
Appendix for “Adversarial Symmetric Variational Autoencoder”

**Yunchen Pu, Weiyao Wang, Ricardo Henao, Liqun Chen, Zhe Gan, Chunyuan Li
and Lawrence Carin**

Department of Electrical and Computer Engineering, Duke University
{yp42, ww109, r.henao, lc267, zg27, cl319, lcarin}@duke.edu

A Proof

Proof of Corollary 1.1 We start from a simple observation $p_\theta(\mathbf{x}) = \int_{\mathbf{z}} p_\theta(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int_{\mathbf{z}} p(\mathbf{z}) p_\theta(\mathbf{x}|\mathbf{z}) d\mathbf{z}$. The second term in (5) of main paper can be rewritten as

$$\mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x}|\mathbf{z}'), \mathbf{z}' \sim p(\mathbf{z}), \mathbf{z} \sim p(\mathbf{z})} [\log(1 - \sigma(f_{\psi_1}(\mathbf{x}, \mathbf{z})))] , \quad (1)$$

$$= \int_{\mathbf{x}} \int_{\mathbf{z}'} \int_{\mathbf{z}} p_\theta(\mathbf{x}|\mathbf{z}') p(\mathbf{z}') p(\mathbf{z}) \log(1 - \sigma(f_{\psi_1}(\mathbf{x}, \mathbf{z}))) d\mathbf{x} d\mathbf{z} d\mathbf{z}' \quad (2)$$

$$= \int_{\mathbf{x}} \int_{\mathbf{z}} \left\{ \int_{\mathbf{z}'} p_\theta(\mathbf{x}|\mathbf{z}') p(\mathbf{z}') d\mathbf{z}' \right\} p(\mathbf{z}) \log(1 - \sigma(f_{\psi_1}(\mathbf{x}, \mathbf{z}))) d\mathbf{x} d\mathbf{z} \quad (3)$$

$$= \int_{\mathbf{x}} \int_{\mathbf{z}} p_\theta(\mathbf{x}) p(\mathbf{z}) \log(1 - \sigma(f_{\psi_1}(\mathbf{x}, \mathbf{z}))) d\mathbf{x} d\mathbf{z} \quad (4)$$

Therefore, the objective function $\mathcal{L}_{A1}(\psi_1)$ in (5) can be expressed as

$$\begin{aligned} & \int_{\mathbf{x}} \int_{\mathbf{z}} q(\mathbf{x}) q_\phi(\mathbf{z}|\mathbf{x}) \log[\sigma(f_{\psi_1}(\mathbf{x}, \mathbf{z}))] d\mathbf{x} d\mathbf{z} + \int_{\mathbf{x}} \int_{\mathbf{z}} p_\theta(\mathbf{x}) p(\mathbf{z}) \log(1 - \sigma(f_{\psi_1}(\mathbf{x}, \mathbf{z}))) d\mathbf{x} d\mathbf{z} \\ &= \int_{\mathbf{x}} \int_{\mathbf{z}} \{ q_\phi(\mathbf{x}, \mathbf{z}) \log[\sigma(f_{\psi_1}(\mathbf{x}, \mathbf{z}))] + p_\theta(\mathbf{x}) p(\mathbf{z}) \log(1 - \sigma(f_{\psi_1}(\mathbf{x}, \mathbf{z}))) \} d\mathbf{x} d\mathbf{z} \end{aligned} \quad (5)$$

This integral of (5) is maximal as a function of $f(\mathbf{x}, \mathbf{z})$ if and only if the integrand is maximal for every (\mathbf{x}, \mathbf{z}) . Note that the problem $\max_x a \log x + b \log(1 - x)$ achieves maximum at $x = a/(a+b)$ and $\sigma(x) = 1/(1 + e^{-x})$. Hence, we have the optimal function of $f_{\psi_1^*}$ at

$$\sigma(f_{\psi_1^*}) = \frac{q_\phi(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{x}, \mathbf{z}) + p_\theta(\mathbf{x}) p(\mathbf{z})} \quad f_{\psi_1^*} = \log q_\phi(\mathbf{x}, \mathbf{z}) + \log p_\theta(\mathbf{x}) p(\mathbf{z}) \quad (6)$$

Similarly, we have $f_{\psi_2^*}(\mathbf{x}, \mathbf{z}) = \log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}) q(\mathbf{x})$

Proof of Proposition 1 If $\{\theta^*, \phi^*, \psi_1^*, \psi_2^*\}$ achieves an equilibrium of (12) of main paper. The Corollary 1.1 indicates that $f_{\psi_1^*} = \log q_\phi(\mathbf{x}, \mathbf{z}) + \log p_\theta(\mathbf{x}) p(\mathbf{z})$ and $f_{\psi_2^*}(\mathbf{x}, \mathbf{z}) = \log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}) q(\mathbf{x})$.

