

# Deep Learning for Computer Vision

113-1/Fall 2024

<https://cool.ntu.edu.tw/courses/41702> (NTU COOL)

<http://vllab.ee.ntu.edu.tw/dlcv.html> (Public website)

Yu-Chiang Frank Wang 王鈺強, Professor

Dept. Electrical Engineering, National Taiwan University

## About Final Projects

- Updates
  - **1:30pm-5pm, Dec 26<sup>th</sup>, Thursday**  
(sorry, no space available on 25<sup>th</sup>)
  - 3~4 people per group  
(team up in mid Nov.)
  - Adapt from latest CVPR/ICCV/ECCV challenges or competitions
  - Poster presentation;  
**code required** for reproduction
  - Intra/inter-group evaluation
  - Snack/drink provided

# FINAL PROJECT

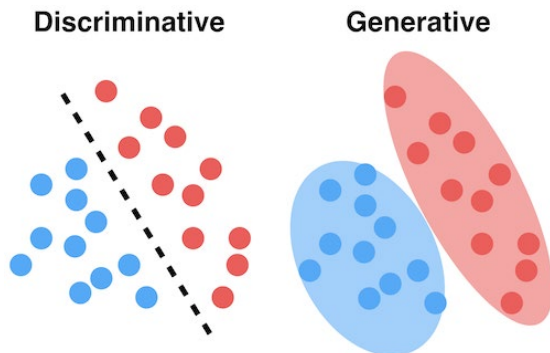
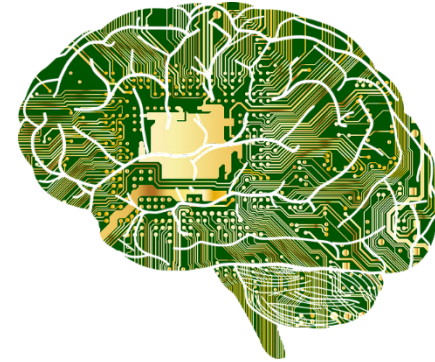


Sp Adobe Spark

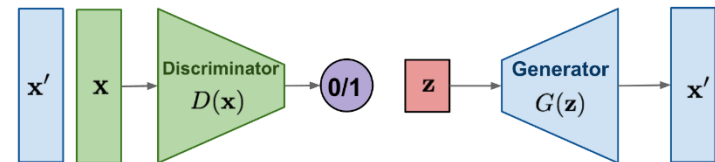


# What's to Be Covered Today...

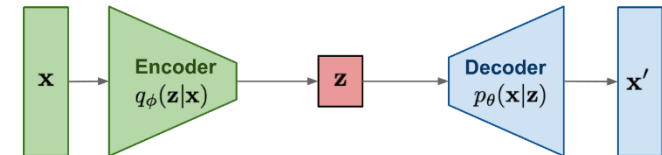
- Generative Models
  - Diffusion Model
    - Conditional Diffusion Model
      - Classifier Guidance
      - Classifier-Free Guidance
      - Text/Image Guidance
    - Personalization via Diffusion Model
  - Generative Adversarial Network
  - HW #2 is out! (due 10/29)



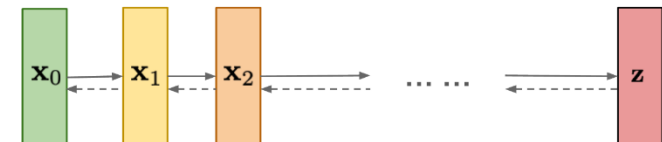
**GAN: Adversarial training**



**VAE: maximize variational lower bound**



**Diffusion models:**  
Gradually add Gaussian noise and then reverse



# A Quick Recap of Generative Models

## Discriminative Model:

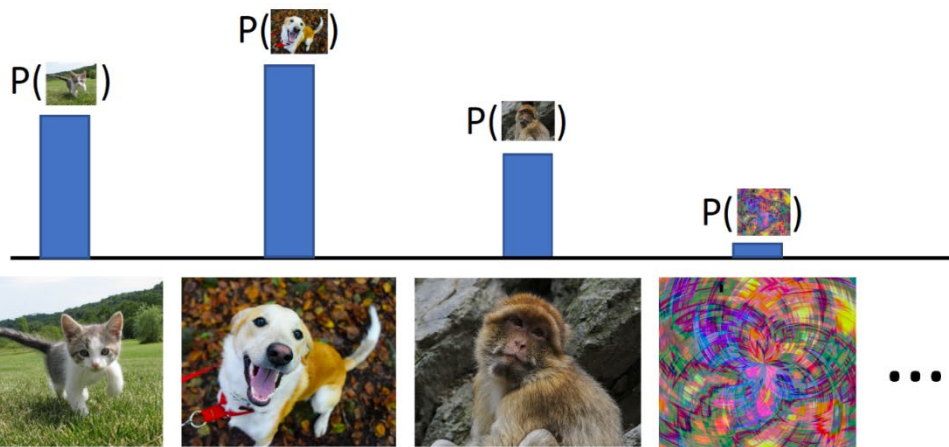
Learn a probability distribution  $p(y|x)$

## Generative Model:

Learn a probability distribution  $p(x)$

## Conditional Generative Model:

Learn  $p(x|y)$

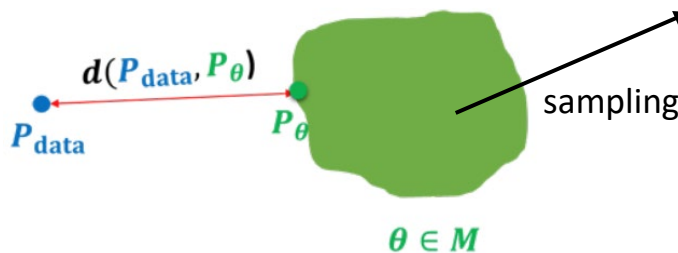


Generative model: All possible images compete with each other for probability mass

Model can “reject” unreasonable inputs by assigning them small values



$$x_i \sim p_{data}$$

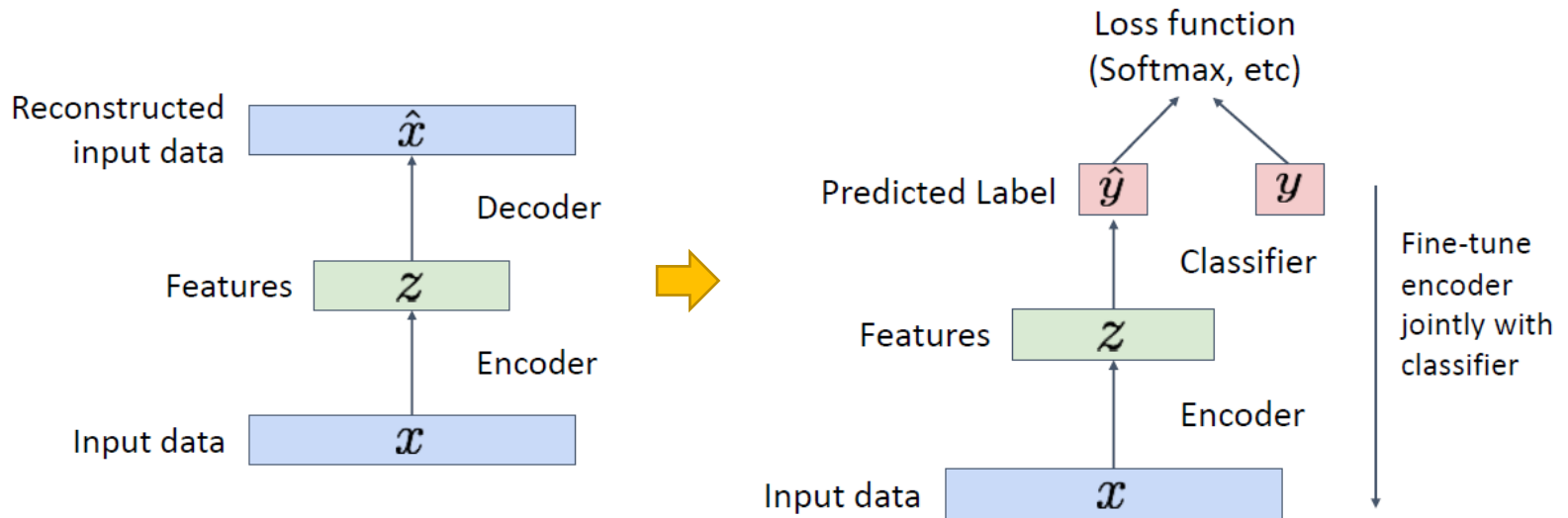


$$\hat{x}_i \sim p_\theta$$



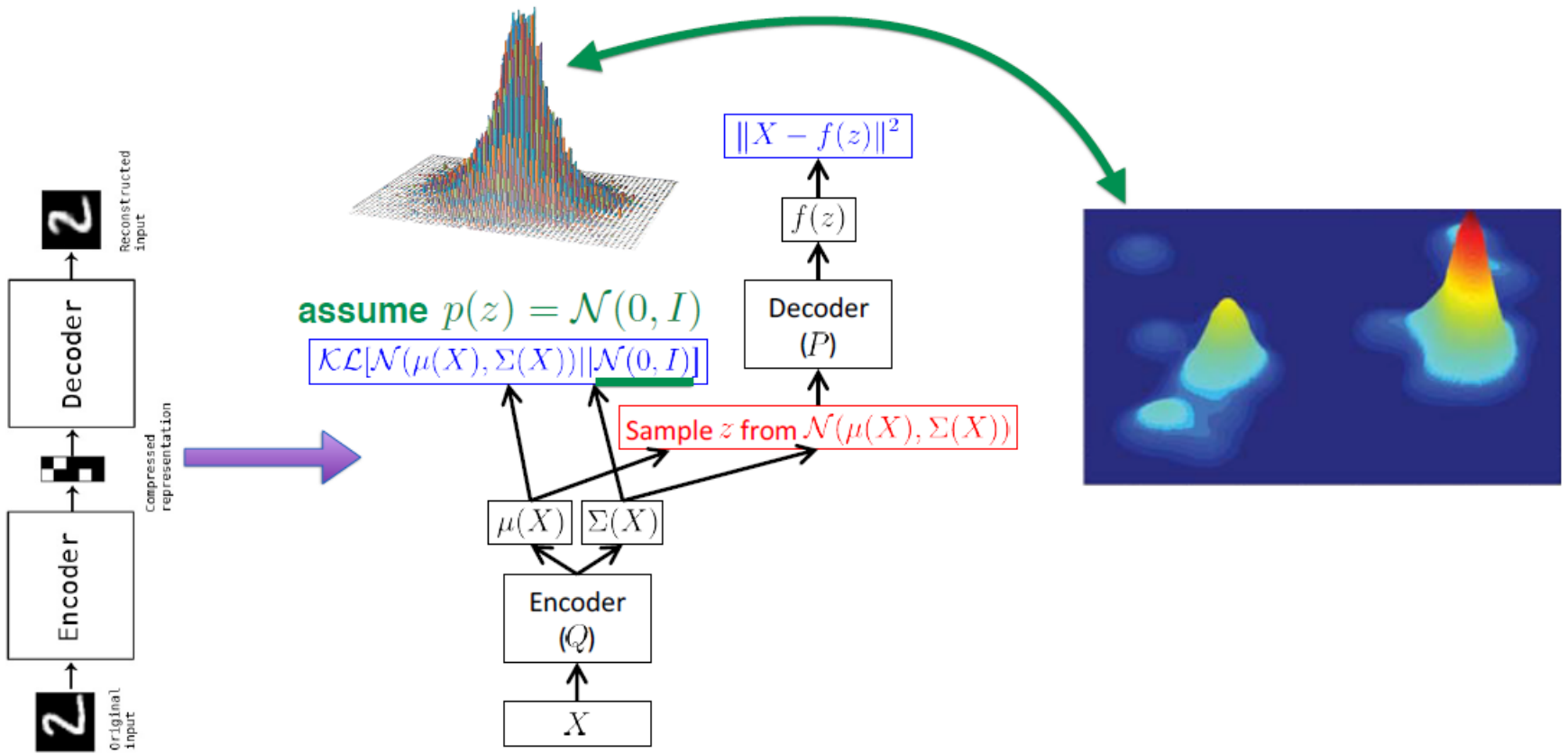
# Autoencoder (AE)

- Unsupervised learning for deriving latent representation
  - Train AE with reconstruction objectives
- Train autoencoder (AE) for downstream tasks
  - After AE training is complete, freeze/finetune the encoder and learn additional modules (e.g., MLP) for downstream tasks
  - E.g., to train a DNN for classification, one can freeze the encoder and learn an additional MLP as the classifier.



# From AE to Variational Autoencoder (VAE)

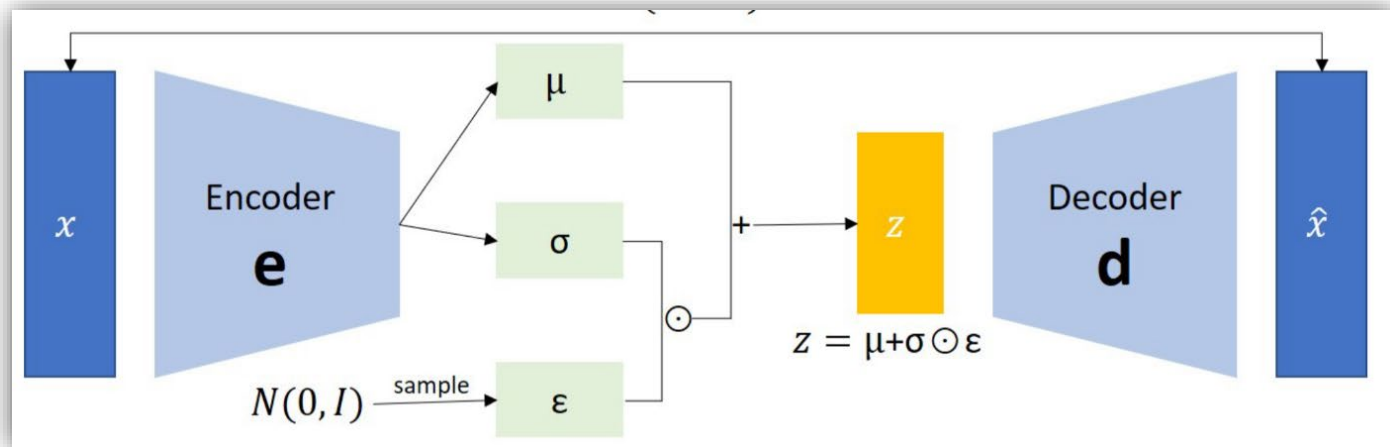
Now is a “*distribution*”, we can assume it to be a distribution easy to sample from, e.g. Gaussian



# Reparameterization Trick in VAE

- Remarks

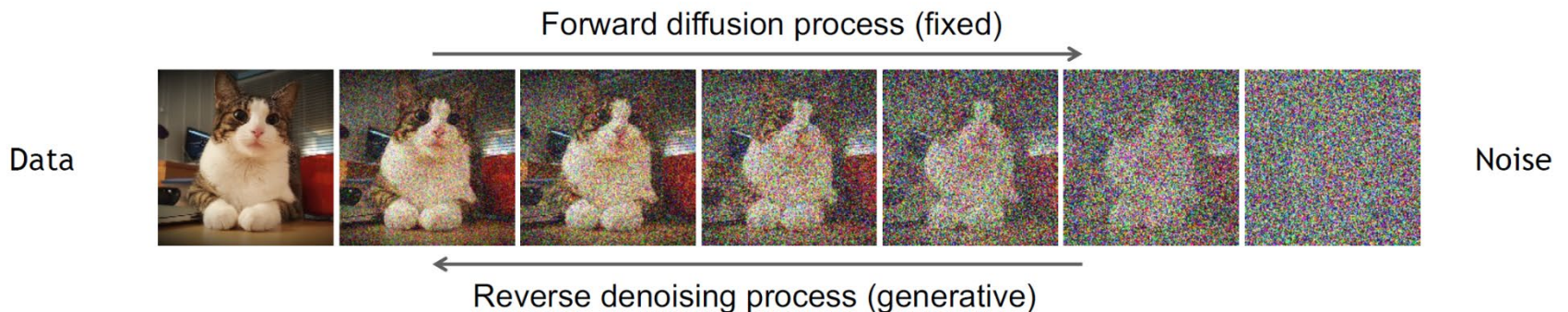
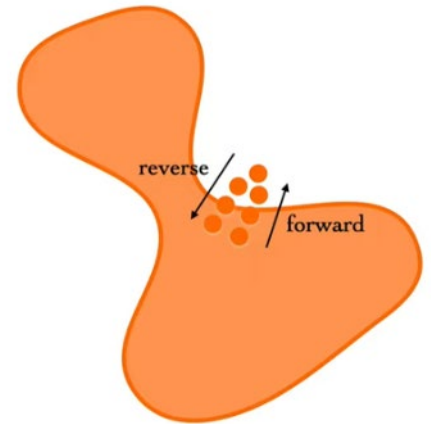
- Given  $x$ , **sample**  $z$  from latent distribution (described by output parameters of encoder), we apply  $z = \mu + \sigma \odot \varepsilon$  ( $\varepsilon$  simply generated by Normal distribution).
- For training, this enables BP gradients in encoder through  $\mu$  and  $\sigma$ ; for inference, this introduces generation stochasticity.



