

# Introduction to Deep Generative Modeling

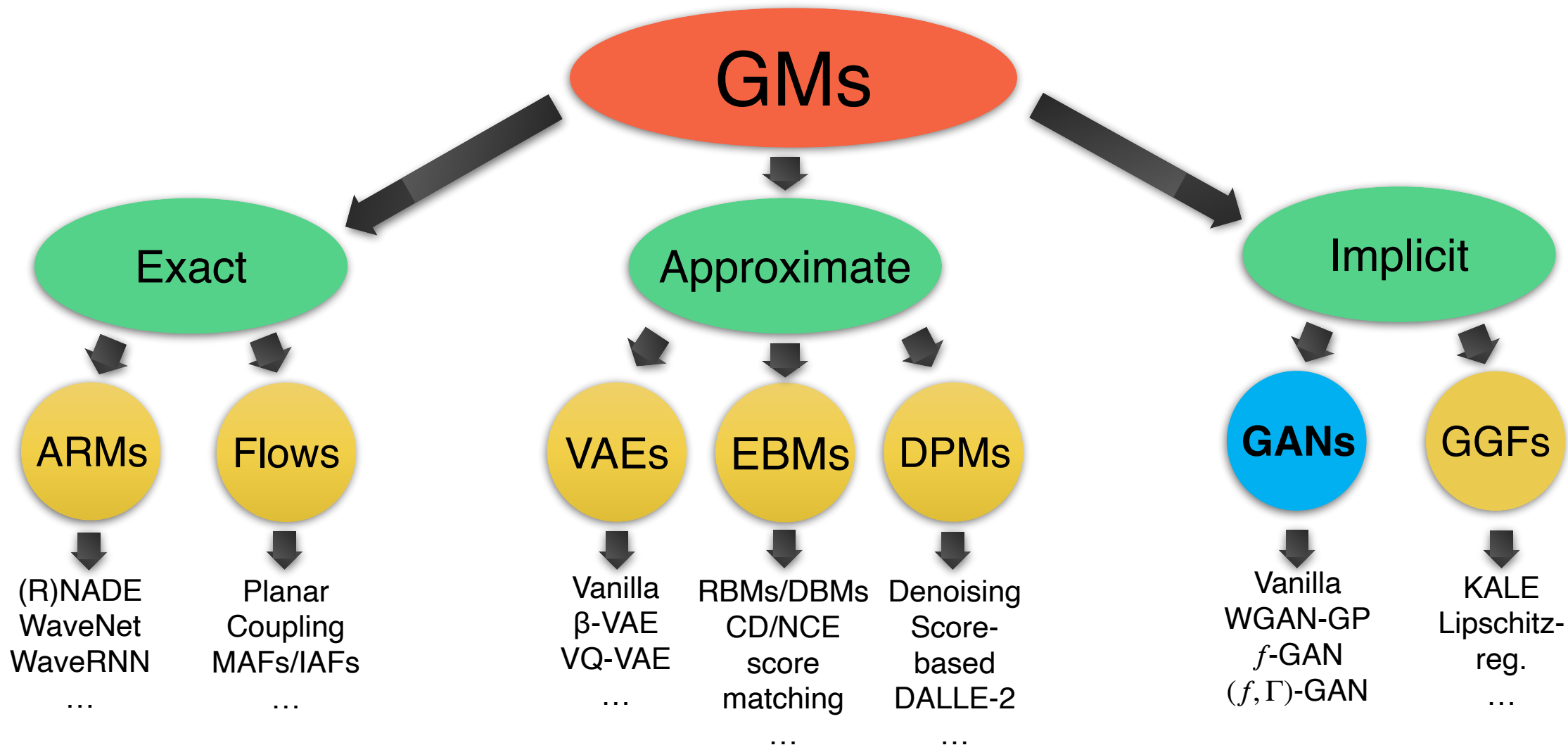
Lecture #15

**HY-673** – Computer Science Dep., University of Crete

Professors: Yannis Pantazis, Yannis Stylianou

Teaching Assistant: Michail Raptakis

# Taxonomy of GMs



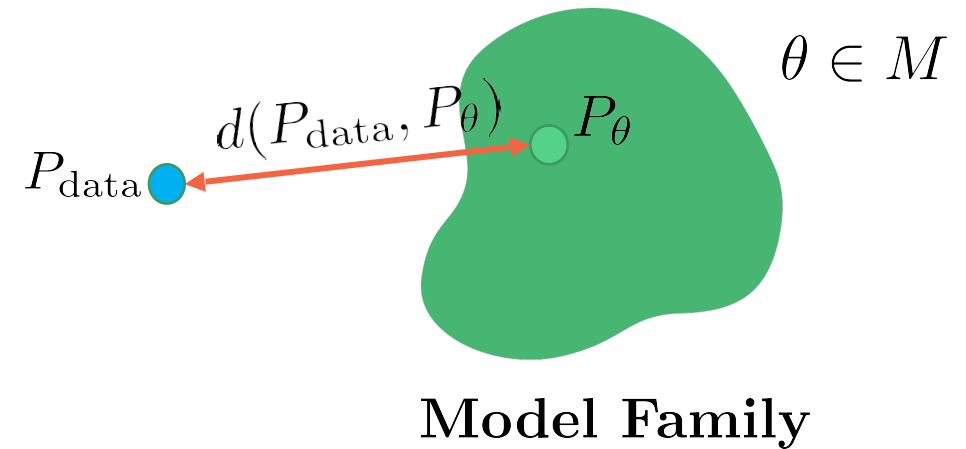
# Generative Adversarial Networks

Lecture #15

Recap:



$$x_i \sim P_{\text{data}} \\ i = 1, 2, \dots, n$$



- Generative Model families

- Autoregressive Models:  $p_{\theta}(x) = \prod_{i=1}^n p_{\theta}(x_i | x_{<i})$ .
- Variational Autoencoders:  $p_{\theta}(x) = \int p_{\theta}(x, z) dz$ .
- Normalizing Flow Models:  $p_X(x; \theta) = p_Z(f_{\theta}^{-1}(x)) \left| \det \left( \frac{\partial f_{\theta}^{-1}(x)}{\partial x} \right) \right|$ .
- Diffusion Probabilistic Models:  $p_{\theta}(x) = p_{\theta}(x | x_1) \prod_{t=2}^T p_{\theta}(x_{t-1} | x_t) p(x_T)$ .

$$\hat{\theta} = \arg \max_{\theta} \sum_{i=1}^n \log p_{\theta}(x_i), \quad x_1, x_2, \dots, x_n \sim p_{\text{data}}(x).$$

- **Optimal statistical efficiency**

- Assume sufficient model capacity, such that there exists a unique  $\theta^* \in \mathcal{M}$  that satisfies  $p_{\theta^*} = p_{\text{data}}$ .
- The convergence of  $\hat{\theta}$  to  $\theta^*$  when  $n \rightarrow \infty$  is the “fastest” among all statistical methods when using maximum likelihood training.

- **Higher likelihood  $\equiv$  better lossless compression.**

- Is the likelihood a good indicator of the quality of samples generated by the model?

