



uc3m

Universidad  
**Carlos III**  
de Madrid



UNIVERSITY OF  
CAMBRIDGE

[07/01/22 GTS Research Seminar]

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

Ignacio Peis<sup>1</sup>, Chao Ma<sup>2,3</sup>, José Miguel Hernández-Lobato<sup>2</sup>

<sup>1</sup>Dept. of Signal Theory and Communications, Universidad Carlos III de Madrid

<sup>2</sup>Dept. of Engineering, University of Cambridge

<sup>3</sup>Microsoft Research Cambridge



# Introduction

## Challenges

- Improve approximate inference in advanced VAEs
- Improve missing data imputation
- Improve predictions under missing data condition
- Improve active information acquisition
- Deal with partial, mixed-type data

# Introduction

## Contributions

- Improved Hierarchical VAE for mixed-type partial data.
- Improved inference via Hamiltonian Monte Carlo with automatic hyperparameter optimization.
- Improved active learning with novel sampling-based active information acquisition technique.

---

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

---

**Ignacio Peis<sup>1</sup>   Chao Ma<sup>2,3</sup>   José Miguel Hernández-Lobato<sup>2</sup>**

<sup>1</sup>Universidad Carlos III de Madrid   <sup>2</sup>University of Cambridge   <sup>3</sup>Microsoft Research Cambridge  
ipeis@tsc.uc3m.es  
{cm905, jmh233}@cam.ac.uk

[\[Preprint\]](#) [\[Code\]](#)

# Variational Autoencoders

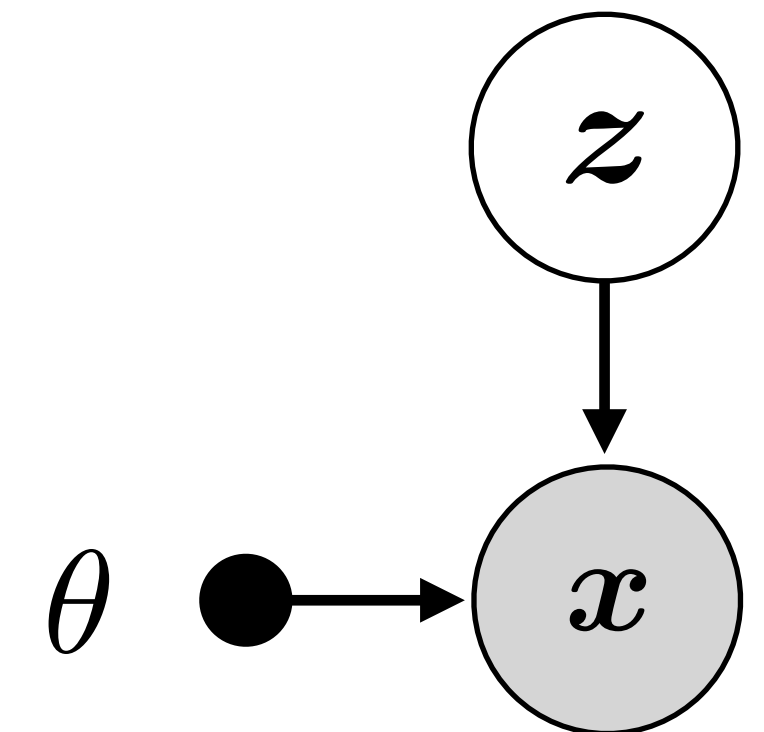
## Definition<sup>1</sup>

- Generative, explicit models with intractable probability.
- **Goal:** optimize parameters to maximize the marginal 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 NNs.
- Posterior distribution:

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



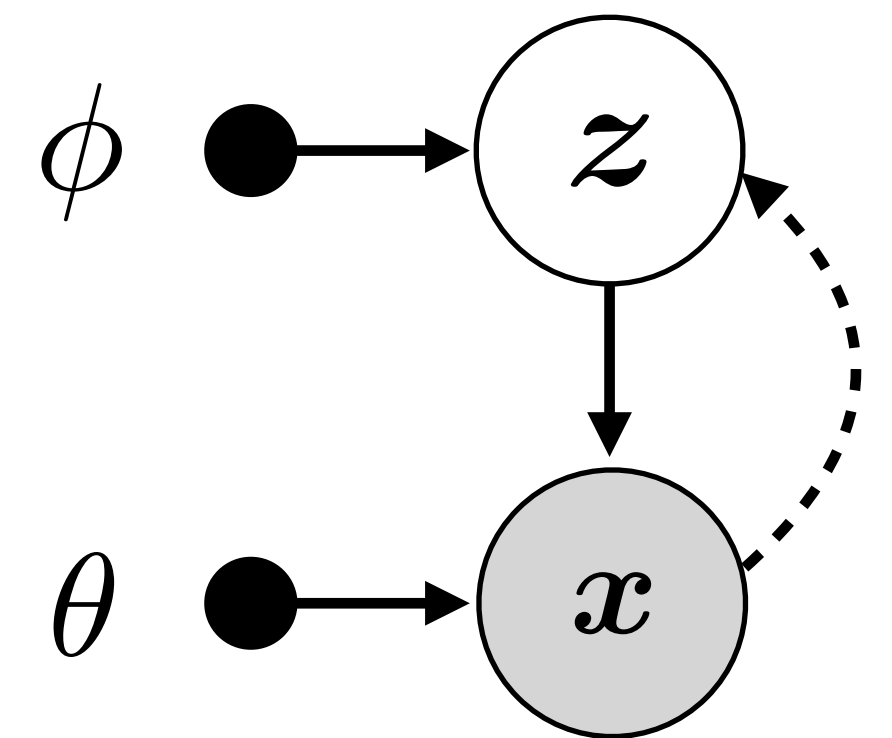
<sup>1</sup> Kigima et al., 2013



# Variational Autoencoders

## Training

- For each batch:
  1. Encode to parameters of the approximate posterior.
  2. Sample from  $q_\phi(\mathbf{z}|\mathbf{x})$ .
  3. Obtain the Monte Carlo approximation of the ELBO.
  4. Optimization step on  $\theta$  and  $\phi$ .



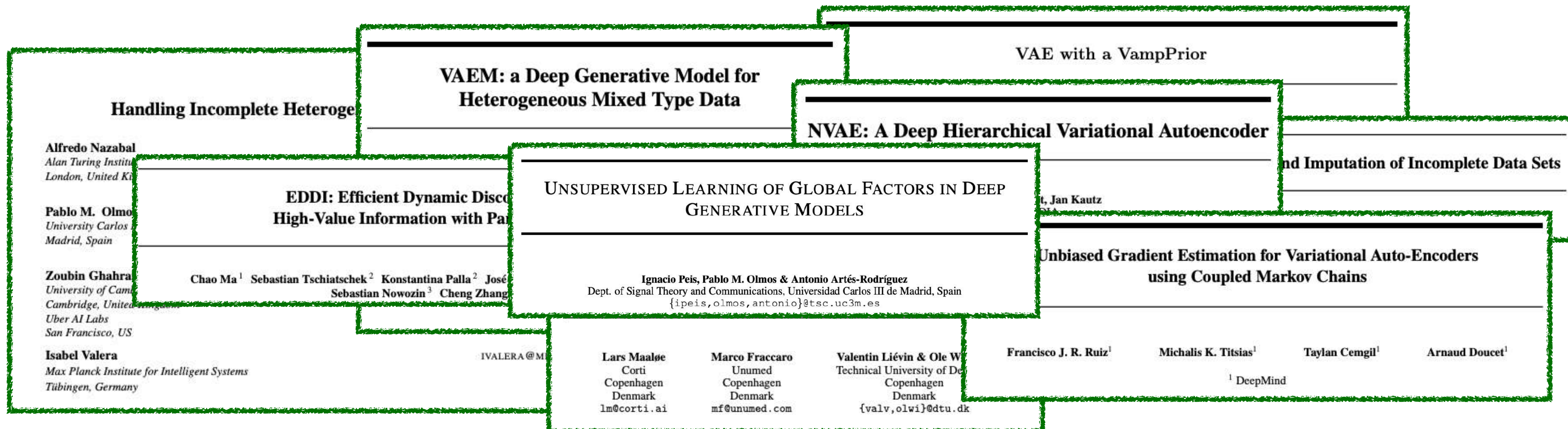
$$\nabla_{(\theta, \phi)} \left( \frac{1}{B} \sum_{i=1}^B [\log p_\theta(\mathbf{x}|\mathbf{z}) - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))] \right)$$



# Variational Autoencoders

## Recent advances

- Handling incomplete and heterogeneous data.
- Improving approximate inference VS increasing flexibility of the prior.



# Hierarchical Variational Autoencoders

## Definition\*

- Add flexibility to the prior with an autoregressive path of latent variables

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

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

- **Inductive bias:** hierarchical flow of information.
- **Potential pitfalls:** *posterior collapse*.

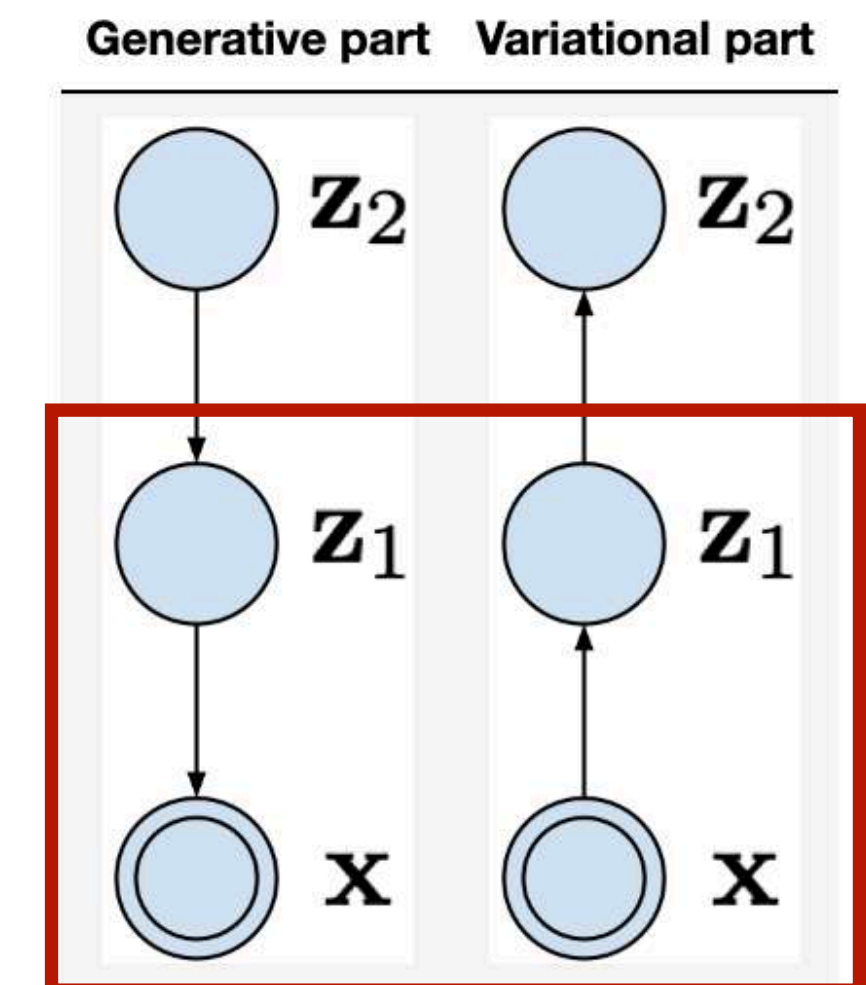


Figure 4. A two-level VAE.

\* [https://jmtomczak.github.io/blog/9/9\\_hierarchical\\_lvm\\_p1.html](https://jmtomczak.github.io/blog/9/9_hierarchical_lvm_p1.html)



# Hierarchical Variational Autoencoders

## Definition

- To avoid posterior collapse, the variational networks learn a residual difference of the posterior from a deterministic bottom-up path<sup>1,2,3</sup>.

$$q(\mathbf{z}_i | \mathbf{x}) = \mathcal{N}(\mu_i + \Delta\mu_i(\mathbf{x}), \sigma_i^2 \Delta\sigma_i^2)$$

- The generative and variational posterior are tightly connected:

$$KL(q(\mathbf{z}^i | x) \| p(\mathbf{z}^i)) = \frac{1}{2} \left( \frac{\Delta\mu_i^2}{\sigma_i^2} + \Delta\sigma_i^2 - \log \Delta\sigma_i^2 - 1 \right)$$

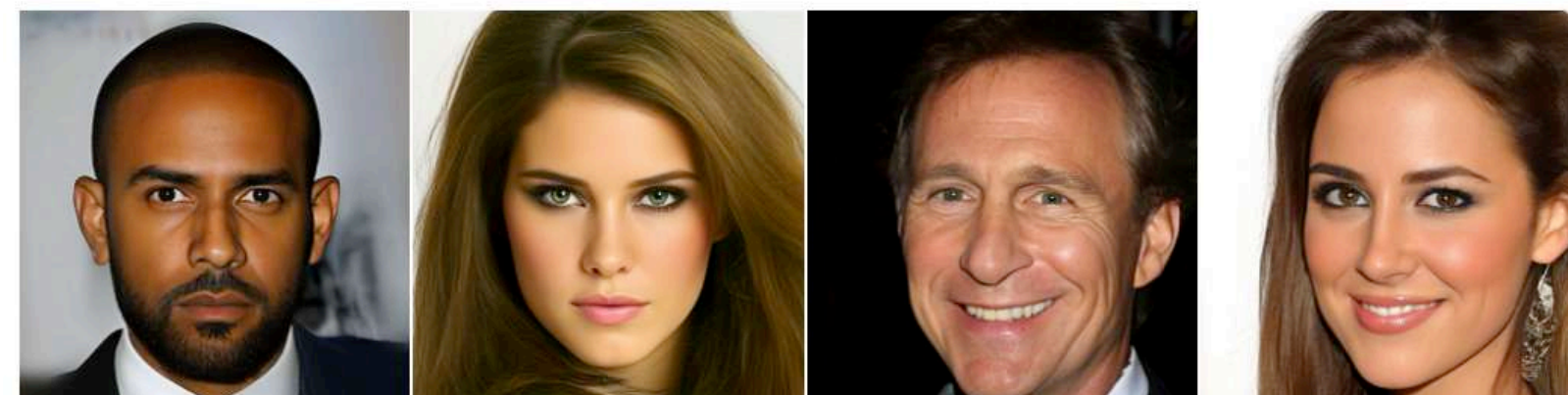


Figure 1: 256×256-pixel samples generated by NVAE, trained on CelebA HQ [28].

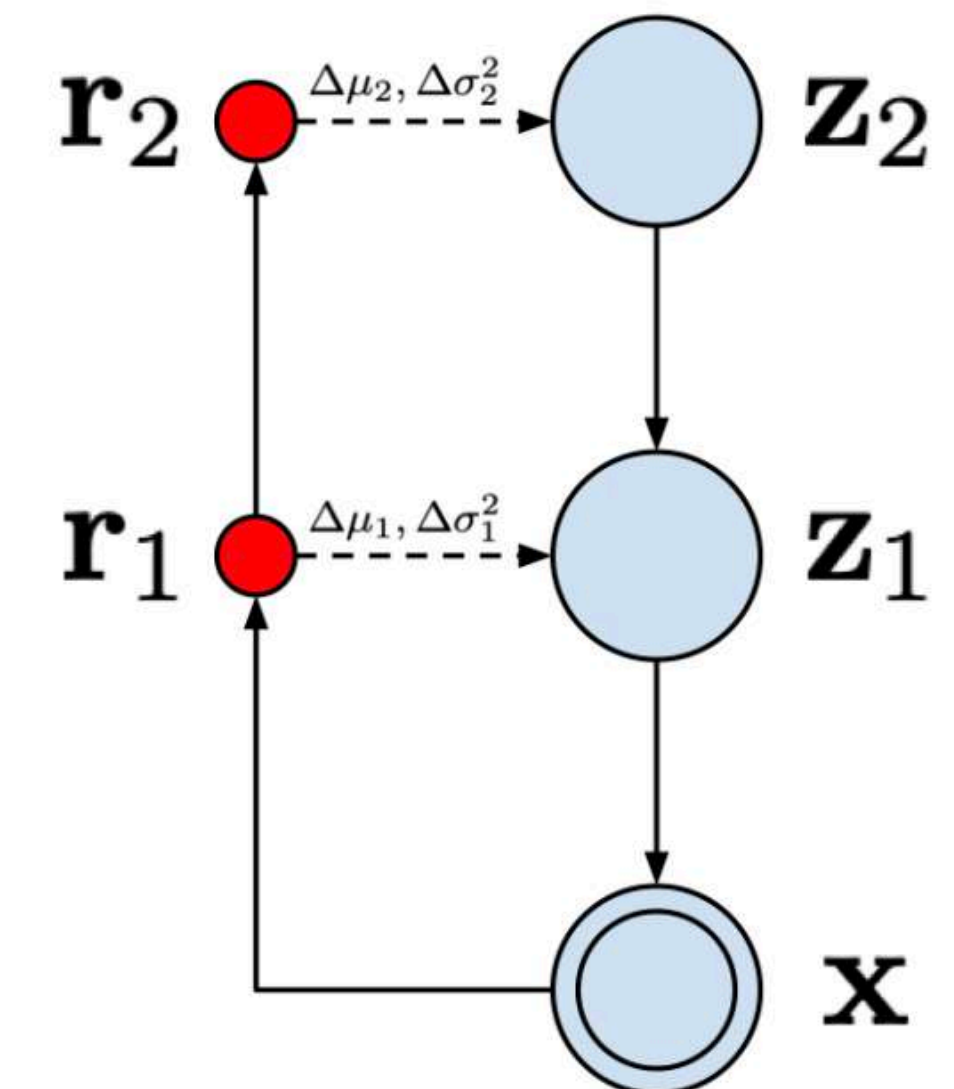


Figure 5. A top-down VAE.