# Denoising Diffusion Probabilistic Models (DDPM)

Learning to generate by denoising

- 2 processes required for training:
  - **Forward diffusion process**
    - gradually add noise to input
  - **Reverse diffusion process**
    - learns to generate/restore data by denoising
    - typically implemented via a **conditional U-net**

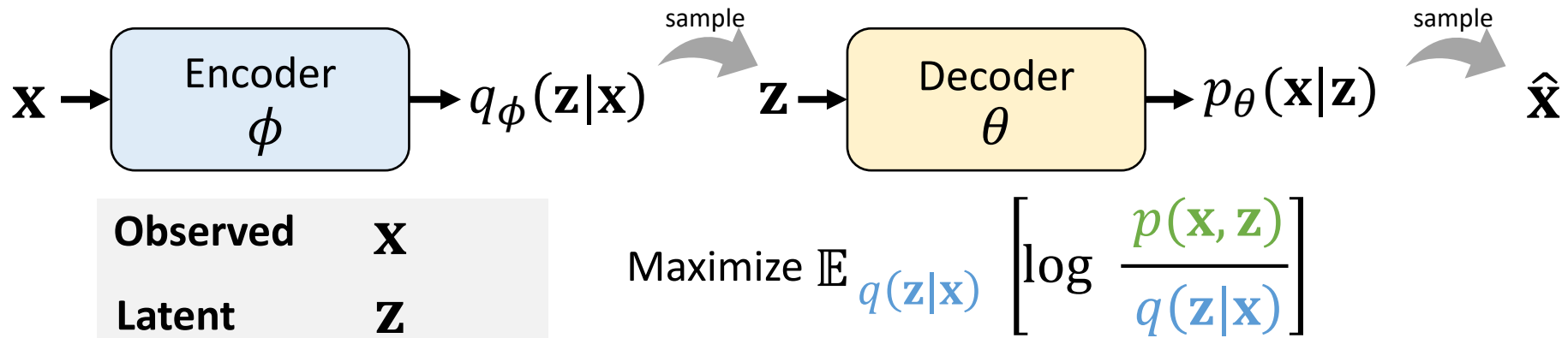


[Ho et al., Denoising Diffusion Probabilistic Models, NeurIPS 2020](#)

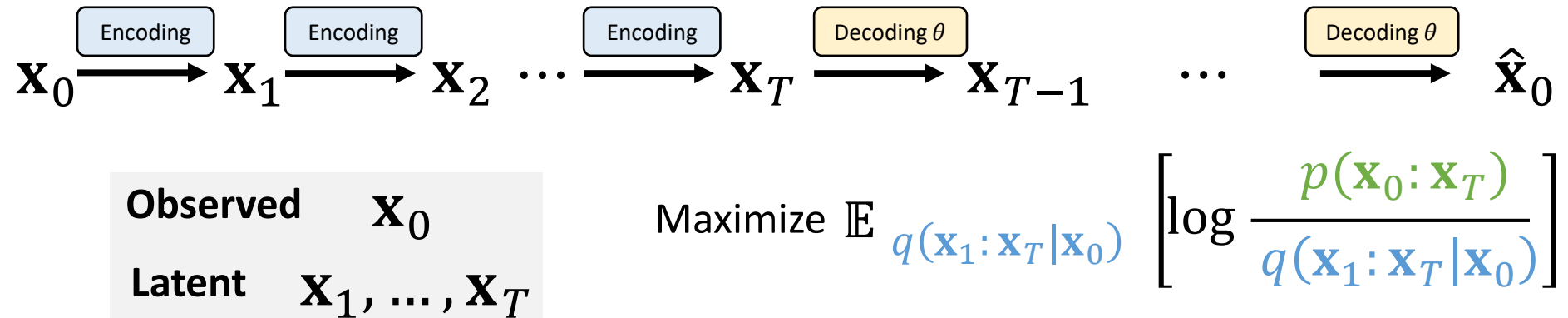
[Song et al., Score-Based Generative Modeling through Stochastic Differential Equations, ICLR 2021](#)

# VAE vs. DDPM

## Variational Autoencoder

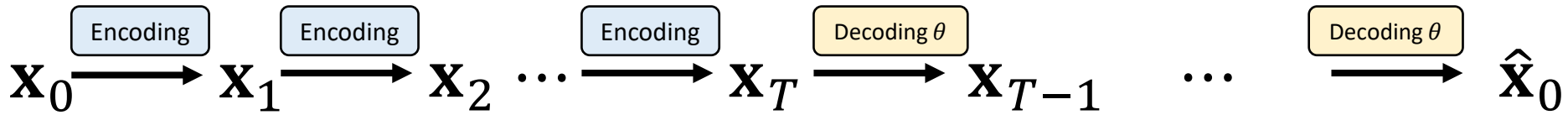


## Diffusion model





# Training DDPM



$$\sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) || p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]$$

$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$

$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$



# From Unconditional to Conditional Generative Models

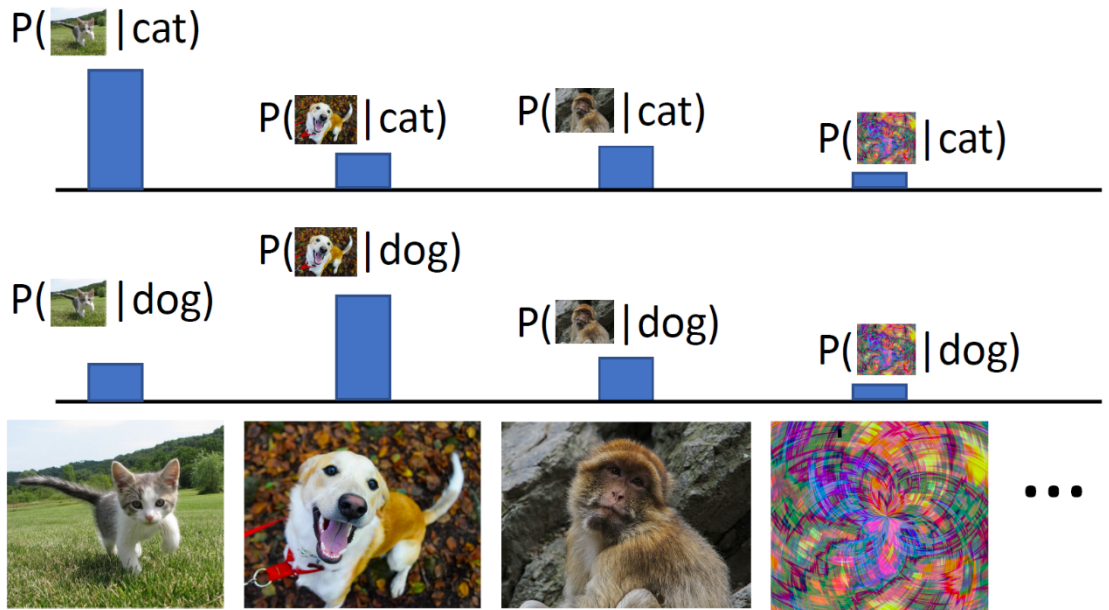
## Discriminative Model:

Learn a probability distribution  $p(y|x)$

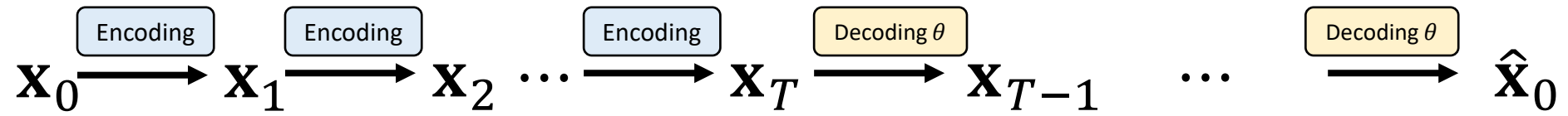
## Generative Model:

Learn a probability distribution  $p(x)$

**Conditional Generative Model:** Learn  $p(x|y)$



Conditional Generative Model: Each possible label induces a competition among all images

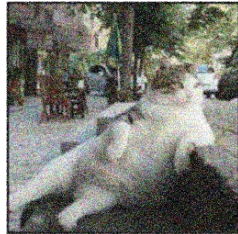


$$\sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) || p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]$$

$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$



$\mathbf{x}_0$

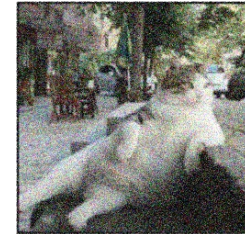


$\mathbf{x}_{t-1}$



$\mathbf{x}_t$

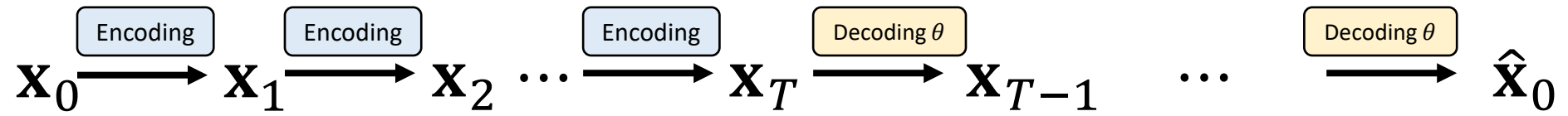
$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$



$\mathbf{x}_{t-1}$



$\mathbf{x}_t$



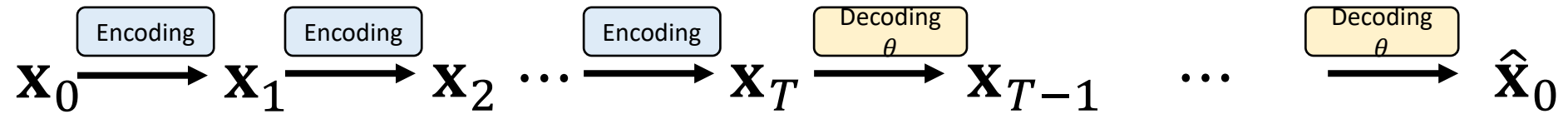
$$\sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) || p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]$$

$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$

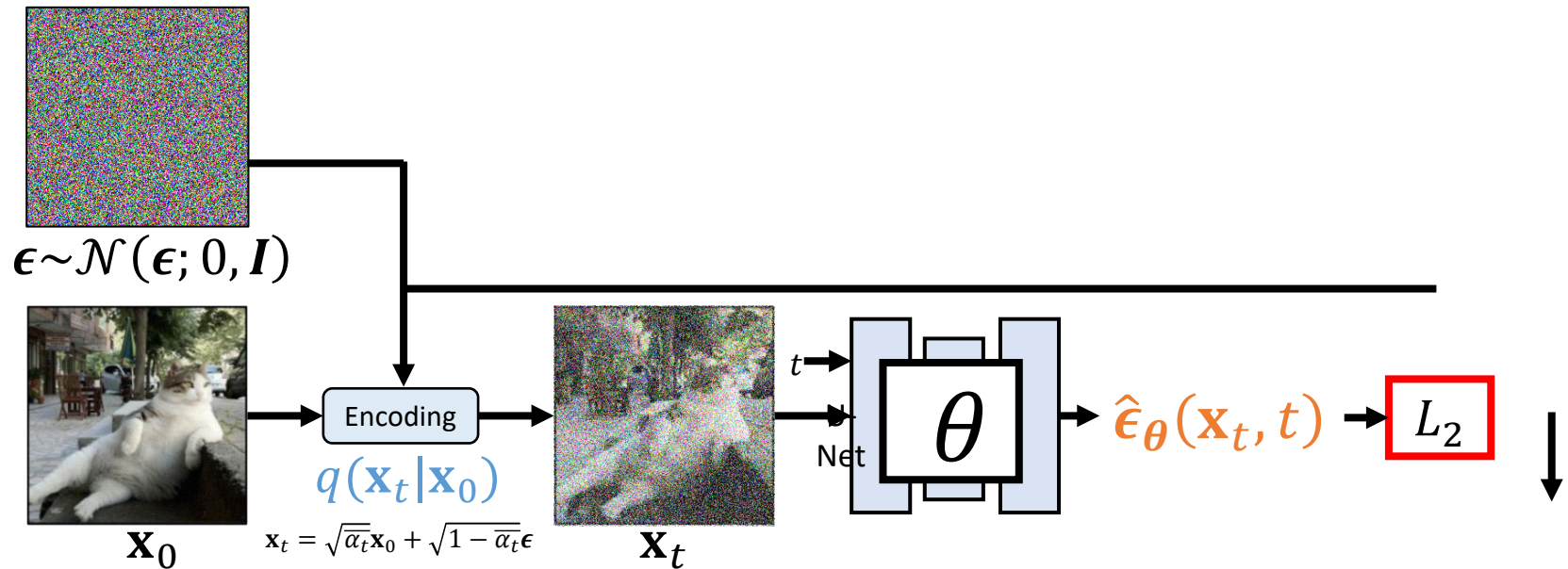
$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$



## Observation #2



$\rightarrow \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [ w(t) [ \| \hat{\epsilon}_{\theta}(\mathbf{x}_t, t) - \epsilon \|_2^2 ] ]$





# Training vs. Inference

- Summary

## Algorithm 1 Training

- 1: **repeat**
- 2:  $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
- 3:  $t \sim \text{Uniform}(\{1, \dots, T\})$
- 4:  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 5: Take gradient descent step on  $\nabla_{\theta} \|\epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|^2$
- 6: **until** converged

