

[80245013 Machine Learning, Fall, 2020]

# **Deep Generative Models**

**Jun Zhu**

dcszj@mail.tsinghua.edu.cn

<http://ml.cs.tsinghua.edu.cn/~jun>

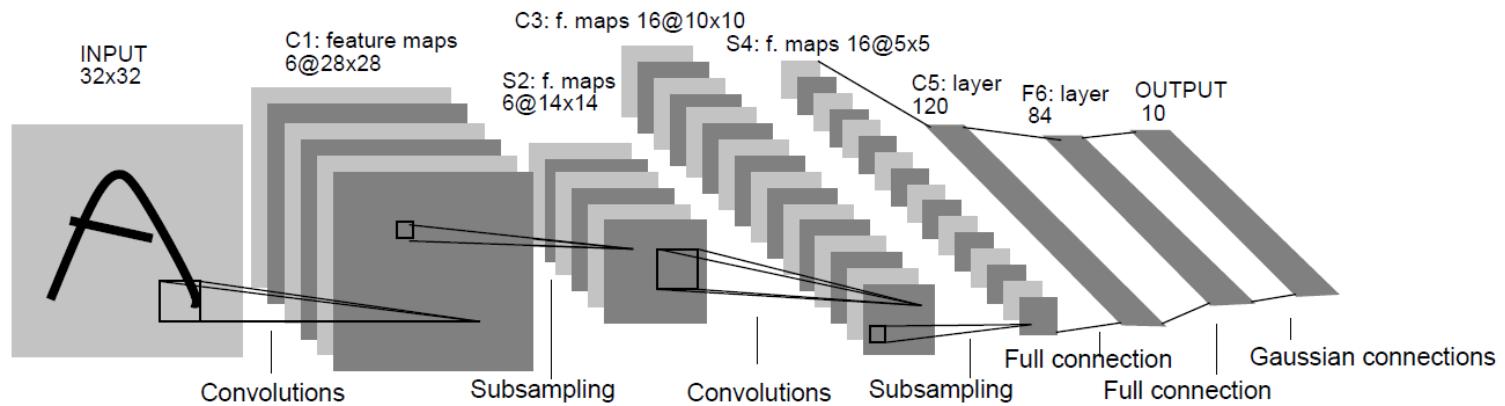
State Key Lab of Intelligent Technology & Systems

Tsinghua University

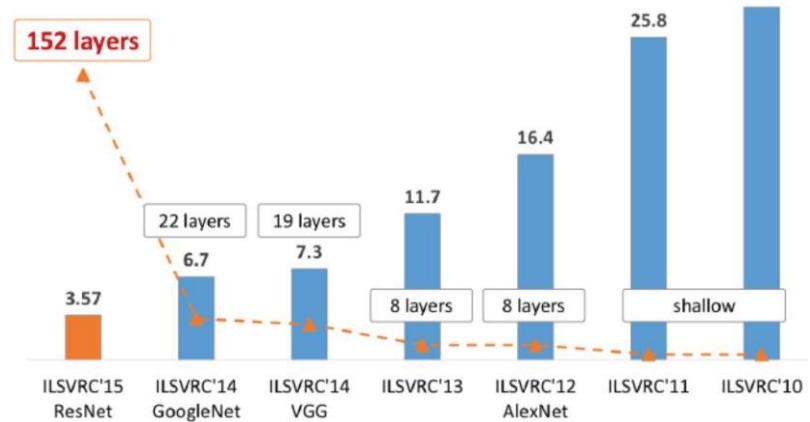
November 3, 2020

# Discriminative Deep Learning

- ◆ Learn a deep NN to map an input to output



- Gradient back-propagation
- Dropout
- Activation functions:
  - rectified linear



# Generative Modeling

- ◆ Have training examples

$$\mathbf{x} \sim p_{\text{data}}(\mathbf{x})$$

- ◆ Want a model that can draw samples:

$$\mathbf{x}' \sim p_{\text{model}}(\mathbf{x})$$

- where  $p_{\text{model}}(\mathbf{x}) \approx p_{\text{data}}(\mathbf{x})$

7	2	1	0	4	1	4	9	5	9
0	6	9	0	1	5	9	7	3	4
9	6	4	5	4	0	7	4	0	1
3	1	3	4	7	2	7	1	2	1
1	7	4	2	3	5	1	2	4	4
6	3	5	5	6	0	4	1	9	5
7	8	9	3	7	4	0	4	3	0
7	0	2	9	1	7	3	2	9	7
7	6	2	7	8	4	7	3	6	1
3	6	9	3	1	4	1	7	6	9

$$\mathbf{x} \sim p_{\text{data}}(\mathbf{x})$$



7	2	1	0	4	9	4	9	5	9
0	6	9	0	1	5	9	7	3	4
9	6	4	5	4	0	7	4	0	1
3	9	3	4	7	2	7	1	2	1
1	7	4	2	3	8	1	2	9	4
6	3	5	5	6	0	4	1	9	9
7	8	9	3	7	9	0	4	8	0
7	0	2	9	1	7	3	2	9	7
9	6	2	9	5	4	7	3	6	1
8	6	7	3	1	4	1	7	6	9

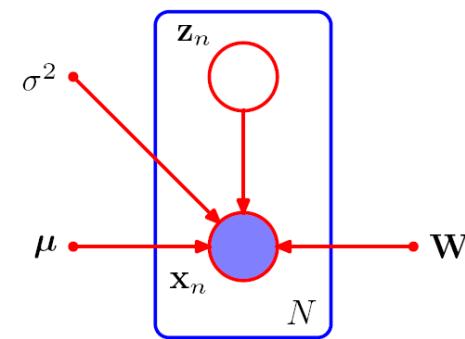
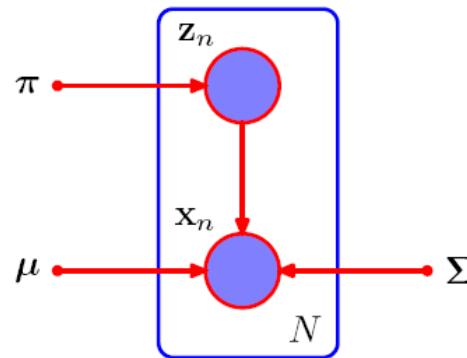
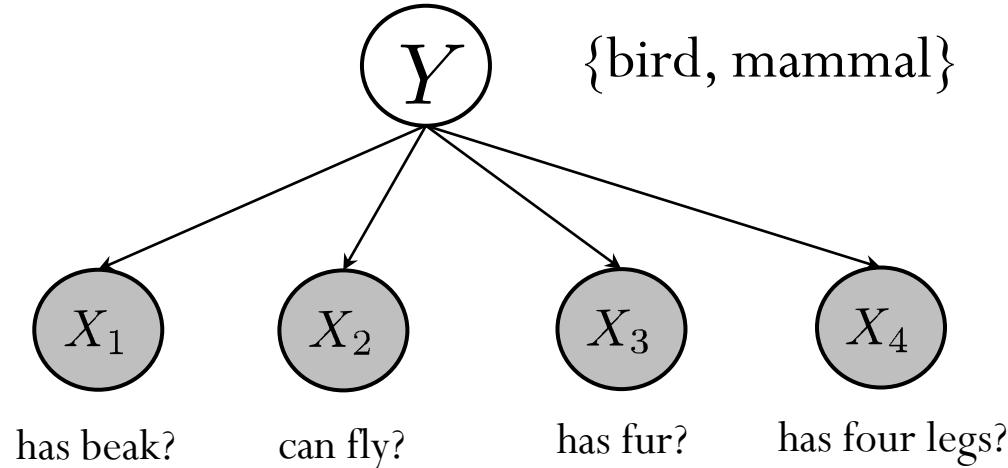
$$\mathbf{x}' \sim p_{\text{model}}(\mathbf{x})$$

# Why generative models?

- ◆ Leverage unlabeled datasets, which are often much larger than labeled ones
  - Unsupervised learning
    - e.g., clustering, density estimation, feature extraction, dimension reduction, data generation
  - Semi-supervised learning
    - e.g., classification, information extraction, learning-to-rank, network analysis, opinion mining
- ◆ Conditional generative models
  - Speech synthesis: Text  $\Rightarrow$  Speech
  - Machine Translation: French  $\Rightarrow$  English
  - Image captioning: Image  $\Rightarrow$  Text

# Generative models are your old friends

- ◆ Naïve Bayes, GMM, Factor Analysis (Probabilistic PCA)



# Deep Generative Models

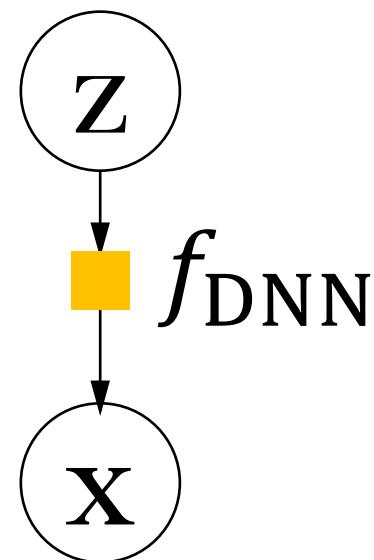
- ◆ More flexible by using differential function mapping between random variables
- ◆ If  $z$  is uniformly distributed over  $(0, 1)$ , then  $y = f(z)$  has the distribution

$$p(y) = p(z) \left| \frac{dz}{dy} \right|$$

- where  $p(z) = 1$
- ◆ This trick is widely used to draw samples from exponential family distributions (e.g., Gaussian, Exponential)

# Deep Generative Models

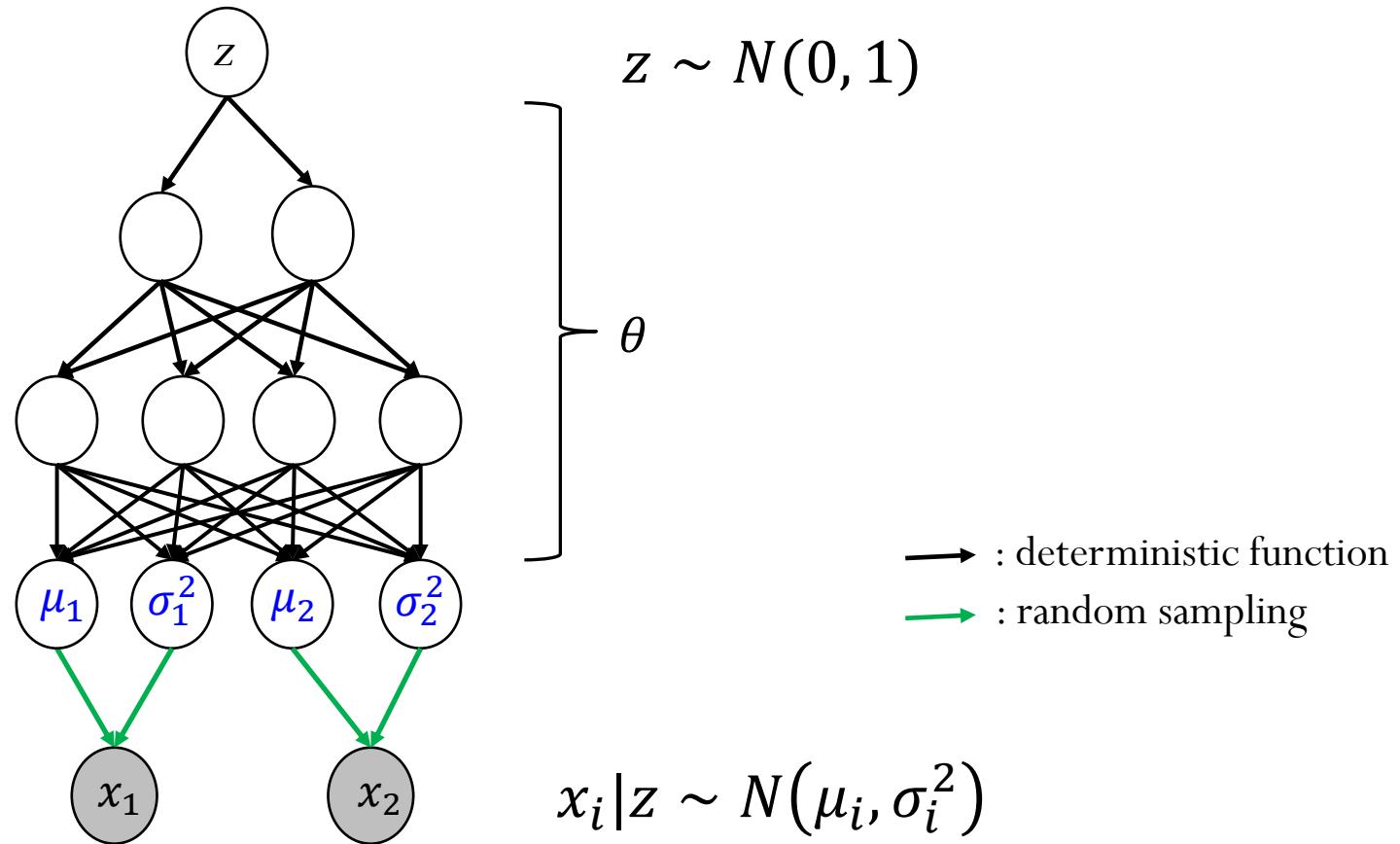
- ◆ More flexible by using differential function mapping between random variables
- ◆ DGMs learn a function transform with deep neural networks



Cause => Disease  
Topics => Docs  
Objects => Images  
Words => Phonemes

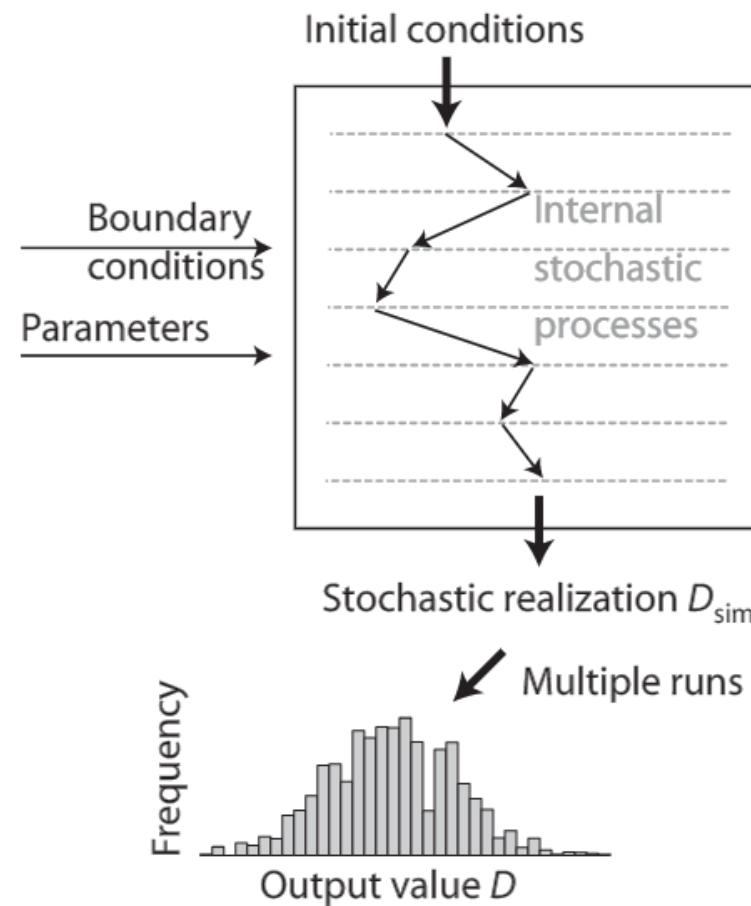
# An example with MLP

- ◆ 1D latent variable  $z$ ; 2D observation  $x$
- ◆ Idea: NN + Gaussian (or Bernoulli) with a diagonal covariance



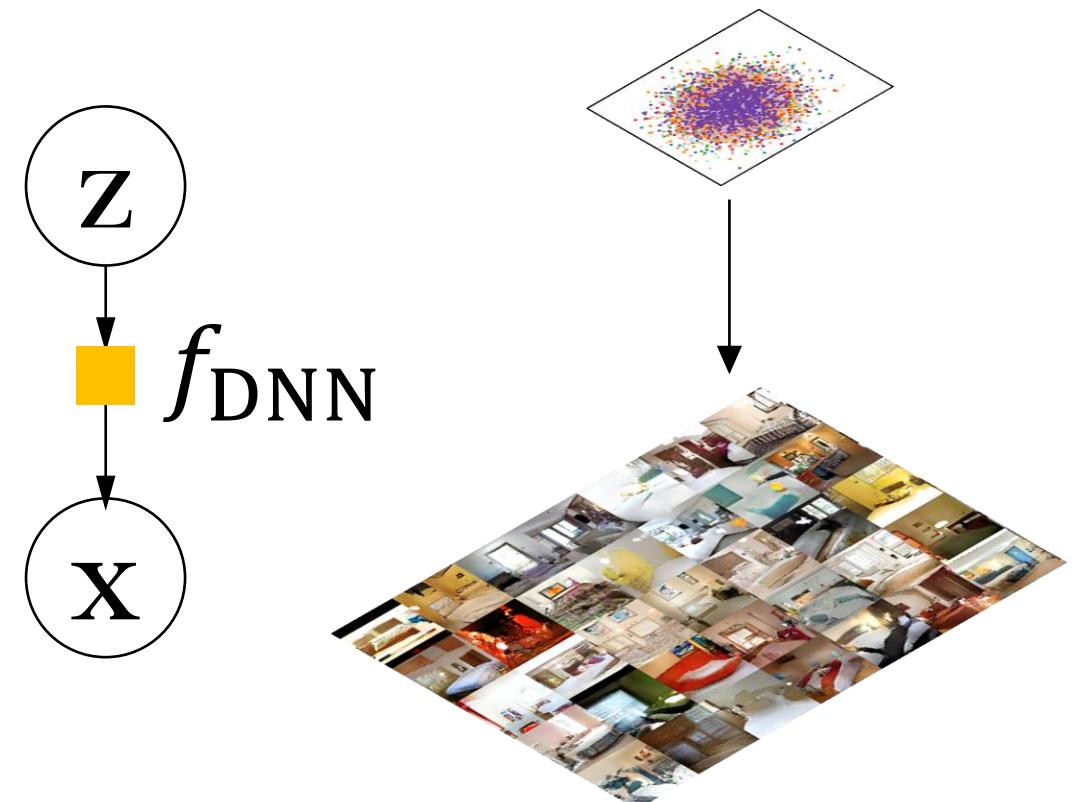
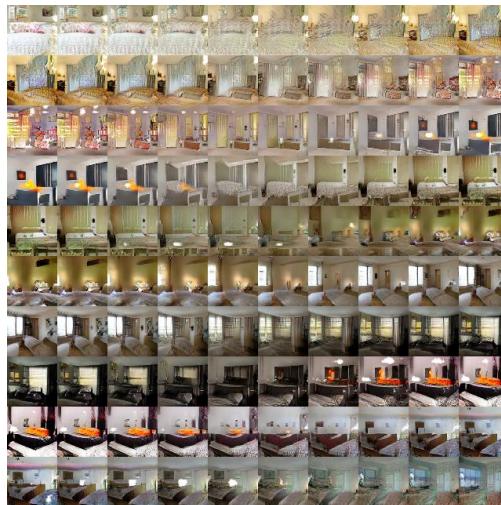
# Implicit Deep Generative Models