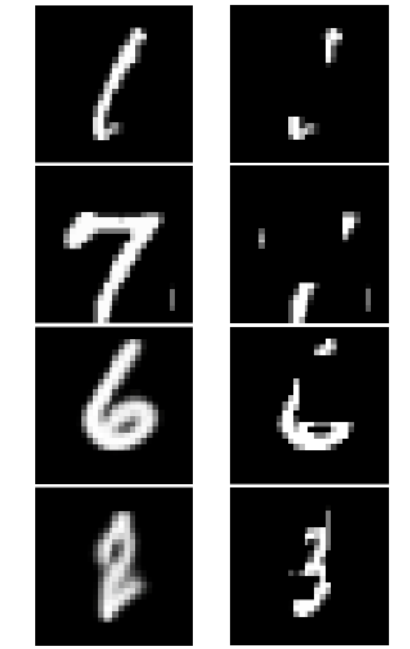
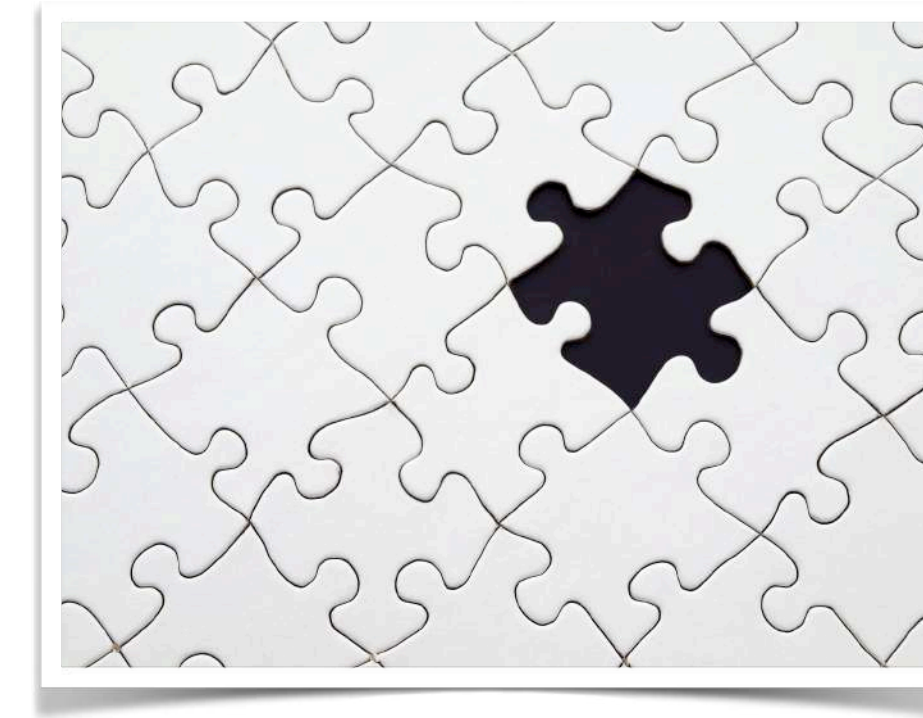
<sup>1</sup> Vahdat et al., 2020

<sup>2</sup> Maaløe et al., 2019

<sup>3</sup> Child, 2019

# Variational Autoencoders

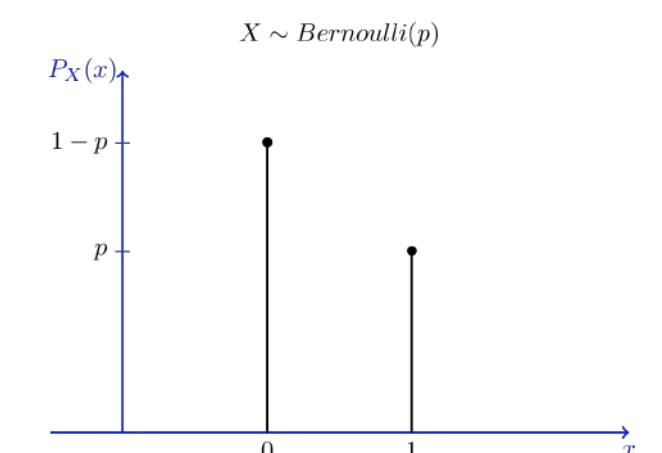
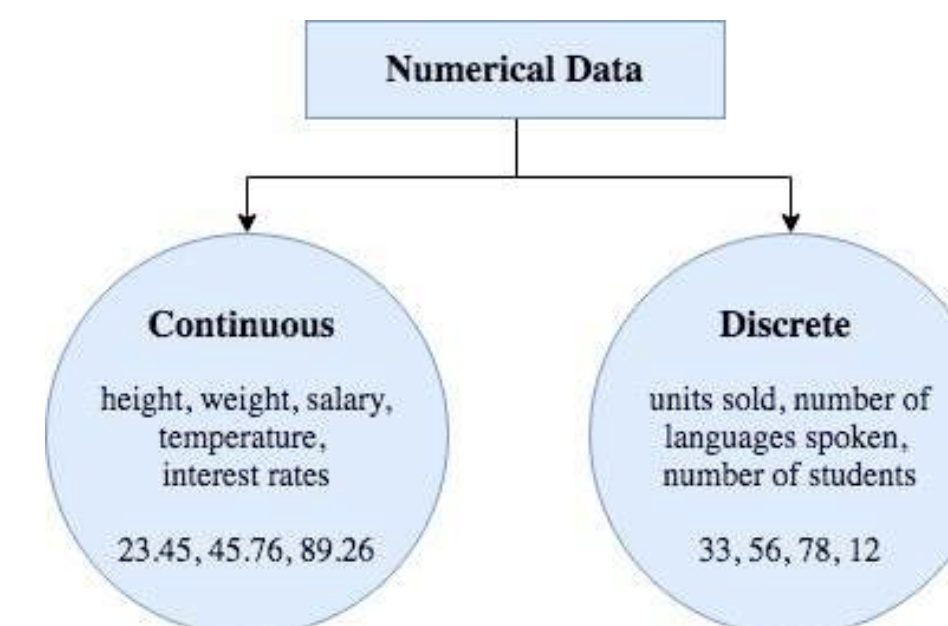
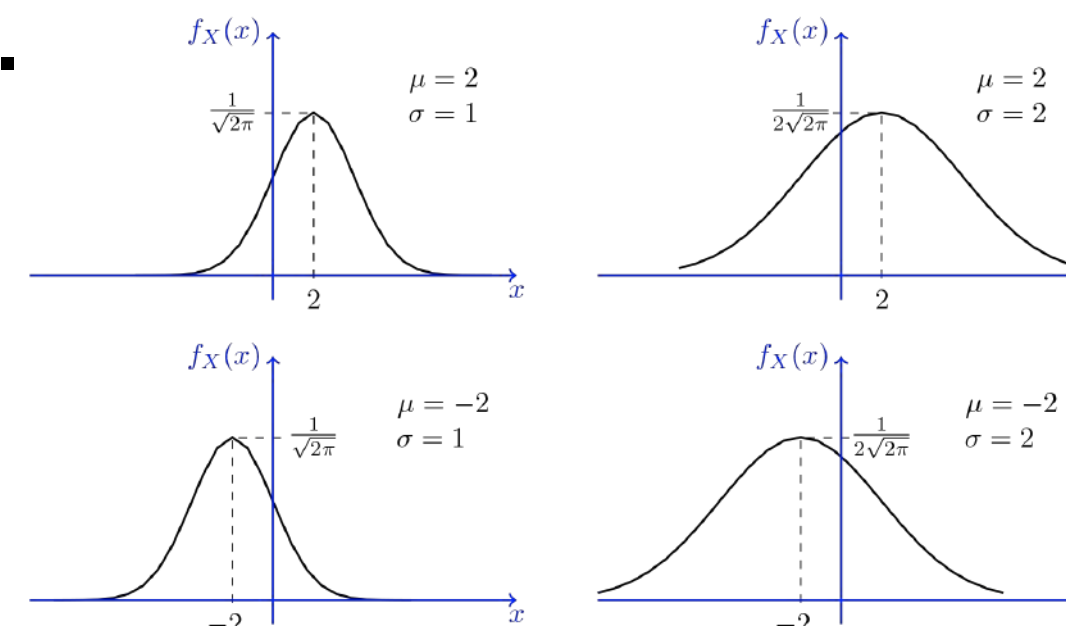
## VAEs for partial, heterogeneous data



- Factorization over dimensions<sup>1,2</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>1</sup>: use a different likelihood per dimension.
- Problem:** unbalanced likelihoods.



<sup>1</sup> Nazabal et al., 2020

<sup>2</sup> Mattei et al., 2019



# Variational Autoencoders

## VAEM for partial heterogeneous data

- **Solution<sup>1</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 and a joint dependency VAE  $(\theta, \phi)$  on the marginally encoded data:

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

- The marginal encodings  $z_d$  are Gaussian, thus, the likelihoods are balanced.
- Interdependencies between heterogeneous variables are better captured by the model.

<sup>1</sup> Ma et al., 2020

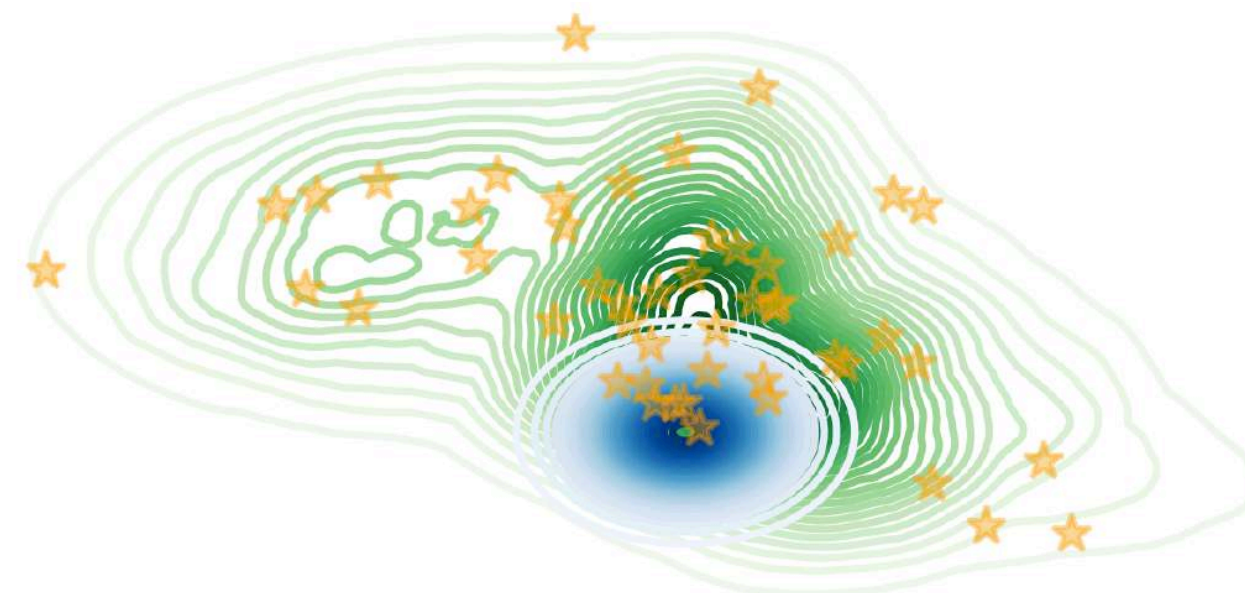
# Approximate Inference in VAEs

## Variational Inference

- VAEs are **restricted** by purely Gaussian approximations of the true posterior.



<sup>1</sup>True posterior (green) vs Approximate posterior (blue)



<sup>2</sup>Samples from the true posterior (orange)

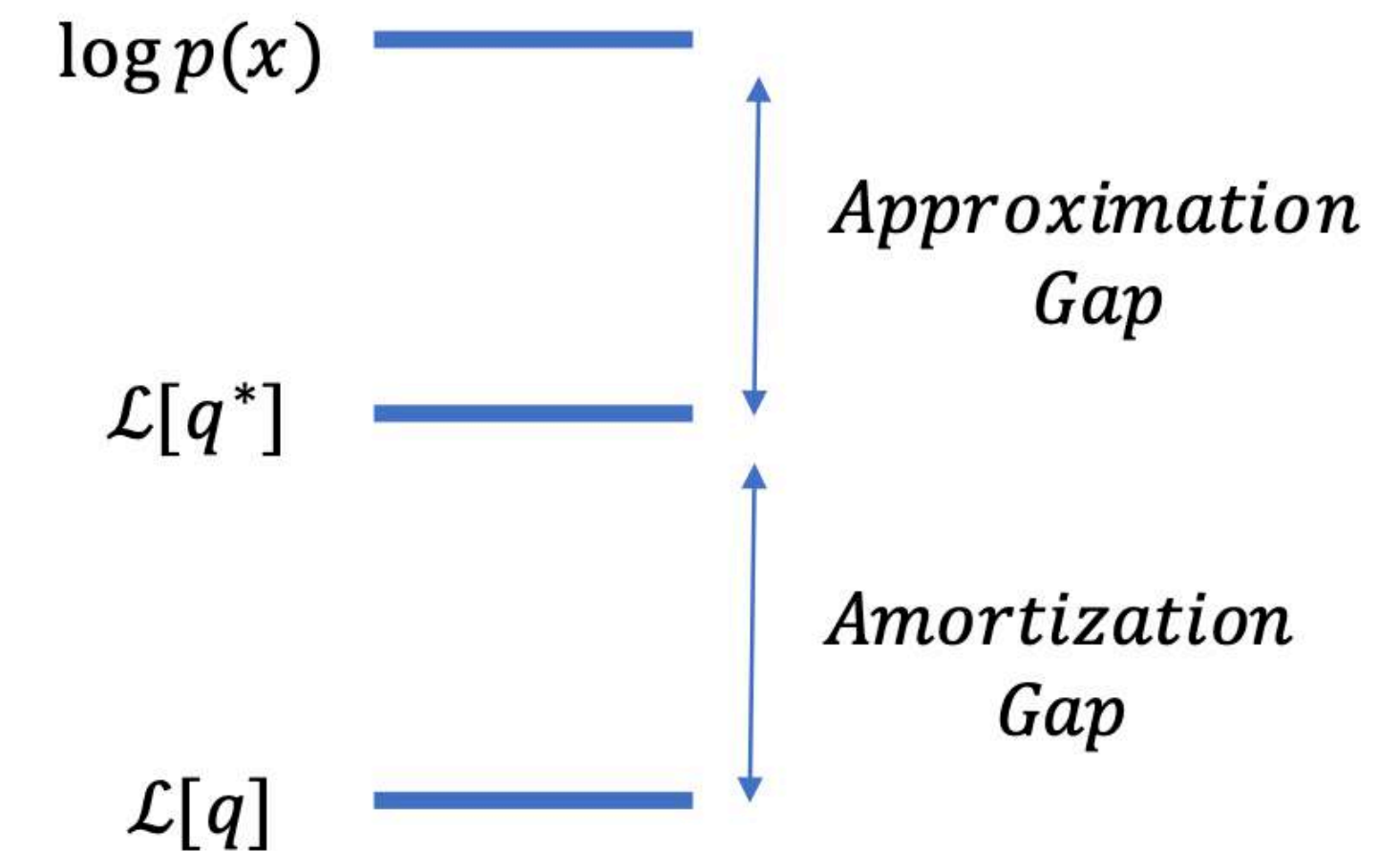


Figure 1. Gaps in Inference

<sup>1</sup> Cremer et al., 2018

<sup>2</sup> Peis et al., 2022



# Bayesian active information acquisition

## Encoder-based

- Reward function as an expected gain of information<sup>1</sup>:

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

- **Intractable**. Solved by transforming into  $\mathbf{z}$  space<sup>2</sup>:

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

- Extension for the VAEM model<sup>3</sup>:

$$\hat{R}_I(\mathbf{x}_i, \mathbf{x}_O) = \mathbb{E}_{p_\theta(\mathbf{x}_i, \mathbf{z}_i, \mathbf{z}_O | \mathbf{x}_O)} \{ \text{KL} [q_\lambda(\mathbf{h} | \mathbf{z}_i, \mathbf{z}_O) || q_\lambda(\mathbf{h} | \mathbf{z}_O)] - \\ \mathbb{E}_{p_\theta(\mathbf{x}_\phi, \mathbf{z}_\phi, | \mathbf{x}_O)} \text{KL} [q_\lambda(\mathbf{h} | \mathbf{z}_\phi, \mathbf{z}_i, \mathbf{z}_O) || q_\lambda(\mathbf{h} | \mathbf{z}_\phi, \mathbf{z}_O)] \}$$

- **Restriction**: this approximation metric still relies on the Gaussian approximation of the encoder.