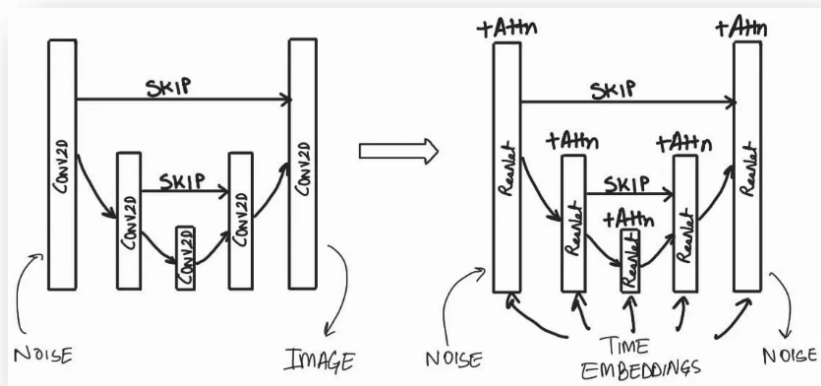
$\mathbf{x}_t$

## Algorithm 2 Sampling

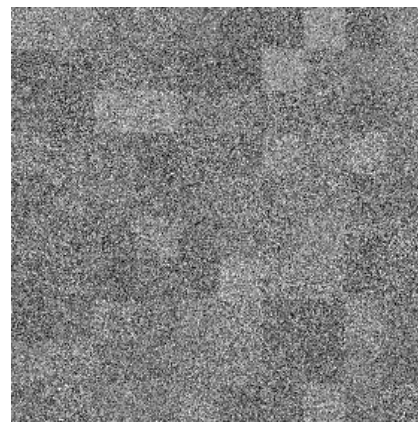
- 1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 2: **for**  $t = T, \dots, 1$  **do**
- 3:  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 4:  $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
- 5: **end for**
- 6: **return**  $\mathbf{x}_0$

allowing generation diversity

$$\mathbf{x}_t \xrightarrow{\text{model}} \epsilon_{\theta}(\mathbf{x}_t, t) \xrightarrow{P(x_t|x_0) \rightarrow P(x_0|x_t, \epsilon_{\theta})} \hat{\mathbf{x}}_0(\mathbf{x}_t, \epsilon_{\theta}) \longrightarrow \mu(\mathbf{x}_t, \hat{\mathbf{x}}_0), \beta_t \xrightarrow{P(x_{t-1}|x_t, x_0)} \hat{\mathbf{x}}_{t-1}$$



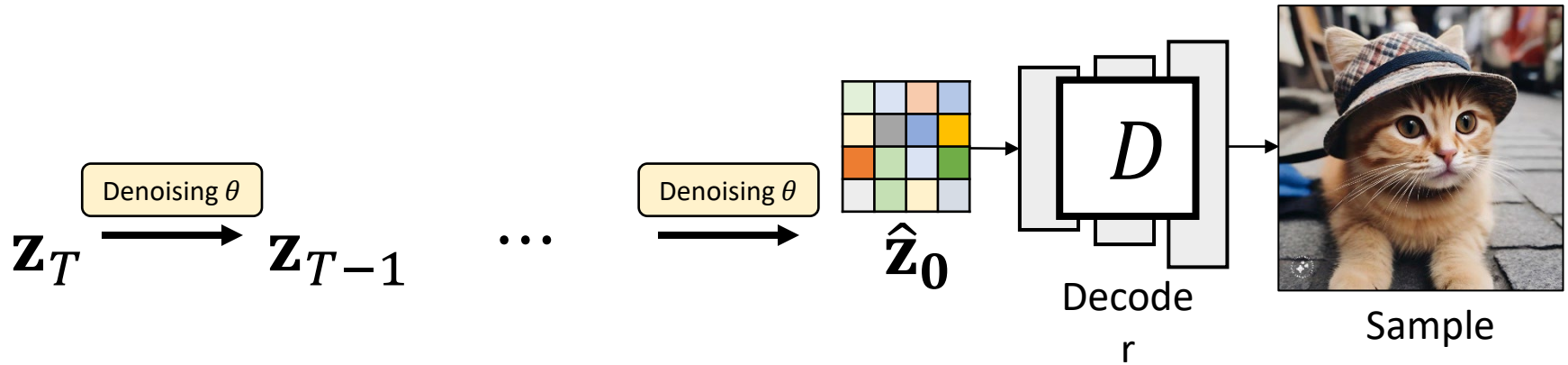
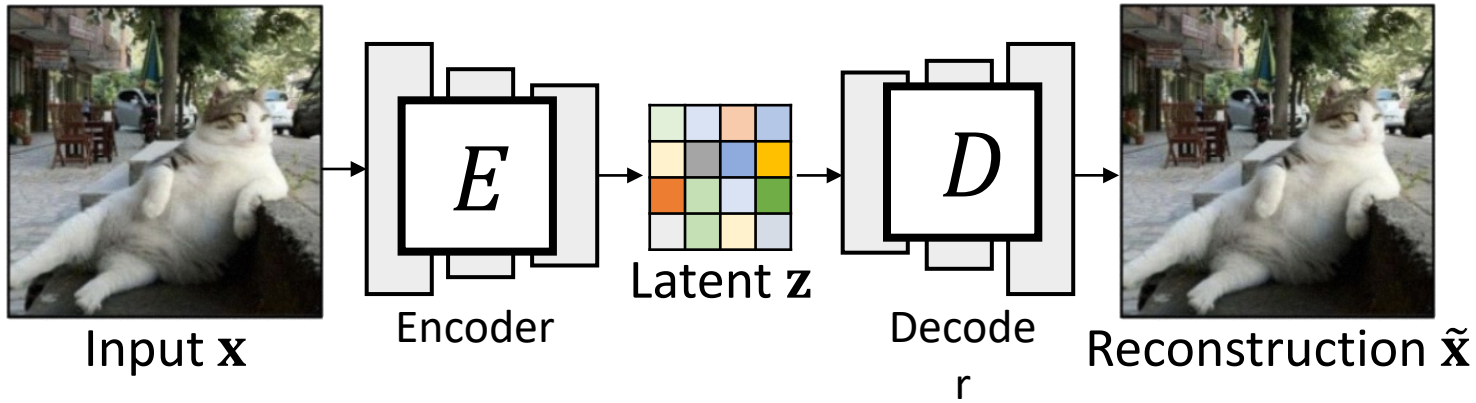
Training U-Net for DDPM noise prediction



MNIST handwritten image data



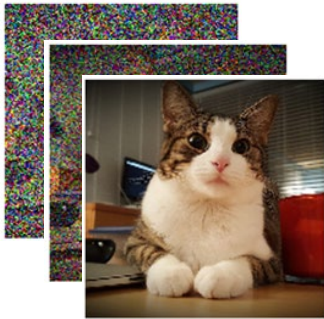
# Latent Diffusion Models



# From DDPM to DDIM: Denoising Diffusion Implicit Models

- DDIM
  - Sampling process for generation

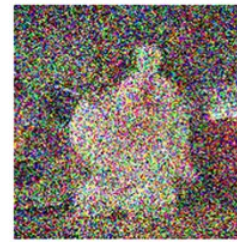
$$\mathbf{x}_{t-1} = \sqrt{\alpha_{t-1}} \left( \underbrace{\frac{\mathbf{x}_t - \sqrt{1 - \alpha_t} \epsilon_{\theta}^{(t)}(\mathbf{x}_t)}{\sqrt{\alpha_t}}}_{\text{“predicted } \mathbf{x}_0\text{”}} \right) + \underbrace{\sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \epsilon_{\theta}^{(t)}(\mathbf{x}_t)}_{\text{“direction pointing to } \mathbf{x}_t\text{”}} + \underbrace{\sigma_t \epsilon_t}_{\text{random noise}}$$



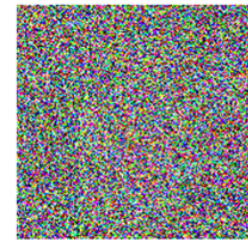
Step 0



Step 123



Step 456

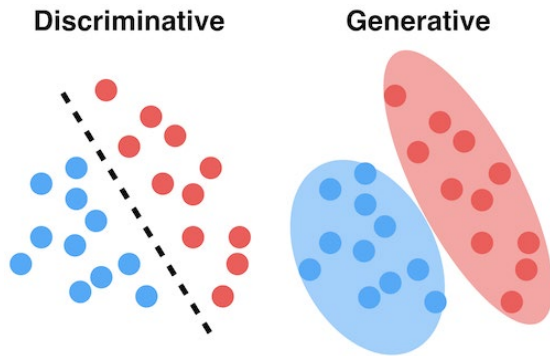
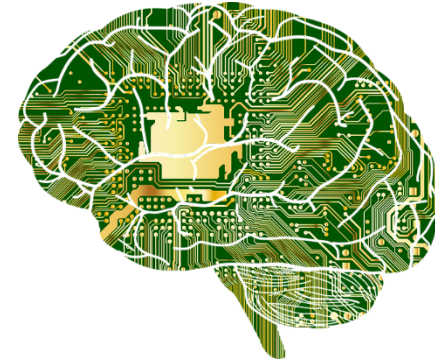


Step 999

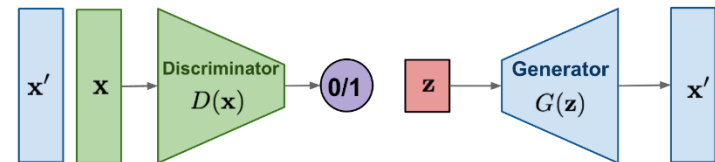
- Additional comment on  $\sigma_t^2$ : *stochastic vs. deterministic* generation process
- Since DDIM and DDPM share **the same objective function**, so one can use a **pretrained** DDPM for DDIM generation.

# What's to Be Covered Today...

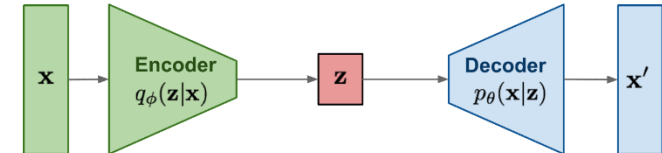
- Generative Models
  - Diffusion Model
    - Conditional Diffusion Model
      - Classifier Guidance
      - Personalization via Diffusion Model
  - Generative Adversarial Network



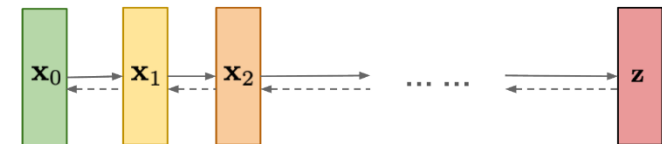
**GAN: Adversarial training**



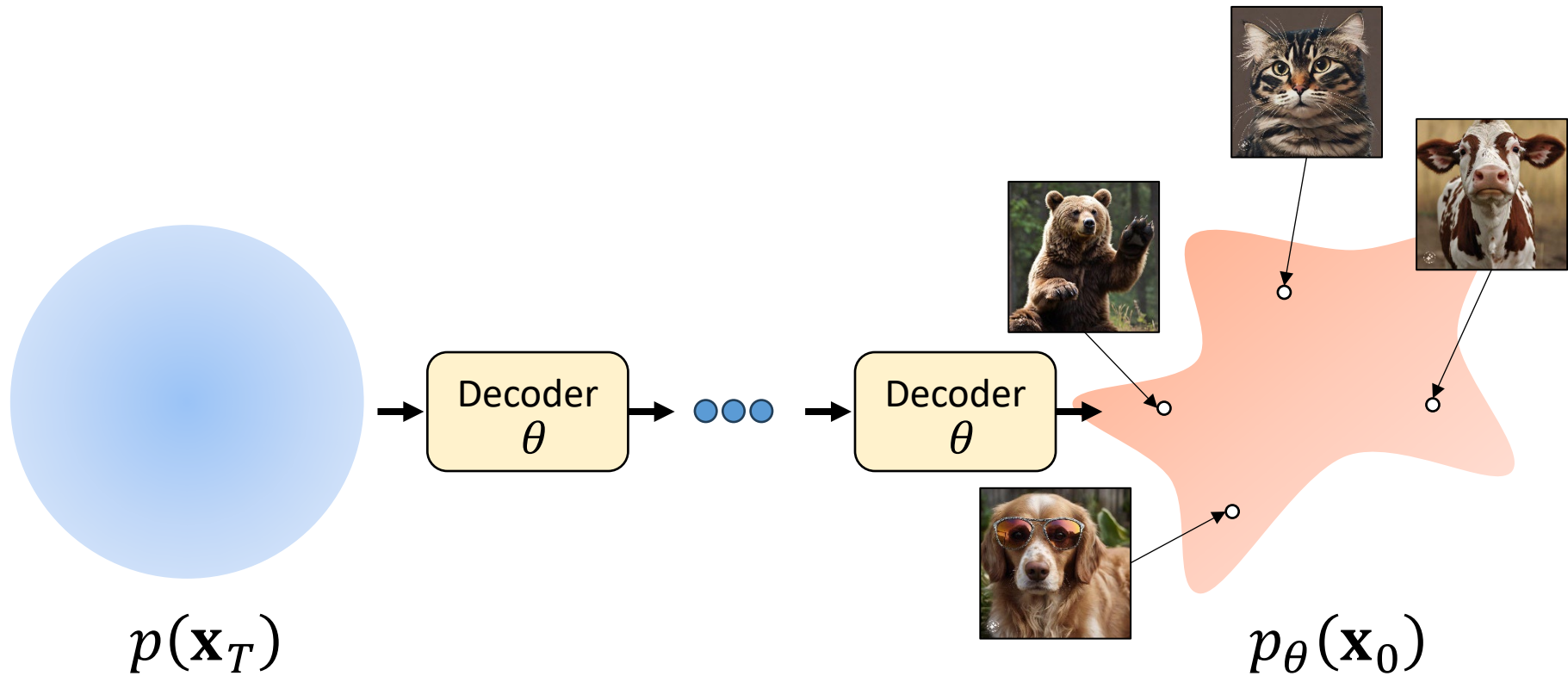
**VAE: maximize variational lower bound**



