# Xingiang Ding<sup>1</sup> and David J. Freedman<sup>1,2</sup> Estimate the Gradient $\nabla_{\theta} \log p_{\theta}(x)$ , which is an Expectation in (1), Using Annealed Importance Sampling/Jarzynski Equality Annealed Importance Sampling Importance Sampling

# Learning Deep Generative Models with Annealed Importance Sampling <sup>1</sup> Department of Neurobiology, The University of Chicago, <sup>2</sup>The Grossman Institute for Neuroscience, Quantitative Biology and Human Behavior, The University of Chicago Learning Deep Generative Models via Maximizing Likelihood **Requires Computing the Gradient of the Log-Likelihood** Given a generative model defined by a joint distribution of observed data variables **x**

and continuous latent variable *z*:

$$p_{\theta}(x,z) = p_{\theta}(z)p_{\theta}(x|z)$$

With training data x, we are interested in learning  $\theta$  by maximizing its marginal likelihood:

$$\log p_{\theta}(x) = \log \int p_{\theta}(x, z) \mathrm{d}$$

Computing log  $p_{\theta}(x)$  is expensive because it involves a high dimensional integral. However, we note that maximizing log  $p_{\theta}(x)$  does not necessarily require computing the value of log  $p_{\theta}(x)$ . What is necessarily required is the gradient of log  $p_{\theta}(x)$  with respect to  $\theta$ :

$$\nabla_{\theta} \log p_{\theta}(x) = \mathbb{E}_{z \sim p_{\theta}(z|x)} \left[ \nabla_{\theta} \log p_{\theta} \right]$$

which can be written as an expectation.

### Several Previous Methods Used for Training Generative Models Can Be Viewed as Devising Ways of Approximating the Expectation in (1)

Variational Auto-Encoder (VAE) In maximizing the ELBO:

$$\mathcal{C}(\theta, \phi) = \mathbb{E}_{z \sim q_{\phi}(z|x)} \left[ \log p_{\theta}(x, z) - \log p_{\theta}(x, z) - \log p_{\theta}(x, z) \right]$$

 $\theta$  follows the gradient:

 $\nabla_{\theta} \mathcal{L}(\theta, \phi) = \mathbb{E}_{z \sim q_{\phi}(z|x)} \left[ \nabla_{\theta} \log p_{\theta}(x, z) \right]$ 

which is an approximation to the expectation in (1).

Importance Weighted Auto-Encoder (IWAE)

$$\mathcal{L}_{K}(\theta,\phi) = \mathbb{E}_{z_{1},...,z_{K}\sim q_{\phi}(z|x)} \left[ \log \frac{1}{K} \sum_{k=1}^{K} \frac{p_{\theta}(x,z_{k})}{q_{\phi}(z_{k}|x)} \right]$$
$$\nabla_{\theta}\mathcal{L}_{K}(\theta,\phi) = \sum_{k=1}^{K} \widetilde{w_{k}} \nabla_{\theta} \log p_{\theta}(x,z_{k}) \qquad (2)$$
$$\widetilde{w_{k}} = w_{k} / \sum_{i=1}^{K} w_{i} \qquad w_{k} = w(x,z_{k},\theta,\phi) = \frac{p_{\theta}(x,z_{k})}{q_{\phi}(z_{k}|x)}$$

Markov Chain Monte Carlo (MCMC) Methods Directly estimate the expectation in (1) by drawing samples from the posterior distribution  $p_{\theta}(z|x)$  with MCMC. It takes a large number of steps for a Markov chain to converge, so it is much slower than variational methods.



## Learn Deep Generative Models with Annealed Importance Sampling

glot			MNIST	
5	K=50	K = 1	K=5	K=50
11 47 94	-103.91 -102.03 -101.64	-86.90 -84.56 -84.14	-85.52 -84.25 -83.78	-84.38 -83.93 -83.63

	$\approx \log p(x)$		$\approx \log p(x)$
MH-HMC (L = 5, K = 5, T = 5) Ours (L = 5, K = 5, T = 5)	-102.57 -102.47	IWAE-DReG (K = 55)	-103.85
MH-HMC (L = 5, K = 5, T = 11) Ours (L = 5, K = 5, T = 11)	-101.32 -101.94	Ours (L = 5, K = 1, T = 11) IWAE-DReG (K = 275)	-102.45
MH-HMC (L = 5, K = 50, T = 5) Ours (L = 5, K = 50, T = 5)	-102.32 -102.03	Ours (L = 5, K = 5, T = 11)	-101.94
MH-HMC (L = 5, K = 50, T = 11) Ours (L = 5, K = 50, T = 11)	-101.25 -101.64	IWAE-DReG (K = 2750) Ours (L = 5, K = 50, T = 11)	-102.40 -101.64

$q(z x)^{1-eta}p_{ heta}(z x)^{eta} \qquad p_{ heta}(z x)$	
$q_{\phi}(z x)$	
m $q(z x)^{1-\beta}p_{\theta}(z x)^{\beta}$ by	
Hamiltonian Monte Carlo (HMC)	
nergy of $U^{\beta}(z) = -\log q(z x)^{1-\beta} p_{\theta}(z x)^{\beta}$	в
$p_{ heta}(z_k^0 x)^eta \qquad p_{ heta}(z_k^1 x)$	
$\frac{\partial q}{\partial x}(x) = \frac{1}{q(z_k^1 x)^{1-\beta}p_{\theta}(z_k^1 x)^{\beta}}$	
	-
ompling $a_{-x}(x x)$ or $a_{-x}(x x)$	

K: number of importance samples L: number of integration steps in HMC T: num of intermediate distributions **case 1**: when importance sampling is used sample  $z_1, ..., z_K \sim q_\phi(z|x)$ set  $w_k = \frac{p_\theta(x, z_k)}{q_\phi(z_k | x)}$  and  $\widetilde{w_k} = w_k / \sum_{k=1}^K w_k$ sample  $z_1, ..., z_K$  and compute  $\widetilde{w_1}, ..., \widetilde{w_K}$ sample  $j \sim Categorical(\widetilde{w_1}, ..., \widetilde{w_K})$ 

### Results of IWAE-DReG, MH-HMC and our approach on the Omniglot dataset with same computational cost.