- ◆ Generate data with a stochastic process whose likelihood function is not explicitly specified (Hartig et al., 2011)

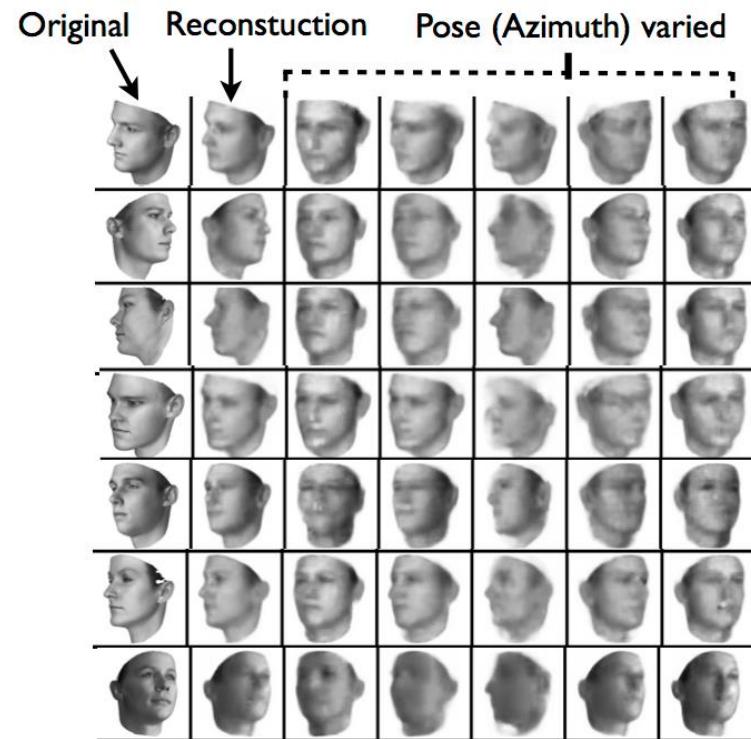
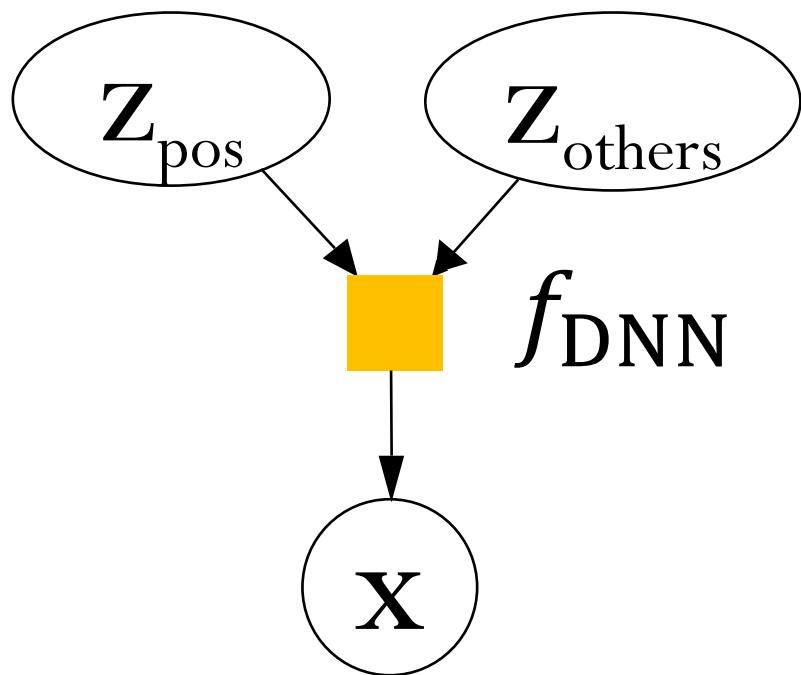


# Deep Generative Models

[Image Generation:  
Generative Adversarial Nets,  
Goodfellow13 & Radford15]



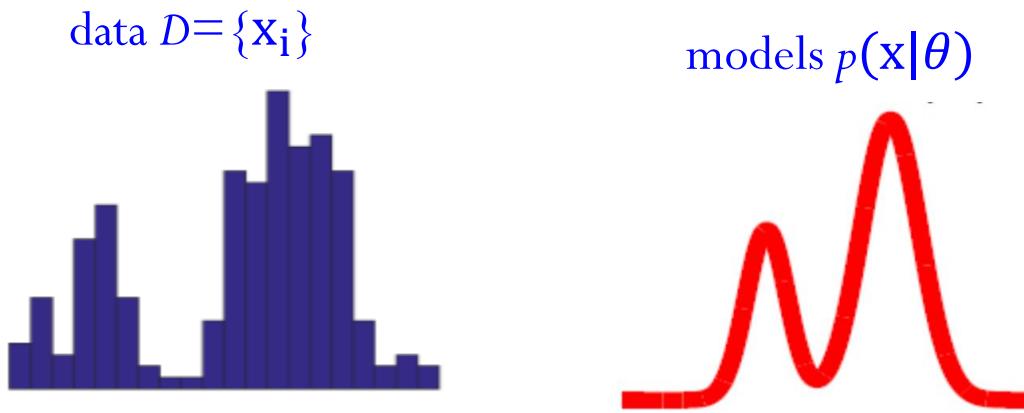
# Deep Generative Models



[Image Understanding: Variational Autoencoders,  
Kingma13 & Tejas15 & Eslami16]

# Learning Deep Generative Models

- ◆ Given a set  $D$  of unlabeled samples, learn the unknown parameters (or a distribution)



- ◆ Find a model that minimizes

$$\mathbb{D}(\text{data } \{x_i\}_{i=1}^n, \text{model } p)$$

# Learning Deep Generative Models

- ◆ Maximum likelihood estimation (MLE):

$$\hat{\theta} = \operatorname{argmax} p(D|\theta)$$

- has an explicit likelihood model
- ◆ Minimax objective (e.g., GAN)
  - a two-player game to reach equilibrium

- ◆ Moment-matching:

- draw samples from  $p$ :  $\widehat{D} = \{y_i\}_{i=1}^M$ , where  $y_i \sim p(x|\theta)$

- Kernel MMD:
$$\mathcal{L}_{MMD^2} = \left\| \frac{1}{N} \sum_{i=1}^N \phi(x_i) - \frac{1}{M} \sum_{j=1}^M \phi(y_j) \right\|_{\mathcal{H}}^2$$

- rich enough to distinguish any two distributions in certain RKHS

# Variational Bayes

- ◆ Consider the log-likelihood of *a single example*

$$\log p(\mathbf{x}; \theta) = \log \int p(\mathbf{z}, \mathbf{x}; \theta) d\mathbf{z}$$

- ◆ Log-integral/sum is annoying to handle directly
- ◆ Derive a variational lower bound  $L(\theta, \phi, \mathbf{x})$

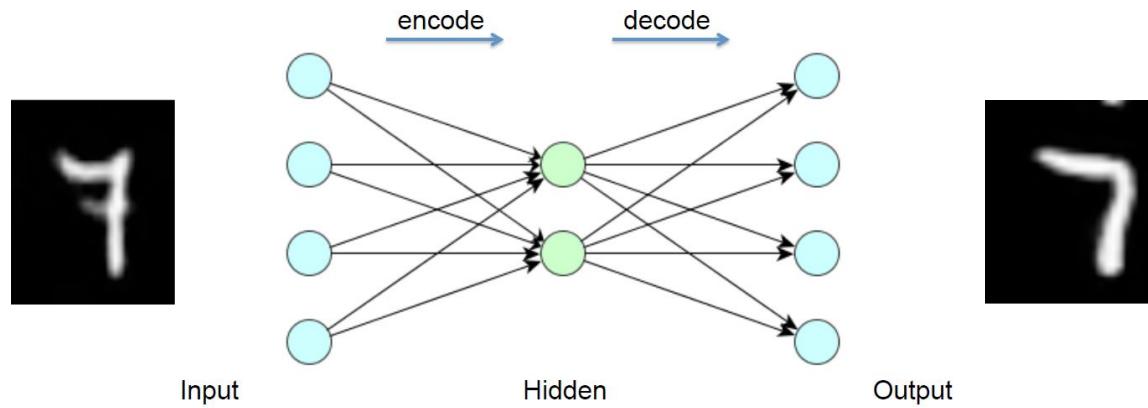
$$\log p(\mathbf{x}; \theta) = L(\theta, \phi, \mathbf{x}) + \text{KL}(q(\mathbf{z}|\mathbf{x}; \phi) \| p(\mathbf{z}|\mathbf{x}; \theta))$$

$$\begin{aligned} L(\theta, \phi, \mathbf{x}) &= \mathbf{E}_{q(\mathbf{z}|\mathbf{x}; \phi)} [\log p(\mathbf{z}, \mathbf{x}; \theta) - \log q(\mathbf{z}|\mathbf{x}; \phi)] \\ &= \mathbf{E}_{q(\mathbf{z}|\mathbf{x}; \phi)} [\log p(\mathbf{x}|\mathbf{z}; \theta) + \log p(\mathbf{z}; \theta) - \log q(\mathbf{z}|\mathbf{x}; \phi)] \\ &= \mathbf{E}_{q(\mathbf{z}|\mathbf{x}; \phi)} [\log p(\mathbf{x}|\mathbf{z}; \theta)] - \text{KL}(q(\mathbf{z}|\mathbf{x}; \phi) \| p(\mathbf{z}; \theta)) \end{aligned}$$

reconstruction term

prior regularization

# Recap: Auto-Encoder

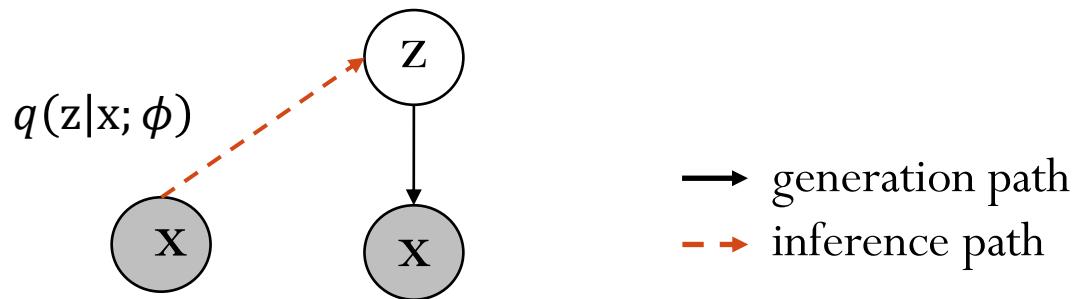


- ◆ Encoder:  $\mathbf{h} = s(W\mathbf{x} + b)$
- ◆ Decoder:  $\mathbf{x}' = s(W'\mathbf{h} + b')$
- ◆ Training: minimize the reconstruction error (e.g., square loss, cross-entropy loss)
- ◆ Denoising AE: randomly corrupted inputs are restored to learn more robust features

# Auto-Encoding Variational Bayes (AEVB)

- ◆ What's unique in AEVB is that *the variational distribution is parameterized by a deep neural network*

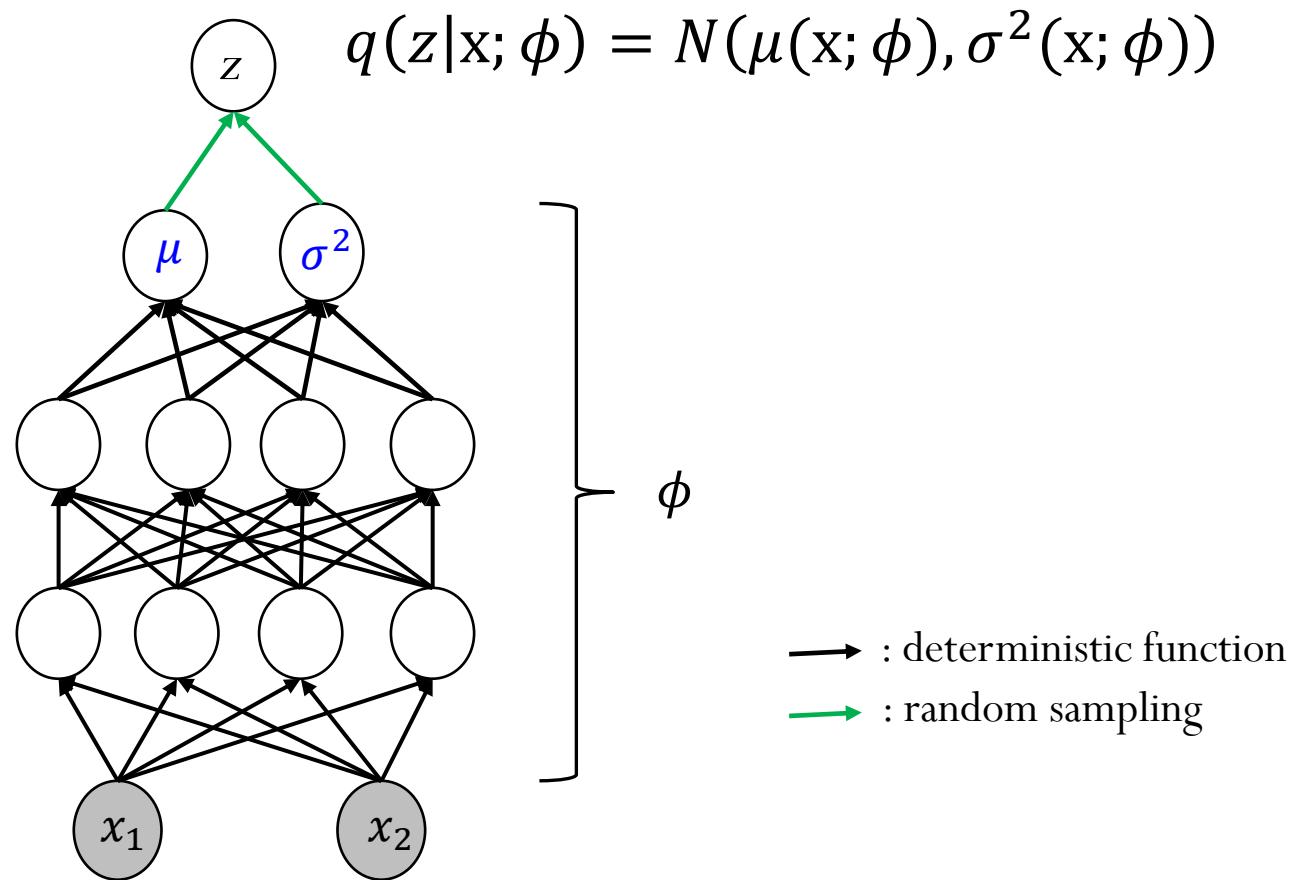
$$q(z|x; \phi) \approx p(z|x; \theta)$$



- We call it an **inference (recognition, encoder) network** or a **Q-network**
- All the parameters are learned jointly via SGD with variance reduction

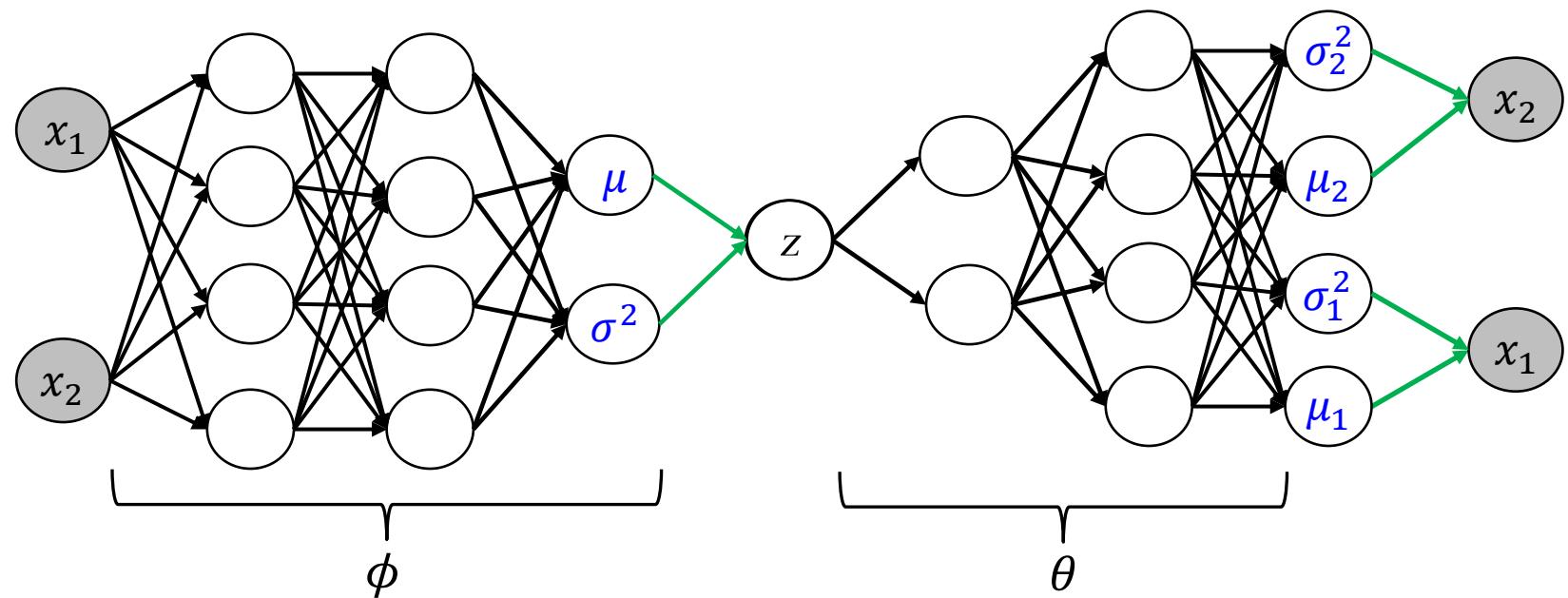
# The Encoder Network

- ◆ A feedforward NN + Gaussian



# The Complete Auto-encoder

- ◆ The Q-P network architecture:



→ : deterministic function

→ : random sampling

# Stochastic Variational Inference

- ◆ Variational lower-bound for *a single example*

$$L(\theta, \phi, \mathbf{x}) = \mathbf{E}_{q(\mathbf{z}|\mathbf{x};\phi)}[\log p(\mathbf{x}|\mathbf{z};\theta)] - \text{KL}(q(\mathbf{z}|\mathbf{x};\phi) \| p(\mathbf{z};\theta))$$

- ◆ Variational lower-bound for *a set of examples*

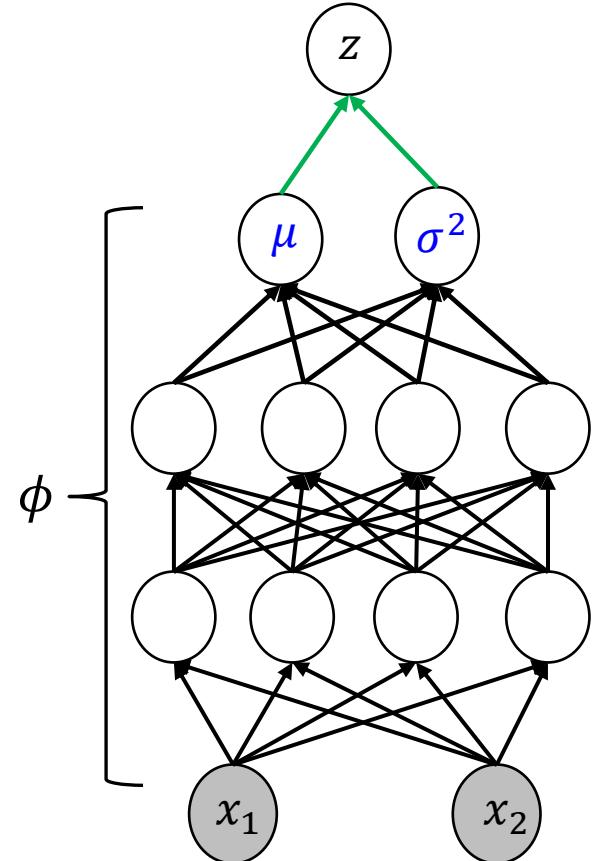
$$L(\theta, \phi, D) = \sum_i \mathbf{E}_{q(\mathbf{z}_i|\mathbf{x}_i;\phi)}[\log p(\mathbf{x}_i|\mathbf{z}_i;\theta)] - \text{KL}(q(\mathbf{z}_i|\mathbf{x}_i;\phi) \| p(\mathbf{z}_i;\theta))$$

- Use stochastic gradient methods to handle large datasets
- Random mini-batch
  - for each  $i$ , infer the posterior  $q(\mathbf{z}_i|\mathbf{x}_i;\phi)$ ; As we parameterize as a neural network, this in fact optimizes  $\phi$
- ◆ *However, calculating the expectation and its gradients is non-trivial, often intractable*

# Example with Gaussian Distributions

- ◆ Use  $N(0, 1)$  as prior for  $z$ ;  $q(z|x; \phi)$  is Gaussian with parameters  $(\mu(x; \phi), \sigma^2(x; \phi))$  determined by NN
  - The KL-divergence

$$-\text{KL}(q(z|x; \phi) \| p(z; \theta)) = \frac{1}{2} (1 + \log \sigma^2 - \mu^2 - \sigma^2)$$



# Example with Gaussian Distributions

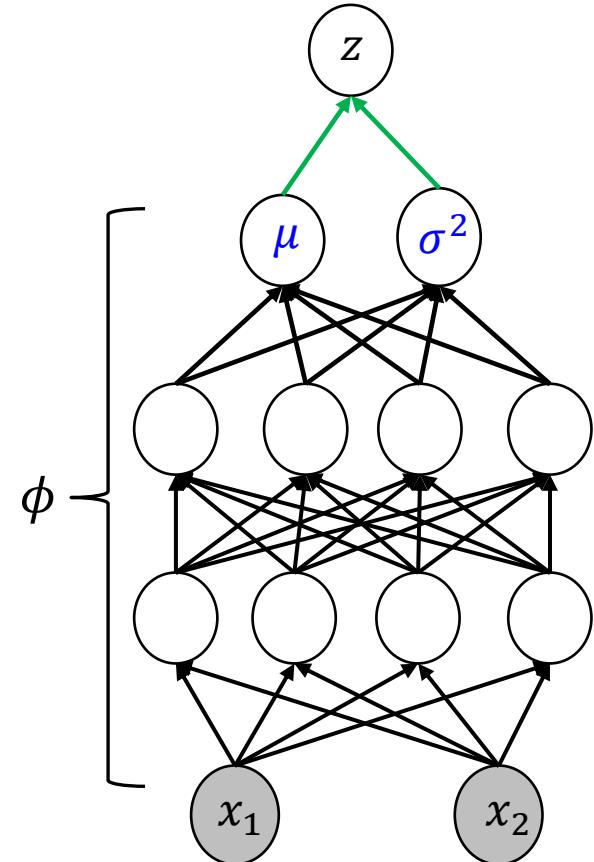
- ◆ Use  $N(0, 1)$  as prior for  $z$ ;  $q(z|x; \phi)$  is Gaussian with parameters  $(\mu(x; \phi), \sigma^2(x; \phi))$  determined by NN
  - The expected log-likelihood

$$\mathbf{E}_{q(z|x;\phi)} [\log p(x|z; \theta)]$$

- If the likelihood is Gaussian

$$-\log p(x_i|z_i) = \sum_j \frac{1}{2} \log \sigma_j^2 + \frac{(x_{ij} - \mu_{xi})^2}{2\sigma_j^2}$$

- *The expectation is still hard to compute because of nonlinearity functions*



# Example with Gaussian Distributions

- ◆ Use  $N(0, 1)$  as prior for  $z$ ;  $q(z|x; \phi)$  is Gaussian with parameters  $(\mu(x; \phi), \sigma^2(x; \phi))$  determined by NN
  - The expected log-likelihood

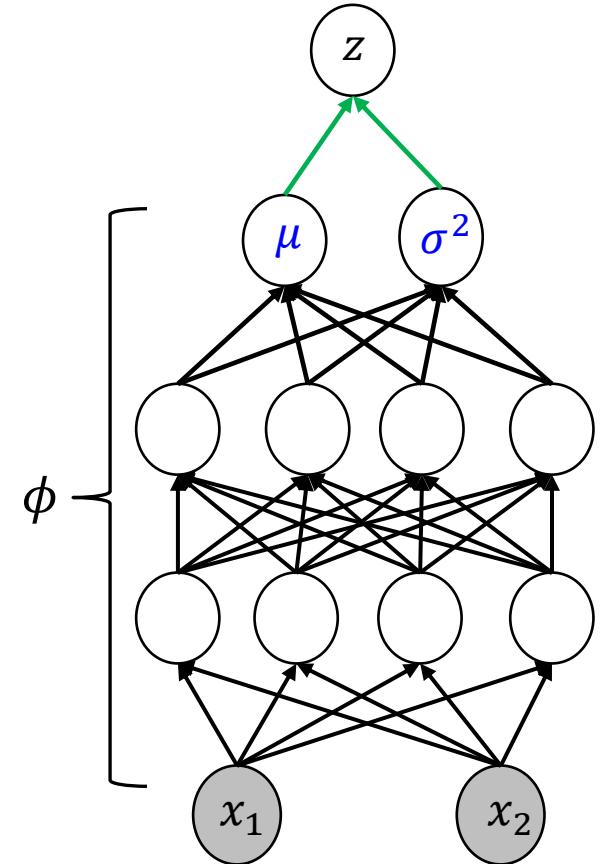
$$\mathbf{E}_{q(z|x;\phi)} [\log p(x|z; \theta)]$$

- Approximate via Monte Carlo methods

$$\mathbf{E}_{q(z|x;\phi)} [\log p(x|z; \theta)] \approx \frac{1}{L} \sum_k \log p(x|z^{(k)})$$

$$z^{(k)} \sim q(z|x; \phi)$$

- An unbiased estimator



# Example with Gaussian Distributions

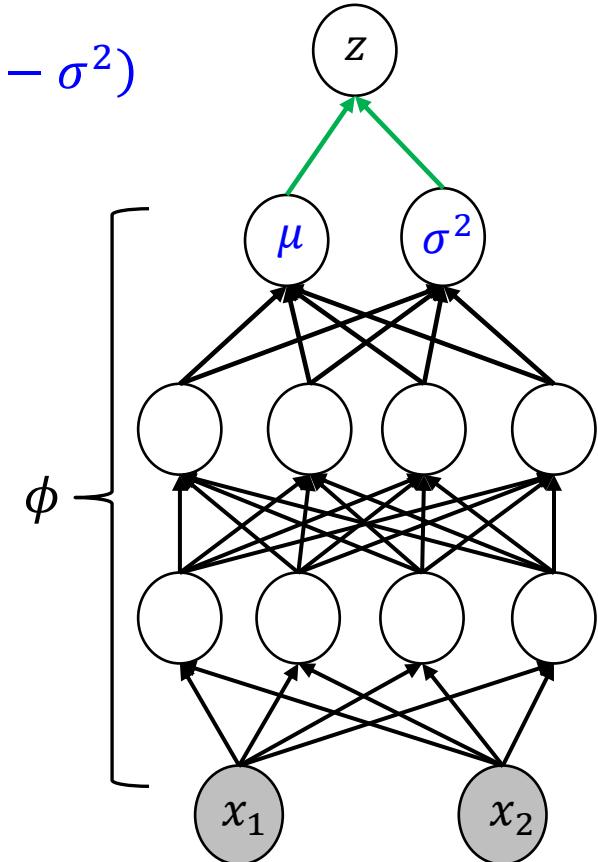
- ◆ The KL-regularization term (closed-form):

$$-\text{KL}(q(z|x; \phi) \| p(z; \theta)) = \frac{1}{2} (1 + \log \sigma^2 - \mu^2 - \sigma^2)$$

- Easy to calculate gradient
- ◆ The expected log-likelihood term (MC estimate)

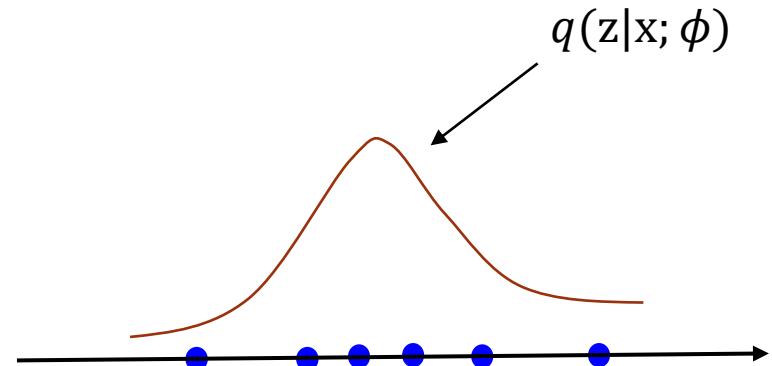
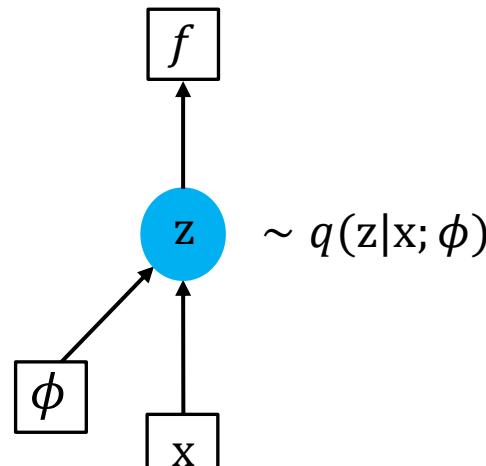
$$\mathbf{E}_{q(z|x;\phi)}[\log p(x|z; \theta)] \approx \frac{1}{L} \sum_k \log p(x|z^{(k)})$$
$$z^{(k)} \sim q(z|x; \phi)$$

- Gradient needs back-propagation!
- *However,  $Z^{(k)}$  is a random variable, we can't take gradient over a randomly drawn number*



# Reparameterization Trick

- ◆ Backpropagation not possible through random sampling



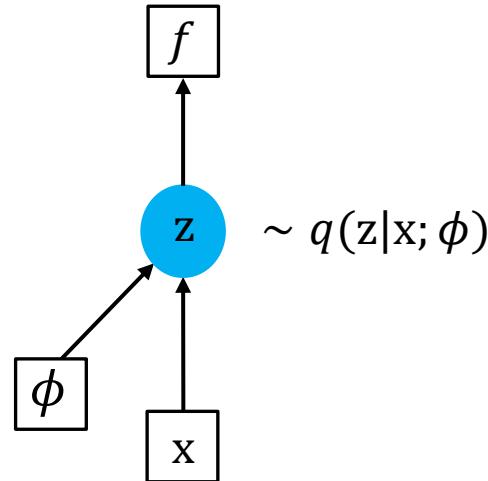
$$z^{(k)} \sim N(\mu(x, \phi), \sigma^2(x, \phi))$$

$$\{-1.5, -0.5, 0.3, 0.6, 1.5, \dots\}$$

Cannot back-propagate through a randomly drawn number

# Reparameterization Trick

- ◆ Backpropagation not possible through random sampling

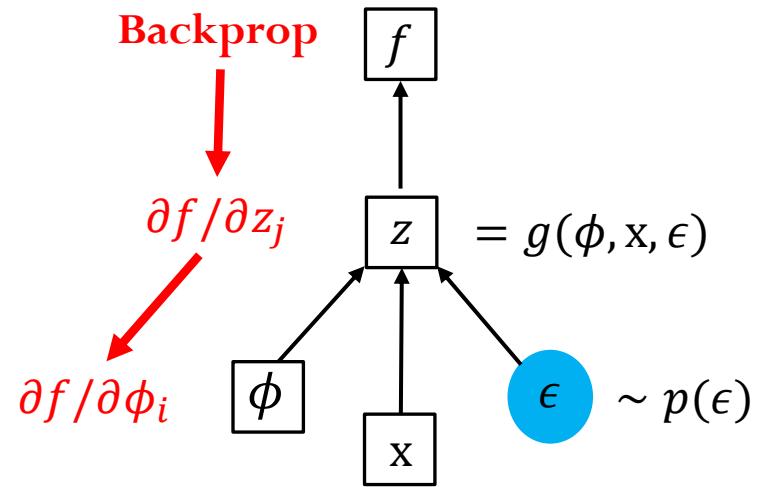


◻ : deterministic node

● : random node

$$z^{(k)} \sim N(\mu(x, \phi), \sigma^2(x, \phi))$$

Cannot back-propagate through a randomly drawn number



$$\epsilon^{(k)} \sim N(0,1)$$

$$z^{(k)} = \mu(x, \phi) + \sigma(x, \phi) \cdot \epsilon^{(k)}$$

$Z$  has the same distribution, but now can back-prop  
Separate into a deterministic part and noise

# The General Form

- ◆ The VAE bound

$$\begin{aligned} L(\theta, \phi, \mathbf{x}) &= \mathbf{E}_{q(z|x;\phi)} [\log p(z, \mathbf{x}; \theta) - \log q(z|x; \phi)] \\ &= \mathbf{E}_{q(z|x;\phi)} \left[ \log \frac{p(z, \mathbf{x}; \theta)}{q(z|x; \phi)} \right] \end{aligned}$$

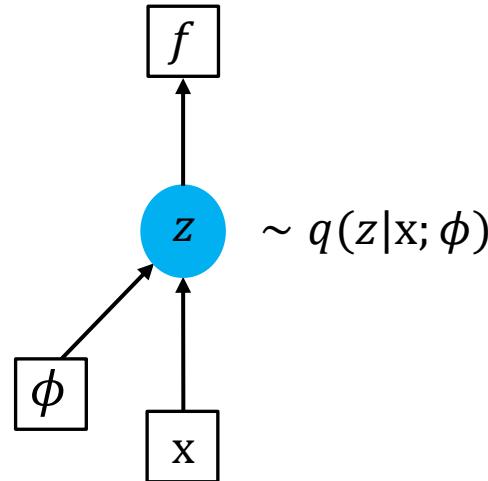
- ◆ Monte Carlo estimate:

$$\begin{aligned} L(\theta, \phi, \mathbf{x}) &\approx \frac{1}{L} \sum_k \log \frac{p(z^{(k)}, \mathbf{x}; \theta)}{q(z^{(k)}|x; \phi)} \\ z^{(k)} &\sim q(z|x; \phi) \end{aligned}$$

