

# Deep Generative Models

CS772A: Probabilistic Machine Learning

Piyush Rai

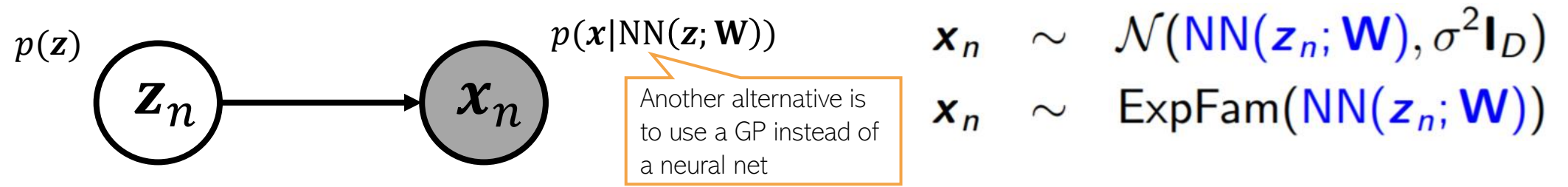
# Plan

- Variational Autoencoders
- Generative Adversarial Networks
- Denoising Diffusion Models



# Constructing Generative Models using Neural Nets<sup>3</sup>

- We can use a neural net to define the mapping from a  $K$ -dim  $\mathbf{z}_n$  to  $D$ -dim  $\mathbf{x}_n$



- If  $\mathbf{z}_n$  has a Gaussian prior, such models are called **deep latent Gaussian models** (DLGM)
- Since NN mapping can be very powerful, DLGM can generate very high-quality data
  - Take the trained network, generate a random  $\mathbf{z}$  from prior, pass it through the model to generate  $\mathbf{x}$

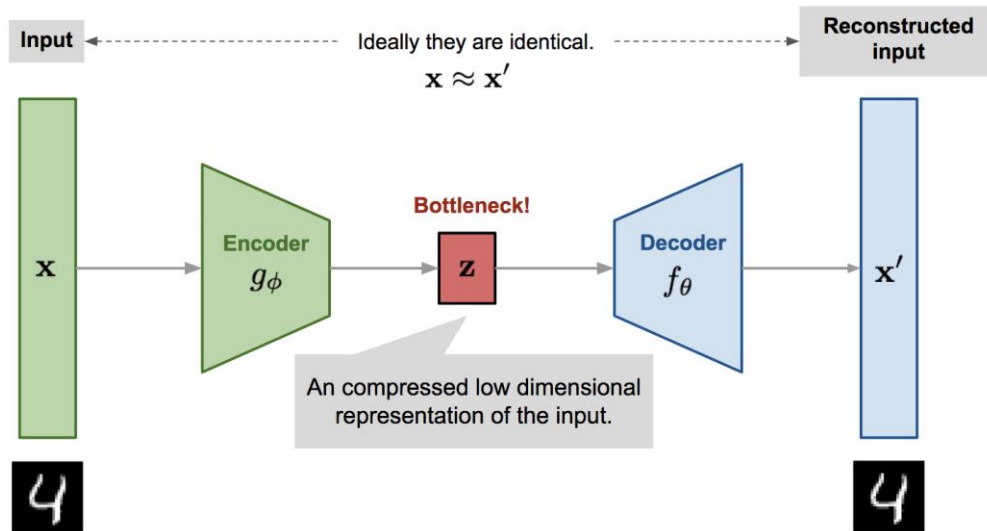


Some sample images generated by Vector Quantized Variational Auto-Encoder (VQ-VAE), a state-of-the-art DLGM



# Variational Autoencoder (VAE)

- VAE\* is a probabilistic extension of autoencoders (AE)



$$L_{AE}(\theta, \phi) = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}^{(i)} - f_{\theta}(g_{\phi}(\mathbf{x}^{(i)})))^2$$