- **Case 1:** Optimal generative model will give best **sample quality** and highest test **log-likelihood**.
- For imperfect models, achieving high log-likelihoods might not always imply good sample quality, and vice-versa (Theis et al., 2016)
- **Case 2:** Great test log-likelihoods, poor samples. E.g., For a discrete noise mixture model  $p_{\theta}(x) = 0.01p_{\text{data}}(x) + 0.99p_{\text{noise}}(x)$ 
  - 99% of the samples are just noise
  - Taking logs, we get a lower bound

$$\log p_{\theta}(x) = \log \left[ 0.01p_{\text{data}}(x) + 0.99p_{\text{noise}}(x) \right] \geq \log 0.01p_{\text{data}}(x) = \log p_{\text{data}}(x) - \log 100$$

- For expected likelihoods, we know that
  - Lower bound
  - Upper bound (via non-negativity of KL)

$$\mathbb{E}_{p_{\text{data}}} [\log p_{\theta}(x)] \geq \mathbb{E}_{p_{\text{data}}} [\log p_{\text{data}}(x)] - \log 100$$

$$\mathbb{E}_{p_{\text{data}}} [\log p_{\text{data}}(x)] \geq \mathbb{E}_{p_{\text{data}}} [\log p_{\theta}(x)]$$

- As we increase the dimension of  $x$ , absolute value of  $\log p_{\text{data}}(x)$  increases proportionally but  $\log 100$  remains constant.

Hence,  $\mathbb{E}_{p_{\text{data}}} [\log p_{\theta}(x)] \approx \mathbb{E}_{p_{\text{data}}} [\log p_{\text{data}}(x)]$  in very high dimensions!

- **Case 3:** Great samples, poor test log-likelihoods. E.g., Memorizing training set
  - Samples look exactly like the training set (cannot do better!)
  - Test set will have zero probability assigned (cannot do worse!)
- The above cases suggest that it might be useful to disentangle likelihoods and samples
- **Likelihood-free learning** consider objectives that do not depend directly on a likelihood function

# Comparing Distributions via Samples



$$S_1 = \{x \sim P\}$$

vs.



$$S_2 = \{x \sim Q\}$$

- Given a finite set of samples from two distributions  $S_1 = \{x \sim P\}$  and  $S_2 = \{x \sim Q\}$ , how can we tell if these samples are from the same distribution? (i.e.,  $P = Q$ ?)

# Two-Sample Tests

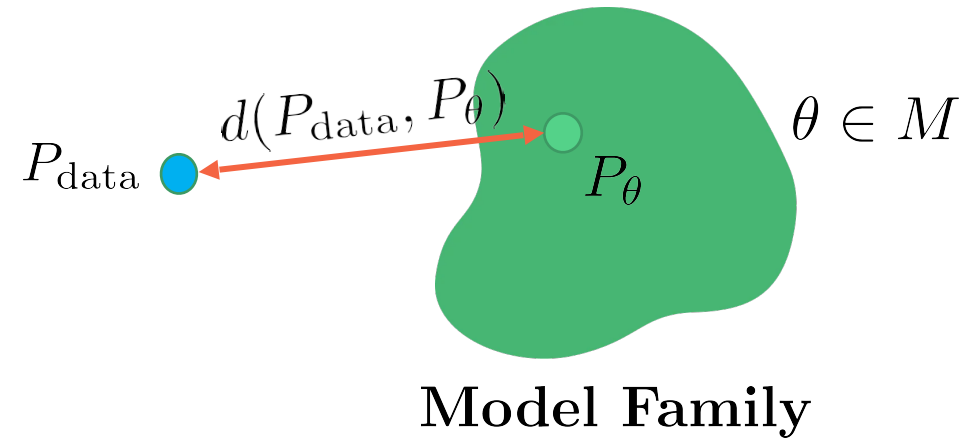
- Given  $S_1 = \{x \sim P\}$  and  $S_2 = \{x \sim Q\}$ , a **two-sample test** considers the following hypotheses
  - Null hypothesis  $H_0 : P = Q$
  - Alternative hypothesis  $H_1 : P \neq Q$
- Test statistic  $T$  compares  $S_1$  and  $S_2$  e.g., difference in means, variances of the two sets of samples.
- If  $T$  is larger than a threshold  $\alpha$ , then reject  $H_0$  otherwise we say  $H_0$  is consistent with observation.
- **Key observation:** Test statistic is **likelihood-free** since it does not involve the densities  $P$  or  $Q$  (only samples)

# Generative Modeling and Two-Sample Tests

Lecture #15



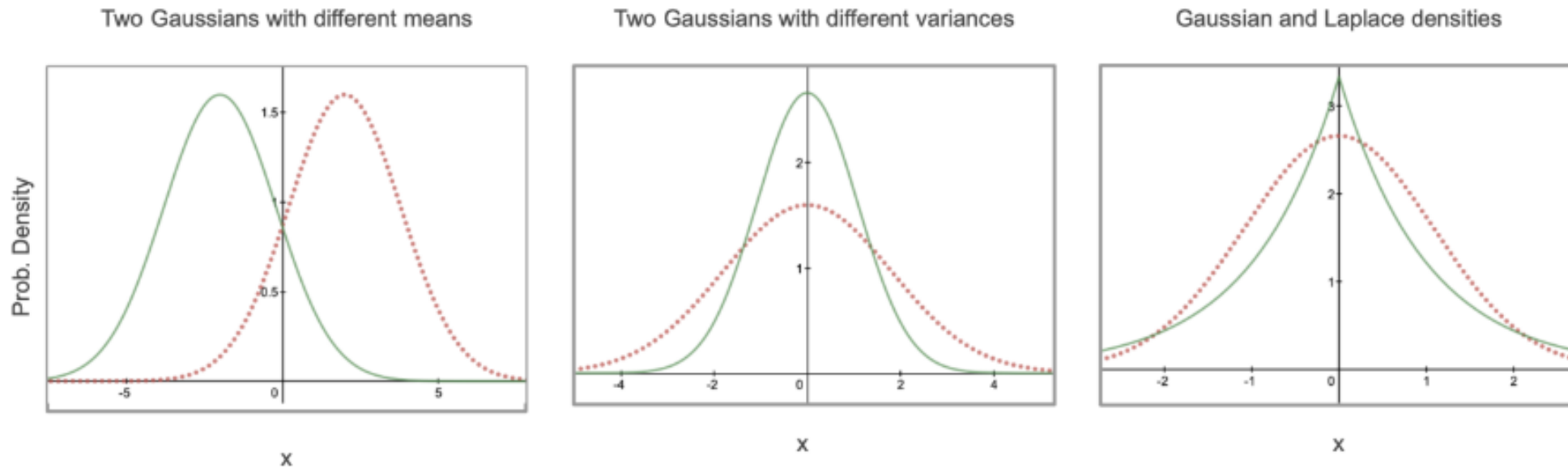
$$x_i \sim P_{\text{data}} \\ i = 1, 2, \dots, n$$



- A priori we assume direct access to  $S_1 = \mathcal{D} = \{x \sim p_{\text{data}}\}$
- In addition, we have a model distribution  $p_{\theta}$
- Assume that the model distribution permits efficient sampling (e.g., directed models). Let  $S_2 = \{x \sim p_{\theta}\}$
- **Alternative notion of distance between distributions:** Train the generative model to minimize a two-sample test objective between  $S_1$  and  $S_2$

# Two-Sample Test via a Discriminator

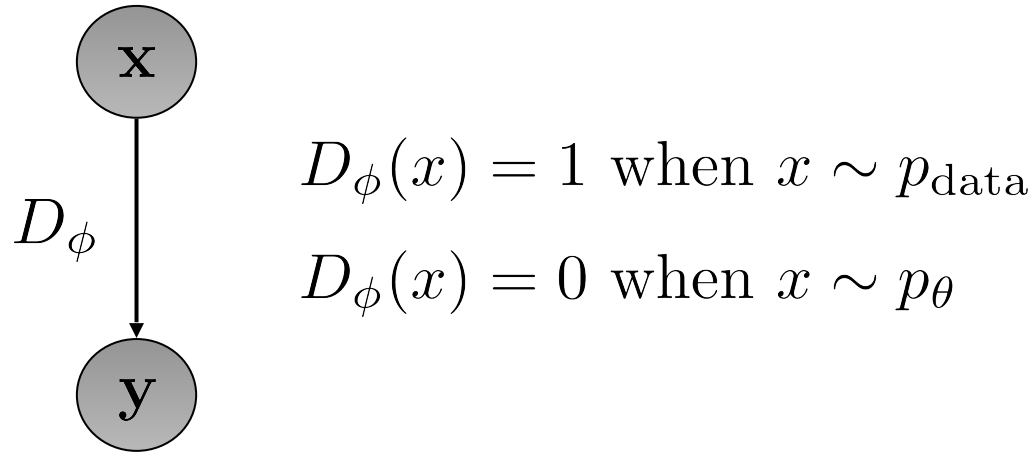
- Finding a two-sample test objective in high dimensions is hard