**Diffusion models:**  
Gradually add Gaussian noise and then reverse

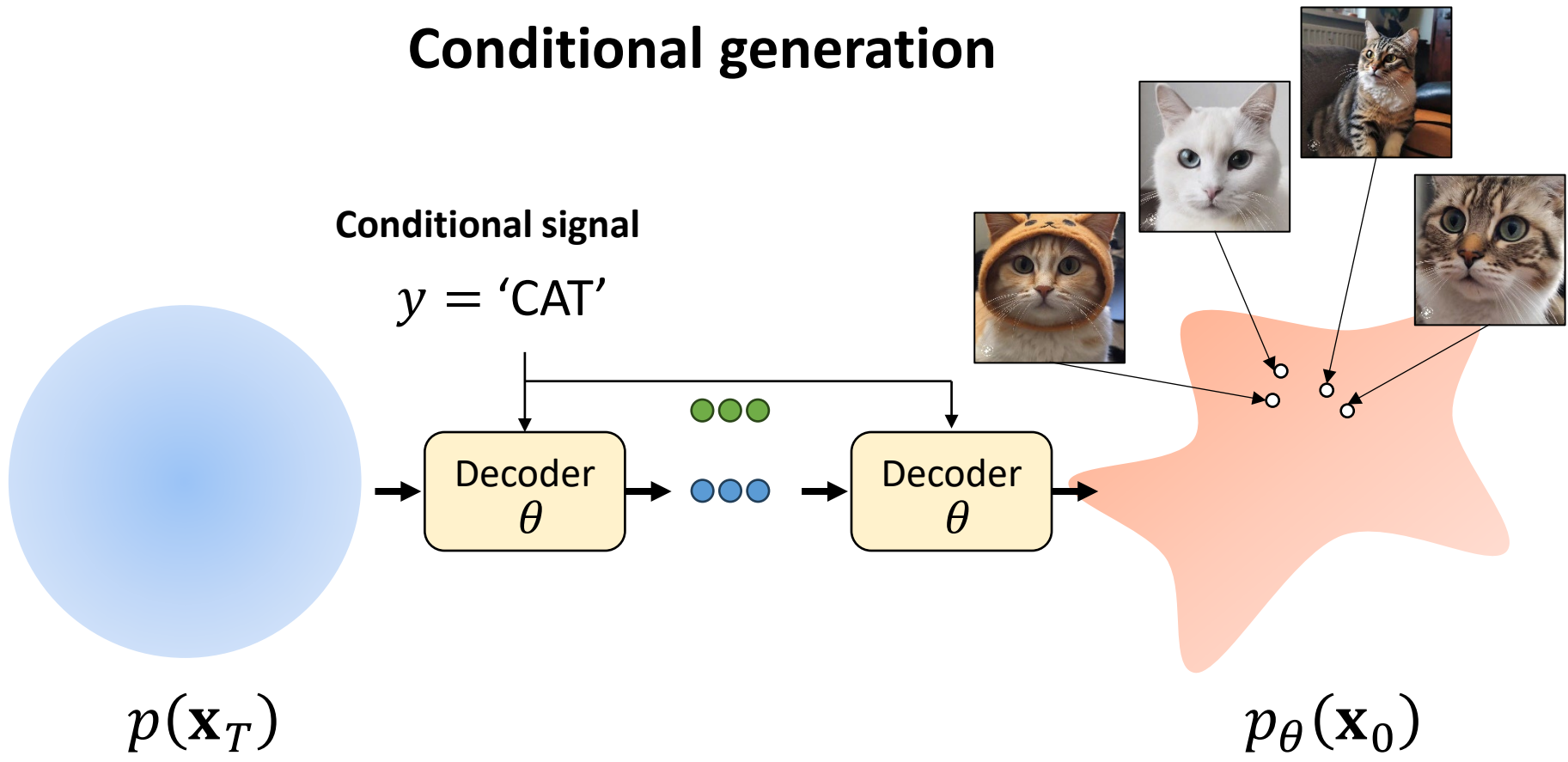


# Unconditional generation





# Conditional generation



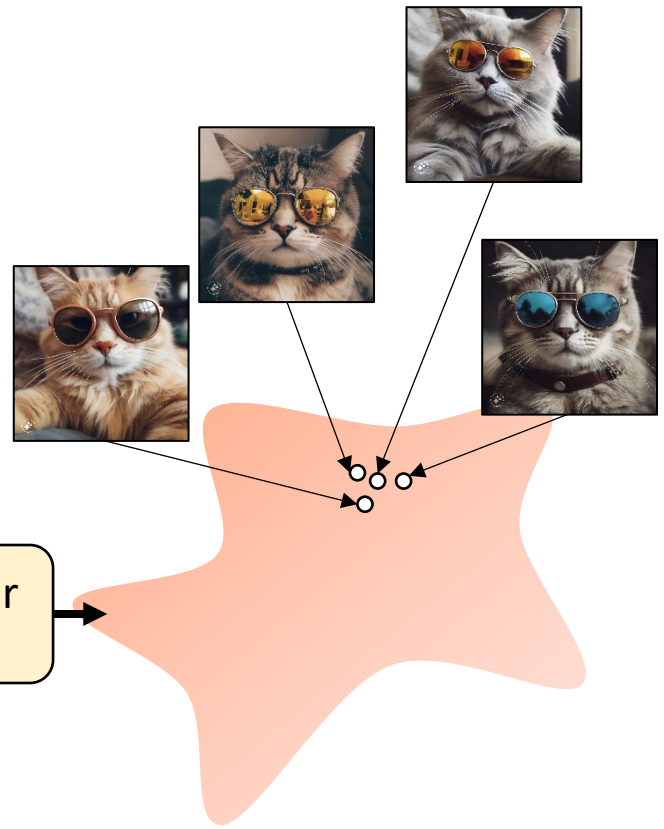
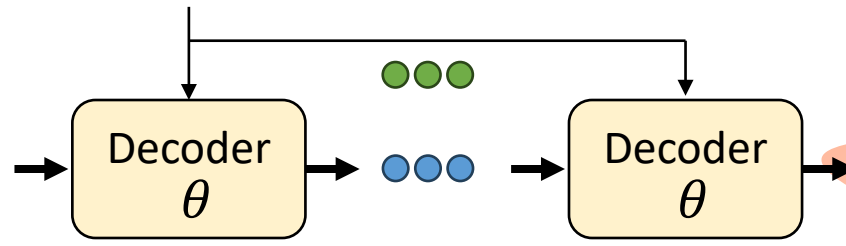
# Conditional generation

Conditional signal

$y = \text{"A cat wearing sunglasses"}$




$p(\mathbf{x}_T)$



$p_\theta(\mathbf{x}_0)$

# Unconditional generation

$\mathbf{x}_t$  

$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$

$s_\theta(\mathbf{x}_t, t) \approx \nabla \log p(\mathbf{x}_t)$   
Unconditional score

$$\mu_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t + \frac{1 - \alpha_t}{\sqrt{\alpha_t}} s_\theta(\mathbf{x}_t, t)$$

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



# Conditional generation

$\mathbf{x}_t$  

$$s_{\theta}(\mathbf{x}_t, t, \mathbf{y}) \approx \nabla \log p(\mathbf{x}_t | \mathbf{y})$$

Conditional score



# Conditional generation

$\mathbf{x}_t$

$$s_{\theta}(\mathbf{x}_t, t, \mathbf{y}) \approx \nabla \log p(\mathbf{x}_t | \mathbf{y})$$

Conditional score

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





# Conditional generation

$\mathbf{x}_t$

$$s_{\theta}(\mathbf{x}_t, t, \mathbf{y}) \approx \nabla \log p(\mathbf{x}_t | \mathbf{y})$$

Conditional score

$$\nabla \log p(\mathbf{x}_t | \mathbf{y}) = \nabla \log \left( \frac{p(\mathbf{x}_t) p(\mathbf{y} | \mathbf{x}_t)}{p(\mathbf{y})} \right)$$

Conditional score

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

# Conditional generation

$\mathbf{x}_t$

$$s_{\theta}(\mathbf{x}_t, t, \mathbf{y}) \approx \nabla \log p(\mathbf{x}_t | \mathbf{y})$$

Conditional score

$$\begin{aligned} \nabla \log p(\mathbf{x}_t | \mathbf{y}) &= \nabla \log p(\mathbf{x}_t) + \nabla \log p(\mathbf{y} | \mathbf{x}_t) \\ &\quad - \nabla \log p(\mathbf{y}) \end{aligned}$$

Conditional score

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

# Conditional generation



**Classifier  
Guidance**

$\mathbf{x}_t$

$$s_{\theta}(\mathbf{x}_t, t, \mathbf{y}) \approx \nabla \log p(\mathbf{x}_t | \mathbf{y})$$

Conditional score

$$\nabla \log p(\mathbf{x}_t | \mathbf{y}) = \nabla \log p(\mathbf{x}_t) + \nabla \log p(\mathbf{y} | \mathbf{x}_t)$$

Conditional score

Unconditional score

Adversarial gradient

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

# Conditional generation



**Classifier  
Guidance**

$\mathbf{x}_t$

$$s_{\theta}(\mathbf{x}_t, t, \mathbf{y}) \approx \nabla \log p(\mathbf{x}_t | \mathbf{y})$$

Conditional score

$$\nabla \log p_{\gamma}(\mathbf{x}_t | \mathbf{y}) = \nabla \log p(\mathbf{x}_t) + \gamma \nabla \log p(\mathbf{y} | \mathbf{x}_t)$$

Conditional score

Unconditional score

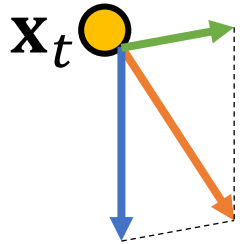
Adversarial gradient

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

# Conditional generation



Classifier  
Guidance



$$s_{\theta}(\mathbf{x}_t, t, \mathbf{y}) \approx \nabla \log p(\mathbf{x}_t | \mathbf{y})$$

Conditional score

$$\nabla \log p_{\gamma}(\mathbf{x}_t | \mathbf{y}) = \nabla \log p(\mathbf{x}_t) + \gamma \nabla \log p(\mathbf{y} | \mathbf{x}_t)$$

Conditional score

Unconditional score

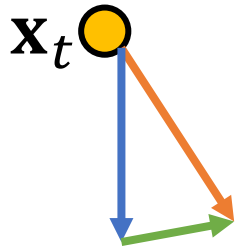
Adversarial gradient

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

# Conditional generation



## Classifier Guidance



$$s_{\theta}(\mathbf{x}_t, t, \mathbf{y}) \approx \nabla \log p(\mathbf{x}_t | \mathbf{y})$$

Conditional score

$$\nabla \log p_{\gamma}(\mathbf{x}_t | \mathbf{y}) = \nabla \log p(\mathbf{x}_t) + \gamma \nabla \log p(\mathbf{y} | \mathbf{x}_t)$$

Conditional score

Unconditional score

Adversarial gradient

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

# Conditional generation



Classifier  
Guidance



Classifier guidance scale = 1



Classifier guidance scale = 10

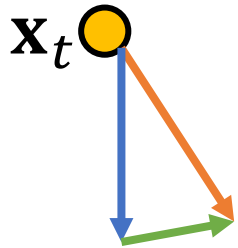
$p_{\text{data}}(\mathbf{x})$   
 $y = \text{"Pembroke Welsh corgi"}$

[Dhariwal & Nichol 2021]

# Conditional generation



**Classifier  
Guidance**



$$s_{\theta}(\mathbf{x}_t, t, \mathbf{y}) \approx \nabla \log p(\mathbf{x}_t | \mathbf{y})$$

Conditional score

$$\nabla \log p_{\gamma}(\mathbf{x}_t | \mathbf{y}) = \nabla \log p(\mathbf{x}_t) + \gamma \nabla \log p(\mathbf{y} | \mathbf{x}_t)$$

Conditional score

Unconditional score

Adversarial gradient

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



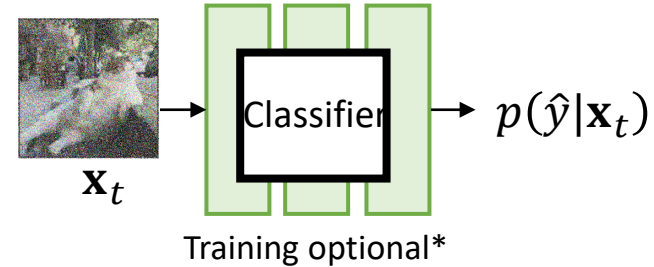
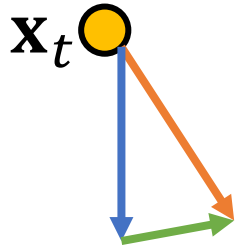
# Conditional generation



## Classifier Guidance

$$\nabla \log p_\gamma(\mathbf{x}_t | \mathbf{y}) = \nabla \log p(\mathbf{x}_t) + \gamma \nabla \log p(\mathbf{y} | \mathbf{x}_t)$$

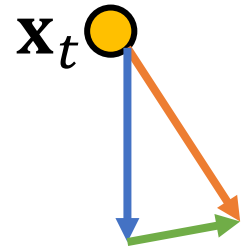
Conditional score      Unconditional score      Adversarial gradient



# Conditional generation

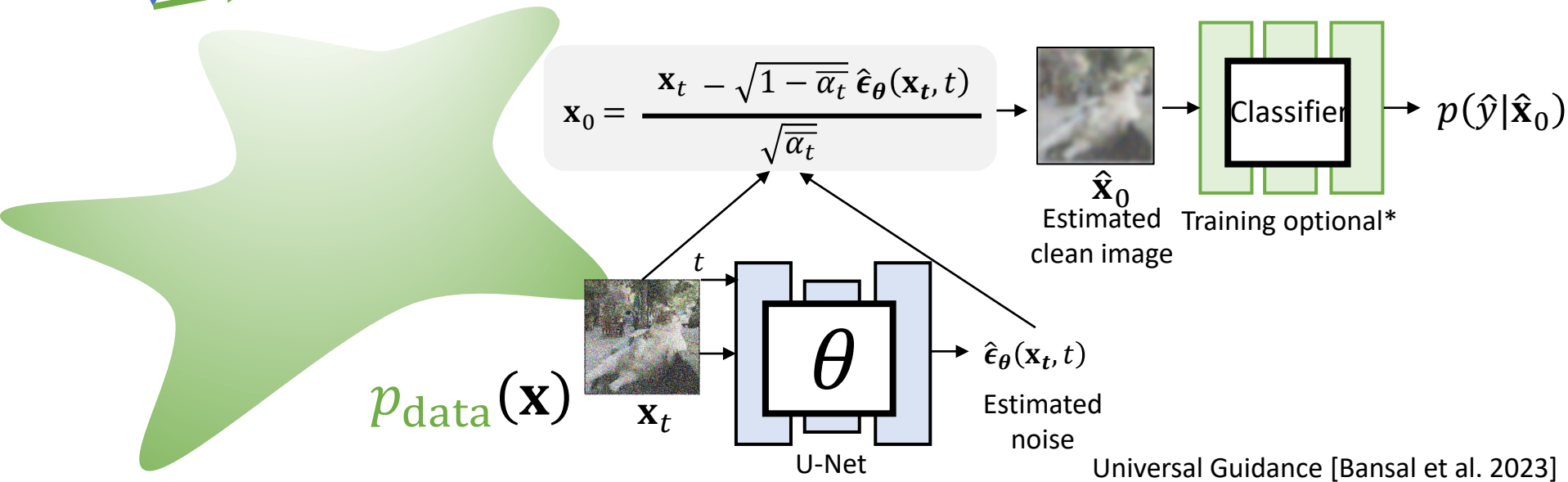


## Classifier Guidance



$$\nabla \log p_\gamma(\mathbf{x}_t | \mathbf{y}) = \nabla \log p(\mathbf{x}_t) + \gamma \nabla \log p(\mathbf{y} | \hat{\mathbf{x}}_0)$$

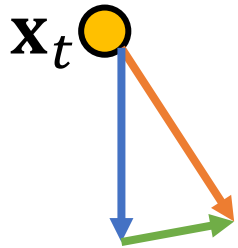
Conditional score      Unconditional score      Adversarial gradient



Universal Guidance [Bansal et al. 2023]