- The basic difference is that VAE assumes a prior  $p(\mathbf{z})$  on the latent code  $\mathbf{z}$ 
  - This enables it to not just compress the data but also generate synthetic data
  - How: Sample  $\mathbf{z}$  from a prior and pass it through the decoder
- Thus VAE can learn good latent representation + generate novel synthetic data
- The name has “Variational” in it since it is learned using VI principles



# Variational Autoencoder (VAE)

- VAE has three main components
  - A prior  $p_{\theta}(\mathbf{z})$  over latent codes
  - A probabilistic decoder/generator  $p_{\theta}(\mathbf{x}|\mathbf{z})$ , modeled by a deep neural net
  - A posterior or probabilistic encoder  $p_{\theta}(\mathbf{z}|\mathbf{x})$  approx. by an “inference network”  $q_{\phi}(\mathbf{z}|\mathbf{x})$

Here  $\theta$  collectively denotes all the parameters of the prior and likelihood

Using the idea of “Amortized Inference” (next slide)

- VAE is learned by maximizing the ELBO

ELBO for a single data point

$$\begin{aligned} \mathcal{L}(\theta, \phi | \mathbf{x}) &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\theta}(\mathbf{z})) \end{aligned}$$

Here  $\phi$  collectively denotes all the parameters that define the inference network

Maximized to find the optimal  $\theta$  and  $\phi$

$q_{\phi}$  should be such that data  $\mathbf{x}$  is reconstruct well from  $\mathbf{z}$  (high log-lik)

$q_{\phi}$  should also be simple (close to the prior)

- The [Reparametrization Trick](#) is commonly used to optimize the ELBO
- Posterior is inferred only over  $\mathbf{z}$ , and usually only point estimate on  $\theta$  and  $\phi$



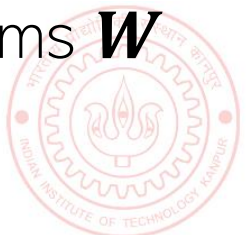
# Amortized Inference

- Latent variable models need to infer the posterior  $p(\mathbf{z}_n | \mathbf{x}_n)$  for each observation  $\mathbf{x}_n$
- This can be slow if we have lots of observations because
  1. We need to iterate over each  $p(\mathbf{z}_n | \mathbf{x}_n)$
  2. Learning the global parameters needs wait for step 1 to finish for all observations
- One way to address this is via Stochastic VI
- Amortized inference is another appealing alternative (used in VAE and other LVMs too)

$$p(\mathbf{z}_n | \mathbf{x}_n) \approx q(\mathbf{z}_n | \phi_n) = q(\mathbf{z}_n | \text{NN}(\mathbf{x}_n; \mathbf{W}))$$

