

---

# Information Acquisition and Distributions of Functions with Deep Generative Models

---

**Ignacio Peis**

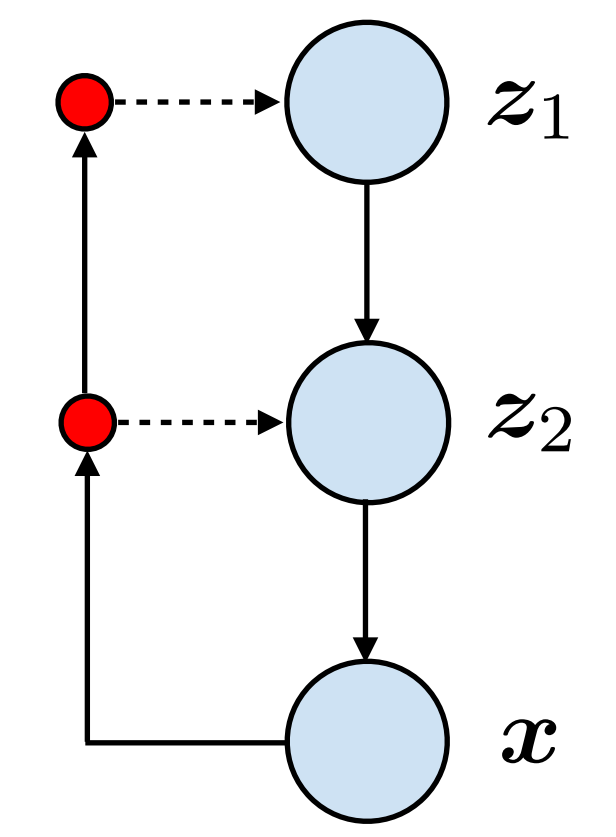
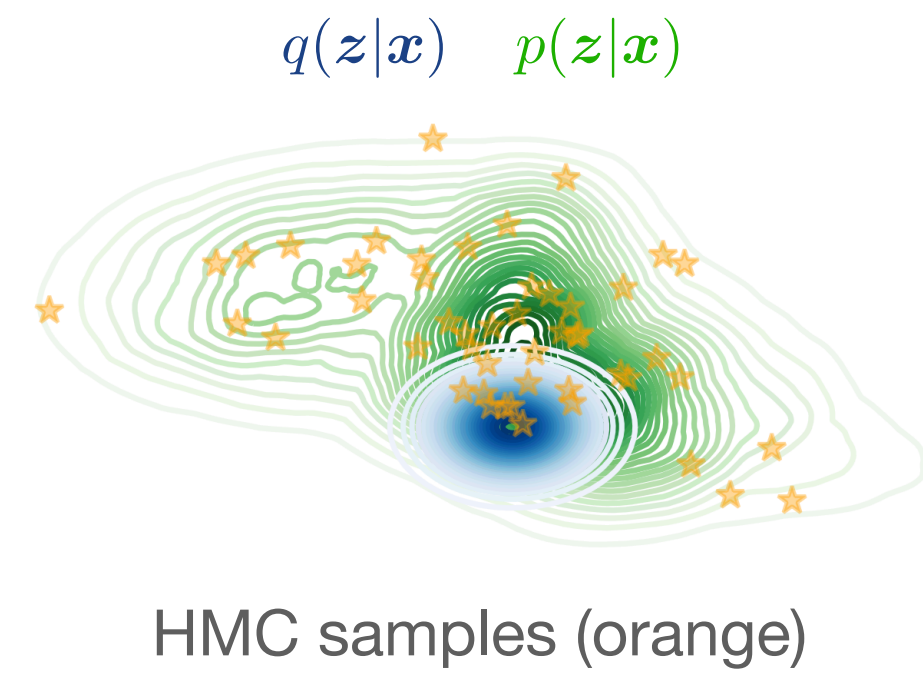
Universidad Carlos III de Madrid  
[ipeis@tsc.uc3m.es](mailto:ipeis@tsc.uc3m.es)



# Introduction

## Research Questions: Part I

- One-layered VAEs approximate inference can be improved via Markov Chain Monte Carlo [1-4].
  1. Could we leverage MCMC methods for Hierarchical VAEs?
  2. If so, could we improve incomplete data handling with MCMC?

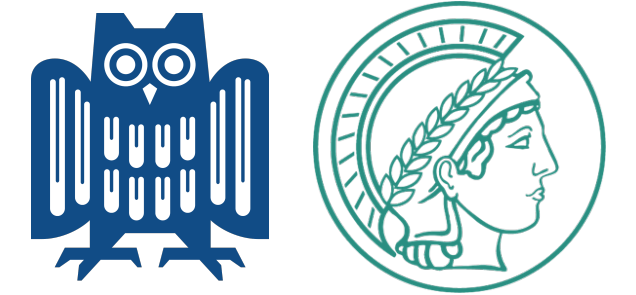
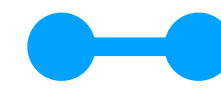


[1] Campbell et al., 2019

[2] Caterini et al., 2018

[3] Salimans et al., 2018

[4] Ruiz et al., 2021



# Introduction

## Research Questions: Part II

- VAEs are successful in generating structured data.
  1. Could we generate non-structured data via **functions** [5,6] using a VAE framework?
  2. Can we encode weights of a Neural Network?

[5] Dupont et al., 2022

[6] Dupont et al., 2022



# Variational Autoencoders

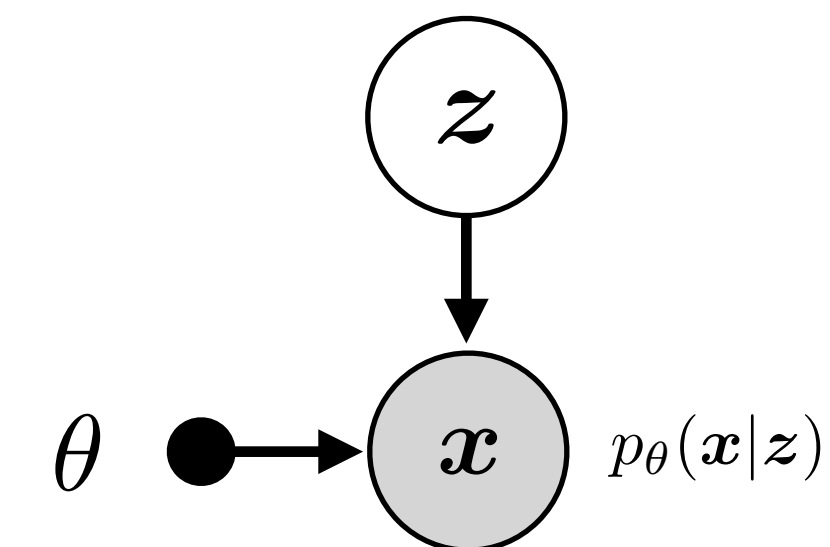
## Definition<sup>[7]</sup>

- Generative, explicit density models with intractable marginal log-likelihood.

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int p_{\theta}(\mathbf{x}|\mathbf{z}) p(\mathbf{z}) d\mathbf{z}$$

- **Intractable**, due to the complexity added by the NNs.
- Posterior distribution:

$$p(\mathbf{z}|\mathbf{x}) = \frac{p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{\underline{p_{\theta}(\mathbf{x})}}$$



<sup>[7]</sup> Kigima et al., 2013





# Variational Autoencoders

## Amortised Variational Inference (I)

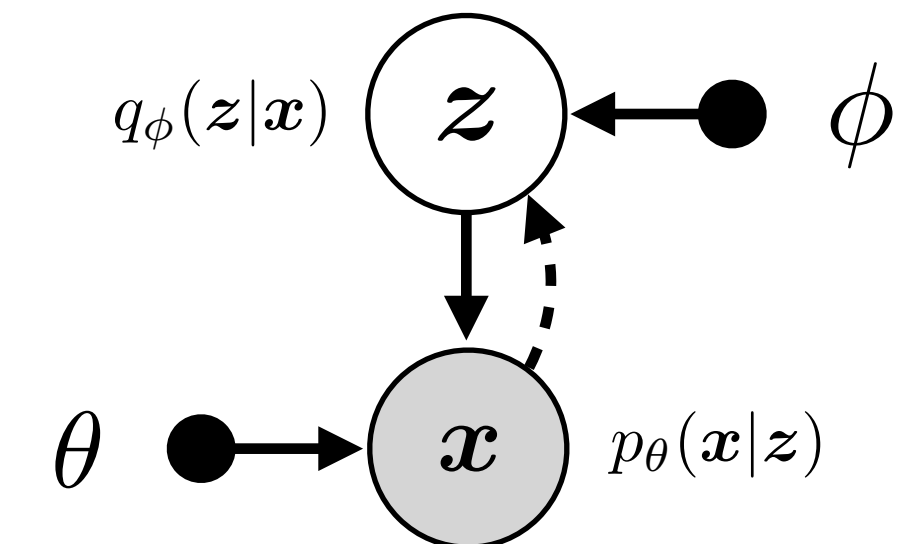
- Learn a **Gaussian** approximation of the posterior using observed data by minimizing

$$D_{KL} (q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}|\mathbf{x}))$$

- Which is equivalent to maximizing

$$\mathcal{L}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})},$$

named **Evidence Lower Bound (ELBO)**.






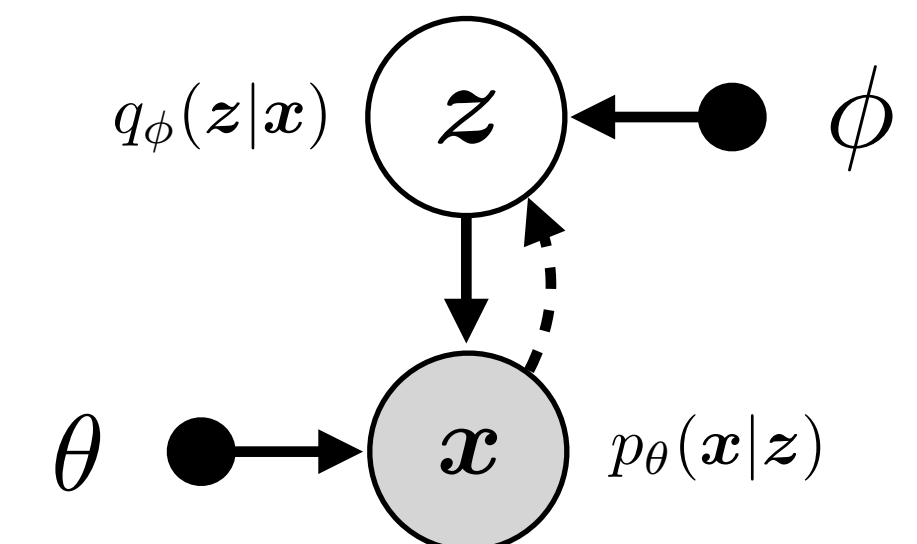
# Variational Autoencoders

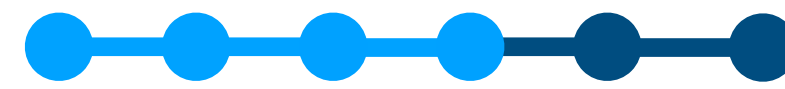
## Amortised Variational Inference (I)

The **ELBO** is typically expressed as

$$\mathcal{L}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \log p_{\theta}(\mathbf{x}|\mathbf{z}) - D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})) \leftarrow \mathcal{N}(\mathbf{0}, \mathbf{I})$$

 **Decoder**      **Encoder**





# Variational Autoencoders

## Amortised Variational Inference (II)

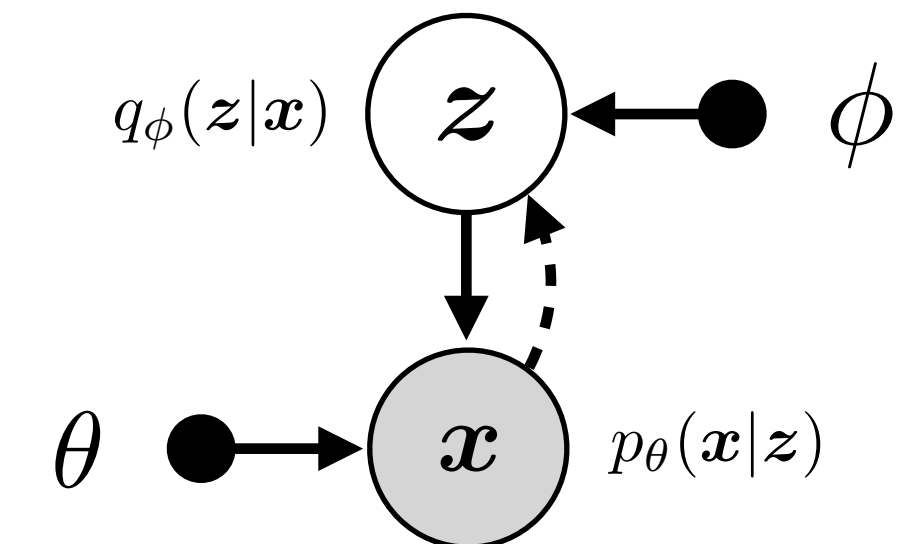
- First **reconstruction** term requires an MC estimator:

$$\hat{\mathcal{L}}(\mathbf{x}) = \frac{1}{S} \sum_{i=1}^S \left( \log p_{\theta}(\mathbf{x} | \mathbf{z}^{(s)}) \right) - D_{KL}(q_{\phi}(\mathbf{z} | \mathbf{x}) || p(\mathbf{z}))$$

Reparameterization trick<sup>[7]</sup>

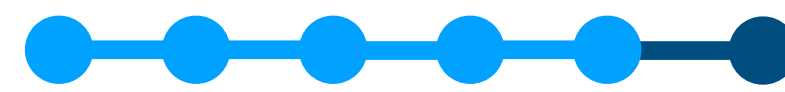
$$\mathbf{z}^{(s)} = f_{\mu}(\mathbf{x}) + f_{\sigma}(\mathbf{x}) \cdot \epsilon^{(s)}$$
$$\epsilon^{(s)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Works reasonably good even for  $S=1$



- Second **regularization** term can be computed in close form.

<sup>[7]</sup> Kigama et al., 2013



# Variational Autoencoders

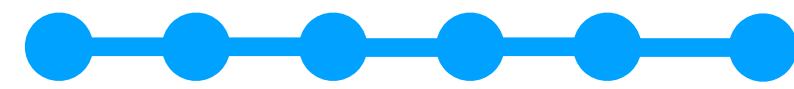
## Training

For each batch of  $B$  samples:

1. Encode to the parameters of the approximate posteriors  $q_\phi(\mathbf{z}|\mathbf{x}_i)$  .
2. Draw a sample from each  $q_\phi(\mathbf{z}|\mathbf{x}_i)$  .
3. Optimization step on  $\theta$  and  $\phi$ .

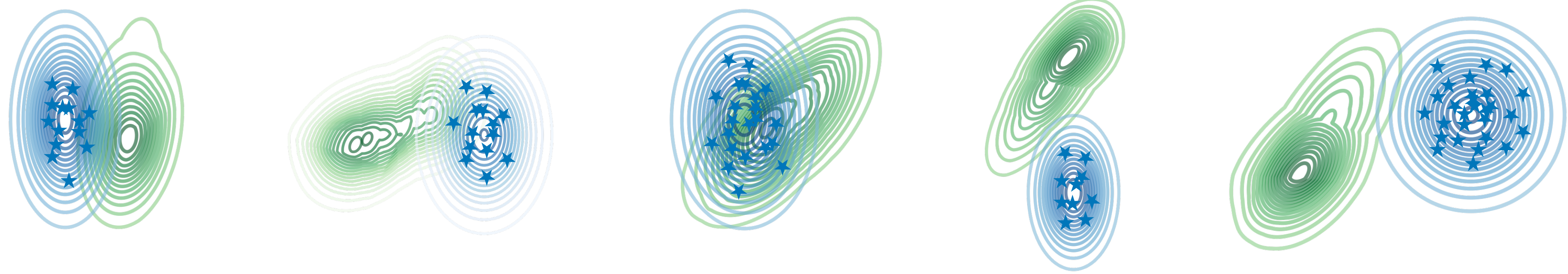
$$\nabla_{(\theta, \phi)} \left( \frac{1}{B} \sum_{i=1}^B (\log p_\theta(\mathbf{x}_i|\mathbf{z}) - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}_i)||p(\mathbf{z}))) \right)$$





# Variational Autoencoders

## Approximate Inference



$q(z|x)$     $p(z|x)$

- Can we get better samples that follow the green contour? ✓

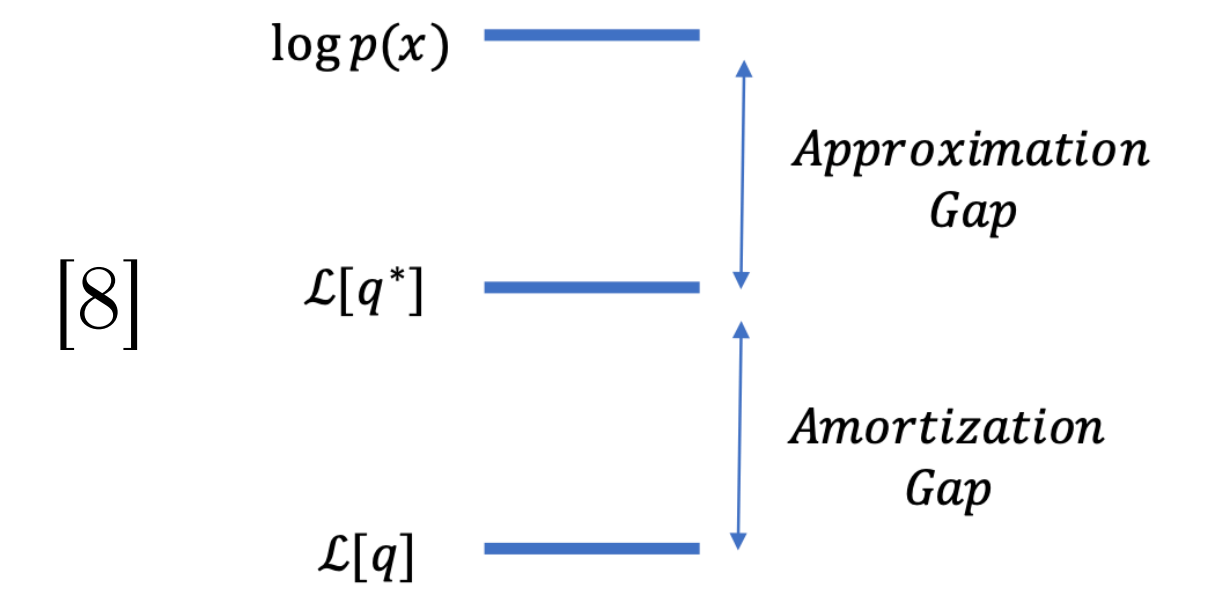
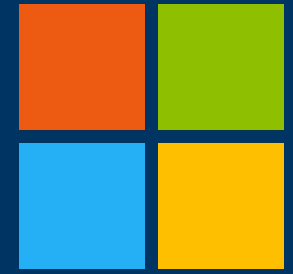


Figure 1. Gaps in Inference

[8] Cremer et al., 2018



## Part I

# Missing Data Imputation and Acquisition with Deep Hierarchical Models and Hamiltonian Monte Carlo

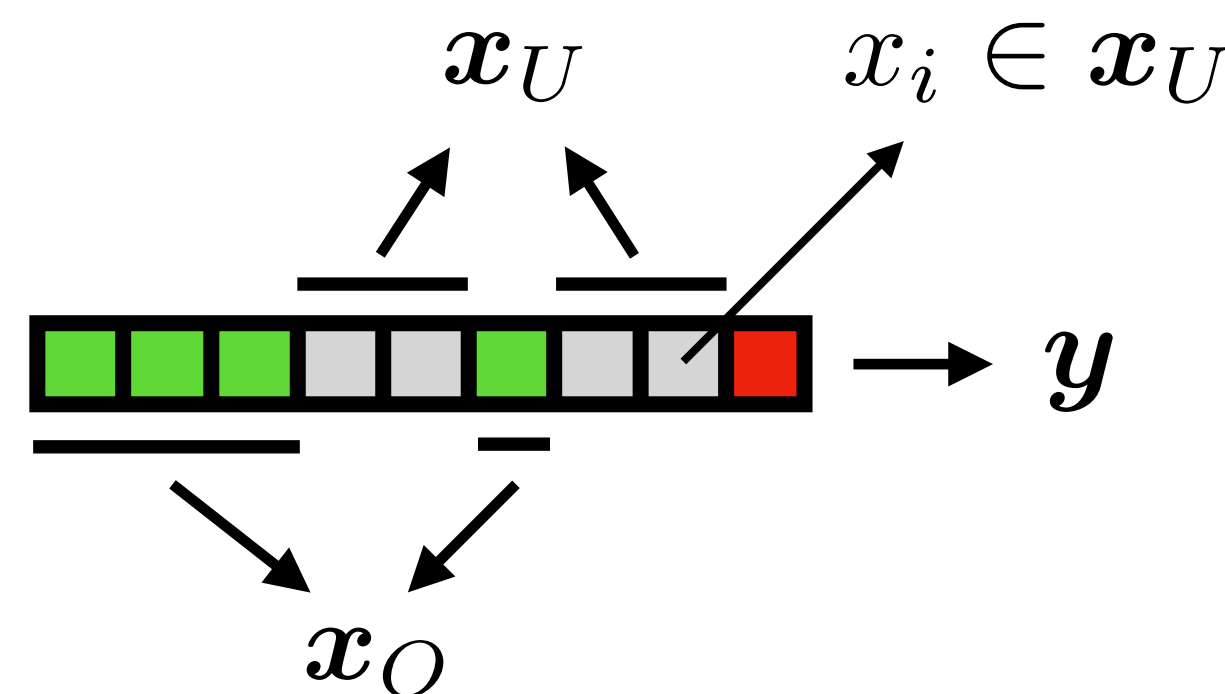


# Challenges

## Enhance information acquisition with VAEs

- Discovery of high-value information.
- Bayesian reward function [9] as an expected gain of information:

$$R(i, \mathbf{x}_O) = \mathbb{E}_{\mathbf{x}_i \sim p(\mathbf{x}_i | \mathbf{x}_O)} D_{\text{KL}} [p(\mathbf{y} | \mathbf{x}_i, \mathbf{x}_O) || p(\mathbf{y} | \mathbf{x}_O)]$$



[9] Bernardo et al., 1979



# Challenges

## Enhance information acquisition with VAEs

- Approximated in [10,11] by transforming the reward into the latent space:

$$\hat{R}(i, \mathbf{x}_o) = \mathbb{E}_{\hat{p}(\mathbf{x}_i|\mathbf{x}_o)} D_{KL} [q(\mathbf{z}|\mathbf{x}_i, \mathbf{x}_o) || q(\mathbf{z}|\mathbf{x}_o)] - \mathbb{E}_{\hat{p}(\mathbf{y}, \mathbf{x}_i|\mathbf{x}_o)} D_{KL} [q(\mathbf{z}|\mathbf{y}, \mathbf{x}_i, \mathbf{x}_o) || q(\mathbf{z}|\mathbf{y}, \mathbf{x}_o)]$$

- These methods are based on **Gaussian** approximations of the true posterior.

[10] Ma et al., 2018    [11] Ma et al., 2020





# Challenges

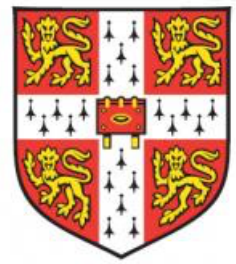
## Improve missing data imputation with VAEs

- Imputation under a VAE framework [10-13]:

$$p(\mathbf{x}_U | \mathbf{x}_O) = \mathbb{E}_{p(\mathbf{z} | \mathbf{x}_O)} [p(\mathbf{x}_U | \mathbf{z})] \approx \mathbb{E}_{q(\mathbf{z} | \mathbf{x}_O)} [p(\mathbf{x}_U | \mathbf{z})]$$

- Also based on **Gaussian** approximations of the true posterior.

[10] Ma et al., 2018   [11] Ma et al., 2020   [12] Nazabal et al., 2020   [13] Mattei et al., 2020



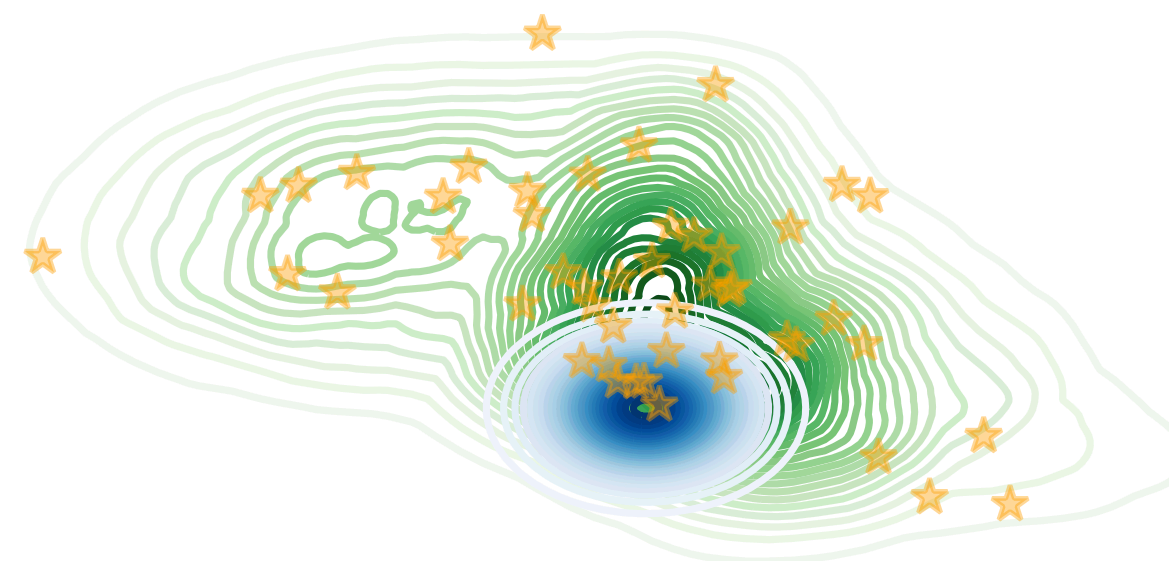
# Challenges

## Sampling-based methods for incomplete data related tasks



Expectations over the intractable posterior should leverage a well-designed MCMC approximation method when compared to a Gaussian-based approximation.

$q(z|x)$   $p(z|x)$



HMC samples (orange)

- Imputation:**  $p(\mathbf{x}_U|\mathbf{x}_O) \approx \mathbb{E}_{q^{(T)}(\epsilon|\mathbf{x}_O)}[p(\mathbf{x}_U|\epsilon)].$
- Prediction:**  $p(\mathbf{y}|\mathbf{x}_O) \approx \mathbb{E}_{q^{(T)}(\epsilon|\mathbf{x}_O)}[p(\mathbf{y}|\epsilon, \mathbf{x}_O, \hat{\mathbf{x}}_U)].$
- Sampling-based active learning.**

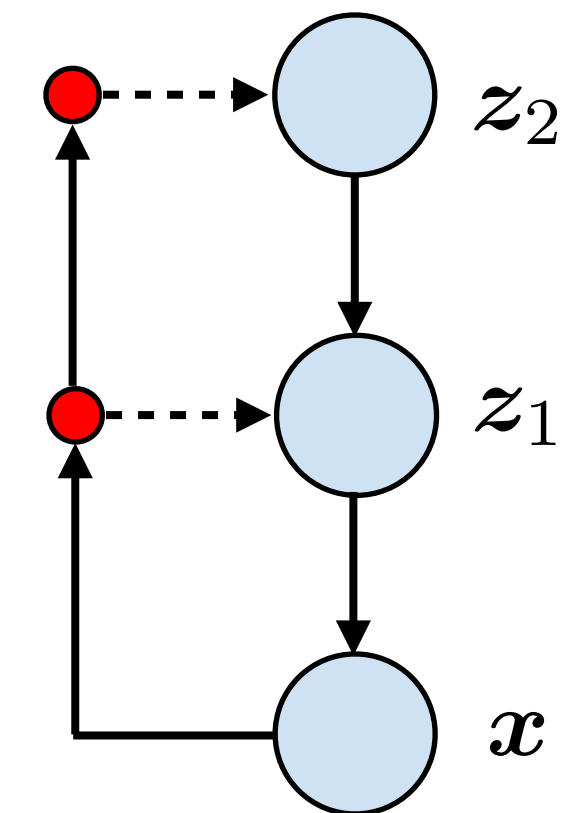


# Challenges Hierarchical VAEs

- Hierarchical VAEs are successful at increasing flexibility.

$$p(\mathbf{z}_L) \prod_{l=1}^{L-1} p(\mathbf{z}_l | \mathbf{z}_{l+1})$$

- The hierarchy allows for modelling:
  - Abstract to specific generative factors.
  - Global to local generative factors.



[14]



[14] Child, 2020

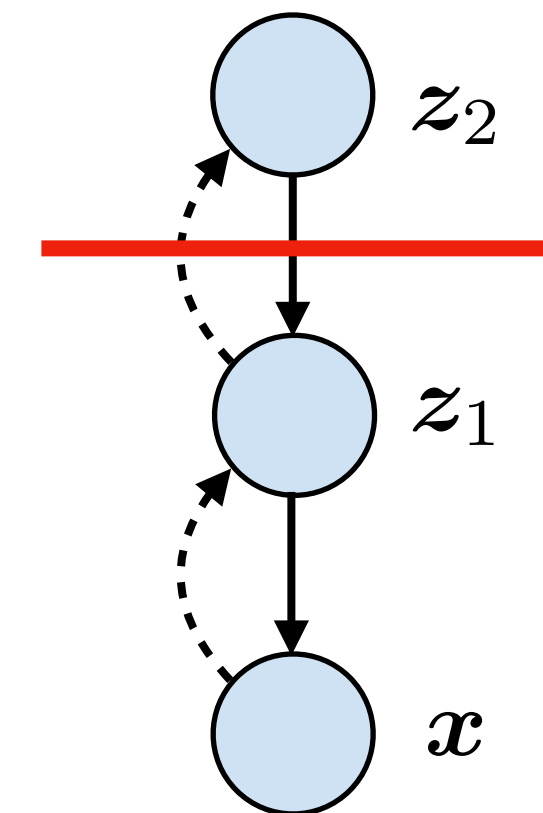


# Challenges Hierarchical VAEs

- **Delicate inference** (posterior collapse).

$$ELBO(\mathbf{x}) = \mathbb{E}_{Q(\mathbf{z}_1, \mathbf{z}_2 | \mathbf{x})} \left[ \ln p(\mathbf{x} | \mathbf{z}_1) - KL[q(\mathbf{z}_1 | \mathbf{x}) || p(\mathbf{z}_1 | \mathbf{z}_2)] - KL[q(\mathbf{z}_2 | \mathbf{z}_1) || p(\mathbf{z}_2)] \right]$$

$$q(\mathbf{z}_2 | \mathbf{z}_1) \approx p(\mathbf{z}_2) \approx \mathcal{N}(0, 1)$$





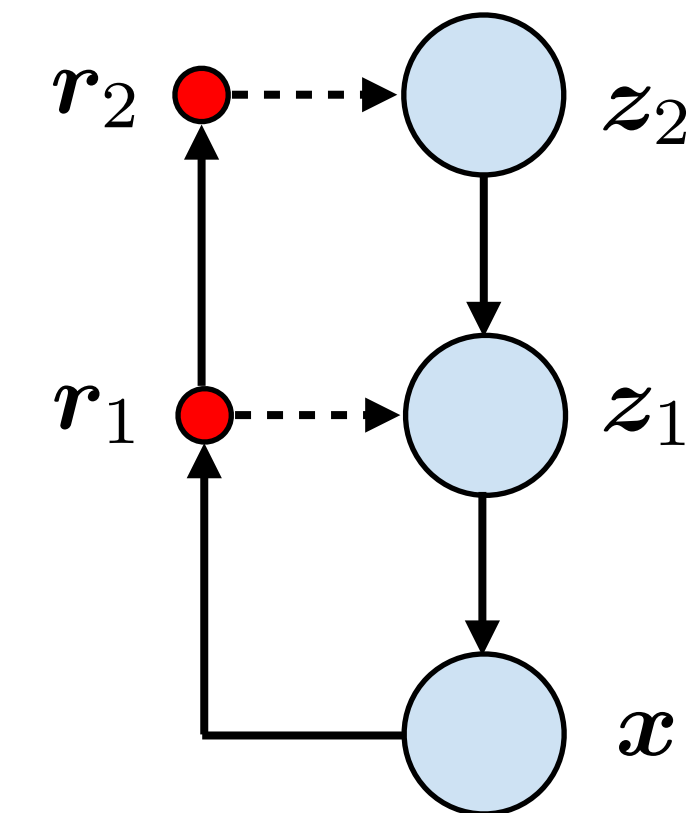


# Challenges Hierarchical VAEs

- Solution: top-down inference.

$$Q(\mathbf{z}_1, \mathbf{z}_2 | \mathbf{x}) = q(\mathbf{z}_1 | \mathbf{z}_2, \mathbf{x})q(\mathbf{z}_2 | \mathbf{x}).$$

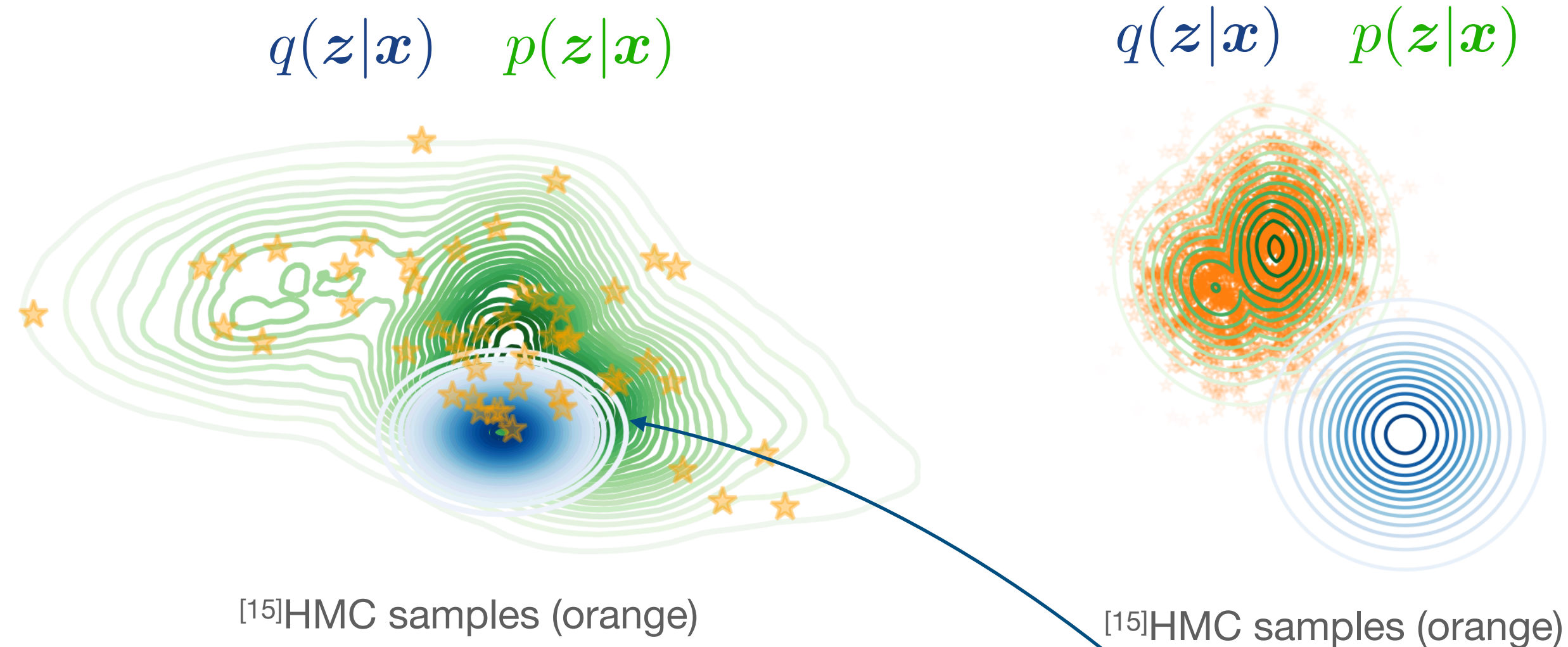
- ✓ Posterior collapse is relaxed.
- ◆ Inference bias worsens with latent dimensionality.