- In the generative model setup, we know that  $S_1$  and  $S_2$  come from different distributions  $p_{\text{data}}$  and  $p_{\theta}$  respectively
- **Key idea:** Learn a statistic that **maximizes** a suitable notion of distance between the two sets of samples  $S_1$  and  $S_2$

# Two-Sample Test via a Discriminator

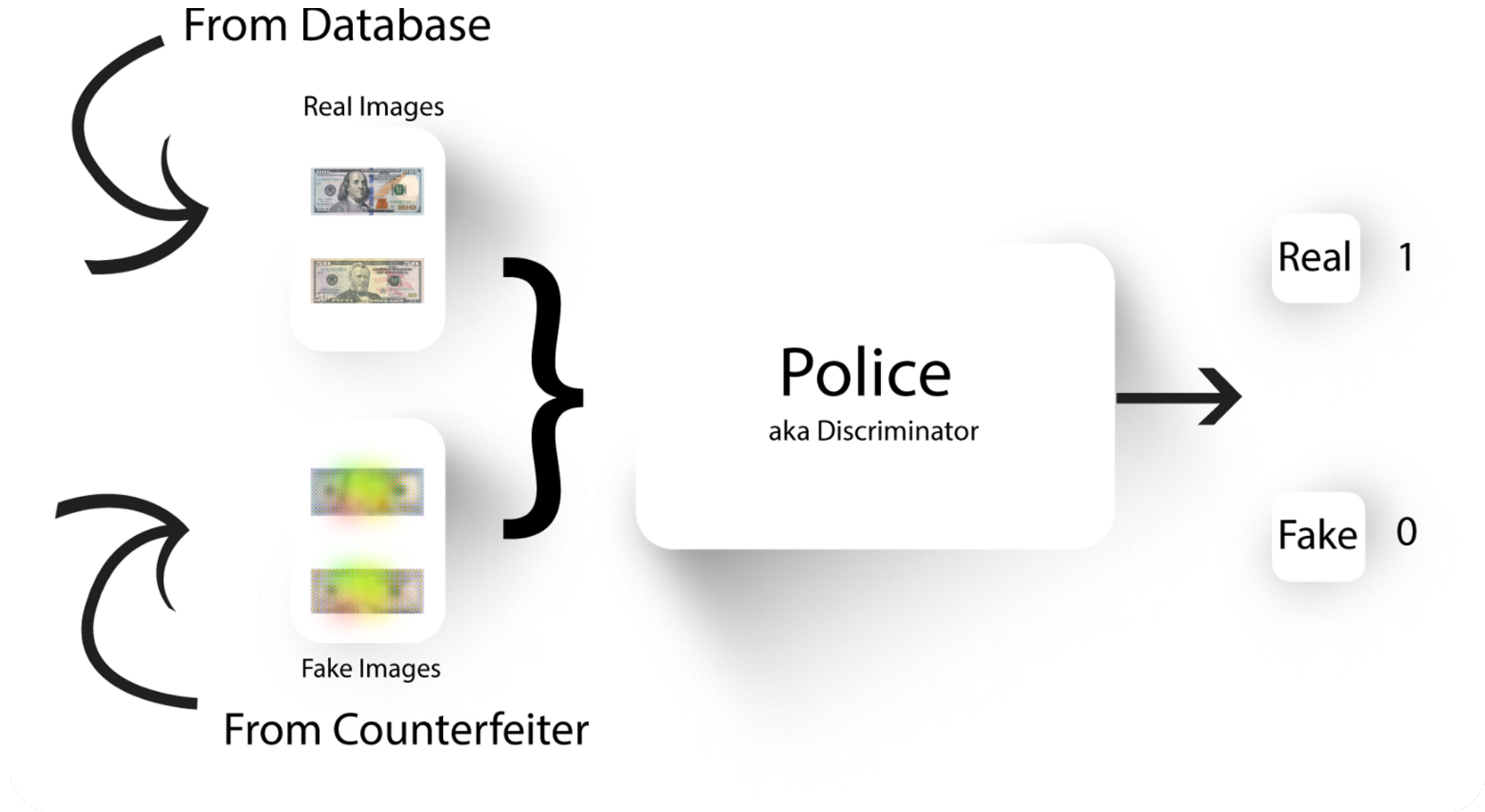
Lecture #15



- **Two-Sample Test via a Discriminator**

- Any function (e.g., neural network) which tries to distinguish “real” samples from the dataset and “fake” samples generated from the model
- Maximizes the two-sample test objective (in support of the alternative hypothesis  $p_{\text{data}} \neq p_\theta$ )

# Two-Sample Test via a Discriminator



# Two-Sample Test via a Discriminator

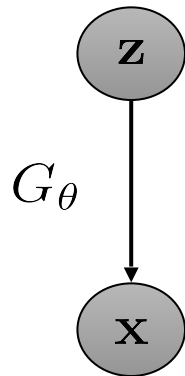
- Training objective for discriminator:

$$\max_D V(D) = \mathbb{E}_{x \sim p_{\text{data}}} [\log D(x)] + \mathbb{E}_{x \sim p_G} [\log(1 - D(x))].$$

- For a fixed generator  $G$ , the discriminator is performing binary classification with the cross entropy objective:
  - Assign probability 1 to true data points  $x \sim p_{\text{data}}$
  - Assign probability 0 to fake samples  $x \sim p_G$

- Optimal Discriminator:  $D_G^*(x) = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_G(x)}.$