# Conditional Generation with Classifier-Guidance

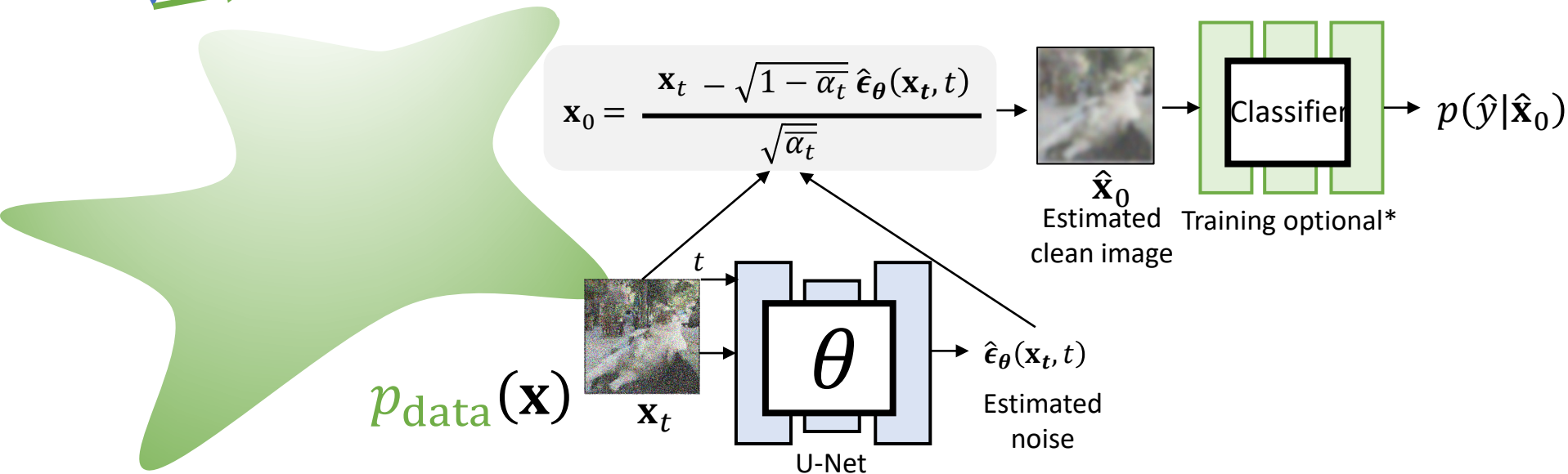
- Any price to pay?



## Classifier Guidance

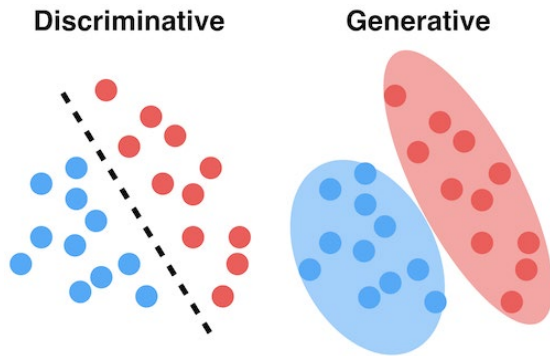
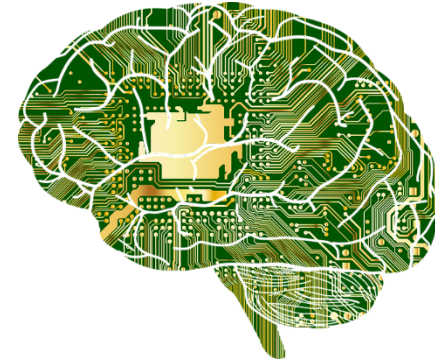
$$\nabla \log p_\gamma(\mathbf{x}_t | \mathbf{y}) = \nabla \log p(\mathbf{x}_t) + \gamma \nabla \log p(\mathbf{y} | \hat{\mathbf{x}}_0)$$

Conditional score
Unconditional score
Adversarial gradient

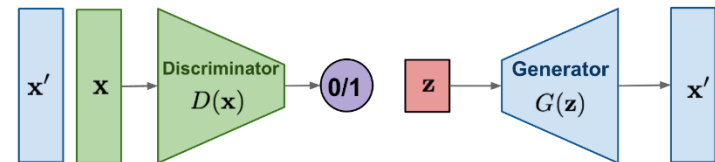


# What's to Be Covered Today...

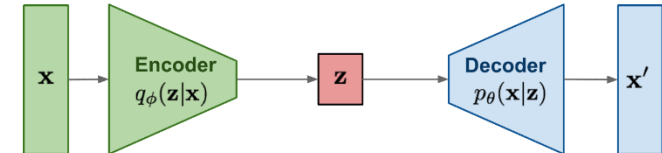
- Generative Models
  - Diffusion Model
    - Conditional Diffusion Model
      - Classifier Guidance
      - Classifier-Free Guidance
    - Personalization via Diffusion Model
  - Generative Adversarial Network



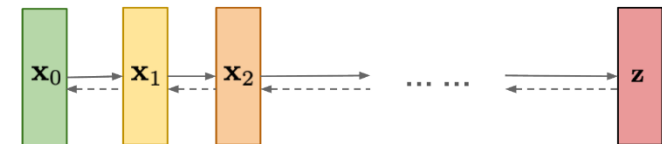
**GAN: Adversarial training**



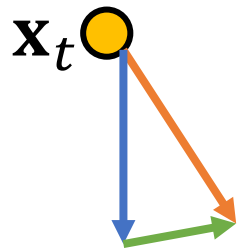
**VAE: maximize variational lower bound**



**Diffusion models:**  
Gradually add Gaussian noise and then reverse



# Conditional generation



$$\nabla \log p_{\gamma}(\mathbf{x}_t | \mathbf{y}) = \nabla \log p(\mathbf{x}_t) + \gamma \nabla \log p(\mathbf{y} | \mathbf{x}_t)$$

Conditional score      Unconditional score      Adversarial gradient



## Classifier Guidance

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

# Conditional generation

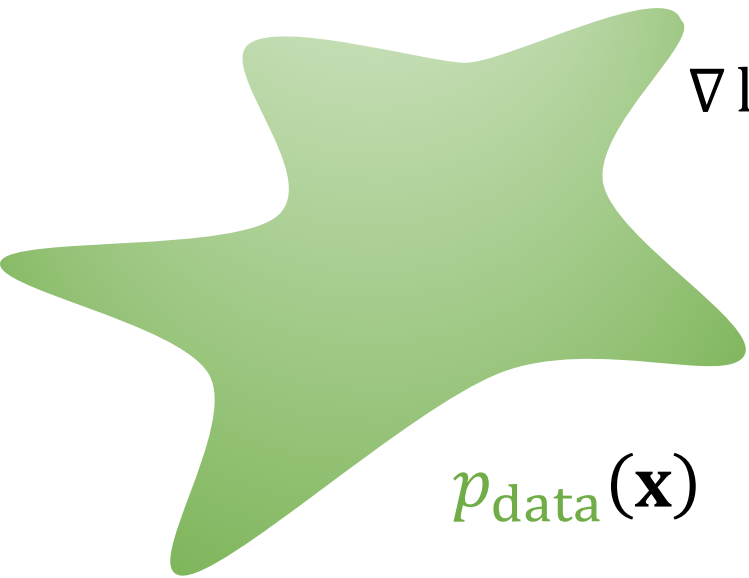
$\mathbf{x}_t$  

$$\nabla \log p_\gamma(\mathbf{x}_t | \mathbf{y}) = \nabla \log p(\mathbf{x}_t) + \gamma \nabla \log p(\mathbf{y} | \mathbf{x}_t)$$

Conditional score

Unconditional score

Adversarial gradient


$$\nabla \log p(\mathbf{y} | \mathbf{x}_t) = \nabla \log \left( \frac{p(\mathbf{x}_t | \mathbf{y}) p(\mathbf{y})}{p(\mathbf{x}_t)} \right)$$

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

# Conditional generation

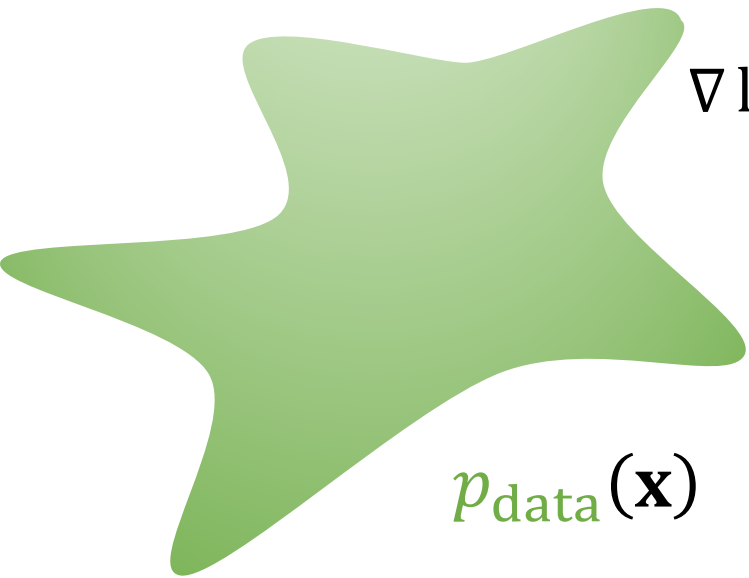
$\mathbf{x}_t$  

$$\nabla \log p_\gamma(\mathbf{x}_t | \mathbf{y}) = \nabla \log p(\mathbf{x}_t) + \gamma \nabla \log p(\mathbf{y} | \mathbf{x}_t)$$

Conditional score

Unconditional score

Adversarial gradient


$$\nabla \log p(\mathbf{y} | \mathbf{x}_t) = \nabla \log p(\mathbf{x}_t | \mathbf{y}) - \nabla \log p(\mathbf{x}_t) + \nabla \log p(\mathbf{y})$$

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

# Conditional generation

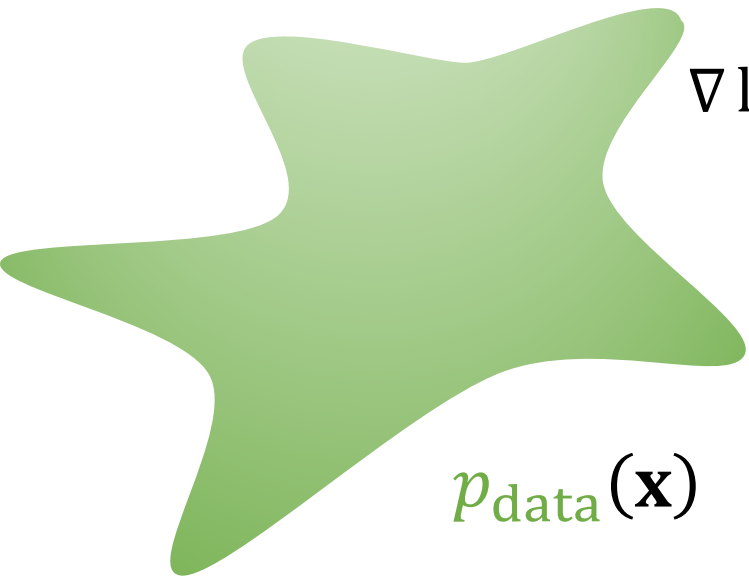
$\mathbf{x}_t$  

$$\nabla \log p_\gamma(\mathbf{x}_t | \mathbf{y}) = \nabla \log p(\mathbf{x}_t) + \gamma \nabla \log p(\mathbf{y} | \mathbf{x}_t)$$

Conditional score

Unconditional score

Adversarial gradient


$$\nabla \log p(\mathbf{y} | \mathbf{x}_t) = \nabla \log p(\mathbf{x}_t | \mathbf{y}) - \nabla \log p(\mathbf{x}_t)$$

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



# Conditional generation

$\mathbf{x}_t$  

$$\nabla \log p_\gamma(\mathbf{x}_t | \mathbf{y}) = \nabla \log p(\mathbf{x}_t) + \gamma(\nabla \log p(\mathbf{x}_t | \mathbf{y}) - \nabla \log p(\mathbf{x}_t))$$

Conditional score

Unconditional score

Adversarial gradient



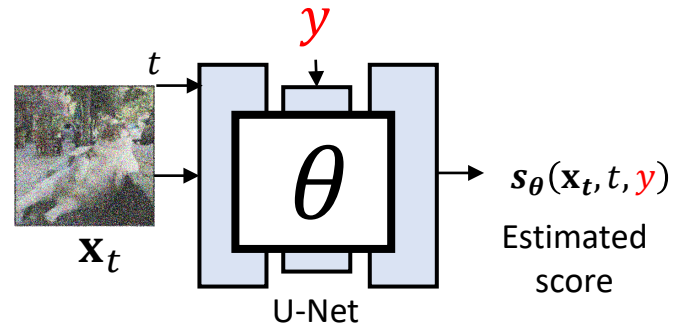
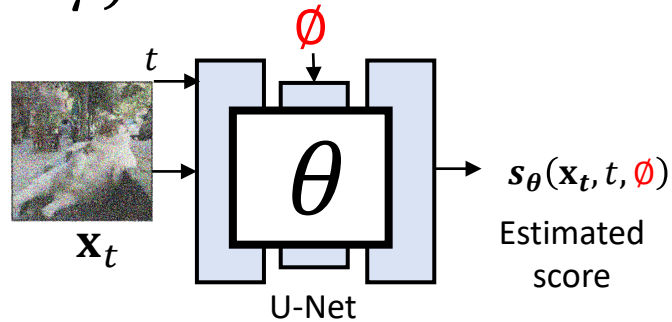
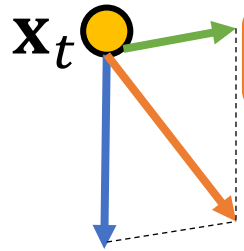
$p_{\text{data}}(\mathbf{x})$

# Conditional generation

## Classifier-free Guidance

$$\nabla \log p_\gamma(\mathbf{x}_t | \mathbf{y}) = (1 - \gamma) \nabla \log p(\mathbf{x}_t) + \gamma \nabla \log p(\mathbf{x}_t | \mathbf{y})$$

Conditional score      Unconditional score      Conditional score



# Conditional generation

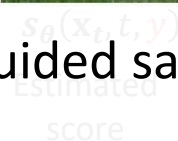
## Classifier-free Guidance



Unguided samples  $(x)$



Guided samples



[Ho and Salimans 2021]

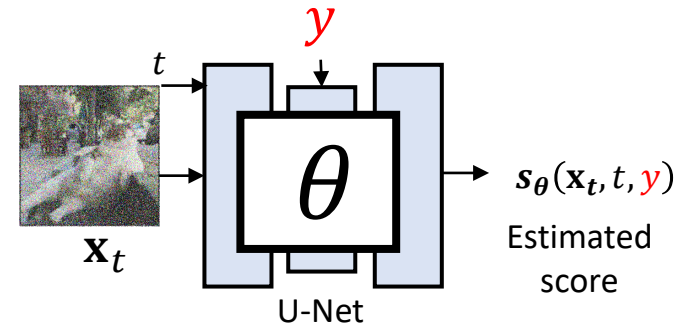
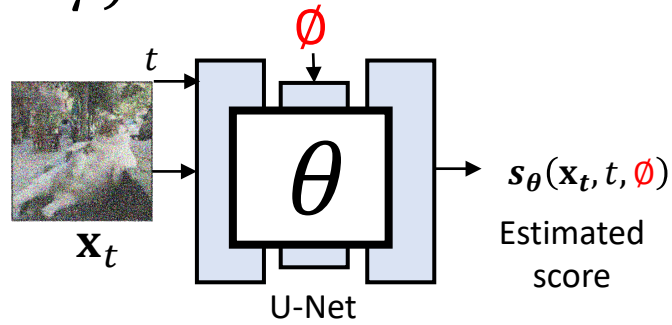
# Conditional Generation with Classifier-Free Guidance

- Any price to pay?