# Challenges

## MCMC for improving inference in VAEs



$$\mathcal{L}(\mathbf{x}; \theta, \phi) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x} | \mathbf{z})] - D_{KL}(q_\phi(\mathbf{z} | \mathbf{x}) || p(\mathbf{z}))$$

$$\approx \frac{1}{S} \sum_{s=1}^S \log p_\theta(\mathbf{x} | \mathbf{z}^{(s)}) - D_{KL}(q_\phi(\mathbf{z} | \mathbf{x}) || p(\mathbf{z}))$$

[15] Peis et al., 2021



# Challenges

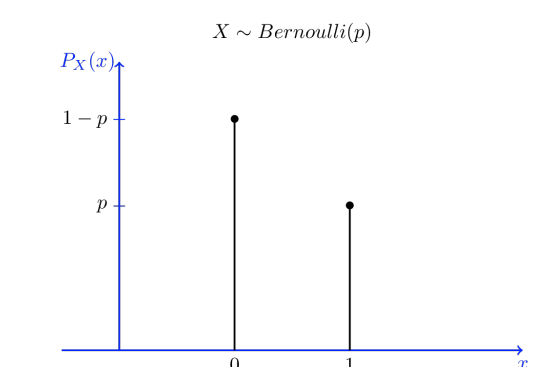
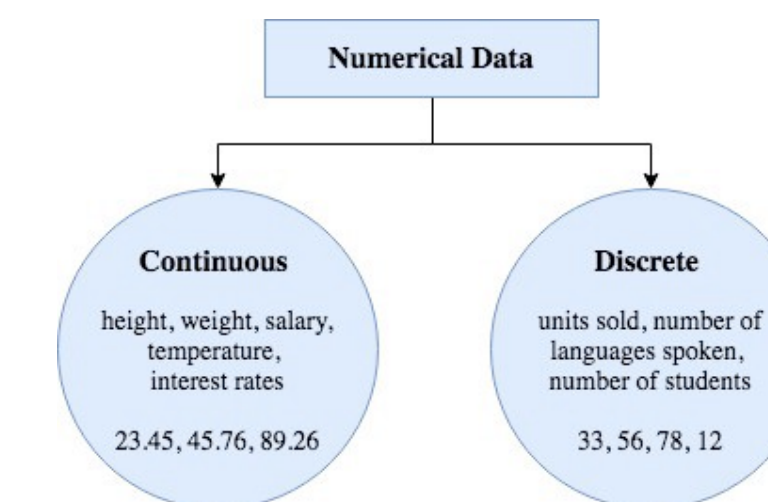
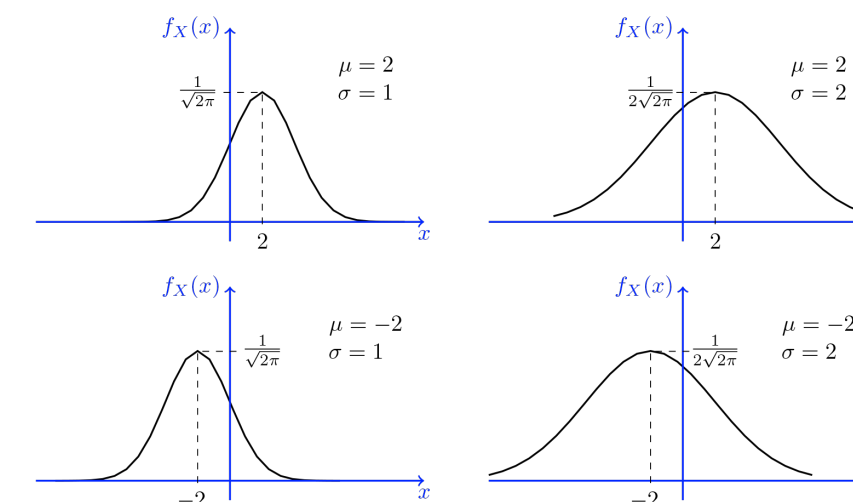
## Handling partial heterogeneous data

- Factorization over dimensions<sup>[10-13]</sup>:

$$\mathcal{L}(\mathbf{x}) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \sum_{d=1}^D \mathbb{I}(x_d \in \mathbf{x}_O) \log p_\theta(x_d|\mathbf{z}) \right] - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}_O)||p(\mathbf{z}))$$

- Naïve approach for heterogeneous<sup>[12,13]</sup>: use a different likelihood per dimension.

- Problem:** imbalanced likelihoods.



[10] Ma et al., 2018    [11] Ma et al., 2020    [12] Nazabal et al., 2020    [13] Mattei et al., 2020



# Challenges

## Handling partial heterogeneous data

- **Solution<sup>[11]</sup>**: learn first  $D$  marginal VAEs  $(\theta_d, \gamma_d)$ :

$$\mathcal{L}_d(x_d; \{\theta_d, \gamma_d\}) = \mathbb{I}(x_d \in \mathbf{x}_O) \mathbb{E}_{q_{\gamma_d}(z_d|x_d)} \log \frac{p_{\theta_d}(x_d, z_d)}{q_{\gamma_d}(z_d|x_d)}$$

and a joint dependency VAE  $(\theta, \phi)$  on the marginally encoded data:

$$\mathcal{L}(\mathbf{z}) = \mathbb{E}_{q_{\phi}(\mathbf{h}|\mathbf{z})} \left[ \sum_{d=1}^D \mathbb{I}(z_d \in \mathbf{z}_O) \mathbb{E}_{q_{\gamma_d}(z_d|x_d)} [\log p_{\theta}(z_d|\mathbf{h})] \right] - D_{KL}(q_{\phi}(\mathbf{h}|\mathbf{z}_O) || p(\mathbf{h}))$$

- Interdependencies between heterogeneous variables are better captured by the dependency VAE.

<sup>[11]</sup> Ma et al., 2020

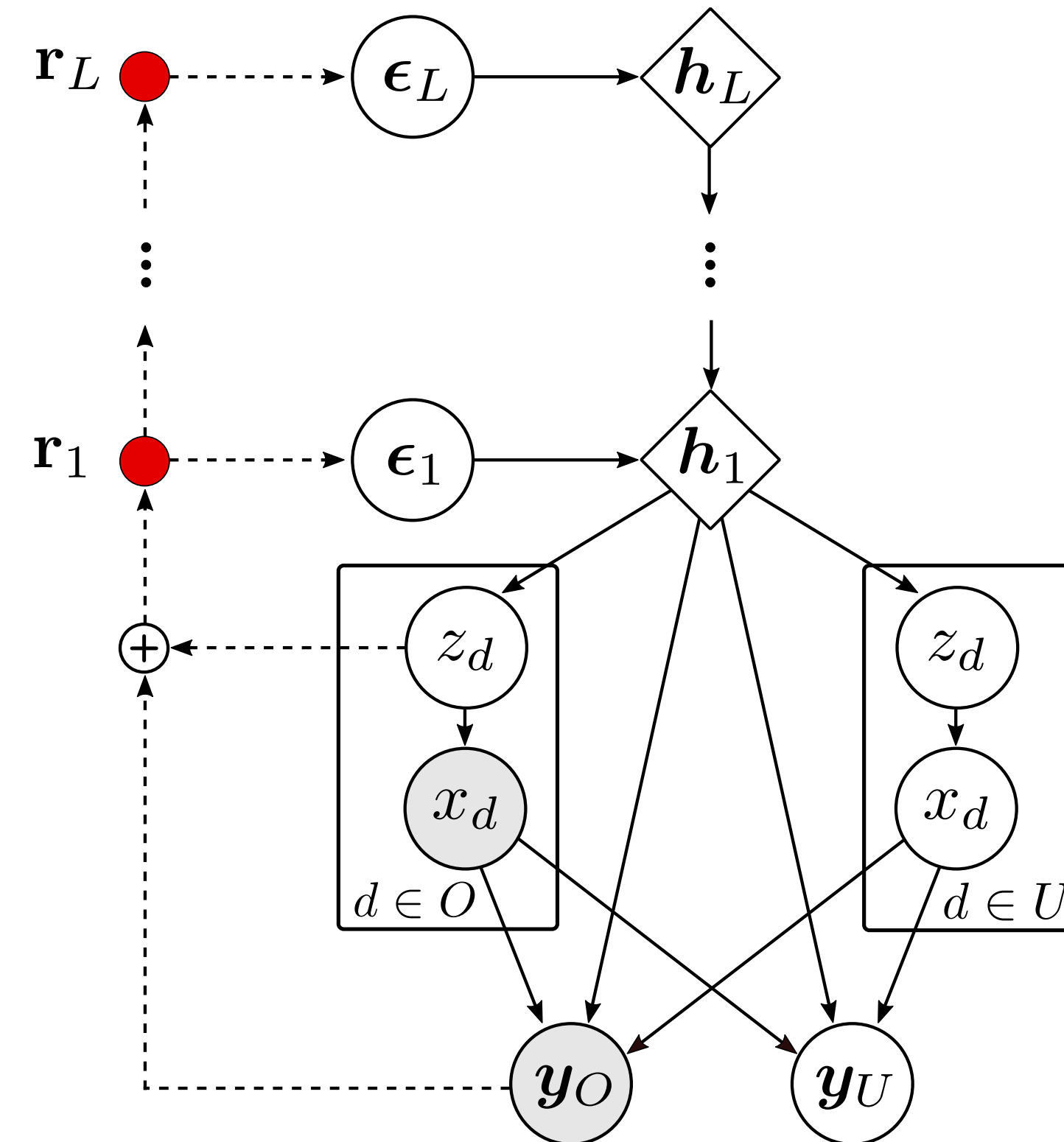


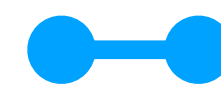


# Contributions

## Hierarchical Hamiltonian VAE for mixed-type incomplete data (HH-VAEM)

- Increased flexibility by using hierarchical latent space.
- Improved inference by means of automatically tuned HMC.
- Reparameterization for well-posed HMC on relaxed posterior.
- Heterogeneous data handling.
- More accurate imputation and prediction.
- More effective information acquisition.





# Contributions (II)

## Sampling-based method for active information acquisition

- **Sampling**-based estimator [16] of the Mutual Information:

$$\begin{aligned} R(i, \mathbf{x}_O) &= D_{\text{KL}} [p(\mathbf{y}, x_i | \mathbf{x}_O) || p(\mathbf{y} | \mathbf{x}_O) p(x_i | \mathbf{x}_O)] = \mathcal{I}(\mathbf{y}; x_i | \mathbf{x}_O) = \\ &= \iint_{x_i, \mathbf{y}} p_{x_i, \mathbf{y} | \mathbf{x}_O}(x_i, \mathbf{y} | \mathbf{x}_O) \log \left( \frac{p_{x_i, \mathbf{y} | \mathbf{x}_O}(x_i, \mathbf{y} | \mathbf{x}_O)}{p_{x_i | \mathbf{x}_O}(x_i | \mathbf{x}_O) p_{\mathbf{y} | \mathbf{x}_O}(\mathbf{y} | \mathbf{x}_O)} \right) \end{aligned}$$

$$\hat{I}(\mathbf{y}; x_i | \mathbf{x}_O) \approx \sum_{ij} p_{x_i, \mathbf{y} | \mathbf{x}_O}(i, j) \log \frac{p_{x_i, \mathbf{y} | \mathbf{x}_O}(i, j)}{p_{x_i | \mathbf{x}_O}(i) p_{\mathbf{y} | \mathbf{x}_O}(j)}$$

- ✓ Avoids the Gaussian approximation
- ✓ Efficient, easy parallelization.

[16] Kraskov et al., 2004



# Method

## Hamiltonian Monte Carlo

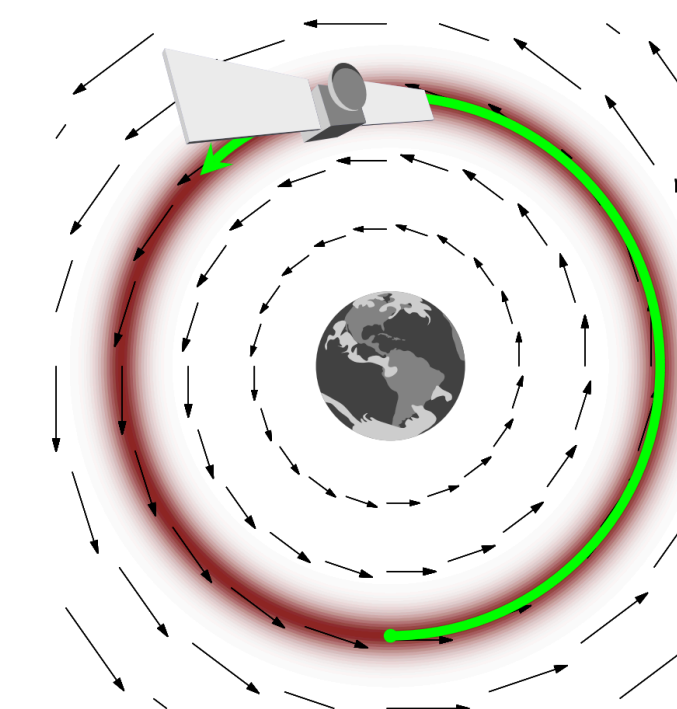
- Sample from complex distributions via unnormalised targets

### 1. Phase and momentum space

$$p(\mathbf{z}, \mathbf{r}) = p(\mathbf{r}|\mathbf{z})p(\mathbf{z})$$

$$H(\mathbf{z}, \mathbf{r}) = -\log p(\mathbf{z}, \mathbf{r}) = -\log p(\mathbf{r}|\mathbf{z}) - \log p(\mathbf{z}) = K(\mathbf{r}, \mathbf{z}) + V(\mathbf{z})$$

$$H(\mathbf{z}, \mathbf{r}) = -\log p^*(\mathbf{z}) + \frac{1}{2}\mathbf{r}^T \mathbf{M}^{-1}\mathbf{r}.$$

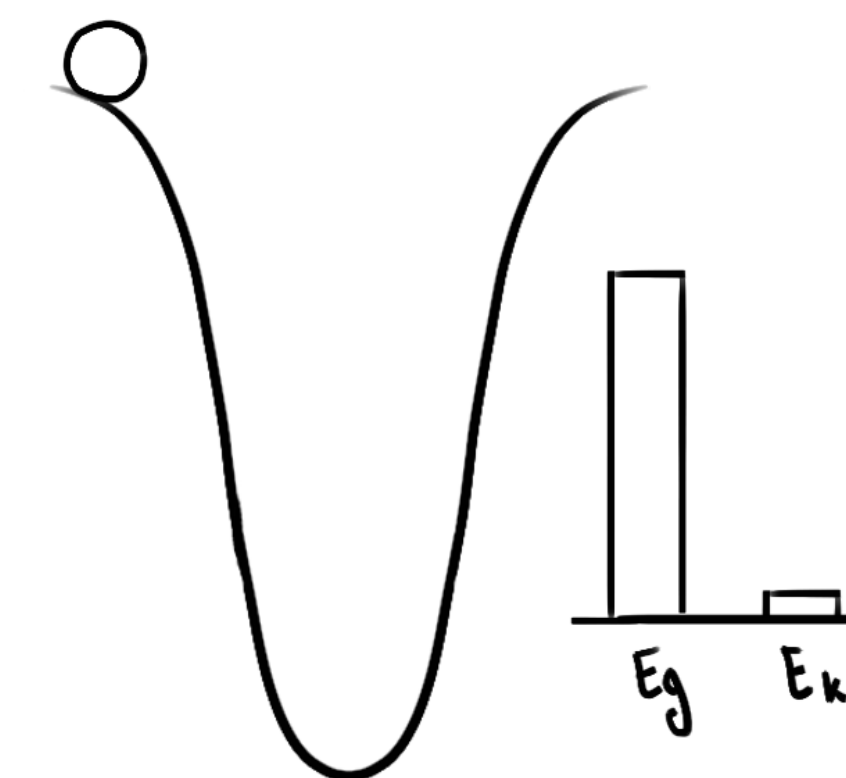


### 2. Hamiltonian equations

$$\begin{aligned} \frac{d\mathbf{z}}{dt} &= +\frac{\partial H}{\partial \mathbf{r}} = \frac{\partial K}{\partial \mathbf{r}} \\ \frac{d\mathbf{r}}{dt} &= -\frac{\partial H}{\partial \mathbf{z}} = -\frac{\partial K}{\partial \mathbf{z}} - \frac{\partial V}{\partial \mathbf{z}} \end{aligned}$$



$$\begin{aligned} \mathbf{r}_{l+\frac{1}{2}} &= \mathbf{r}_l + \frac{1}{2}\boldsymbol{\phi} \odot \nabla_{\mathbf{z}_l} \log p^*(\mathbf{z}_l), \\ \mathbf{z}_{l+1} &= \mathbf{z}_k + \mathbf{r}_{l+\frac{1}{2}} \odot \boldsymbol{\phi} \odot \frac{1}{\mathbf{M}}, \\ \mathbf{r}_{l+1} &= \mathbf{r}_{l+\frac{1}{2}} + \frac{1}{2}\boldsymbol{\phi} \odot \nabla_{\mathbf{z}_{l+1}} \log p^*(\mathbf{z}_{l+1}), \end{aligned}$$





# Method

## Hamiltonian Monte Carlo

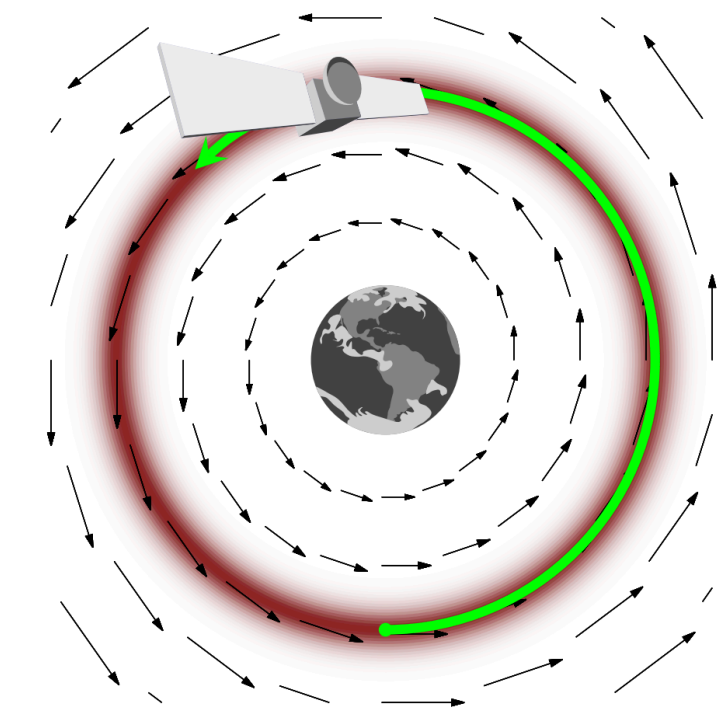
- Sample from complex distributions via unnormalised targets

### 1. Phase and momentum space

$$p(\mathbf{z}, \mathbf{r}) = p(\mathbf{r}|\mathbf{z})p(\mathbf{z})$$

$$H(\mathbf{z}, \mathbf{r}) = -\log p(\mathbf{z}, \mathbf{r}) = -\log p(\mathbf{r}|\mathbf{z}) - \log p(\mathbf{z}) = K(\mathbf{r}, \mathbf{z}) + V(\mathbf{z})$$

$$H(\mathbf{z}, \mathbf{r}) = -\log p^*(\mathbf{z}) + \frac{1}{2}\mathbf{r}^T \mathbf{M}^{-1}\mathbf{r}.$$

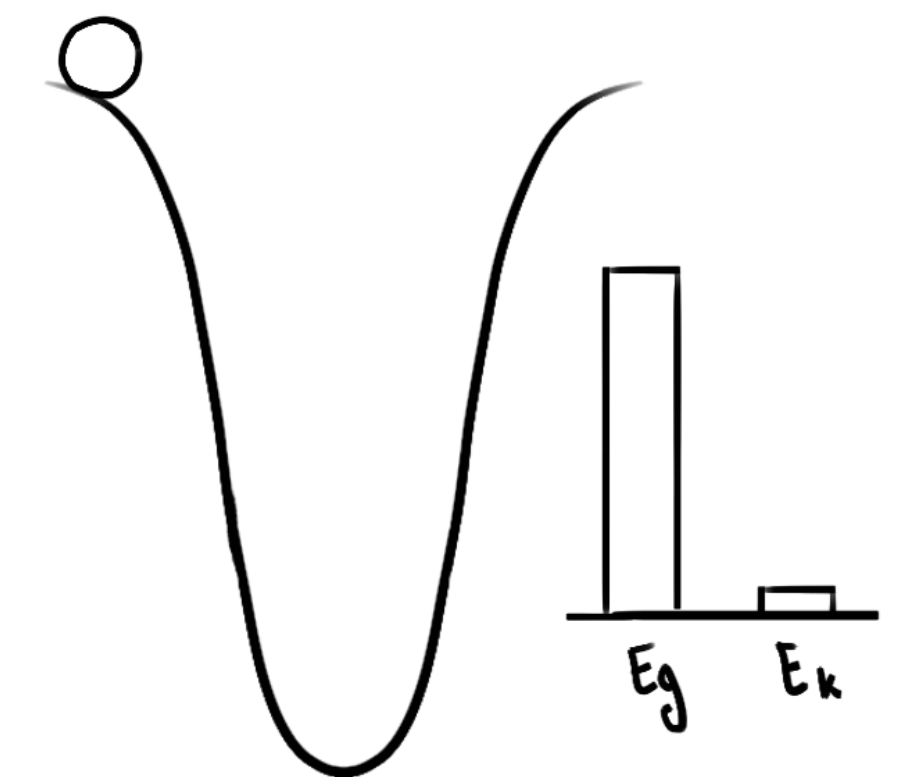


### 2. Hamiltonian equations

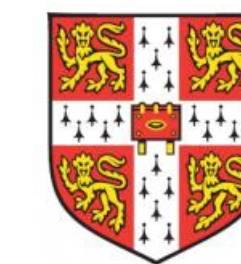
$$\begin{aligned} \frac{d\mathbf{z}}{dt} &= +\frac{\partial H}{\partial \mathbf{r}} = \frac{\partial K}{\partial \mathbf{r}} \\ \frac{d\mathbf{r}}{dt} &= -\frac{\partial H}{\partial \mathbf{z}} = -\frac{\partial K}{\partial \mathbf{z}} - \frac{\partial V}{\partial \mathbf{z}} \end{aligned}$$



$$\begin{aligned} \mathbf{r}_{l+\frac{1}{2}} &= \mathbf{r}_l + \frac{1}{2}\boldsymbol{\phi} \odot \nabla_{\mathbf{z}_l} \log p^*(\mathbf{z}_l), \\ \mathbf{z}_{l+1} &= \mathbf{z}_k + \mathbf{r}_{l+\frac{1}{2}} \odot \boldsymbol{\phi} \odot \frac{1}{\mathbf{M}}, \\ \mathbf{r}_{l+1} &= \mathbf{r}_{l+\frac{1}{2}} + \frac{1}{2}\boldsymbol{\phi} \odot \nabla_{\mathbf{z}_{l+1}} \log p^*(\mathbf{z}_{l+1}), \end{aligned}$$







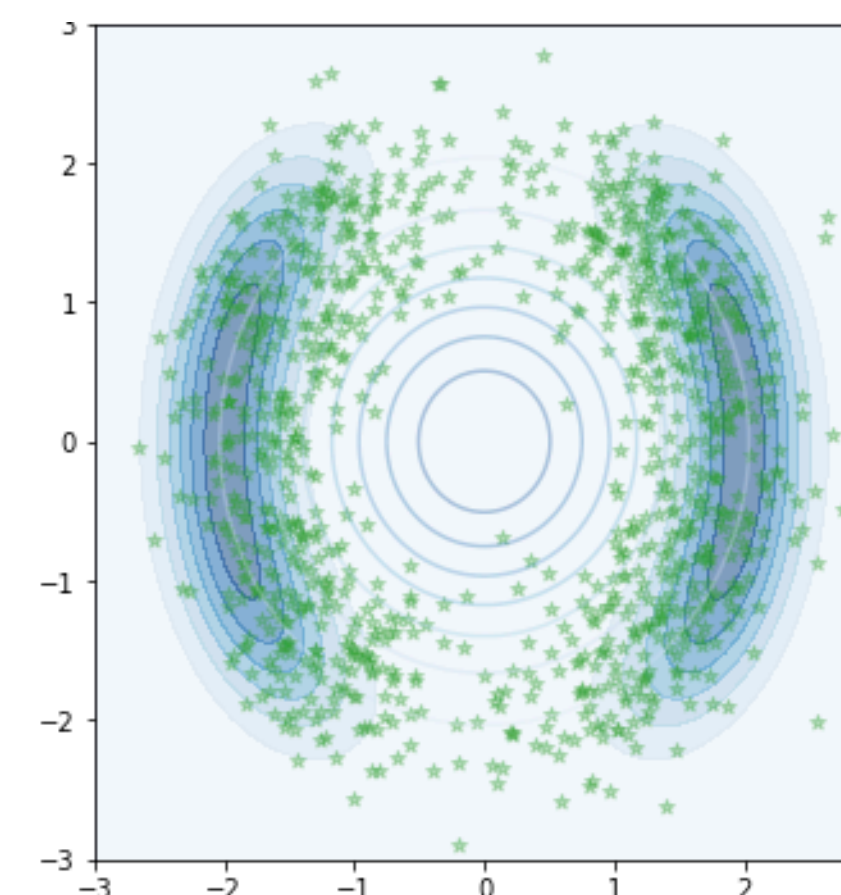
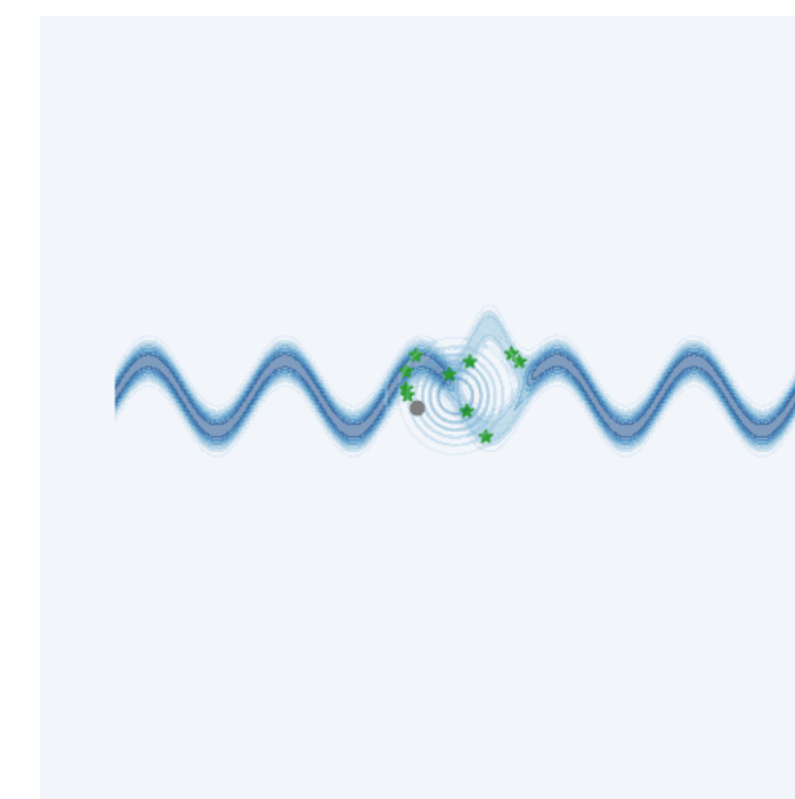
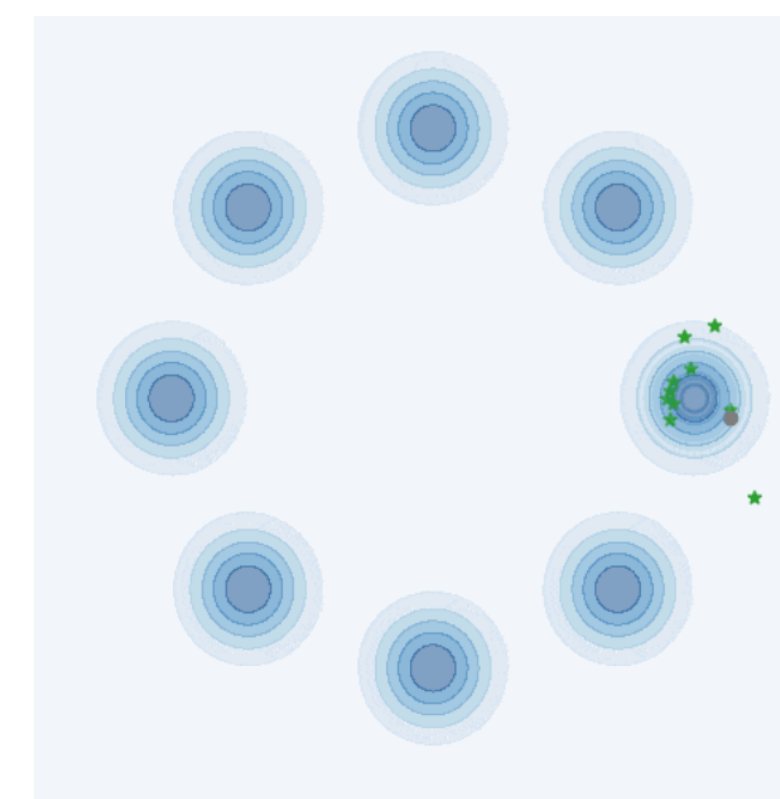
# Method

## Hamiltonian Monte Carlo

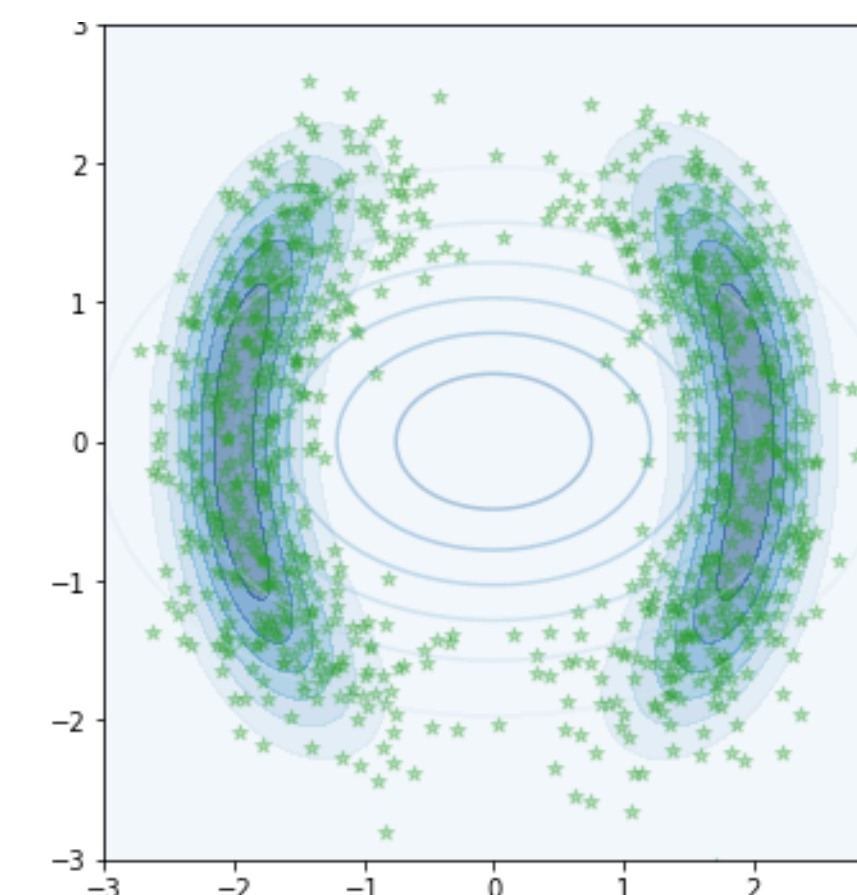
- Discrete trajectories (*chains*) of  $T$  updates, ending in:

$$q^{(T)}(\mathbf{z}|\mathbf{x}) \approx p(\mathbf{z}|\mathbf{x})$$

- Target: true posterior density.
- Needed:
  - 1. Good initial proposal (encoder).
  - 2. Well-defined hyperparameters.



Random hyperparameters



Tuned hyperparameters





# Method

## HMC Hyperparameter tuning [1]

- Tuning the hyperparameters via Variational Inference:

$$\phi^* = \operatorname{argmax}_{\phi} \mathbb{E}_{q_{\phi}^{(T)}(z)} [\log p^*(z)] + H [q_{\phi}^{(T)}(z)]$$

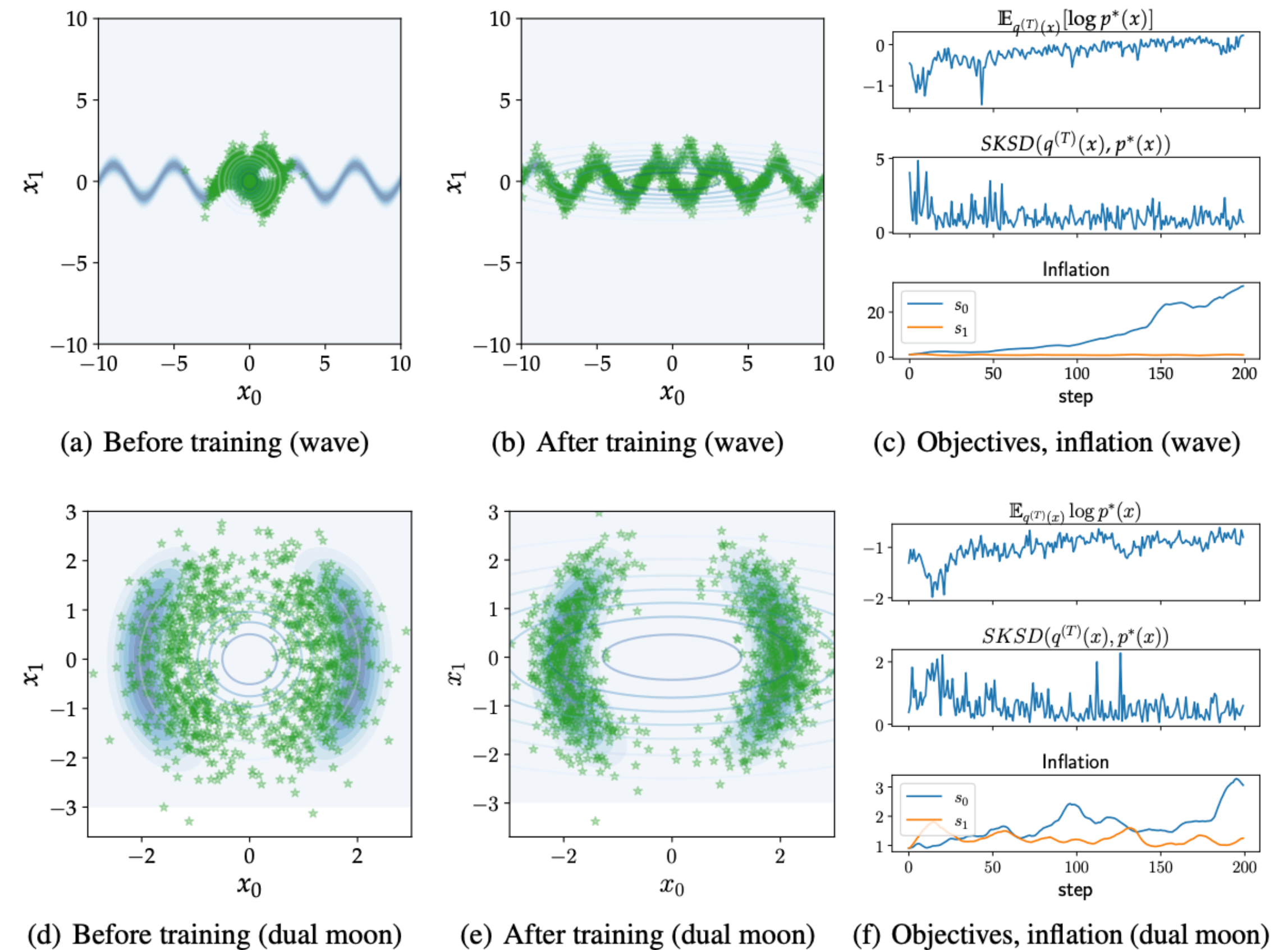
- Add an **inflation** parameter for scaling the proposal [17]

$$s^* = \operatorname{argmin}_s \operatorname{SKSD}(z^{(T)}, \nabla_z \log p^*(z))$$

- Code available at:



<https://github.com/ipeis/HMCTuning>



[1] Campbell et al., 2021 [17] Gong et al., 2020



# Method

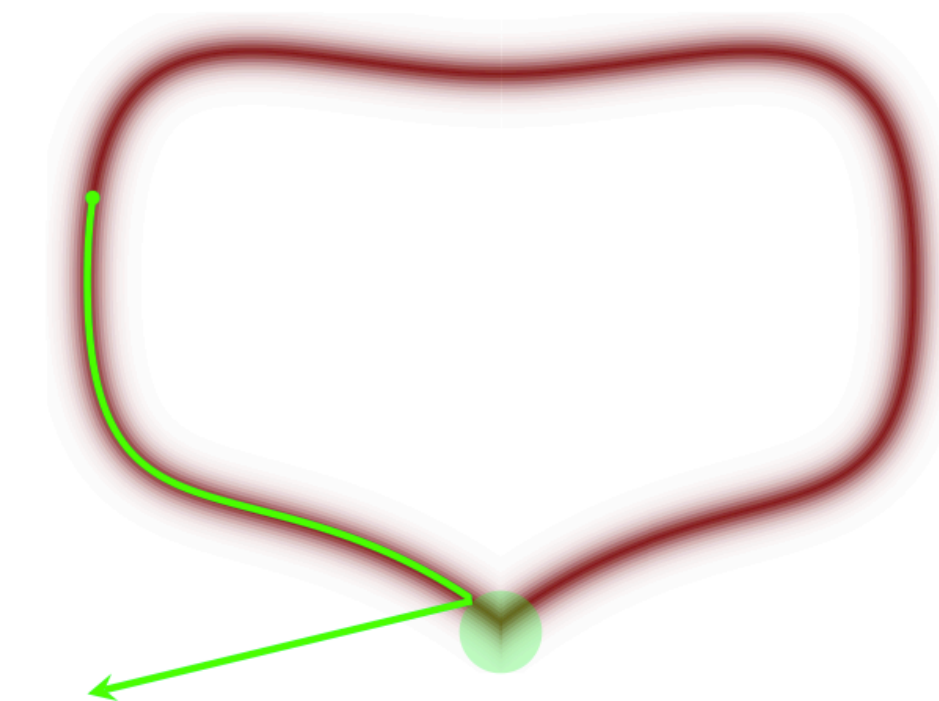
## HMC is ill-posed for Hierarchical VAEs

- Hierarchical dependencies lead to **huge gradients** [18, 19]

$$\mathbf{r}_{l+\frac{1}{2}} = \mathbf{r}_l + \frac{1}{2} \phi \odot \nabla_{z_l} \log p^*(z_l),$$

$$\mathbf{z}_{l+1} = \mathbf{z}_k + \mathbf{r}_{l+\frac{1}{2}} \odot \phi \odot \frac{1}{M},$$

$$\mathbf{r}_{l+1} = \mathbf{r}_{l+\frac{1}{2}} + \frac{1}{2} \phi \odot \nabla_{z_{l+1}} \log p^*(z_{l+1}),$$



- Samples can diverge due to integrator overflow issues.