- A two player minimax game between a **generator** and a **discriminator**

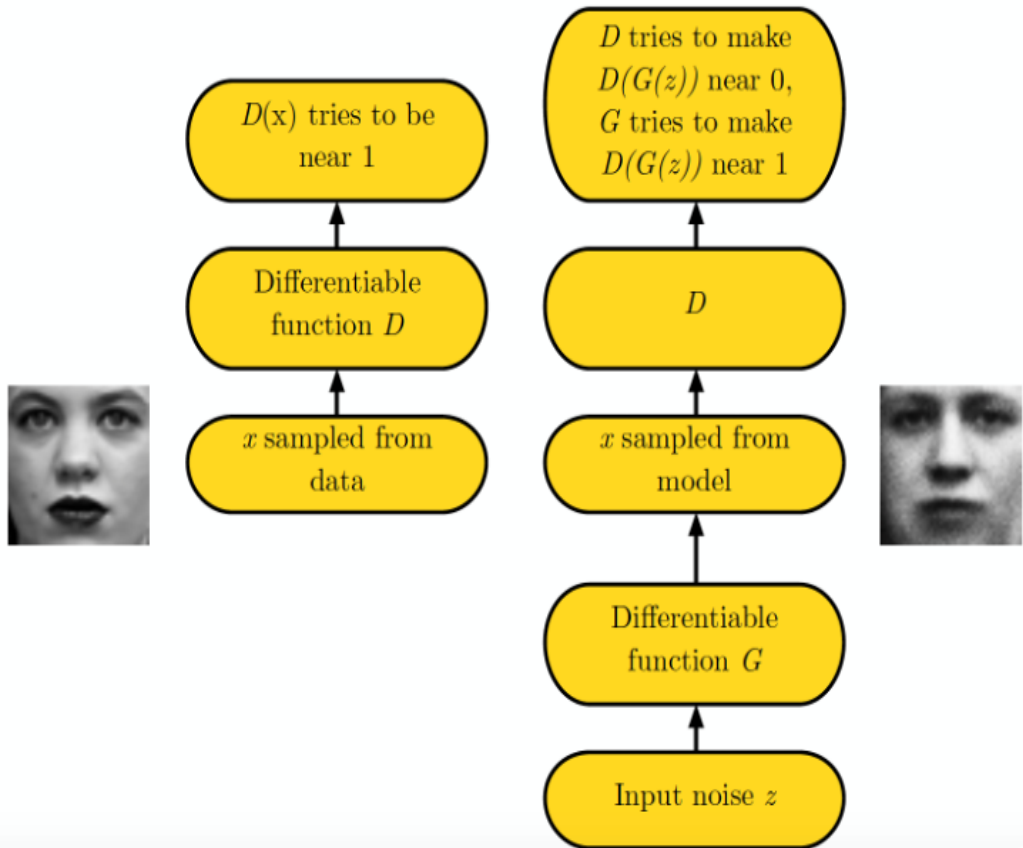


- **Generator**
  - Directed, latent variable model with a deterministic mapping between  $z$  and  $x$  given by  $G_\theta$
  - Minimizes a two-sample test objective (in support of the null hypothesis  $p_{\text{data}} = p_\theta$ )

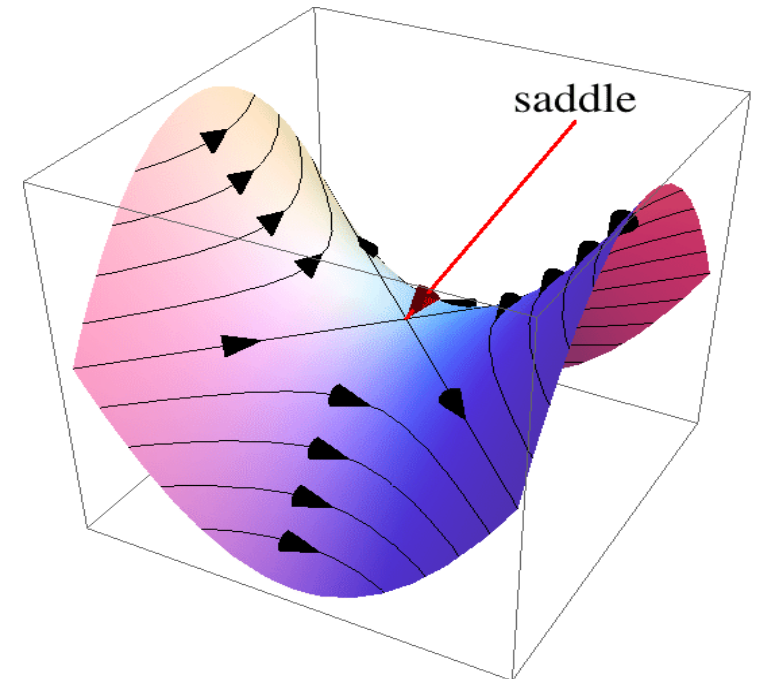
# Example of GAN Objective

- Training objective for both generator and discriminator:

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



The joint optimum  $(G^*, D^*)$  is a saddle point.



# The GAN Training Algorithm

- Sample minibatch of  $m$  training points  $x^{(1)}, x^{(2)}, \dots, x^{(m)}$  from  $\mathcal{D}$
- Sample minibatch of  $m$  noise vectors  $z^{(1)}, z^{(2)}, \dots, z^{(m)}$  from  $p_Z$
- Update the discriminator parameters  $\phi$  by stochastic gradient **ascent**

$$\nabla_{\phi} V(G_{\theta}, D_{\phi}) = \frac{1}{m} \nabla_{\phi} \sum_{i=1}^m \left[ \log D_{\phi}(x^{(i)}) + \log(1 - D_{\phi}(G_{\theta}(z^{(i)}))) \right].$$

- Update the generator parameters  $\theta$  by stochastic gradient **descent**

$$\nabla_{\theta} V(G_{\theta}, D_{\phi}) = \frac{1}{m} \nabla_{\theta} \sum_{i=1}^m \log(1 - D_{\phi}(G_{\theta}(z^{(i)}))).$$

- Repeat for fixed number of iterations

# Example of GAN Objective

- For the optimal discriminator  $D_G^*(\cdot)$  and fixed generator  $G(\cdot)$ , we have

$$\begin{aligned} V(G, D_G^*(x)) &= \mathbb{E}_{x \sim p_{\text{data}}} \left[ \log \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_G(x)} \right] + \mathbb{E}_{x \sim p_G} \left[ \log \frac{p_G(x)}{p_{\text{data}}(x) + p_G(x)} \right] \\ &= \mathbb{E}_{x \sim p_{\text{data}}} \left[ \log \frac{p_{\text{data}}(x)}{\frac{1}{2}(p_{\text{data}}(x) + p_G(x))} \right] + \mathbb{E}_{x \sim p_G} \left[ \log \frac{p_G(x)}{\frac{1}{2}(p_{\text{data}}(x) + p_G(x))} \right] - \log 4 \\ &= \underbrace{D_{\text{KL}} \left( p_{\text{data}} \parallel \frac{p_{\text{data}} + p_G}{2} \right) + D_{\text{KL}} \left( p_G \parallel \frac{p_{\text{data}} + p_G}{2} \right)}_{2 \times \text{Jensen-Shannon Divergence (JSD)}} - \log 4 \\ &= 2D_{\text{JS}}(p_{\text{data}}, p_G) - \log 4. \end{aligned}$$

# Jenson-Shannon Divergence

- Also called as the symmetric KL divergence

$$D_{\text{JS}}(p, q) = \frac{1}{2} \left( D_{\text{KL}} \left( p \parallel \frac{1}{2}(p + q) \right) + D_{\text{KL}} \left( q \parallel \frac{1}{2}(p + q) \right) \right).$$