## Classifier-free Guidance

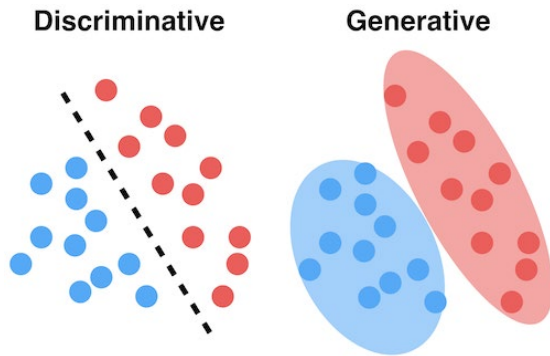
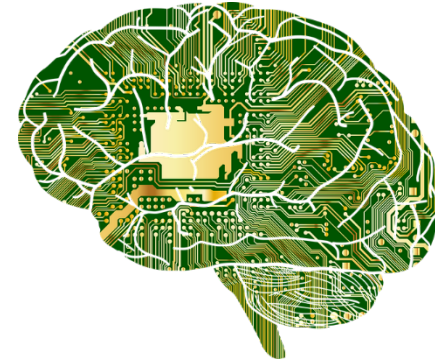
$$\nabla \log p_\gamma(\mathbf{x}_t | \mathbf{y}) = (1 - \gamma) \nabla \log p(\mathbf{x}_t) + \gamma \nabla \log p(\mathbf{x}_t | \mathbf{y})$$

Conditional score
Unconditional score
Conditional score

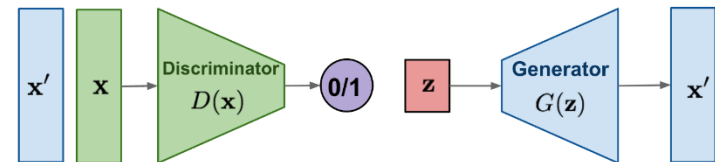


# What's to Be Covered Today...

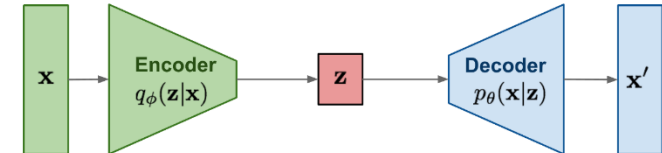
- Generative Models
  - Diffusion Model
    - Conditional Diffusion Model
      - Classifier Guidance
      - Classifier-Free Guidance
      - Text/Image Guidance
    - Personalization via Diffusion Model
  - Generative Adversarial Network



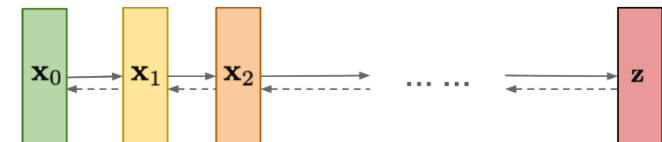
**GAN: Adversarial training**



**VAE: maximize variational lower bound**

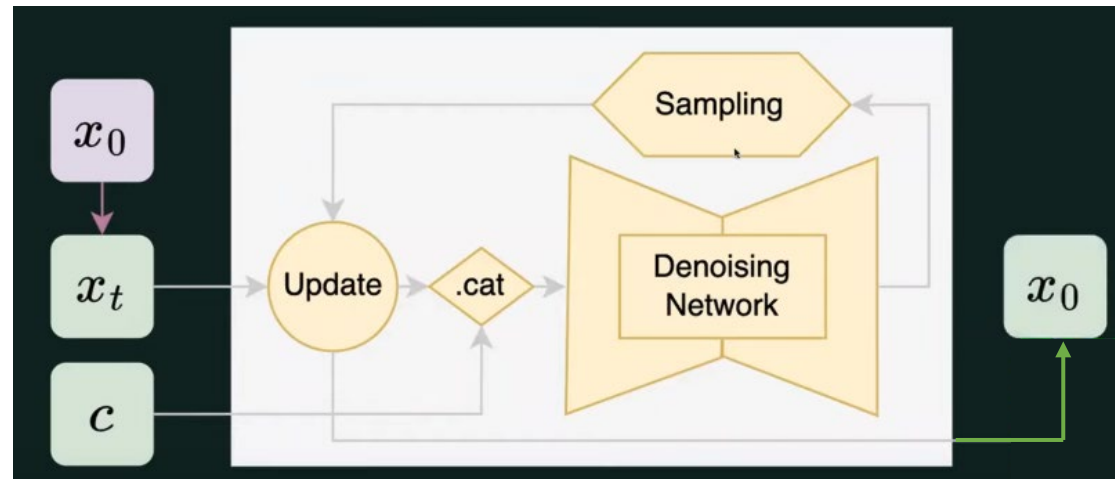
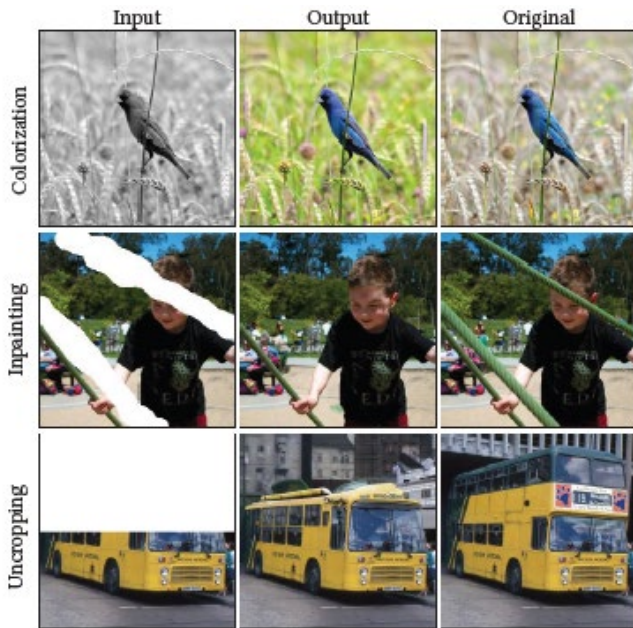


**Diffusion models:**  
Gradually add Gaussian noise and then reverse



# Conditional Generation with Image Guidance

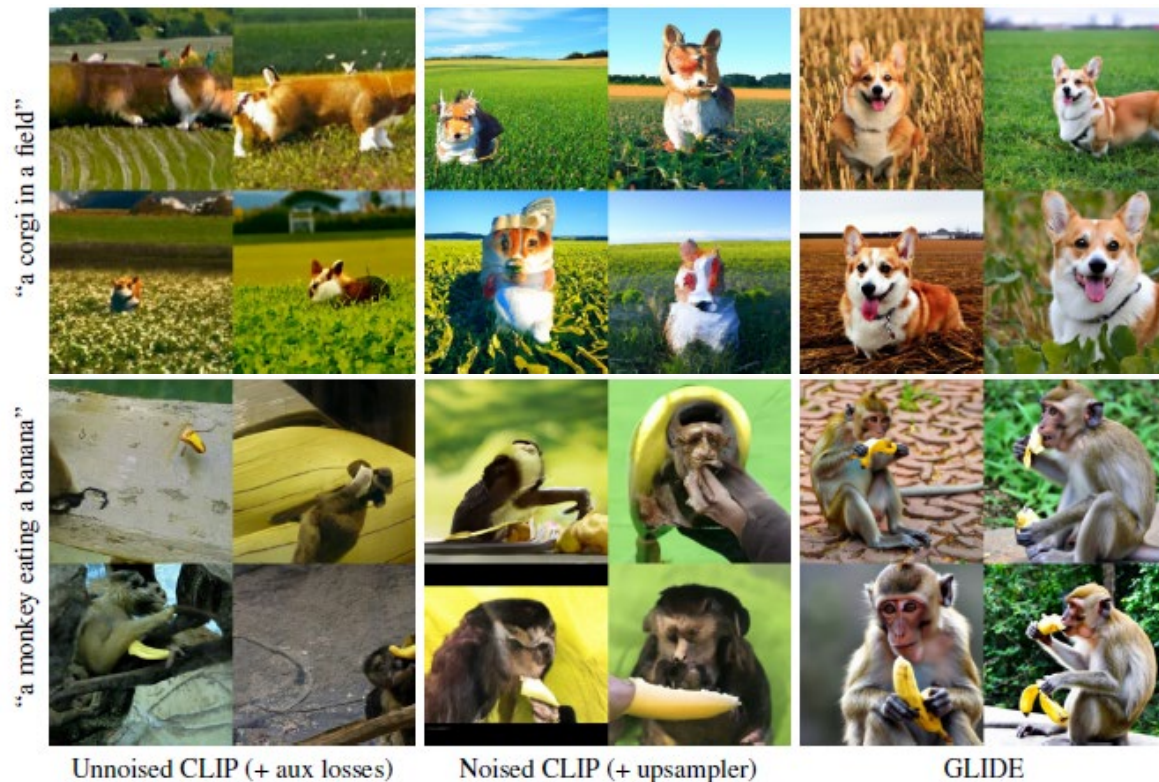
- *Palette: Image-to-Image Diffusion Models*, Google Research, arXiv 2022
  - Applications: colorization, inpainting, outpainting, etc.
  - Input image as condition (via concatenation)





# Conditional Generation with Text Guidance

- *GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models*, OpenAI, arXiv 2022
  - CLIP (Contrastive Language-Image Pretraining) is previously proposed to measure alignment between text and image inputs
  - Classifier guidance -> CLIP guidance (not training-free)
  - What is CLIP? (see next slide)



# CLIP: Contrastive Language-Image Pretraining

- *Learning Transferable Visual Models From Natural Language Supervision*, OpenAI, NeurIPS WS 2021 (w/ 22000+ citations)
- Why not just CNN?
  - Require **annotated data** for training image classification
  - **Domain gap** between closed and open-world domain data
  - Lack of ability for **zero-shot classification**



let



geNet V2



ImageNet Rendition



ObjectNet



ImageNet Sketch



ImageNet Adversarial



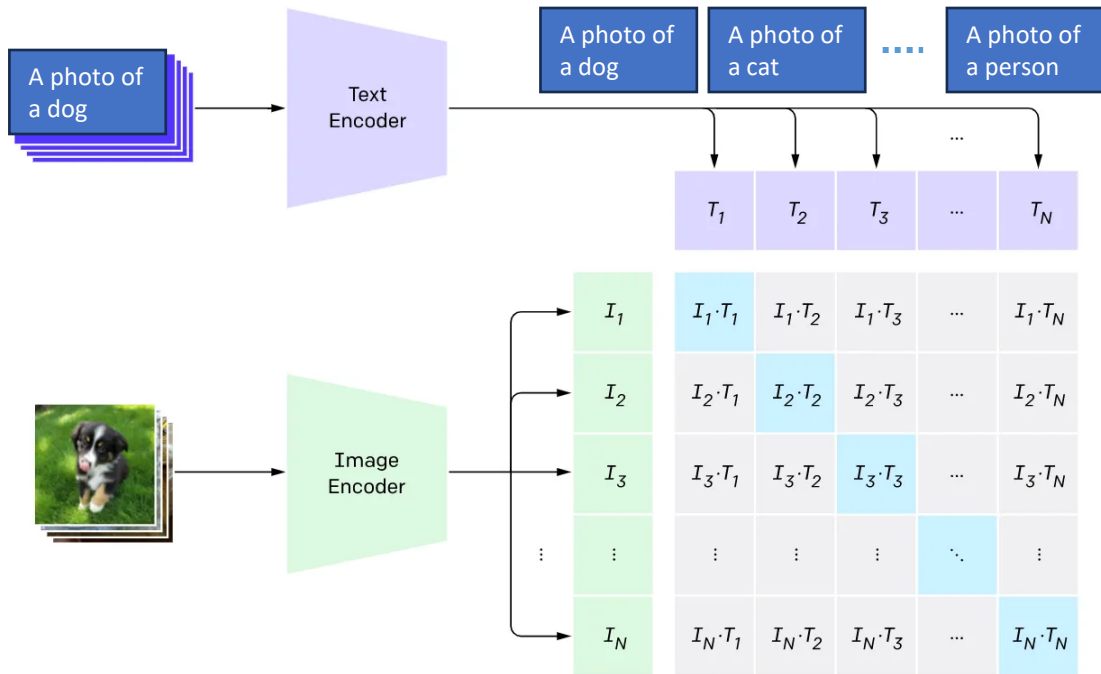


# CLIP (cont'd)

- Objectives

- Cross-domain contrastive learning from large-scale image-language data
- Next-token prediction (**what's this & why?**); will talk more about this for the lecture of Transformer

## 1. Contrastive pre-training

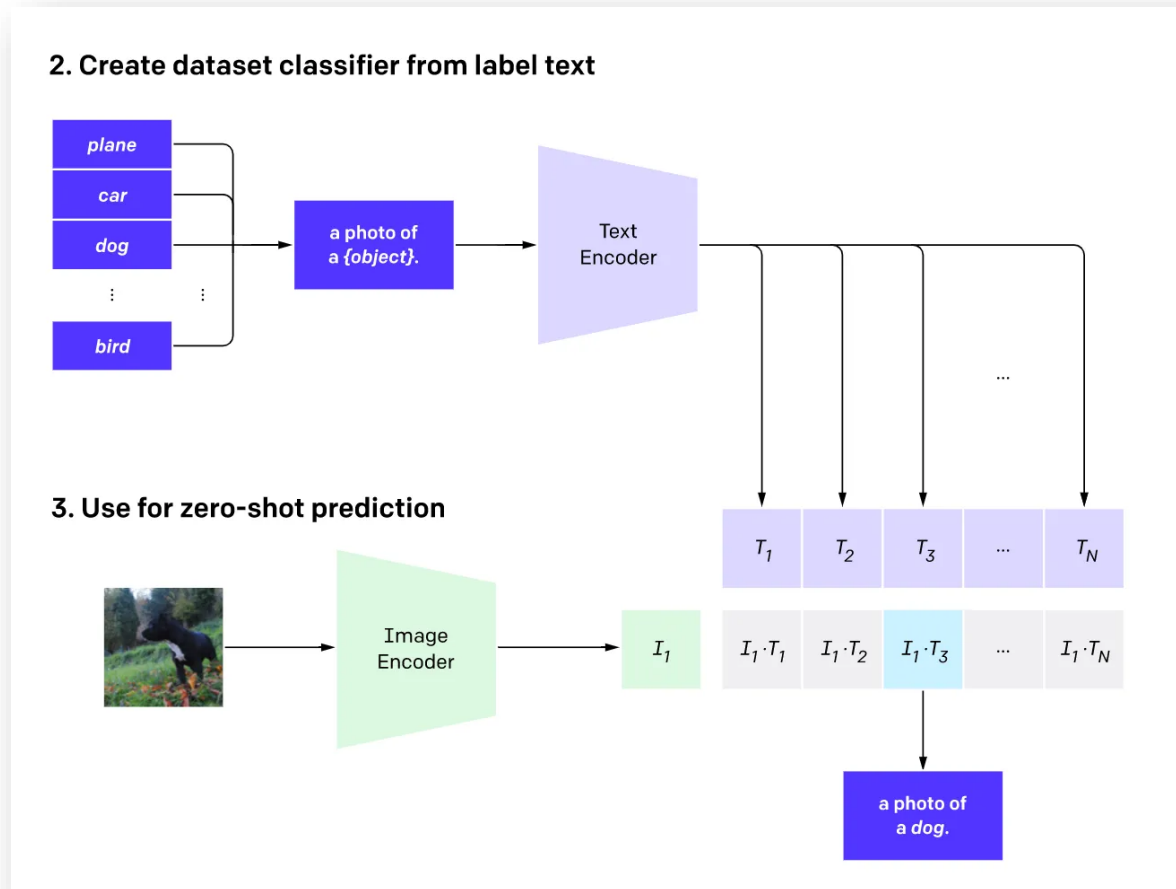


## 2. Next-token prediction

e.g., a \_\_\_\_;  
a photo \_\_\_\_;  
a photo of \_\_\_\_\_, etc.