[18] Betancourt et al., 2017 [19] Betancourt et al., 2015



# Method

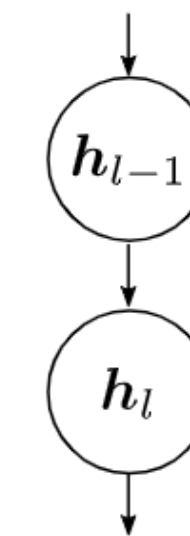
## HMC is ill-posed for Hierarchical VAEs

- **Solution:**

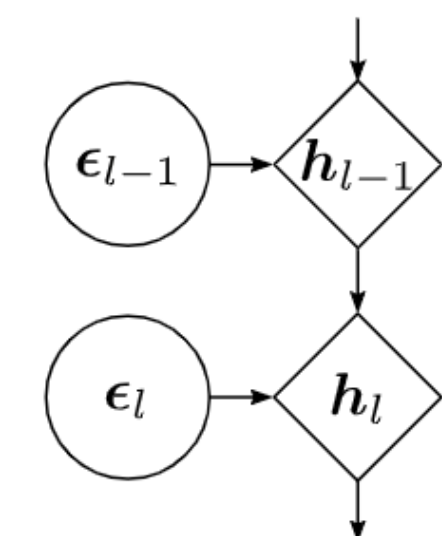
✓ Reparameterization for relaxed posterior:

$$\mathbf{h}_l = f_{\mu_l}(\mathbf{h}_{l+1}) + f_{\sigma_l}(\mathbf{h}_{l+1}) \cdot \epsilon_l$$

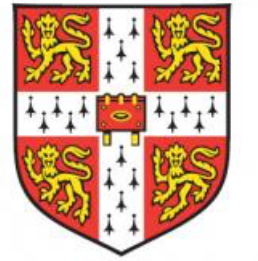
NNs with parameters  $\theta_{\mu_l} \rightarrow f_{\mu_l}, \theta_{\sigma_l} \rightarrow f_{\sigma_l}$



(a) AR hierarchy



(b) Reparameterization



# Method

## Ill-posed for HMC

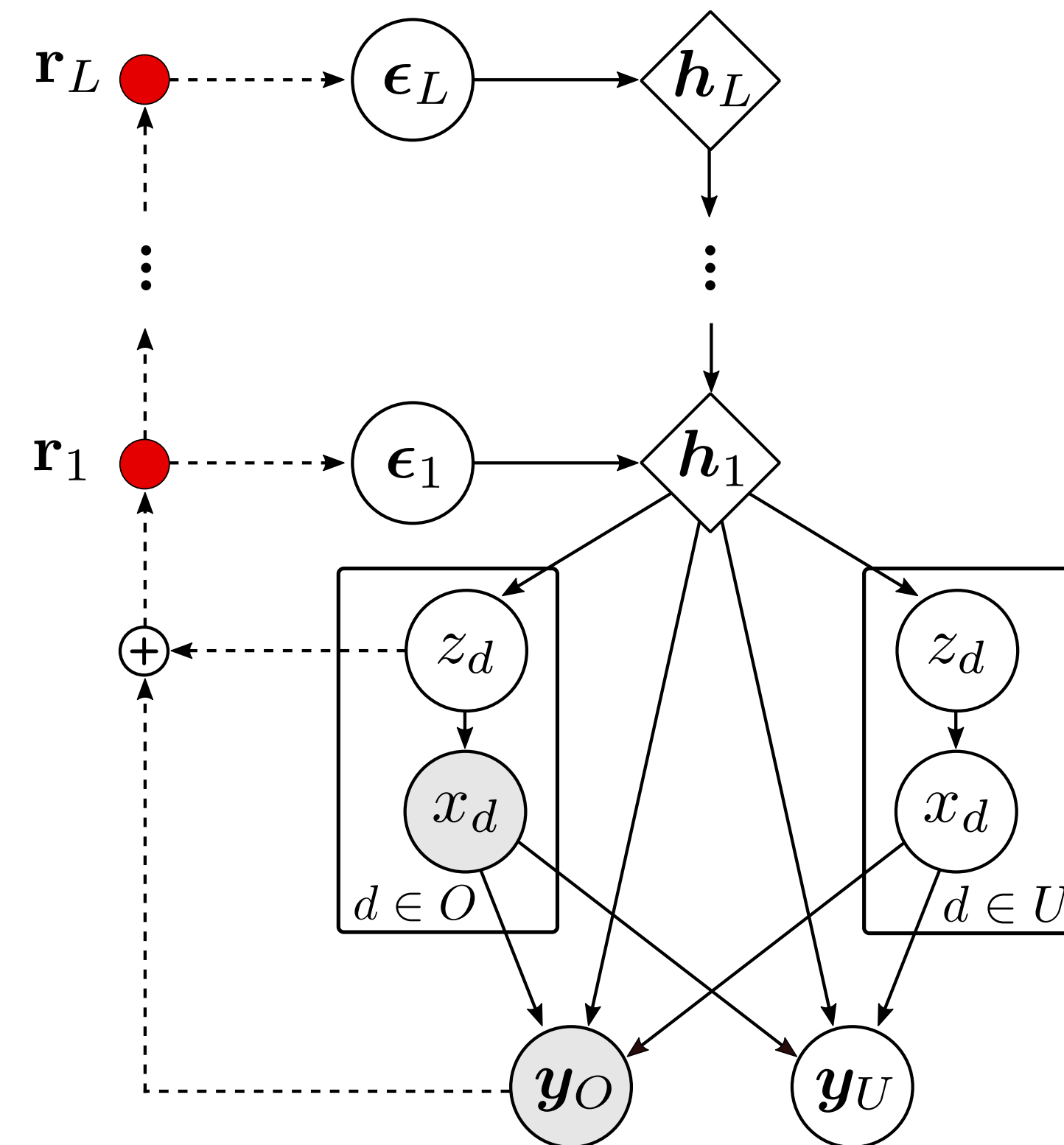
- **Solution:**

- ✓ Reparameterization for relaxed posterior:

$$\mathbf{h}_l = f_{\mu_l}(\mathbf{h}_{l+1}) + f_{\sigma_l}(\mathbf{h}_{l+1}) \cdot \epsilon_l$$

NNs with parameters  $\theta_{\mu_l} \rightarrow f_{\mu_l}, \theta_{\sigma_l} \rightarrow f_{\sigma_l}$

- ✓ Perform inference on  $\epsilon = \{\epsilon_1, \dots, \epsilon_L\}$  with standard Gaussian prior.







# Method

## Ill-posed for HMC

- **Solution:**

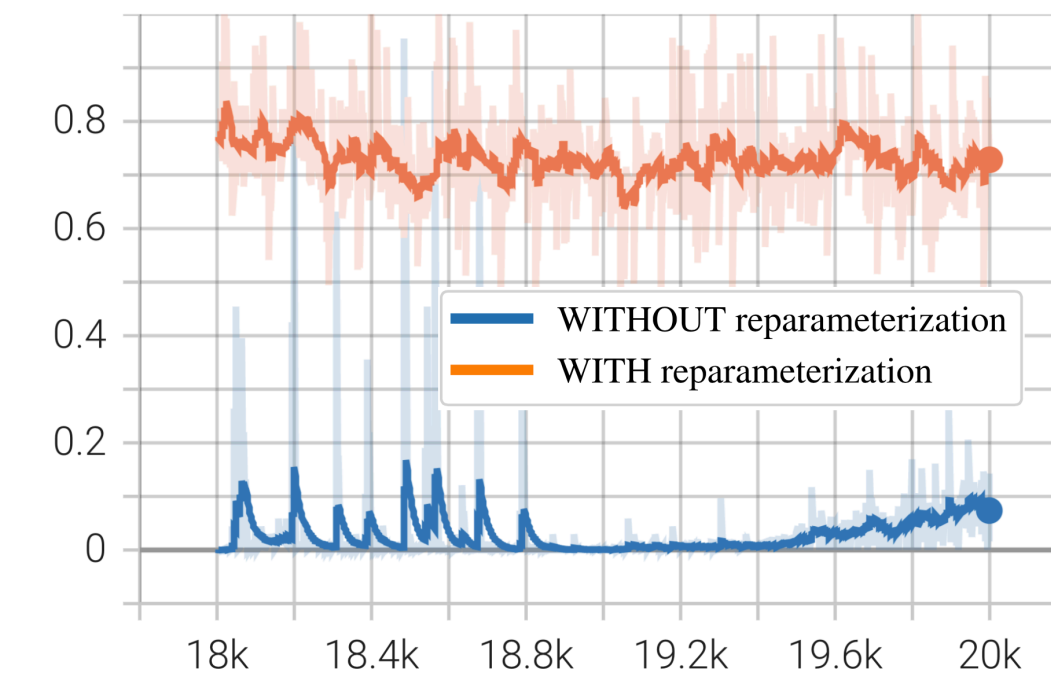
- ✓ Reparameterization for relaxed posterior:

$$\mathbf{h}_l = f_{\mu_l}(\mathbf{h}_{l+1}) + f_{\sigma_l}(\mathbf{h}_{l+1}) \cdot \epsilon_l$$

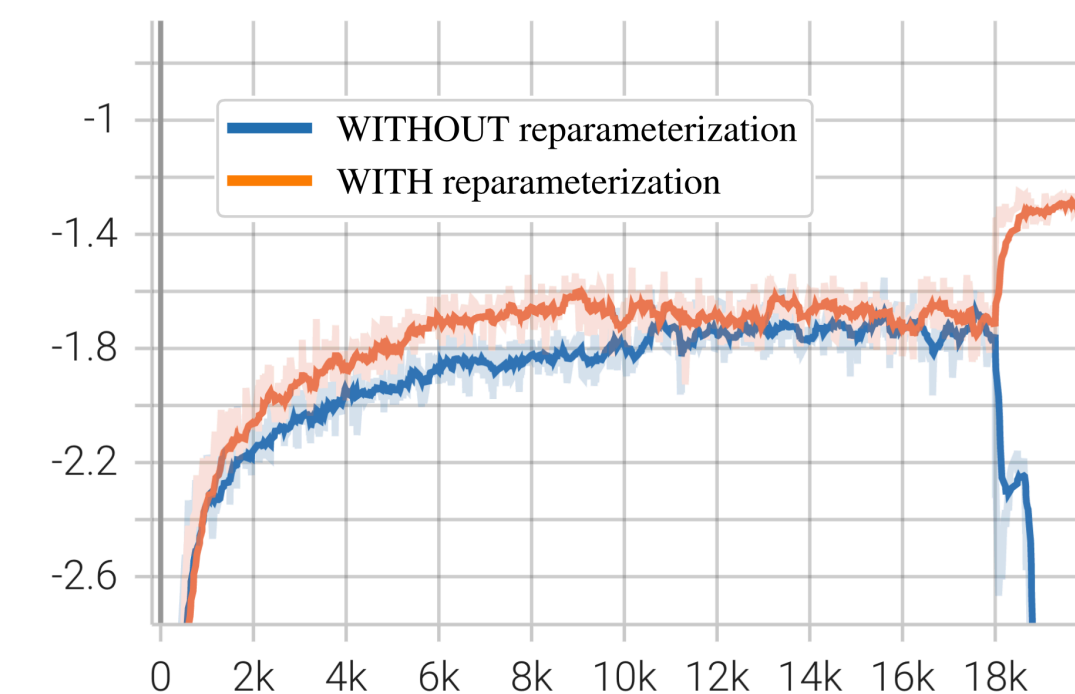
NNs with parameters  $\theta_{\mu_l} \rightarrow f_{\mu_l}, \theta_{\sigma_l} \rightarrow f_{\sigma_l}$

- ✓ Perform inference on  $\epsilon = \{\epsilon_1, \dots, \epsilon_1\}$  with standard Gaussian prior.

- ✓ No need to increase complexity of the HMC method.



(a) Mean acceptance rate  $\bar{p}_a$



(b)  $\log p(\mathbf{x}_U | \mathbf{x}_O)$





# Method

## Optimization algorithm

---

**Algorithm 1** Training algorithm for HH-VAEM

---

**Input:** data  $(\mathbf{x}_O^{(1:N)}, \mathbf{y}_O^{(1:N)})$ , steps:  $T_d, T_{VI}, T_{HMC}$

**Parameters:**  $\gamma, \theta, \psi, \phi, s$

STAGE 1: MARGINAL VAES

**for**  $d = 1$  **to**  $D$  **do**

  Initialize marginal VAE  $\{\theta_d, \gamma_d\}$

**for**  $t = 1$  **to**  $T_d$  **do**

$\gamma_d^{t+1}, \theta_d^{t+1} \leftarrow \text{Adam}_{\gamma_d^t, \theta_d^t}(\mathcal{L}_d)$

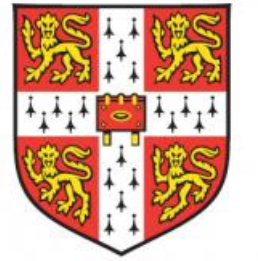
**end for**

**end for**

---

- 1. Train marginal VAEs using:

$$\mathcal{L}_d(x_d; \{\theta_d, \gamma_d\}) = \mathbb{I}(x_d \in \mathbf{x}_O) \mathbb{E}_{q_{\gamma_d}(z_d|x_d)} \log \frac{p_{\theta_d}(x_d, z_d)}{q_{\gamma_d}(z_d | x_d)}$$



# Method

## Optimization algorithm

---

**Algorithm 1** Training algorithm for HH-VAEM

---

**Input:** data  $(\mathbf{x}_O^{(1:N)}, \mathbf{y}_O^{(1:N)})$ , steps:  $T_d, T_{VI}, T_{HMC}$

**Parameters:**  $\gamma, \theta, \psi, \phi, s$

STAGE 1: MARGINAL VAES

**for**  $d = 1$  **to**  $D$  **do**

  Initialize marginal VAE  $\{\theta_d, \gamma_d\}$

**for**  $t = 1$  **to**  $T_d$  **do**

$\gamma_d^{t+1}, \theta_d^{t+1} \leftarrow \text{Adam}_{\gamma_d^t, \theta_d^t}(\mathcal{L}_d)$

**end for**

**end for**

STAGE 2: DEPENDENCY VAE

**for**  $t = 1$  **to**  $T_{VAE}$  **do**

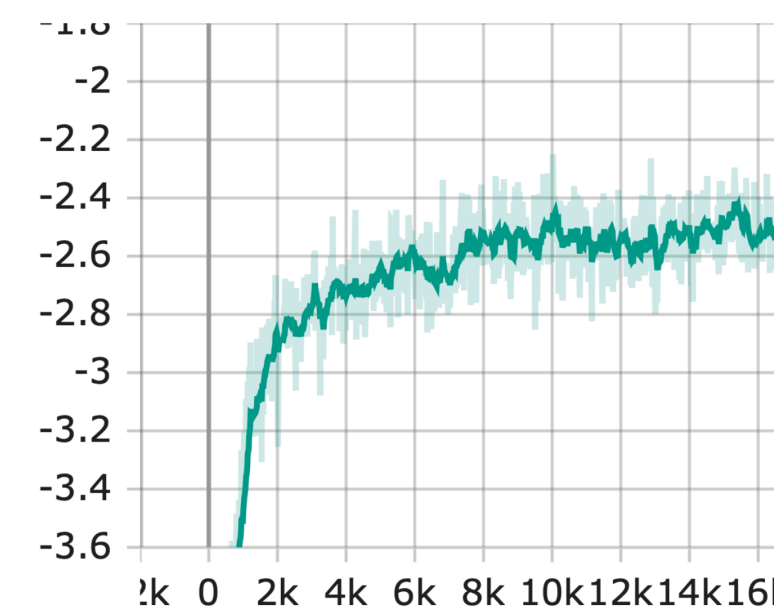
$\theta^{t+1}, \psi^{t+1} \leftarrow \text{Adam}_{\theta^t, \psi^t}(\mathcal{L}_{VI})$

**end for**

---

- 2. Training Hierarchical VAE using the ELBO:

$$\mathcal{L}_{VI}(\mathbf{x}_O, \mathbf{y}_O; \{\theta, \psi\}) = \mathbb{E}_{q_\psi} [\log p_\theta(\mathbf{z}_O | \mathbf{h}_1) + \log p_\theta(\mathbf{y}_O | \hat{\mathbf{x}}, \mathbf{h}_1)] - \sum_{l=1}^L D_{\text{KL}}(q_\psi(\epsilon_l | \mathbf{x}_O, \mathbf{y}_O) \| p(\epsilon_l))$$



(a)  $\log p(\mathbf{x}_U | \mathbf{x}_O)$



# Method

## Optimization algorithm

---

### Algorithm 1 Training algorithm for HH-VAEM

---

**Input:** data  $(\mathbf{x}_O^{(1:N)}, \mathbf{y}_O^{(1:N)})$ , steps:  $T_d, T_{VI}, T_{HMC}$

**Parameters:**  $\gamma, \theta, \psi, \phi, s$

STAGE 1: MARGINAL VAES

**for**  $d = 1$  **to**  $D$  **do**

  Initialize marginal VAE  $\{\theta_d, \gamma_d\}$

**for**  $t = 1$  **to**  $T_d$  **do**

$\gamma_d^{t+1}, \theta_d^{t+1} \leftarrow \text{Adam}_{\gamma_d^t, \theta_d^t}(\mathcal{L}_d)$

**end for**

**end for**

STAGE 2: DEPENDENCY VAE

**for**  $t = 1$  **to**  $T_{VAE}$  **do**

$\theta^{t+1}, \psi^{t+1} \leftarrow \text{Adam}_{\theta^t, \psi^t}(\mathcal{L}_{VI})$

**end for**

STAGE 3: JOINTLY OPTIMIZING VAE + HMC

**for**  $t = 1$  **to**  $T_{HMC}$  **do**

$\psi^{t+1} \leftarrow \text{Adam}_{\psi^t}(\mathcal{L}_{VI})$

$\theta^{t+1}, \phi^{t+1} \leftarrow \text{Adam}_{\theta^t, \phi^t}(\mathcal{L}_{HMC})$

$s^{t+1} \leftarrow \text{Adam}_{s^t}(\mathcal{L}_{SKSD})$

**end for**

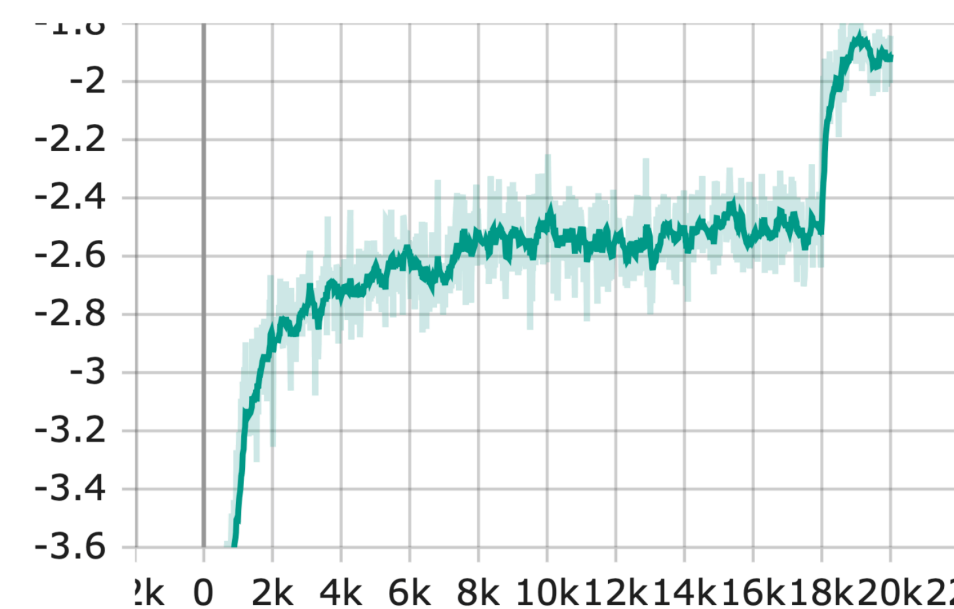
---

- Train a) encoder using ELBO, b) HMC hyperparams, decoder and predictor parameters using HMC objective and c) scale using SKSD.

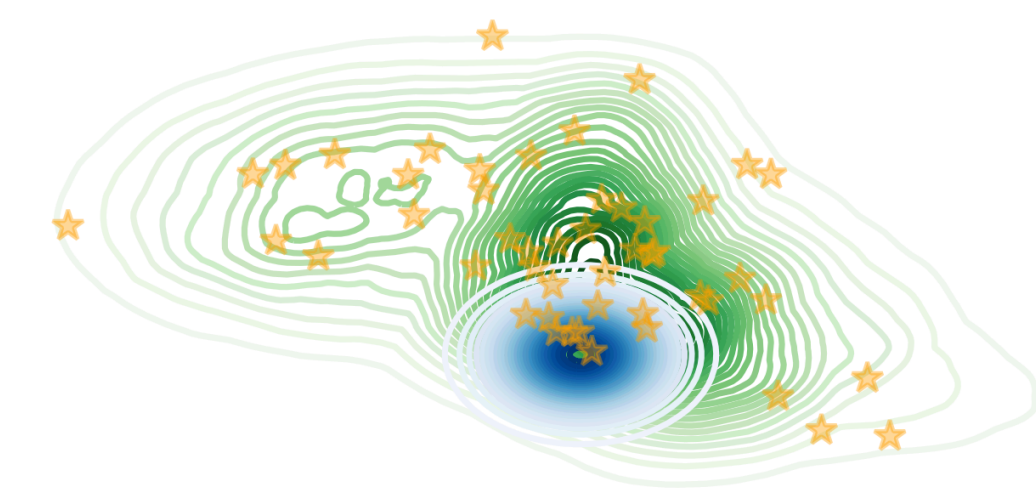
$$\mathcal{L}_{HMC}(\mathbf{z}_O, \mathbf{y}_O; \{\theta, \psi, \phi\}) = \mathbb{E}_{q_\phi^{(T)}(\epsilon)} [\log p_\theta(\mathbf{z}_O | \mathbf{h}_1) + \log p_\theta(\mathbf{y}_O | \hat{\mathbf{x}}, \mathbf{h}_1) + \sum_{l=1}^L p(\epsilon_l^{(T)})]$$

$$\mathcal{L}_{SKSD}(\mathbf{x}_O, \mathbf{y}_O; s) = \text{SKSD} \left( q_\phi^{(T)}(\epsilon | \mathbf{z}_O, \mathbf{x}_O, \mathbf{y}_O; s), p(\epsilon | \mathbf{z}_O, \mathbf{x}_O, \mathbf{y}_O) \right)$$

$q(z|x)$      $p(z|x)$



(a)  $\log p(\mathbf{x}_U | \mathbf{x}_O)$



$\epsilon^2$ HMC samples (orange)



# Experiments

## Set up

- **HH-VAEM** with 2 layers of latent variables.
- Baseline models:
  1. **VAEM**: Gaussian-based, 1 layer.
  2. **MIWAEM**: Gaussian-based, 1 layer, importance weighted.
  3. **HMC-VAEM**: HMC-based, 1 layer.
  4. **H-VAEM**: Gaussian-based, 2 layers.
- **Training**: missing features and target with a probability sampled from  $U(0.01, 0.99)$  each batch.
- **Test**: 50% missing features, fully-missing target.





# Experiments

## Missing data imputation and target prediction

$$\log p(\mathbf{x}_U | \mathbf{x}_O) = \log \mathbb{E}_{\epsilon \sim q^{(T)}(\epsilon | \mathbf{x}_O)} [p(\mathbf{x}_U | \epsilon)] \approx \log \frac{1}{k} \sum_i^k p(\mathbf{x}_U | \epsilon_i)$$

	Bank	Insurance	Avocado	Naval	Yatch	Diabetes	Concrete	Wine	Energy	Boston
VAEM	2.84 ± 0.07	1.81 ± 0.03	1.89 ± 0.01	0.55 ± 0.05	3.15 ± 0.28	2.78 ± 0.16	2.45 ± 0.26	3.01 ± 0.61	2.09 ± 0.10	2.01 ± 0.23
MIWAEM	2.74 ± 0.05	1.88 ± 0.04	1.92 ± 0.04	0.57 ± 0.03	2.66 ± 0.11	2.55 ± 0.09	2.34 ± 0.51	2.76 ± 0.48	2.06 ± 0.14	1.94 ± 0.23
H-VAEM	2.82 ± 0.06	1.80 ± 0.04	1.89 ± 0.01	0.48 ± 0.06	3.06 ± 0.31	2.74 ± 0.09	2.42 ± 0.21	2.85 ± 0.56	1.72 ± 0.11	1.89 ± 0.24
HMC-VAEM	2.69 ± 0.05	1.77 ± 0.06	1.89 ± 0.02	0.49 ± 0.07	<b>2.21 ± 0.24</b>	2.72 ± 0.20	2.28 ± 0.29	2.83 ± 0.46	1.73 ± 0.05	1.83 ± 0.16
<b>HH-VAEM</b>	<b>2.63 ± 0.04</b>	<b>1.75 ± 0.03</b>	<b>1.88 ± 0.05</b>	<b>0.40 ± 0.05</b>	2.47 ± 0.27	<b>2.54 ± 0.13</b>	<b>2.28 ± 0.09</b>	<b>1.90 ± 0.17</b>	<b>1.71 ± 0.04</b>	<b>1.83 ± 0.11</b>

Table 1: Test negative log likelihood of the unobserved features for our model and baselines.

$$\log p(\mathbf{y} | \mathbf{x}_O) = \log \mathbb{E}_{\epsilon \sim q^{(T)}(\epsilon | \mathbf{x}_O)} [p(\mathbf{y} | \epsilon)] \approx \log \frac{1}{k} \sum_i^k p(\mathbf{y} | \epsilon_i),$$

	Bank	Insurance	Avocado	Naval	Yatch	Diabetes	Concrete	Wine	Energy	Boston
VAEM	0.56 ± 0.06	1.20 ± 0.03	1.18 ± 0.02	2.69 ± 0.01	0.61 ± 0.02	1.59 ± 0.19	1.07 ± 0.09	0.28 ± 0.09	0.61 ± 0.14	0.85 ± 0.21
MIWAEM	0.51 ± 0.03	1.15 ± 0.03	1.15 ± 0.03	2.70 ± 0.01	0.60 ± 0.03	<b>1.36 ± 0.10</b>	0.95 ± 0.22	0.28 ± 0.13	0.54 ± 0.12	0.80 ± 0.21
H-VAEM	0.50 ± 0.03	1.06 ± 0.02	1.18 ± 0.02	2.68 ± 0.01	0.60 ± 0.02	1.71 ± 0.14	1.02 ± 0.09	0.26 ± 0.11	0.46 ± 0.14	0.90 ± 0.22
HMC-VAEM	0.52 ± 0.02	1.00 ± 0.03	1.12 ± 0.03	2.71 ± 0.01	<b>0.52 ± 0.15</b>	1.55 ± 0.29	0.95 ± 0.26	0.28 ± 0.09	0.41 ± 0.07	0.71 ± 0.13
<b>HH-VAEM</b>	<b>0.49 ± 0.03</b>	<b>0.93 ± 0.06</b>	<b>1.10 ± 0.01</b>	<b>2.62 ± 0.01</b>	0.56 ± 0.02	1.38 ± 0.18	<b>0.95 ± 0.08</b>	<b>0.20 ± 0.04</b>	<b>0.32 ± 0.05</b>	<b>0.55 ± 0.04</b>

Table 2: Test negative log likelihood of the predicted target for our model and baselines.





# Experiments (MNIST datasets)

	VAE	MIWAE	H-VAE	HMC-VAE	<b>HH-VAE</b>
MNIST	0.124 ± 0.001	0.121 ± 0.001	0.119 ± 0.001	0.101 ± 0.004	<b>0.094 ± 0.003</b>
F-MNIST	0.162 ± 0.002	0.160 ± 0.002	0.156 ± 0.002	0.150 ± 0.002	<b>0.144 ± 0.002</b>

Table 3: Test negative log likelihood of the unobserved features for the MNIST datasets.



	VAE	MIWAE	H-VAE	HMC-VAE	<b>HH-VAE</b>
MNIST	0.153 ± 0.009	0.151 ± 0.007	0.146 ± 0.006	0.067 ± 0.007	<b>0.056 ± 0.019</b>
F-MNIST	0.501 ± 0.012	0.496 ± 0.008	0.494 ± 0.007	0.357 ± 0.060	<b>0.337 ± 0.069</b>

Table 4: Test negative log likelihood of the predicted target for the MNIST datasets.

	VAE	MIWAE	H-VAE	HMC-VAE	<b>HH-VAE</b>
MNIST	0.953 ± 0.004	0.953 ± 0.003	0.953 ± 0.003	0.978 ± 0.003	<b>0.981 ± 0.005</b>
F-MNIST	0.824 ± 0.005	0.824 ± 0.004	0.824 ± 0.004	0.869 ± 0.015	<b>0.876 ± 0.017</b>

Table 5: Test accuracy of the predicted digits for the MNIST datasets.



# Experiments

## Sequential Active Information Acquisition (SAIA)

- Sequentially acquiring high-value information by selecting features that maximize **our proposed sampling-based reward**:

$$\hat{I}(\mathbf{y}; x_i | \mathbf{x}_O) \approx \sum_{ij} p_{x_i, \mathbf{y} | \mathbf{x}_O}(i, j) \log \frac{p_{x_i, \mathbf{y} | \mathbf{x}_O}(i, j)}{p_{x_i | \mathbf{x}_O}(i) p_{\mathbf{y} | \mathbf{x}_O}(j)}$$

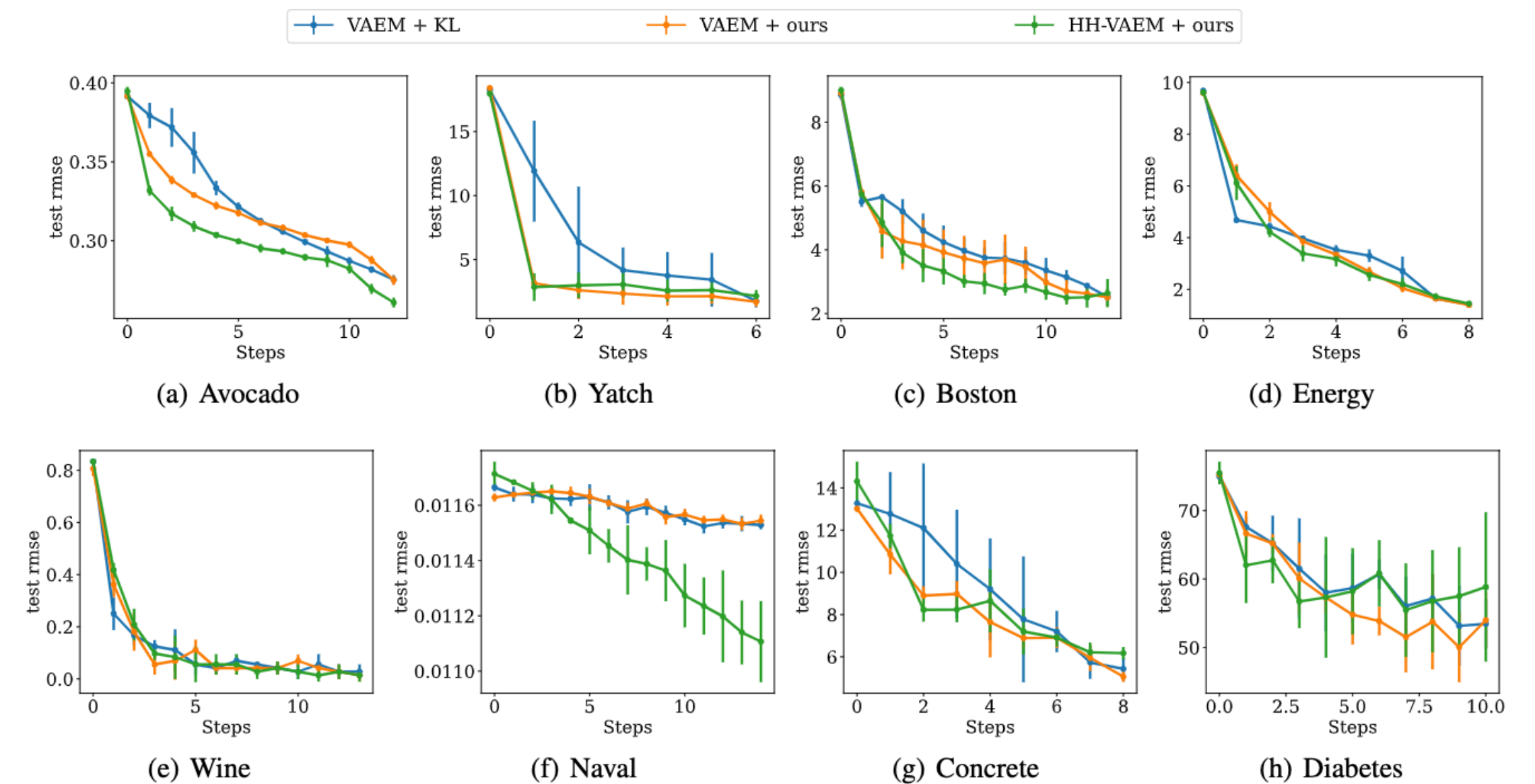
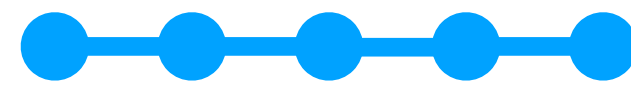


Figure 5: SAIA metric curves. Horizontal axis shows acquisition steps (number of discovered features). Vertical axis is the RMSE.





# Experiments

## Conditional image inpainting

1. Encode to  $q_{\phi}^{(0)}(\epsilon | z_O, x_O, y_O)$
2. Using HMC, sample from  $q_{\phi}^{(T)}(\epsilon | z_O, x_O, y_O)$
3. Decode to  $p(x_U | \epsilon^{(T)})$



# Conclusion

- We presented:
  1. **HH-VAEM**: novel Hierarchical VAE improved with HMC with automatic hyperparameter optimization.
  2. Novel **sampling-based technique** based on the Mutual Information estimation for efficient information acquisition.
- Based on the provided experiments, we demonstrate that our methods:
  - ✓ Improve approximate inference in hierarchical VAEs wrt to the Gaussian approximation.
  - ✓ Improve missing data imputation task.
  - ✓ Improve prediction task.
  - ✓ Improve active information acquisition task.



# Further details



---

## Missing Data Imputation and Acquisition with Deep Hierarchical Models and Hamiltonian Monte Carlo

---

**Ignacio Peis**  
Universidad Carlos III de Madrid  
Madrid, Spain  
ipeis@tsc.uc3m.es

**Chao Ma**  
Microsoft Research  
University of Cambridge  
Cambridge, UK  
cm905@cam.ac.uk

**José Miguel Hernández-Lobato**  
University of Cambridge  
Cambridge, UK  
jmh233@cam.ac.uk

[\[Code\]](#)

[\[Automatic HMC code\]](#)



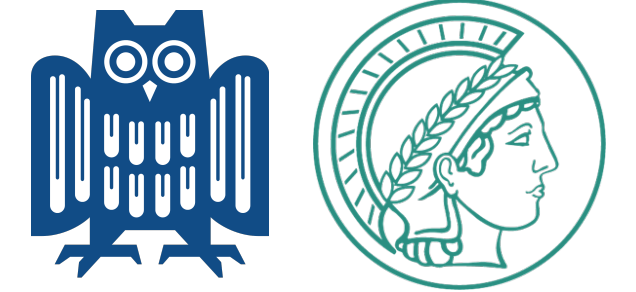
[\[Paper\]](#)





## Part II

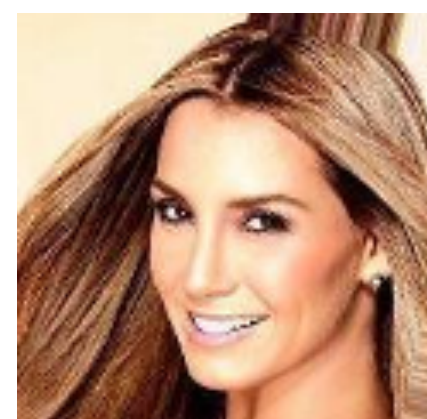
# Variational Mixture of HyperGenerators for Learning Distributions over Functions



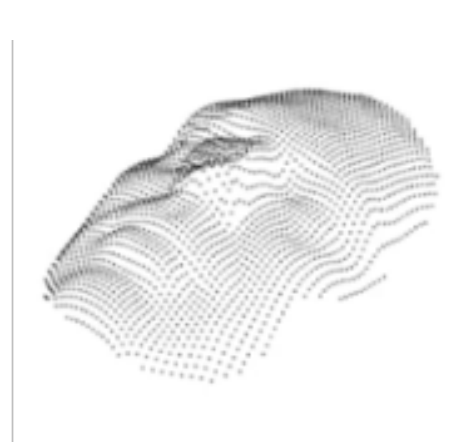
# Motivation

- We typically deal with discretized versions of data that are continuous in nature.