- Properties

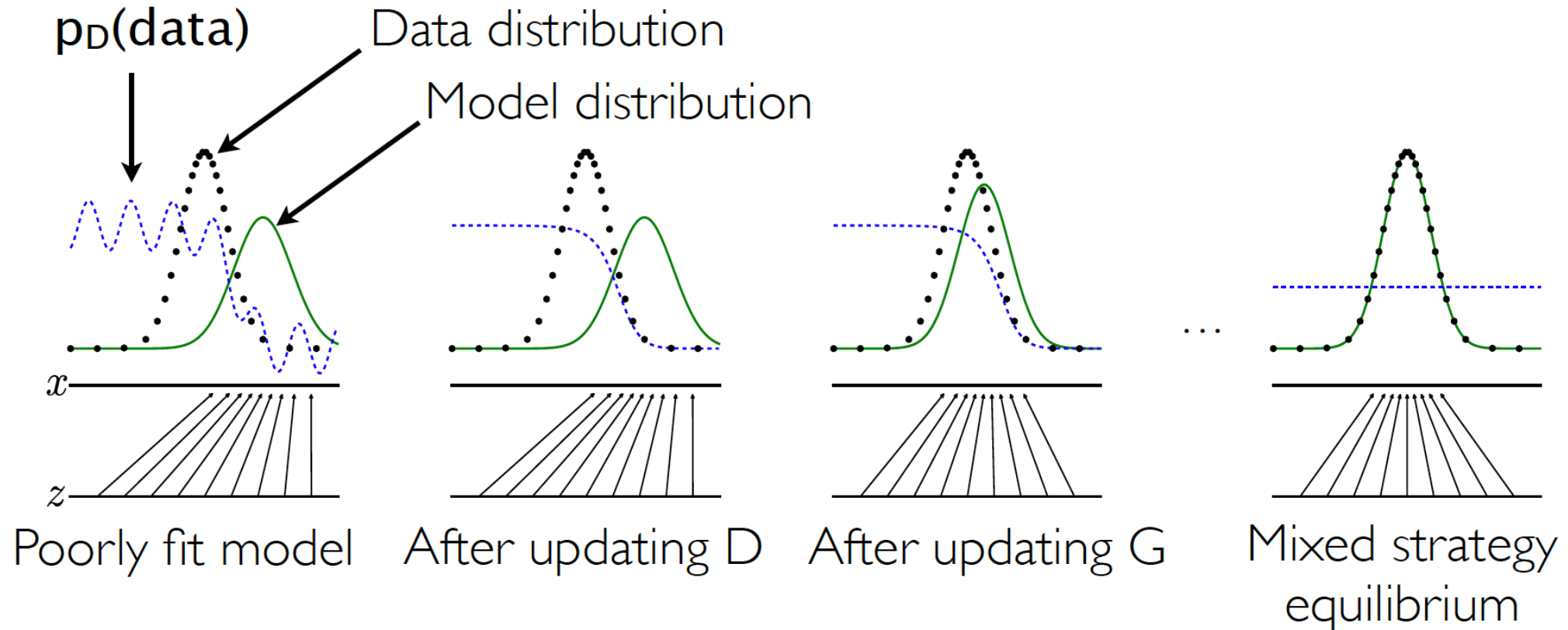
- $D_{\text{JS}}(p, q) \geq 0$
- $D_{\text{JS}}(p, q) = 0$  iff  $p = q$
- $D_{\text{JS}}(p, q) = D_{\text{JS}}(q, p)$
- $\sqrt{D_{\text{JS}}(p, q)}$  satisfies triangle inequality  $\Rightarrow$  Jenson-Shannon Distance

- Optimal generator for the JSD/Negative Cross Entropy GAN  $p_{G^*} = p_{\text{data}}$

- For the optimal discriminator  $D_{G^*}^*$  and generator  $G^*$ , we have  $V(G^*, D_{G^*}^*) = -\log 4$ .

# Alternating Optimization in GANs

$$\min_{\theta} \max_{\phi} V(G_{\theta}, D_{\phi}) = \mathbb{E}_{x \sim p_{\text{data}}} [\log D_{\phi}(x)] + \mathbb{E}_{z \sim p_Z} [\log(1 - D_{\phi}(G_{\theta}(z)))].$$



# Frontiers in GAN Research

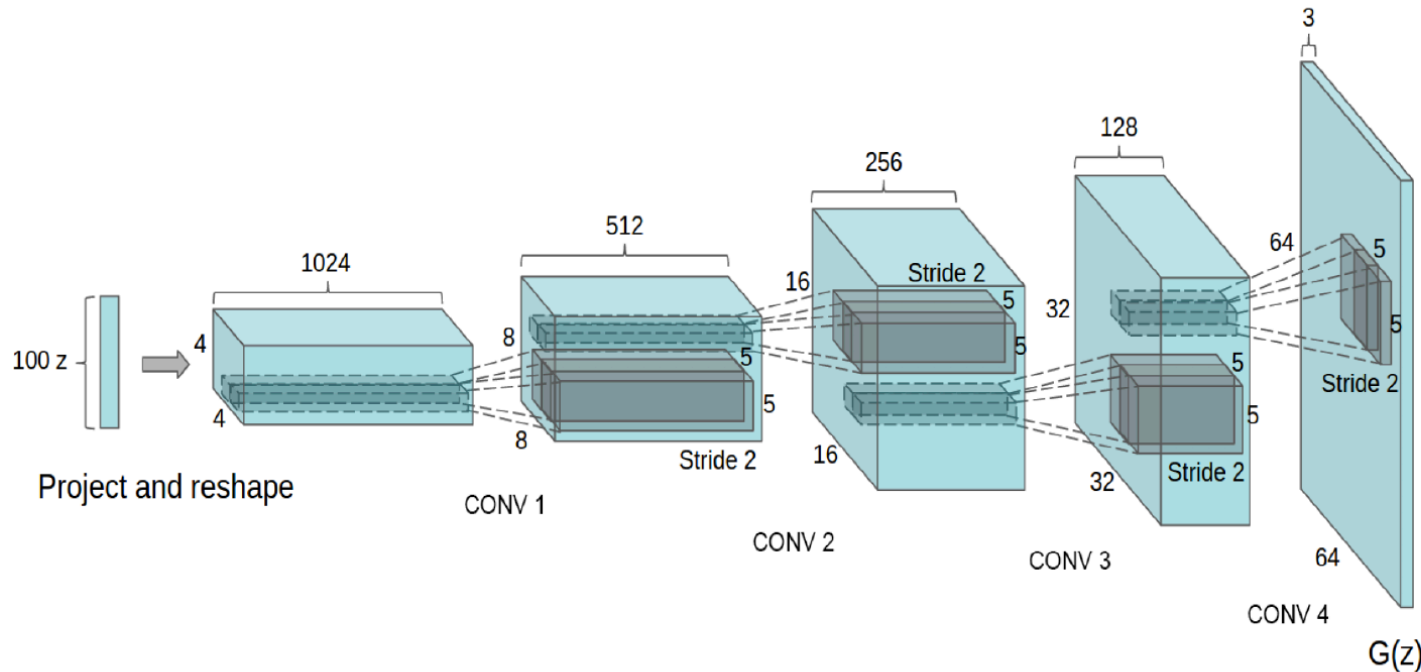
Lecture #15



- GANs have been successfully applied to several domains and tasks
- However, working with GANs can be very challenging in practice
  - Unstable optimization
  - Mode collapse
  - Performance evaluation
- Many tricks have been proposed to successfully train GANs

Image Source: Ian Goodfellow. Samples from Goodfellow et al., 2014, Radford et al., 2015, Liu et al., 2016, Karras et al., 2017, Karras et al., 2018

## Generator Architecture



### Key ideas:

- Replace FC hidden layers with Convolutions
  - **Generator:** Fractional-Strided convolutions
- Use Batch Normalization after each layer
- **Inside Generator**
  - Use ReLU for hidden layers
  - Use Tanh for the output layer

# DCGAN Example – LSUN bedrooms

Lecture #15



Radford, Alec, Luke Metz, and Soumith Chintala. "Unsupervised representation learning with deep convolutional generative adversarial networks." arXiv:1511.06434 (2015).