- Again, we cannot back-prop through the randomly drawn numbers

# Reparameterization Trick

- ◆ Backpropagation not possible through random sampling



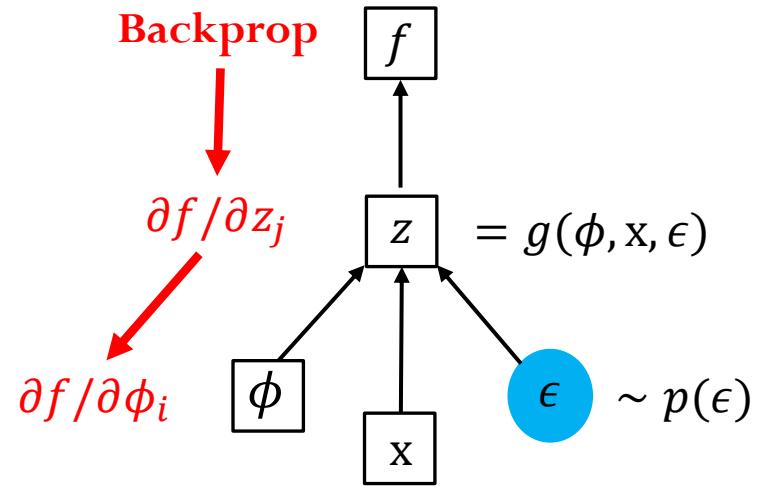
$$\sim q(z|x; \phi)$$

◻ : deterministic node

● : random node

$$z^{(k)} \sim q(z|x; \phi)$$

Cannot back-propagate through a randomly drawn number



$$= g(\phi, x, \epsilon) \quad \sim p(\epsilon)$$

$$\epsilon^{(k)} \sim p(\epsilon)$$

$$z^{(k)} = g(\phi, x, \epsilon^{(k)})$$

$Z$  has the same distribution, but now can back-prop  
Separate into a deterministic part and noise

# Reparam-Trick Summary

## ◆ The VAE bound

$$L(\theta, \phi, x) = \mathbf{E}_{q(z|x; \phi)} \left[ \log \frac{p(z, x; \theta)}{q(z|x; \phi)} \right]$$

□ Reparameterized as

$$L(\theta, \phi, x) = \mathbf{E}_{p(\epsilon)} \left[ \log \frac{p(g(x, \epsilon, \phi), x; \theta)}{q(g(x, \epsilon, \phi)|x; \phi)} \right]$$

- where  $\epsilon$  is a simple distribution (e.g., standard normal) and  $g$  is a deep NN

## ◆ The gradients are

$$\nabla_{\theta} L(\theta, \phi, x) = \mathbf{E}_{p(\epsilon)} \left[ \nabla_{\theta} \log \frac{p(g(x, \epsilon, \phi), x; \theta)}{q(g(x, \epsilon, \phi)|x; \phi)} \right]$$

- Back-prop is applied over the deep NN
- Similar for  $\phi$

# Importance Weighted Auto-Encoder (IWAE)

- ◆ The VAE lower bound of log-likelihood

$$L(\theta, \phi, \mathbf{x}) = \mathbf{E}_{q(z|x; \phi)} \left[ \log \frac{p(z, \mathbf{x}; \theta)}{q(z|x; \phi)} \right]$$

- ◆ A better variational lower bound (IWAE)

$$L_K(\theta, \phi, \mathbf{x}) = \mathbf{E}_{q(z|x; \phi)} \left[ \log \left( \frac{1}{K} \sum_{k=1:K} \frac{p(z^{(k)}, \mathbf{x}; \theta)}{q(z^{(k)}|x; \phi)} \right) \right]$$

where  $z^{(k)} \sim q(z|x; \phi)$

- This is a lower-bound of the log-likelihood
- When  $K=1$ , recovers the VAE bound
- When  $K = \infty$ , recovers the log-likelihood
- A monotonic sequence:

$$L_K(\theta, \phi, \mathbf{x}) \leq L_{K+1}(\theta, \phi, \mathbf{x}), \quad \forall \theta, \phi, \mathbf{x}$$

# Reparametrization Trick

- ◆ The IWAE bound:

$$L_K(\theta, \phi, \mathbf{x}) = \mathbf{E}_{q(z|x;\phi)} \left[ \log \left( \frac{1}{K} \sum_{k=1:K} w(z^{(k)}, \mathbf{x}; \theta) \right) \right]$$

$$\text{where } z^{(k)} \sim q(z|x; \phi) \quad w(z^{(k)}, \mathbf{x}; \theta, \phi) = \frac{p(z^{(k)}, \mathbf{x}; \theta)}{q(z^{(k)}|x; \phi)}$$

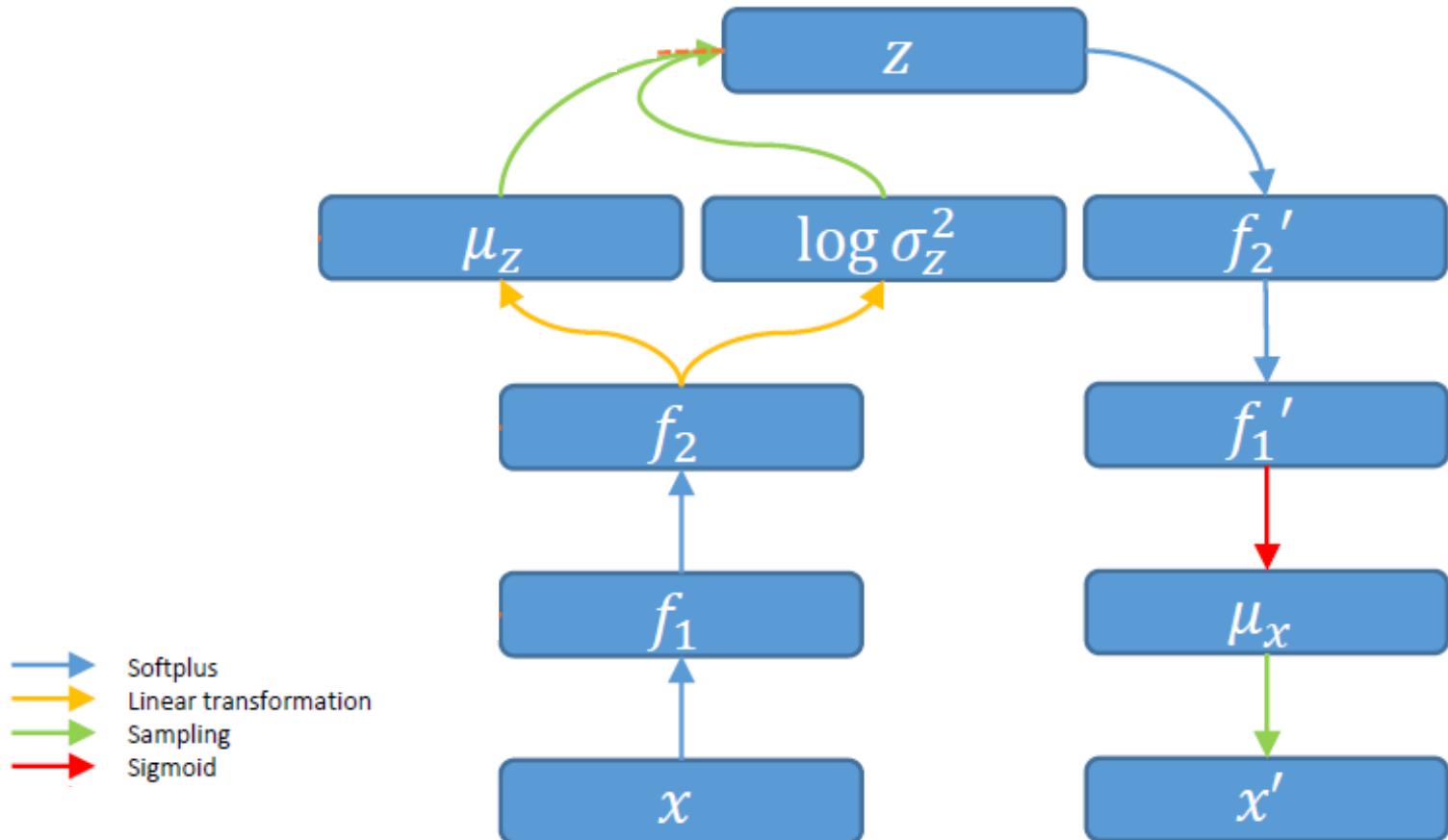
- ◆ Reparameterization form:

$$L_K(\theta, \phi, \mathbf{x}) = \mathbf{E}_{p(\epsilon)} \left[ \log \left( \frac{1}{K} \sum_{k=1:K} w(g(\epsilon^{(k)}, \mathbf{x}, \phi), \mathbf{x}; \theta) \right) \right]$$

$$\text{where } \epsilon^{(k)} \sim p(\epsilon)$$

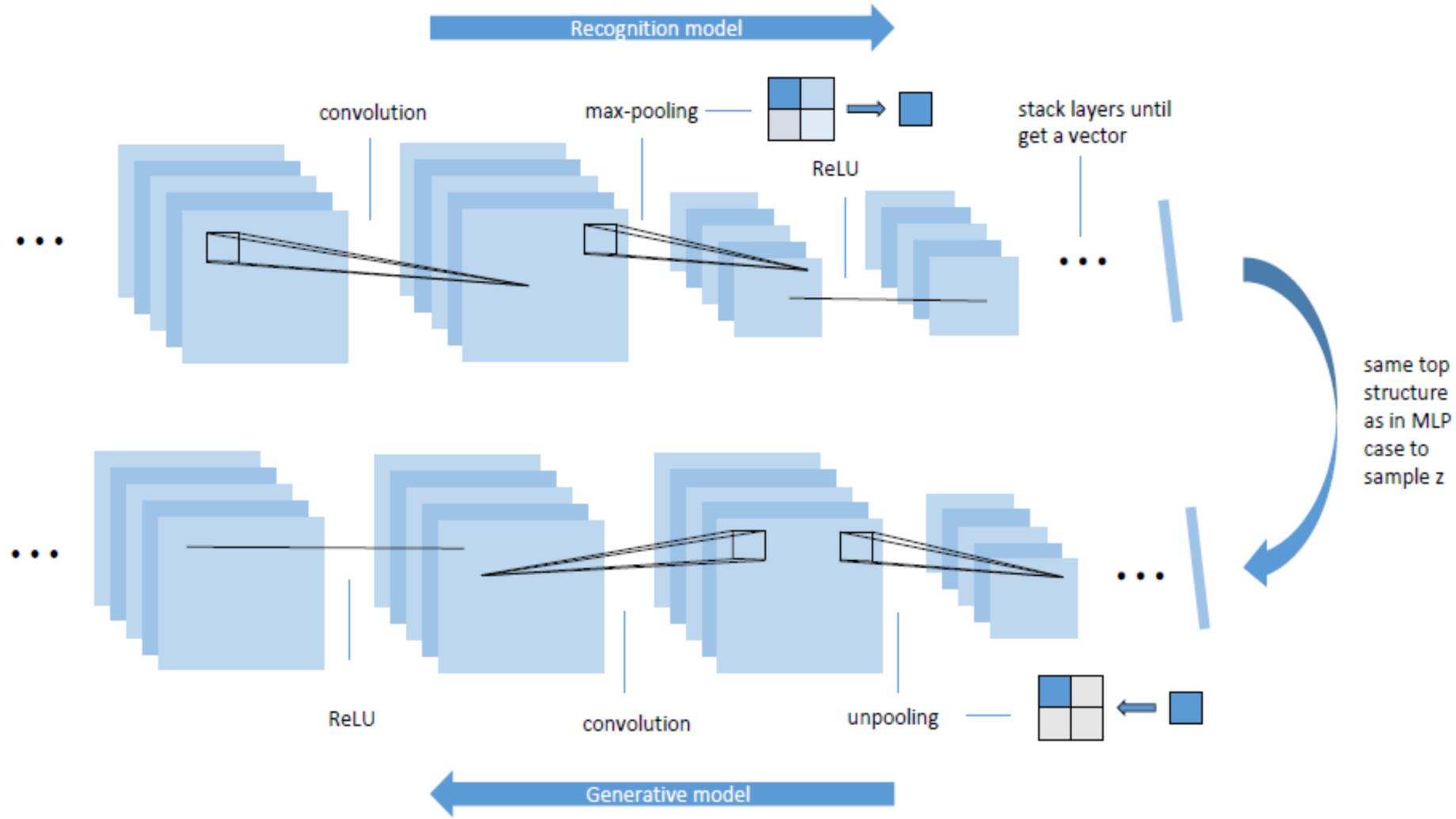
- The gradient can be calculated as in VAE

## 2-Layer MLP: Q-P network architecture



\*Same as in Auto-Encoding Variational Bayes (VA) [Kingma & Welling, 2014]

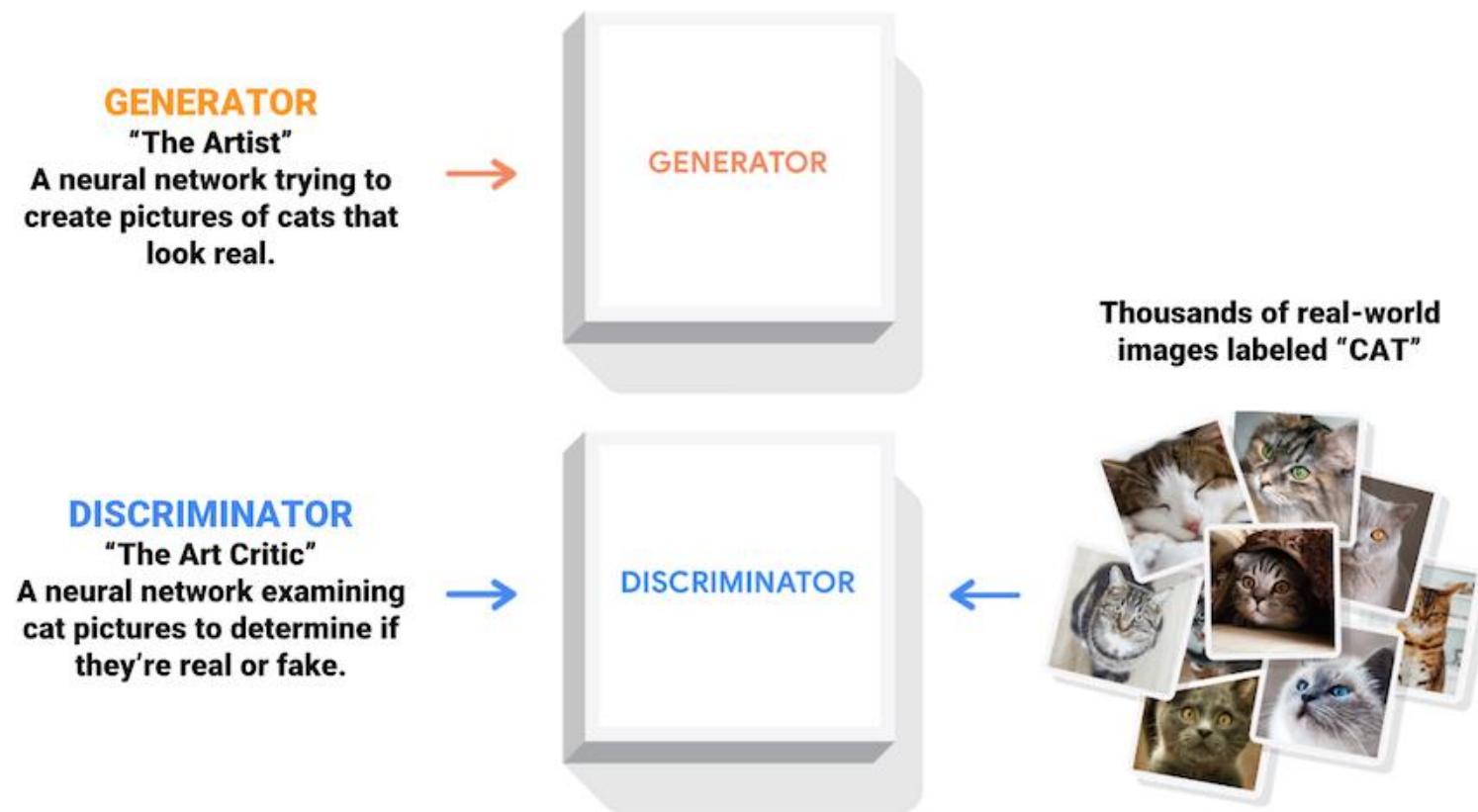
# 5-Layer CNN: Q-P network architecture



# Generative Adversarial Networks (GAN)

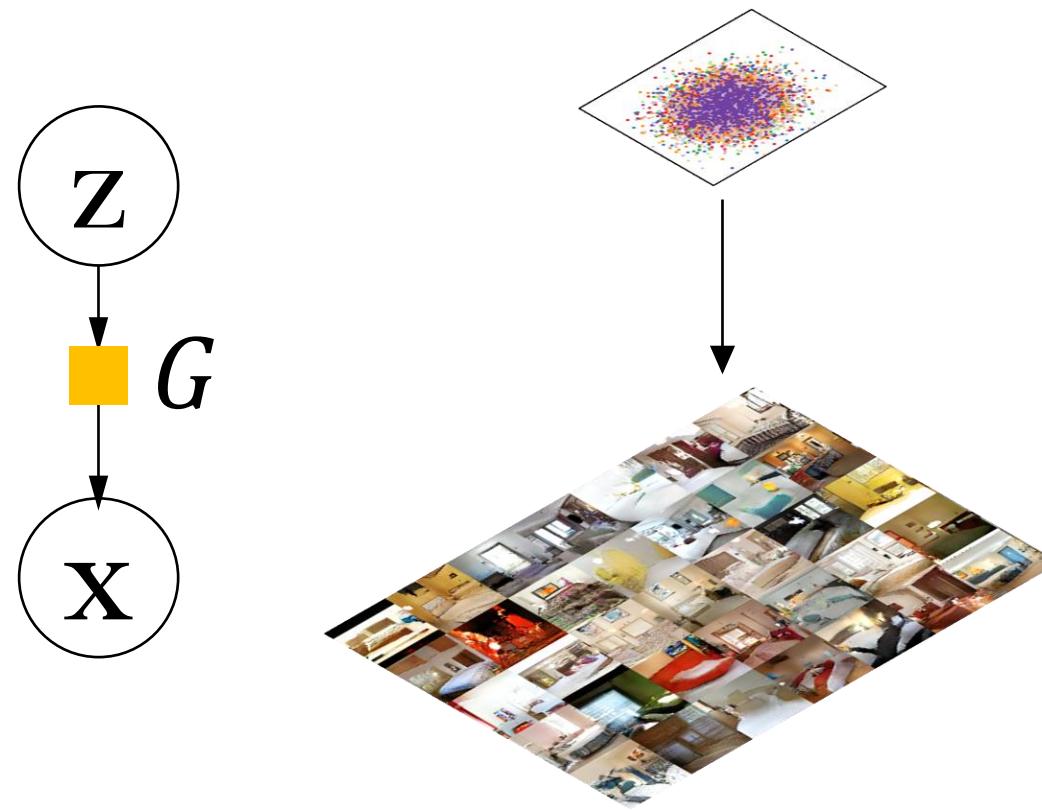
- ◆ A game between two players:
  - A discriminator  $D$
  - A generator  $G$
- ◆  $D$  tries to discriminate between:
  - A sample from the data distribution.
  - And a sample from the generator  $G$ .
- ◆  $G$  tries to “trick”  $D$  by generating samples that are hard for  $D$  to distinguish from data.

