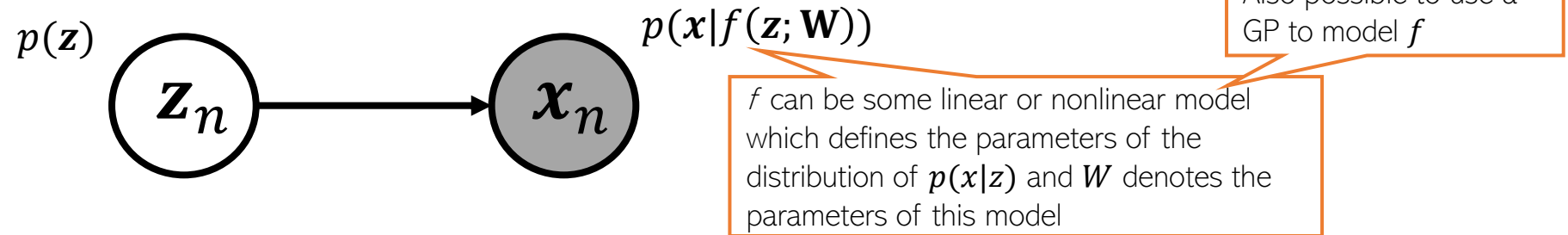


# VAE and GAN

# Latent Variable Models for Generation Tasks

- Assume a  $K$ -dim latent variable  $\mathbf{z}_n$  is transformed to generate to  $D$ -dim observation  $\mathbf{x}_n$

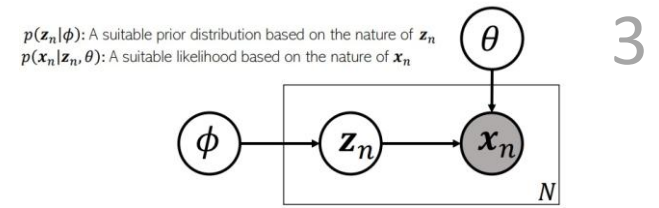


- It is common to use a Gaussian prior for  $\mathbf{z}_n$  (though other priors can be used)
- If we use a neural net or GP, such models can generate very high-quality data
  - Take the trained network, generate a random  $\mathbf{z}$  from prior, pass it through the model to generate  $\mathbf{x}$



Some sample images generated by Vector Quantized Variational Auto-Encoder (VQ-VAE), a state-of-the-art latent variable model for generation

# Factor Analysis and Probabilistic PCA



- FA and PPCA assume  $f$  to be a linear model
- In FA/PPCA, latent variables  $\mathbf{z}_n \in \mathbb{R}^K$  typically assumed to have a Gaussian prior
  - If we want sparse latent variable, can use Laplace prior on  $\mathbf{z}_n$
  - More complex extensions of FA/PPCA use a mixture of Gaussians prior on  $\mathbf{z}_n$
- Assumption: Observations  $\mathbf{x}_n \in \mathbb{R}^D$  typically assumed to have a Gaussian likelihood
  - Other likelihood models (e.g., exp-family) can also be used if data not real-valued
- Relationship between  $\mathbf{z}_n$  and  $\mathbf{x}_n$  modeled by a noisy linear mapping

$$\mathbf{x}_n = \mathbf{W}\mathbf{z}_n + \epsilon_n = \sum_{k=1}^K \mathbf{w}_k z_{nk} + \epsilon_n$$

Zero-mean and diagonal or spherical Gaussian noise

Linear combination of the columns of  $\mathbf{W}$

$$p(\mathbf{z}_n) = \mathcal{N}(\mathbf{z}_n | \mathbf{0}, \mathbf{I})$$
$$p(\mathbf{x}_n | \mathbf{z}_n) = \mathcal{N}(\mathbf{x}_n | \mathbf{W}\mathbf{z}_n, \Psi)$$

Diagonal for FA, spherical for PPCA

- Linear Gaussian Model.  $\mathbf{W}$ ,  $\mathbf{z}_n$ 's, and  $\Psi$  can be learned (e.g, using EM, VI, MCMC)