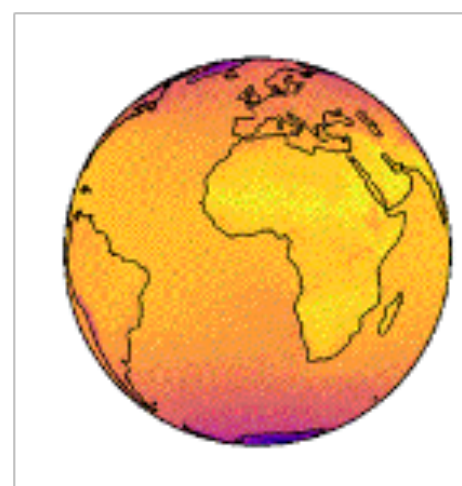
2D Images



3D Images

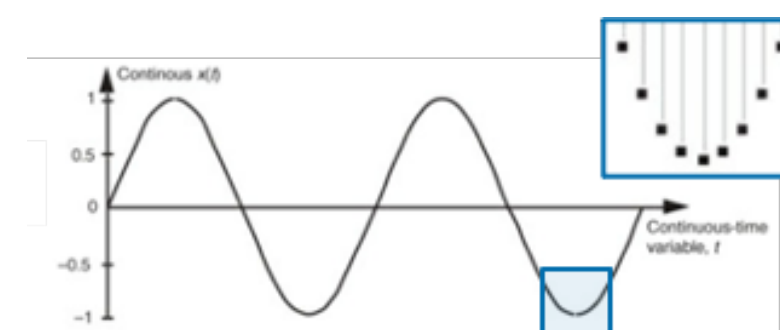


Polar data



Spatial

Time series

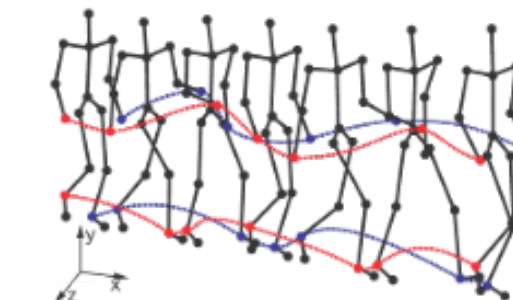


Temporal

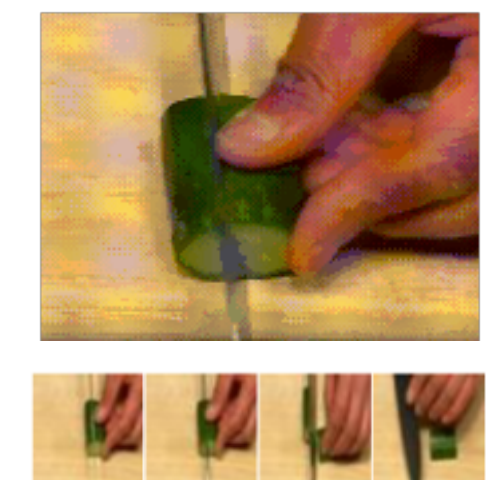
Audio



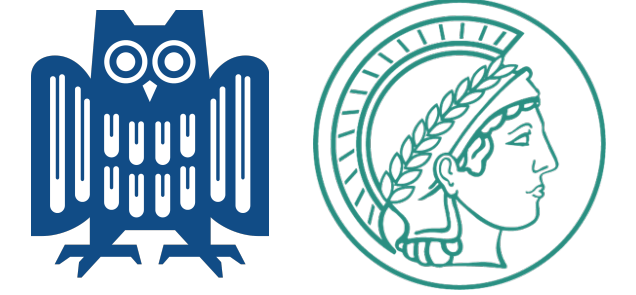
Motion sequences



Video

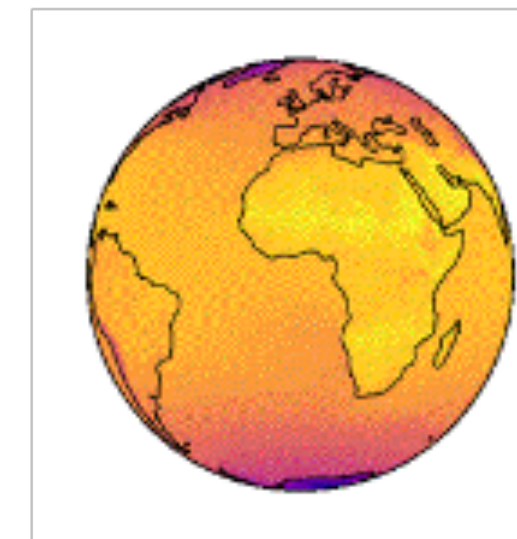
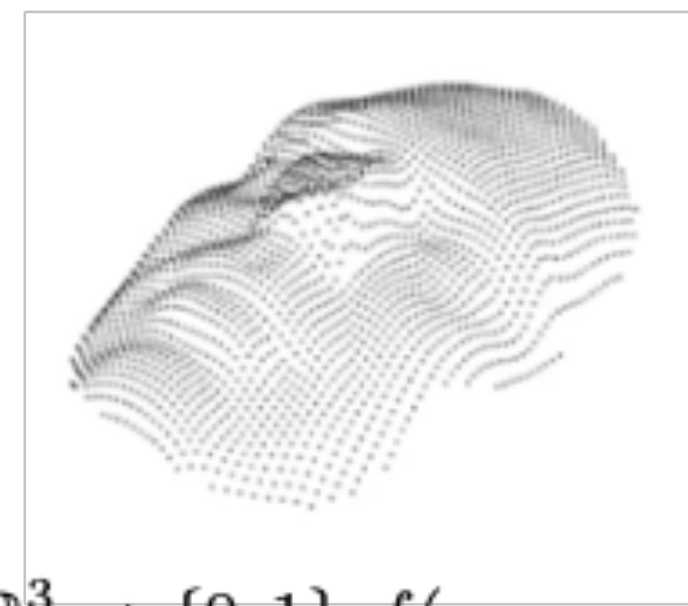


Spatio-temporal



# Motivation

- Data can be expressed as a function over continuous coordinate systems.



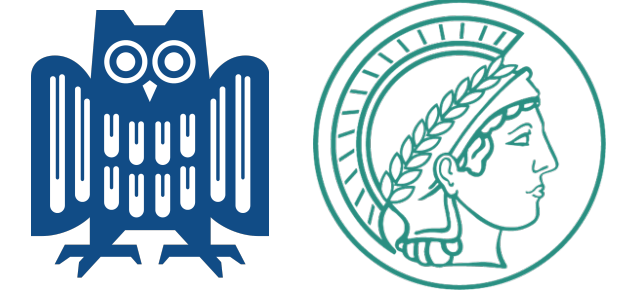
$$f : \mathbb{R}^2 \rightarrow \mathbb{R}^3, f(x_1, x_2) = (r, g, b)$$

$$f : \mathbb{R}^3 \rightarrow \{0, 1\}, f(x_1, x_2, x_3) = p$$

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}, f(\varphi, \lambda) = T$$

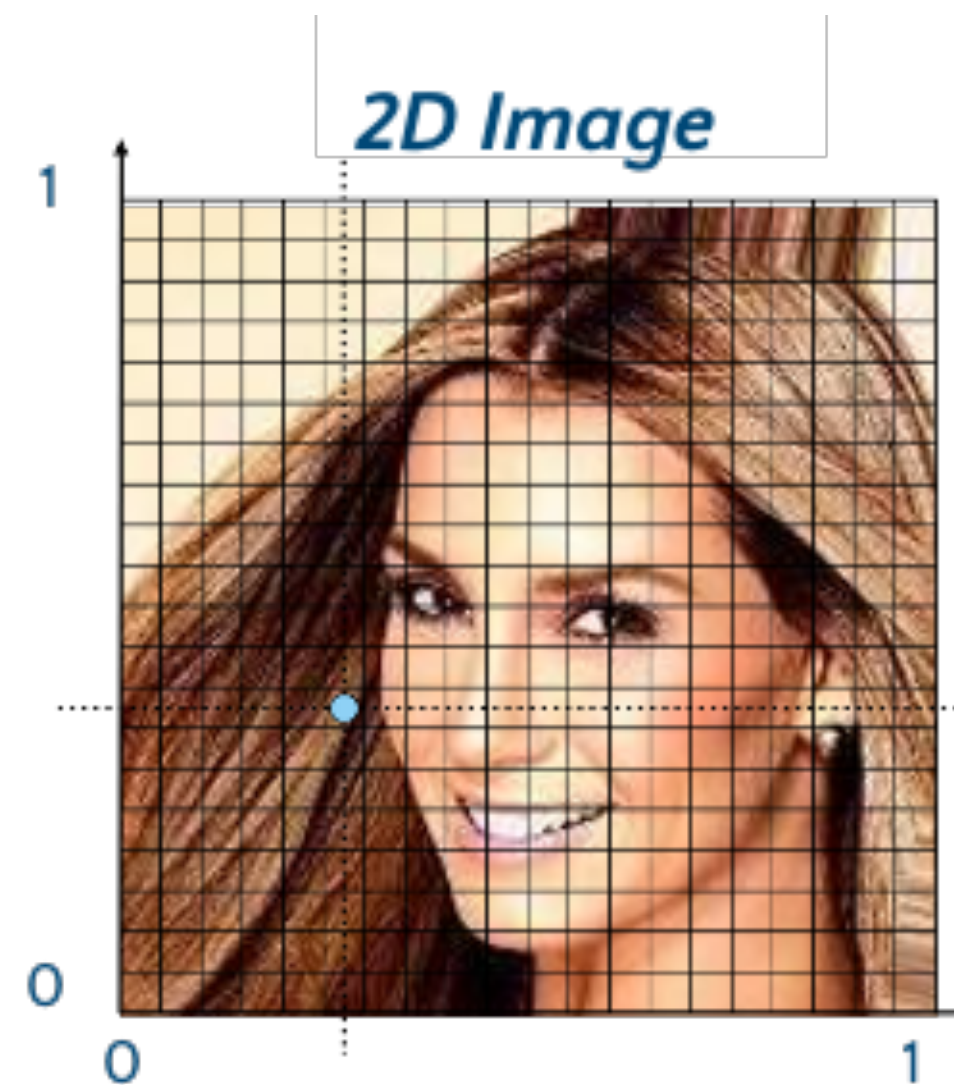
$$f : \mathbb{R}^3 \rightarrow \mathbb{R}^3, f(x_1, x_2, t) = (r, g, b)$$





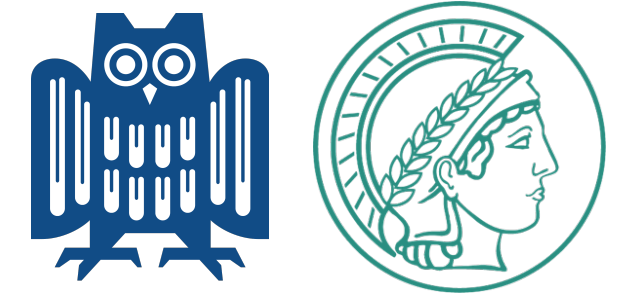
# Motivation

- Focusing on images:



- Each pixel is now a pair  $\{\mathbf{x}_d, \mathbf{y}_d\}$  where  $x_d \in \mathbb{R}^2$ ,  $y_d \in \mathbb{R}^3$
- Full image is a pair of sets  $\mathbf{X} = \{\mathbf{x}_d\}_{d=1}^D$ ,  $\mathbf{Y}_d = \{\mathbf{y}_d\}_{d=1}^D$
- Generator function  $f : \mathbf{X} \rightarrow \mathbf{Y}$  creates this specific image with the mapping  $f(\mathbf{x}_d) = \mathbf{y}_d$ ,  $d \in [1, \dots, D]$

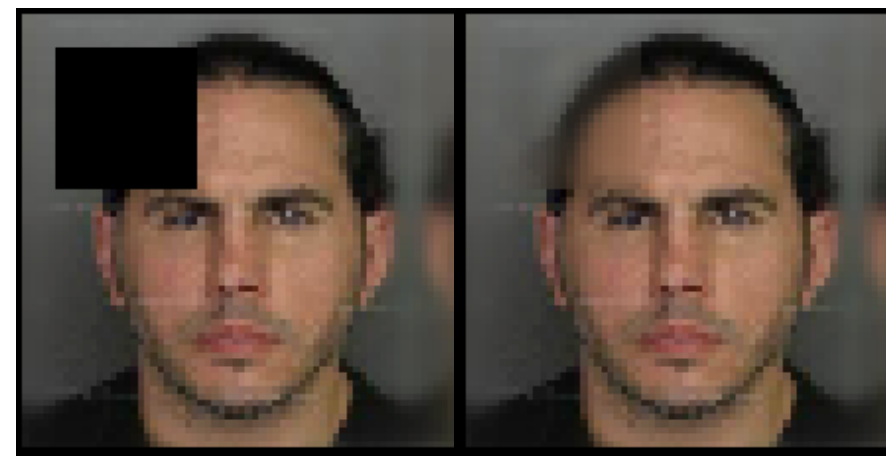




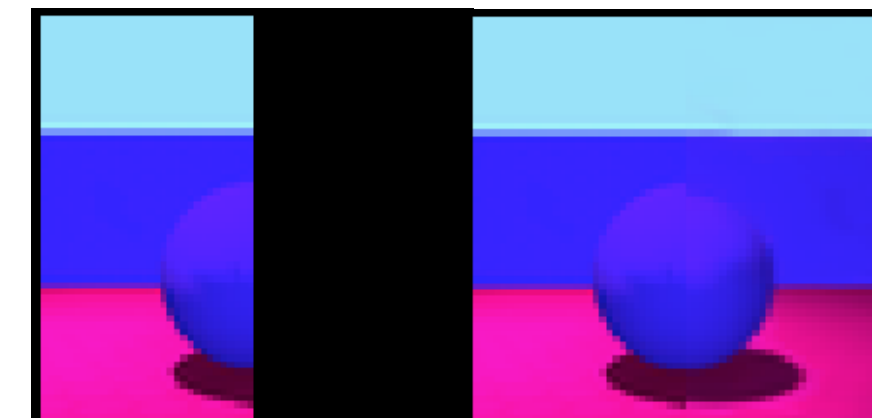
# Motivation

- This approach will easily allow for:

Inpainting [6,7]



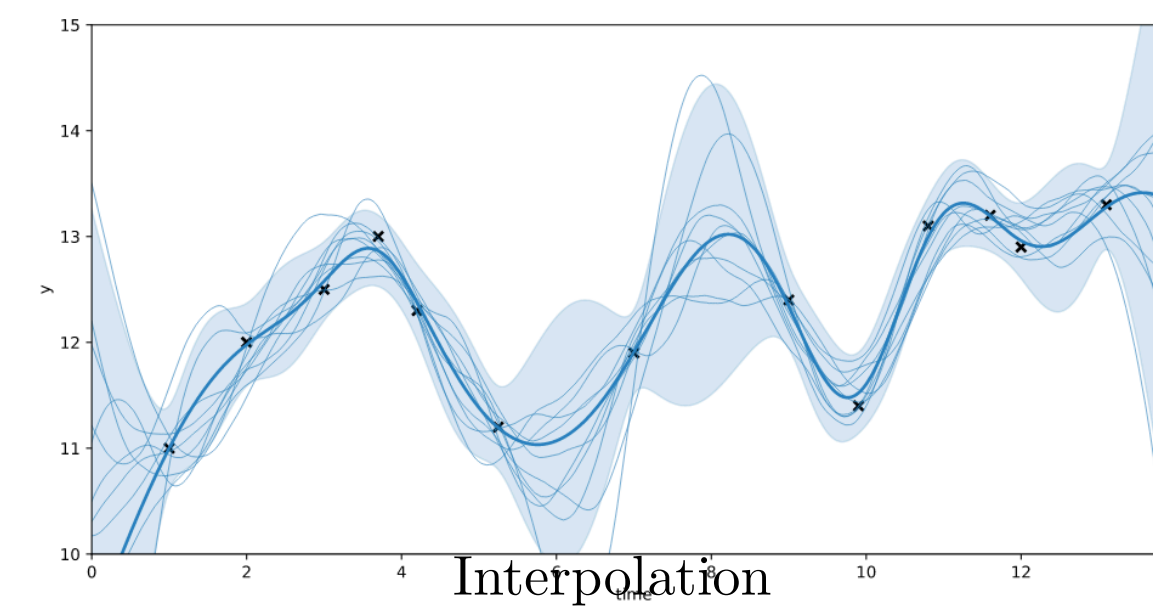
Outpainting [6,7]

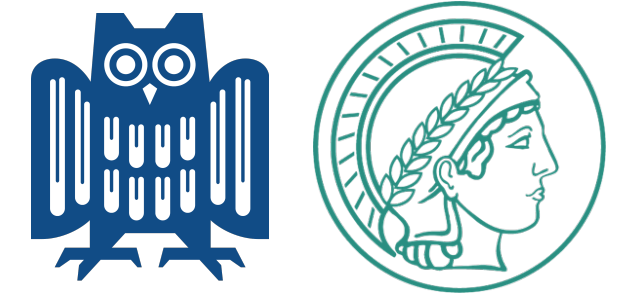
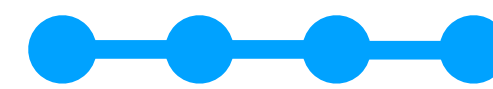


Superresolution [6,7]




Conditional Generation

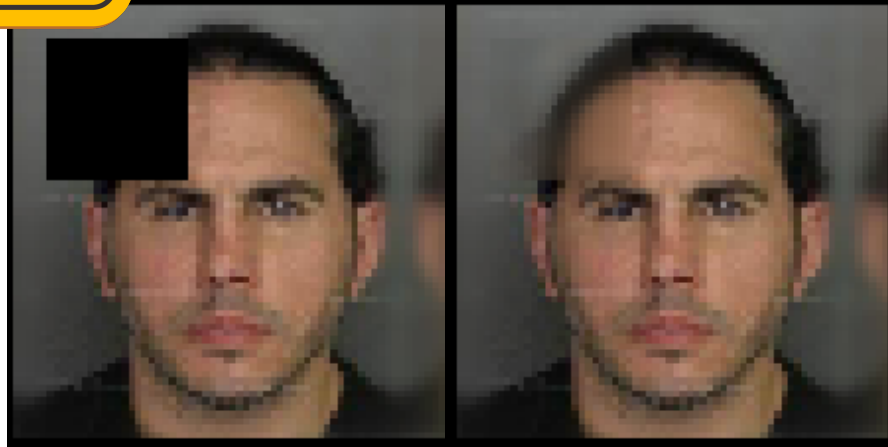





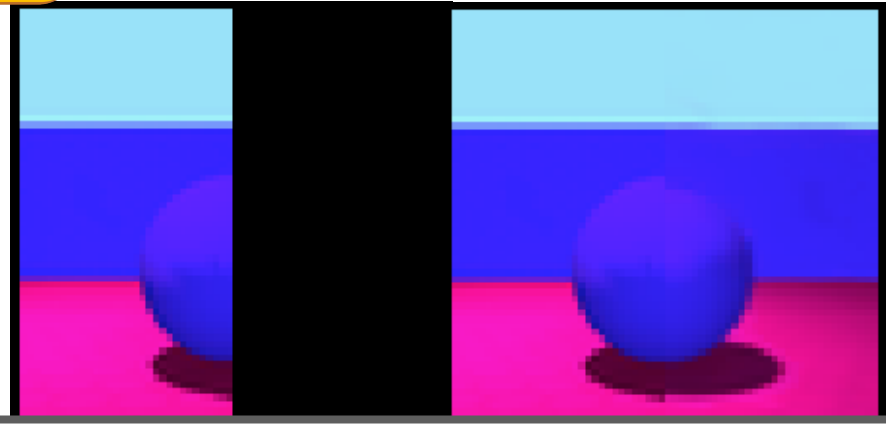
# Motivation


- This approach will easily allow for:


 Inpainting [6,7]




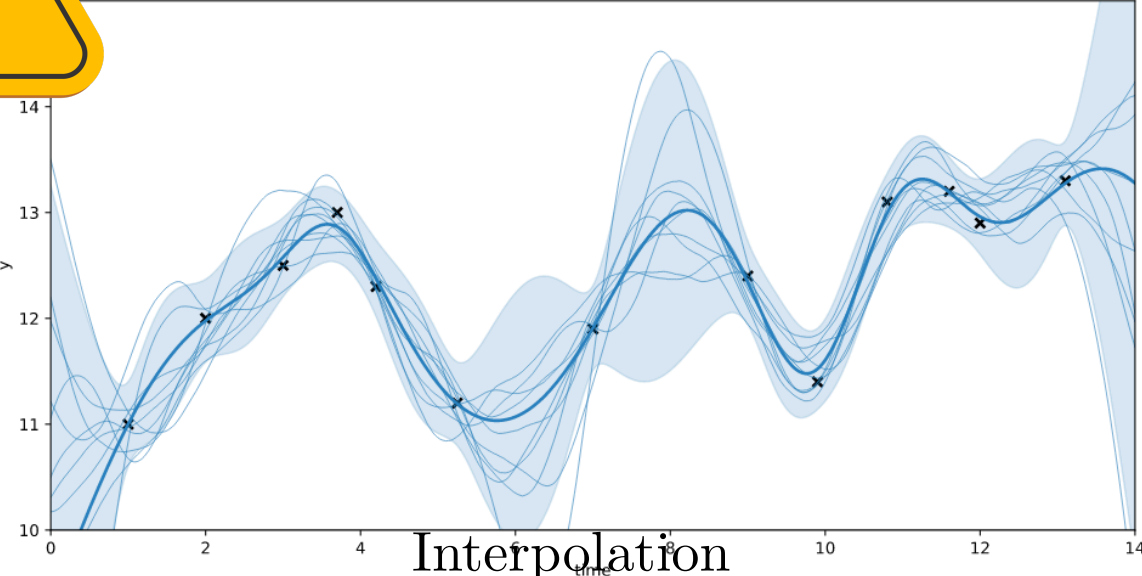
 Outpainting [6,7]



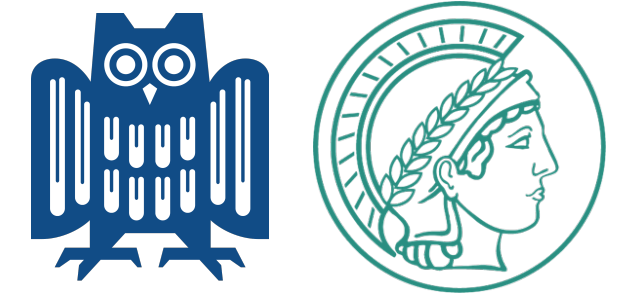
 Superresolution [6,7]



 Conditional Generation

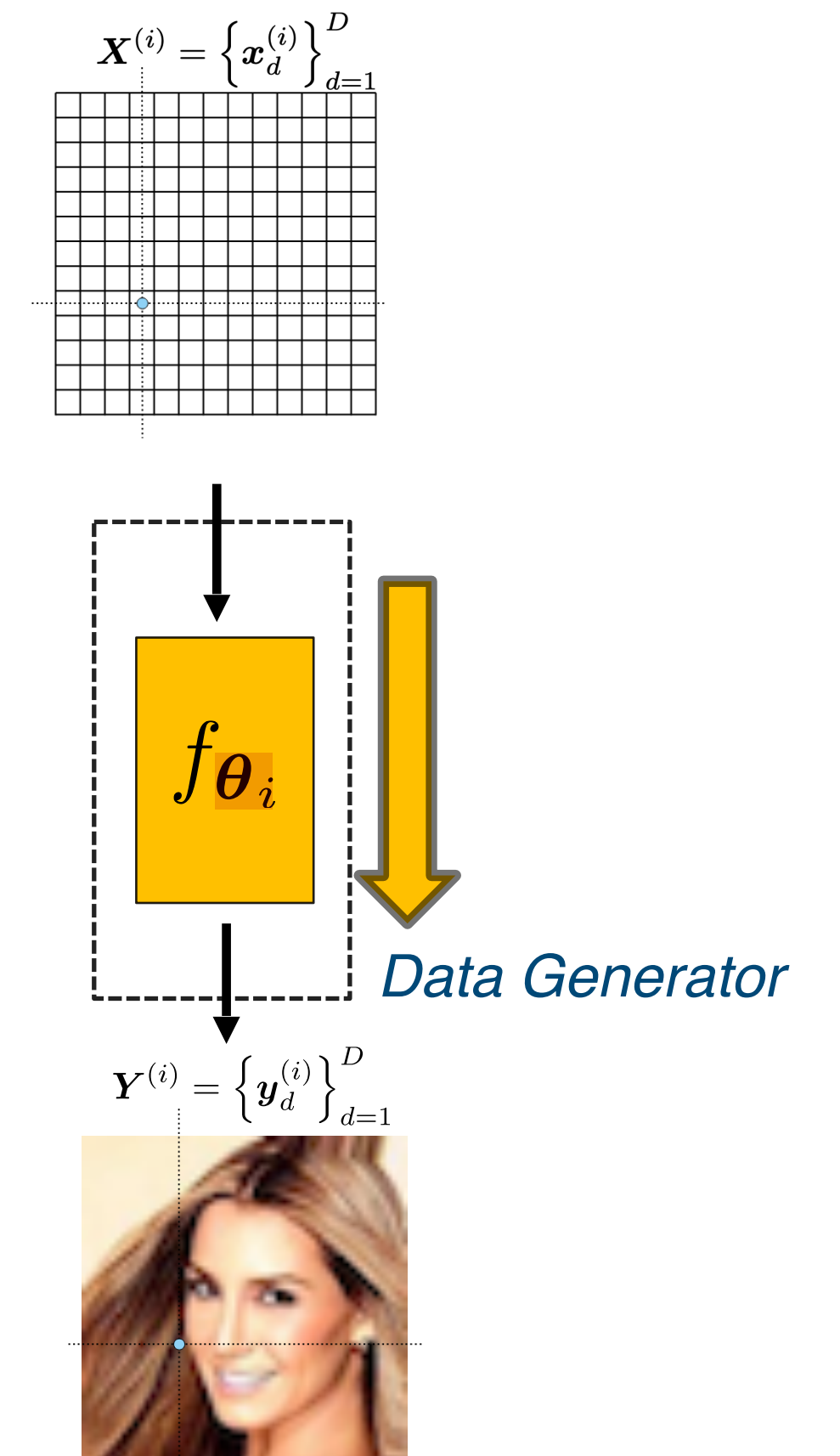


These ones require inference!



# Implicit Neural Representations

## INRs [20-22]



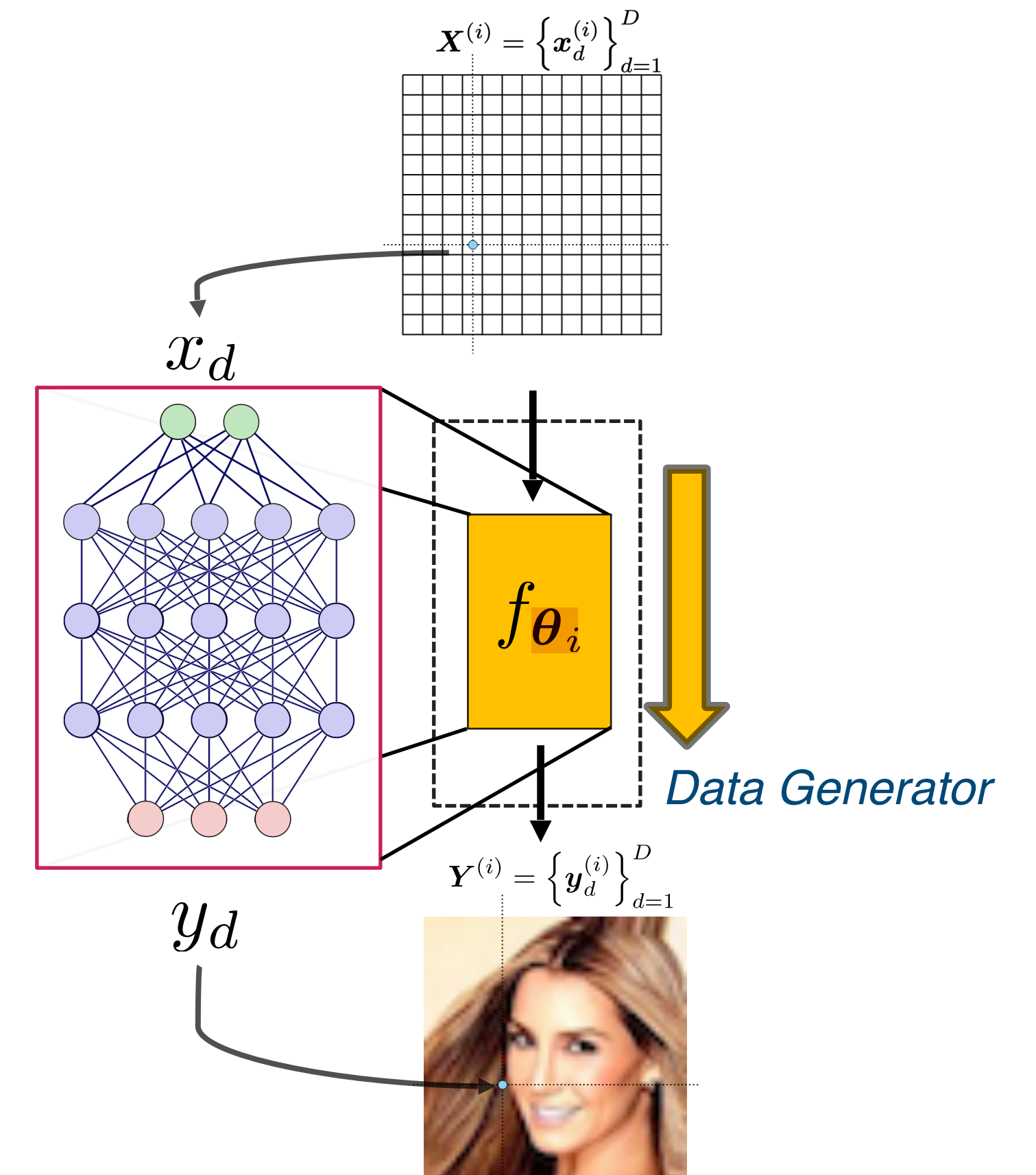
[20] Sitzmann et al., 2020

[21] Mescheder et al., 2019

[20] Sitzmann et al., 2019

# Implicit Neural Representations

## INRs [20-22]



[20] Sitzmann et al., 2020

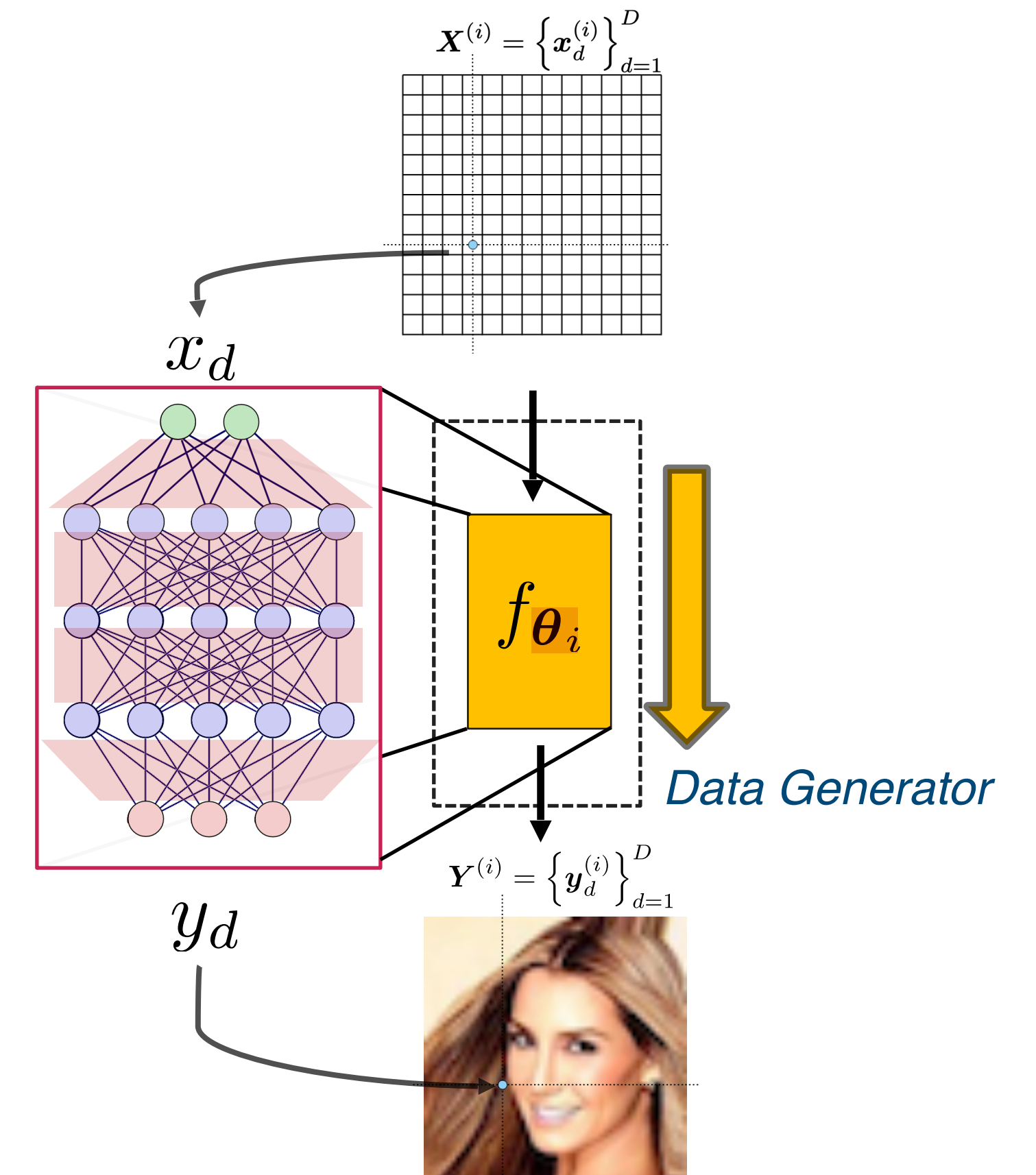
[21] Mescheder et al., 2019

[20] Sitzmann et al., 2019



# Implicit Neural Representations

## INRs [20-22]



[20] Sitzmann et al., 2020

[21] Mescheder et al., 2019

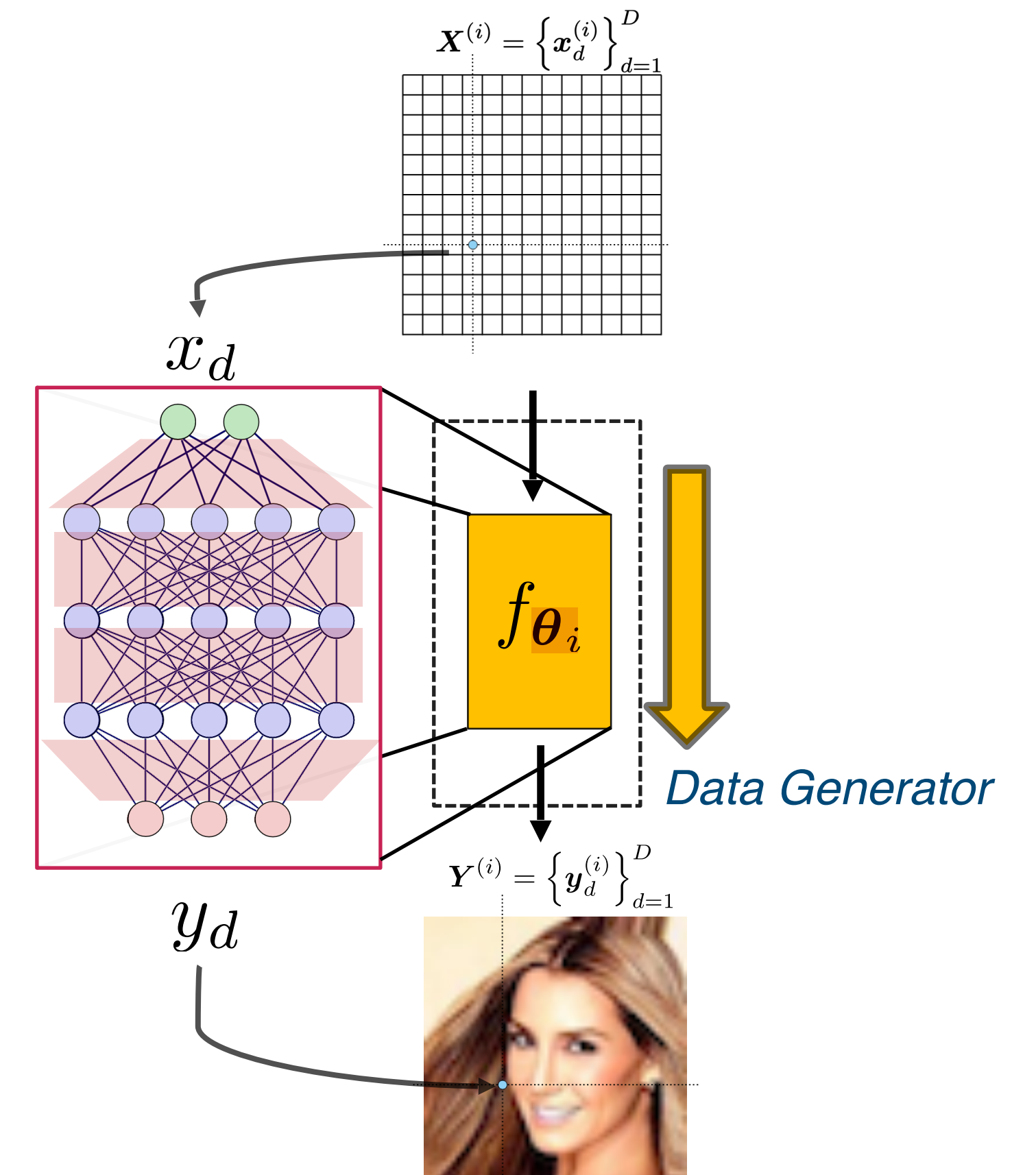
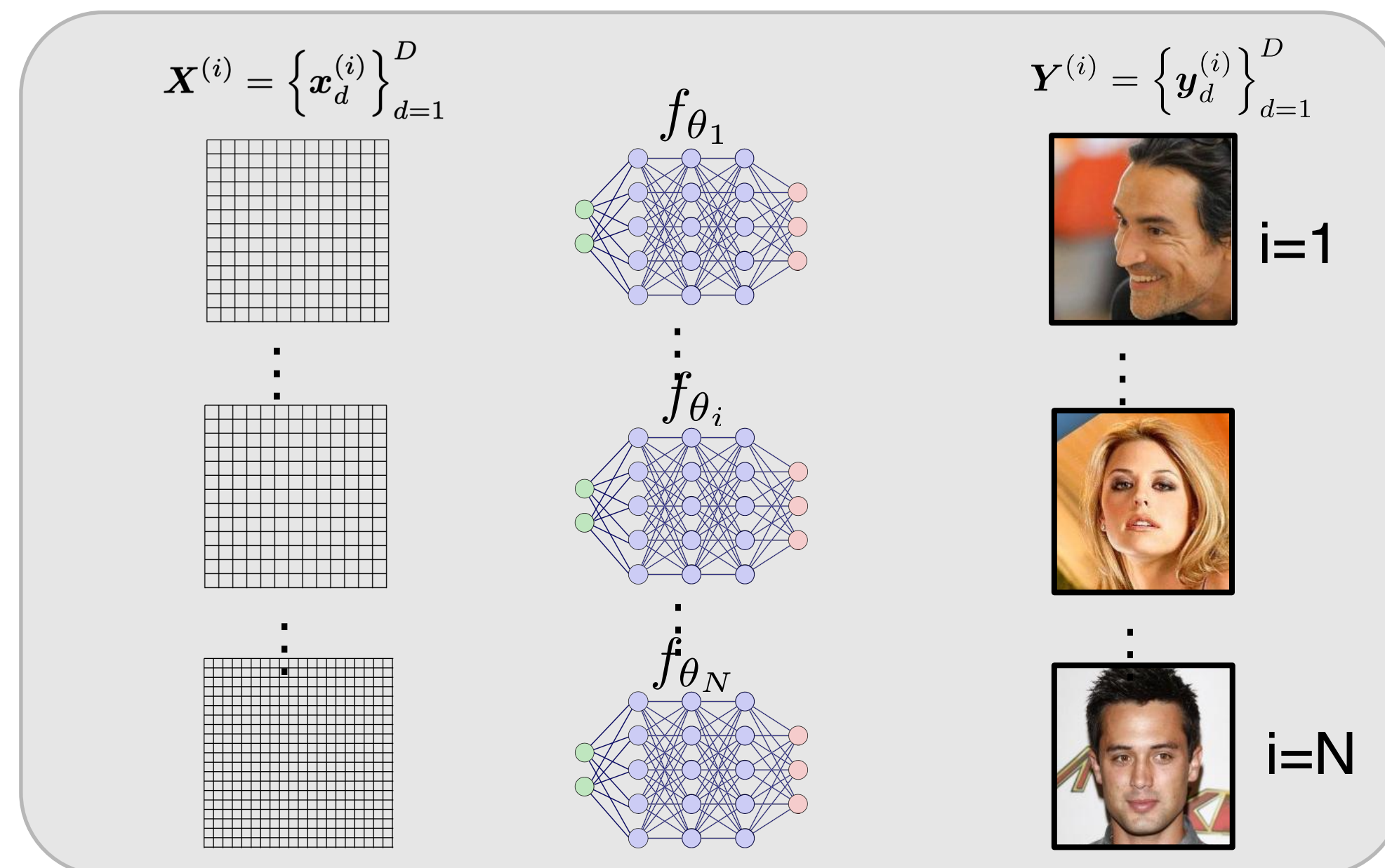
[20] Sitzmann et al., 2019

# Implicit Neural Representations

## INRs [20-22]

- Learning distributions functions within a VAE framework.

