

# Langevin Autoencoders for Learning Deep Latent Variable Models

Shohei Taniguchi, Yusuke Iwasawa, Wataru Kumagai, Yutaka Matsuo



# Deep Latent Variable Models

## Definition

**x** : observation  
**z** : latent variable  
 **$\theta$**  : model parameter

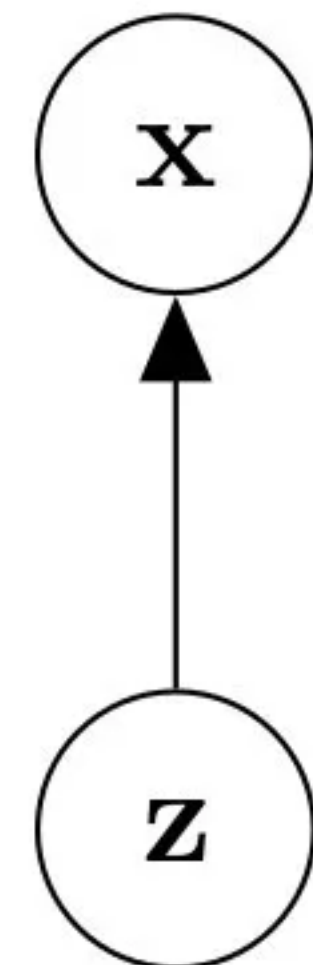
$$p(\mathbf{x}; \boldsymbol{\theta}) = \int p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta}) p(\mathbf{z}) d\mathbf{z}$$

- $p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta})$  is typically constructed using a deep neural network  $f_{\mathbf{x}|\mathbf{z}}$

e.g.,

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, I)$$

$$p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta}) = \mathcal{N}\left(\mathbf{x}; f_{\mathbf{x}|\mathbf{z}}(\mathbf{z}; \boldsymbol{\theta}), \sigma^2 I\right)$$



# Deep Latent Variable Models

## Training via Maximum Likelihood

**x** : observation  
**z** : latent variable  
**θ** : model parameter

$$\nabla_{\theta} \mathbb{E}_{\hat{p}_{\text{data}}(\mathbf{x})} [\log p(\mathbf{x}; \theta)] \approx \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{p(\mathbf{z}^{(i)} | \mathbf{x}^{(i)}; \theta)} \left[ \nabla_{\theta} \log p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)}; \theta) \right]$$

- $\hat{p}_{\text{data}}$ : Empirical distribution defined by training set  
 $n$ : Minibatch size
- Due to the intractability of the posterior  $p(\mathbf{z} | \mathbf{x}; \theta)$ , approximation is needed



# Langevin Dynamics

**$x$**  : observation  
 **$z$**  : latent variable  
 **$\theta$**  : model parameter

- Langevin dynamics (LD) is an MCMC based on the following Langevin equation

$$dz = -\nabla_z U(x, z; \theta) dt + \sqrt{2} dB$$

$$U(x, z; \theta) = -\log p(x, z; \theta)$$

- This stochastic differential equation has the posterior as a stationary distribution

# Langevin Dynamics

**$\mathbf{x}$**  : observation  
 **$\mathbf{z}$**  : latent variable  
 **$\boldsymbol{\theta}$**  : model parameter  
 **$\eta$**  : stepsize

- By simulating the dynamics, samples asymptotically approach to the posterior

for  $t = 1, \dots, T$

$$\mathbf{z}_{t+1} \sim \mathcal{N} \left( \mathbf{z}_{t+1}; \mathbf{z}_t - \eta \nabla_{\mathbf{z}_t} U(\mathbf{x}, \mathbf{z}_t; \boldsymbol{\theta}), 2\eta \mathbf{I} \right)$$

- Optionally, Metropolis-Hastings rejection steps can be added for calibrating discretization error