If  $q$  is Gaussian then the NN will output a mean and a variance

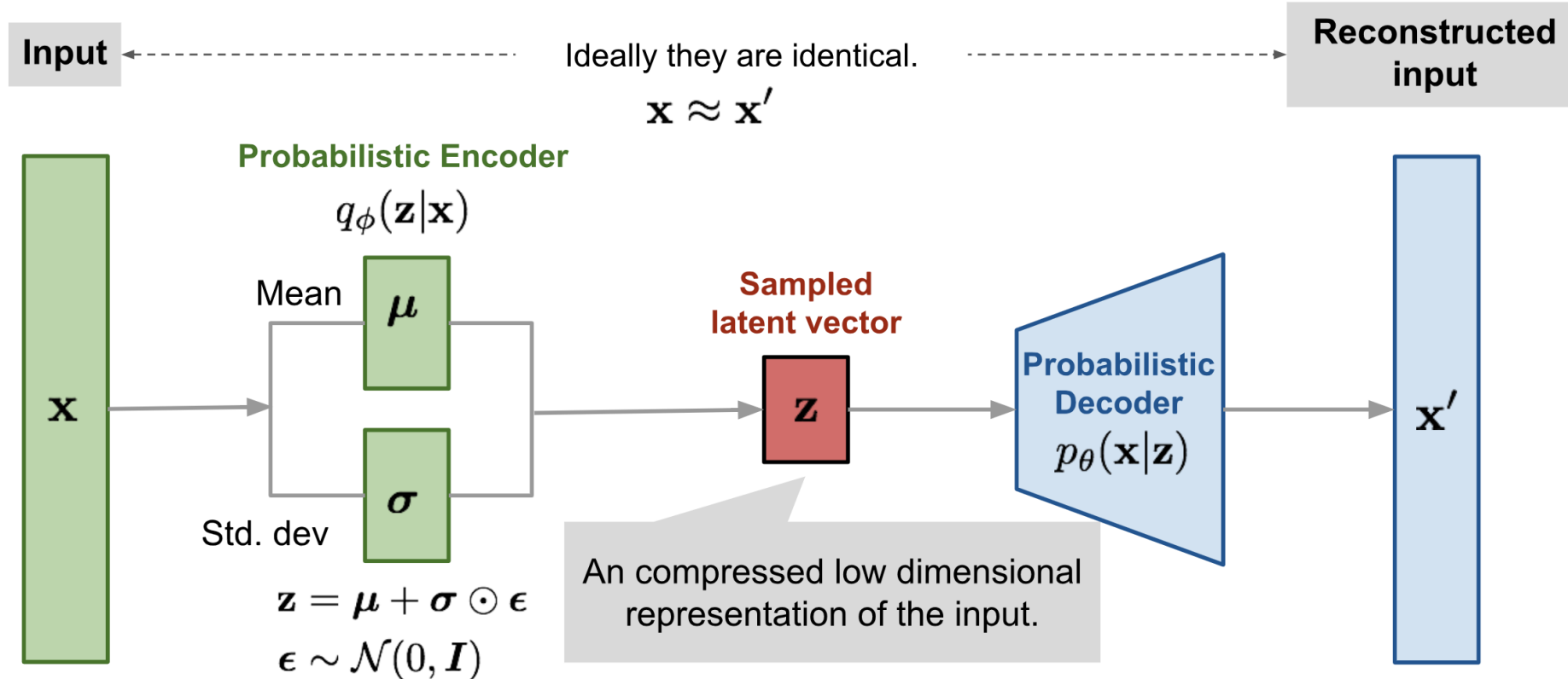
- Thus no need to learn  $\phi_n$ 's (one per data point) but just a single NN with params  $\mathbf{W}$ 
  - This will be our “encoder network” for learning  $\mathbf{z}_n$
  - Also very efficient to get  $p(\mathbf{z}_* | \mathbf{x}_*)$  for a new data point  $\mathbf{x}_*$



# Variational Autoencoder: The Complete Pipeline

- Both probabilistic encoder and decoder learned jointly by maximizing the ELBO

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}, \phi | \mathbf{x}) &= \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z} | \mathbf{x})] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} [\log p_{\boldsymbol{\theta}}(\mathbf{x} | \mathbf{z})] - \text{KL}(q_{\phi}(\mathbf{z} | \mathbf{x}) || p_{\boldsymbol{\theta}}(\mathbf{z}))\end{aligned}$$





# VAE and Posterior Collapse

- VAEs may suffer from **posterior collapse**

Decoder is a neural net and can be arbitrarily powerful making this term very large

Consequently, KL will become close to zero collapsing posterior to the prior

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

- Thus, due to posterior collapse, reconstruction will still be good but the code  $\mathbf{z}$  may be garbage (not useful as a representation for  $\mathbf{x}$ )
- Several ways to prevent posterior collapse, e.g.,