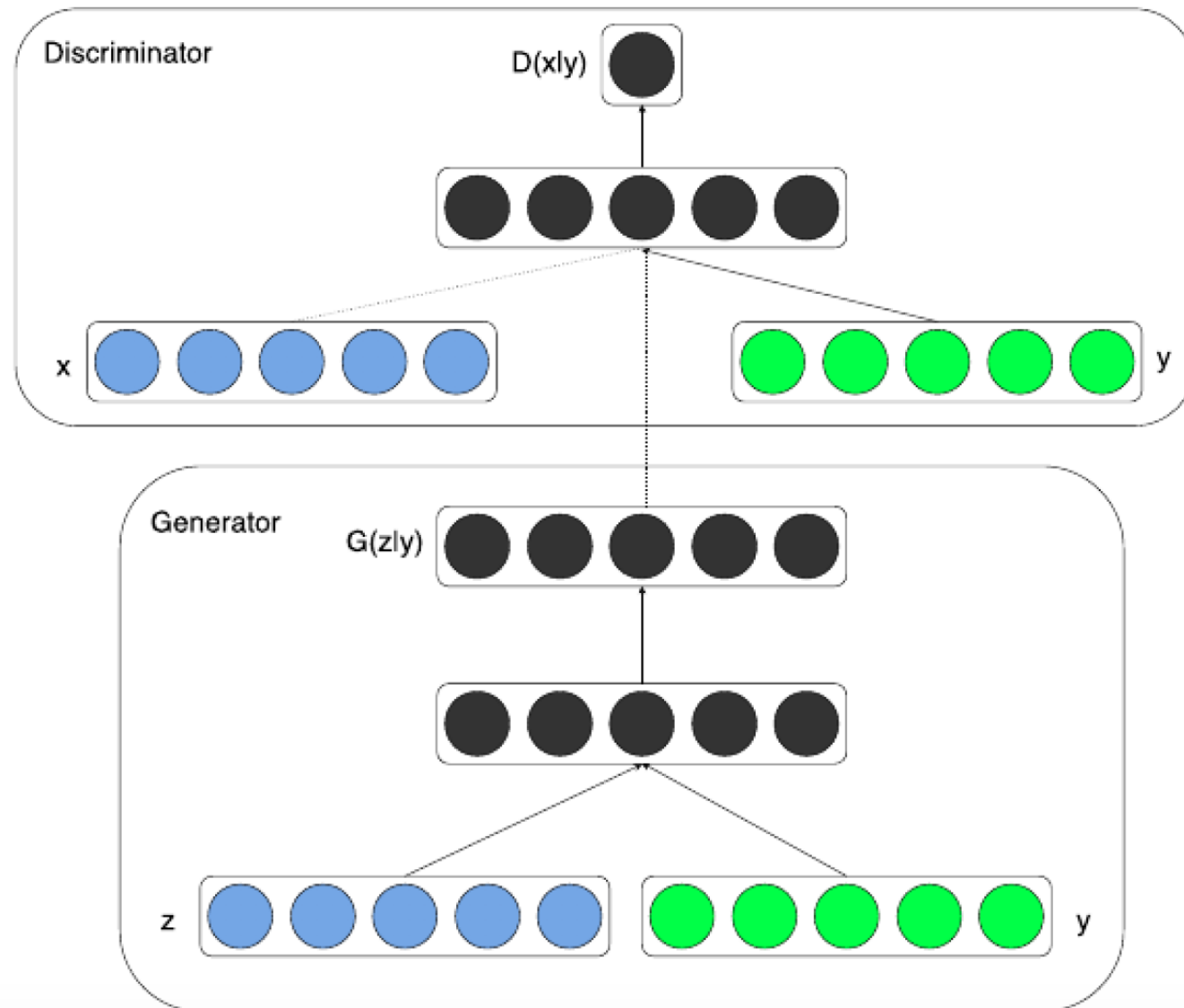
- GAN is too free. How to add some constraints?
- Add conditional variables  $\mathbf{y}$  into  $G$  and  $D$

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

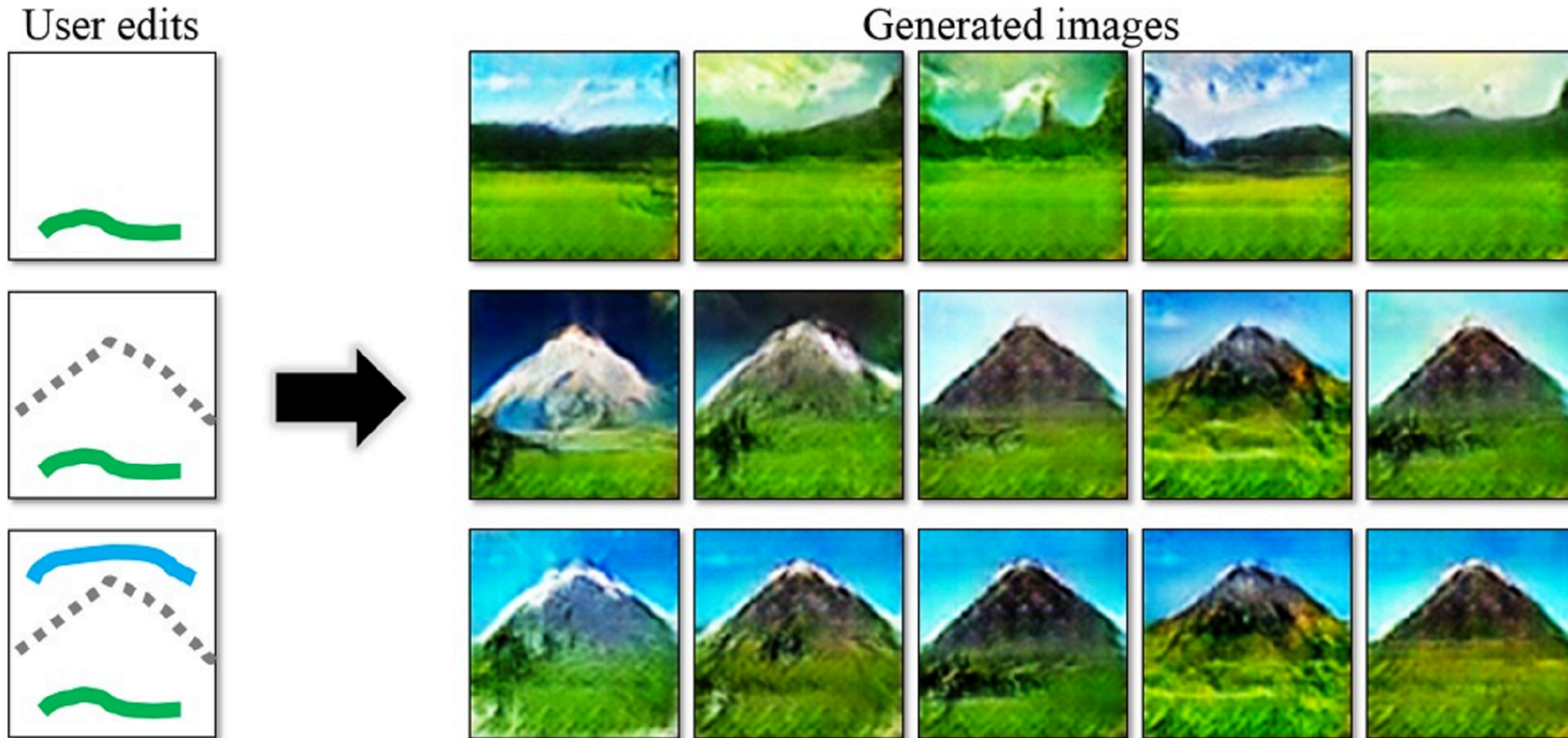


$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x}|\mathbf{y})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z}|\mathbf{y})))].$$