Note that

$$\mathcal{L}_{\text{VAEx}}(\theta, \phi) = \mathbb{E}_{q(\mathbf{x})} \log p_\theta(\mathbf{x}) - \text{KL}(q_\phi(\mathbf{x}, \mathbf{z}) \| p_\theta(\mathbf{x}, \mathbf{z})) \quad (7)$$

$$= \mathbb{E}_{q(\mathbf{x})} \log q(\mathbf{x}) - \text{KL}(q_\phi(\mathbf{x}, \mathbf{z}) \| p_\theta(\mathbf{x}, \mathbf{z})) - \text{KL}(q_\phi(\mathbf{x}) \| p_\theta(\mathbf{x})) \quad (8)$$

and

$$\mathcal{L}_{\text{VAEz}}(\theta, \phi) = \mathbb{E}_{p(\mathbf{z})} \log q_\phi(\mathbf{z}) - \text{KL}(p_\theta(\mathbf{x}, \mathbf{z}) \| q_\phi(\mathbf{x}, \mathbf{z})) \quad (9)$$

$$= \mathbb{E}_{p(\mathbf{z})} \log p(\mathbf{z}) - \text{KL}(p_\theta(\mathbf{x}, \mathbf{z}) \| q_\phi(\mathbf{x}, \mathbf{z})) - \text{KL}(p_\theta(\mathbf{z}) \| q_\phi(\mathbf{z})) \quad (10)$$

where $\mathbb{E}_{p(z)} \log p(z)$ and $\mathbb{E}_{q(x)} \log q(x)$ can be considered as constant. Therefore, maximize $\mathcal{L}_{\text{VAE}xz}$ is equivalent to minimize

$$\text{KL}(p_\theta(x, z) \| q_\phi(x, z)) + \text{KL}(q_\phi(x, z) \| p_\theta(x, z)) + \text{KL}(p_\theta(z) \| q_\phi(z)) + \text{KL}(q_\phi(x) \| p_\theta(x))$$

The minimum of first two terms is achieved if and only if $p_\theta(x, z) = q_\phi(x, z)$ while the minimums of last two terms are achieved at $p_\theta(x) = q_\phi(x)$ and $p(z) = q_\phi(z)$, respectively. Note that the joint match $p_\theta(x, z) = q_\phi(x, z)$ is achieved, the marginals also matches which indicates the optimal (θ^*, ϕ^*) is achieved if and only if $p_{\theta^*}(x, z) = q_{\phi^*}(x, z)$.

B Model Architecture

The model architectures are shown as following. For $f_{\psi_1}(x, z)$ and $f_{\psi_2}(x, z)$, we use the same architecture but the parameters are not shared.