# The “Artist” and “Critic” Argument



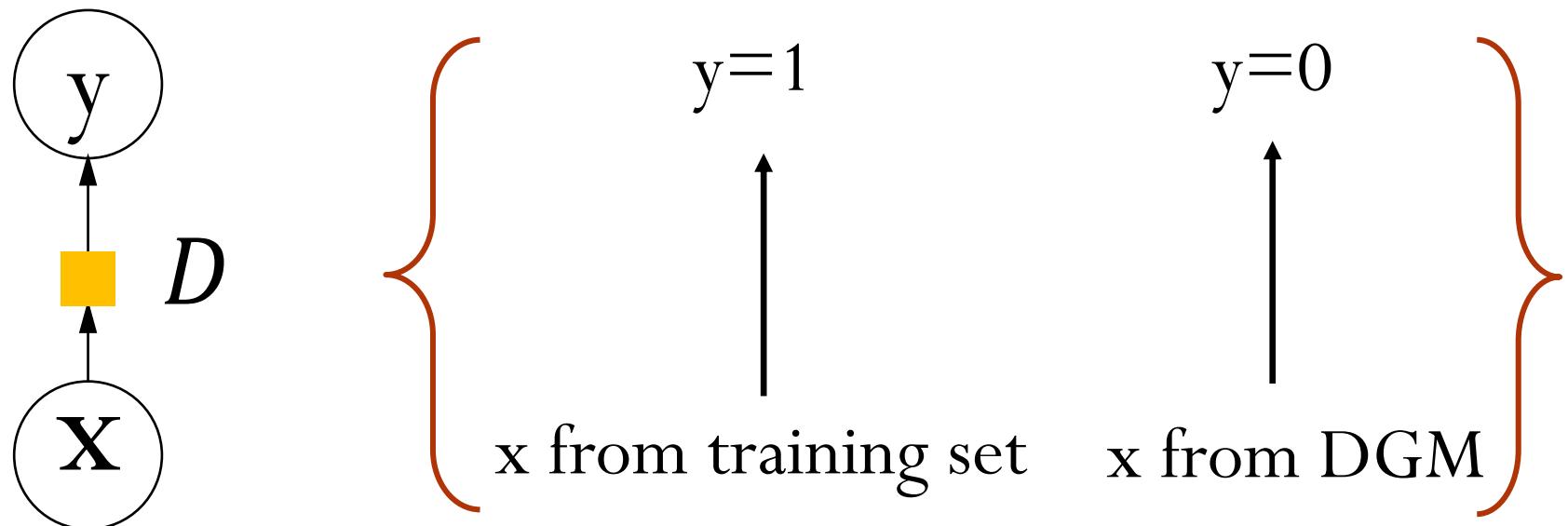
# GAN – architecture

- ◆ The generator-network  $G$  is a DGM
  - It generates samples from random noise



# GAN – architecture

- ◆ The discriminator-network  $D$  is a binary classifier
  - It aims to assign the correct label to both training samples and the samples from  $G$



- This is a supervised learning task (binary classification)!

# GAN - objective

- ◆ The discriminator-network  $D$  is a binary classifier
    - It aims to assign the correct label to both training samples and the samples from  $G$
  - ◆ Maximum likelihood estimation (MLE) is the natural choice!

$$\max_D \mathbf{E}_{p_{\text{data}}(\mathbf{x})}[\log D(\mathbf{x})] + \mathbf{E}_{p(\mathbf{z})} \left[ \log (1 - D(G(\mathbf{z}))) \right]$$


  
 x from training set       $G(\mathbf{z})$  from generator

- $D(x) = p(y=1 \mid x)$
  - aka. cross-entropy loss minimization

# GAN - objective

- ◆ The generator aims to fool the discriminator  $D$ 
  - Generated samples should be identified as “real” by  $D$
  - Maximize the likelihood of being real:

$$\max_G \mathbf{E}_{p(z)} [\log(D(G(z)))]$$

↑

$G(z)$  is a sample from generator

- Or minimize the likelihood of being fake:

$$\min_G \mathbf{E}_{p(z)} [\log(1 - D(G(z)))]$$

# GAN – objective

- ◆ Minimax objective function

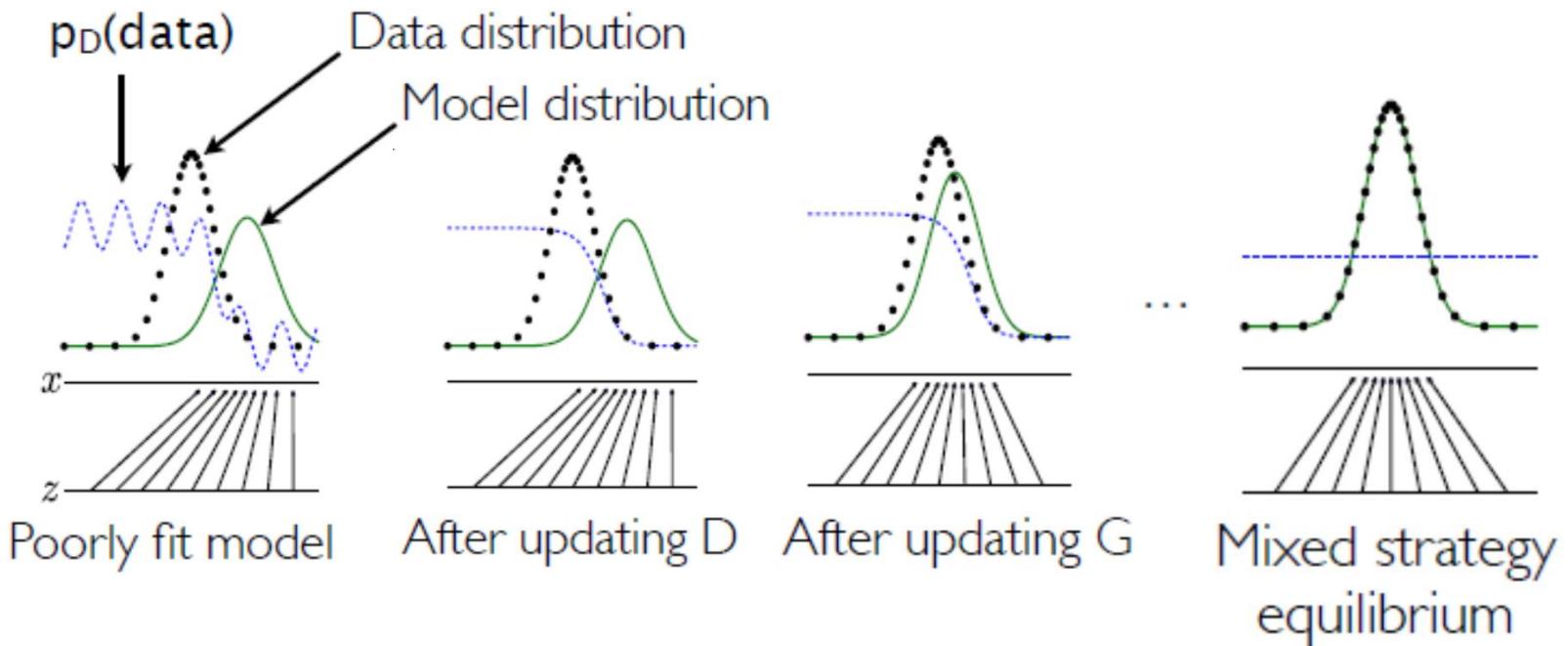
$$\min_G \max_D \mathbf{E}_{p_{\text{data}}(x)}[\log D(x)] + \mathbf{E}_{p(z)} [\log (1 - D(G(z)))]$$

- Optimal strategy of the discriminator for any  $p_{\text{model}}(x)$  is

$$D(x) = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_{\text{model}}(x)}$$

- Assume infinite data, infinite model capacity, direct updating generator's distribution
  - Unique global optimum
  - Optimum corresponds to data distribution

# GAN – learning process



- An adversarial pair near convergence:  $D$  is partially accurate
- Update discriminator distribution  $D$  according to the optimal strategy
- Update generator distribution to be more similar to data distribution
- Iterate until convergence

# Practice

- ◆ Early stopping of discriminator training to avoid overfitting
  - Multiple updates of  $D$
  - One update of  $G$
- ◆ Gradient saturation
- ◆ Implement VAE and GAN on Tensorflow at the Review session.

# Flow-based Models

- ◆ A generative process with invertible (or bijective) function  $g$ :

$$z \sim p(z)$$

$$x = g(z)$$

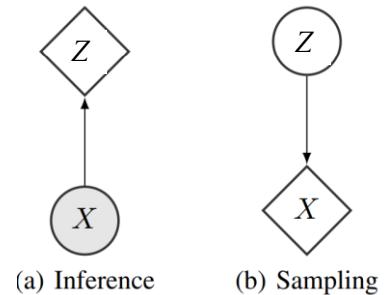
- ◆ The inference network for the latent variable is:

$$z = f(x) = g^{-1}(x)$$

- ◆ The density:

$$p(x) = p(f(x)) \left| \det \frac{\partial f(x)}{\partial x} \right|$$

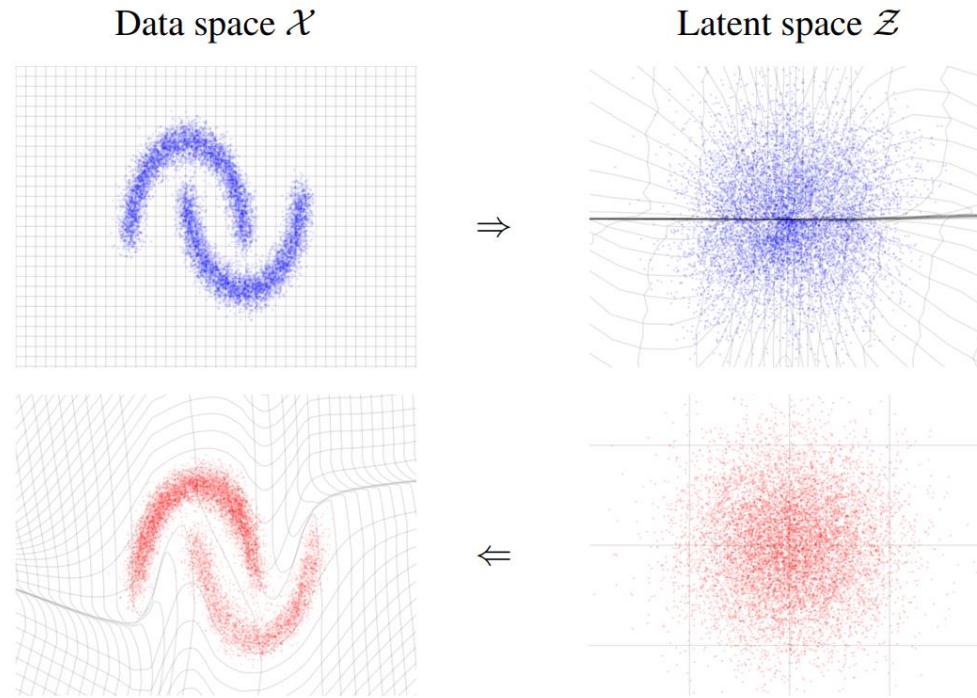
- ◆ Desirae of the function: easy determinant of Jacobian, easy inverse



# An illustrative example

**Inference**  
 $x \sim \hat{p}_X$   
 $z = f(x)$

**Generation**  
 $z \sim p_Z$   
 $x = f^{-1}(z)$



# An Example of the Invertible Function

- ◆ A simple idea: split  $x$  into two parts and define

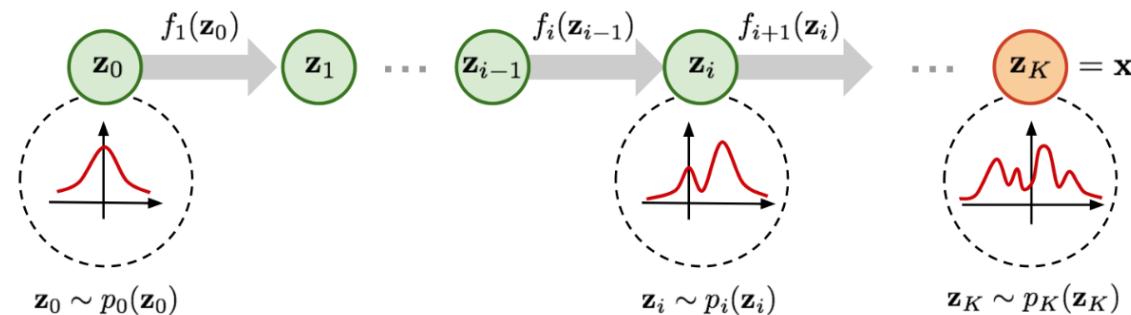
$$y_1 = x_1$$

$$y_2 = x_2 + m(x_1)$$

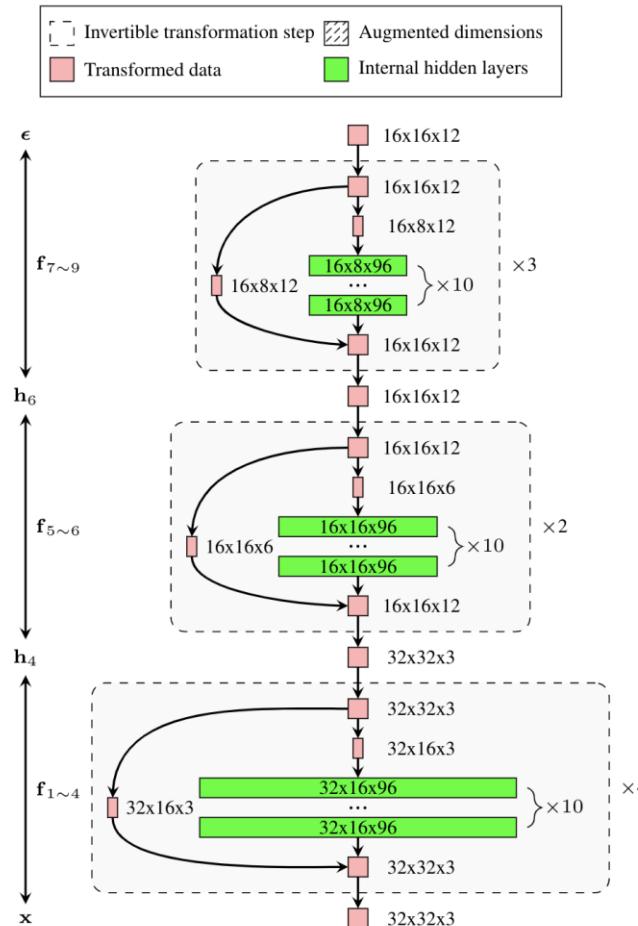
◆ Then the determinant-Jacobian 1!

◆ The inverse is  $x_1 = y_1$   
 $x_2 = y_2 - m(y_1)$

- $m$  is an arbitrarily complex function (e.g., ReLU networks)
- ◆ Define the invertible transform function as a composition of simple functions
  - Thus, the variable change is a sequence, called a *flow* (the sequence on density change is a *normalizing flow*)



# Bottleneck Problem in Flow++



Invertibility imposes constraints on architectures!

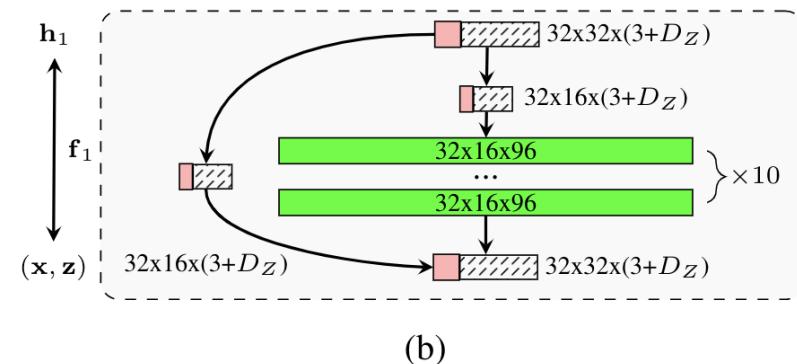
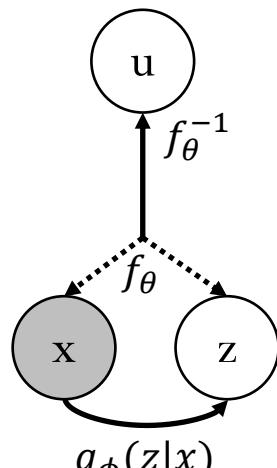


Figure 1. (a) Bottleneck problem in a Flow++ (Ho et al., 2019) for CIFAR-10. Dimensionality of the transformed data (red) limits the model capacity. (b) Our solution VFlow, where  $D_Z$  is the dimensionality of the augmented random variable. Only the transformation step  $f_1$  is shown due to space constraint.