<sup>1</sup> Bernardo et al., 1979

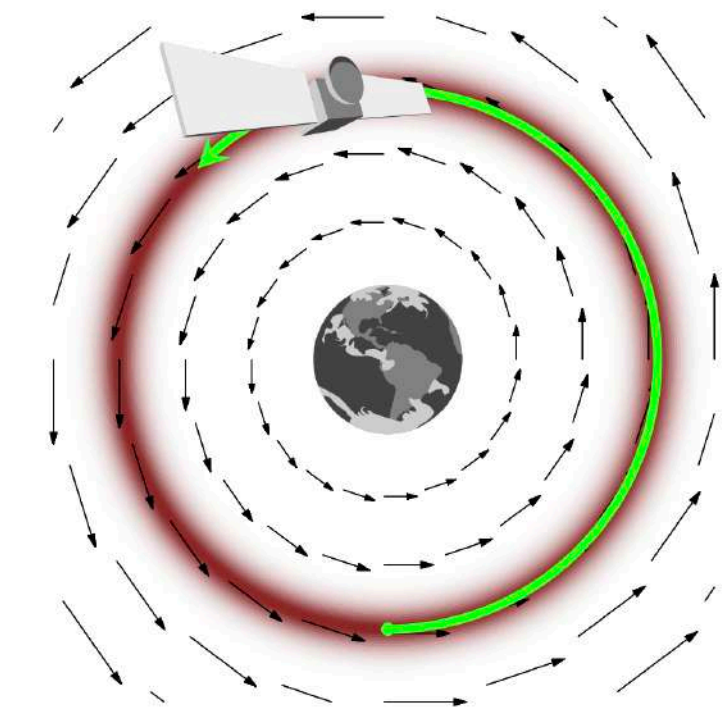
<sup>2</sup> Ma et al., 2018

<sup>2</sup> Ma et al., 2020

# Hamiltonian Monte Carlo<sup>1,2</sup>

## Definition

- Sample from complex distributions via unnormalised targets  $p(\mathbf{z}) = \frac{1}{Z} p^*(\mathbf{z})$
- Highly efficient exploration using differential geometry and conservative dynamics



1. Expand to phase space with momentum variable:

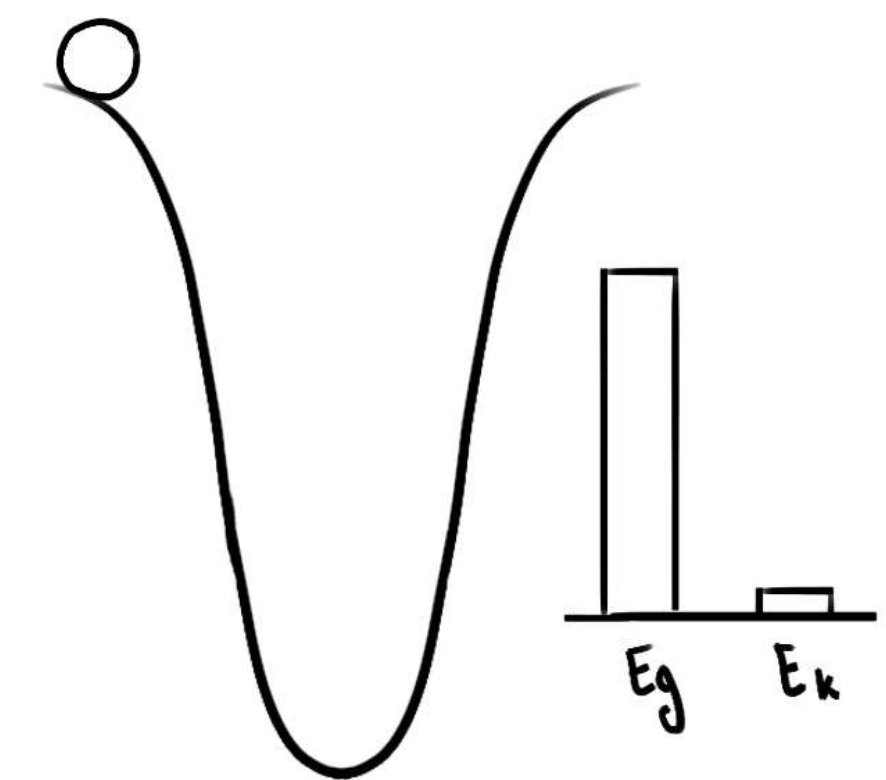
$$(\mathbf{z}) \rightarrow (\mathbf{z}, \mathbf{r}) \quad p(\mathbf{z}, \mathbf{r}) = p(\mathbf{r}|\mathbf{z})p(\mathbf{z})$$

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

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

2. Hamiltonian equations

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



<sup>1</sup> Betancourt, 2017

<sup>2</sup> Neal., 2017



# Hamiltonian Monte Carlo<sup>1</sup>

## HMC in practice

- Leapfrog integrator for Hamiltonian equations.
- Discrete trajectories (*chains*) of  $T$  updates.
- From an initial proposal, each update consists on  $L$  cyclic Leapfrog steps:

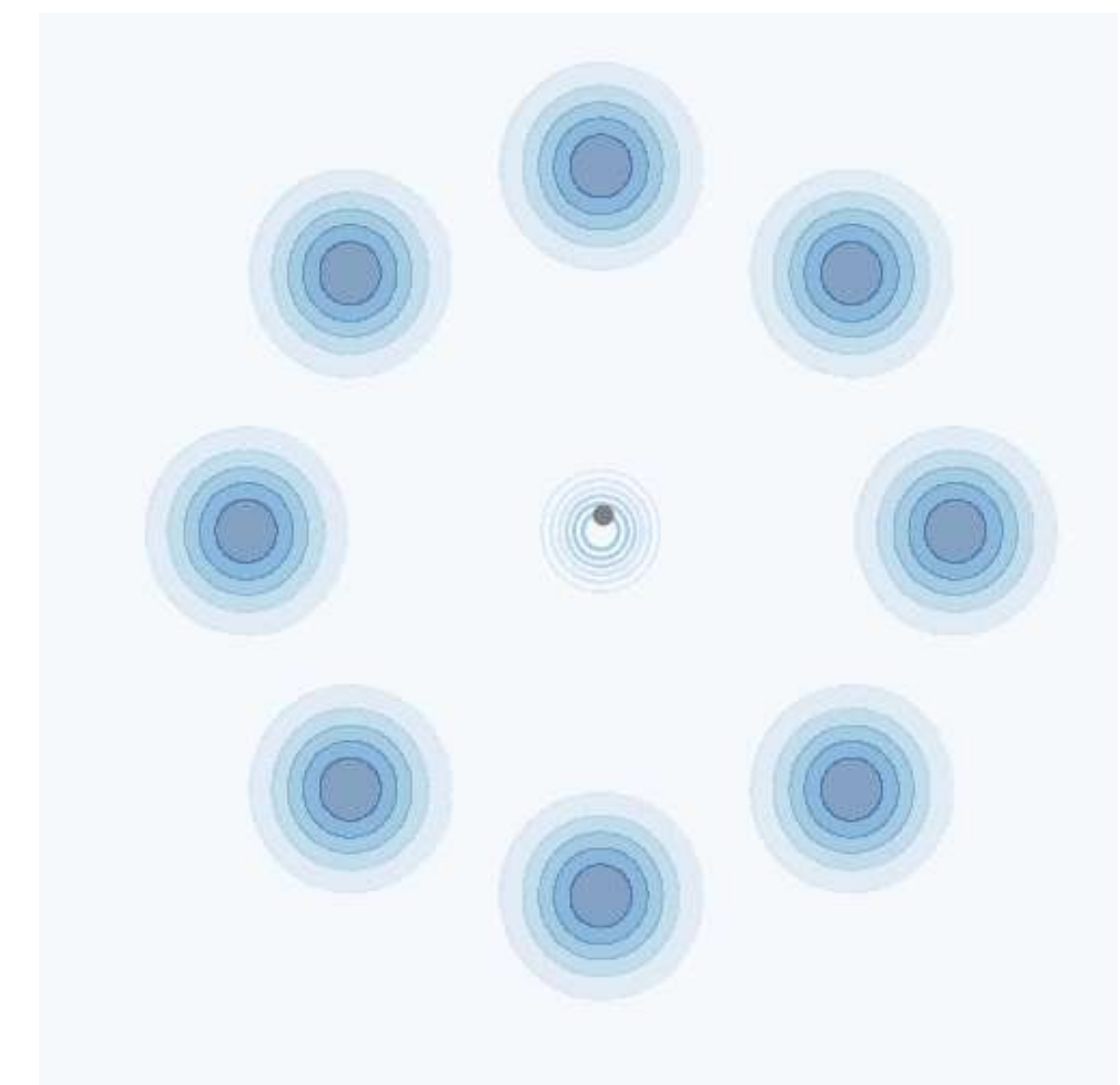
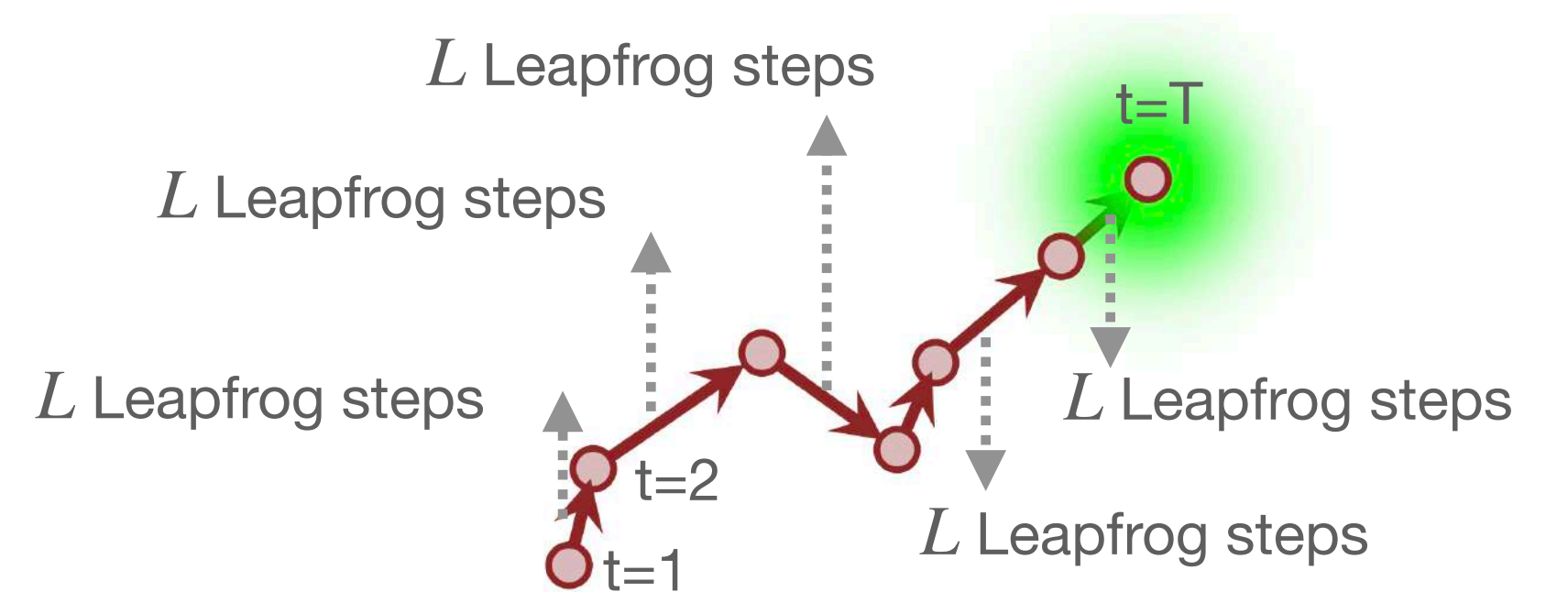
$$\mathbf{r}_{l+\frac{1}{2}} = \mathbf{r}_l + \frac{1}{2} \boldsymbol{\phi} \odot \nabla_{\mathbf{z}_l} \log p^*(\mathbf{z}_l),$$

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

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

- Ending in a new proposal  $(\mathbf{z}', \mathbf{r}')$ , which is accepted with probability:

$$\min [1, \exp(-H(\mathbf{z}', \mathbf{r}') + H(\mathbf{z}, \mathbf{r}))]$$

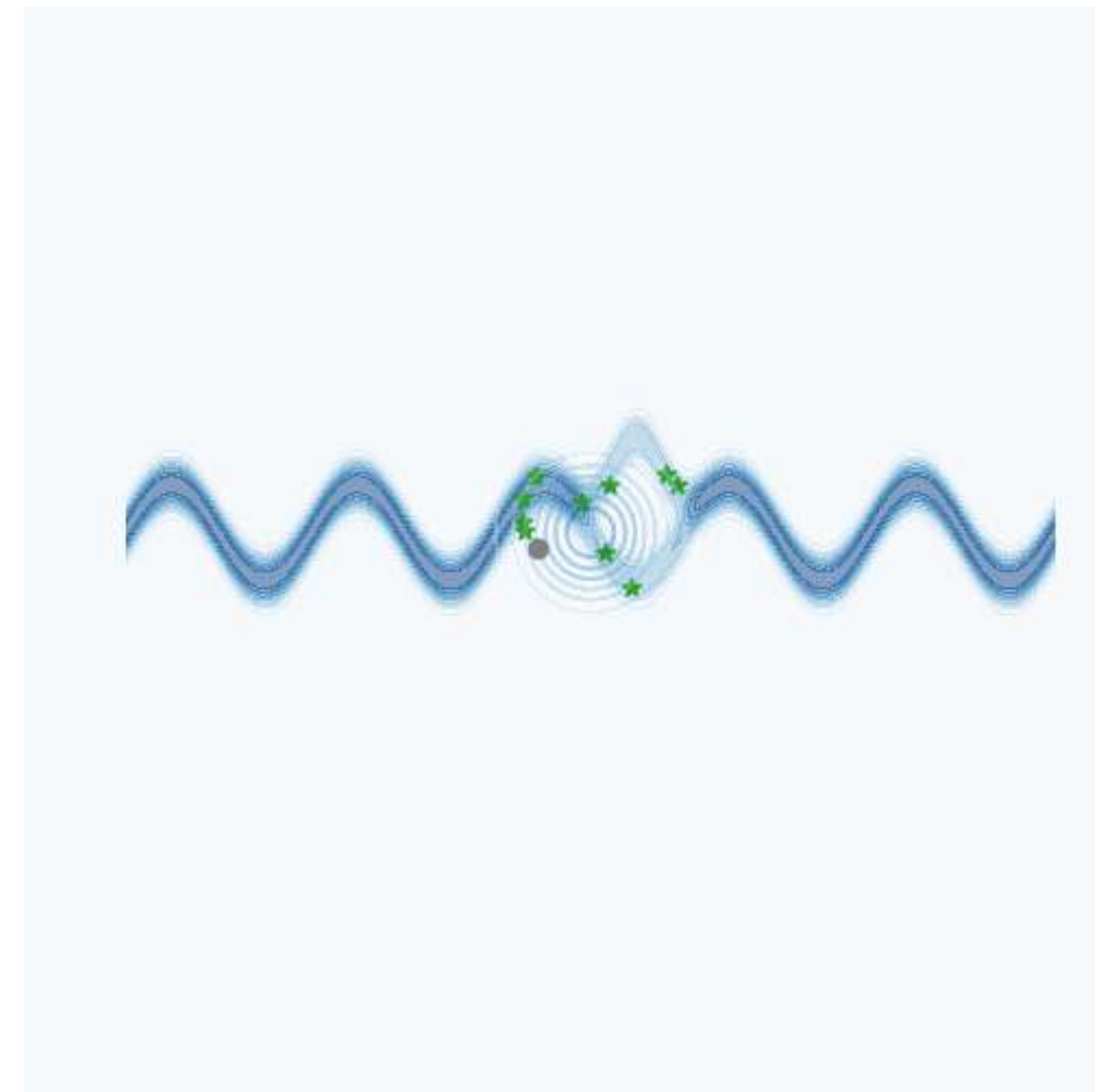
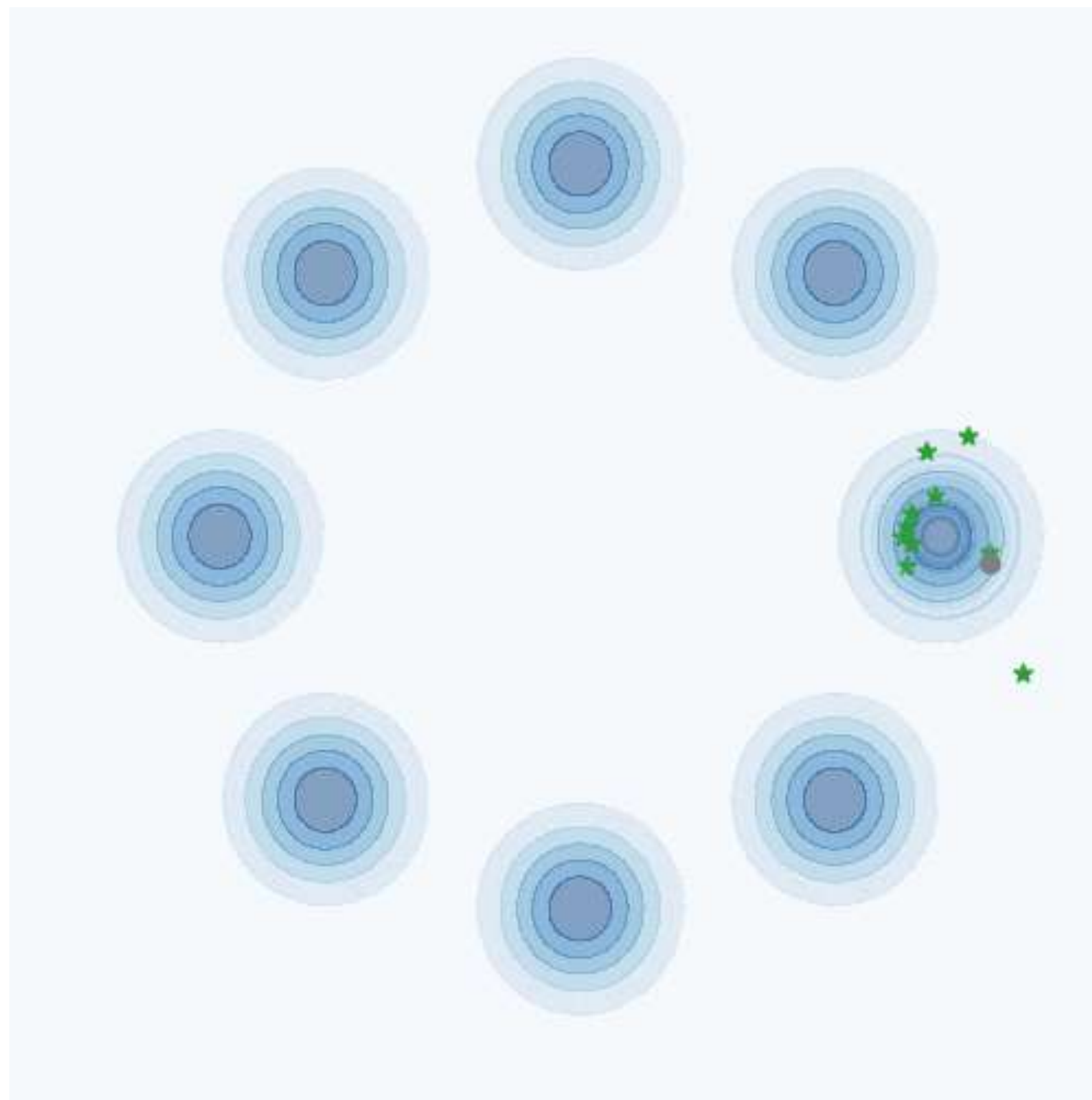


<sup>1</sup> Betancourt et al., 2017

# Training Hamiltonian Monte Carlo [\[code\]](#)



## Effect of the hyperparameter choice





# Training Hamiltonian Monte Carlo [code]

## Gradient-based Optimization<sup>1</sup>

- Tuning the **hyperparameters** via Variational Inference:

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

Implicit

- Add an **inflation** parameter,  $\mathbf{s}$ , for scaling the proposal  $q^{(0)}(\mathbf{z}) = \mathcal{N}(\boldsymbol{\mu}_0, \mathbf{s}^2 \boldsymbol{\Sigma}_0)$

$$\mathbf{s}^* = \operatorname{argmin}_{\mathbf{s}} d(q_{\phi}^{(T)}(\mathbf{z}), p(\mathbf{z}))$$

- Sliced Kernelized Stein Discrepancy<sup>2</sup> (SKSD) measures discrepancy between  $p(\mathbf{z})$  and  $q(\mathbf{z})$  using:

✓ Samples of the approximated distribution

✓ Gradients of the true target

$$\mathbf{s}^* = \operatorname{argmin}_{\mathbf{s}} \text{SKSD}(\mathbf{z}^{(T)}, \nabla_{\mathbf{z}} \log p^*(\mathbf{z}))$$

✓ Robust in high-dimensional spaces.

<sup>1</sup> Campbell et al., 2021

<sup>2</sup> Gong et al., 2020

# Training Hamiltonian Monte Carlo [\[code\]](#)

## Gradient-based Optimization

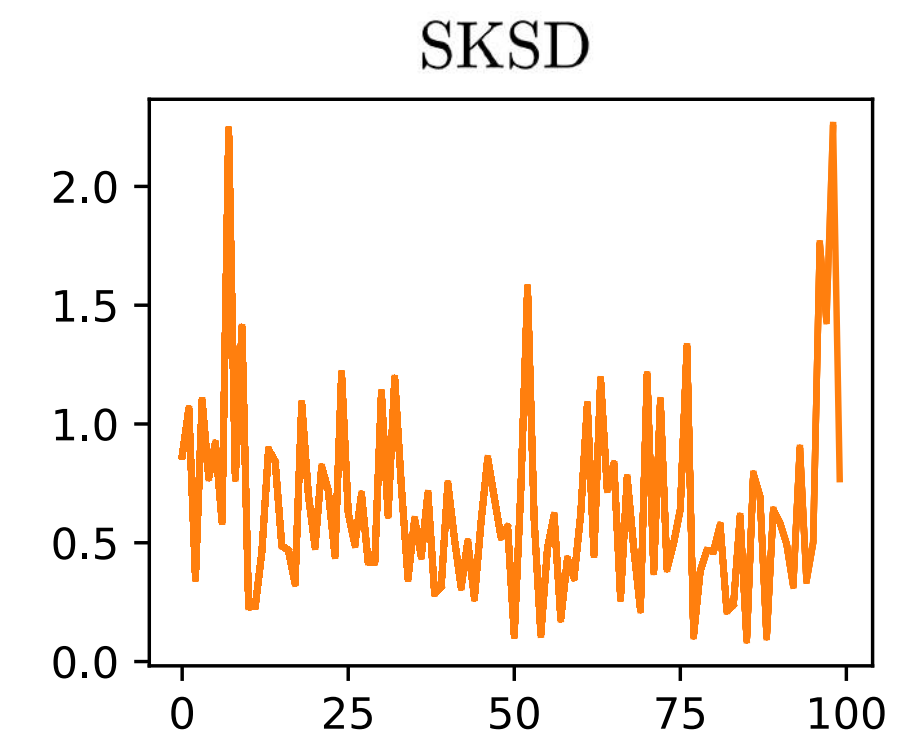
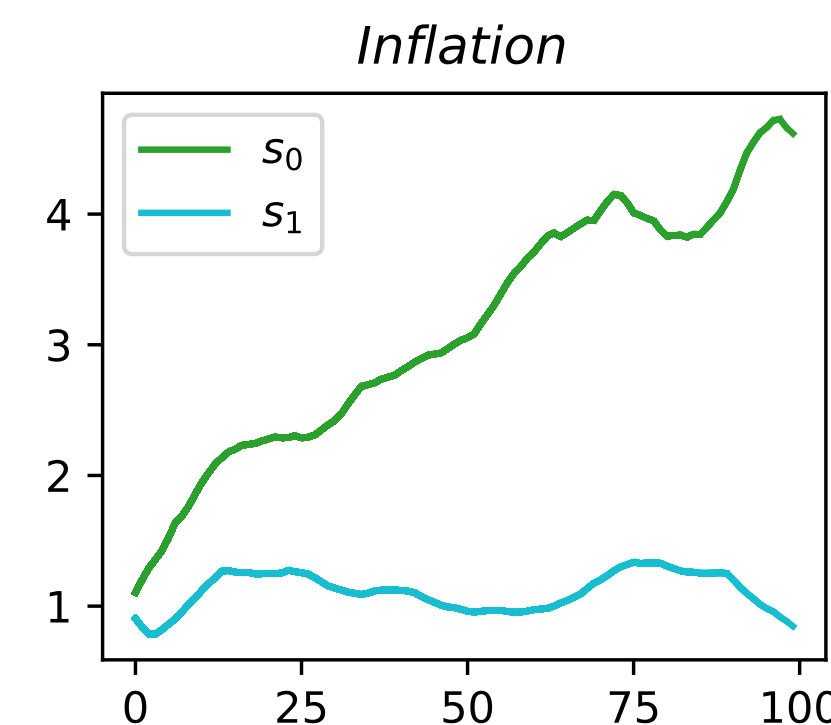
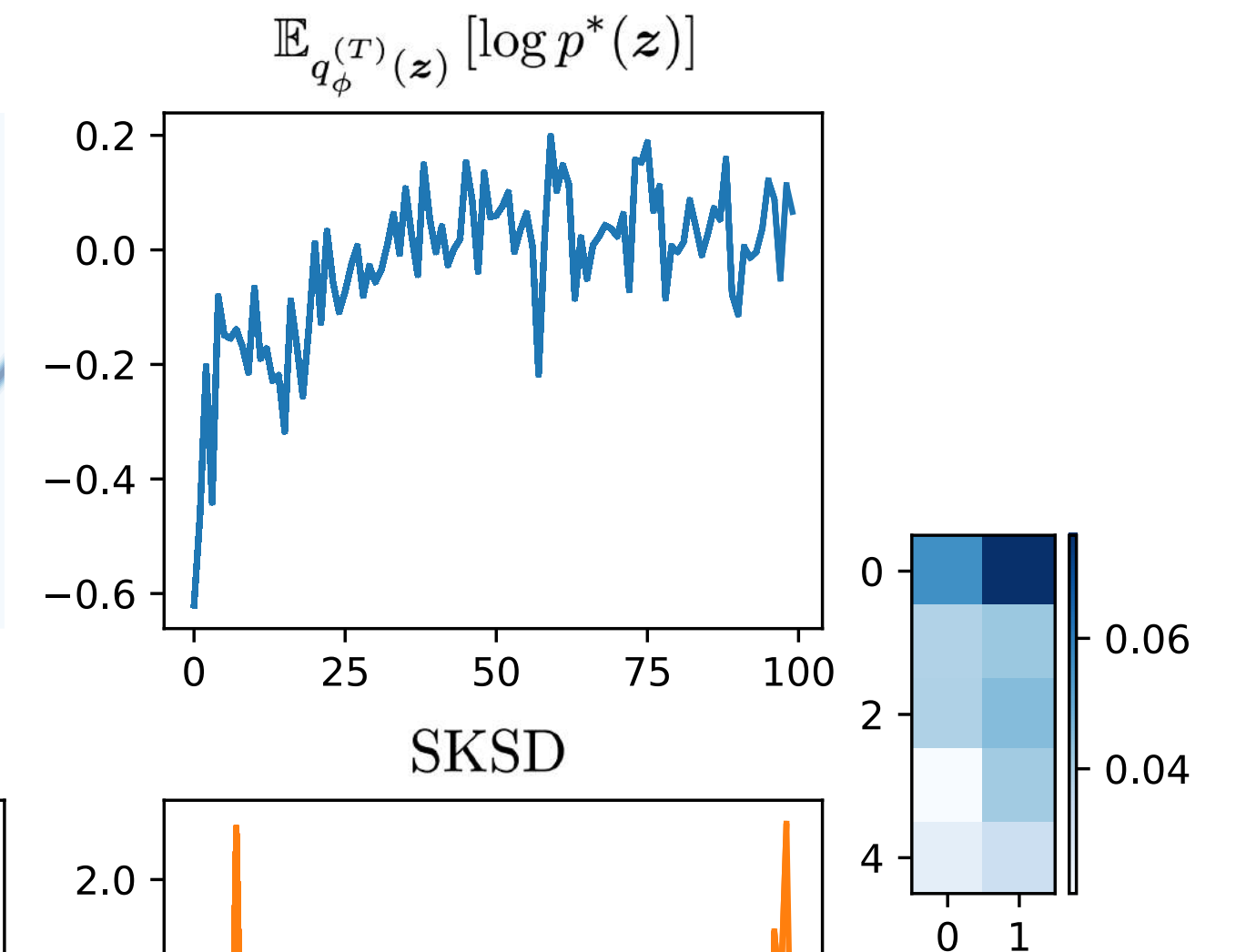
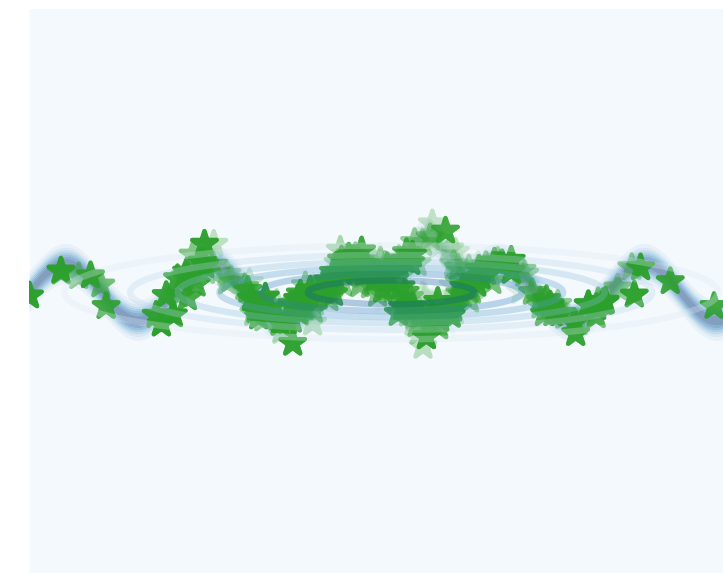
- For each step:

1. Update HMC hyperparameters:

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

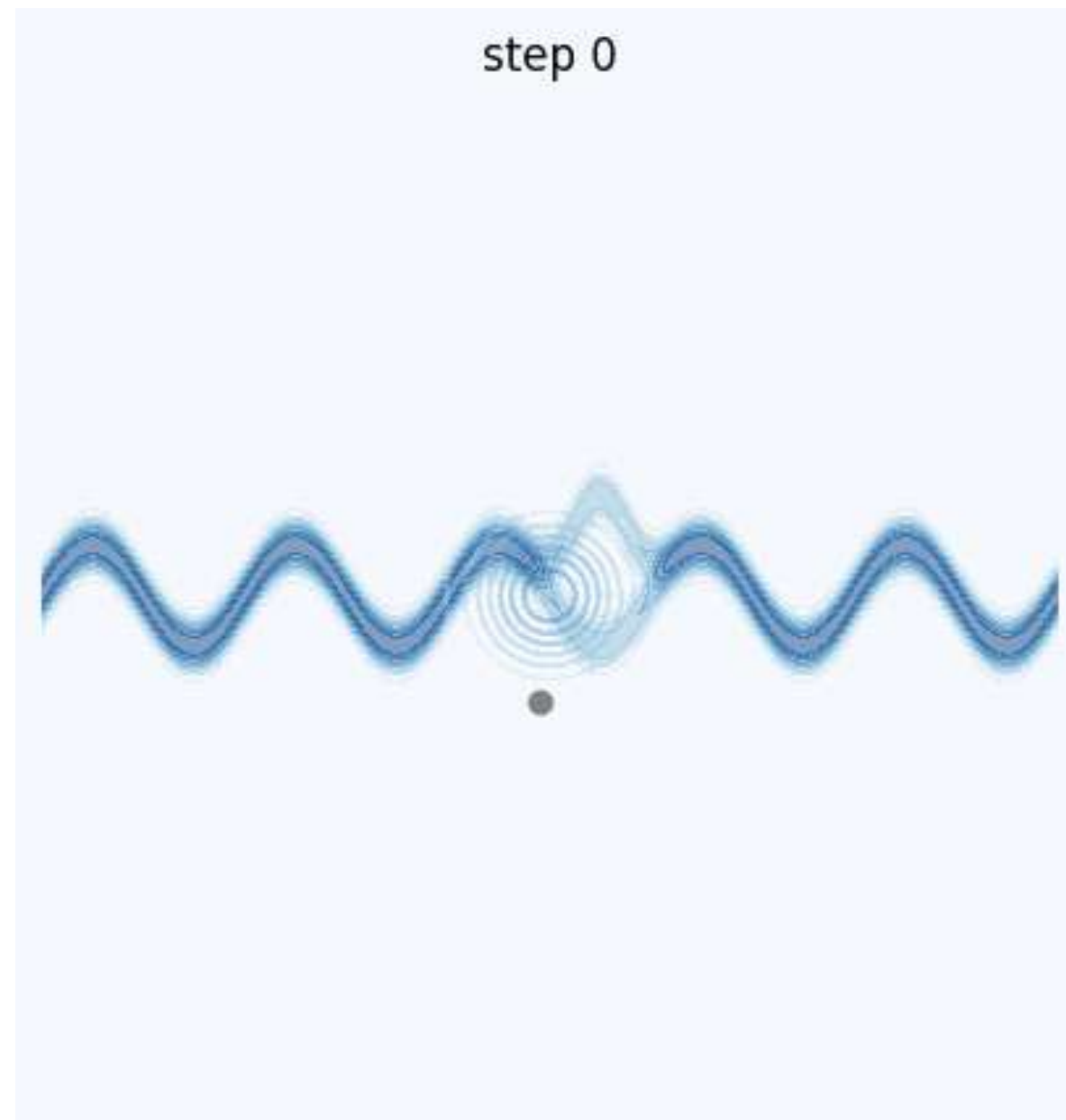
2. Update scaling factor:

$$\mathbf{s}^* = \operatorname{argmin}_{\mathbf{s}} \operatorname{SKSD}(\mathbf{z}^{(T)}, \nabla_{\mathbf{z}} \log p^*(\mathbf{z}))$$