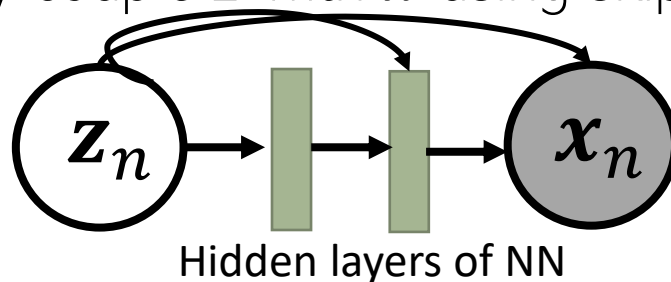
- Use KL annealing

A carefully tuned value between 0 and 1

$$\mathcal{L}(\theta, \phi | \mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} [\log p_{\theta}(\mathbf{x} | \mathbf{z})] - \beta \text{KL}(q_{\phi}(\mathbf{z} | \mathbf{x}) || p_{\theta}(\mathbf{z}))$$

For example, keep the variance of  $q$  as fixed

- Avoid KL from becoming 0 using some  $q$  that doesn't collapse to the prior
- More tightly couple  $\mathbf{z}$  with  $\mathbf{x}$  using skip-connections (Skip-VAE)



Besides these, MCMC (sometimes used for inference in VAE), or improved VI techniques can also help in preventing posterior collapse in VAEs

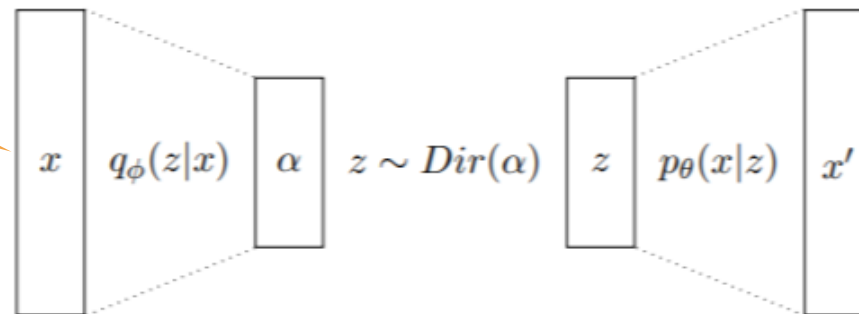




# VAE: Some Comments

- One of the state-of-the-art latent variable models
- Useful for both generation as well as representation learning
- Many improvements and extensions, e.g.,
  - For text data and sequences (VAE for topic models or “neural topic models”)

Document (e.g., as a vector of word counts)



(a) Dirichlet VAE (DVAE)

- VAE-style models with more than one layer of latent variables (Sigmoid Belief Networks, hierarchical VAE, Ladder VAE, Deep Exponential Families, etc)