Data generator  $f_{\theta_i}$  is unique to each image



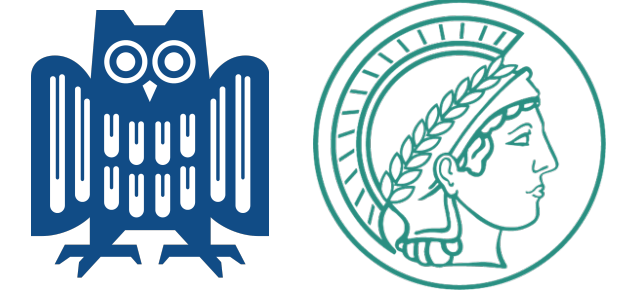
[20] Sitzmann et al., 2020

[21] Mescheder et al., 2019

[20] Sitzmann et al., 2019

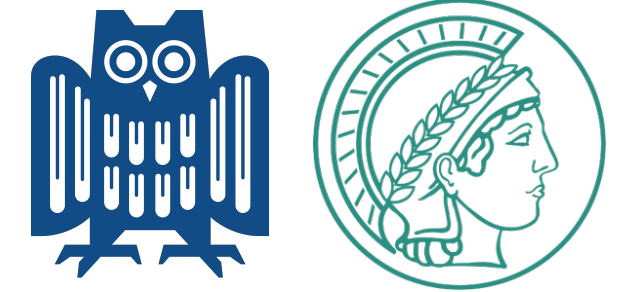


# Deep Generative Models of INRs

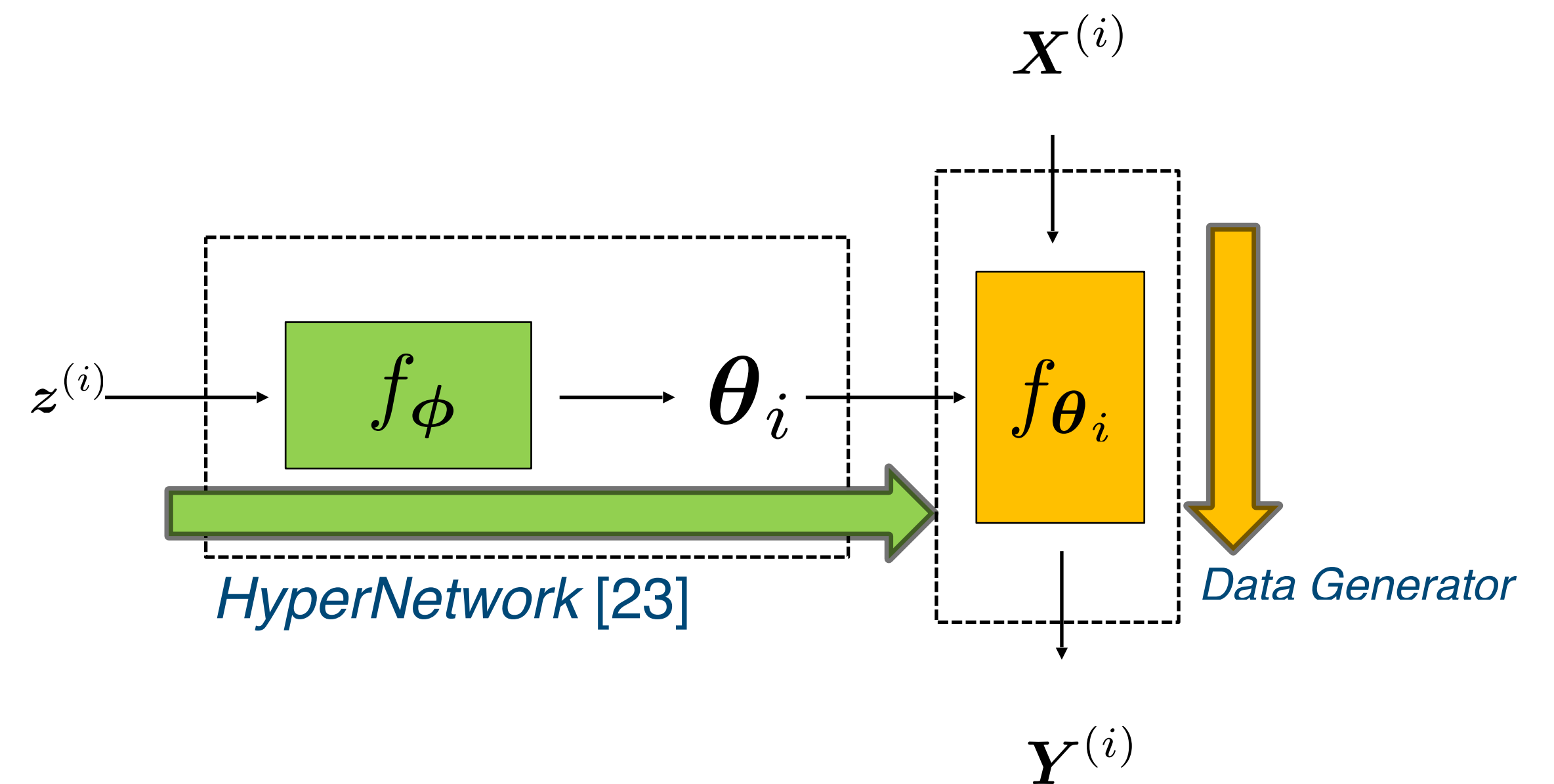


How to scale to large datasets?

How to map a latent representation to an INR?



# Deep Generative Models of INRs



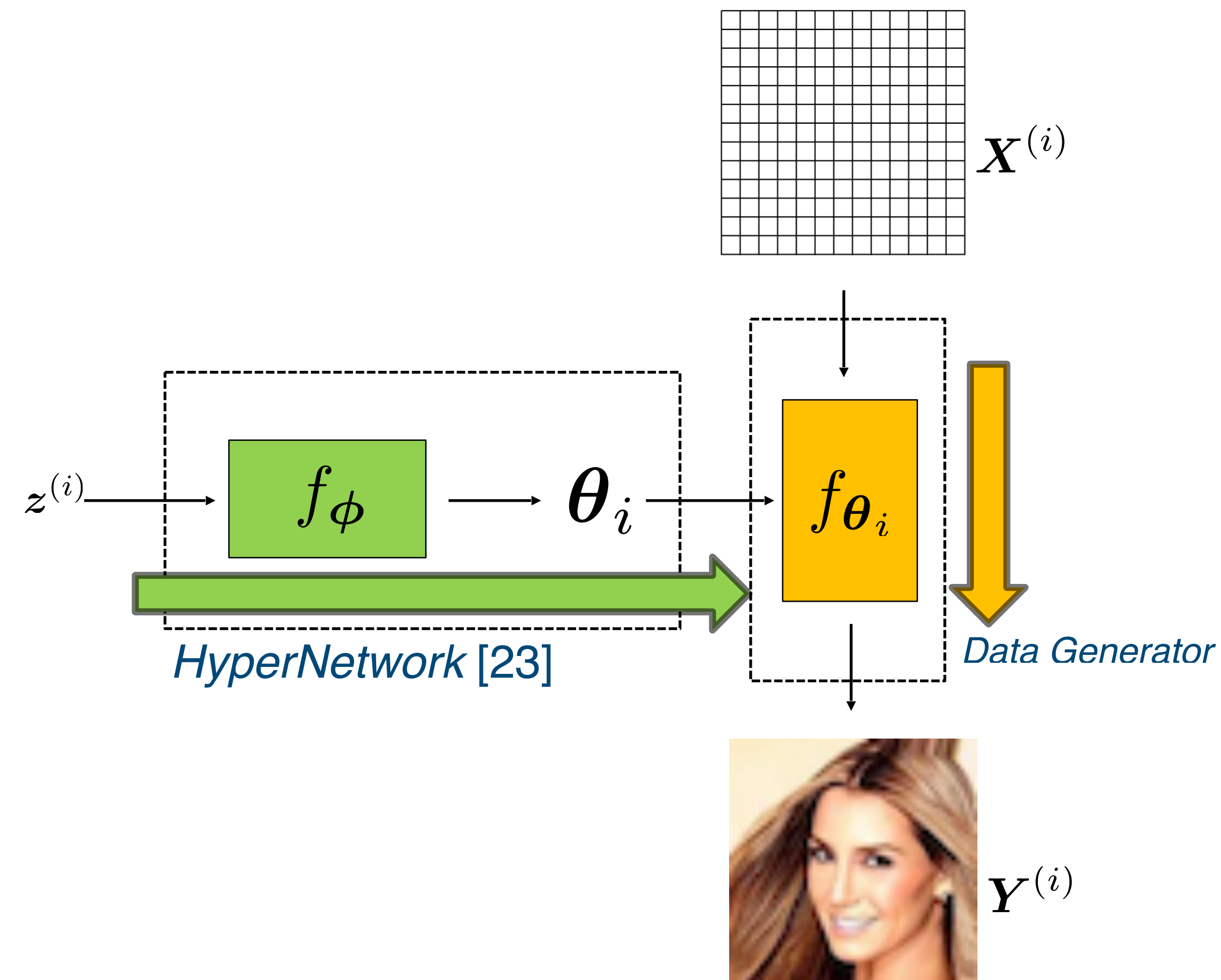
[23] Ha et al., 2017



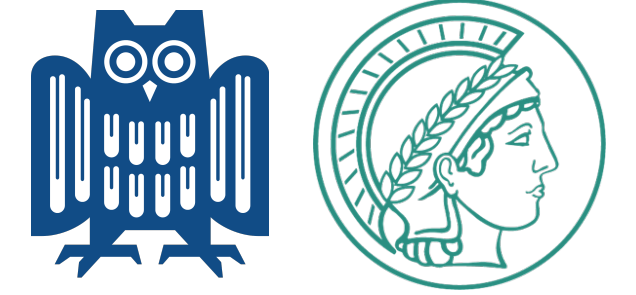


# Deep Generative Models of INRs

Have  $z^{(i)}$ , a summary representation of image.



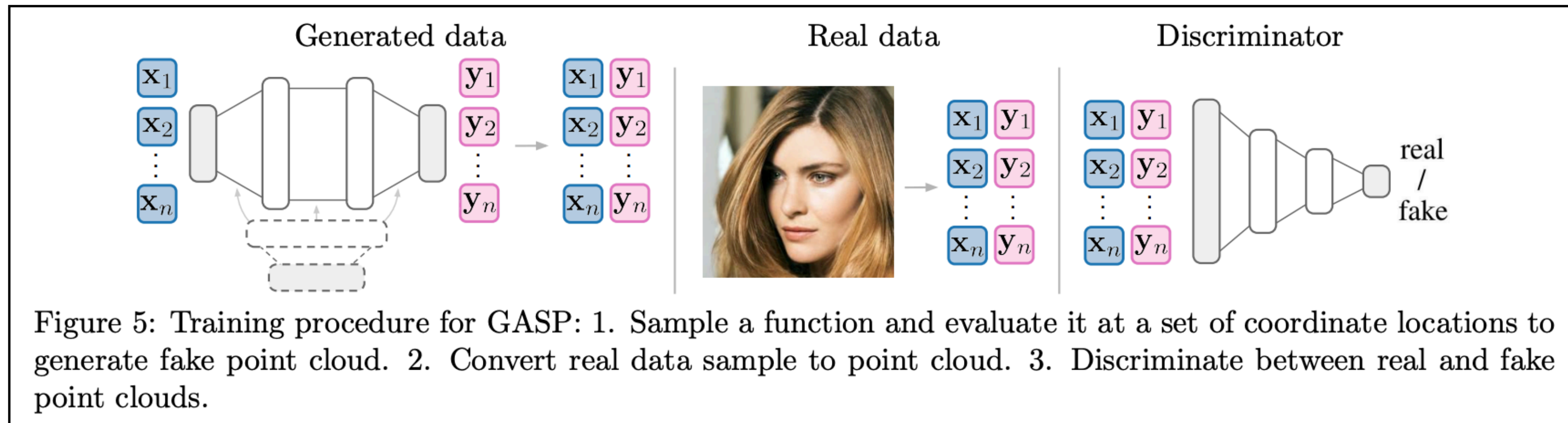
[23] Ha et al., 2017



# Previous work

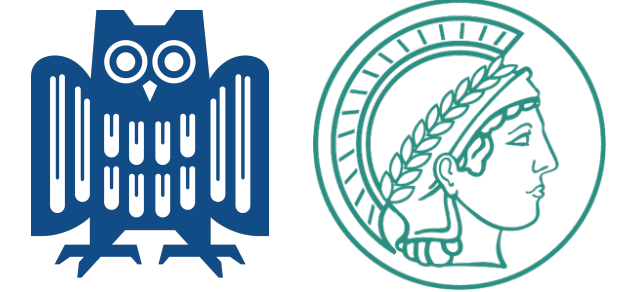
## GASP<sup>[5]</sup>

- Adversarial training:



**✗** Can't tackle inference related tasks.

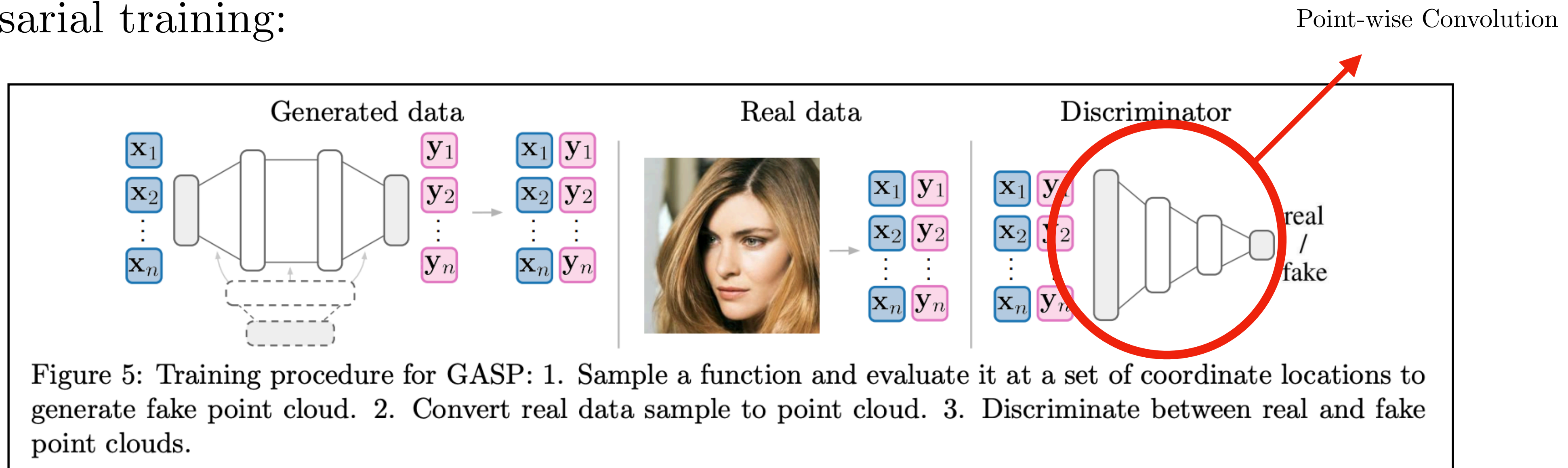
<sup>[5]</sup> Dupont et al., 2020



# Previous work

## GASP<sup>[5]</sup>

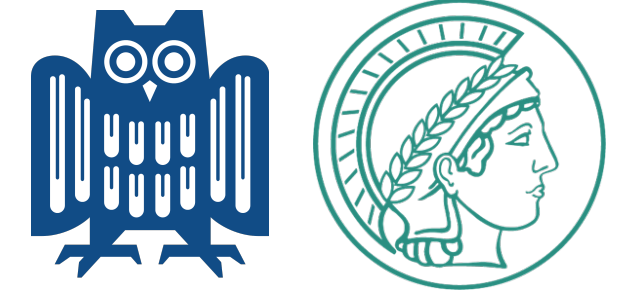
- Adversarial training:



✘ Can't tackle inference related tasks.

<sup>[5]</sup> Dupont et al., 2020



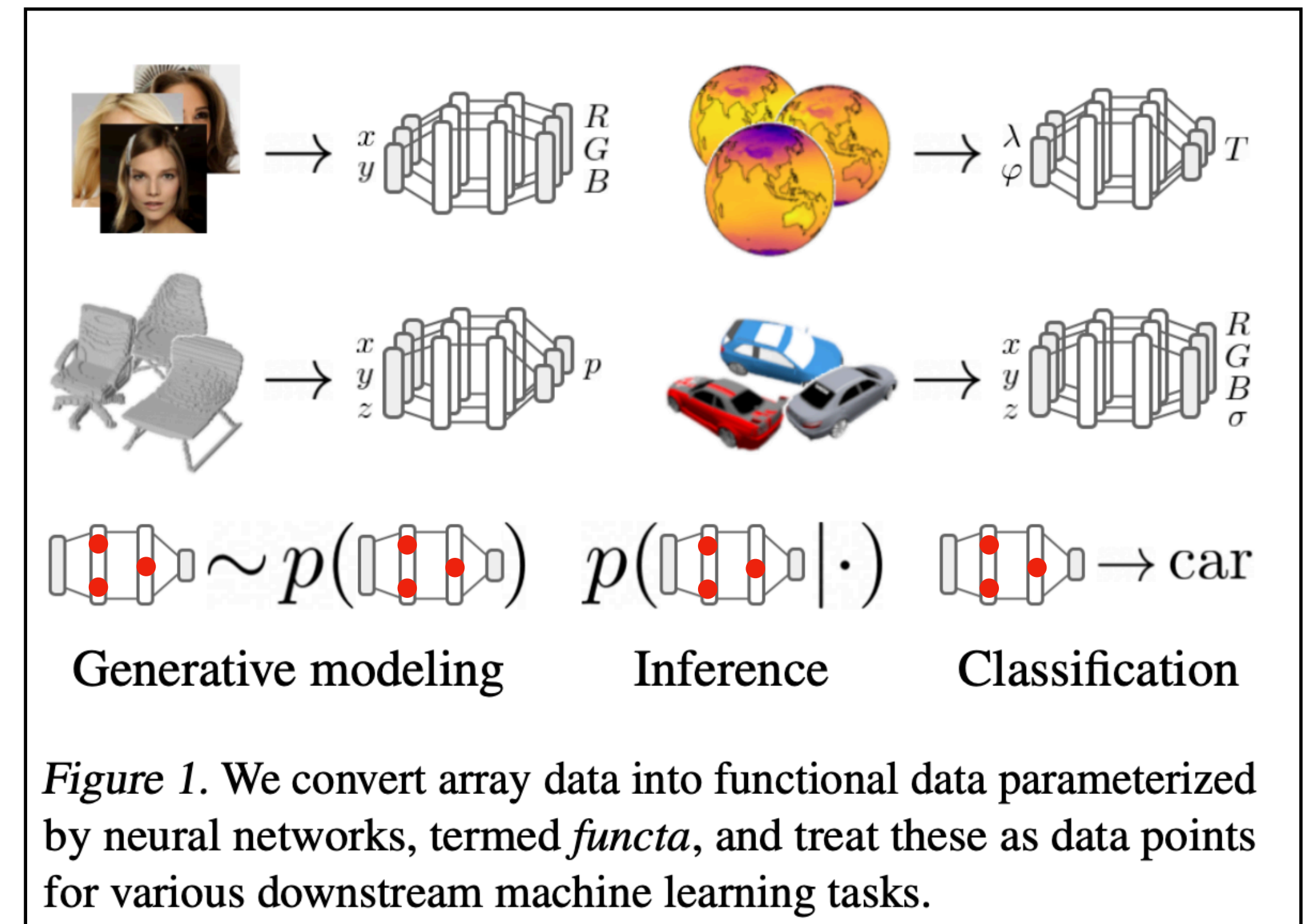


# Previous work

## Functas<sup>[6]</sup>

- Decoupled training:
  1. Fit an INR per datapoint using SIREN<sup>[20]</sup> and **modulation vectors**, named **functas**.
  2. Train any generative model on the functa dataset of vectors.

✘ Computationally expensive inference.

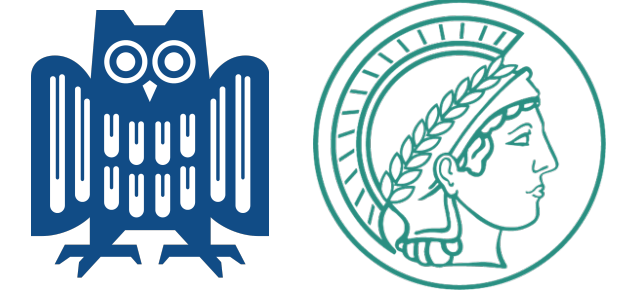


[6] Dupont et al., 2022

[20] Sitzmann et al., 2020

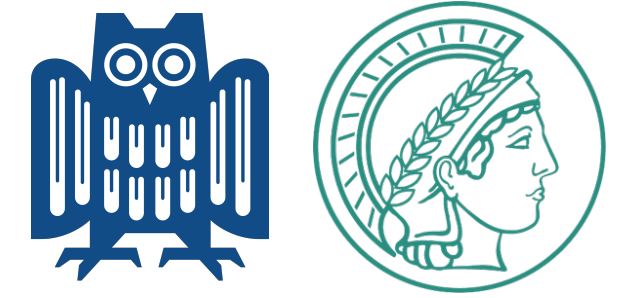


# Deep Generative Models of INRs



How to infer the latent representation  $\mathbf{z}$ ?

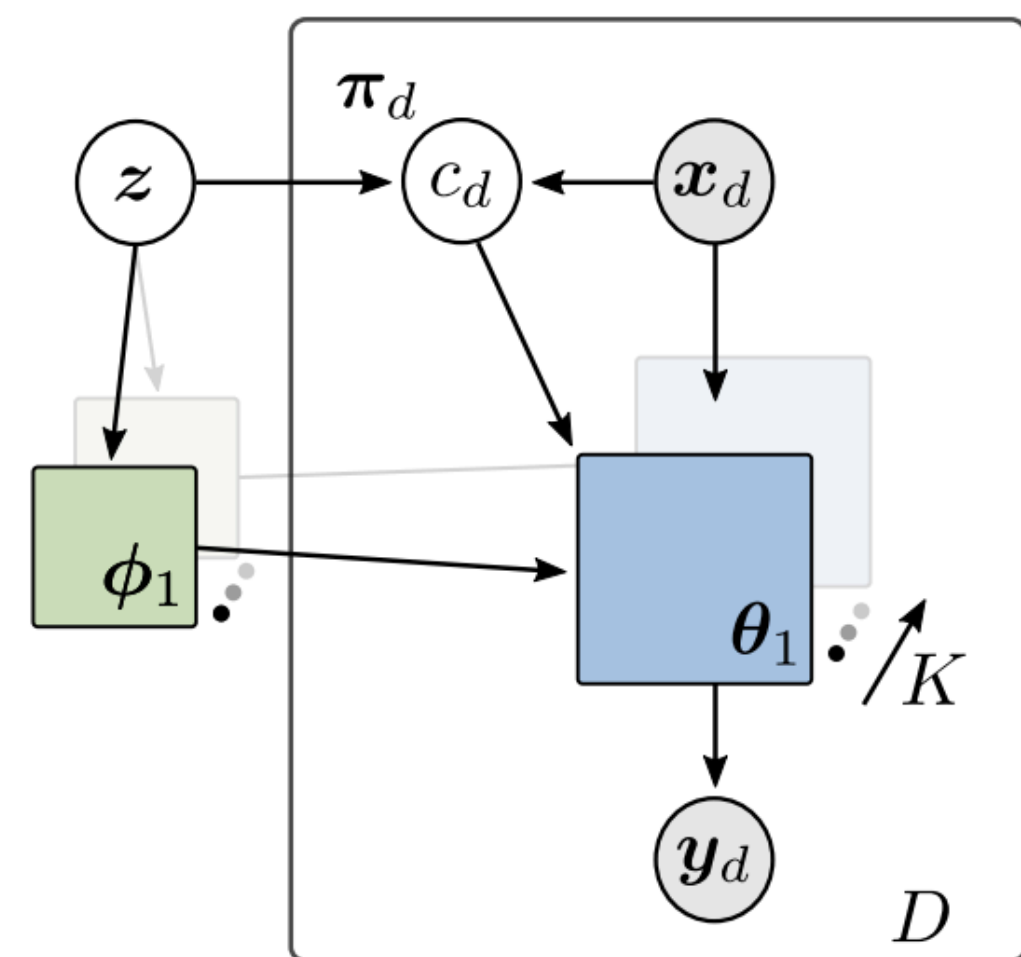
$$q_{\gamma}(\mathbf{z}|\mathbf{Y}, \mathbf{X}) \quad p_{\psi}(\mathbf{z})$$



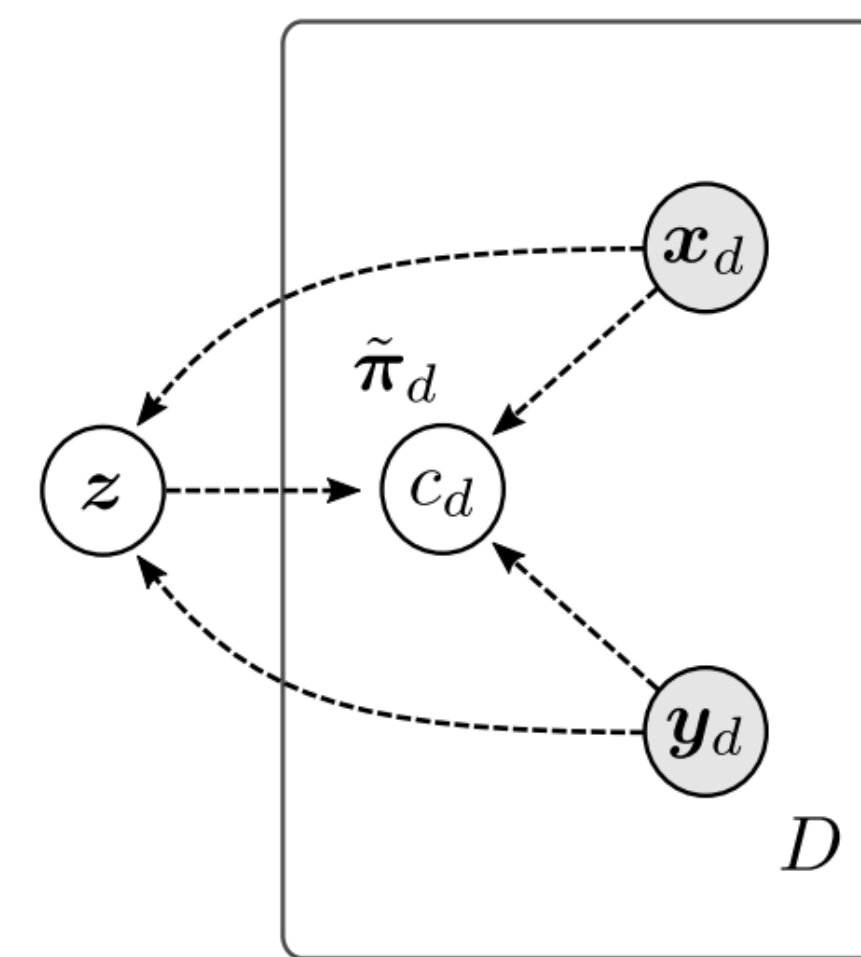
# Proposed method: VAMoH

## Variational Mixture of HyperGenerators

HyperNetwork
  Data generator



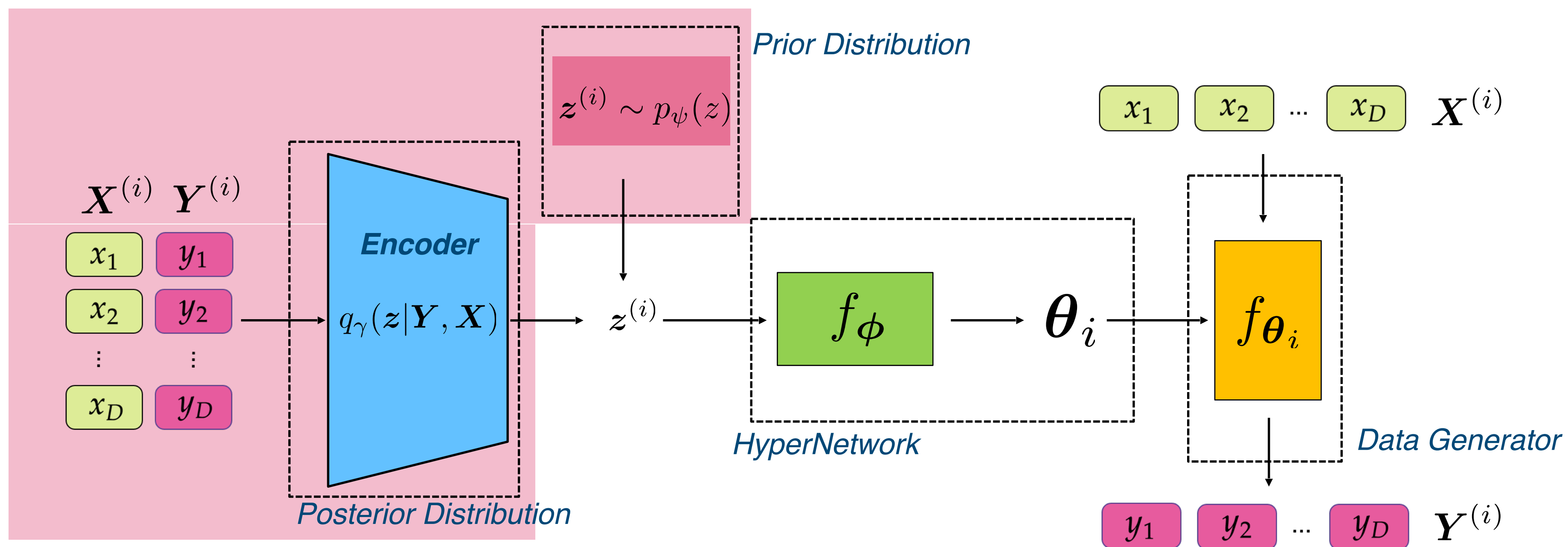
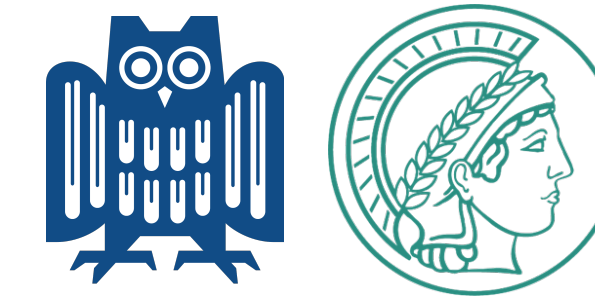
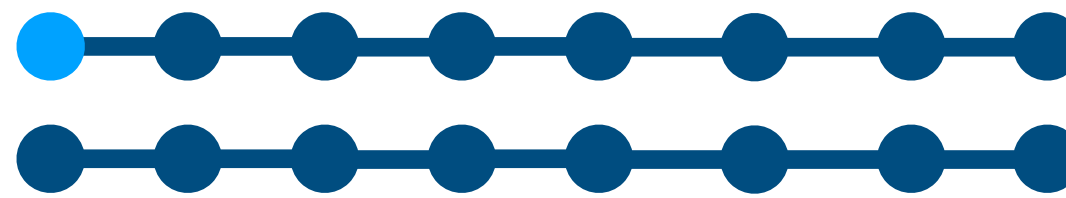
(a) Generative model



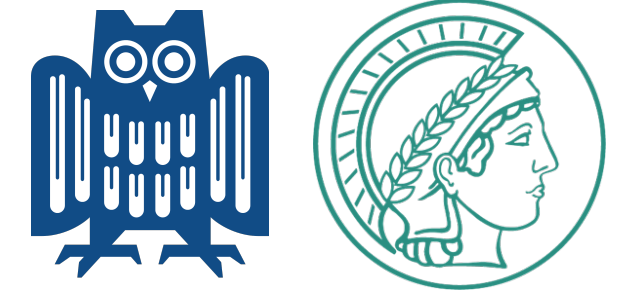
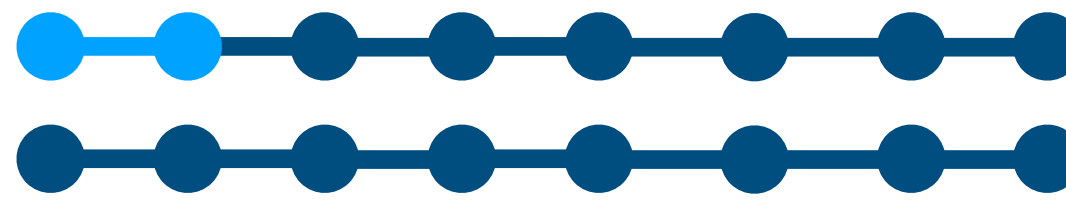
(b) Inference model

# VAMoH

## Encoder



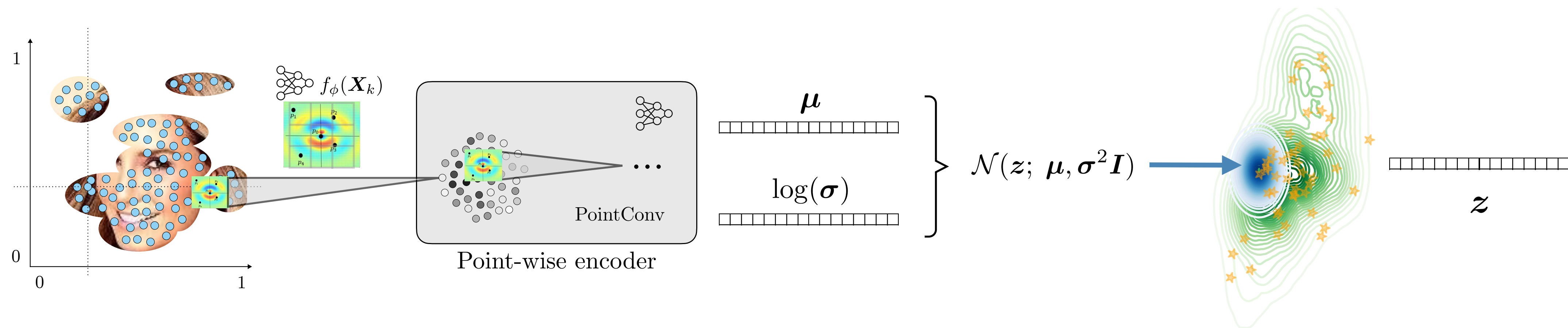
$z^{(i)}$ : Latent Variable



# VAMoH

## Encoder

- PointConv<sup>[21]</sup> encoder for point clouds.