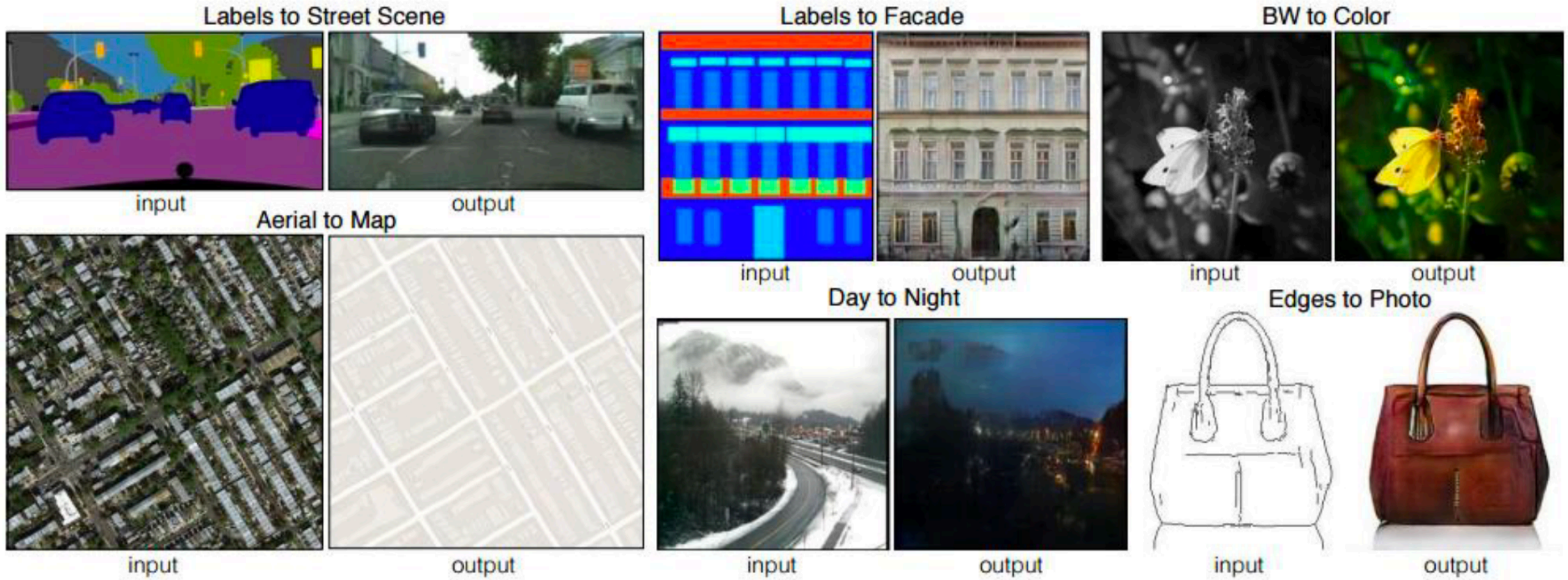
# Conditional GAN



# Conditional GAN - Examples



# Conditional GAN - Examples

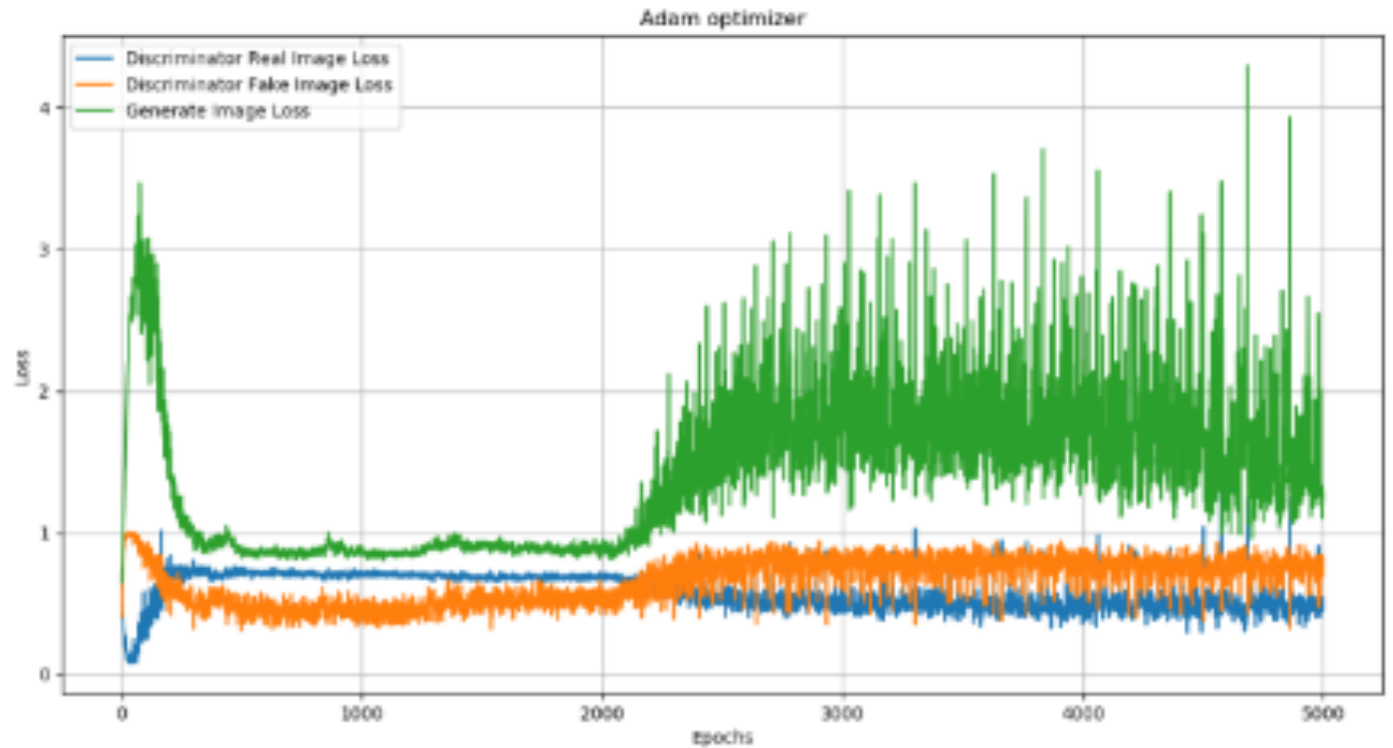
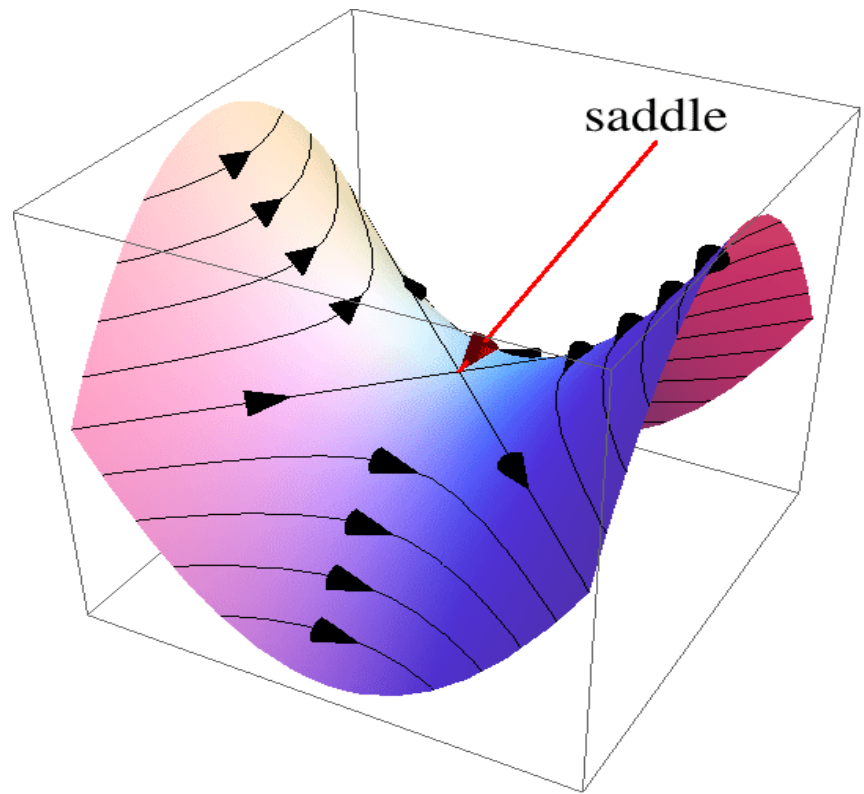