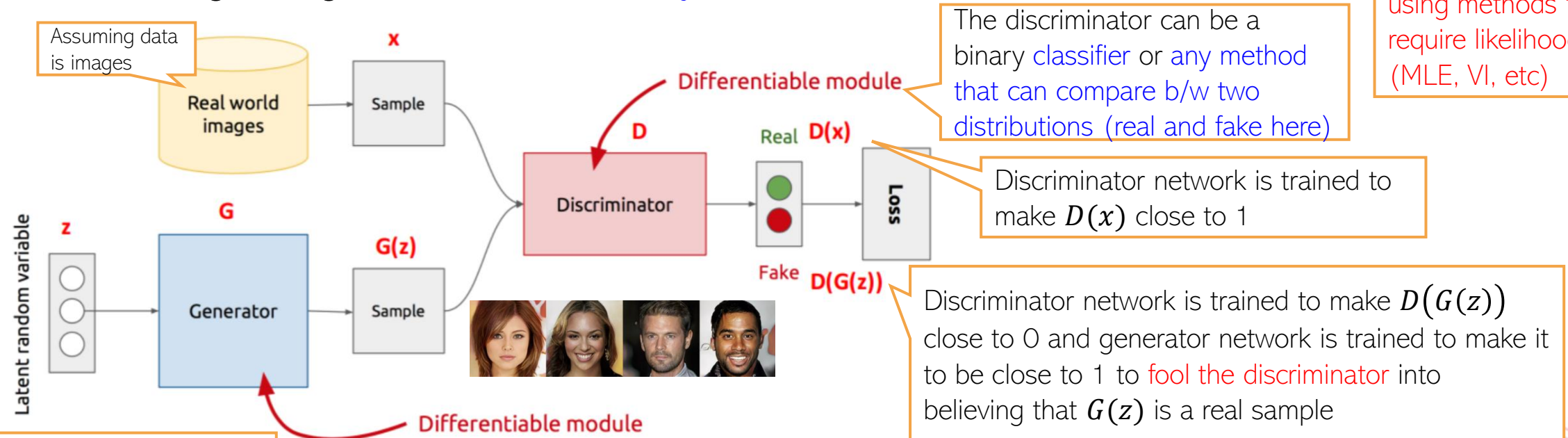


# Generative Adversarial Network (GAN)

- GAN is an implicit generative latent variable model
- Can generate from it but can't compute  $p(\mathbf{x})$  - the model doesn't define it explicitly
- GAN is training using an **adversarial way** (Goodfellow et al, 2013)

Unlike VAE, no explicit parametric likelihood model  $p(\mathbf{x}|\mathbf{z})$

Thus can't train using methods that require likelihood (MLE, VI, etc)



Min-max optimization

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$

# Generative Adversarial Network (GAN)

- The GAN training criterion was

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

- With  $G$  fixed, the optimal  $D$  (exercise)

$$D_G^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}$$

Distribution of real data
Distribution of synthetic data

- Given the optimal  $D$ , The optimal generator  $G$  is found by minimizing

$$V(D_G^*, G) = \mathbb{E}_{\mathbf{x} \sim p_{data}} \left[ \log \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p_g} \left[ \log \frac{p_g(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})} \right]$$

$$= \text{KL} \left[ p_{data}(x) \left\| \frac{p_{data}(x) + p_g(x)}{2} \right. \right] + \text{KL} \left[ p_g(x) \left\| \frac{p_{data}(x) + p_g(x)}{2} \right. \right] - \log 4$$

Jensen-Shannon divergence between  $p_{data}$  and  $p_g$ .  
Minimized when  $p_g = p_{data}$

Thus GAN can learn the true data distribution if the generator and discriminator have enough modeling power



# GAN Optimization

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

- The GAN training procedure can be summarized as

- 1 Initialize  $\theta_g, \theta_d$ ;
 

$\theta_g$  and  $\theta_d$  denote the params of the deep neural nets defining the generator and discriminator, respectively
- 2 **for** *each training iteration* **do**

In practice, for stable training, we run  $K > 1$  steps of optimizing w.r.t.  $D$  and 1 step of optimizing w.r.t.  $G$

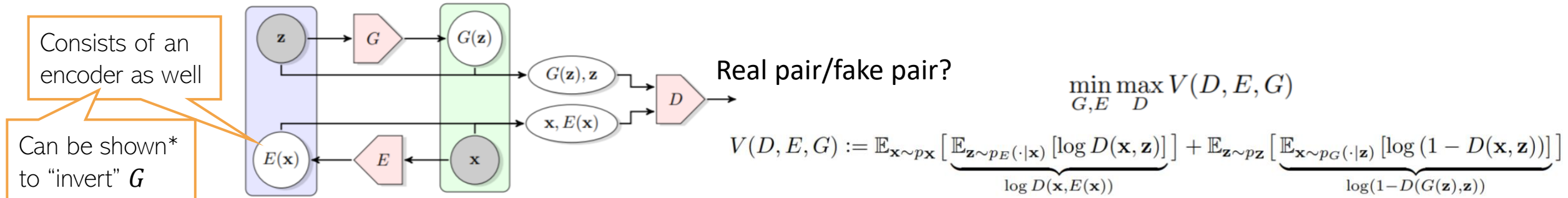
  - 3 **for**  $K$  steps **do**
  - 4     Sample minibatch of  $M$  noise vectors  $\mathbf{z}_m \sim q_z(\mathbf{z})$ ;
  - 5     Sample minibatch of  $M$  examples  $\mathbf{x}_m \sim p_D$ ;
  - 6     Update the discriminator by performing stochastic gradient *ascent* using this gradient:
 