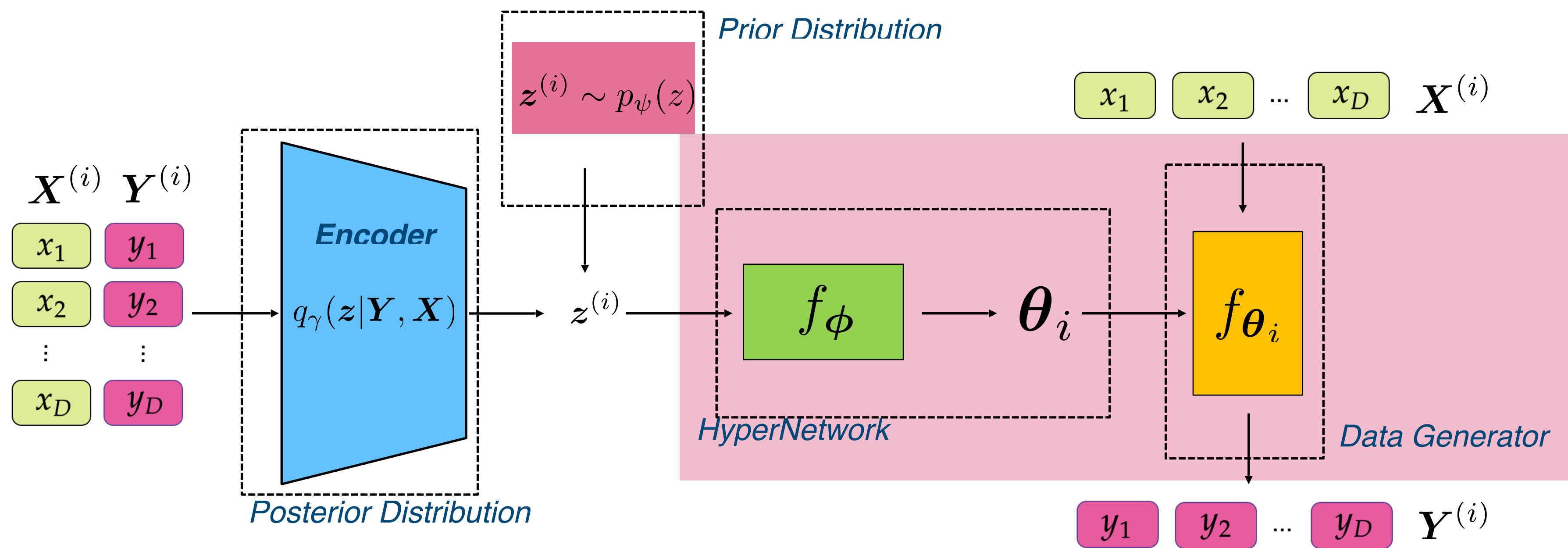
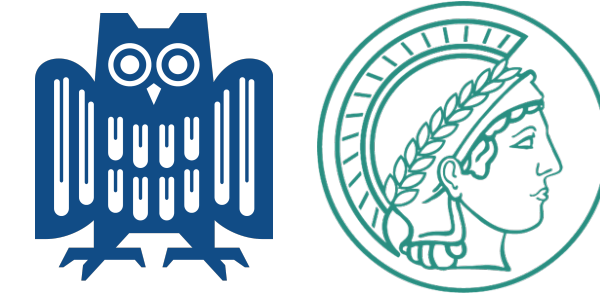
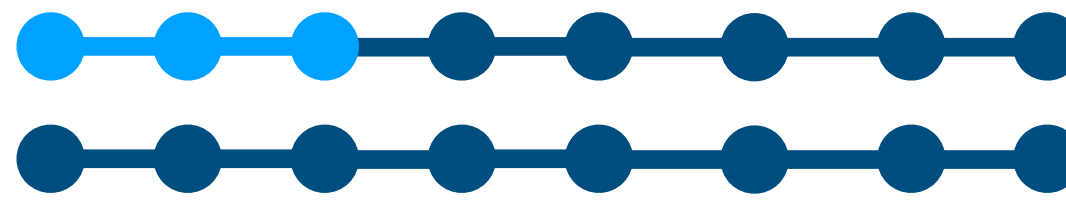


[21] Wu et al., 2019



# VAMoH

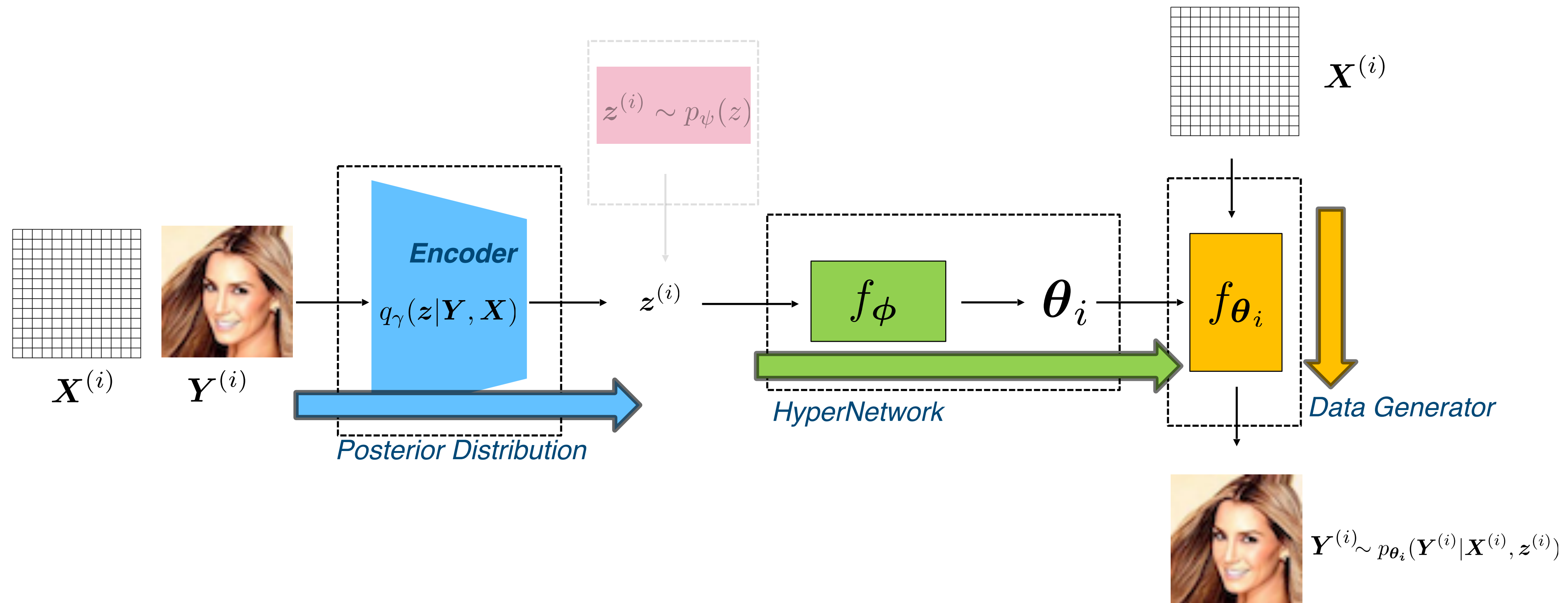
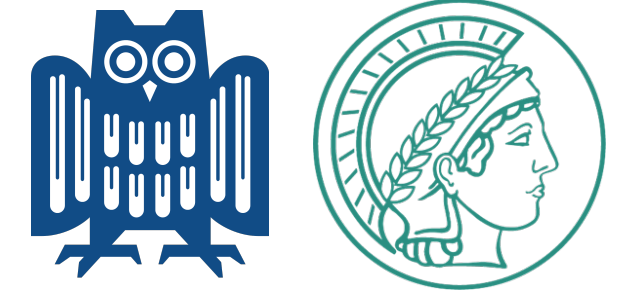
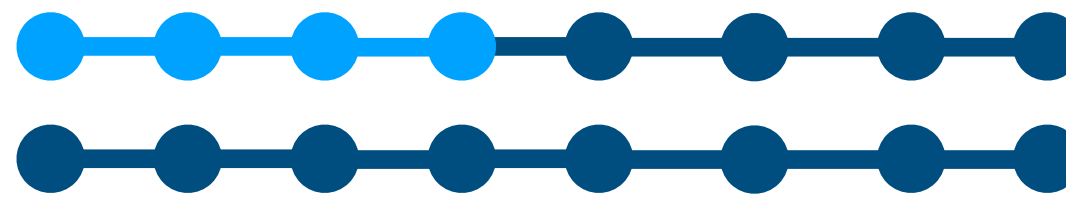
## Decoder



$z^{(i)}$ : Latent Variable

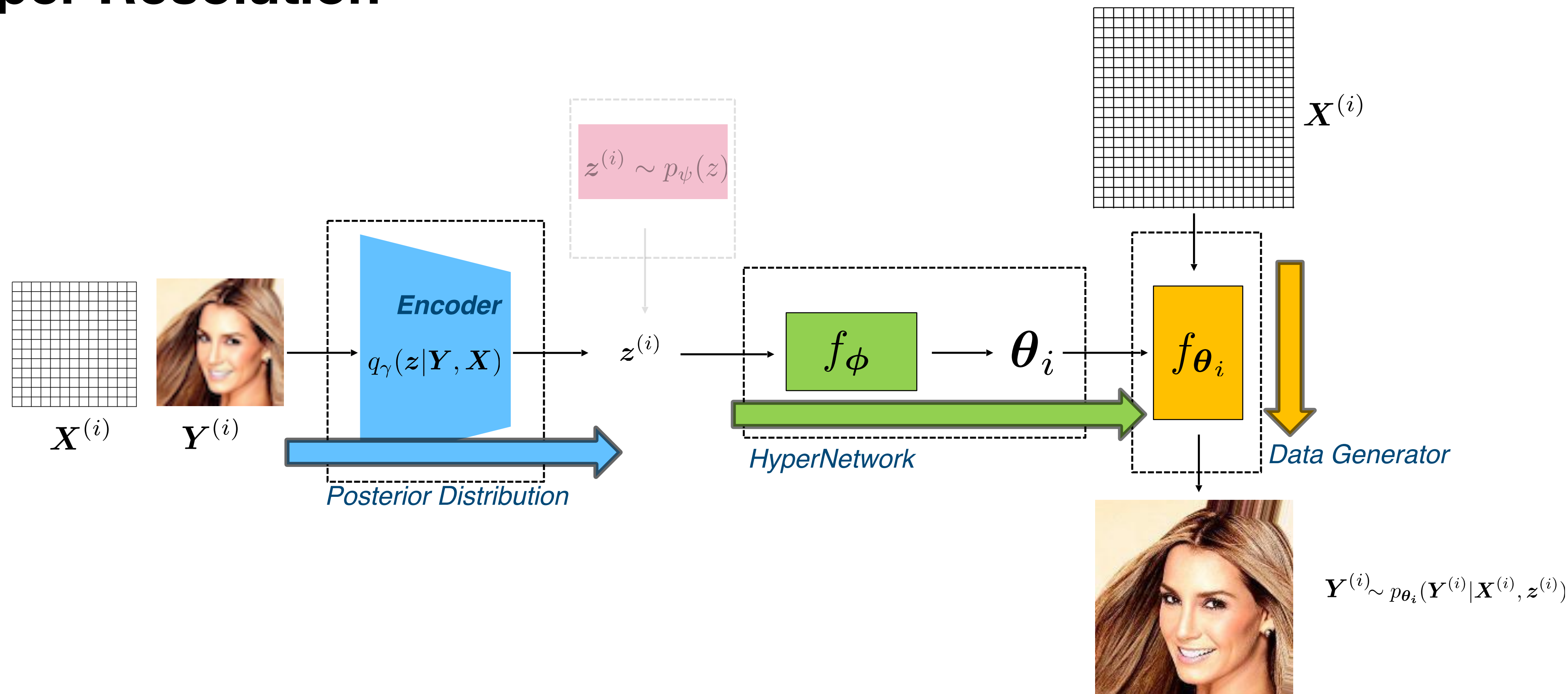
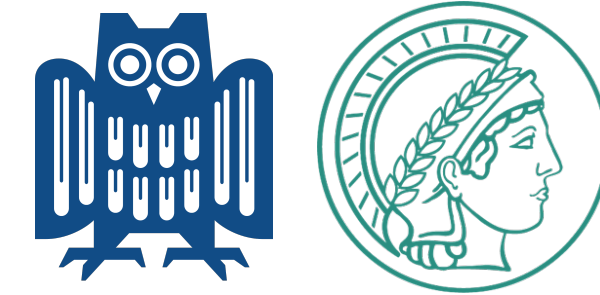
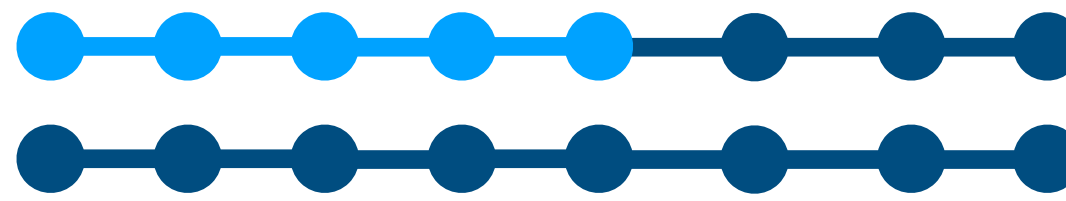
# VAMoH

## Reconstruction



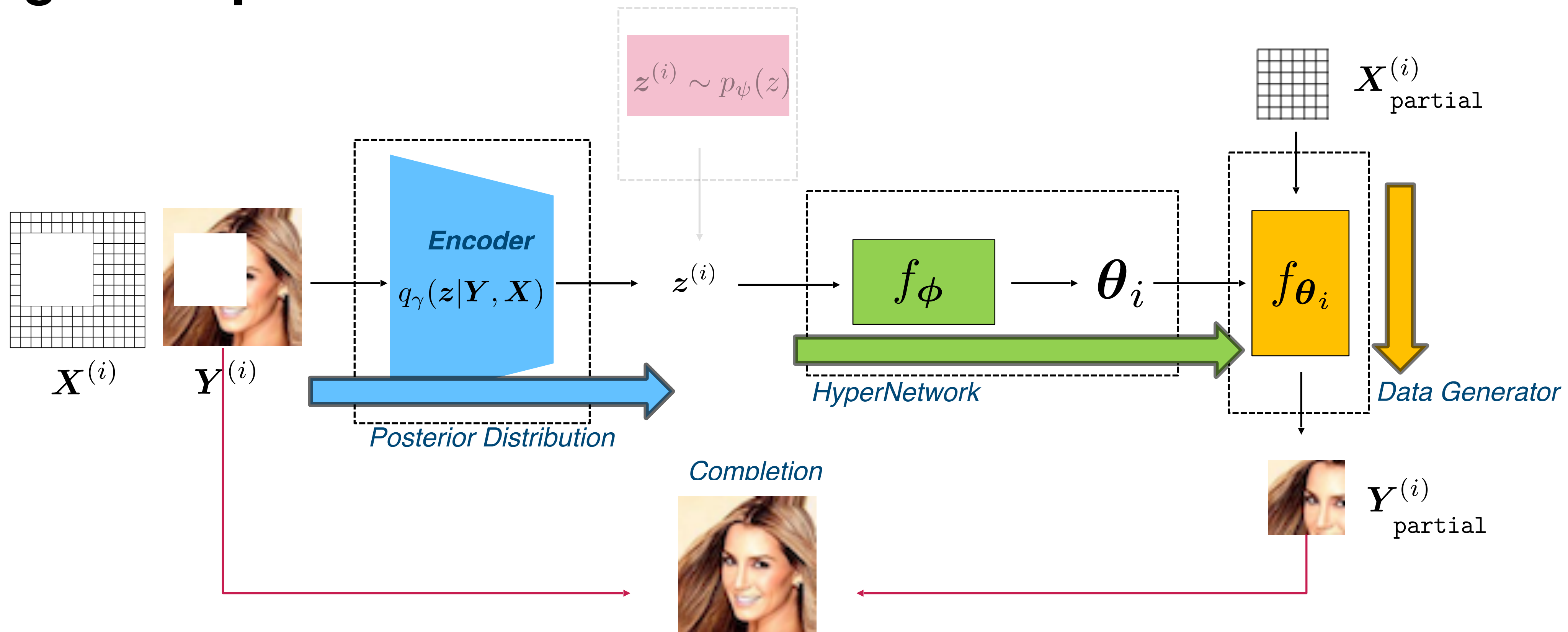
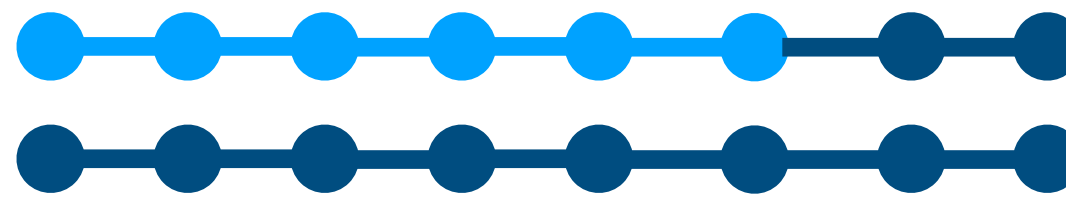
# VAMoH

## Super Resolution



# VAMoH

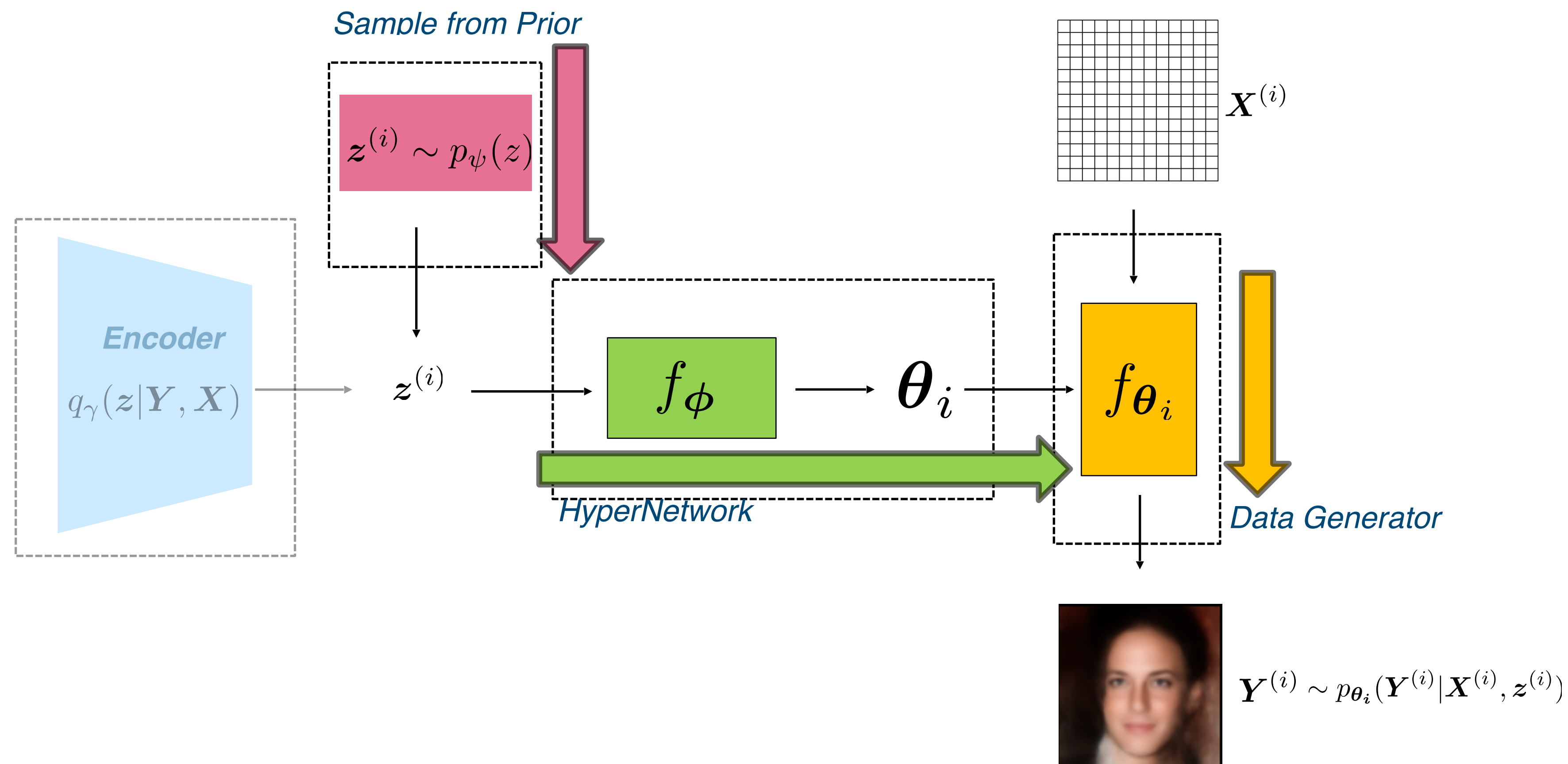
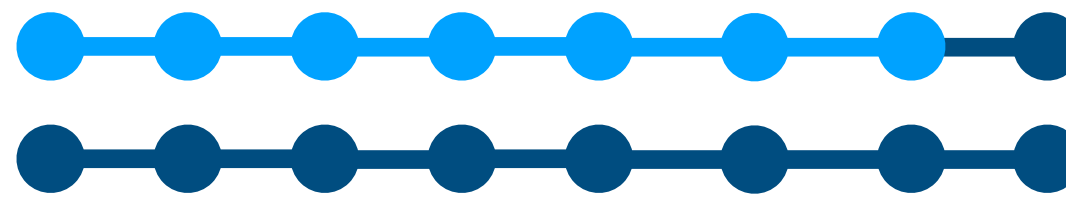
## Image Completion





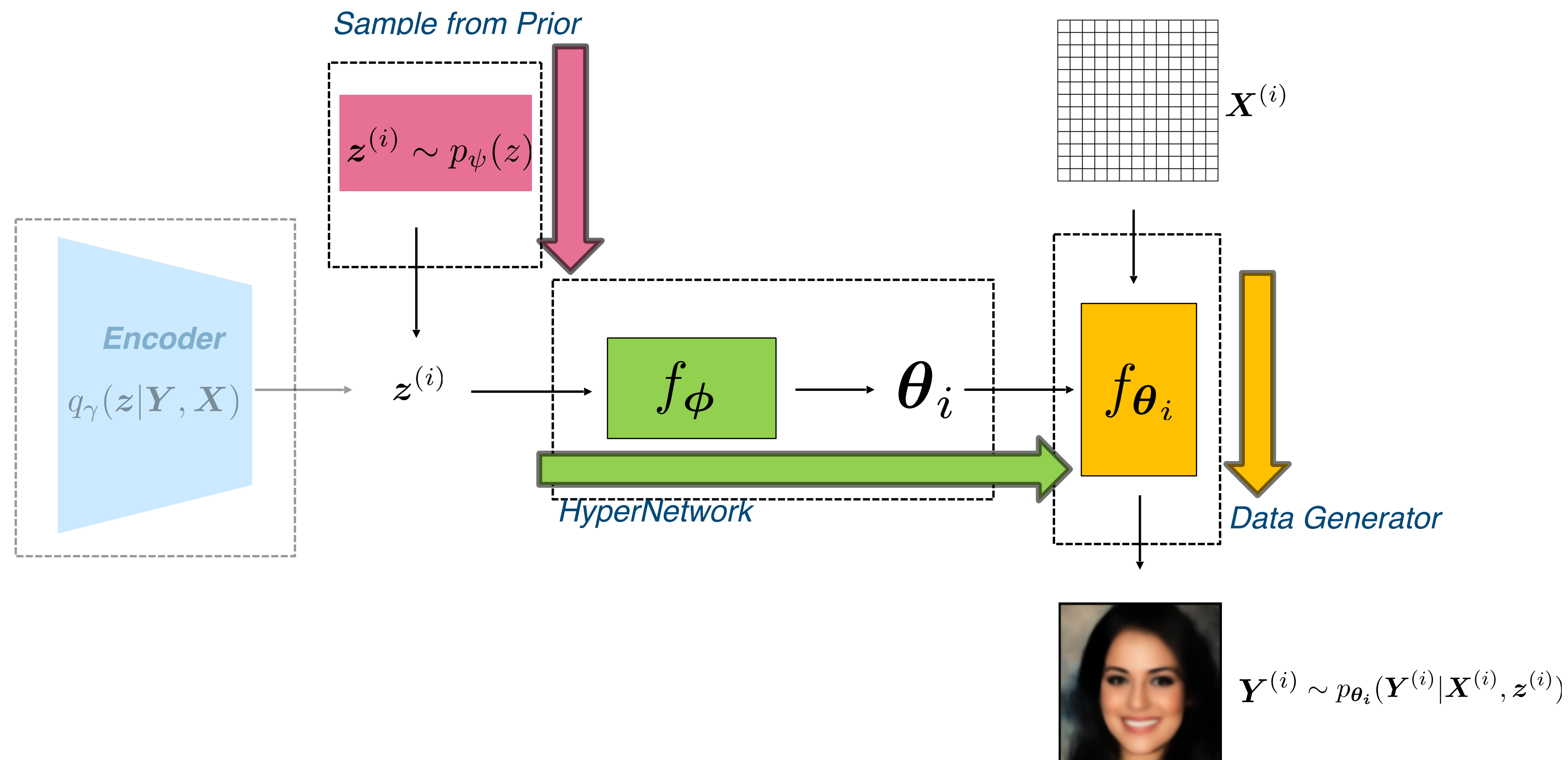
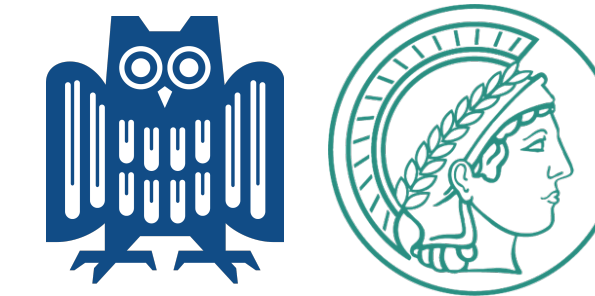
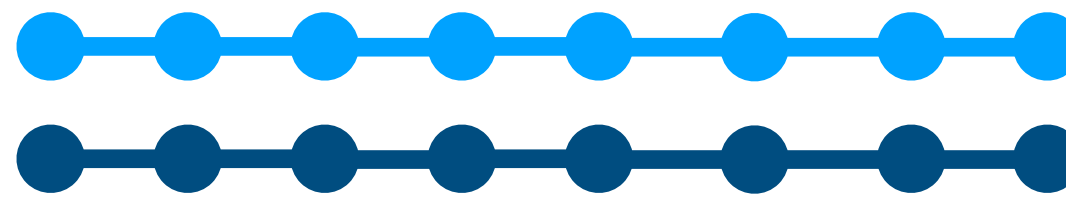
# VAMoH

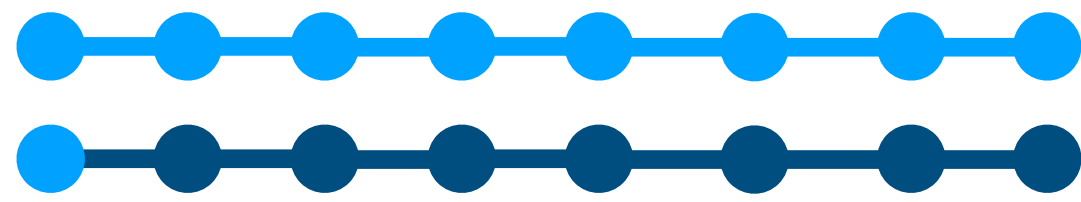
## Image Generation



# VAMoH

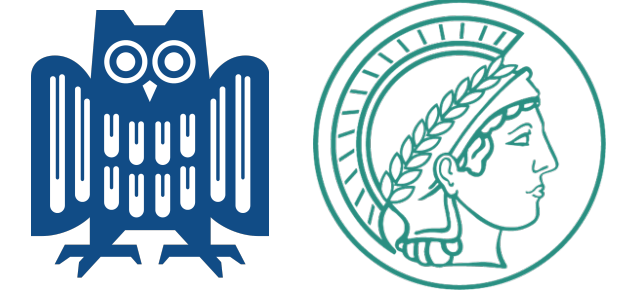
## Image Generation





# VAMoH

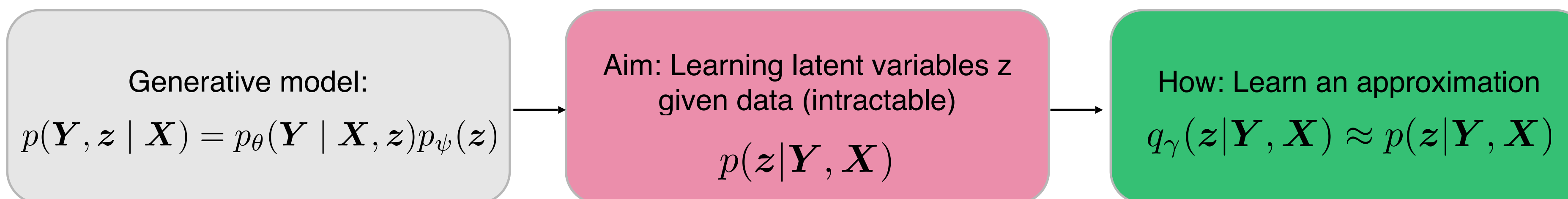
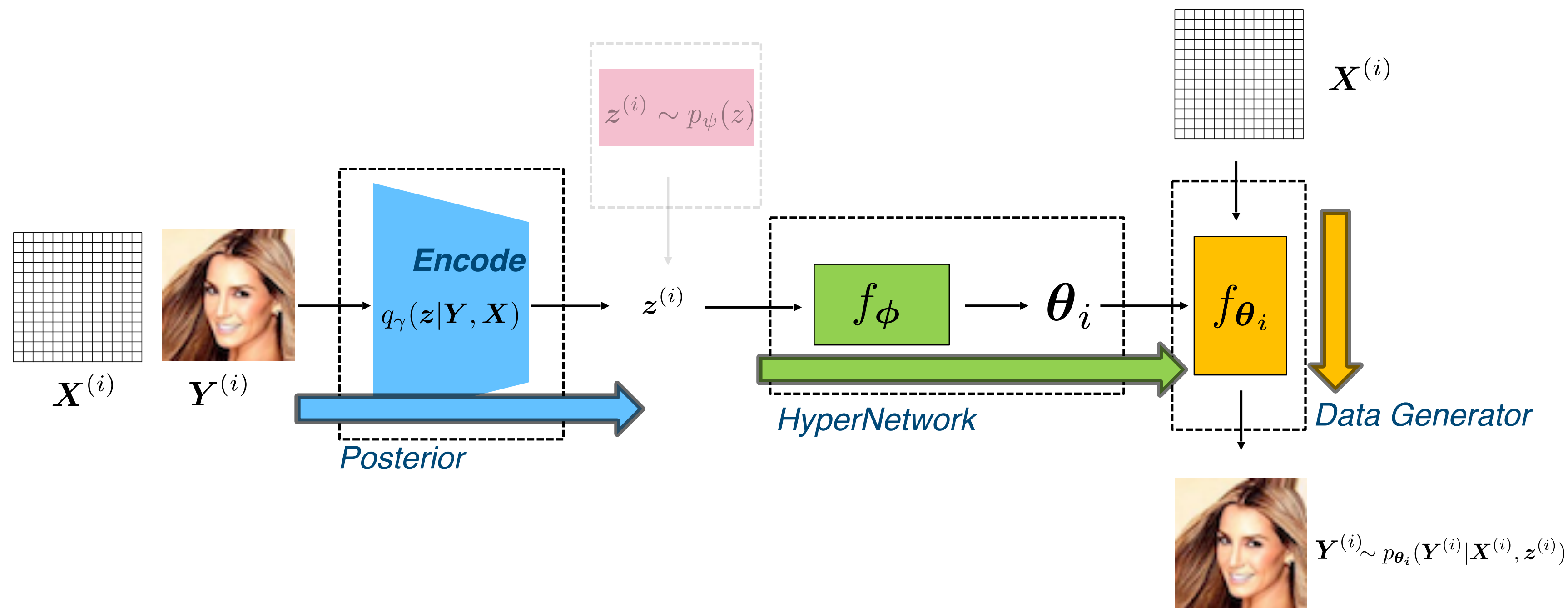
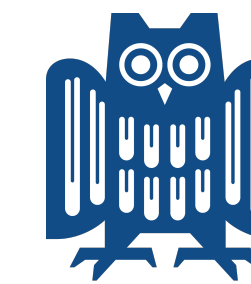
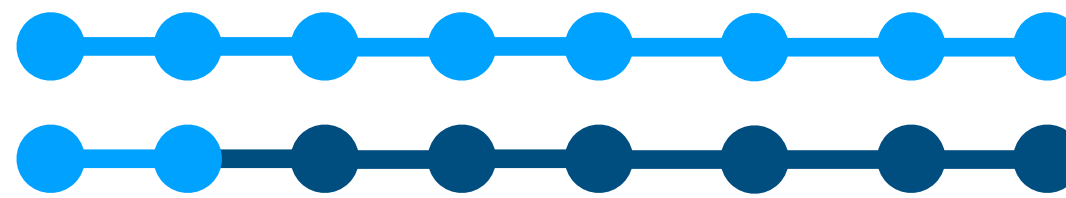
## Optimization



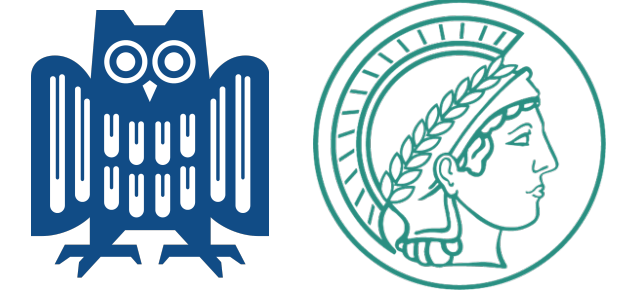
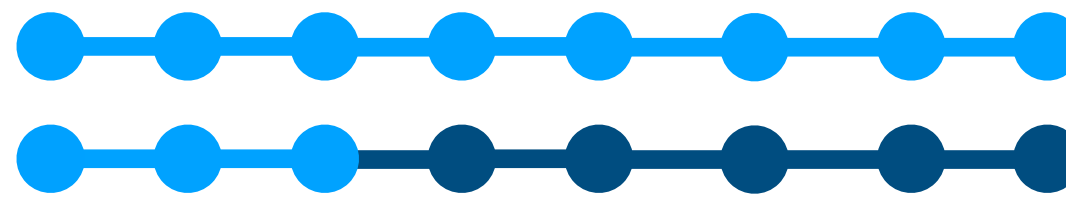
How to learn all these steps end-to-end from data?