# Optimization Challenges

- **Theorem (informal):** If the generator updates are made in function space and discriminator is optimal at every step, then the generator is guaranteed to converge to the data distribution
- **Unrealistic assumptions!**
- In practice, the generator and discriminator loss keeps oscillating during GAN training
- No robust stopping criteria in practice (unlike MLE)

# Optimization Challenges

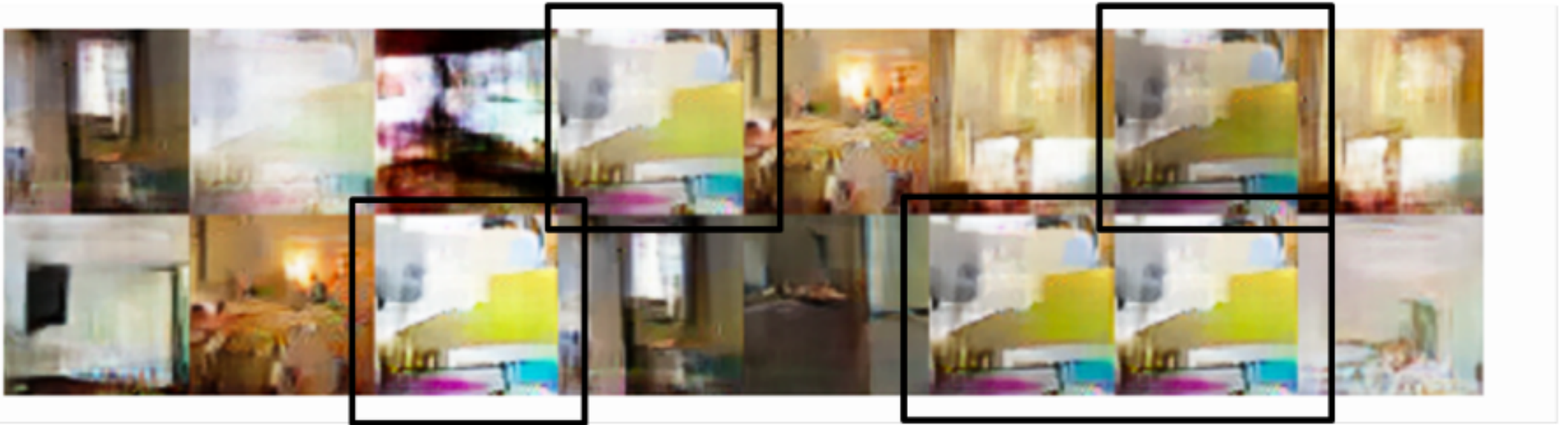
Lecture #15



Source: Mirantha Jayathilaka

# Mode Collapse

- GANs are notorious for suffering from **mode collapse**
- Intuitively, this refers to the phenomena where the generator of a GAN collapses to one of few samples (dubbed as “modes”)

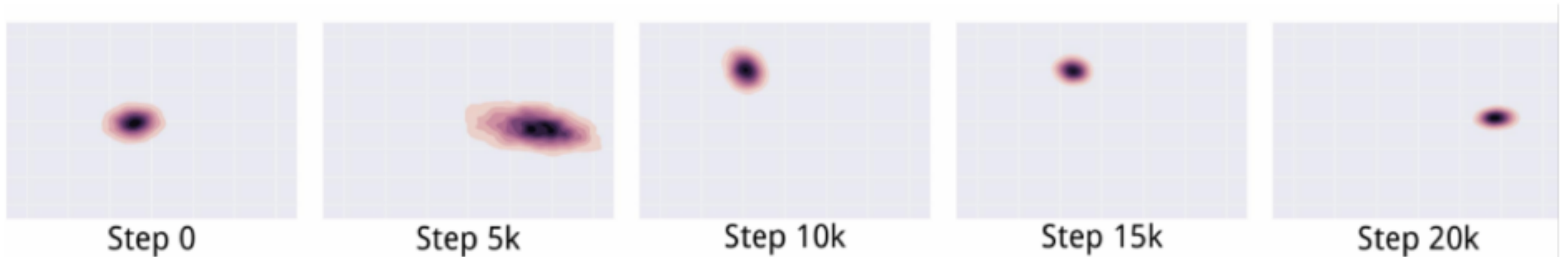


# Mode Collapse

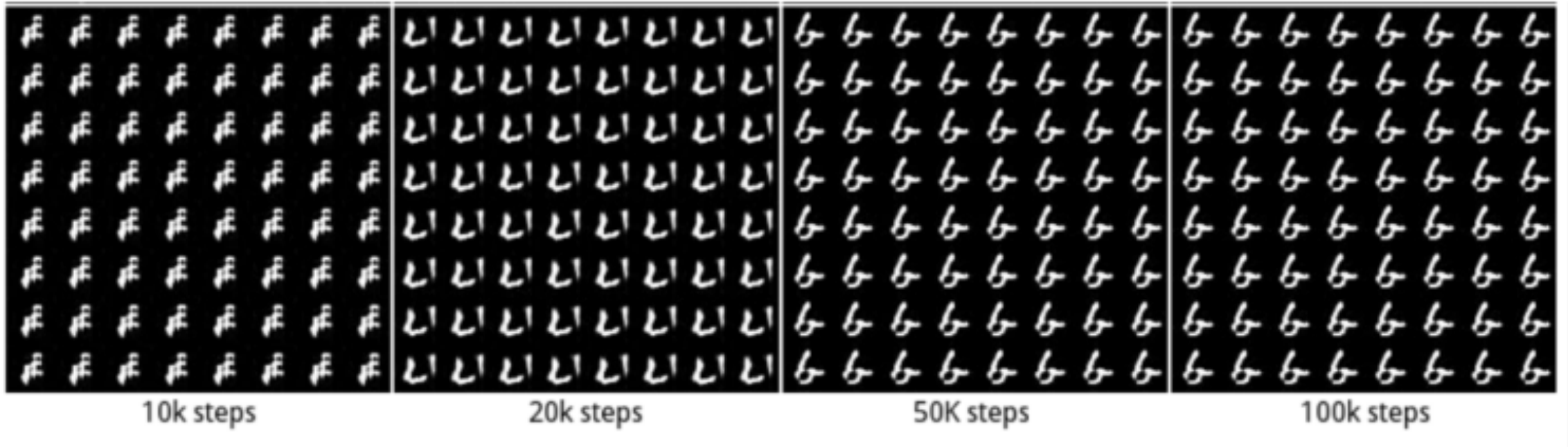
- True distribution is a mixture of Gaussians



- The generator distribution keeps oscillating between different modes



# Mode Collapse



Source: Metz et al., 2017

- Fixes to mode collapse are mostly empirically driven: alternative architectures, alternative GAN loss, adding regularization terms, etc.
- How to Train a GAN? Tips and tricks to make GANs work by Soumith Chintala et al. <https://github.com/soumith/ganhacks>

# Beauty Lies in the Eyes of the Discriminator

Lecture #15



GAN generated art auctioned at Christie's.  
**Expected Price:** \$7,000 – \$10,000  
**True Price:** \$432,500

Source: Robbie Barrat, Obvious

# Introduction to Deep Generative Modeling

Lecture #15

**HY-673** – Computer Science Dep., University of Crete

Professors: Yannis Pantazis, Yannis Stylianou

Teaching Assistant: Michail Raptakis