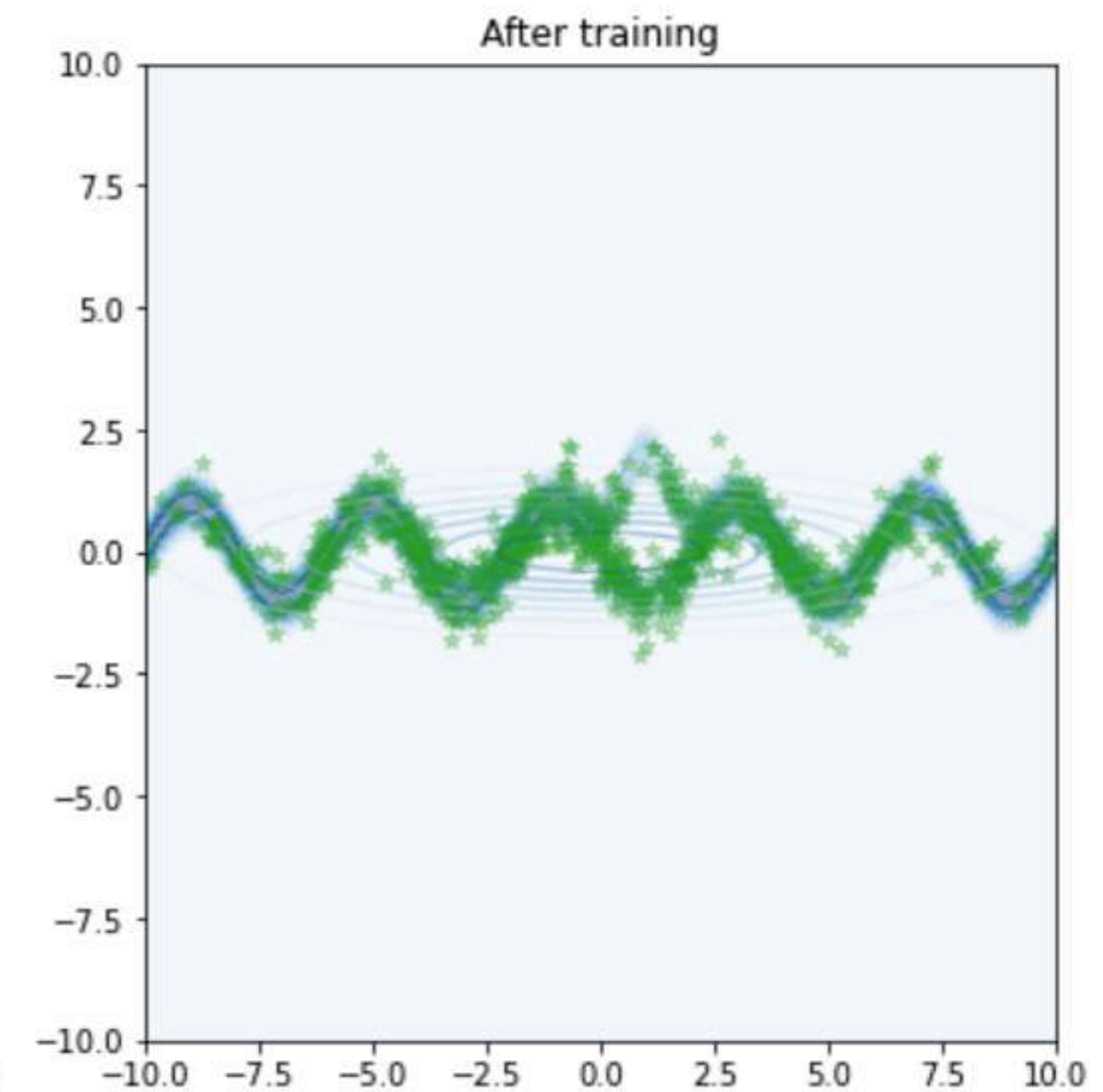
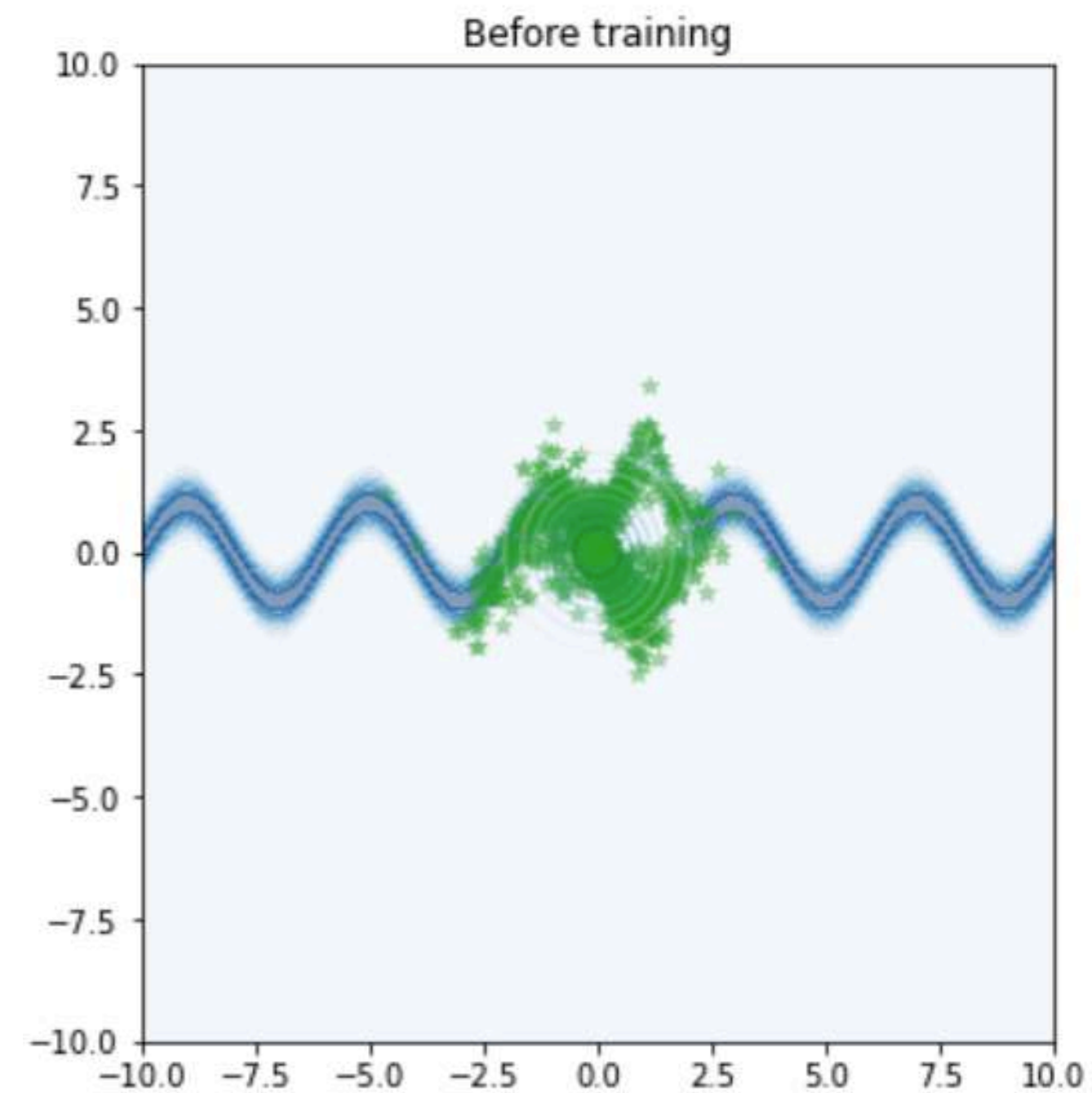


# Training Hamiltonian Monte Carlo [\[code\]](#)

## Gradient-based Optimization



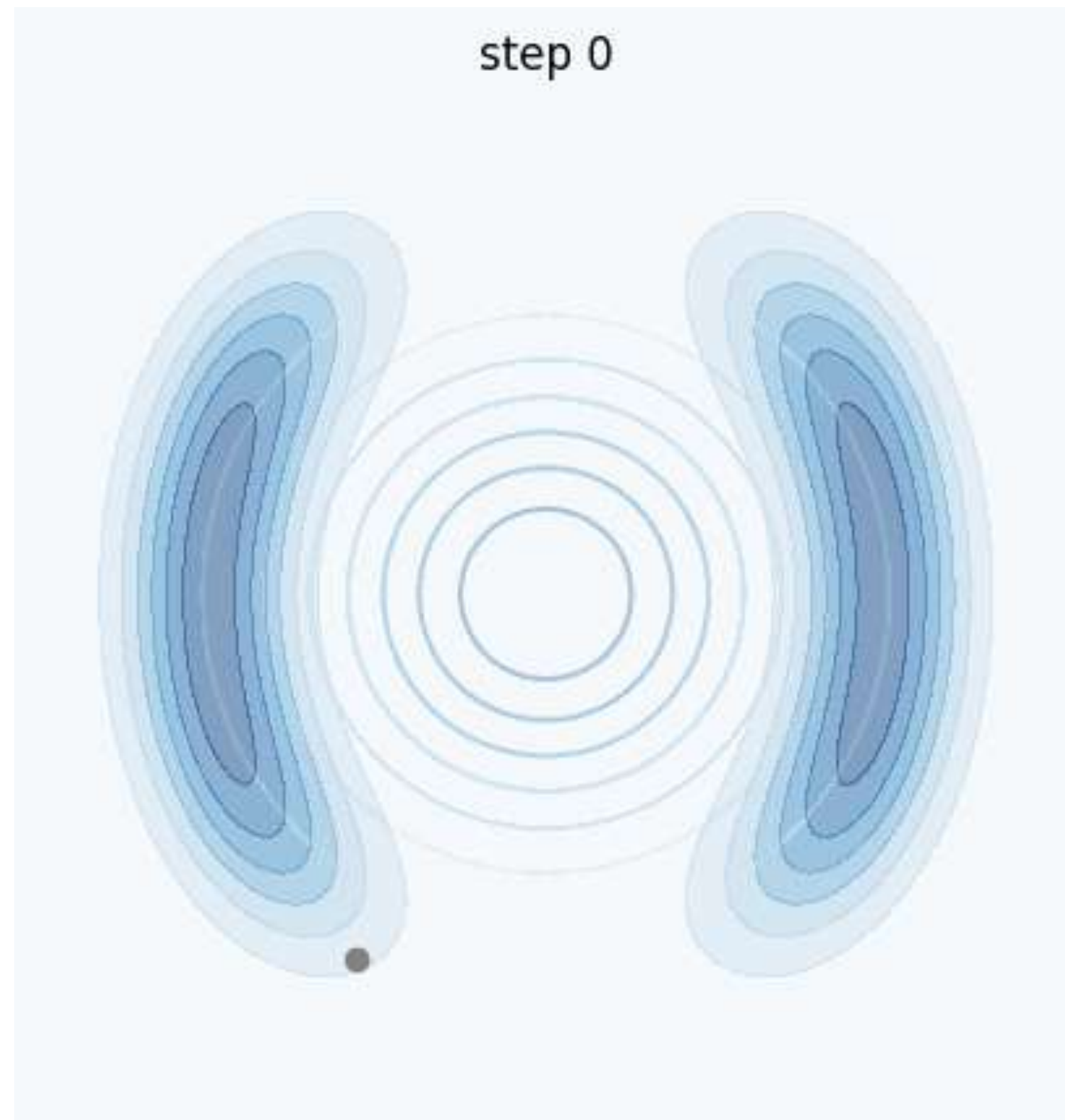
Optimization (gif)



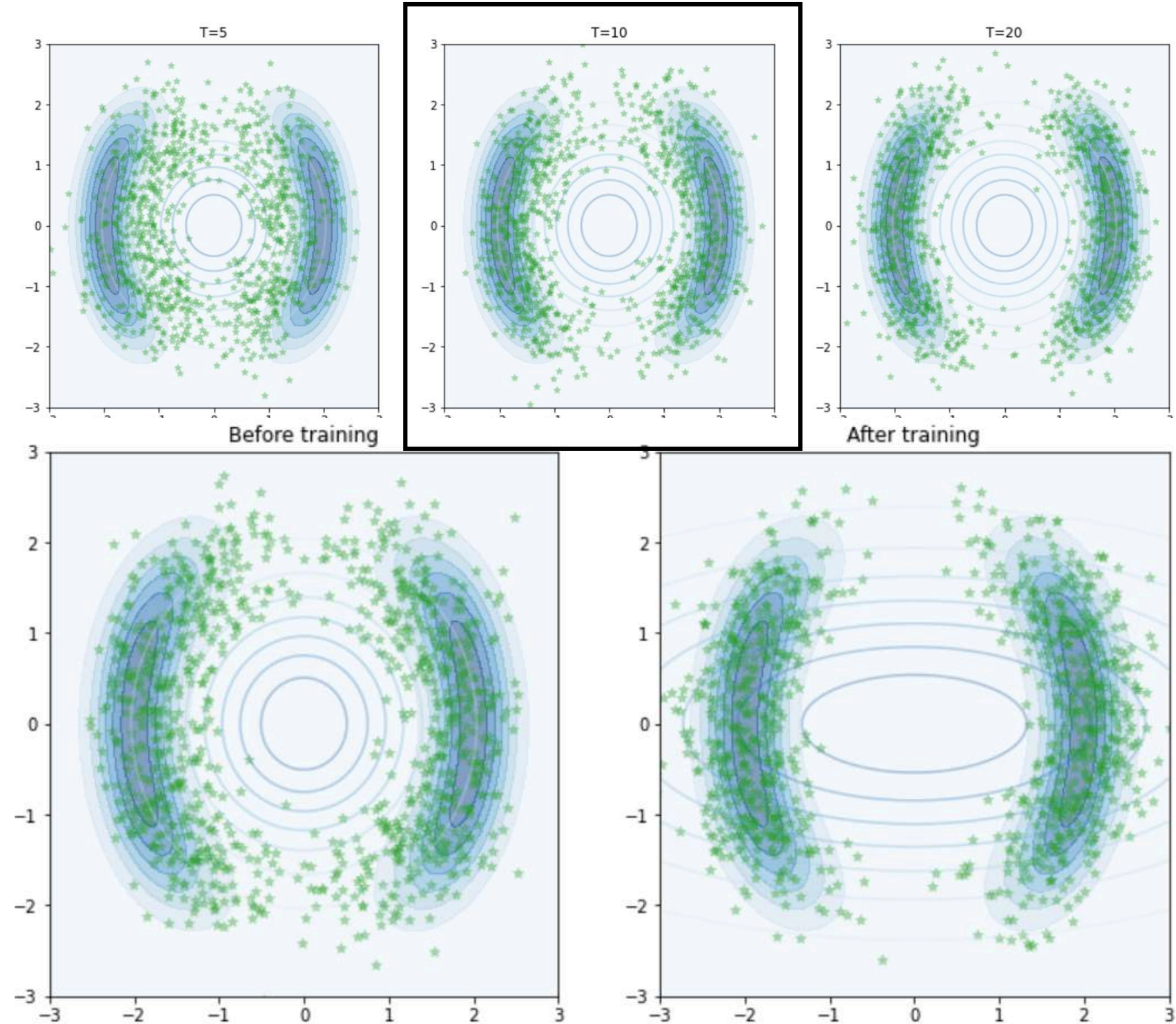


# Training Hamiltonian Monte Carlo [\[code\]](#)

## Gradient-based Optimization



Optimization (gif)





# Hamiltonian Monte Carlo [code]

## Application to VAEs<sup>1,2</sup>

- Approximation of  $p(\mathbf{z} | \mathbf{x})$  improved with  $q^{(T)}(\mathbf{z} | \mathbf{x})$ , using as initial proposal  $q^{(0)}(\mathbf{z} | \mathbf{x})$  given by the encoder.