# VAMoH

## Optimization







# VAMoH

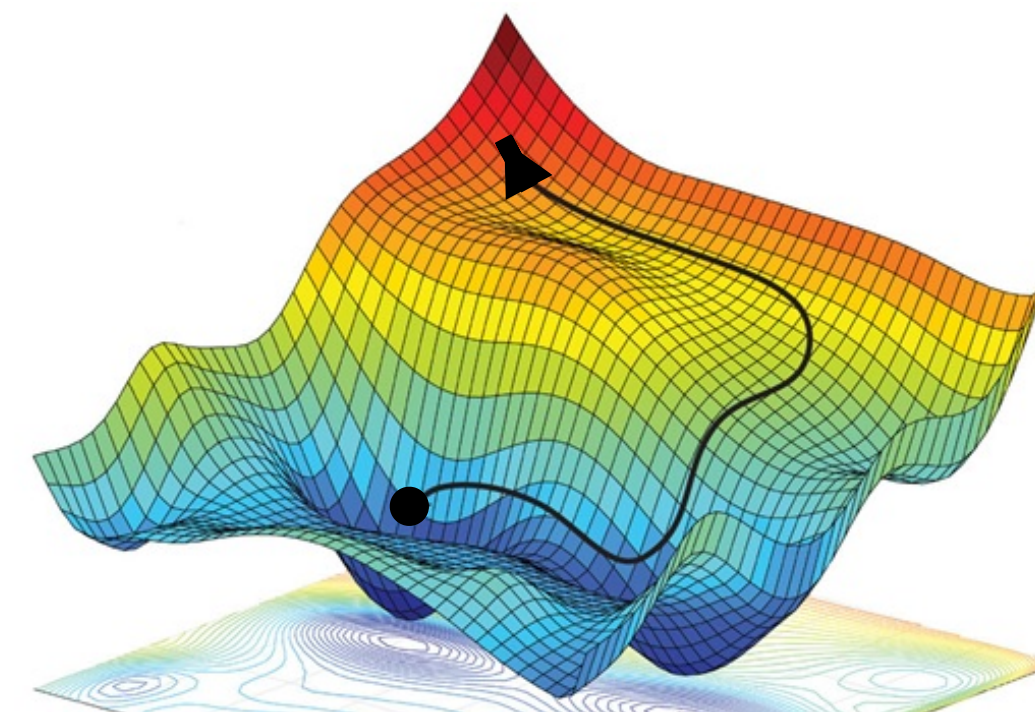
## Optimization

- For a single data sample  $(\mathbf{X}, \mathbf{Y})$

$$\max_{\phi, \gamma} \mathcal{L}(\phi, \gamma; \mathbf{Y}, \mathbf{X}) = \underbrace{\max_{\phi, \gamma} \mathbb{E}_{q_{\gamma}(z|\mathbf{Y}, \mathbf{X})} [\log p_{\theta}(\mathbf{Y} | \mathbf{X}, z)]}_{\text{Reconstruction}} - \underbrace{D_{KL}(q_{\gamma}(z|\mathbf{Y}, \mathbf{X}) || p_{\psi}(z))}_{\text{Regularization}}$$

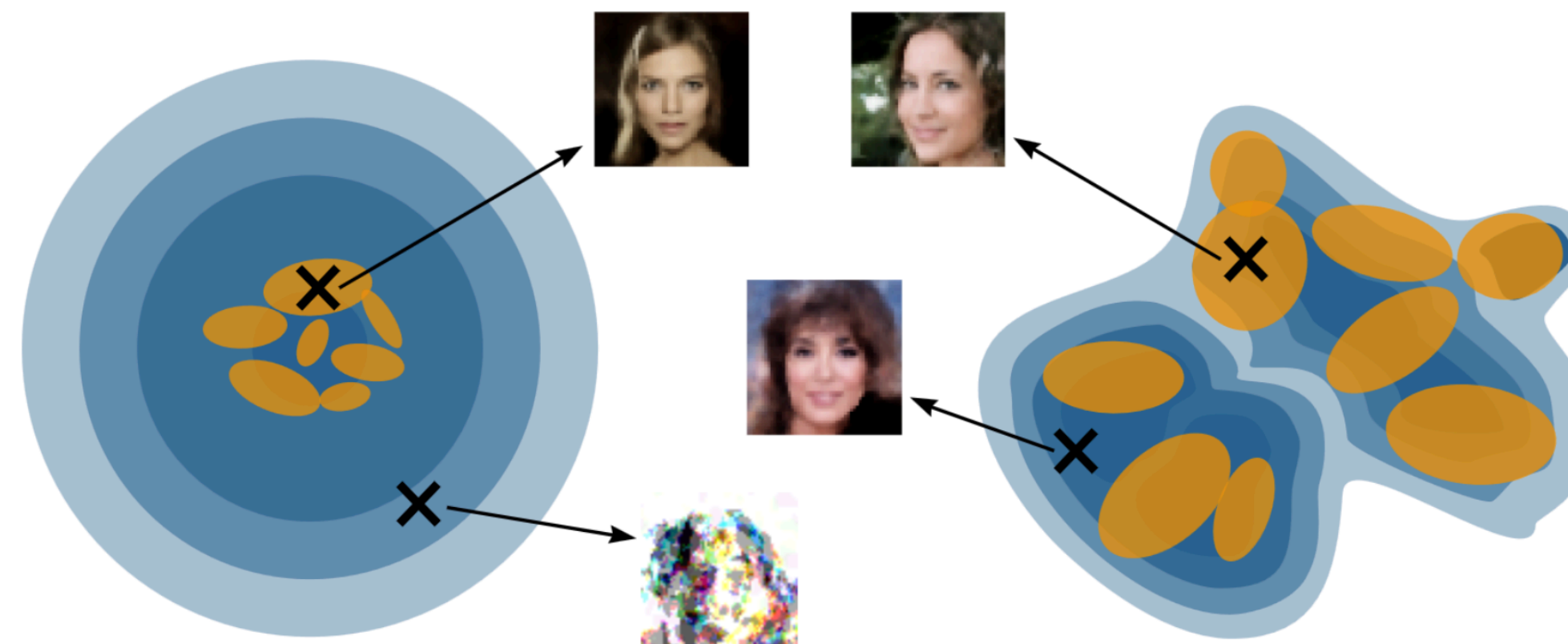
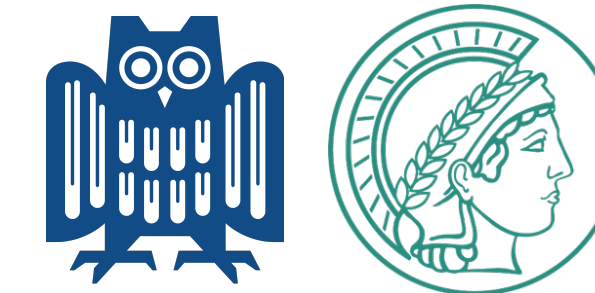
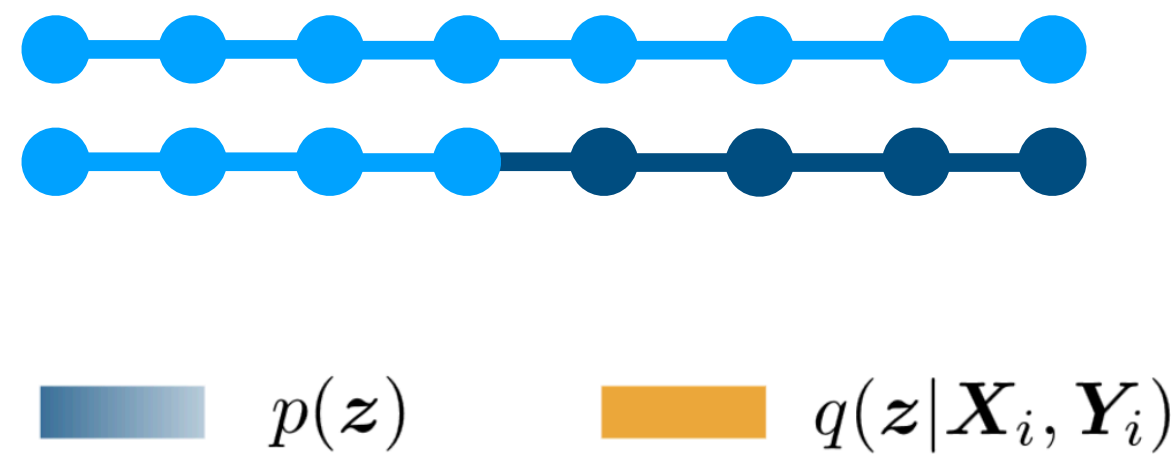
- For all samples in our dataset  $(\mathbf{X}^{(i)}, \mathbf{Y}^{(i)})$ ,  $i \in [N]$

$$\max_{\phi, \gamma} \sum_{i=1}^N \mathcal{L}(\phi, \gamma; \mathbf{Y}^{(i)}, \mathbf{X}^{(i)})$$



# VAMoH

## 'Holes' problem



Regularization Term:

$$\min_{\gamma} D_{KL}(q_{\gamma}(z | \mathbf{Y}, \mathbf{X}) || p_{\psi}(z))$$

We need to align the approximate posterior with the prior.

$$p_{\psi}(z) \quad q_{\gamma}(z)$$

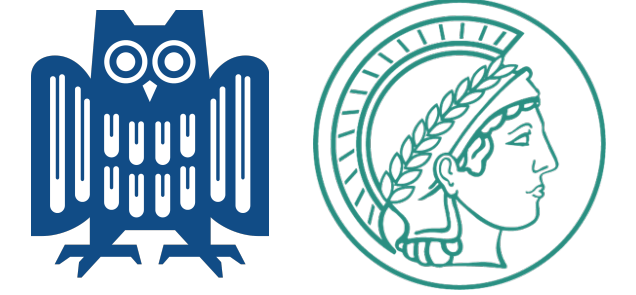
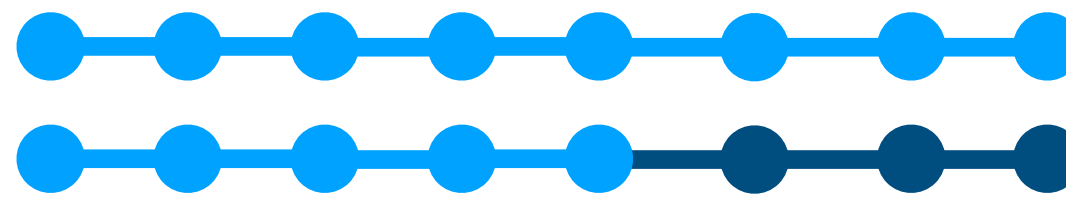
Problem:

If the prior is too simple, it hinders generation quality.

Solution:

Learn a more complex  $p_{\psi}(z)$

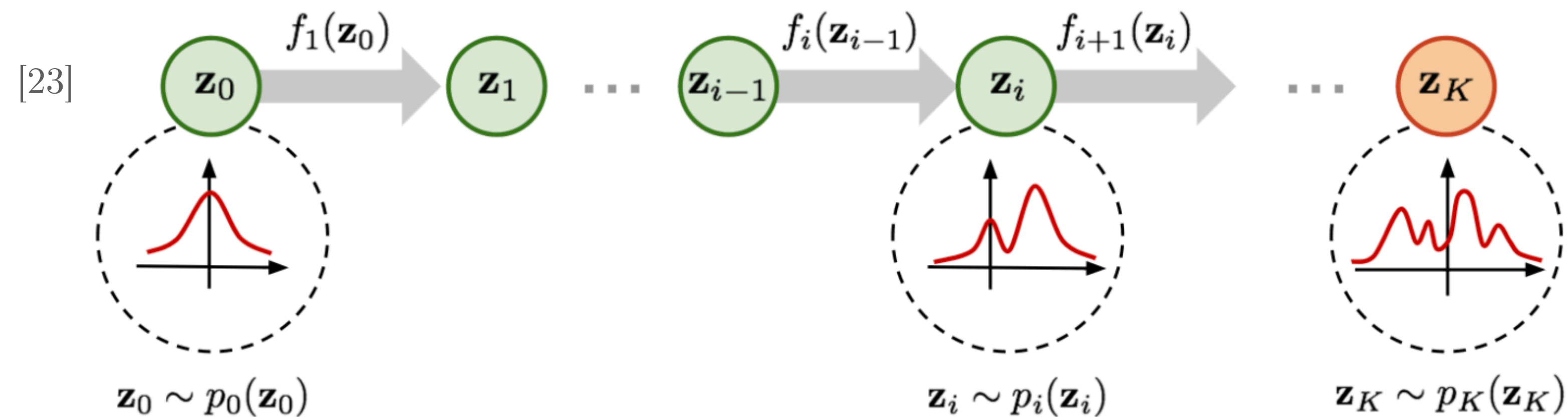
$$\min_{\gamma, \psi} D_{KL}(q_{\gamma}(z | \mathbf{Y}, \mathbf{X}) || p_{\psi}(z))$$



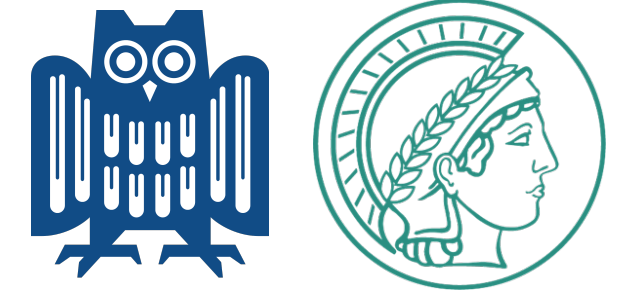
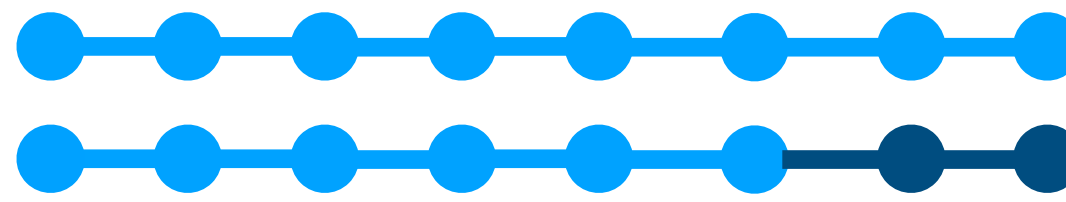
# VAMoH

## Flow-based prior

- More expressive prior using RealNVP (Real-valued, Non-Volume Preserving) Flow.



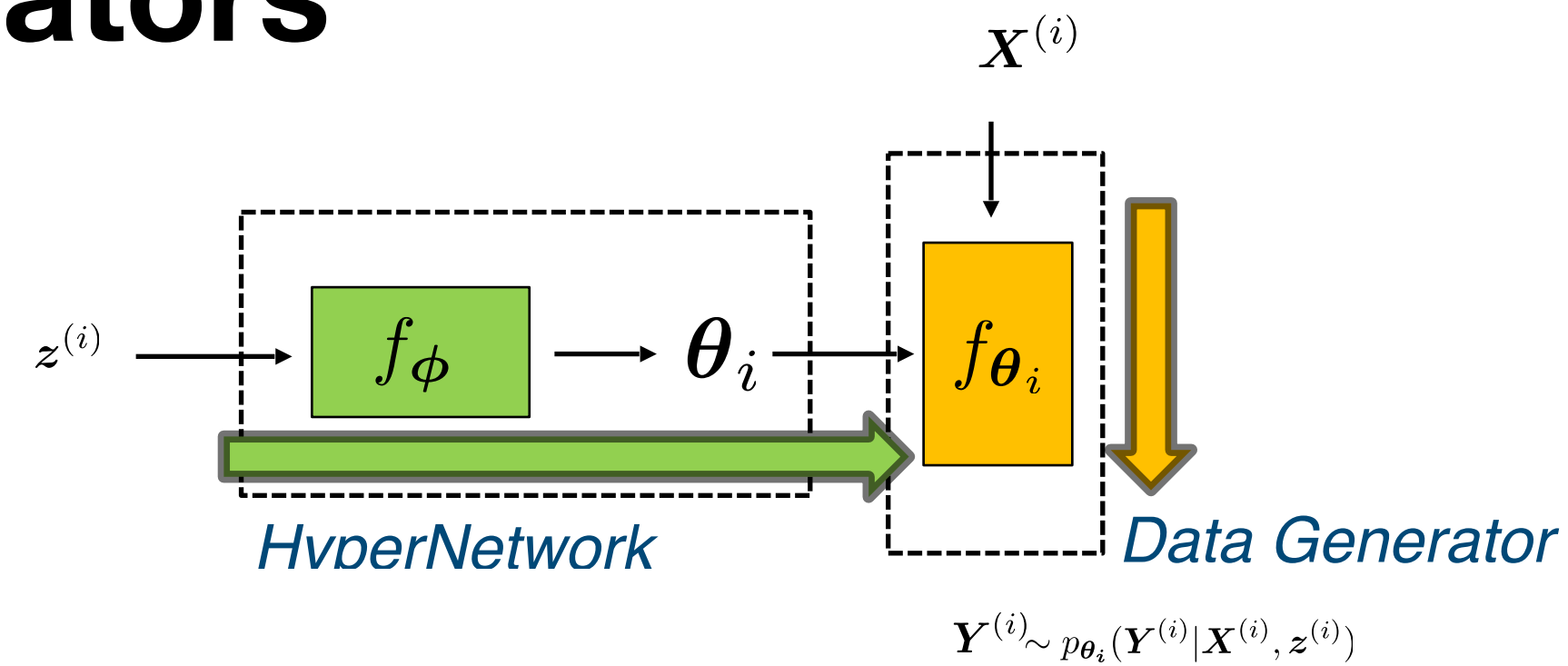
$$z^{(i)} \sim p_\psi(z)$$



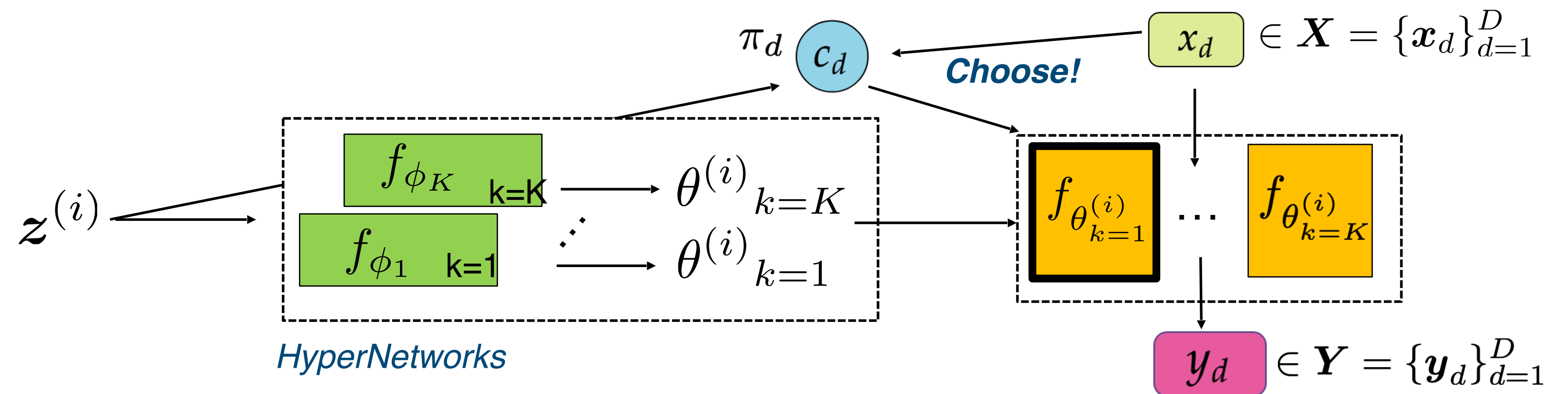
# VAMoH

## Mixture of HyperGenerators

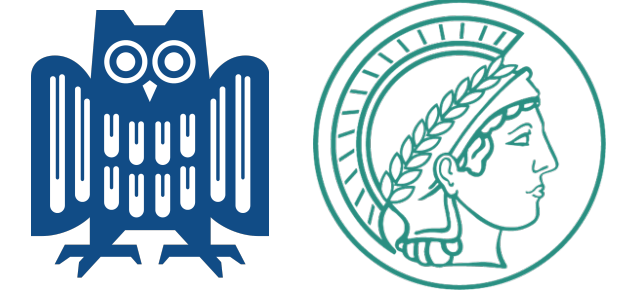
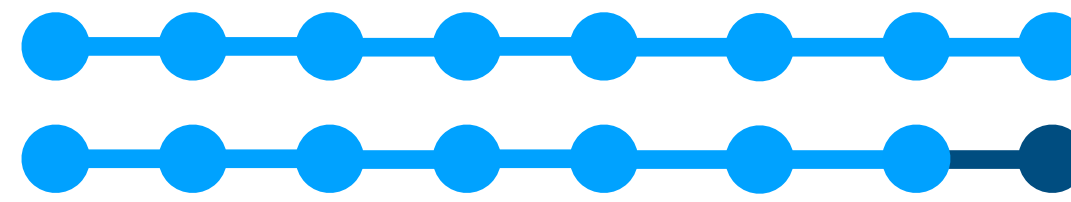
Single HyperGenerator



Mixture of HyperGenerators







# VAMoH

## Mixture of HyperGenerators

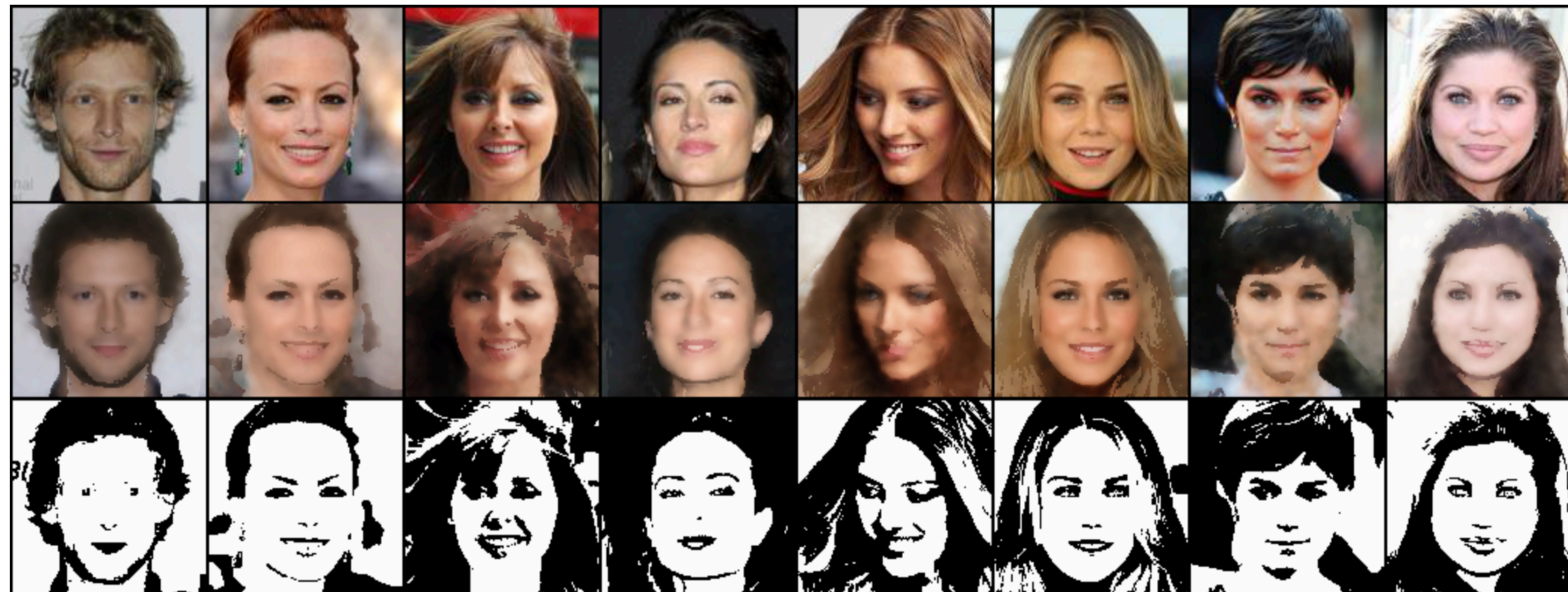
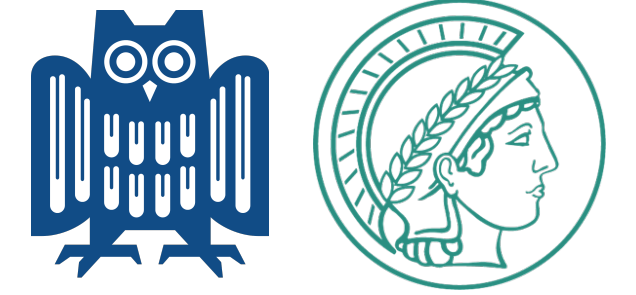
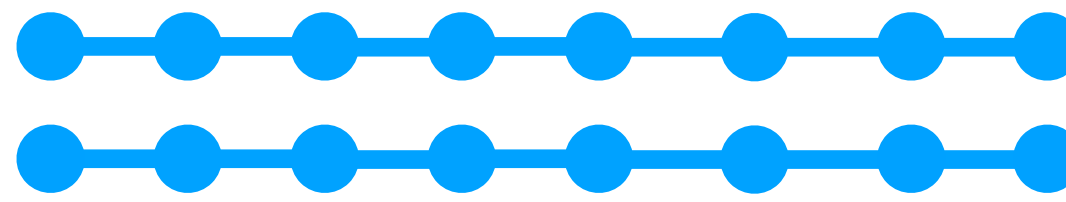


Image Reconstruction with Mixture of HyperGenerators



# VAMoH


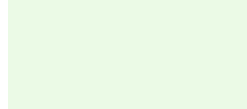
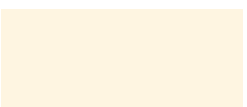
## ELBO

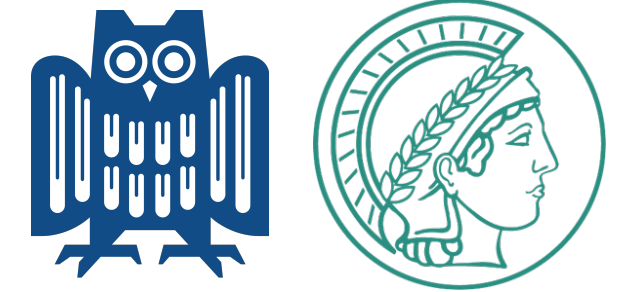
- For a single data sample  $(\mathbf{X}, \mathbf{Y})$

$$\mathcal{L}(\mathbf{Y}, \mathbf{X}; \psi, \phi, \gamma) = \sum_{d=1}^D \mathbb{E}_{q_{\gamma_z}(z|\mathbf{Y}, \mathbf{X})} \left[ \sum_{k=1}^K \log p_{\theta_k}(\mathbf{y}_d | \mathbf{x}_d) \cdot \pi_{dk} \right] - D_{KL}(q_{\gamma_z}(z | \mathbf{X}, \mathbf{Y}) \| p_{\psi_z}(z)) - D_{KL}(q_{\gamma_c}(\mathbf{C} | z, \mathbf{X}, \mathbf{Y}) \| p_{\psi_c}(\mathbf{C} | z, \mathbf{X}))$$

- For all samples in our dataset  $(\mathbf{X}^{(i)}, \mathbf{Y}^{(i)})$ ,  $i \in [N]$

$$\max_{\phi, \gamma, \psi} \sum_{i=1}^N \mathcal{L}(\phi, \gamma, \psi; \mathbf{Y}^{(i)}, \mathbf{X}^{(i)})$$

-  Reconstruction the features of the observed pixels
-  KL of the continuous latent variable
-  KL of the discrete latent variable

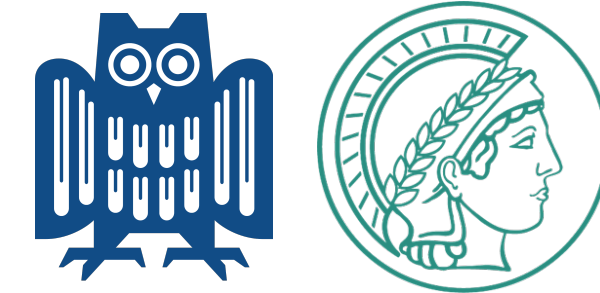


# Experiments

## Baselines


Model	Approach	Training Procedure	Generation	Reconstruction, Imputation, Super Resolution
GASP (2021) [5]	GAN	Minimax	Forward Pass	✗



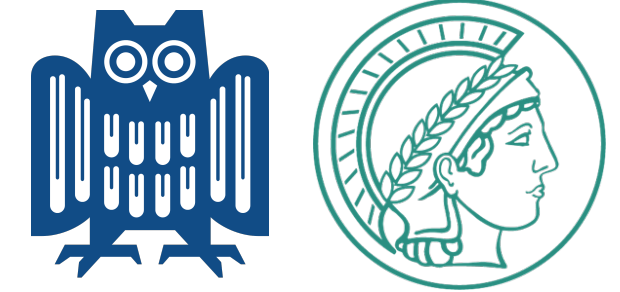


# Experiments

## Baselines


Model	Approach	Training Procedure	Generation	Reconstruction, Imputation, Super Resolution
GASP (2021) [5]	GAN	Minimax	Forward Pass	✗
Functa (2022) [6]	Flow-based	Bilevel optimization	+ Extra Generative Model	Optimization procedure(s) per sample 

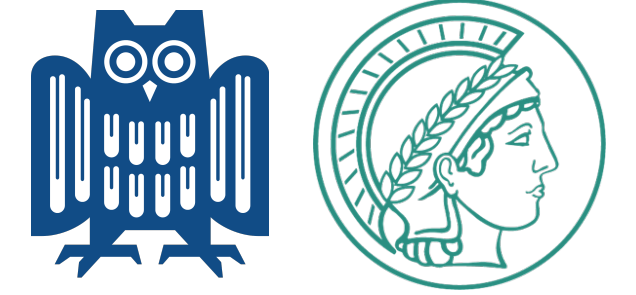




# Experiments



## Baselines

Model	Approach	Training Procedure	Generation	Reconstruction, Imputation, Super-Resolution
GASP (2021) [5]	GAN	Minimax	Forward Pass	$\min_{\phi} -\log p(\phi) + \lambda \sum_{i \in \mathcal{I}} \ f_{\phi}(\mathbf{x}_i) - \mathbf{f}_i\ _2^2$
Functa (2022) [6]	Flow-based	Bilevel optimization	+ Extra Generative Model	Optimization procedure(s) per sample 

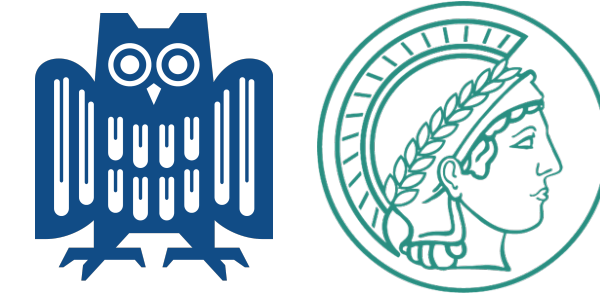


# Experiments

## Baselines

Model	Approach	Training Procedure	Generation	Reconstruction, Imputation, Super Resolution
GASP (2021) [5]	GAN	Minimax	Forward Pass	✗
Functa (2022) [6]	Flow-based	Bilevel optimization	+ Extra Generative Model	Optimization procedure(s) per sample 
VaMoH (ours)	VAE-based	<b>Single optimization</b>	<b>Forward Pass</b>	<b>Forward pass</b> 

VAMoH provides a probabilistic generative model that is efficient, robust, and expressive for modeling distribution over functions.



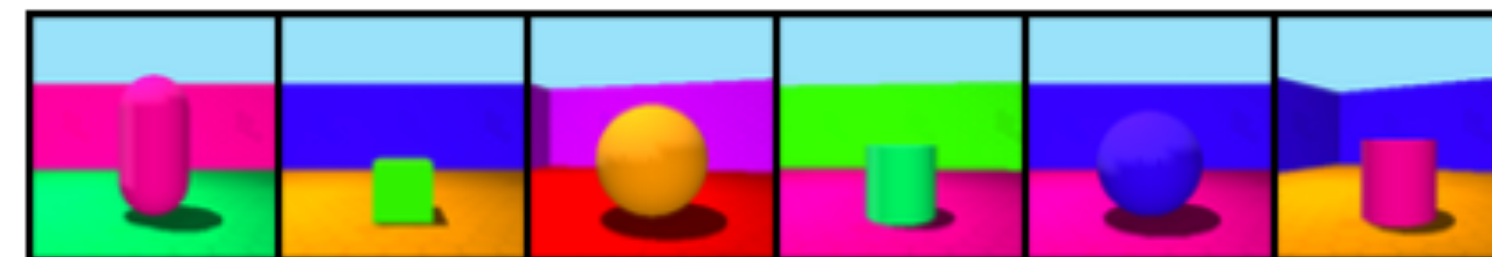
# Experiments

## Datasets

*PolyMNIST*  
(28x28)



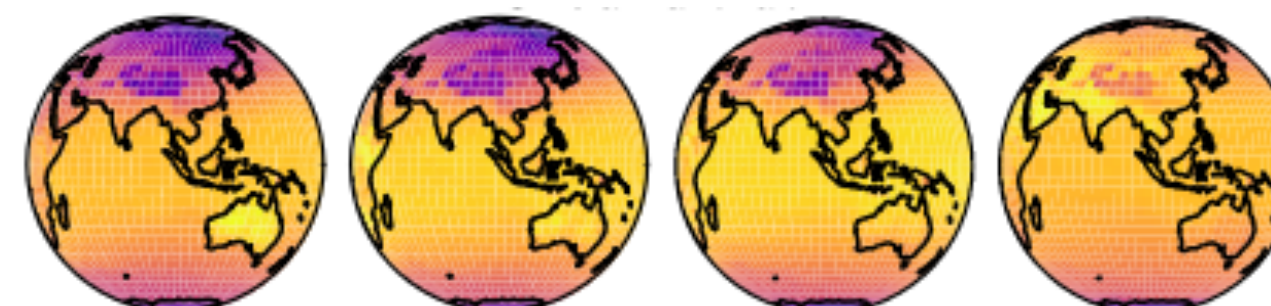
*Shapes3D* (64x64)



*CelebA-HQ* (64x64)