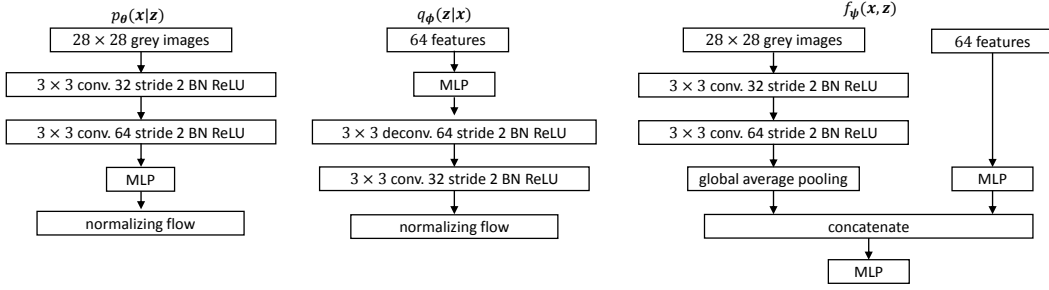


Figure 1: Model architecture for MNIST

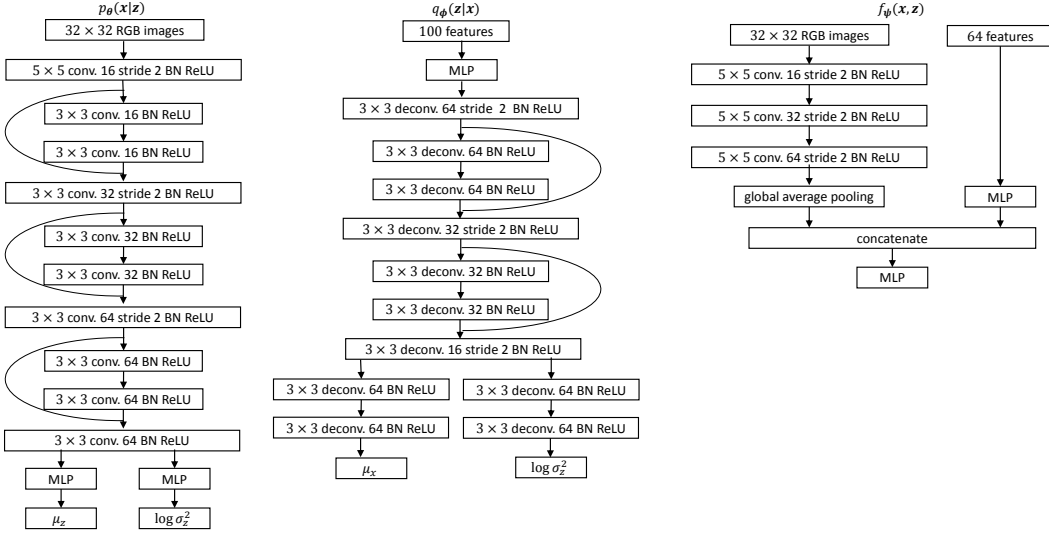


Figure 2: Model architecture for CIFAR

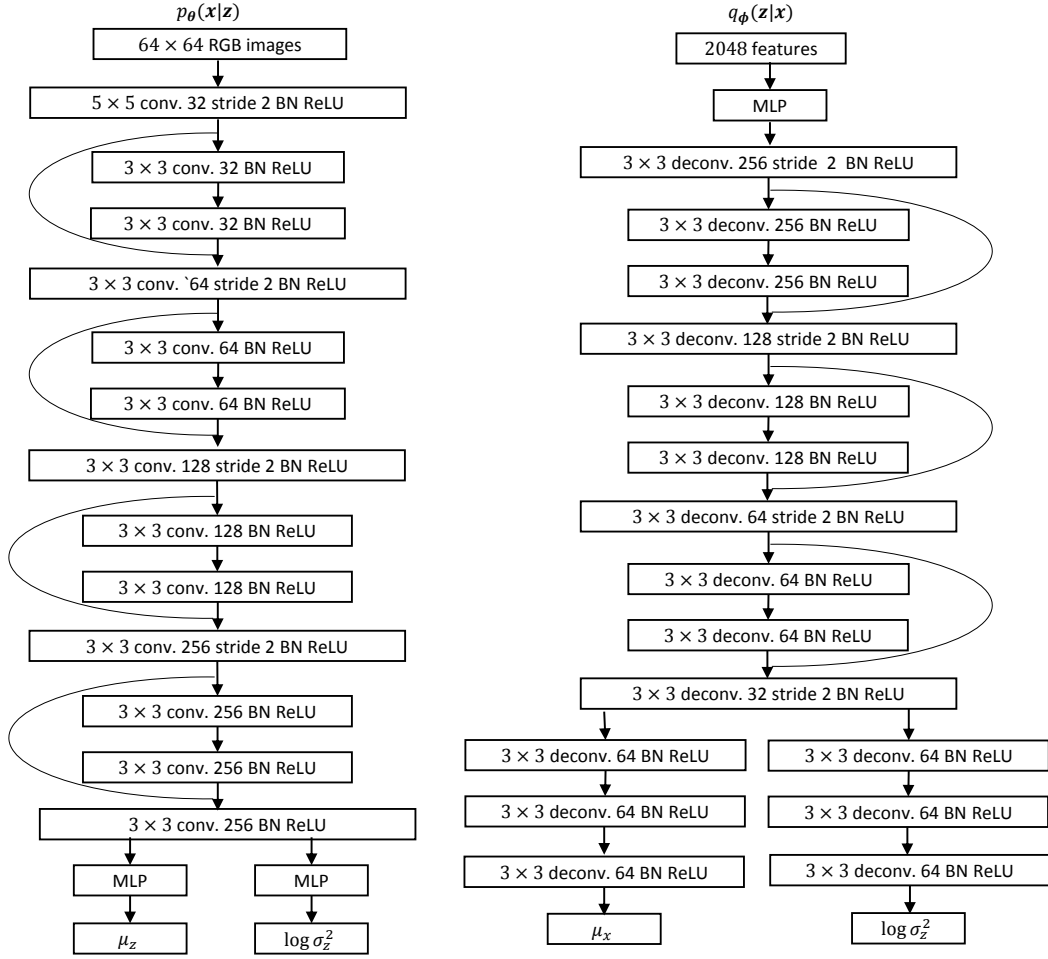


Figure 3: Encoder and decoder for ImageNet

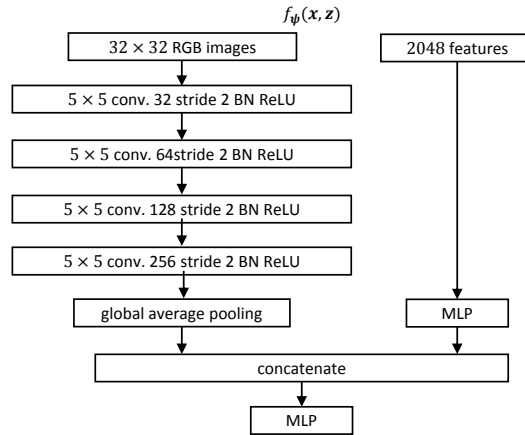


Figure 4: Discriminator for ImageNet

C Additional Results



Figure 5: Generated samples trained on CIFAR-10.



Figure 6: Generated samples trained on ImageNet.