### Stage 1:

- ★ Pretrain VAE  $(\theta, \psi)$  using ELBO.

### Stage 2:

- ★ Keep training encoder parameters  $\psi$  using ELBO:

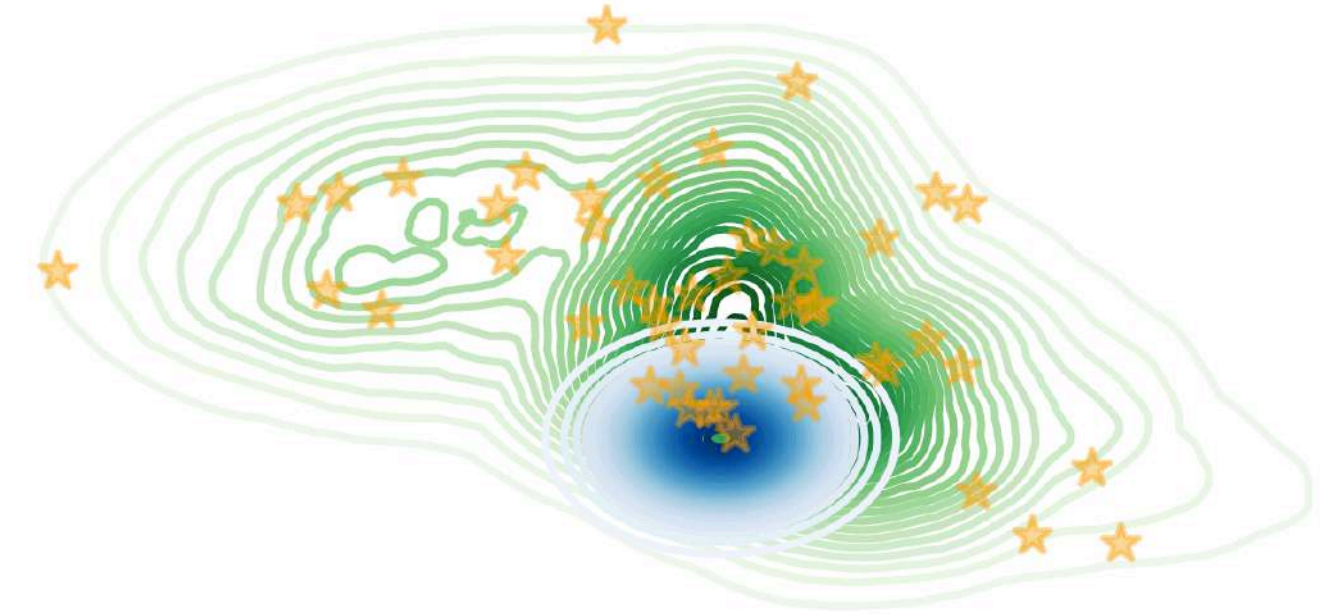
$$\psi^* \leftarrow \text{Adam}_{\psi} \left( \log p_{\theta}(\mathbf{x} | \mathbf{z}^{(0)}) - D_{KL} \left( q_{\psi}^{(0)}(\mathbf{z} | \mathbf{x}) \| p(\mathbf{z}) \right) \right)$$

- ★ Train decoder parameters  $\theta$  and HMC hyperparameters  $\phi$ :

$$(\theta, \phi)^* \leftarrow \text{Adam}_{(\theta, \phi)} \left( \mathbb{E}_{q^{(T)}(\mathbf{z} | \mathbf{x})} [\log p(\mathbf{z}, \mathbf{x})] \right)$$

- ★ Train inflation parameter using:

$$\mathbf{s}^* \leftarrow \text{Adam}_{\mathbf{s}} \left( \text{SKSD} \left( \mathbf{z}^{(T)}, \nabla_{\mathbf{z}} \log p(\mathbf{z}, \mathbf{x}) \right) \right)$$



<sup>2</sup>Samples from the true posterior (orange)

Model	MNIST			Fashion MNIST		
	Scale	Mean	SE	Scale	Mean	SE
VAE	-	-85.08	0.22	-	-108.54	0.60
DReG-IWAE	-	-83.73	0.21	-	-104.48	0.58
maxELT $\alpha = 0$	1.0	-83.48	0.21	1.0	-104.08	0.58
maxELT $\alpha = 1$	1.0	-82.46	0.21	1.0	-103.57	0.58
maxELT $\alpha = 0$ SKSD	6.79	-81.91	0.20	5.58	-103.18	0.58
maxELT $\alpha = 1$ SKSD	3.90	-81.94	0.20	3.59	<b>-102.29</b>	0.57
Hoffman	-	<b>-81.74</b>	0.20	-	-103.04	0.58
Ruiz & Titsias	-	-82.45	0.21	-	-105.13	0.59
Salimans et al.	-	-81.94	-	-	-104.44	0.59
Caterini et al.	-	-82.62	-	-	-104.26	0.58

Approximated test  $\log p(\mathbf{x})$

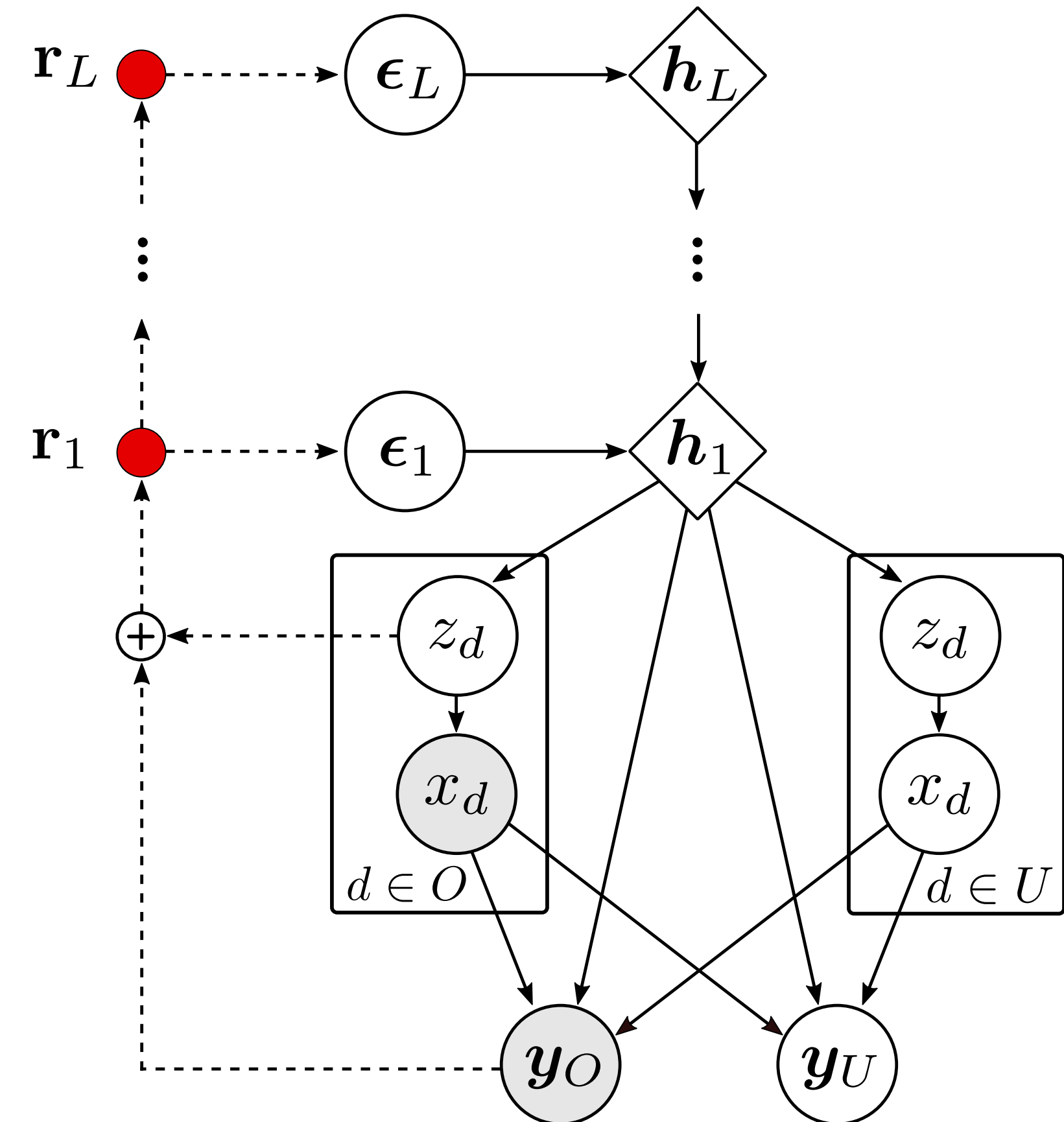
<sup>1</sup> Campbell et al., 2021

<sup>2</sup> Peis et al., 2022

# HH-VAEM

## The Hierarchical Hamiltonian VAE for Mixed-type incomplete data

- Generates data and predictions.
- Models heterogeneous, incomplete data.
- Flexible hierarchical latent space.
- Improved inference via tuned HMC.





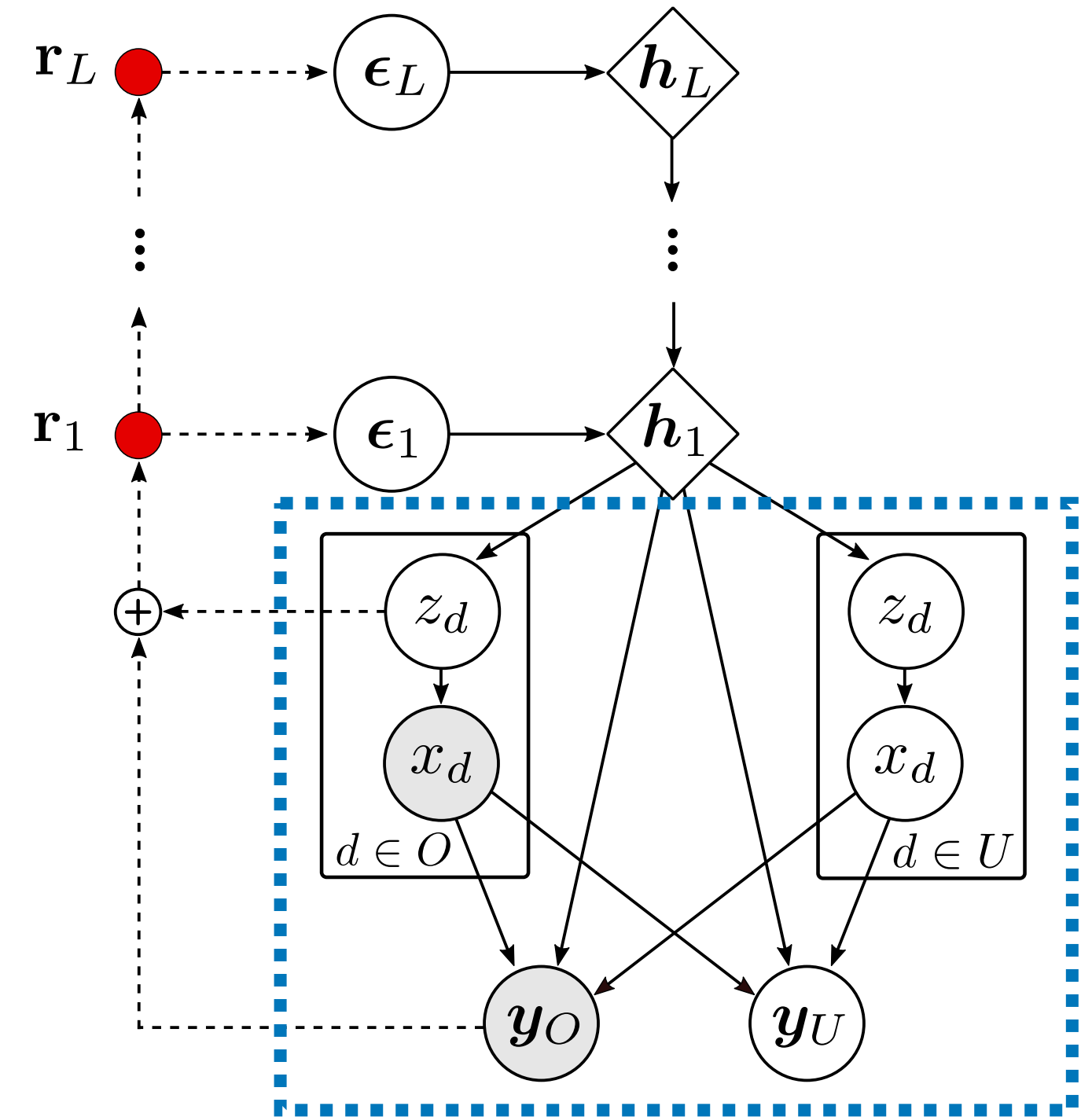
# HH-VAEM

## Mixed-type incomplete data

- Marginal VAEs  $(\theta_d, \gamma_d)$  are pretrained independently on each dimension, with different likelihoods:

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

- Dependency VAE over Gaussian factored dimensions allows dealing with partial heterogeneous data and capture dependencies from balanced likelihoods.



# HH-VAEM

## Hierarchical latent space

- Hierarchical latent space with  $L$  variables  $\mathbf{h} = \{\mathbf{h}_1, \dots, \mathbf{h}_L\}$ .
- **Problem<sup>1</sup>**: HMC fails in densities with huge correlations (AR variables).

$$\nabla_{\mathbf{h}_{1:L}} \log p^*(\mathbf{h}) \uparrow\uparrow$$

- **Solution:**

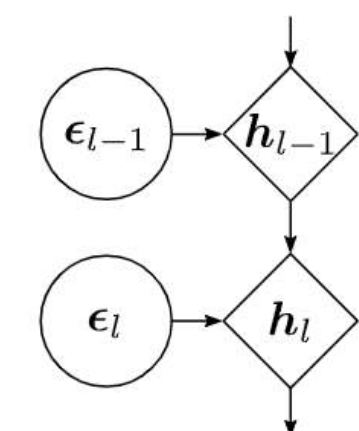
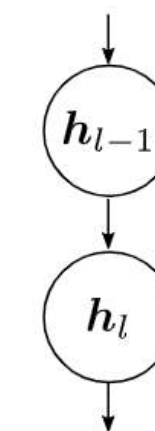
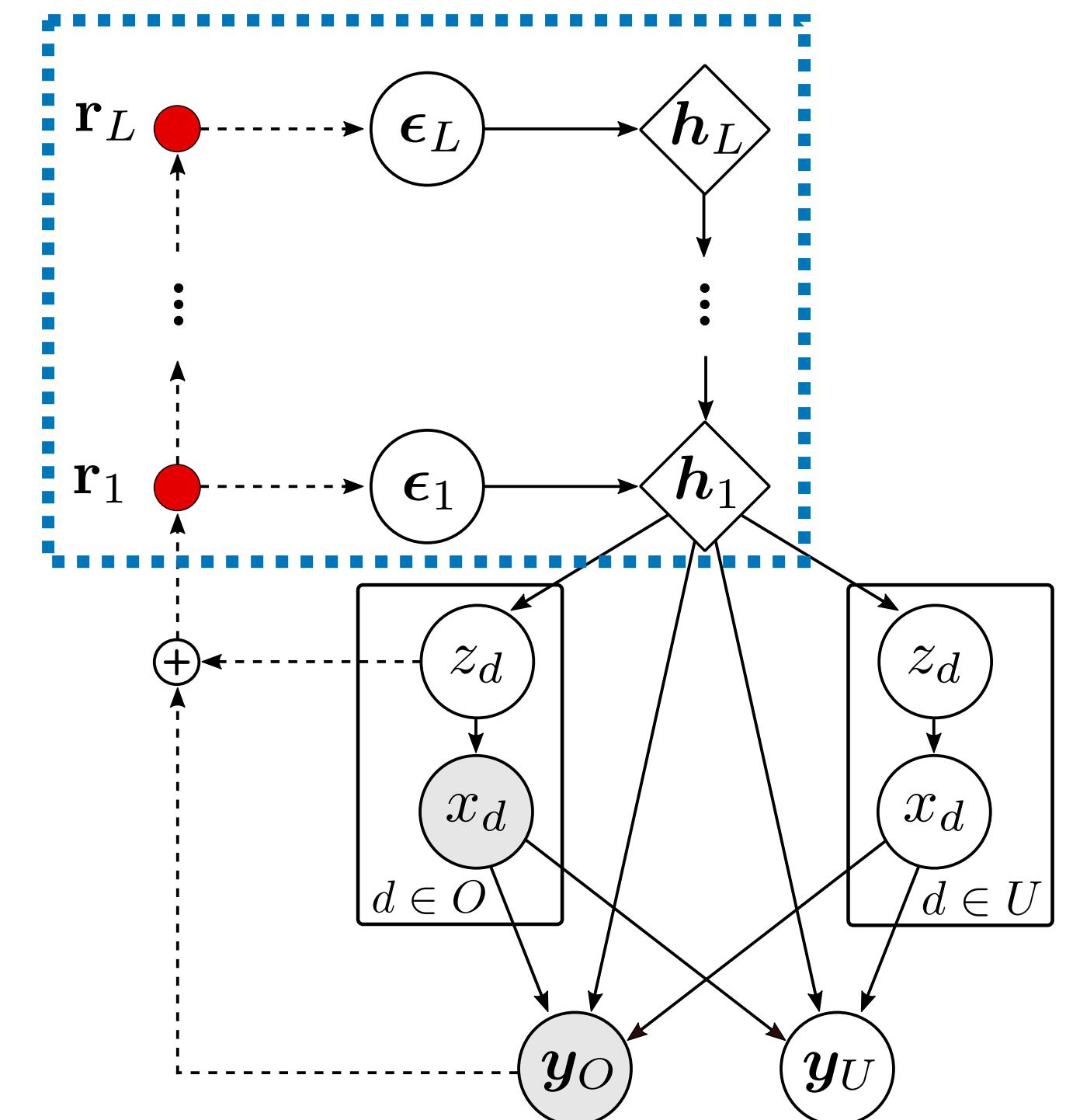
✓ Reparameterization, deterministic hierarchy, relaxed posterior

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

NNs with parameters  $\theta_{\mu_l} \rightarrow f_{\mu_l}$ ,  $\theta_{\sigma_l} \rightarrow f_{\sigma_l}$ . Then  $\theta_l = \{\theta_{\mu_l}, \theta_{\sigma_l}\}$

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

✓ No need to increase complexity of the HMC method.