$$\nabla_{\theta_d} \frac{1}{M} \sum_{m=1}^M [\log D(\mathbf{x}_m) + \log(1 - D(G(\mathbf{z}_m)))] . ;$$
  - 7     Sample minibatch of  $M$  noise vectors  $\mathbf{z}_m \sim q_z(\mathbf{z})$ ;
  - 8     Update the generator by performing stochastic gradient *descent* using this gradient:
 
$$\nabla_{\theta_g} \frac{1}{M} \sum_{m=1}^M \log(1 - D(G(\mathbf{z}_m))) . ;$$

In practice, in this step, instead of minimizing  $\log(1 - D(G(\mathbf{z})))$ , we **maximize**  $\log(D(G(\mathbf{z})))$

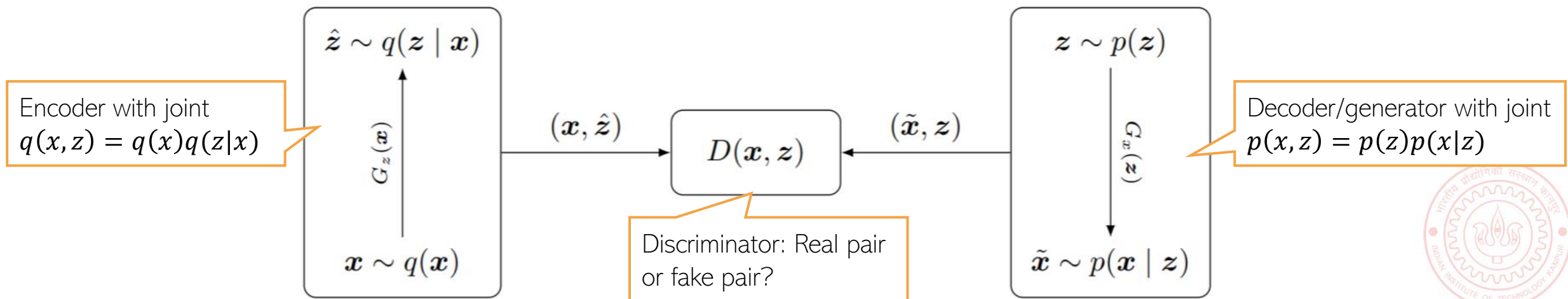
Reason: Generator is bad initially so discriminator will always predict correctly initially and  $\log(1 - D(G(\mathbf{z})))$  will saturate
- 9 Return  $\theta_g, \theta_d$

# GANs that also learn latent representations

- The standard GAN can only generate data. Can't learn the latent  $\mathbf{z}$  from  $\mathbf{x}$
- Bidirectional GAN\* (BiGAN) is a GAN variant that allows this



- Adversarially Learned Inference# (ALI) is another variant that can learn representations



# Evaluating GANs

- Two measures that are commonly used to evaluate GANs
  - Inception score (IS): Evaluates the distribution of generated data
  - Frechet inception distance (FID): Compared the distribution of real data and generated data
- Inception Score defined as  $\exp(\mathbb{E}_{x \sim p_g} [\text{KL}(p(y|x) || p(y))])$  will be high if
  - Very few high-probability classes in each sample  $x$ : Low entropy for  $p(y|x)$
  - We have diverse classes across samples: Marginal  $p(y)$  is close to uniform (high entropy)
- FID uses extracted features (using a deep neural net) of real and generated data
  - Usually from the layers closer to the output layer
- These features are used to estimate two Gaussian distributions

