2013

ICASSP 2013

ICRA 2014

TASLP 2014

MLSP 2014

WASPAA 2015

IJRR 2015

CVIU 2016

ECCV-W 2016

TPAMI 2016

ICASSP 2016

TASLP 2016

ICCV-W 2017

IROS 2017

ICASSP 2017

ECCV 2018

ICASSP 2018

IWAENC 2018

JSTSP 2019

SPL 2019

TPAMI 2019

TPAMI 2020

CVPR 2020

TASLP 2020

ICASSP 2020

ICPR 2021

ICPR-W 2021

TNNLS 2021

ICASSP 2022

IJCV 2022

| ICMI 2011 | EUSIPCO 2012 | Humanoids 2012 | Humanoids 2012 | JMUI 2013 | WASPAA 2013 | Humanoids 2013 | ICASSP 2013 | ICRA 2014 |

**MIAI Grenoble Alpes**

**SPRING**

**RobotLearn**

| TPAMI 2020 | CVPR 2020 | TASLP 2020 | ICASSP 2020 | ICPR 2021 | ICPR-W 2021 | TNNLS 2021 | ICASSP 2022 | IJCV 2022 |

Its the not the Destination,
it's the journey.

*Merci beaucoup, Radu !*