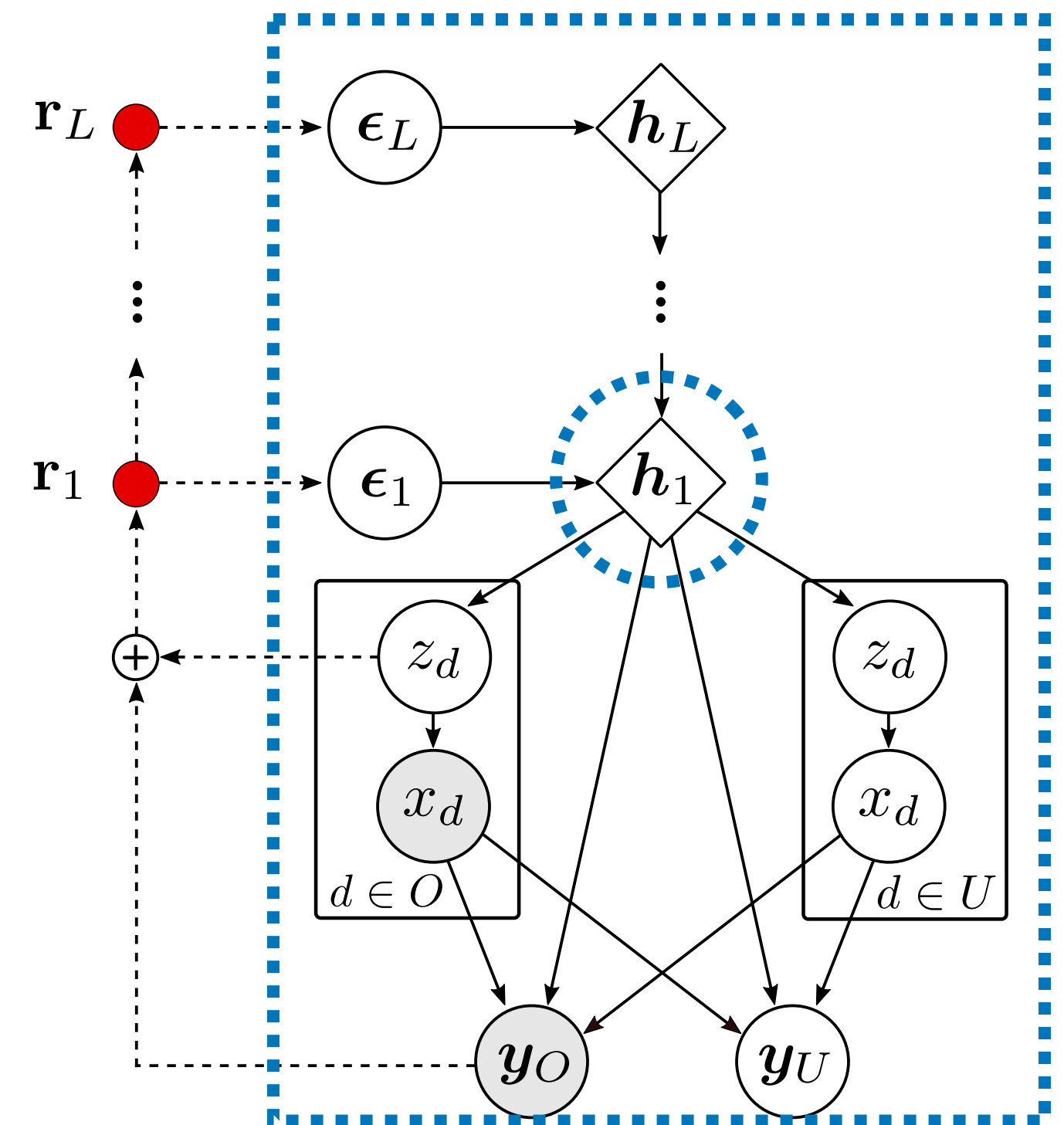
(a) AR hierarchy    (b) Reparameterization

<sup>1</sup> Betancourt et al., 2015

# HH-VAEM

## Generative model

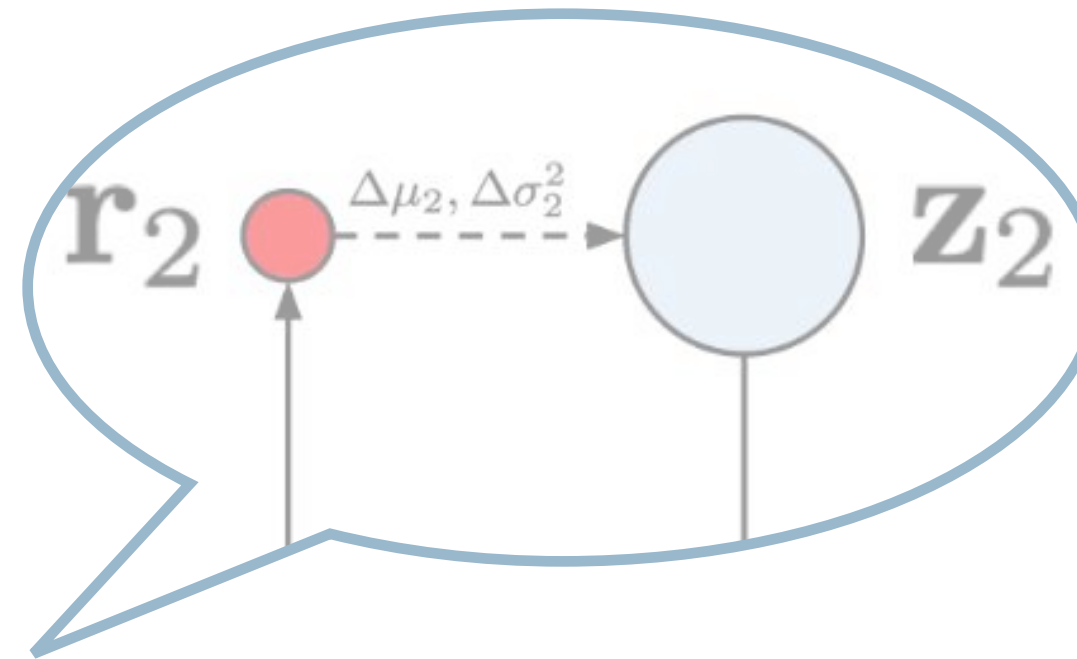
- Generate from the superficial layer:
  1. Data:  $p(z|h_1)$  as NN with parameters  $\theta_z$
  2. Predictions:  $p(y|\hat{x}, h_1)$  as NN with parameters  $\theta_y$
- $\hat{x}$  is the concatenation  $[x_O, \hat{x}_U]$  with the imputed missing part.
- The **predictor** parameters are jointly trained with the model.
- Generative parameters:  $\theta = \{\theta_z, \theta_y, \theta_1, \dots, \theta_L\}$





# HH-VAEM

## Hierarchical encoder



- Bottom-up path, but no sharing parts needed! <sup>1,2</sup>

$$\mathbf{r}_0 = \{\mathbf{x}_O, \mathbf{y}_O\}, \quad \mathbf{r}_l = f_r(\mathbf{r}_{l-1})$$

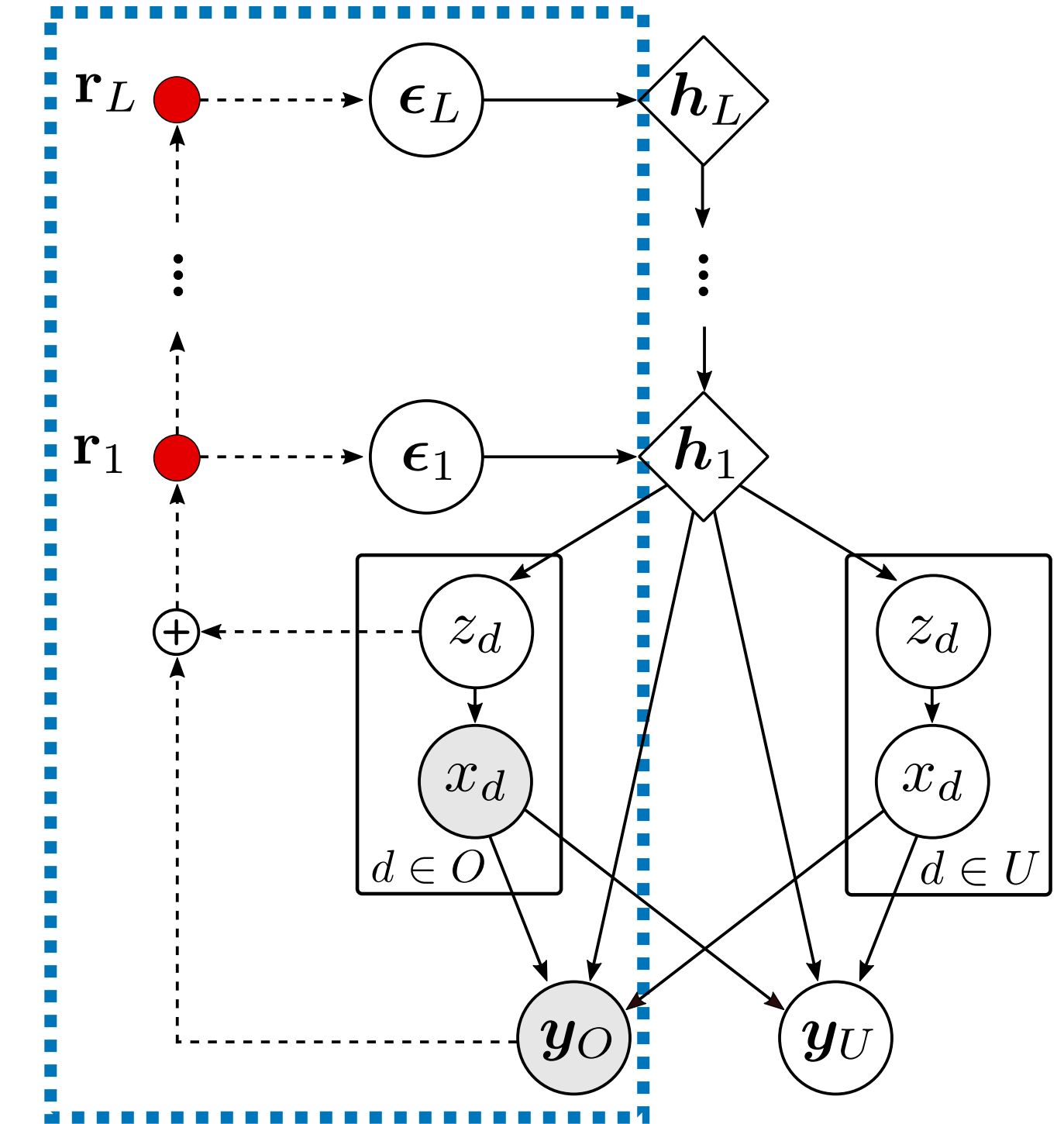
NNs with parameters  $\psi_{r_l} \rightarrow f_{r_l}$ .

- Posterior approximation for each layer:

$$q_{\psi_l}(\epsilon_l | \mathbf{x}_O, \mathbf{y}_O) = \mathcal{N}(g_{\mu_l}(\mathbf{r}_l), g_{\sigma_l}(\mathbf{r}_l))$$

as NNs with parameters  $\psi_{\mu_l} \rightarrow g_{\mu_l}$ ,  $\psi_{\sigma_l} \rightarrow g_{\sigma_l}$

- Encoder parameters:  $\psi = \{\psi_{r_l}, \psi_{\mu_l}, \psi_{\sigma_l}\}_{l=1}^L$



<sup>1</sup> Vahdat et al., 2020

<sup>1</sup> Maaløe et al., 2019



# HH-VAEM

## Training HMC on the posterior distribution

- Include HMC with hyper parameters  $\phi$ .
- Define HMC target as the unnormalised posterior for sampling  $\epsilon = \{\epsilon_1, \dots, \epsilon_L\}$  from the posterior:

$$p^*(\epsilon_1, \dots, \epsilon_L, \mathbf{z}_O, \mathbf{y}_O) = p_\theta(\mathbf{z}_O | \mathbf{h}_1) p_\theta(\mathbf{y}_O | \hat{\mathbf{x}}, \mathbf{h}_1) \prod_{i=1}^L p(\epsilon_i)$$

- Define the HMC objective as

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

- Tune the scale factor  $s$  using the SKSD discrepancy:

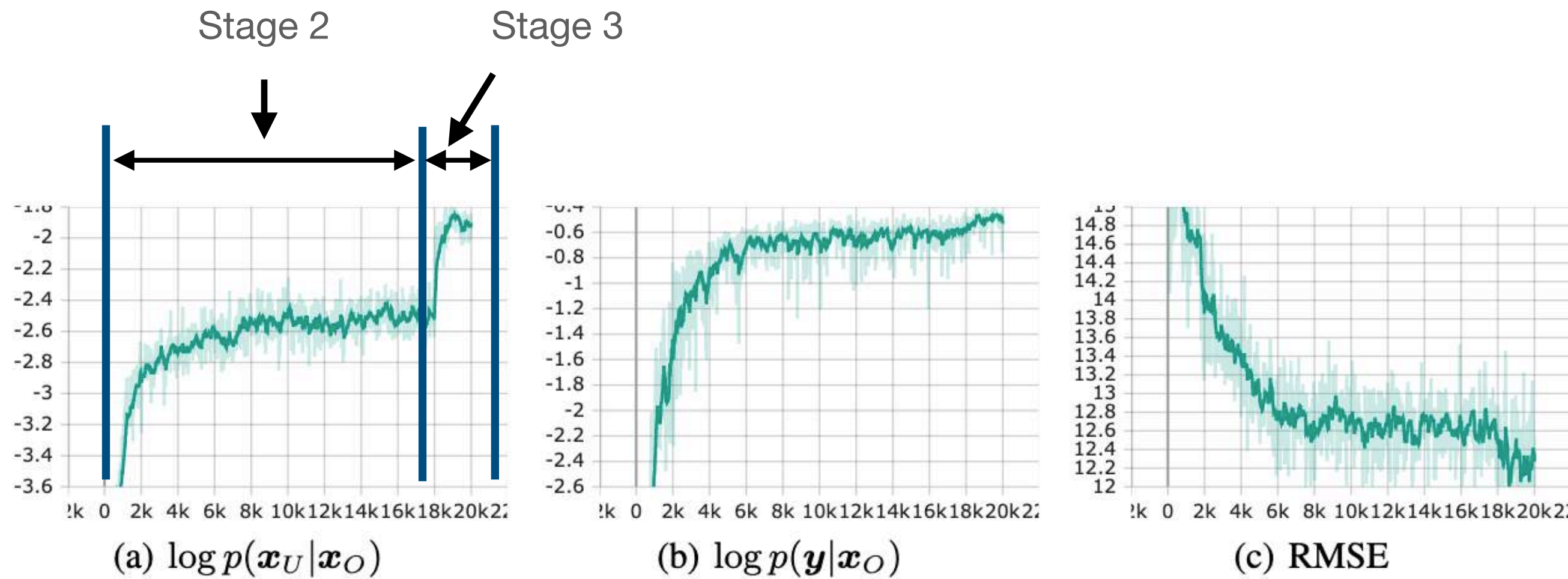
$$\mathcal{L}_{SKSD}(\mathbf{x}_O, \mathbf{y}_O; s) = \text{SKSD} \left( q_\phi^{(T)}(\epsilon | \mathbf{z}_O, \mathbf{x}_O, \mathbf{y}_O; s), p(\epsilon | \mathbf{z}_O, \mathbf{x}_O, \mathbf{y}_O) \right)$$



# HH-VAEM

## Training algorithm

- **Stage 1:** train marginal VAEs on each dimension.
- **Stage 2:** pretrain using ELBO.
- **Stage 3:** jointly train VAE + HMC.