# Langevin Dynamics

## Pros

- Samples are asymptotically unbiased
  - High approximation power

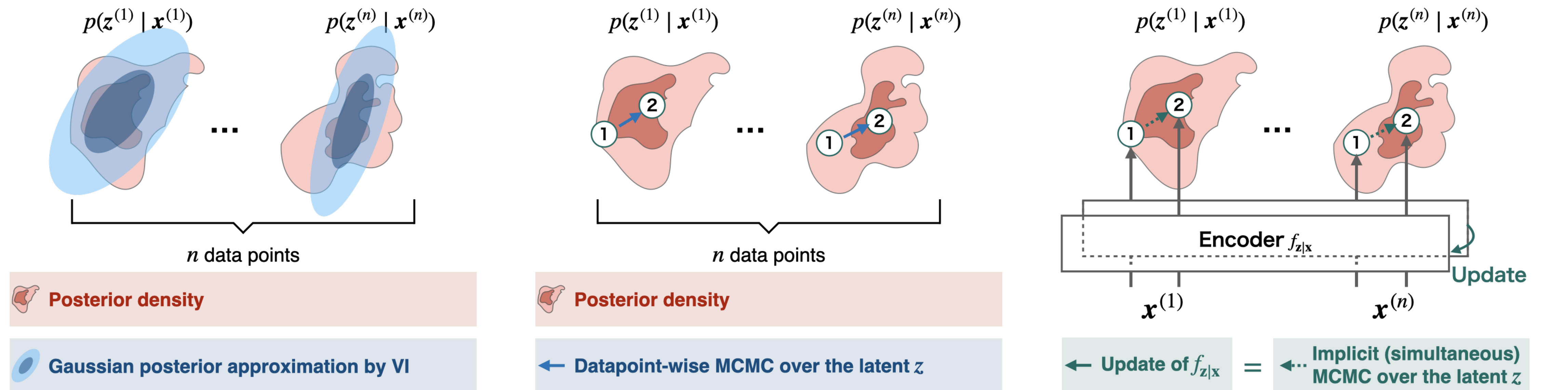
## Cons

- Datapoint-wise sampling is costly
  - Need to run MCMC independently for all minibatch data  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$

# Method

# Amortized Langevin Dynamics

## Overview



Variational Inference

Langevin Dynamics

Amortized Langevin Dynamics (ours)

**MCMC on the latent space is replaced by MCMC on the encoder's parameter space**



# Amortized Langevin Dynamics

## Formulation

$\mathbf{x}$  : observation  
 $\mathbf{z}$  : latent variable  
 $\theta$  : model parameter  
 $f_{\mathbf{z}|\mathbf{x}}$  : encoder

- Consider an SDE on the parameter of an encoder  $f_{\mathbf{z}|\mathbf{x}}$  that maps  $\mathbf{x}$  into  $\mathbf{z}$

$$d\boldsymbol{\phi} = -\nabla_{\boldsymbol{\phi}} V(\boldsymbol{\phi}) dt + \sqrt{2} dB$$

$$V(\boldsymbol{\phi}) = \sum_{i=1}^n U\left(\mathbf{x}^{(i)}, f_{\mathbf{z}|\mathbf{x}}(\mathbf{x}^{(i)}; \boldsymbol{\phi}); \theta\right)$$

# Amortized Langevin Dynamics

## Algorithm

$\mathbf{x}$  : observation  
 $\mathbf{z}$  : latent variable  
 $\theta$  : model parameter  
 $f_{\mathbf{z}|\mathbf{x}}$  : encoder

- By simulating the SDE, samples of  $\phi$  are collected
- MCMC over the latents is implicitly performed by collecting the mapping by  $f_{\mathbf{z}|\mathbf{x}}$

for  $t = 1, \dots, T$

$$\phi_{t+1} \sim \mathcal{N} \left( \phi_{t+1}; \phi_t - \eta \nabla_{\phi} V(\phi), 2\eta \mathbf{I} \right)$$

for  $i = 1, \dots, n$

$$\mathbf{z}_{t+1}^{(i)} = f_{\mathbf{z}|\mathbf{x}} \left( \mathbf{x}^{(i)}; \phi_{t+1} \right)$$