# CLIP (cont'd)

- (Zero-shot) Inference:



- Potential concerns/disadvantages of CLIP?

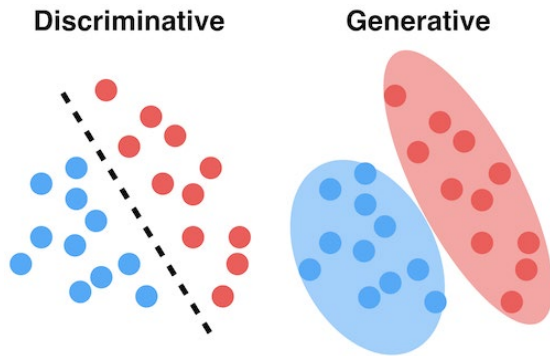
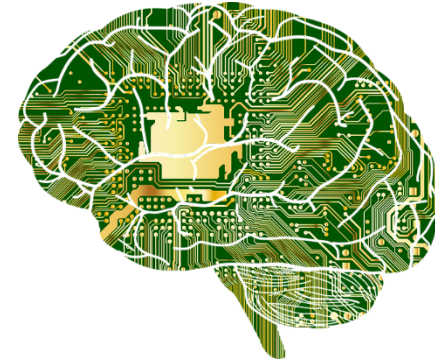
# Questions for Image Generation

- How to evaluate your unconditional image generation results?
- How to evaluate your conditional image generation results?
- Any objective/subjective and quantitative/qualitative evaluation?

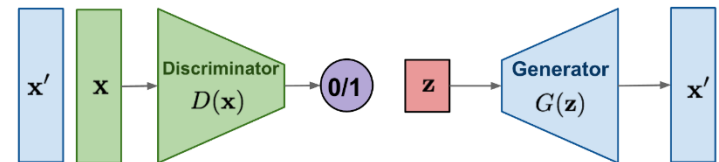


# What's to Be Covered Today...

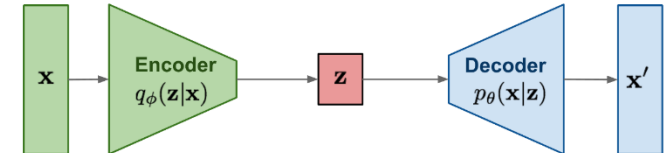
- Generative Models
  - Diffusion Model
    - Conditional Diffusion Model
    - Personalization via Diffusion Model
  - Generative Adversarial Network



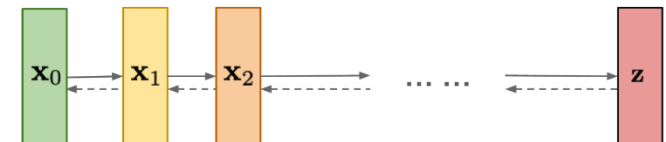
**GAN: Adversarial training**



**VAE: maximize variational lower bound**



**Diffusion models:**  
Gradually add Gaussian noise and then reverse



# Diffusion Model for Personalization (1): Textual Inversion

- Proposed by NV Research, ICLR 2023
- Goal: Learn a special token (e.g.,  $S_*$ ) to represent the concept of interest

Input samples  $\xrightarrow{\text{invert}}$  " $S_*$ "

→

“An oil painting of  $S_*$ ”

“App icon of  $S_*$ ”

“Elmo sitting in the same pose as  $S_*$ ”

“Crochet  $S_*$ ”

Input samples  $\xrightarrow{\text{invert}}$  " $S_*$ "

→

“Painting of two  $S_*$  fishing on a boat”

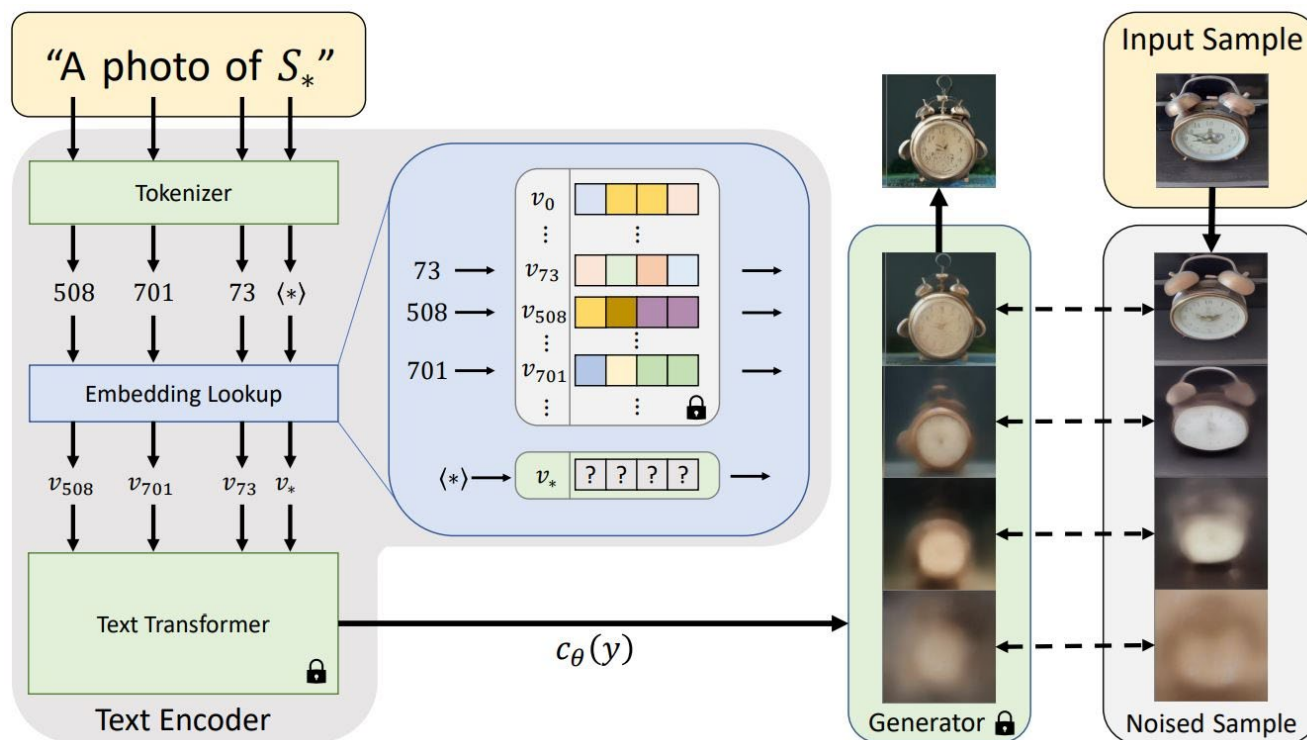
“A  $S_*$  backpack”

“Banksy art of  $S_*$ ”

“A  $S_*$  themed lunchbox”

# Diffusion Model for Personalization (1): Textual Inversion (cont'd)

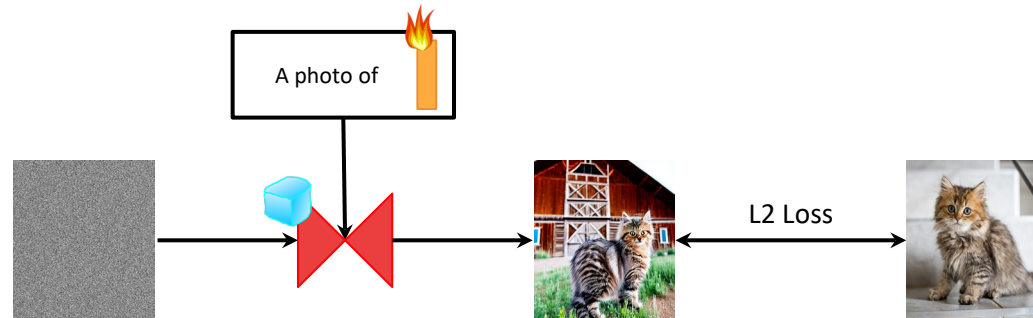
- Learning of special token  $S^*$ 
  - Pre-train and fix text encoder & diffusion model (i.e., generator)
  - Randomly initialize a token as the text encoder input
  - Optimize this token via image reconstruction objectives



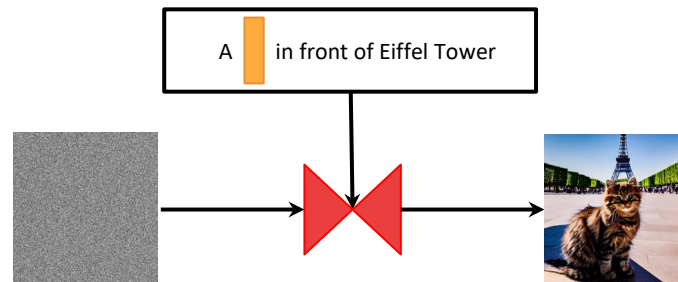
# Diffusion Model for Personalization (1): Textual Inversion (cont'd)

- Learning of special token  $S^*$ 
  - Pre-train and fix text encoder & diffusion model (i.e., generator)
  - Randomly initialize a token as the text encoder input
  - Optimize this token via image reconstruction objectives

- Training:



- Inference:

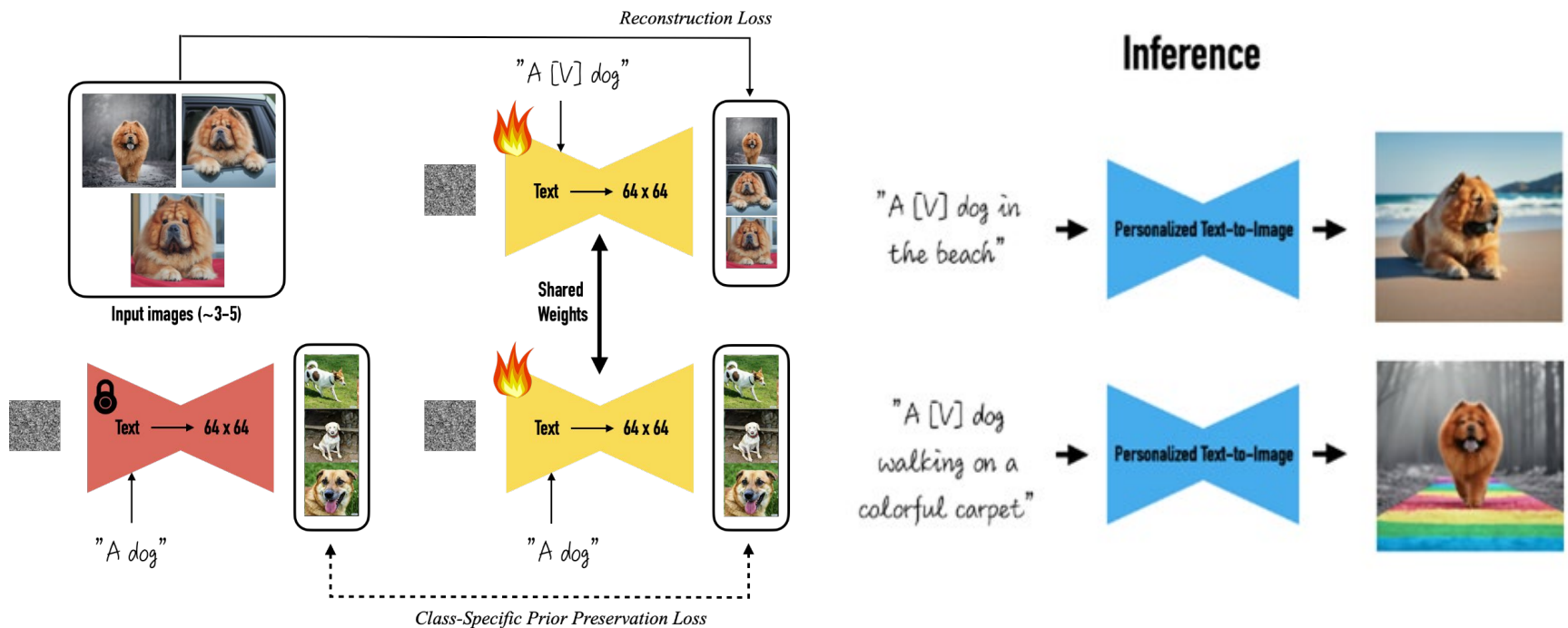


- Any potential concern?