# Some Variants of FA/PPCA

4

## ■ Gamma-Poisson latent factor model

Popular for modeling count-valued data (in text analysis, recommender systems, etc)

Non-negative priors often give a nice interpretability to such latent variable models (will see some more examples of such models shortly)

- Assumes  $K$ -dim non-negative latent variable  $\mathbf{z}_n$  and  $D$ -dim count-valued observations  $\mathbf{x}_n$
- An example: Each  $\mathbf{x}_n$  is the word-count vector representing a document

$$p(\mathbf{z}_n) = \prod_{k=1}^K \text{Gamma}(z_{nk} | a_k, b_k)$$
$$p(\mathbf{x}_n | \mathbf{z}_n) = \prod_{d=1}^D \text{Poisson}(x_{nd} | f(\mathbf{w}_d, \mathbf{z}_n))$$

This is the rate of the Poisson. It should be non-negative,  $\exp(\mathbf{w}_d^T \mathbf{z}_n)$ , or simply  $\mathbf{w}_d^T \mathbf{z}_n$  if  $\mathbf{w}_d$  is also non-negative (e.g., using a gamma/Dirichlet prior on it)

- This can be thought of as a probabilistic non-negative matrix factorization model

## ■ Dirichlet-Multinomial/Multinoulli PCA

- Assumes  $K$ -dim non-negative latent variable  $\mathbf{z}_n$  and  $D$  categorical obs  $\mathbf{x}_n = \{x_{nd}\}_{d=1}^D$
- An example: Each  $\mathbf{x}_n$  is a document with  $D$  words in it (each word is a categorical value)

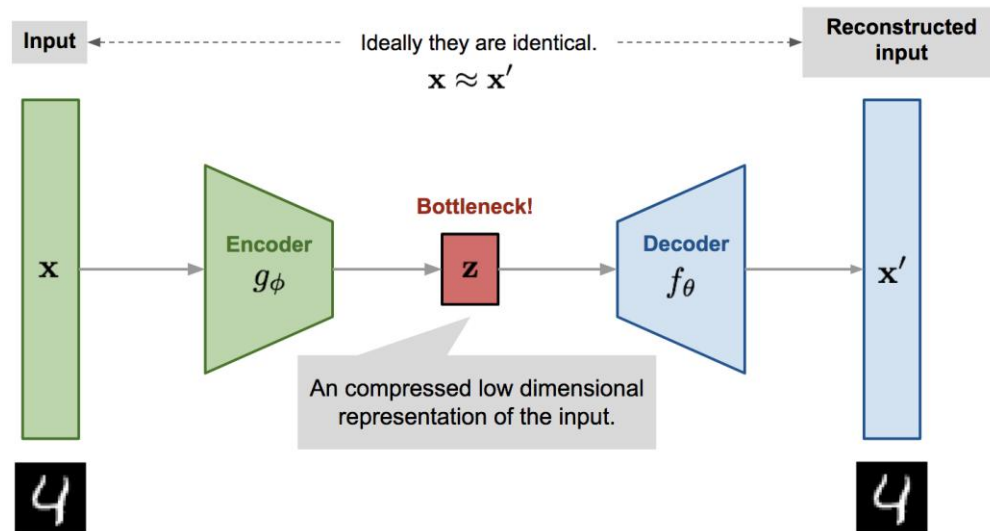
Also sums to 1

$$p(\mathbf{z}_n) = \text{Dirichlet}(\mathbf{z}_n | \boldsymbol{\alpha})$$
$$p(\mathbf{x}_n | \mathbf{z}_n) = \prod_{d=1}^D \text{Multinoulli}(x_{nd} | f(\mathbf{w}_d, \mathbf{z}_n))$$

This should give the probability vector of the multinoulli over  $x_{nd}$ . It should be non-negative and should sum to 1

# A Deep Generative Model: Variational Auto-encoder (VAE)

- VAE\* is a probabilistic extension of autoencoders (AE). An AE is shown below



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

- 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

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

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

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

$q_{\phi}$  should be such that data  $\mathbf{x}$  is reconstructed 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  $\theta$