# Amortized Langevin Dynamics

## Theoretical Analysis

$\mathbf{x}$  : observation  
 $\mathbf{z}$  : latent variable  
 $\theta$  : model parameter  
 $f_{\mathbf{z}|\mathbf{x}}$  : encoder

- When the following conditions are satisfied, ALD has the true posterior as a stationary distribution
  1. Encoder takes the form of  $f_{\mathbf{z}|\mathbf{x}}(\mathbf{x}; \Phi) = \Phi g(\mathbf{x})$
  2. Rank of  $\mathbf{G}$  is  $n$ , where  $\mathbf{G}$  is a matrix with  $g(\mathbf{x}^{(i)})$  in row  $\mathbf{G}_{i,:}$
- **ALD is valid as MCMC under mild assumptions**

# Amortized Langevin Dynamics

## Remarks

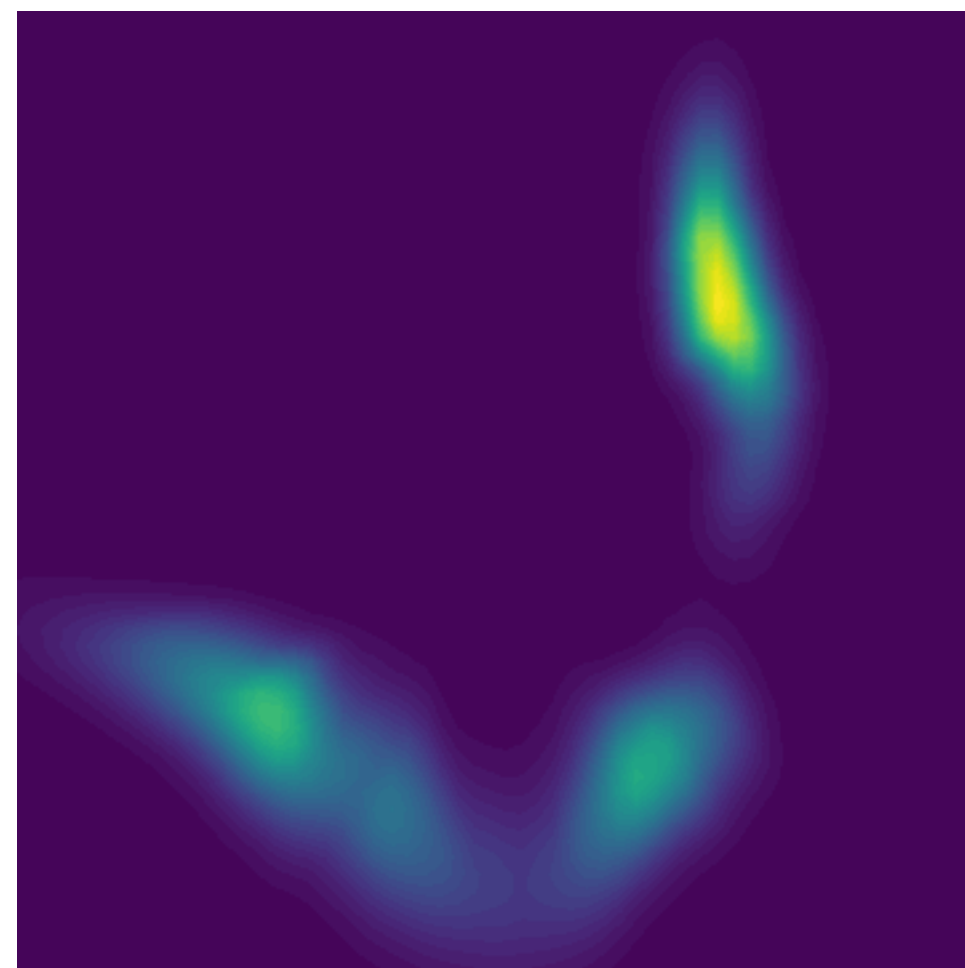
$\mathbf{x}$  : observation  
 $\mathbf{z}$  : latent variable  
 $\theta$  : model parameter  
 $f_{\mathbf{z}|\mathbf{x}}$  : encoder