High IS and low FID is desirable

Both IS and FID measure how realistic the generated data is

Using real data  $\mathcal{N}(\mu_R, \Sigma_R)$

$\mathcal{N}(\mu_G, \Sigma_G)$  Using generated data

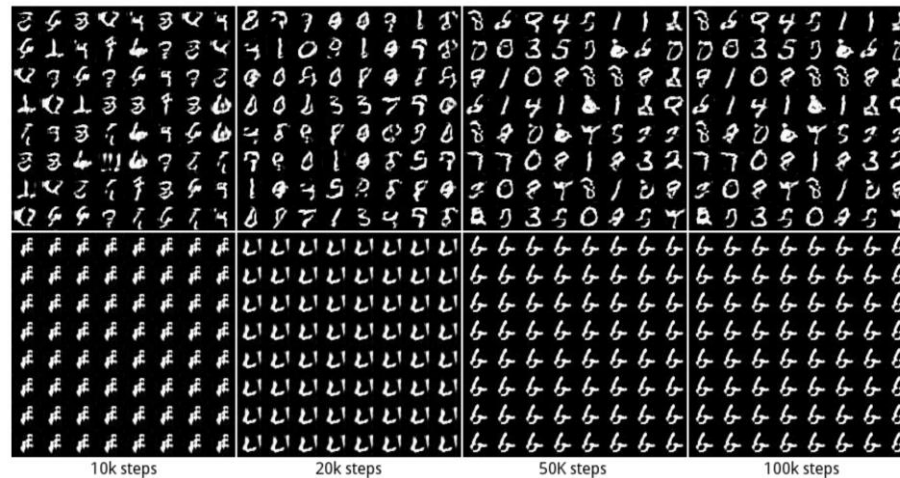
- FID is then defined as  $\text{FID} = |\mu_G - \mu_R|^2 + \text{trace}(\Sigma_G + \Sigma_R - (\Sigma_G \Sigma_R)^{1/2})$





# GAN: Some Issues/Comments

- GAN training can be hard and the basic GAN suffers from several issues
- Instability of training procedure
- Mode Collapse problem: Lack of diversity in generated samples
  - Generator may find some data that can easily fool the discriminator
  - It will stuck at that mode of the data distribution and keep generating data like that



GAN 1: No mode collapse (all 10 modes captured in generation)

GAN 2: Mode collapse (stuck on one of the modes)

- Some work on addressing these issues (e.g., [Wasserstein GAN](#), [Least Squares GAN](#), etc)

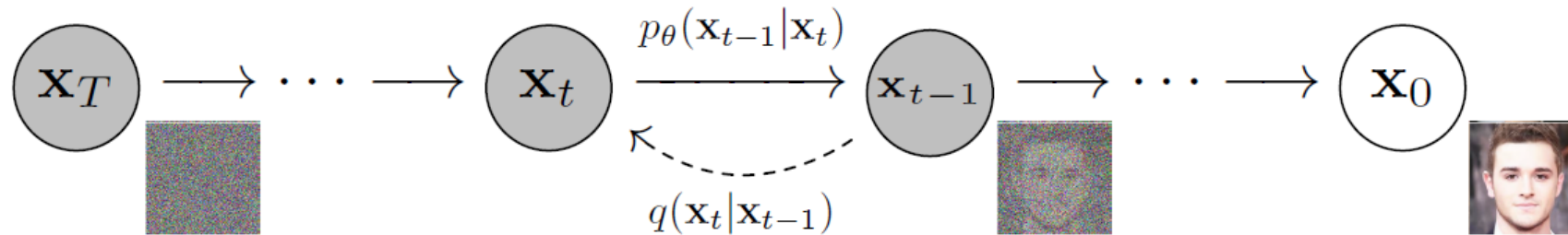




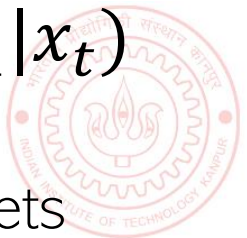
# Denoising Diffusion Models

After learning the model, can use the reverse process to generate data from random noise