# Stochastic Variational 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  $N$
  3.  $\mathcal{L}(\phi) = \mathbb{E}_{q_\phi(\mathbf{z})} [\log p(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z})]$  summation over all data points  $\mathbf{x}_i$  is costly for large  $N$
- One way to address this is via Stochastic VI
  - Use mini-batches of data and stochastic gradients to optimize the ELBO.

$$\widehat{\nabla_\phi \mathcal{L}} = \frac{N}{M} \sum_{i \in \text{batch}} \nabla_\phi \mathbb{E}_{q_\phi} [\log p(x_i, z_i) - \log q_\phi(z_i)]$$

# Amortized Inference

- Amortized inference is another appealing alternative (used in VAE and other LVMs too)
- 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  $\mathbf{x}_*$

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

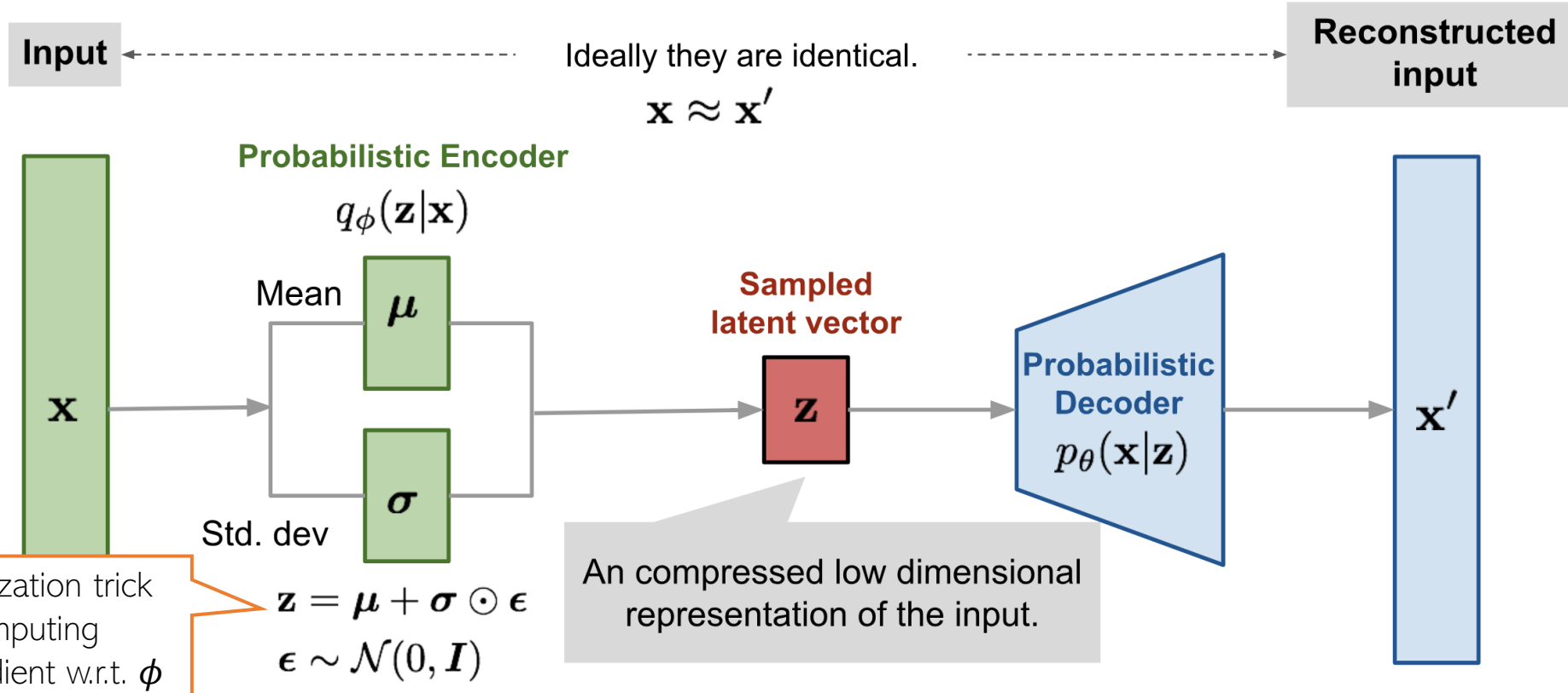
If  $q$  is Gaussian then the NN will output a mean and a variance