1. ALD completely removes datapoint-wise iterations
  2. ALD is valid as an MCMC algorithm
  3. Encoder may accelerate the convergence of MCMC
- We name the learning algorithm of DLVMs using ALD the ***Langevin autoencoder***

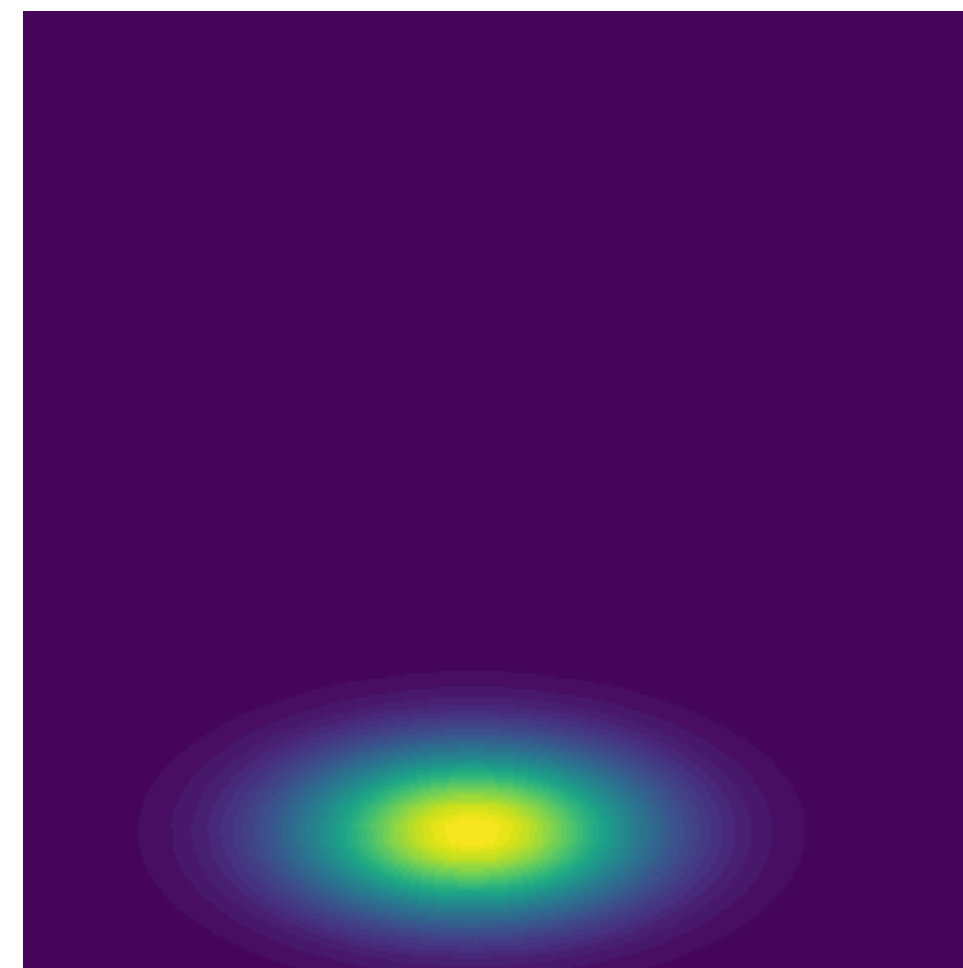
# Experiments

## Toy Example

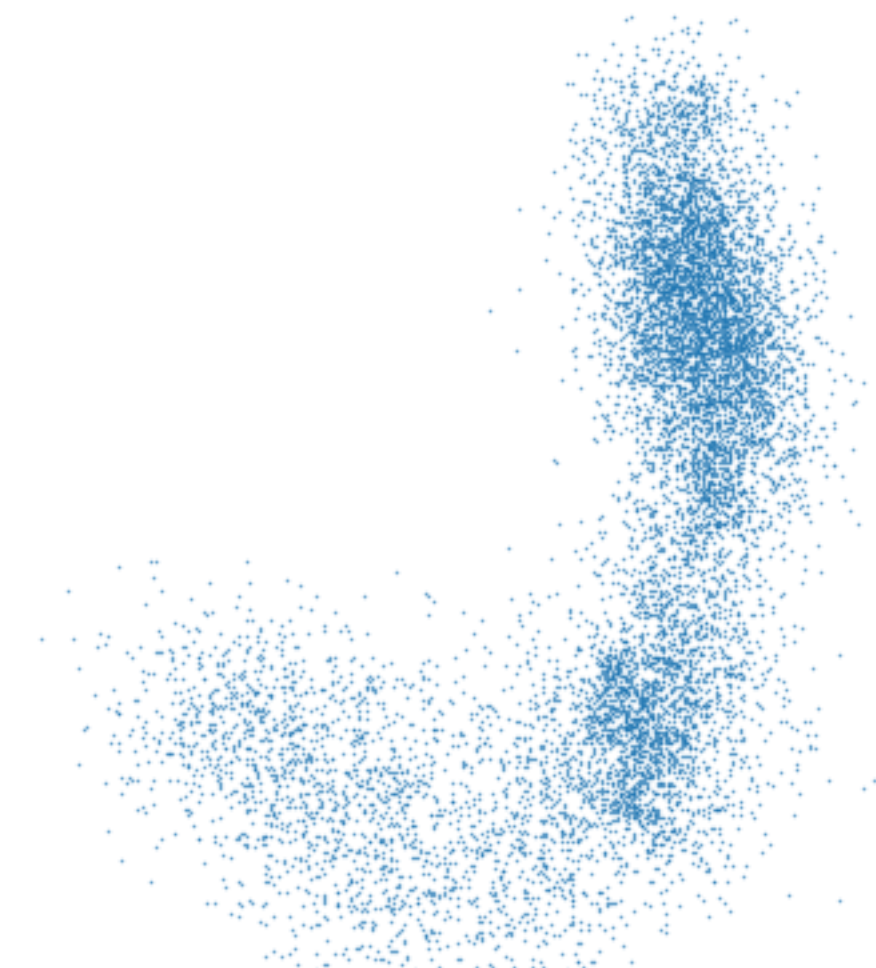
### Visualization of Posterior Approximation



Ground Truth



Variational Inference



ALD (ours)

Our ALD can accurately capture the multimodality

# Experiments

## Image Generation

### Test Likelihood Comparison

	MNIST	SVHN	CIFAR-10	CelebA
VAE	$1.189 \pm 0.002$	$4.442 \pm 0.003$	$4.820 \pm 0.005$	$4.671 \pm 0.001$
VAE-flow	$1.183 \pm 0.001$	$4.454 \pm 0.016$	$4.828 \pm 0.005$	$4.667 \pm 0.005$
Hoffman [2017]	$1.189 \pm 0.002$	$4.440 \pm 0.007$	$4.831 \pm 0.005$	$4.662 \pm 0.011$
LAE (ours)	<b><math>1.177 \pm 0.001</math></b>	<b><math>4.412 \pm 0.002</math></b>	<b><math>4.773 \pm 0.003</math></b>	<b><math>4.636 \pm 0.003</math></b>

Our LAE consistently outperforms VAE and existing LD-based method

# Future Works

- Extend to more sophisticated MCMC (e.g., Hamiltonian Monte Carlo)
- Completely remove the bias of gradient estimation using unbiased MCMC method

# Conclusion

- To train a deep latent variable model (DLVM), we need to approximate intractable posterior distribution
- Langevin dynamics (LD) can be a possible choice, but it is too slow due to datapoint-wise sampling process
- We propose **amortized Langevin dynamics** (ALD), which alleviate the problem by introducing an encoder
- ALD-based learning algorithm of DLVM named the **Langevin autoencoder** empirically outperforms existing methods (e.g., VAE and other LD-based method)