# VFlow: more Expressive Generative Flows through Variational Data Augmentation

learn a generative flow  $p(x, z)$  in the augmented data space jointly with the augmented data distribution



$$u \sim \mathcal{N}(0, I)$$

$$y = concat([x, z])$$

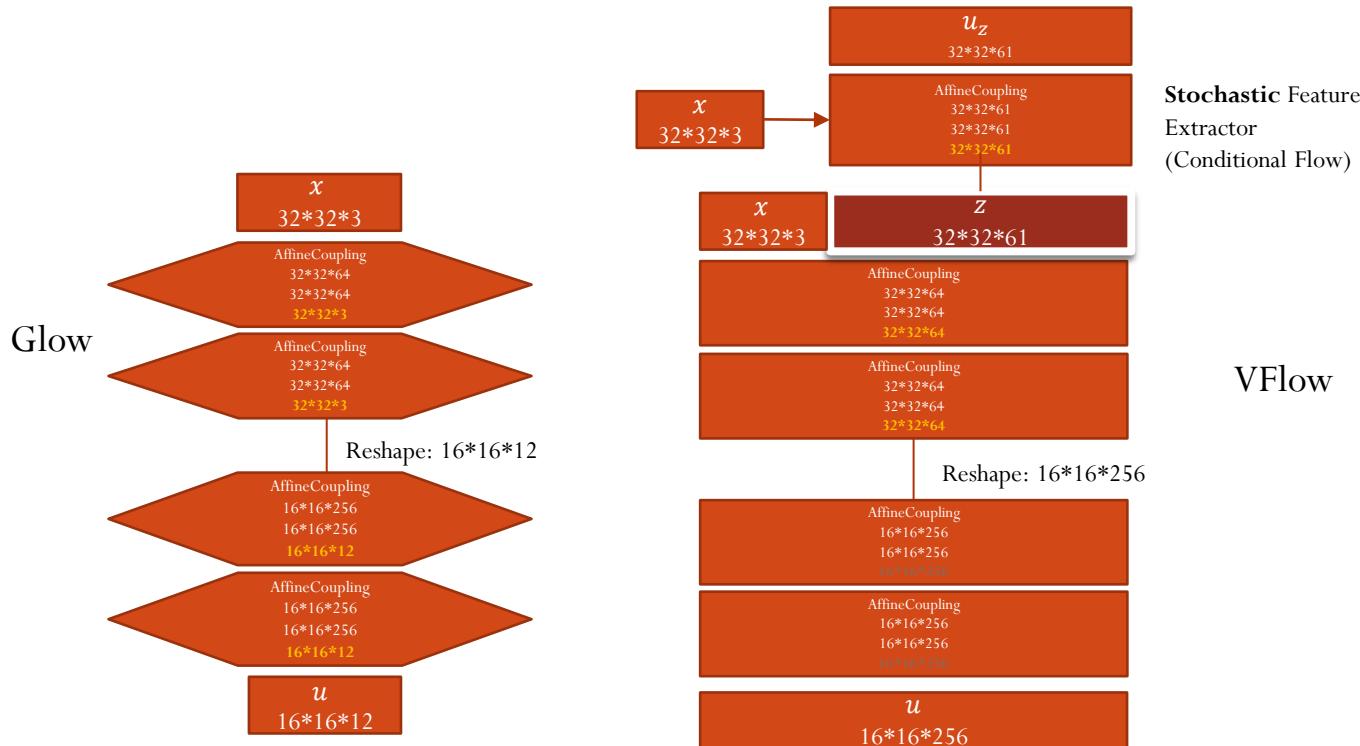
$$\log p(x, z) = \log p(u) + \log \det\left(\frac{du}{dy}\right)$$

$$\log p(x) \geq \mathbb{E}_{q(z|x)} [\log p(x, z) - \log q(z|x)]$$

$$\max_{\theta, \phi} \mathbb{E}_{\hat{p}(\mathbf{x})q(\mathbf{z}|\mathbf{x}; \phi)} [\log p(\mathbf{x}, \mathbf{z}; \theta) - \log q(\mathbf{z}|\mathbf{x}; \phi)]$$

Provably better than generative flows! Also subsumes VAEs as a special case.

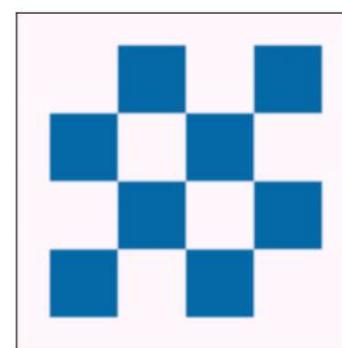
# VFlow



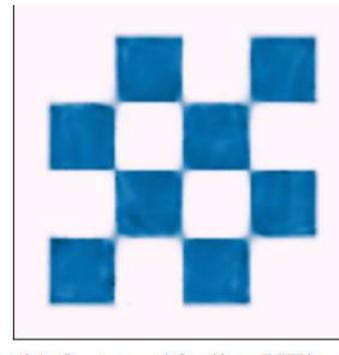
[Chen, Lu, Chenli, Zhu, Tian. *VFlow: More Expressive Generative Flows with Variational Data Augmentation*. ICML 2020]

# Results on Toy 2D Data

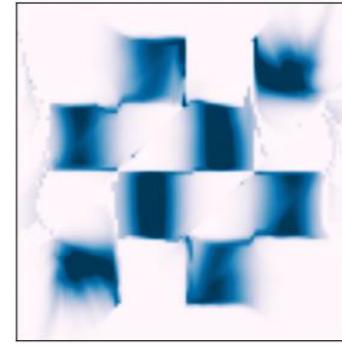
- ◆ VFlow significantly outperforms Glow under similar model size (e.g., 3-step)
- ◆ The 3-layer, 10-dimensional VFlow even outperforms a much larger 20-step Glow



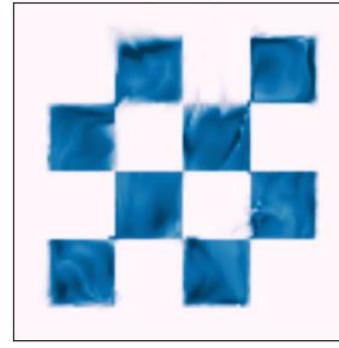
(a) Data (-3.47)



(b) 3-step, 10-dim VFlow  
(-3.51)



(c) 3-step Glow (-3.66)



(d) 20-step Glow (-3.52)

# Results on Toy 2D Data

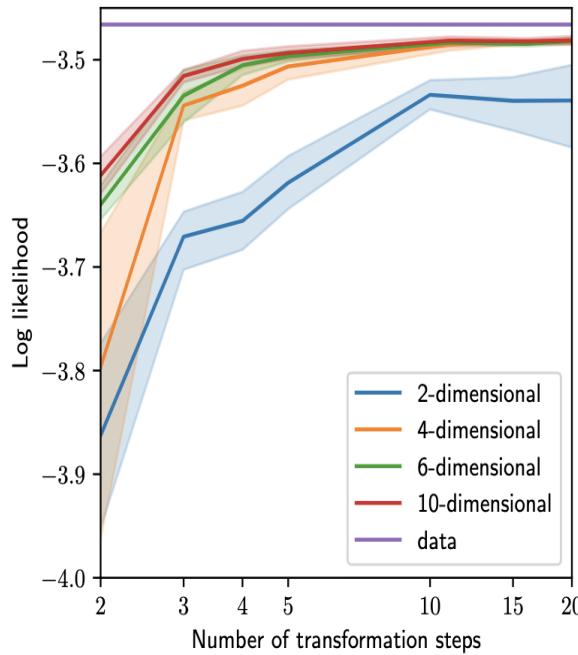


Figure 4. Impact of the dimensionality on the toy dataset.

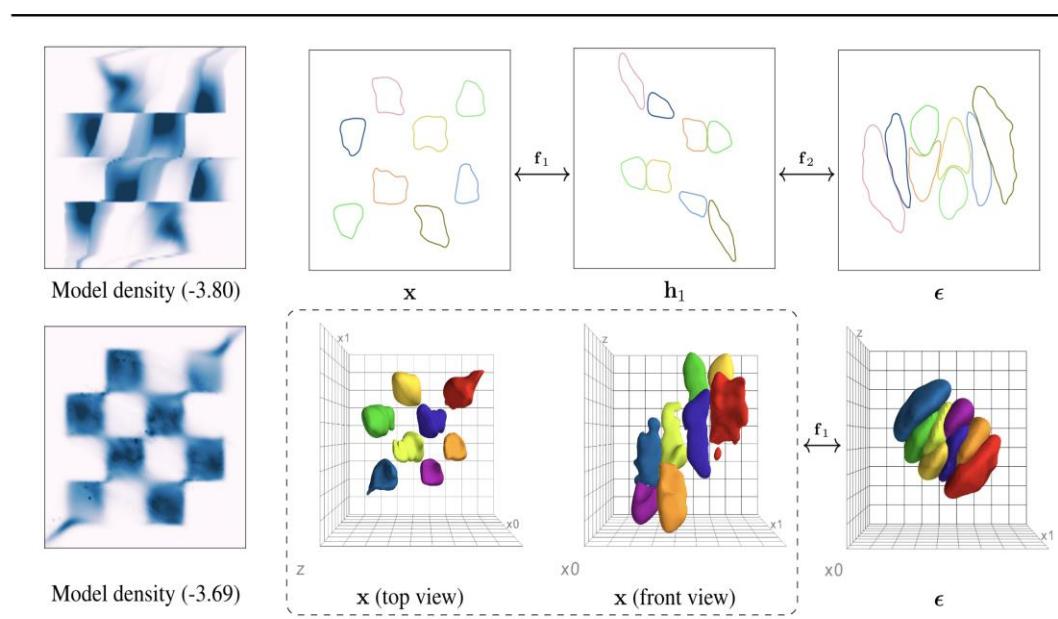


Figure 3. Visualization of learnt transformation on toy data. Top row: 2-step Glow. Bottom row: 2-step, 3-dimensional VFlow. Log-likelihood is shown in parenthesis. We sample  $\epsilon$  and visualize the transformed density in  $x$ ,  $h_1$  and  $\epsilon$  space. The density is estimated from samples by kernel density estimation, and we show the 50% probability contour / isosurface for each mode in different color.

# Results on CIFAR-10

Table 1. Density modeling result on the CIFAR-10 dataset.

Model	bpd
Glow (Kingma & Dhariwal, 2018)	3.35
FFJORD (Grathwohl et al., 2019)	3.40
Residual Flow (Chen et al., 2019)	3.28
MintNet (Song et al., 2019)	3.32
Flow++ (Ho et al., 2019)	3.08
VFlow	<b>2.98</b>

Table 3. Parameter efficiency.

Model	bpd	Parameters	$D_H$	$B$
3-channel Flow++	3.08	31.4M	96	10
6-channel VFlow	<b>2.98</b>	37.8M	96	10
6-channel VFlow	3.03	16.5M	64	10
6-channel VFlow	3.08	<b>11.9M</b>	56	10

Table 2. Impact of dimensionality on the CIFAR-10 dataset.

Model	bpd	Parameters	$D_H$	$B$
3-channel Flow++	3.21	4.02M	32	13
4-channel VFlow	3.15	4.03M	32	11
6-channel VFlow	<b>3.12</b>	<b>4.01M</b>	32	10

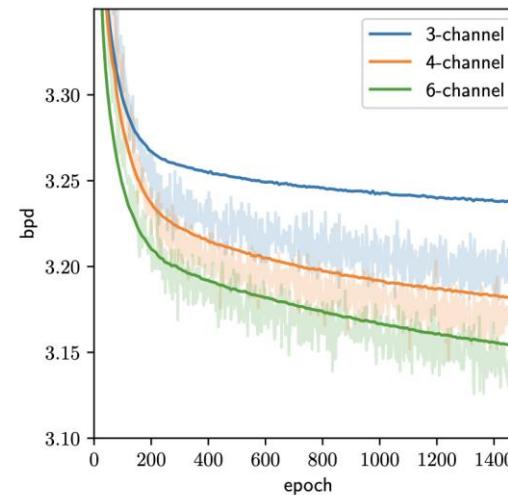


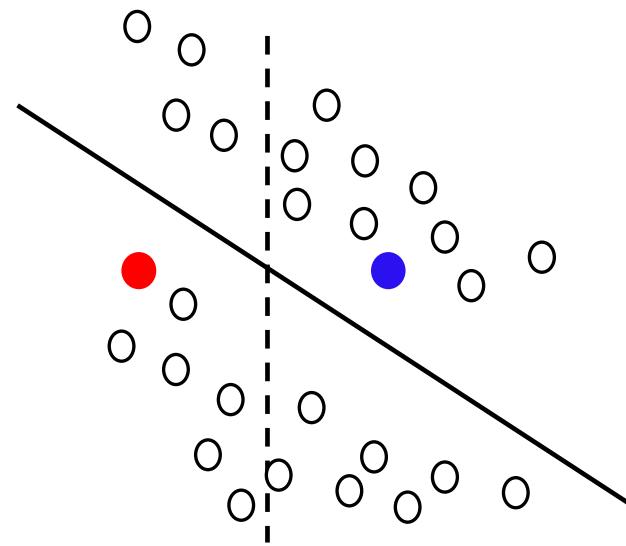
Figure 6. Bpd on training (light) and testing (dark) dataset of Flow++ and VFlow under a 4-million parameter budget. Here bpd is only a upper bound because we evaluate it with ELBO as Eq. (7) instead of the marginal likelihood.

# **Semi-supervised Learning with GANs**

(how to learn efficiently from a limited number of labels?)

# Semi-supervised Learning

- ◆ A toy example



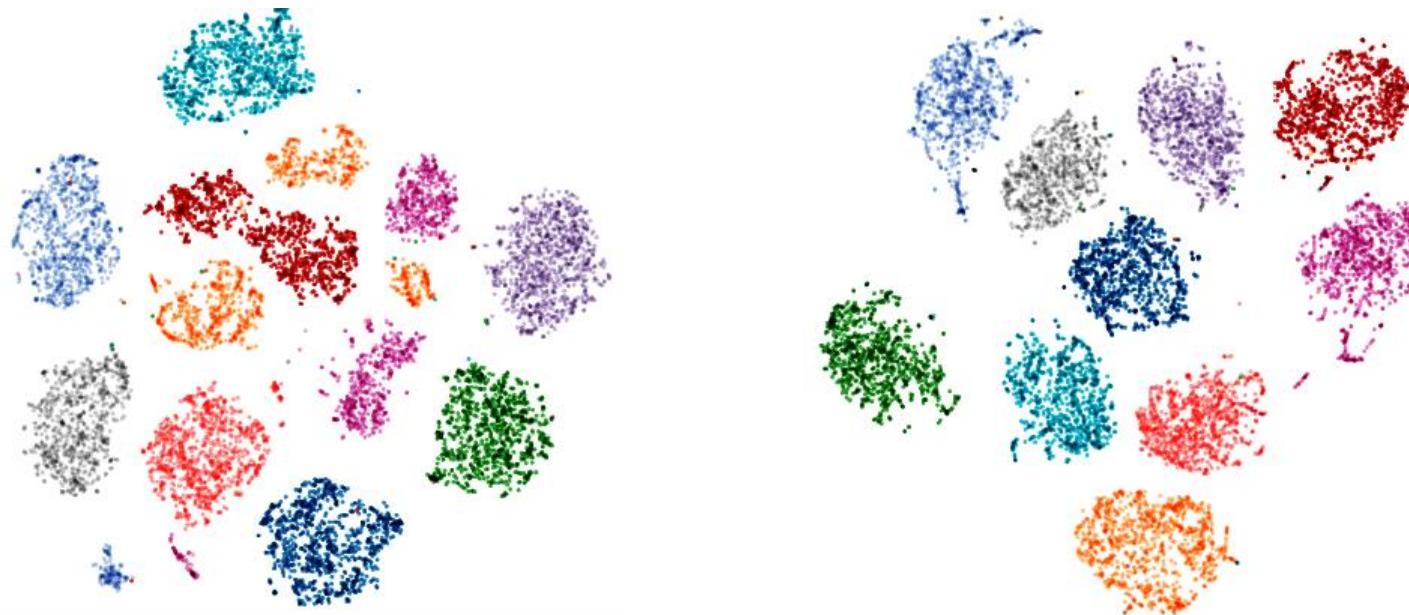
# Representation Matters

- ◆ PCA embedding of learned representations by different DGM models on CIFAR10



# Representation Matters

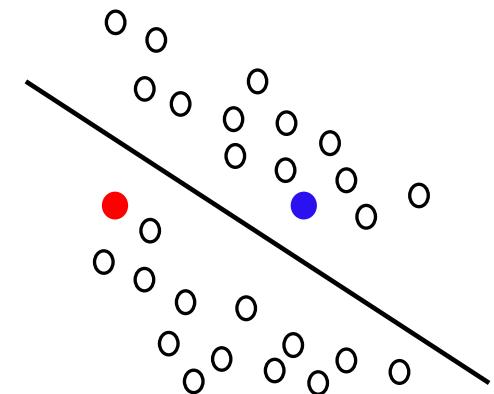
- ◆ t-SNE embedding of learned representations by different DGM models on CIFAR10