---

### 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**

---

# HH-VAEM

## Computational cost

- The  $L$  Leapfrog steps are executed each cycle in  $t = 1 : T$

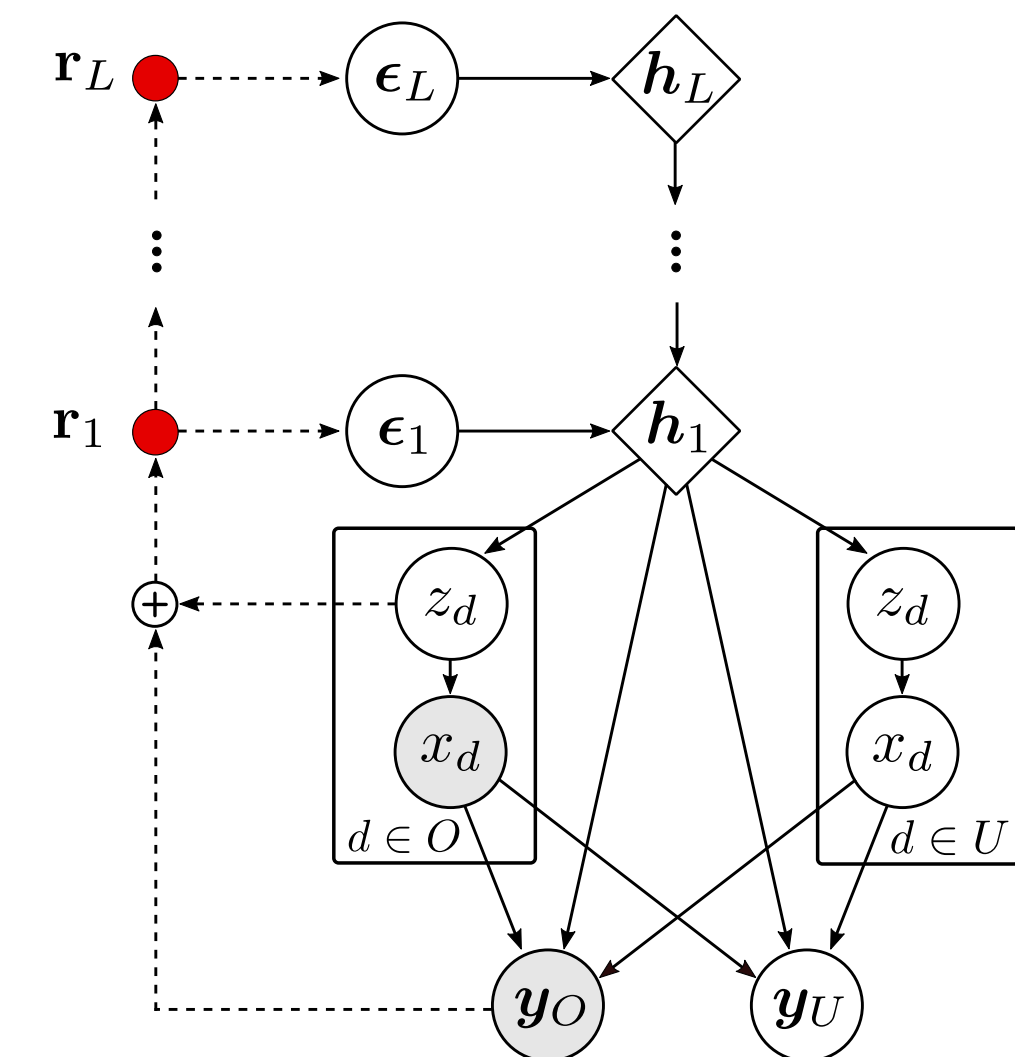
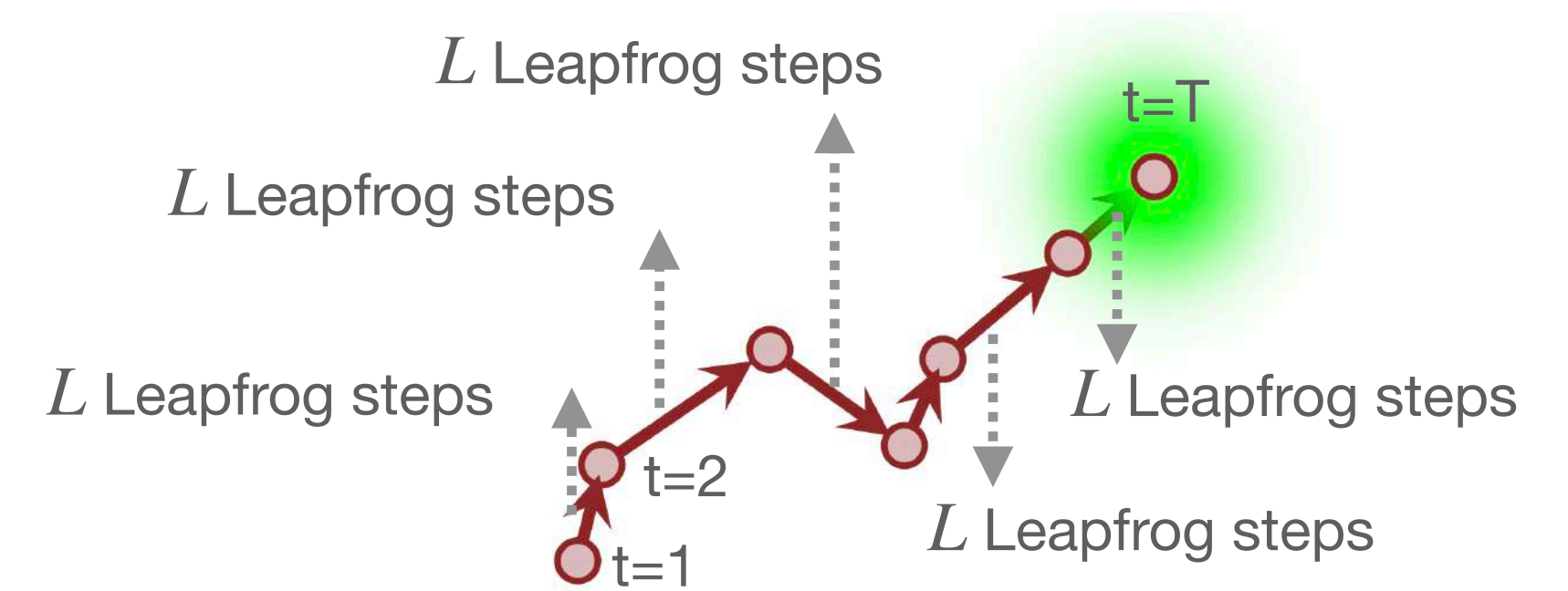
(batch\_size, parallel\_chains, latent\_dimension)

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

$$\epsilon_{i+1} = \epsilon_i + \mathbf{r}_{i+\frac{1}{2}} \odot \phi \odot \frac{1}{M},$$

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

- Where  $\log p^*(\epsilon) = \log p(\epsilon, \mathbf{z}, \mathbf{y}) = \log p(\mathbf{z} | \mathbf{h}_1) + p(\mathbf{y} | \mathbf{h}_1, \hat{\mathbf{x}}) + p(\epsilon)$
- Computing the **gradients** requires:
  - Obtaining likelihood parameters (decoder for  $p(\mathbf{z} | \mathbf{h}_1)$ , predictor for  $p(\mathbf{y} | \mathbf{h}_1, \hat{\mathbf{x}})$ ).
  - Evaluating and perform the automatic differentiation.
- The extra computational cost is approximately a factor of  $2TL$ .



# Bayesian active feature acquisition

## Proposed method: sampling-based

- Bayesian reward<sup>1</sup> can also be expressed in terms of the Mutual Information:

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

- Sampling-based estimator of the Mutual Information<sup>2</sup>:

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

✓ More flexible than the encoder-based method.

✓ Efficient, easy parallelization.

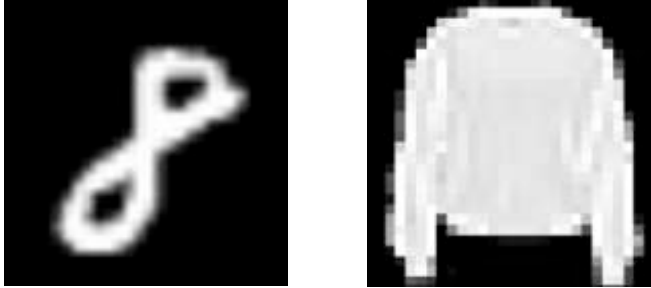
<sup>1</sup> Bernardo et al., 1979

<sup>2</sup> Kraskov et al., 2004



# Experiments

## Set up

- **HH-VAEM** with 2 layers of latent variables.
- Baseline models:
  - **VAEM**: The VAEM<sup>1</sup> strategy.
  - **MIWAEM**: VAEM combined with the importance weighted estimation proposed in MIWAE<sup>2</sup>.
  - **H-VAEM**: A Hierarchical VAEM with two layers of latent variables and a Gaussian encoder.
  - **HMC-VAEM**: A VAEM that includes a tuned HMC sampler for the true posterior.
- Datasets:
  - MNIST, Fashion-MNIST -> without marginal VAEs. 
  - 10 UCI Datasets with mixed-type data: Bank, Insurance, Avocado, Naval, Yatch, Diabetes, Concrete, Wine, Energy, Boston.
- Configuration: manually introduced missing features and target with a probability sampled from  $U(0.01, 0.99)$  each batch.

<sup>1</sup> Ma et al., 2020    <sup>2</sup> Mattei et al., 2019

# Experiments

## Mixed-type data

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

	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	$0.124 \pm 0.001$	$0.121 \pm 0.001$	$0.119 \pm 0.001$	$0.101 \pm 0.004$	<b><math>0.094 \pm 0.003</math></b>
F-MNIST	$0.162 \pm 0.002$	$0.160 \pm 0.002$	$0.156 \pm 0.002$	$0.150 \pm 0.002$	<b><math>0.144 \pm 0.002</math></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	$0.153 \pm 0.009$	$0.151 \pm 0.007$	$0.146 \pm 0.006$	$0.067 \pm 0.007$	<b><math>0.056 \pm 0.019</math></b>
F-MNIST	$0.501 \pm 0.012$	$0.496 \pm 0.008$	$0.494 \pm 0.007$	$0.357 \pm 0.060$	<b><math>0.337 \pm 0.069</math></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	$0.953 \pm 0.004$	$0.953 \pm 0.003$	$0.953 \pm 0.003$	$0.978 \pm 0.003$	<b><math>0.981 \pm 0.005</math></b>
F-MNIST	$0.824 \pm 0.005$	$0.824 \pm 0.004$	$0.824 \pm 0.004$	$0.869 \pm 0.015$	<b><math>0.876 \pm 0.017</math></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  $\hat{I}(\mathbf{y}; x_i | \mathbf{x}_O)$

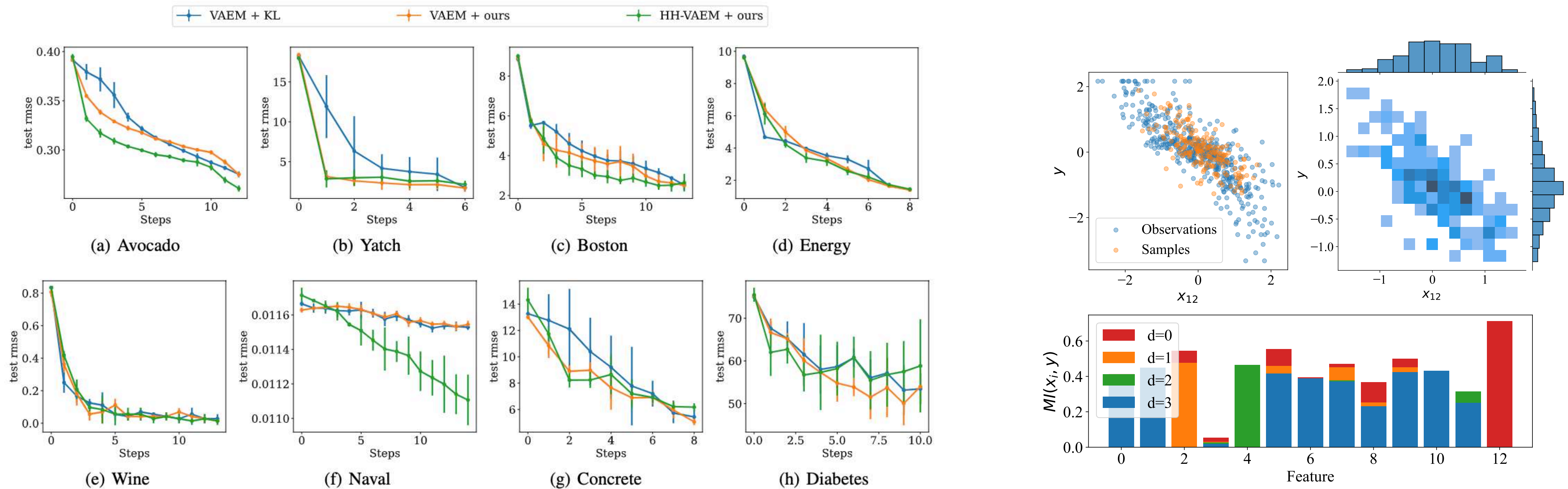


Figure 5: SAIA metric curves. Horizontal axis shows acquisition steps (number of discovered features). Vertical axis is the RMSE.

# 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 learning task.

# References

- 📄 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.
- 📄 Kingma, D. P., & Welling, M. (2013). Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.
- 📄 Nazabal, A., Olmos, P. M., Ghahramani, Z., & Valera, I. (2020). Handling incomplete heterogeneous data using vaes. *Pattern Recognition, 107*, 107501.
- 📄 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.
- 📄 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*.
- 📄 Vahdat, A., & Kautz, J. (2020). NVAE: A deep hierarchical variational autoencoder. *Advances in Neural Information Processing Systems, 33*, 19667-19679.
- 📄 Maaløe, L., Fraccaro, M., Liévin, V., & Winther, O. (2019). Biva: A very deep hierarchy of latent variables for generative modeling. *Advances in neural information processing systems, 32*.



# References

- Child, R. (2020). Very deep vaes generalize autoregressive models and can outperform them on images. *arXiv preprint arXiv:2011.10650*.
- Tomczak, J., & Welling, M. (2018, March). VAE with a VampPrior. In *International Conference on Artificial Intelligence and Statistics* (pp. 1214-1223). PMLR.
- 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.
- Cremer, C., Li, X., & Duvenaud, D. (2018, July). Inference suboptimality in variational autoencoders. In *International Conference on Machine Learning* (pp. 1078-1086). PMLR.
- 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.
- Bernardo, J. M. (1979). Expected information as expected utility. *the Annals of Statistics*, 686-690.
- Betancourt, M. (2017). A conceptual introduction to Hamiltonian Monte Carlo. *arXiv preprint arXiv:1701.02434*.
- Neal, R. M. (2011). MCMC using Hamiltonian dynamics. *Handbook of markov chain monte carlo*, 2(11), 2.

# References

- 📄 Betancourt, M., & Girolami, M. (2015). Hamiltonian Monte Carlo for hierarchical models. *Current trends in Bayesian methodology with applications*, 79(30), 2-4.
- 📄 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.
- 📄 Gong, W., Li, Y., & Hernández-Lobato, J. M. (2020). Sliced kernelized Stein discrepancy. *arXiv preprint arXiv:2006.16531*.
- 📄 Kraskov, A., Stögbauer, H., & Grassberger, P. (2004). Estimating mutual information. *Physical review E*, 69(6), 066138.