# Variational Autoencoder: The Complete Pipeline

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

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

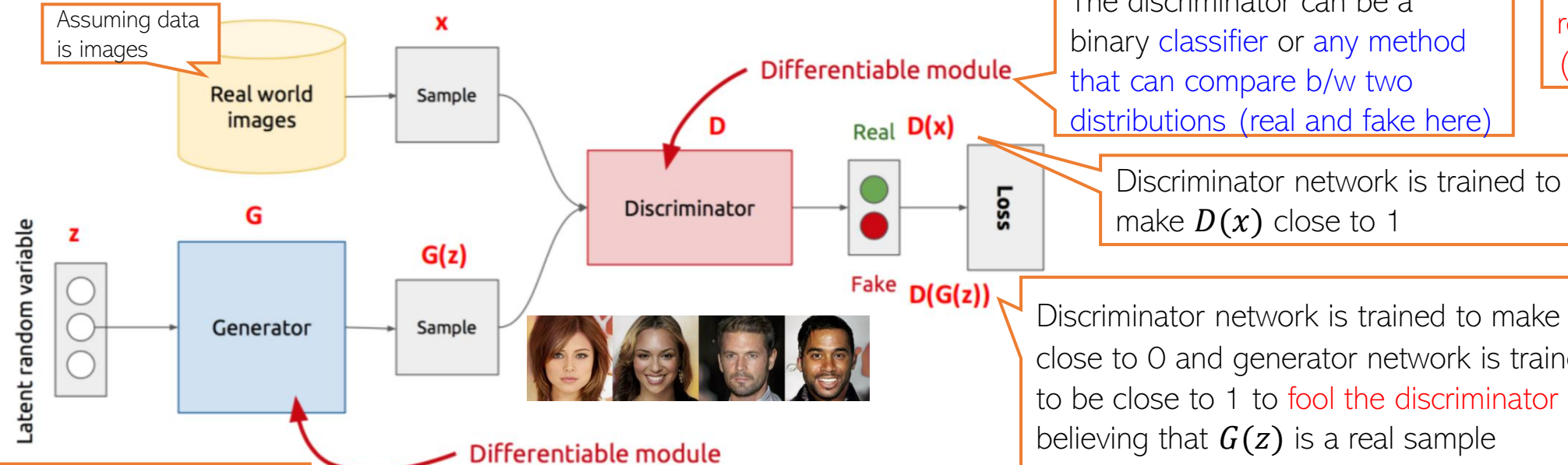


# Generative Adversarial Network (GAN)

- 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 trained 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

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

# 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

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

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

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:

$$\nabla_{\theta_g} \frac{1}{M} \sum_{m=1}^M \log(1 - D(G(\mathbf{z}_m))) ;$$

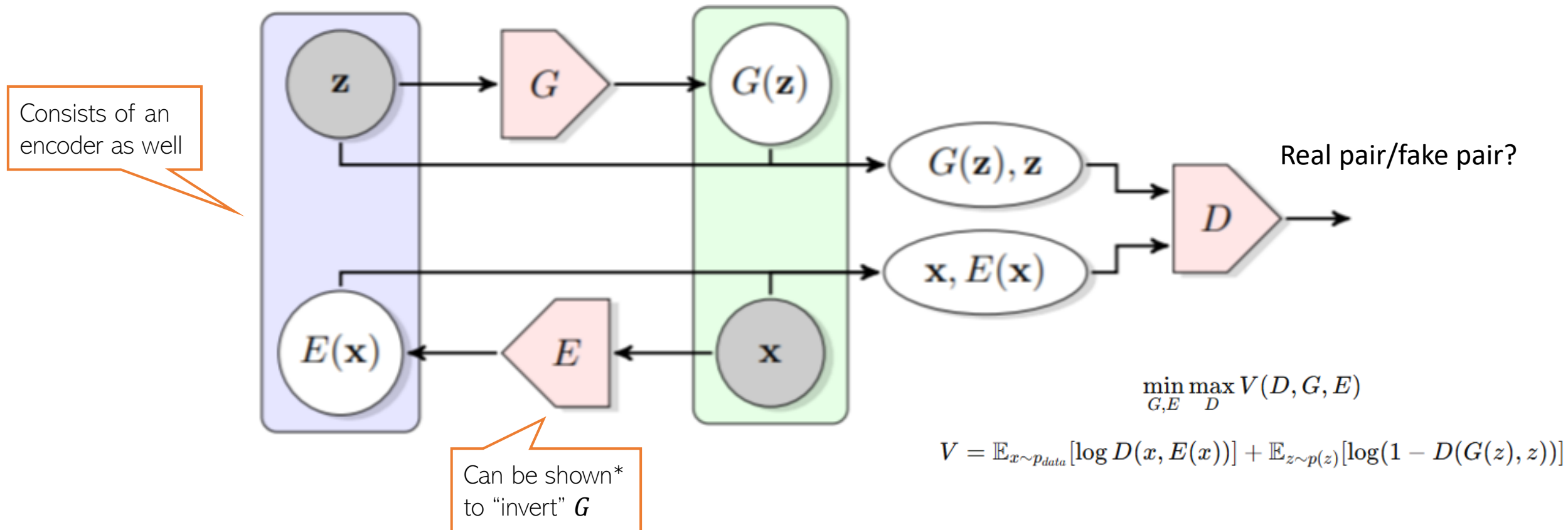
9 Return  $\theta_g, \theta_d$

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

# 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



# 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  $\mathcal{N}(\mu_R, \Sigma_R)$

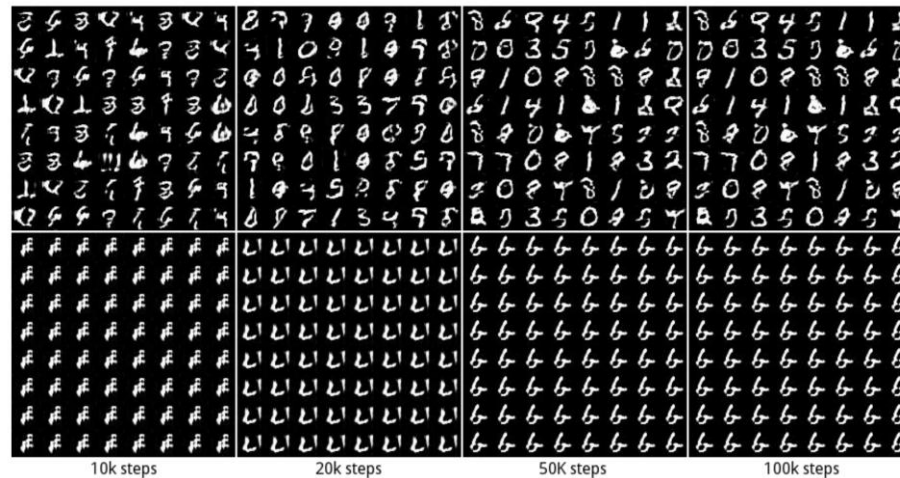
$\mathcal{N}(\mu_G, \Sigma_G)$  Using generated data

- FID is then defined as  $\text{FID} = |\mu_G - \mu_R|^2 + \text{trace}(\Sigma_G + \Sigma_R - (\Sigma_G \Sigma_R)^{1/2})$
- These measures can also be used for evaluating other deep gen models like VAE



# 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)

# References

- Chapter 21 and 26, Kevin Murphy, [Probabilistic Machine Learning: Advanced Topics](#)
- Chapter 19.2 and 17, Christopher N. Bishop, Deep Learning: Foundations and Concepts