# Triple Generative Adversarial Nets

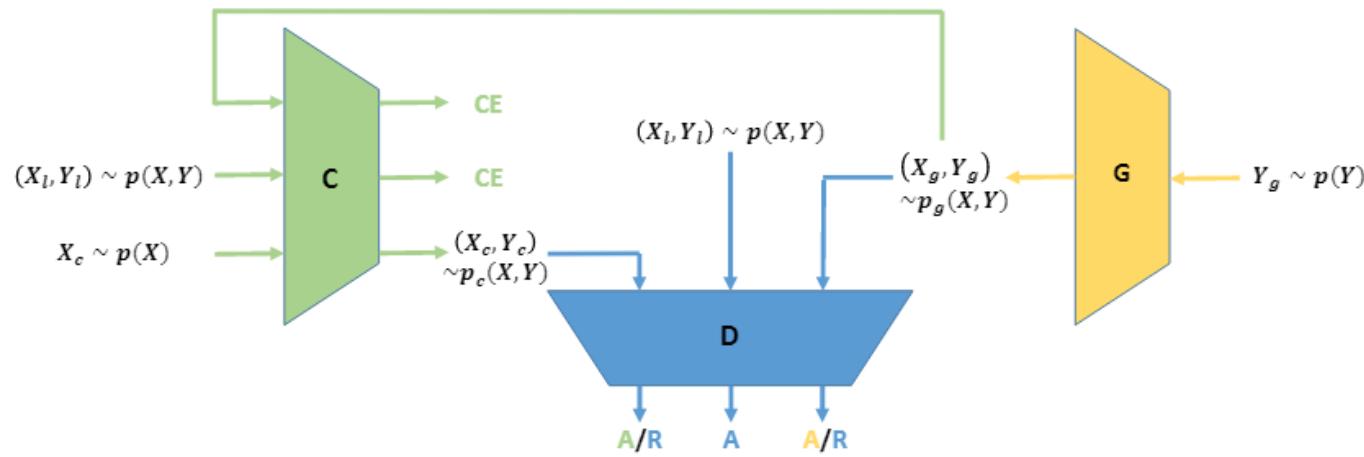
- ◆ A minimax game for semi-supervised learning
  - GAN is for unsupervised learning  $p_{\text{model}}(x) = p_{\text{data}}(x)$
  - We aim to learn the joint distribution  $p_{\text{model}}(x, y) = p_{\text{data}}(x, y)$
- ◆ We need three players
  - factorization form with conditionals
$$\begin{aligned} p(x, y) &= p(x)p(y|x) \\ &= p(y)p(x|y) \end{aligned}$$

A classifier                            A class-conditional generator
  - Two generators to generate  $(x, y)$
  - A discriminator to distinguish fake  $(x, y)$



# Triple-GAN

- ◆ The network architecture



- Both **C** and **G** are generators
- **D** is the discriminator
- **CE**: cross-entropy loss for learning classifier

# A minimax game

- ◆ The optimization problem

$$\min_{C,G} \max_D U(C, G, D) = E_p[\log D(x, y)] + \\ \alpha E_{p_c}[\log(1 - D(x, y))] + (1 - \alpha) E_{p_g}[\log(1 - D(x, y))]$$

- The hyper-parameter  $\alpha$  is often set at  $1/2$
- The

$$\min_{C,G} \max_D \tilde{U}(C, G, D) = U(C, G, D) + E_p[-\log p_c(y|x)]$$

# Major theoretical results

## Theorem

*The equilibrium of  $\tilde{U}(C, G, D)$  is achieved if and only if  $p(x, y) = p_g(x, y) = p_c(x, y)$  with  $D_{C,G}^*(x, y) = \frac{1}{2}$  and the optimum value is  $-\log 4$ .*

## Lemma

*For any fixed  $C$  and  $G$ , the optimal discriminator  $D$  is:*

$$D_{C,G}^*(x, y) = \frac{p(x, y)}{p(x, y) + p_\alpha(x, y)},$$

*where  $p_\alpha(x, y) := (1 - \alpha)p_g(x, y) + \alpha p_c(x, y)$ .*

# Some Practical Tricks for SSL

- Pseudo discriminative loss: using  $(x, y) \sim p_g(x, y)$  as labeled data to train  $C$ 
  - Explicit loss, equivalent to  $KL(p_g(x, y) || p_c(x, y))$
  - Complementary to the implicit regularization by  $D$
- Collapsing to the empirical distribution  $p(x, y)$ 
  - Sample  $(x, y) \sim p_c(x, y)$  as true data for  $D$
  - Biased solution: target shifting towards  $p_c(x, y)$
- Unlabeled data loss on  $C$ 
  - Confidence (Springenberg [2015])
  - Consistence (Laine and Aila [2016])

# Some Results

## ◆ Semi-supervised classification

Table 1: Error rates (%) on partially labeled MNIST, SHVN and CIFAR10 datasets. The results with  $\dagger$  are trained with more than 500,000 extra unlabeled data on SVHN.

Algorithm	MNIST $n = 100$	SVHN $n = 1000$	CIFAR10 $n = 4000$
<i>MI+M2</i> [11]	3.33 ( $\pm 0.14$ )	36.02 ( $\pm 0.10$ )	
VAT [18]	2.33		24.63
<i>Ladder</i> [23]	1.06 ( $\pm 0.37$ )		20.40 ( $\pm 0.47$ )
<i>Conv-Ladder</i> [23]	<b>0.89</b> ( $\pm 0.50$ )		
<i>ADGM</i> [17]	0.96 ( $\pm 0.02$ )	22.86 $\dagger$	
<i>SDGM</i> [17]	1.32 ( $\pm 0.07$ )	16.61( $\pm 0.24$ ) $\dagger$	
<i>MMCVA</i> [15]	1.24 ( $\pm 0.54$ )	<b>4.95</b> ( $\pm 0.18$ ) $\dagger$	
<i>CatGAN</i> [26]	1.39 ( $\pm 0.28$ )		19.58 ( $\pm 0.58$ )
<i>Improved-GAN</i> [25]	0.93 ( $\pm 0.07$ )	8.11 ( $\pm 1.3$ )	18.63 ( $\pm 2.32$ )
<i>ALI</i> [5]		7.3	18.3
<i>Triple-GAN (ours)</i>	<b>0.91</b> ( $\pm 0.58$ )	<b>5.77</b> ( $\pm 0.17$ )	<b>16.99</b> ( $\pm 0.36$ )

# Some Results

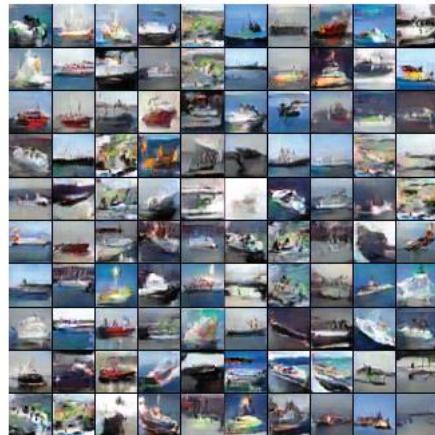
- ◆ Class-conditional generation



(d) Dog



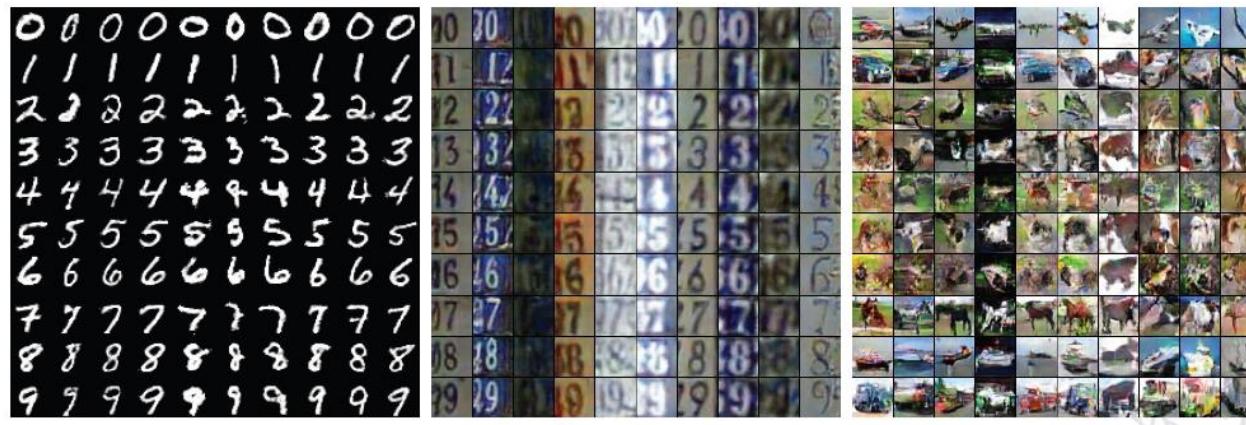
(e) Horse



(f) Ship

# Some Results

- ◆ Disentangle class and style



**Figure:** Same  $y$  for each row. Same  $z$  for each column.

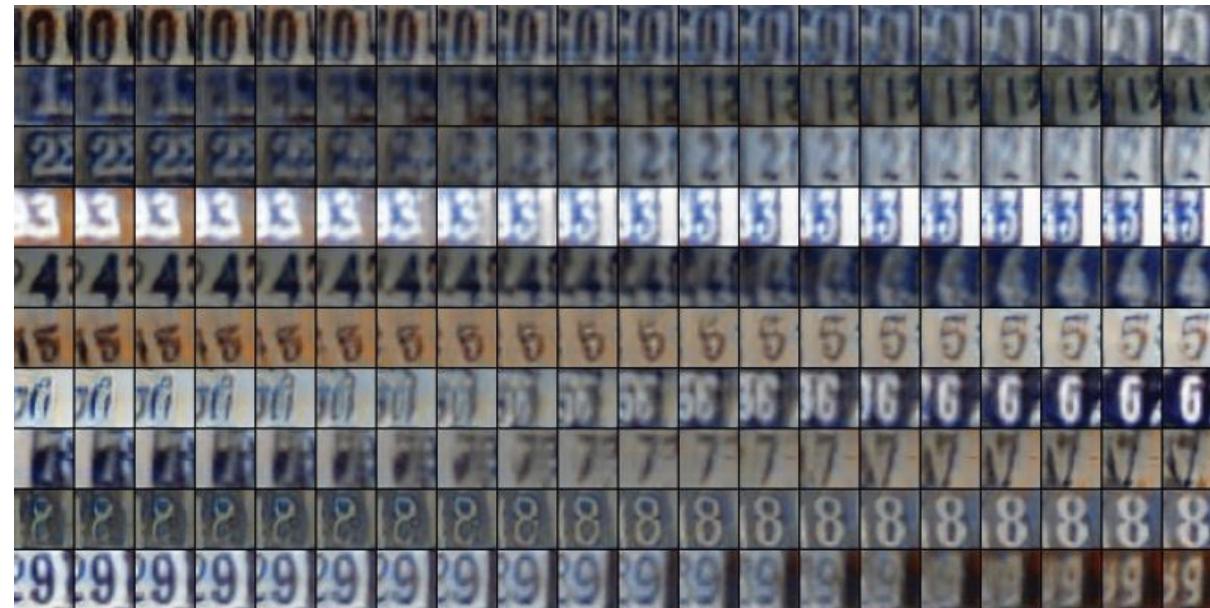
# Some Results

- ◆ Latent space interpolation on MNIST



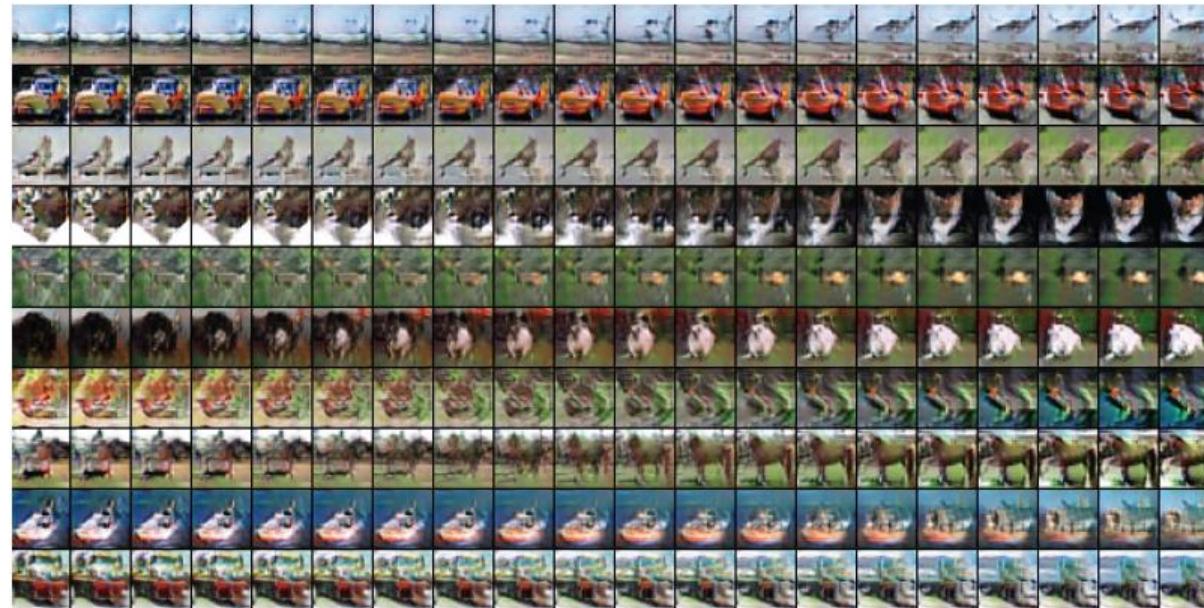
# Some Results

- ◆ Latent space interpolation on SVHN



# Some Results

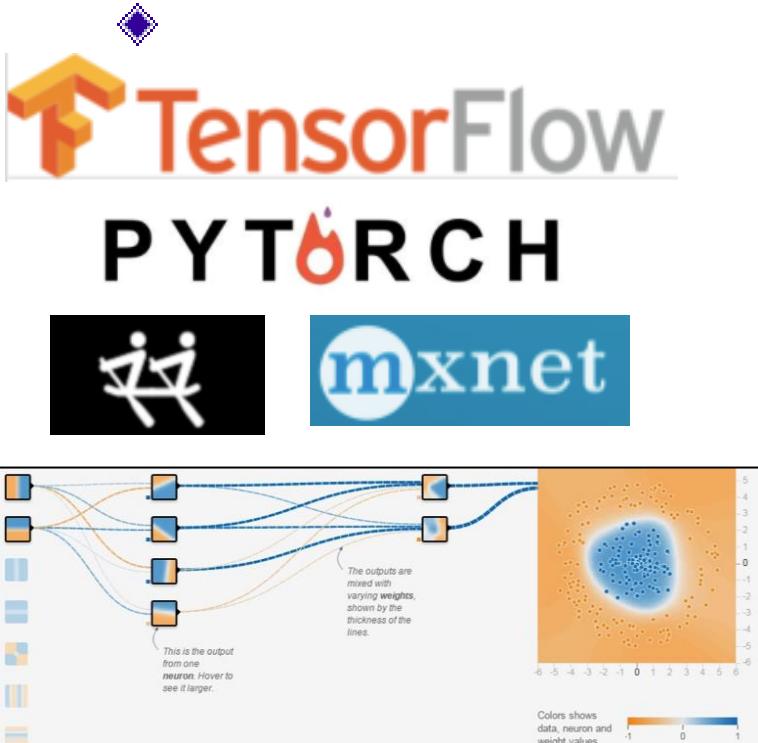
- ◆ Latent space interpolation on CIFAR10



# **Probabilistic Programming Library**

(how to auto/semi-auto implement Bayesian deep learning models?)

# DNN Programming $\rightarrow$ Deep Probabilistic Programming (Deep PPL)



Learning a mapping from  $x$  to  $y$

The screenshot shows the ZhuSuan documentation homepage. It features a header with the ZhuSuan logo and navigation links for 'Docs' and 'Edit on GitHub'. The main content area is titled 'Welcome to ZhuSuan' and includes several small diagrams of neural network graphs. Below this is a detailed description of ZhuSuan: 'ZhuSuan is a python probabilistic programming library for Bayesian deep learning, which conjoins the complimentary advantages of Bayesian methods and deep learning. ZhuSuan is built upon [Tensorflow](#). Unlike existing deep learning libraries, which are mainly designed for deterministic neural networks and supervised tasks, ZhuSuan provides deep learning style primitives and algorithms for building probabilistic models and applying Bayesian inference. The supported inference algorithms include:

- Variational inference with programmable variational posteriors, various objectives and advanced gradient estimators (SGVB, REINFORCE, VIMCO, etc.).
- Importance sampling for learning and evaluating models, with programmable proposals.
- Hamiltonian Monte Carlo (HMC) with parallel chains, and optional automatic parameter tuning.