- Based on a forward (adding noise) process and a reverse (denoising) process



- Steps of the forward process are defined by a **fixed** Gaussian  $q(\mathbf{x}_t|\mathbf{x}_{t-1})$ 
  - The f.p. starts with the clean image  $\mathbf{x}_0$  and adds zero-mean Gaussian noise at each step
  - The f.p. distribution is defined as  $q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t|\sqrt{(1-\beta_t)}\mathbf{x}_{t-1}, \beta_t I)$   $\beta_t \in (0,1)$
  - Eventually as  $T \rightarrow \infty$ , we get  $\mathbf{x}_T$  which is isotropic Gaussian noise
  - Can show:  $q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t|\sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1-\bar{\alpha}_t)I)$  where  $\alpha_t = 1 - \beta_t$  and  $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$
- Steps of the reverse process are defined by a **learnable** Gaussian  $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ 
  - $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$  is an approximation of the **reverse diffusion**  $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$
  - $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$  modeled as  $\mathcal{N}(\mathbf{x}_{t-1}|\mu_\theta(\mathbf{x}_t), \Sigma_\theta(\mathbf{x}_t))$  where  $\mu_\theta$  and  $\Sigma_\theta$  are neural nets



# Denoising Diffusion Models: Training

- The model is trained by minimizing the following objective

$$\mathbb{E}[-\log p_\theta(x_0)] \leq \mathbb{E} \left[ -\log \frac{p_\theta(x_{0:T})}{q(x_{1:T}|x_0)} \right] := \mathcal{L}$$

Upper bound on the negative log-likelihood (negative of the ELBO)

$$p_\theta(\mathbf{x}_{0:T}) := p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) := p(\mathbf{x}_T) \prod_{t=1}^T \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t))$$

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) := \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}) := \prod_{t=1}^T \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$$

$$\mathcal{L} = L_0 + L_1 + L_2 + \dots + L_{T-1} + L_T$$

$$L_0 = -\log p_\theta(x_0|x_1)$$

This is also a Gaussian

$$L_{t-1} = D_{KL}(q(x_{t-1}|x_t, x_0) || p_\theta(x_{t-1}|x_t))$$

$$L_T = D_{KL}(q(x_T|x_0) || p(x_T))$$

Overall loss is just a sum of several KL divergences between Gaussians, and thus available in closed form

- In some ways, denoising diffusion models are similar to VAEs

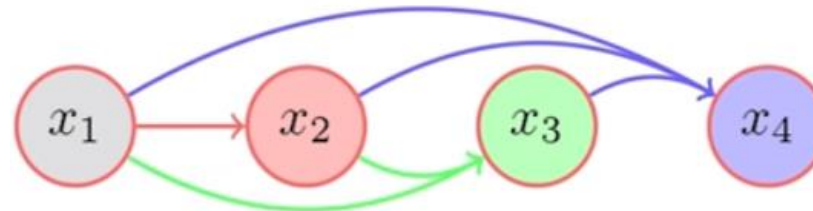


# Summary

- Looked at various methods for generative modeling for unsupervised learning
  - Classical methods (FA, PPCA, other latent factor models, topic models, etc)
  - Deep generative models (VAE, GAN, Denoising Diffusion Models)
- Many of these methods can also be extended to model data other than images
- There are also generative models that do not use latent variables
  - Can still be used to generate data and learn the underlying data distribution

Assuming each observation is n-dimensional

$$\prod_{i=1}^n p(x_i | \mathbf{x}_{<k})$$



An auto-regressive model

Can use a neural network to learn (parameters of) each of these distributions

An example: Neural Autoregressive Density Estimator (NADE)