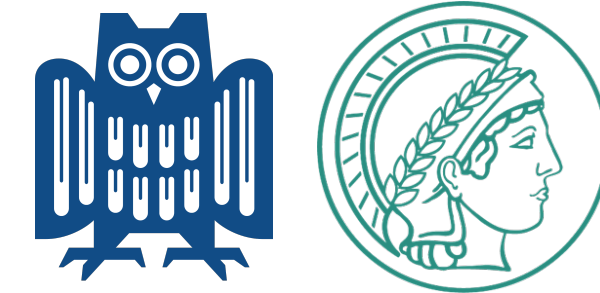
*ERA5 (Polar)*



*ShapeNET (Voxels)*





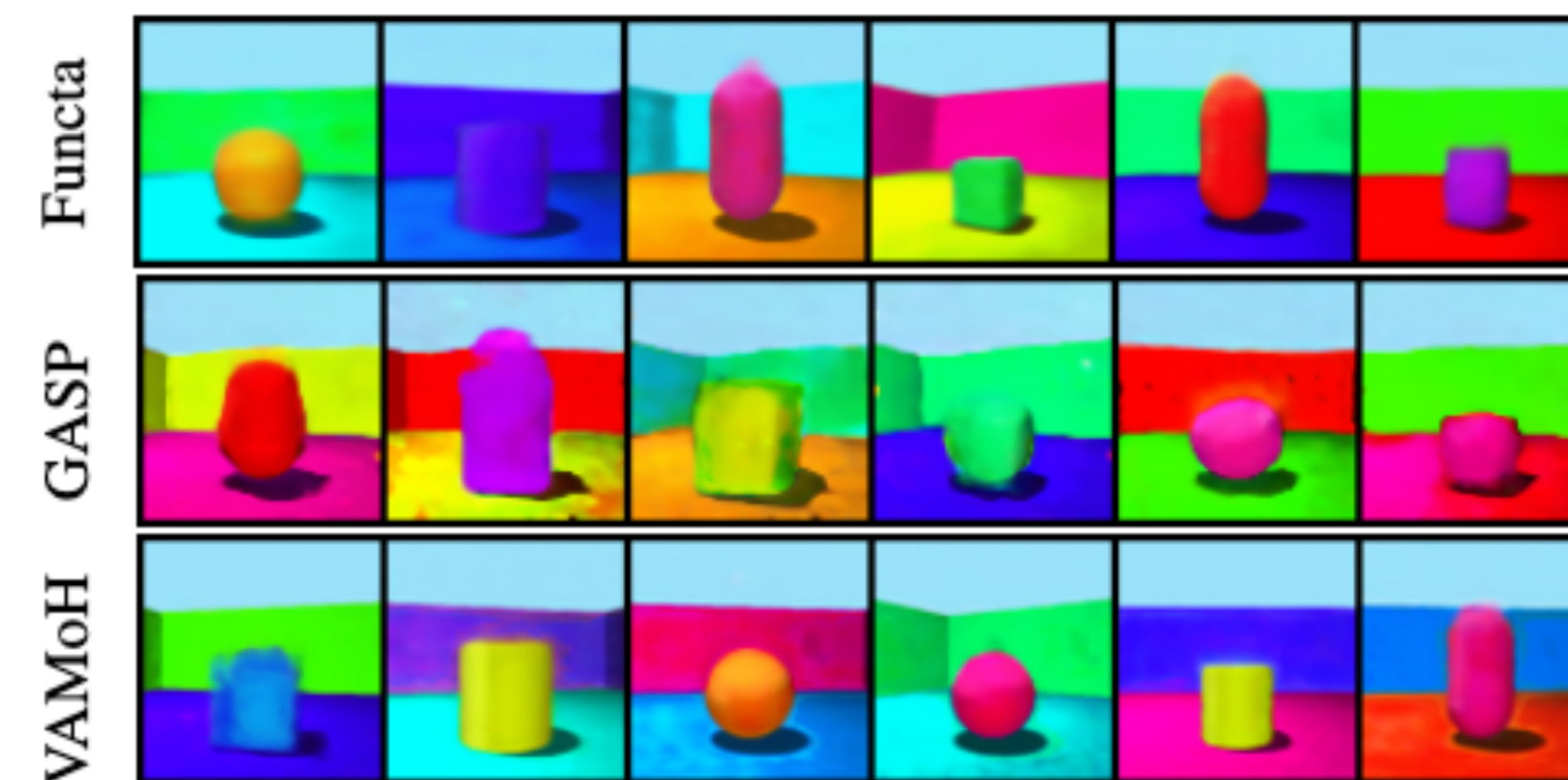


# Experiments

## Generation

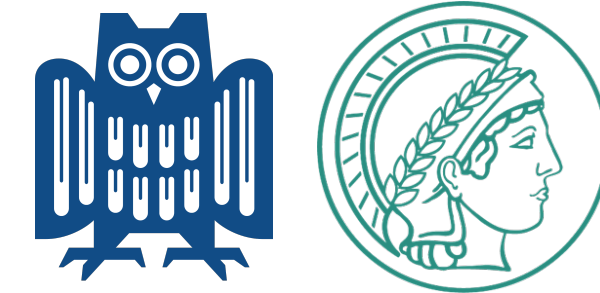


*CelebA-HQ*



*Shapes3D*



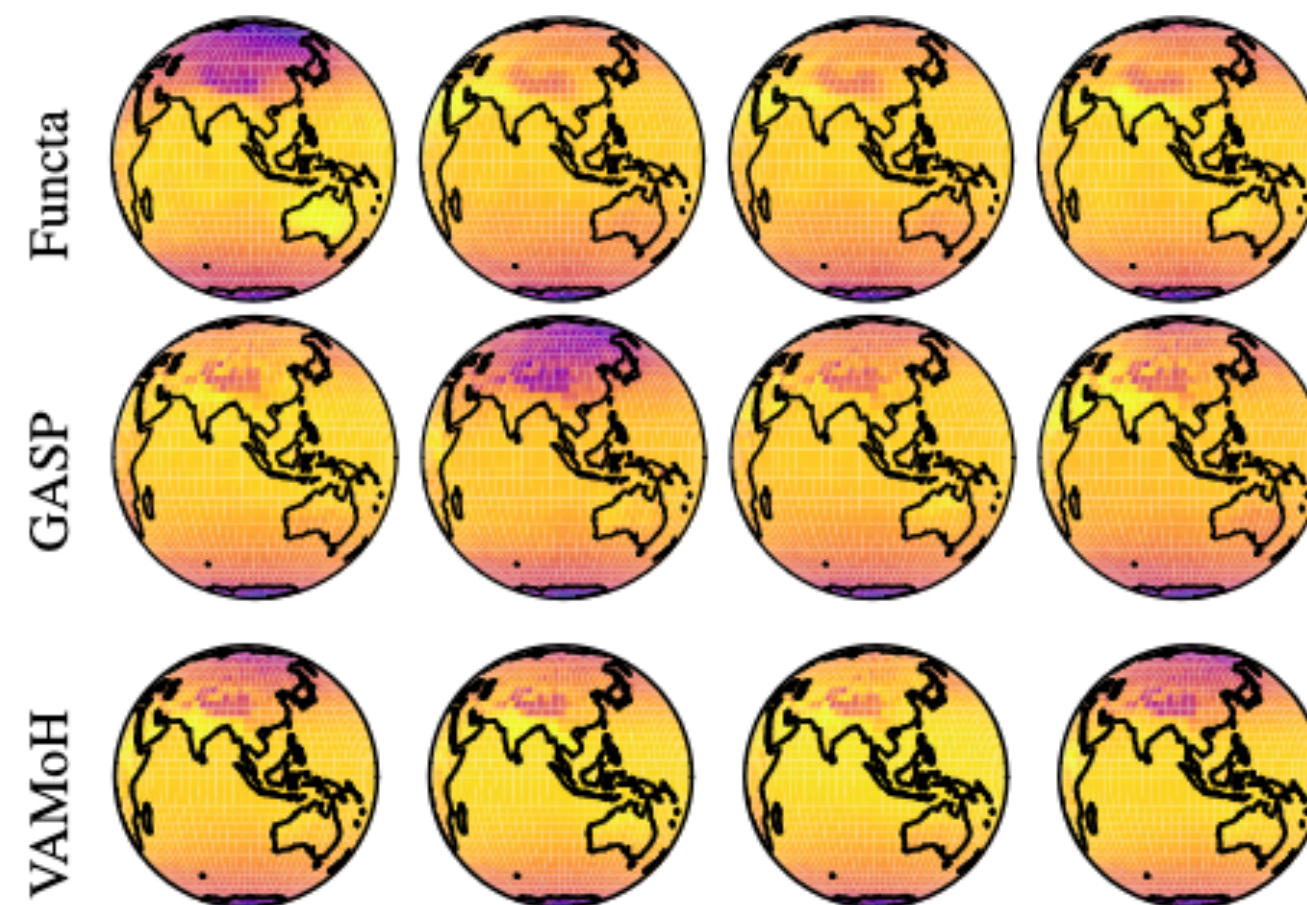


# Experiments

## Generation



*PolyMNIST*



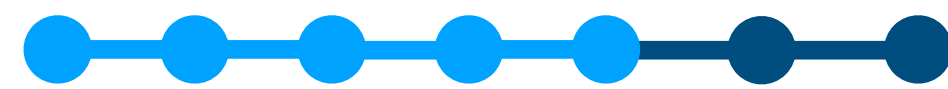
*ERA5*

VAMoH GASP



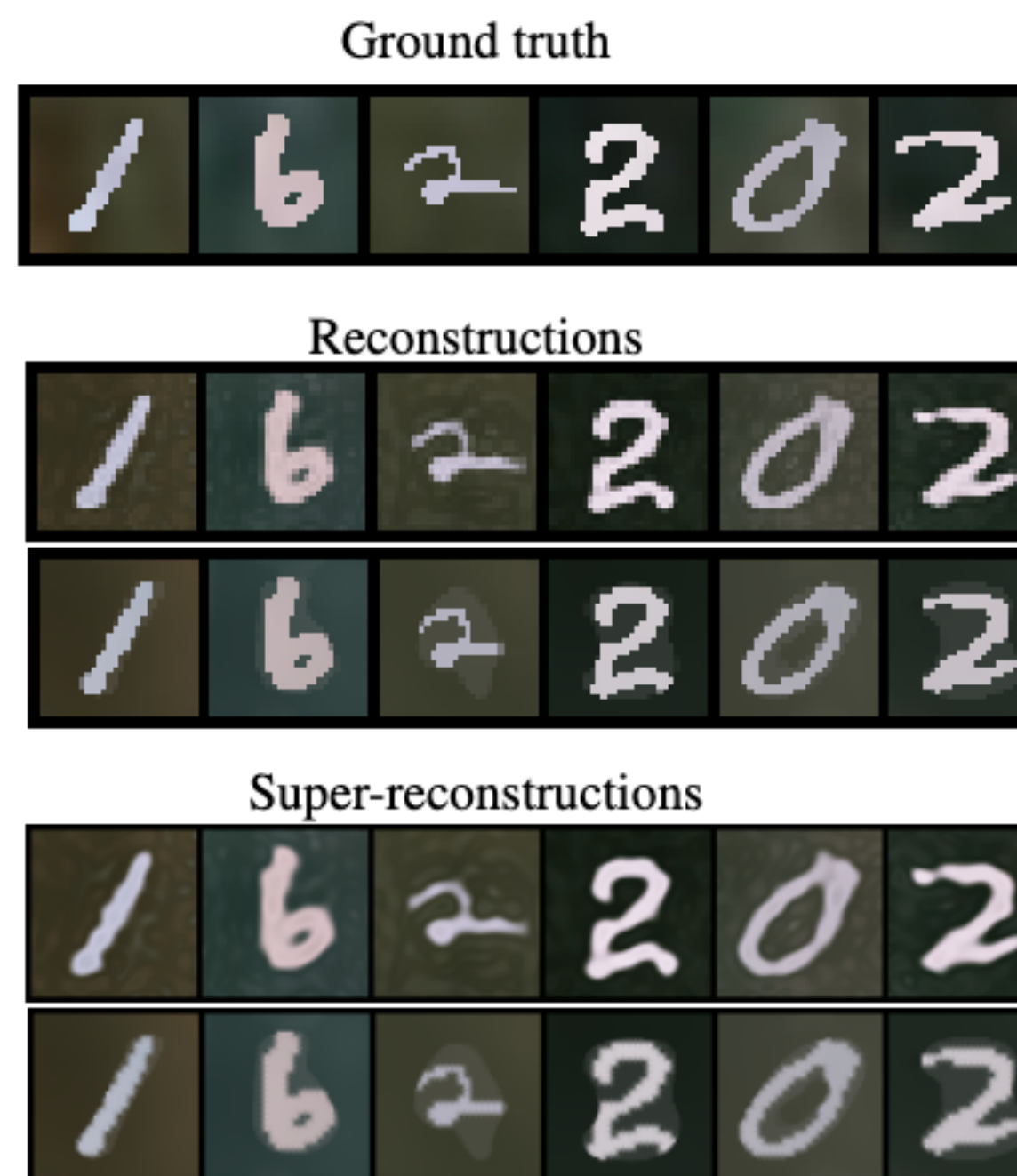
*ShapeNET*



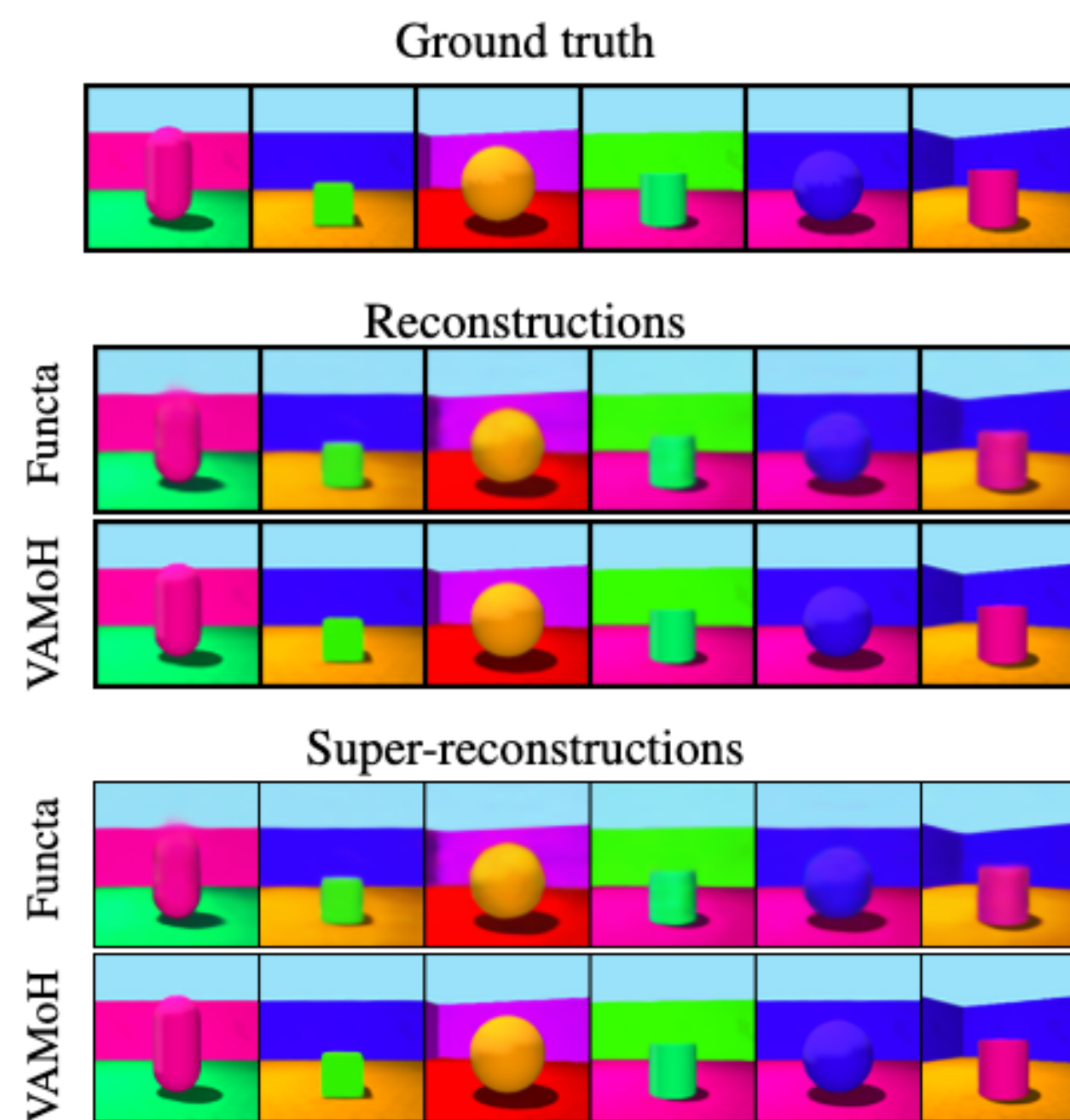


# Experiments

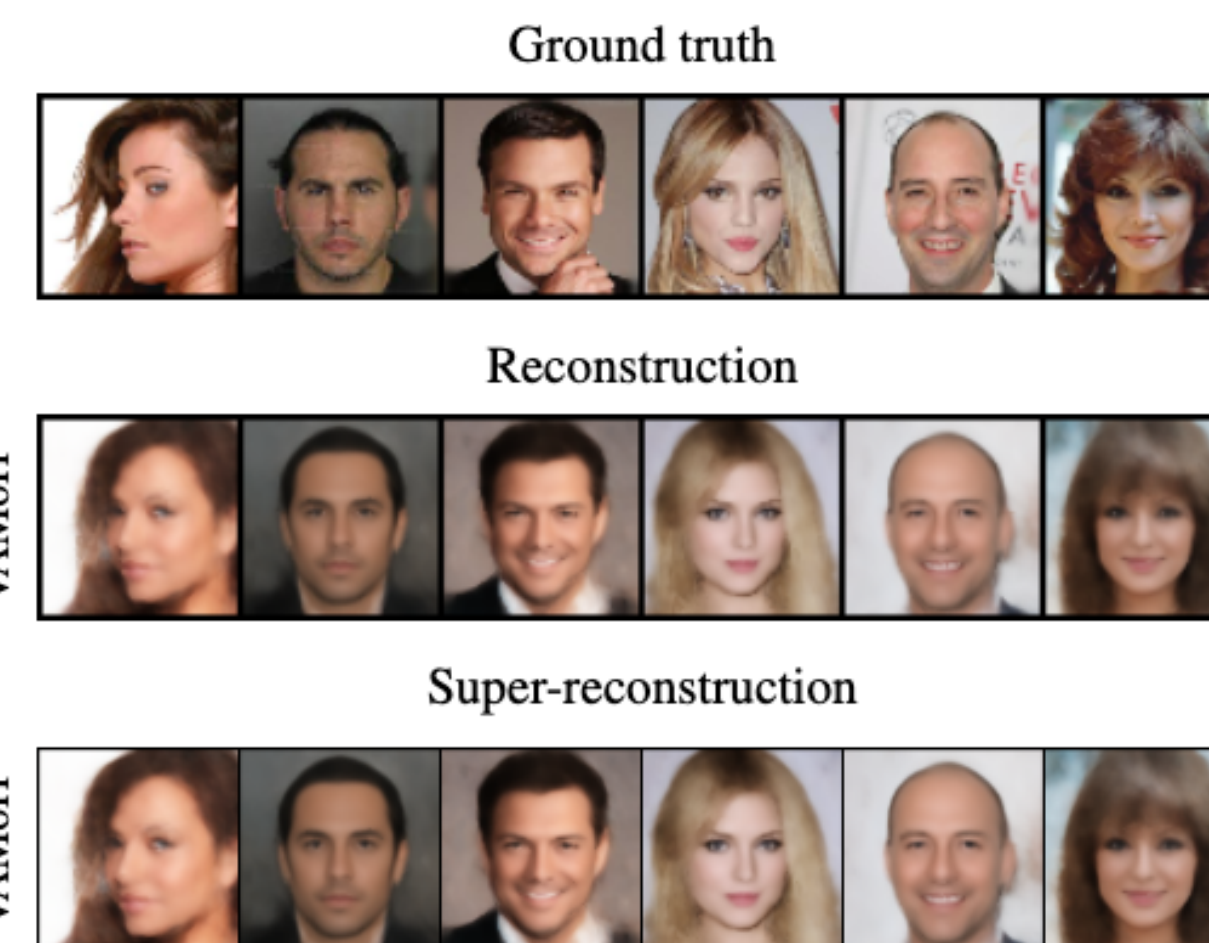
## Reconstructions



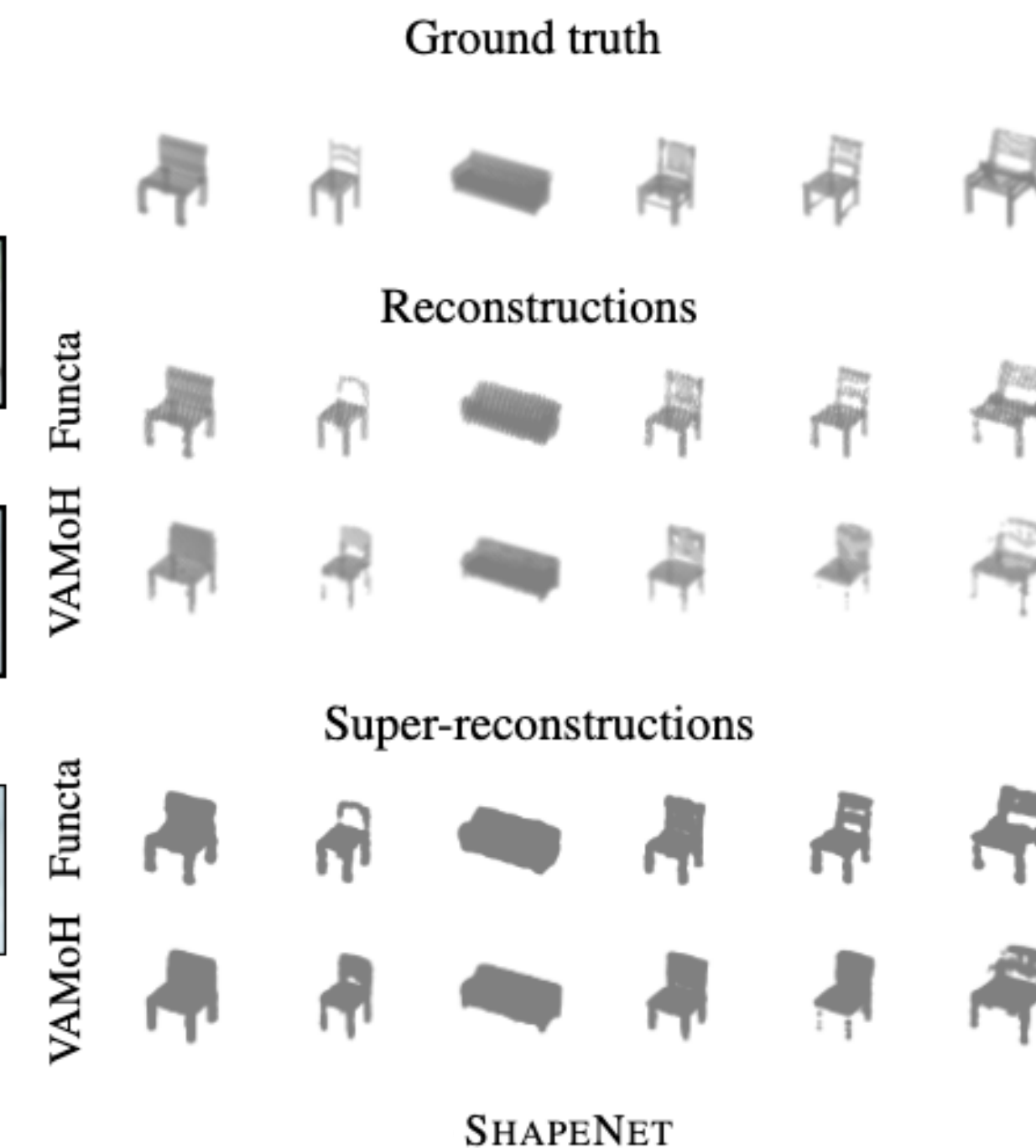
(c) POLYMNIST



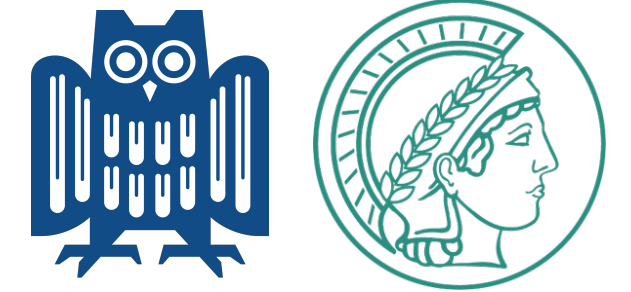
SHAPES3D



CELEBA HQ



SHAPENET



# Experiments

## Inference times

Table 2: Comparison of inference time (seconds) for reconstruction task of VaMoH and Functa. On the right-most two columns, we show the speed improvement of VaMoH compared to Functa (3) which is trained with 3 gradient steps as suggested in the original paper [Dupont et al., 2022b] and Functa (10) which is trained with 10 gradient step to obtain the results of Functa depicted in Figures 16,17. Please note that these experiments are run on the same GPU device.

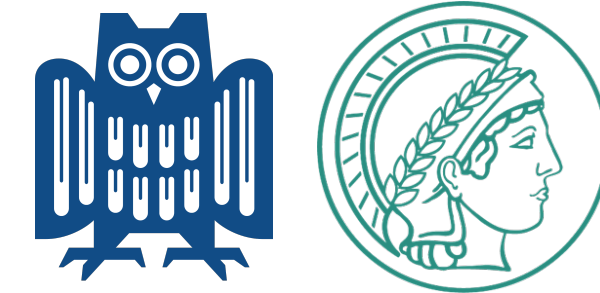
Dataset	Model Inference Time (secs)			Speed Improvement	
	VaMoH	Functa (3)	Functa (10)	vs. Functa (3)	vs. Functa (10)
POLYMNIST	<b>0.00453</b>	0.01648	0.05108	<b>x 3.64</b>	<b>x 11.28</b>
SHAPES3D	<b>0.00536</b>	0.01759	0.05480	<b>x 3.28</b>	<b>x 10.22</b>
CELEBA HQ	<b>0.00757</b>	0.01733	0.05381	<b>x 2.29</b>	<b>x 7.11</b>
ERA5	<b>0.00745</b>	0.01899	0.05932	<b>x 2.55</b>	<b>x 7.96</b>
SHAPENET	<b>0.00689</b>	0.02095	0.06576	<b>x 3.04</b>	<b>x 9.54</b>

*Reconstruction*

Dataset	Model Inference Time (secs)			Speed Improvement	
	VaMoH	Functa (3)	Functa (10)	vs. Functa (3)	vs. Functa (10)
POLYMNIST	<b>0.00455</b>	0.01649	0.05109	<b>x 3.62</b>	<b>x 11.23</b>
SHAPES3D	<b>0.00544</b>	0.01768	0.05489	<b>x 3.25</b>	<b>x 10.09</b>
CELEBA HQ	<b>0.00833</b>	0.01729	0.05377	<b>x 2.08</b>	<b>x 6.46</b>
ERA5	<b>0.00790</b>	0.01997	0.06030	<b>x 2.53</b>	<b>x 7.63</b>
SHAPENET	<b>0.01440</b>	0.02089	0.06569	<b>x 1.45</b>	<b>x 4.56</b>

*Super-reconstruction*



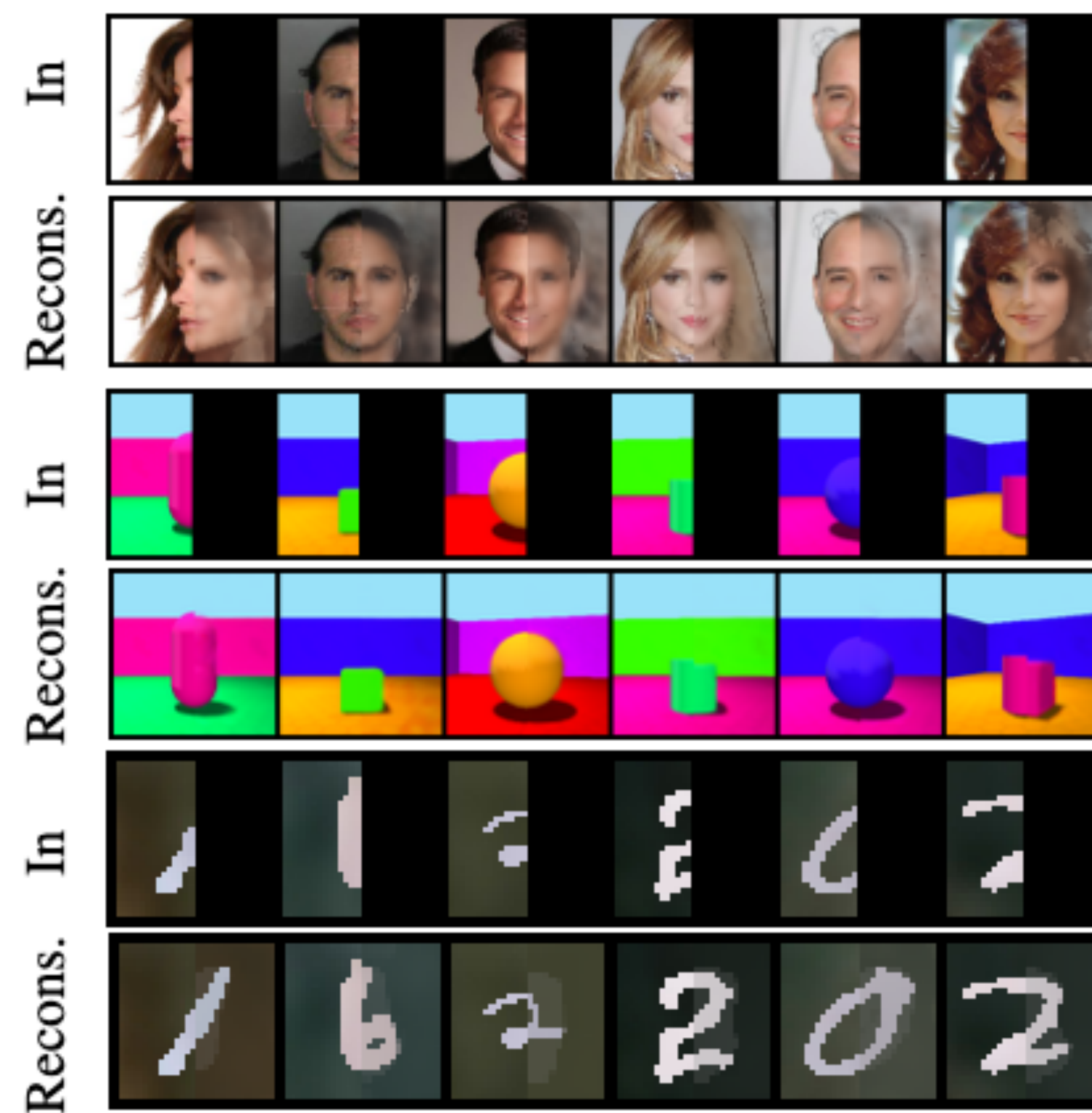


# Experiments

## Image completion



Missing a patch (in-painting)



Missing half of the image

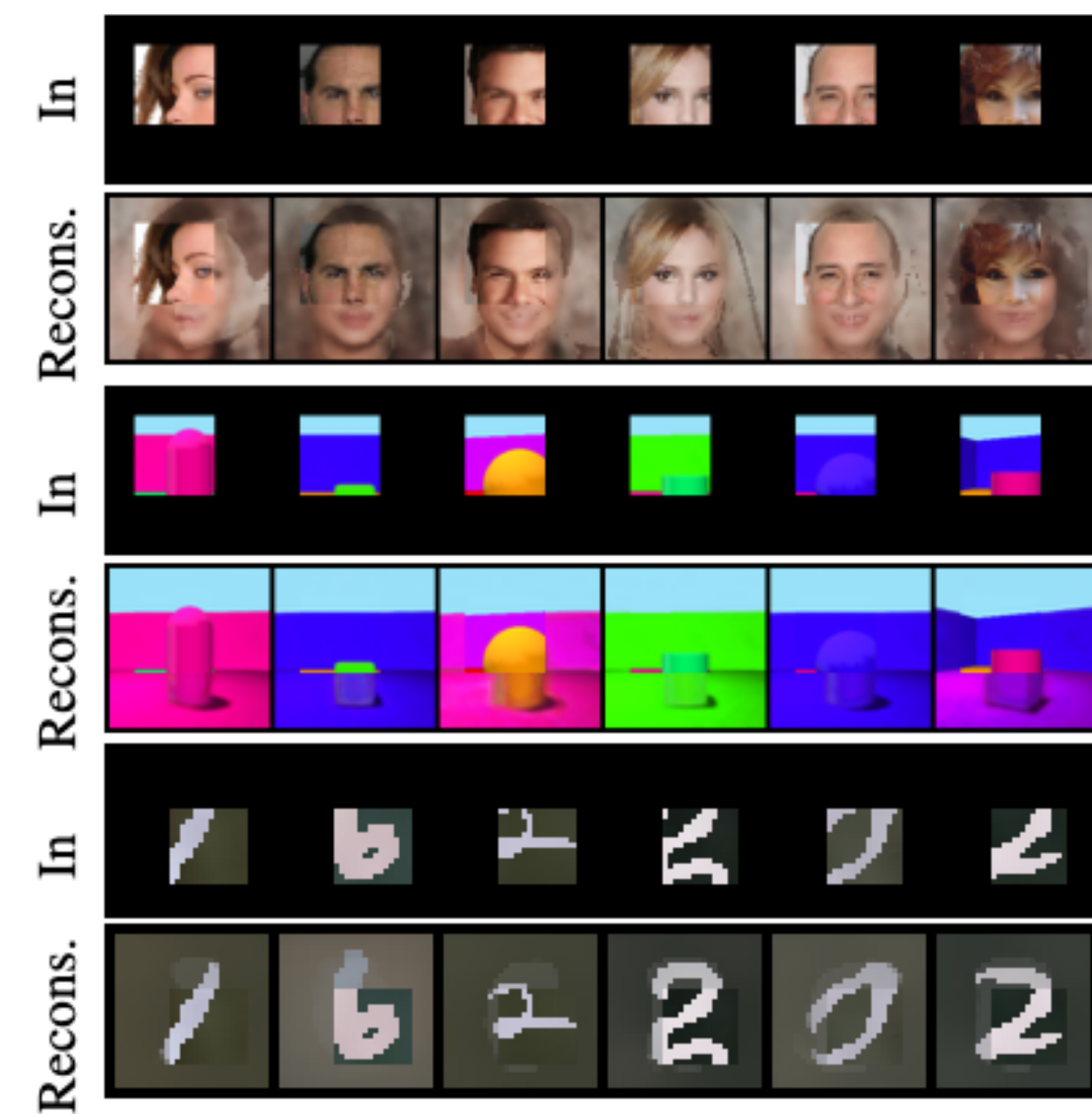


Image out-painting



# Conclusion



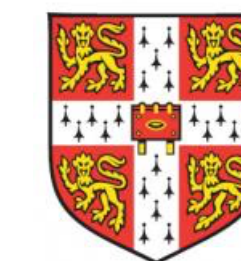
Thanks to learning distributions of functions, our proposed VAMoH can easily perform:

- Generation.
- Reconstruction.
- Conditional generation.
- Super resolution (interpolation).

While being:

- ✓ Robust to partially observed data.
- ✓ Expressive for generating high-quality data.
- ✓ Efficient in terms of inference.

# Further details



---

## VARIATIONAL MIXTURE OF HYPERGENERATORS FOR LEARNING DISTRIBUTIONS OVER FUNCTIONS

---

**Batuhan Koyuncu\***  
Saarland University  
Saarbrücken, Germany

**Pablo Sánchez-Martín**  
Max Planck Institute for Intelligent Systems  
Tübingen, Germany

**Ignacio Peis**  
Universidad Carlos III de Madrid  
Madrid, Spain

**Pablo M. Olmos**  
Universidad Carlos III de Madrid  
Madrid, Spain

**Isabel Valera**  
Saarland University  
Saarbrücken, Germany



[Paper]

[25] Koyuncu et al., 2023

# References

- [1] Campbell, A., Chen, W., Stimper, V., Hernandez-Lobato, J. M., & Zhang, Y. (2021, July). A gradient based strategy for hamiltonian monte carlo hyperparameter optimization. In *International Conference on Machine Learning* (pp. 1238-1248). PMLR.
- [2] Caterini, A. L., Doucet, A., & Sejdinovic, D. (2018). Hamiltonian variational auto-encoder. *Advances in Neural Information Processing Systems*, 31.
- [3] Salimans, T., Kingma, D., & Welling, M. (2015, June). Markov chain monte carlo and variational inference: Bridging the gap. In *International conference on machine learning* (pp. 1218-1226). PMLR.
- [4] Ruiz, F. J., Titsias, M. K., Cemgil, T., & Doucet, A. (2021, December). Unbiased gradient estimation for variational auto-encoders using coupled Markov chains. In *Uncertainty in Artificial Intelligence* (pp. 707-717). PMLR.
- [5] Dupont, E., Whye Teh, Y. & Doucet, A.. (2022). Generative Models as Distributions of Functions. *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics, in Proceedings of Machine Learning Research* 151:2989-3015.
- [6] Dupont, E., Kim, H., Eslami, S. A., Rezende, D. J., & Rosenbaum, D. (2022, June). From data to functa: Your data point is a function and you can treat it like one. In *International Conference on Machine Learning* (pp. 5694-5725). PMLR.

# References



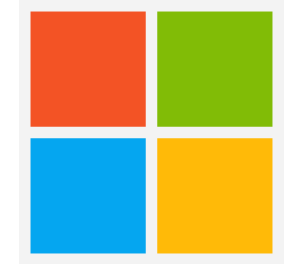
- [7] Kingma, D. P., & Welling, M. (2013). Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.
- [8] Cremer, C., Li, X., & Duvenaud, D. (2018, July). Inference suboptimality in variational autoencoders. In *International Conference on Machine Learning* (pp. 1078-1086). PMLR.
- [9] Bernardo, J. M. (1979). Expected information as expected utility. *the Annals of Statistics*, 686-690.
- [10] Ma, C., Tschitschek, S., Palla, K., Hernández-Lobato, J. M., Nowozin, S., & Zhang, C. (2018). Eddi: Efficient dynamic discovery of high-value information with partial vae. *arXiv preprint arXiv:1809.11142*.
- [11] Ma, C., Tschitschek, S., Turner, R., Hernández-Lobato, J. M., & Zhang, C. (2020). VAEM: a deep generative model for heterogeneous mixed type data. *Advances in Neural Information Processing Systems*, 33, 11237-11247.
- [12] Child, R. (2020). Very deep vaes generalize autoregressive models and can outperform them on images. *arXiv preprint arXiv:2011.10650*.



# References

- 📄 [13] Nazabal, A., Olmos, P. M., Ghahramani, Z., & Valera, I. (2020). Handling incomplete heterogeneous data using vaes. *Pattern Recognition*, 107, 107501.
- 📄 [14] Mattei, P. A., & Frelsen, J. (2019, May). MIWAE: Deep generative modelling and imputation of incomplete data sets. In *International conference on machine learning* (pp. 4413-4423). PMLR.
- 📄 [15] Peis, I., Ma, C., & Hernández-Lobato, J. M. (2022). Missing Data Imputation and Acquisition with Deep Hierarchical Models and Hamiltonian Monte Carlo. arXiv preprint arXiv:2202.04599.
- 📄 [16] Kraskov, A., Stögbauer, H., & Grassberger, P. (2004). Estimating mutual information. *Physical review E*, 69(6), 066138.
- 📄 [17] Gong, W., Li, Y., & Hernández-Lobato, J. M. (2020). Sliced kernelized Stein discrepancy. *arXiv preprint arXiv:2006.16531*.
- 📄 [18] Betancourt, M. (2017). A conceptual introduction to Hamiltonian Monte Carlo. *arXiv preprint arXiv:1701.02434*.

# References



- [19] Betancourt, M., & Girolami, M. (2015). Hamiltonian Monte Carlo for hierarchical models. *Current trends in Bayesian methodology with applications*, 79(30), 2-4.
- [20] Sitzmann, V., Martel, J., Bergman, A., Lindell, D., & Wetzstein, G. (2020). Implicit neural representations with periodic activation functions. *Advances in Neural Information Processing Systems*, 33, 7462-7473.
- [21] Mescheder, L., Oechsle, M., Niemeyer, M., Nowozin, S., & Geiger, A. (2019). Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 4460-4470).
- [22] Sitzmann, V., Zollhöfer, M., & Wetzstein, G. (2019). Scene representation networks: Continuous 3d-structure-aware neural scene representations. *Advances in Neural Information Processing Systems*, 32.
- [23] Ha, D., Dai, A. M., & Le, Q. V. HyperNetworks. In International Conference on Learning Representations.
- [24] Wu, W., Qi, Z., & Fuxin, L. (2019). Pointconv: Deep convolutional networks on 3d point clouds. In *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition* (pp. 9621-9630).
- [25] Koyuncu, B., Sanchez-Martin, P., Peis, I., Olmos, P. M., & Valera, I. (2023). Variational Mixture of HyperGenerators for Learning Distributions Over Functions. In Proceedings of the 40th International Conference on Machine Learning, 2023.

# Thank you!



[ipeis@tsc.uc3m.es](mailto:ipeis@tsc.uc3m.es)