**Installation**

ZhuSuan is still under development. Before the first stable release (1.0), please clone the [GitHub repository](#) and run

```
pip install .
```

[github.com/thu-ml/zhusuan](https://github.com/thu-ml/zhusuan)

# “ZhuSuan” Deep Probabilistic Programming



ZHUSUAN

Improved modeling language

```
def build_model(...):
```

1

Model build function

```
    bn = zs.BayesianNet()
```

2

Create a Bayesian Net

```
    x = bn.normal("x", x_mean,  
std=...)
```

3

Add stochastic node, e.g., a Gaussian  
variable x

```
    y = bn.deterministic("y",  
tf.any_op(x))
```

4

Add deterministic nodes, can use any  
Tensorflow operation

```
Or y = tf.any_op(x)
```

.....

```
return bn
```

5

Return the built Bayesian Net

# “ZhuSuan” Deep Probabilistic Programming



ZHUSUAN

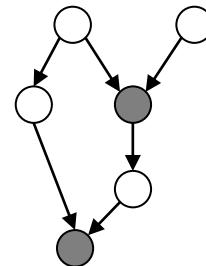
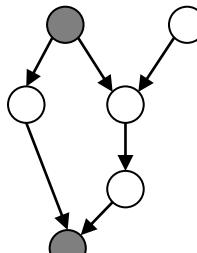
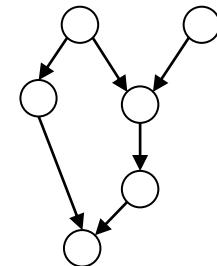
New model reuse strategy

Just need a decorator `@zs.meta_bayesian_net` to the model-building function

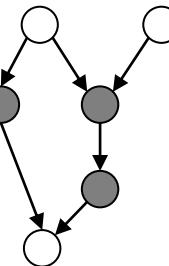
```
@zs.meta_bayesian_net(scope="model", reuse_variables=True)
```

```
def build_model(...):
```

The returned is a reusable `zs.MetaBayesianNet` object, just call `observe()` to assign **observed values** to any subset of random variables



`zs.BayesianNet`



`zs.MetaBayesianNet`

```
bn = meta_bn.observe(x=x_obs, y=y_obs)
```

[github.com/thu-ml/zhusuan](https://github.com/thu-ml/zhusuan)

# “ZhuSuan” Deep Probabilistic Programming



ZHUSUAN

## New Algorithms: Stochastic gradient MCMC

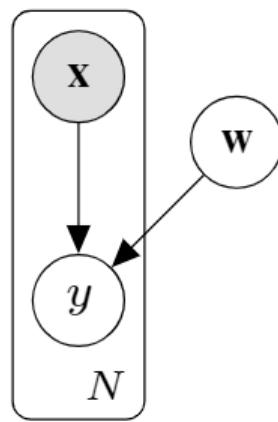
变分推理	梯度估计器	实现
<code>elbo()</code>	SGVB	<code>.sgvb()</code>
<code>k1pq()</code>	Reinforce	<code>.reinforce()</code>
<code>iw_objective()</code>	Importance Sampling	<code>.importance()</code>
	SGVB	<code>.sgvb()</code>
	VIMCO	<code>.vimco()</code>
马尔科夫链蒙特卡洛		实现
	Hamiltonian Monte Carlo	<code>HMC()</code>
Stochastic Gradient MCMC	(Preconditioned) SGLD	<code>SGLD()</code> , <code>PSGLD()</code>
	SGHMC	<code>SGHMC()</code>
	SGNHT	<code>SGNHT()</code>

[github.com/thu-ml/zhusuan](https://github.com/thu-ml/zhusuan)

# Examples: Bayesian Logistic Regression



ZHUSUAN



(a)

```
import zhusuan as zs
import tensorflow as tf

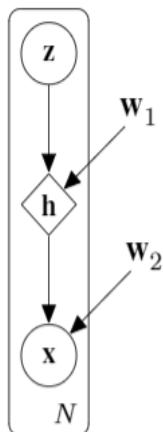
@zs.meta_bayesian_net(scope="blr")
def build_blr(x, alpha, D):
    bn = zs.BayesianNet()
    w = bn.normal('w', mean=tf.zeros([D]),
                  std=alpha, group_ndims=1)
    y_logit = tf.reduce_sum(
        tf.expand_dims(w, 0)*x, axis=1)
    bn.bernoulli('y', y_logit)
    return bn
```

(b)

# Example: Variational Auto-Encoder

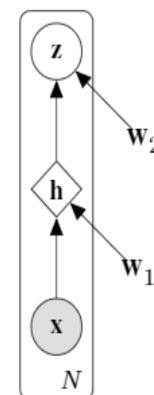


ZHUSUAN



(a)

```
@zs.meta_bayesian_net("vae", reuse_variables=True)
def build_vae(N, D):
    bn = zs.BayesianNet()
    z_mean = tf.zeros([N, D])
    z = bn.normal('z', z_mean, std=1., group_ndims=1)
    h = tf.layers.dense(z, 500,
                        activation=tf.nn.relu)
    x_logits = tf.layers.dense(h, 784)
    bn.bernoulli('x', x_logits, group_ndims=1)
    return bn
```



(a)

```
@zs.reuse_variables("vae_q")
def build_q_net(h_dim, z_dim):
    bn = zs.BayesianNet()
    h = tf.layers.dense(x, h_dim,
                        activation=tf.nn.relu)
    z_mean = tf.layers.dense(h, z_dim)
    z_logstd = tf.layers.dense(h, z_dim)
    bn.normal('z', z_mean, logstd=z_logstd,
              group_ndims=1)
    return bn
```

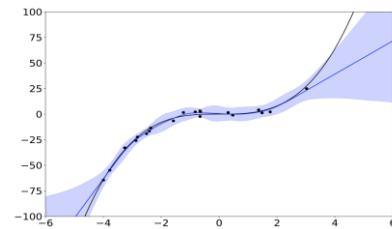
(b)

# “ZhuSuan” Deep Probabilistic Programming

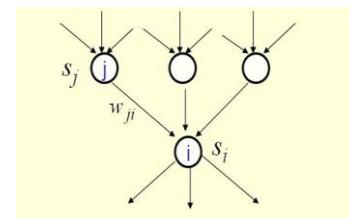
More models

$$D \begin{matrix} \text{document} \rightarrow \text{word} \\ V \end{matrix} = D \begin{matrix} \Theta \\ K \end{matrix} K \begin{matrix} \Phi \\ V \end{matrix}$$

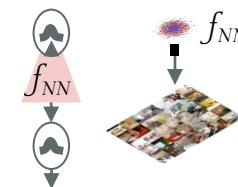
Topic Models; Probabilistic Matrix Factorization  
for text analysis, Recommendor Sys.



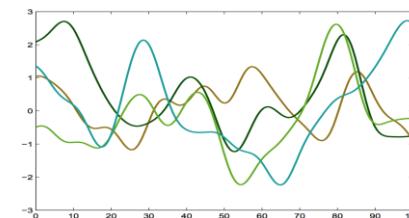
Bayesian Neural Networks  
for uncertainty in deep nets



Deep Belief Nets  
Unsupervised pre-training  
for DNNs



Variational Auto-encoder  
for structured data generation, discrete representation  
learning, semi-supervised learning, etc.



Gaussian Processes  
Nonlinear regression problems

[github.com/thu-ml/zhusuan](https://github.com/thu-ml/zhusuan)

# “ZhuSuan” Deep Probabilistic Programming

- Open-sourced in GitHub:

<https://github.com/thu-ml/zhusuan>



The screenshot shows the GitHub repository page for 'zhusuan'. At the top is the repository logo, a stylized orange 'Z' made of dots. Below it is the repository name 'ZHUSUAN' in a large, bold, white font. Underneath is a 'latest' tag. A search bar labeled 'Search docs' is present. The main content area has a dark background with white text. On the left is a sidebar with the following menu items:

- TUTORIALS
  - Tutorial slides
  - Variational Autoencoders
    - Build the model
    - Reuse the model
    - Inference and learning
    - Generate images
    - Run gradient descent
  - Basic Concepts in ZhuSuan
  - Bayesian Neural Networks
  - Logistic Normal Topic Models
- API DOCS
  - zhusuan.distributions
  - zhusuan.framework
  - zhusuan.variational
  - zhusuan.hmc
  - zhusuan.sgmcmc
  - zhusuan.evaluation
  - zhusuan.transform
  - zhusuan.diagnostics
  - zhusuan.utils
  - zhusuan.legacy
- COMMUNITY
  - Contributing

Docs » Variational Autoencoders

Edit on GitHub

## Variational Autoencoders

Variational Auto-Encoders (VAE) [VAEKW13] is one of the most widely used deep generative models. In this tutorial, we show how to implement VAE in ZhuSuan step by step. The full script is at [examples/variational\\_autoencoders/vae.py](#).

The generative process of a VAE for modeling binarized [MNIST](#) data is as follows:

$\begin{aligned} & \text{import } \text{zhusuan} \text{ as } \text{zs} \\ & \text{def build\_gen}(\text{x\_dim}, \text{z\_dim}, \text{n}, \text{n\_particles}=1): \\ & \quad \text{bn} = \text{zs}.\text{BayesianNet}() \end{aligned}$

This generative process is a stereotype for deep generative models, which starts with a latent representation ( $\text{(}\langle\text{z}\rangle\text{)}$ ) sampled from a simple distribution (such as standard Normal). Then the samples are forwarded through a deep neural network ( $\text{(}\langle\text{f\_NN}\rangle\text{)}$ ) to capture the complex generative process of high dimensional observations such as images. Finally, some noise is added to the output to get a tractable likelihood for the model. For binarized MNIST, the observation noise is chosen to be Bernoulli, with its parameters output by the neural network.

### Build the model

In ZhuSuan, a model is constructed using `BayesianNet`, which describes a directed graphical model, i.e., Bayesian networks. The suggested practice is to wrap model construction into a function (we shall see the meanings of these arguments soon):

```
import zhusuan as zs
def build_gen(x_dim, z_dim, n, n_particles=1):
    bn = zs.BayesianNet()
```

Following the generative process, first we need a standard Normal distribution to generate the latent representations ( $\text{(}\langle\text{z}\rangle\text{)}$ ). As presented in our graphical model, the data is generated in batches with batch size `n`, and for each data, the latent representation is of dimension `z_dim`. So we add a stochastic node by `bn.normal` to generate samples of shape `[n, z_dim]`:

```
# z ~ N(z|0, I)
z_mean = tf.zeros([n, z_dim])
z = bn.normal("z", z_mean, std=1., group_ndims=1, n_samples=n_particles)
```

The method `bn.normal` is a helper function that creates a `Normal` distribution and adds a stochastic node that follows this distribution to the `BayesianNet` instance. The returned `z` is a `StochasticTensor`, which is Tensor-like and can be mixed with Tensors and fed into almost any Tensorflow primitives.

#### Note

To learn more about `Distribution` and `BayesianNet`, Please refer to [Basic Concepts in ZhuSuan](#).

## ZhuSuan: A Library for Bayesian Deep Learning

[J. Shi, J. Chen, J. Zhu, S. Sun, Y. Luo, Y. Gu, Y. Zhou](#)

arXiv preprint, arXiv:1709.05870 , 2017

## Online Documents:

- <http://zhusuan.readthedocs.io/>

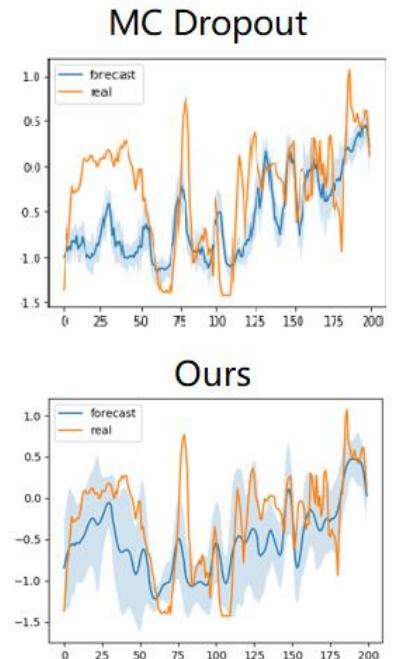
# **Examples**

(what problems can be solved by Bayesian deep learning?)

# Examples: Air Quality Prediction

- ◆ Bayesian inference for a baseline model with LSTM and attention
  - Calculate the uncertainty of prediction
  - Decrease the prediction error of the baseline model

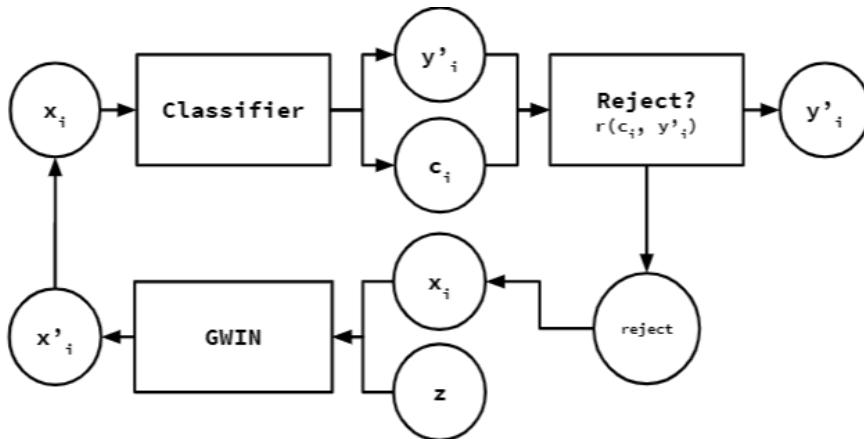
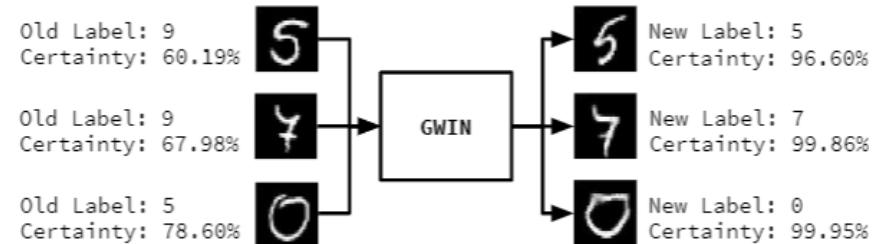
MSE ( $\text{NO}_2$ )	BNN	Baseline NN
+1 hr	<b>0.145</b>	0.160
+7 hr	<b>0.371</b>	0.423
+16 hr	<b>0.389</b>	0.508



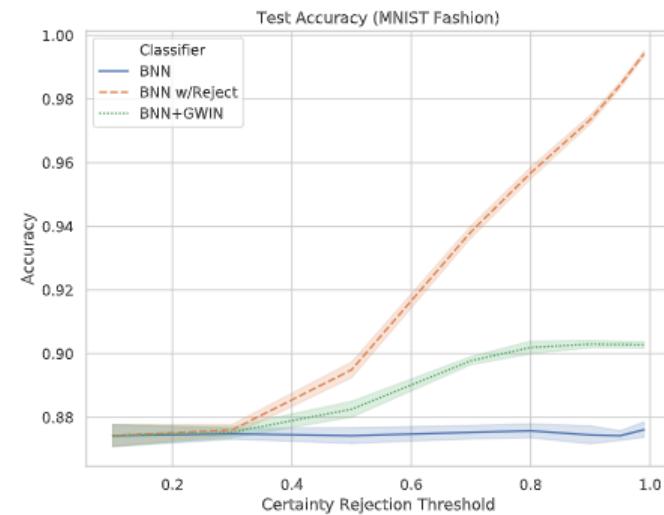
[Function space particle optimization for Bayesian neural networks. Wang et al., ICLR 2019]

# Example: Learning with Rejection

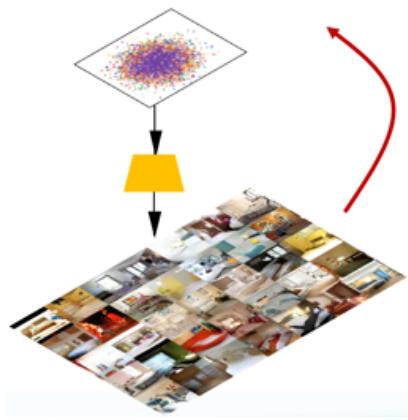
- ◆ Leverage deep generative models to transfer “rejected” samples into high-confidence region



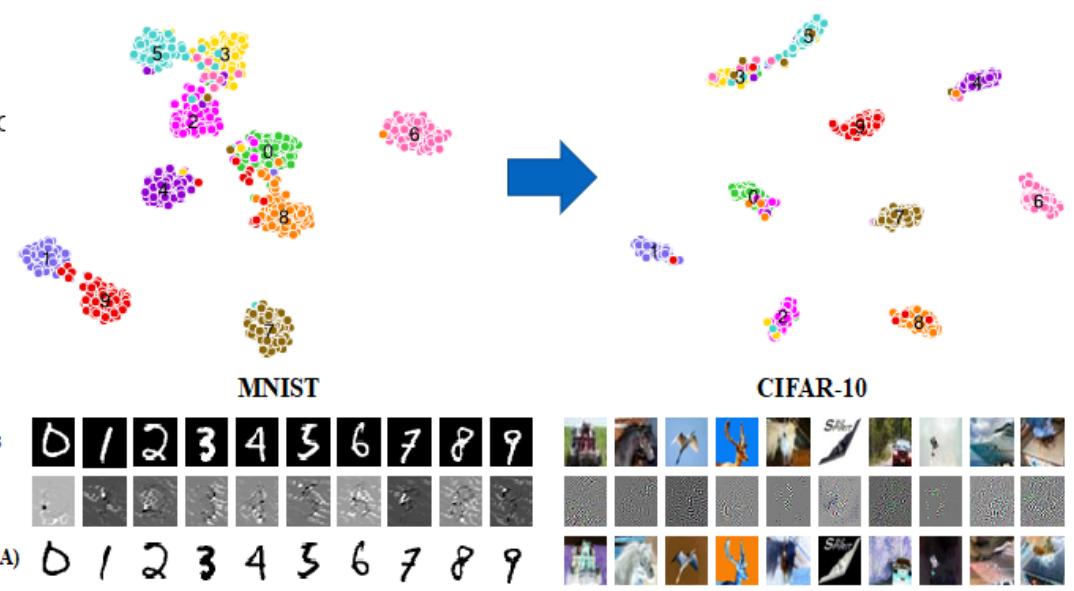
Generative Well-intentioned Networks (GWIN)  
(Cosentino & Zhu, NeurIPS 2019)



# Examples: Adversarial Robustness



Learn a deep network to  
map the data  
into a simple MoG  
distribution



- Theoretically show optimal robustness
- Algorithmically train MM-LDA network using SGD to minimize cross-entropy loss

No extra cost, as compared to vanilla DNN

# Thanks!



**ZhuSuan: A Library for Bayesian Deep Learning**  
J. Shi, et al., arXiv preprint, arXiv:1709.05870 , 2017

**Online Documents:** <http://zhusuan.readthedocs.io/>