# Diffusion Model for Personalization (2): DreamBooth

- Proposed by Google Research, CVPR 2023
- Finetune the diffusion model w/ a fixed token to represent the image concept
  - Determine and fix a rare token (e.g., [V])
  - Finetune the diffusion model for image restoration objectives
  - Enforce a class-specific prior (**why?**)

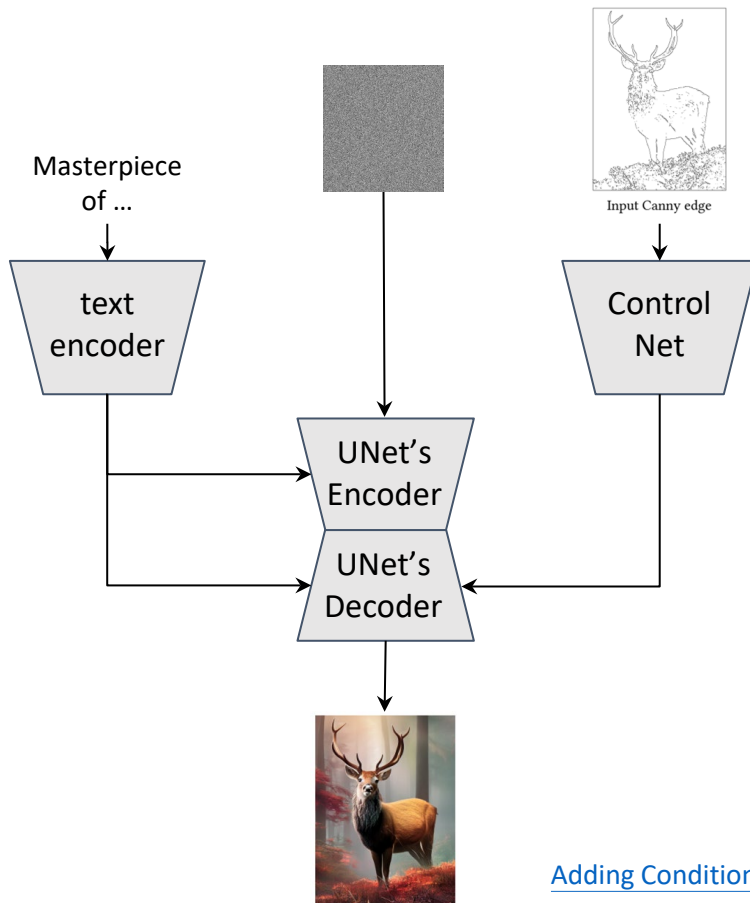


- **Any concern?**



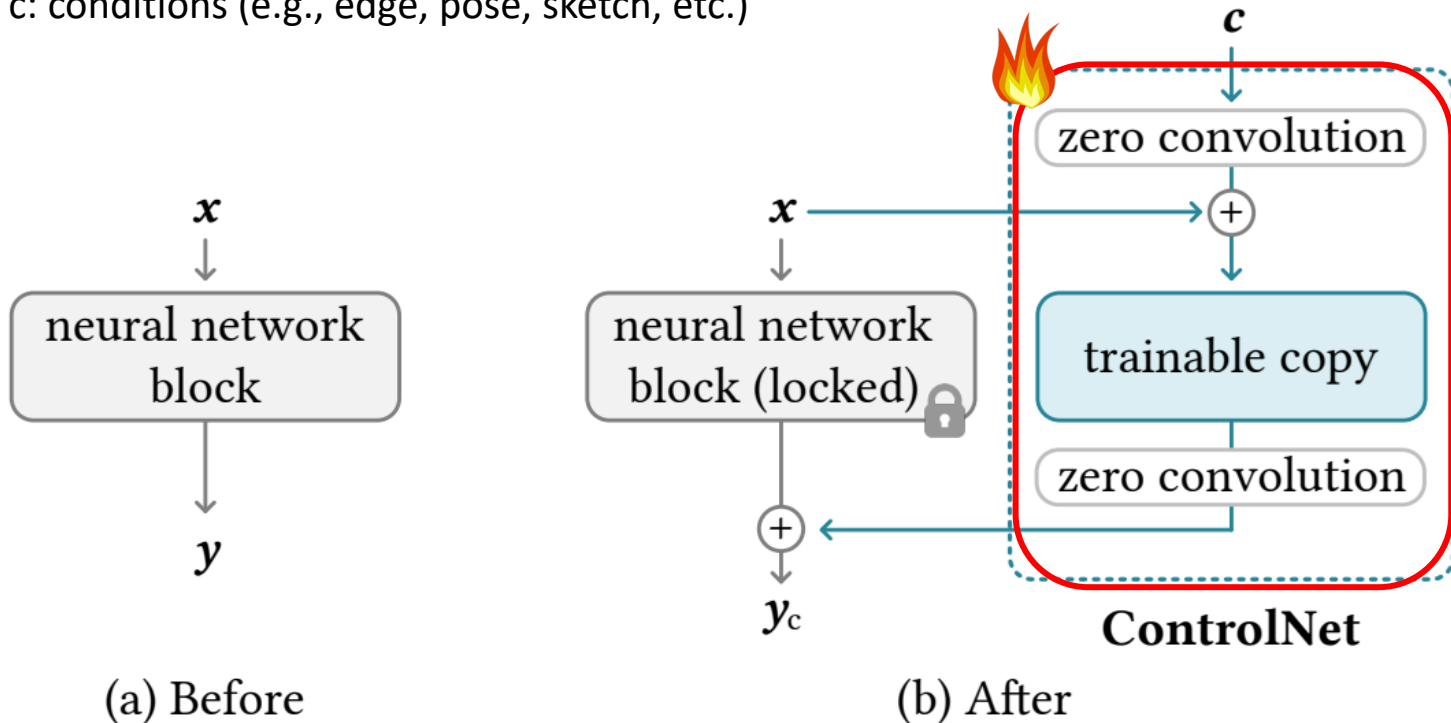
# Diffusion Model for Personalization (3): ControlNet

- Proposed by Stanford, ICCV 2023
- Goal: personalization via user-determined condition



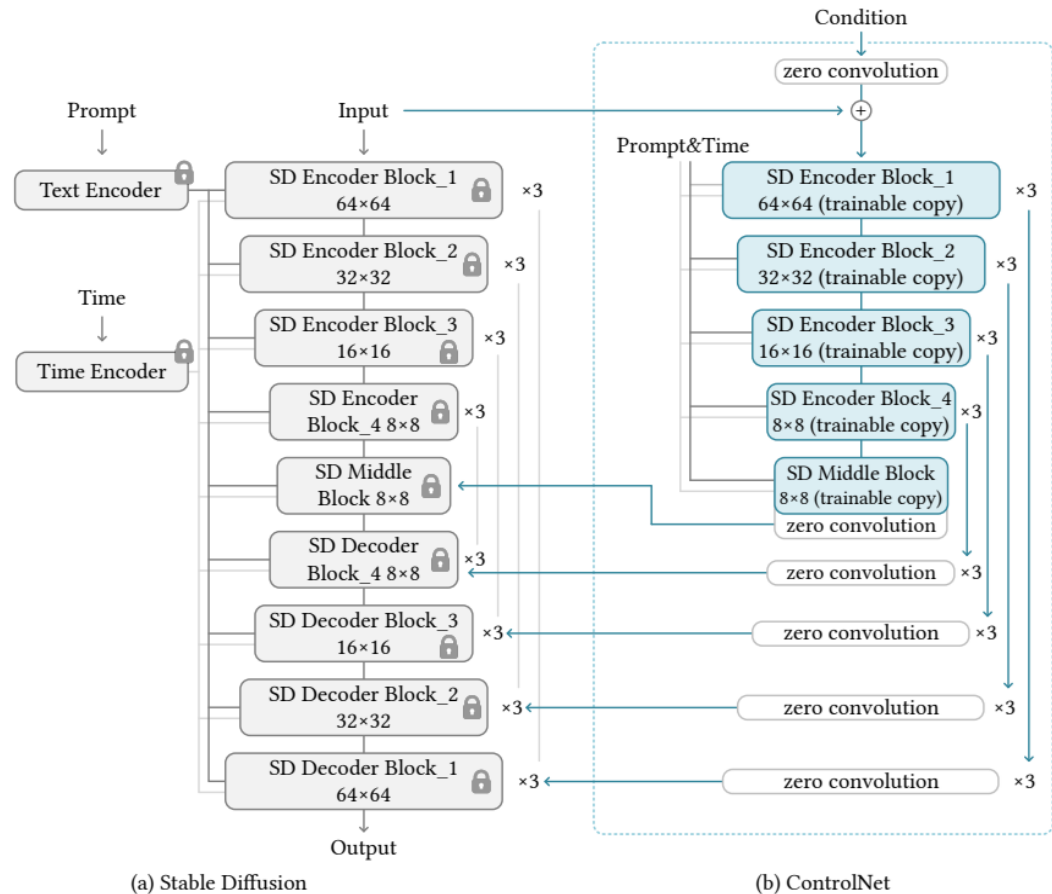
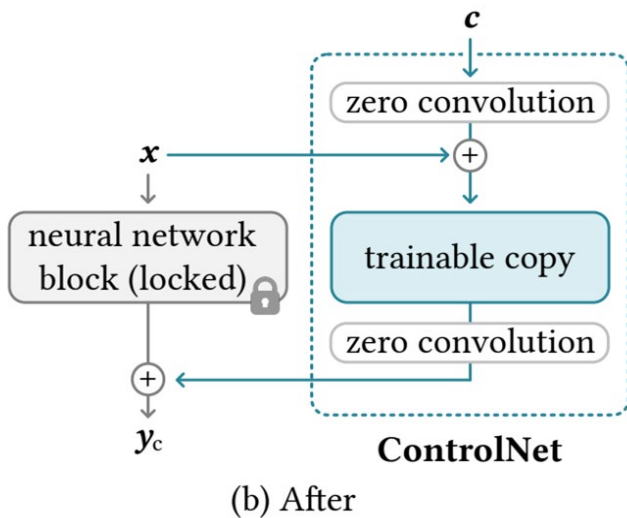
# Diffusion Model for Personalization (3): ControlNet

- Initialized from UNet's encoder
- Notations:
  - $x$ : input noise of each layer
  - $y$ : output noise of each layer
  - $c$ : conditions (e.g., edge, pose, sketch, etc.)



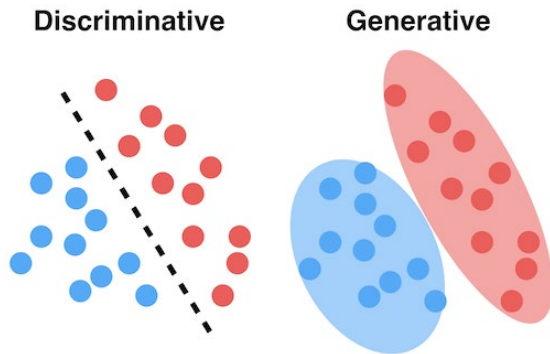
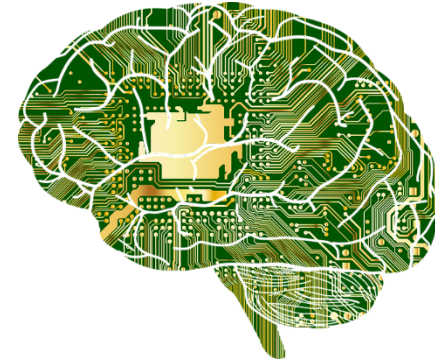
# Diffusion Model for Personalization (3): ControlNet

- Initialized from UNet's encoder
- Notations:
  - $x$ : input noise of each layer
  - $y$ : output noise of each layer
  - $c$ : conditions (e.g., edge,, etc.)

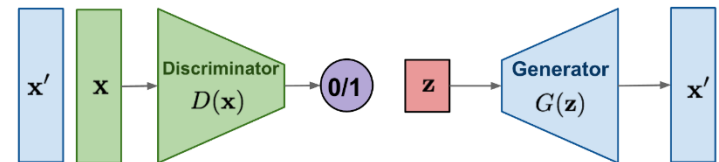


# What's to Be Covered Today...

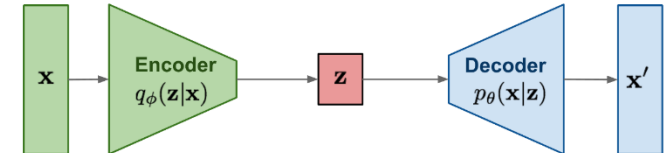
- Generative Models
  - Diffusion Model
    - Conditional Diffusion Model
      - Classifier Guidance
      - Classifier-Free Guidance
      - Text/Image Guidance
    - Personalization via Diffusion Model
  - Generative Adversarial Network



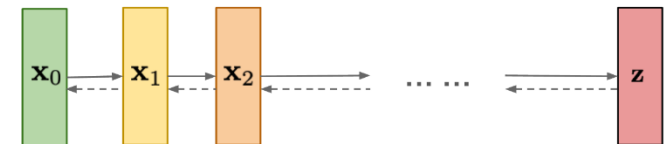
**GAN: Adversarial training**



**VAE: maximize variational lower bound**

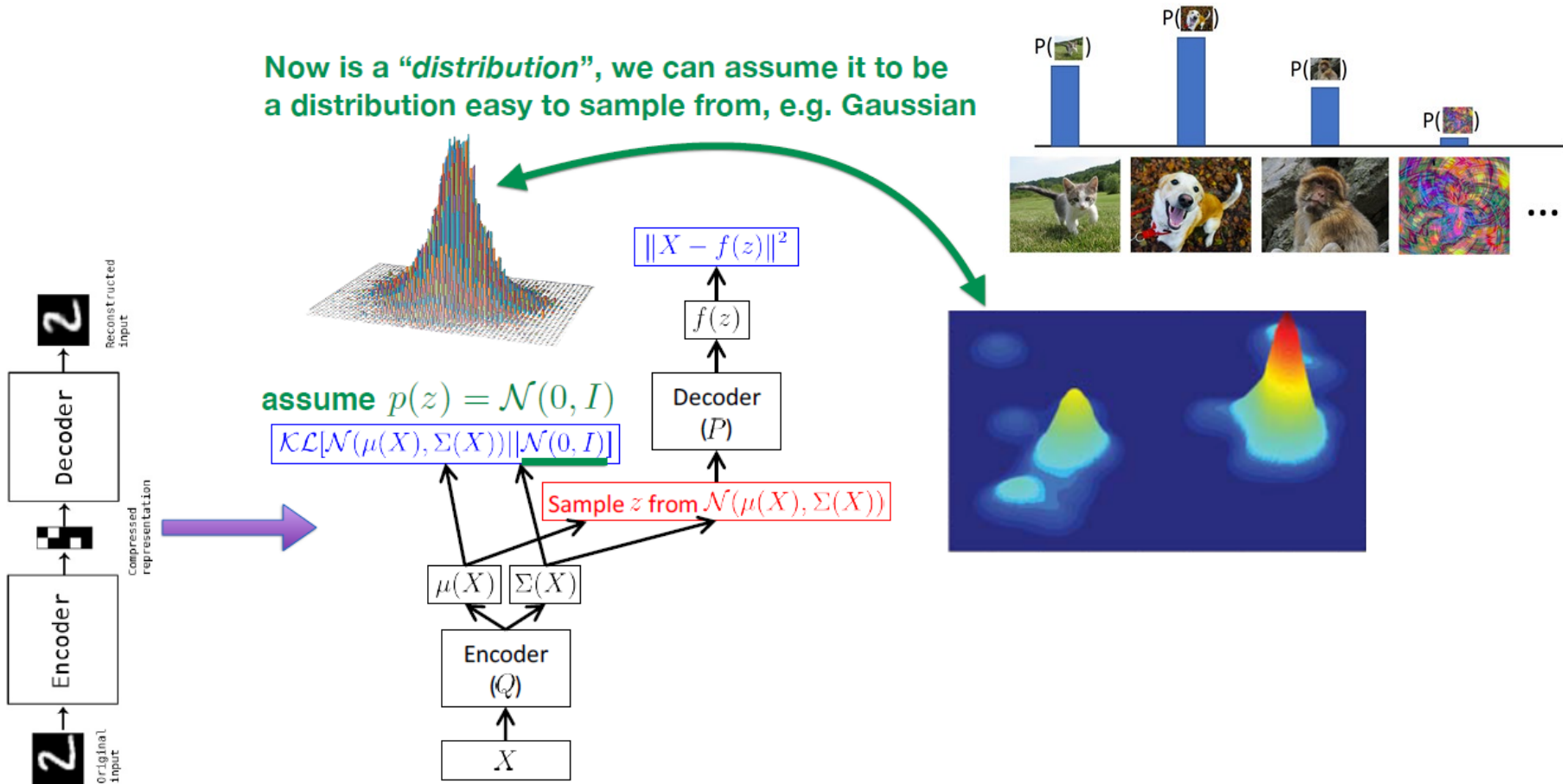


**Diffusion models:**  
Gradually add Gaussian noise and then reverse



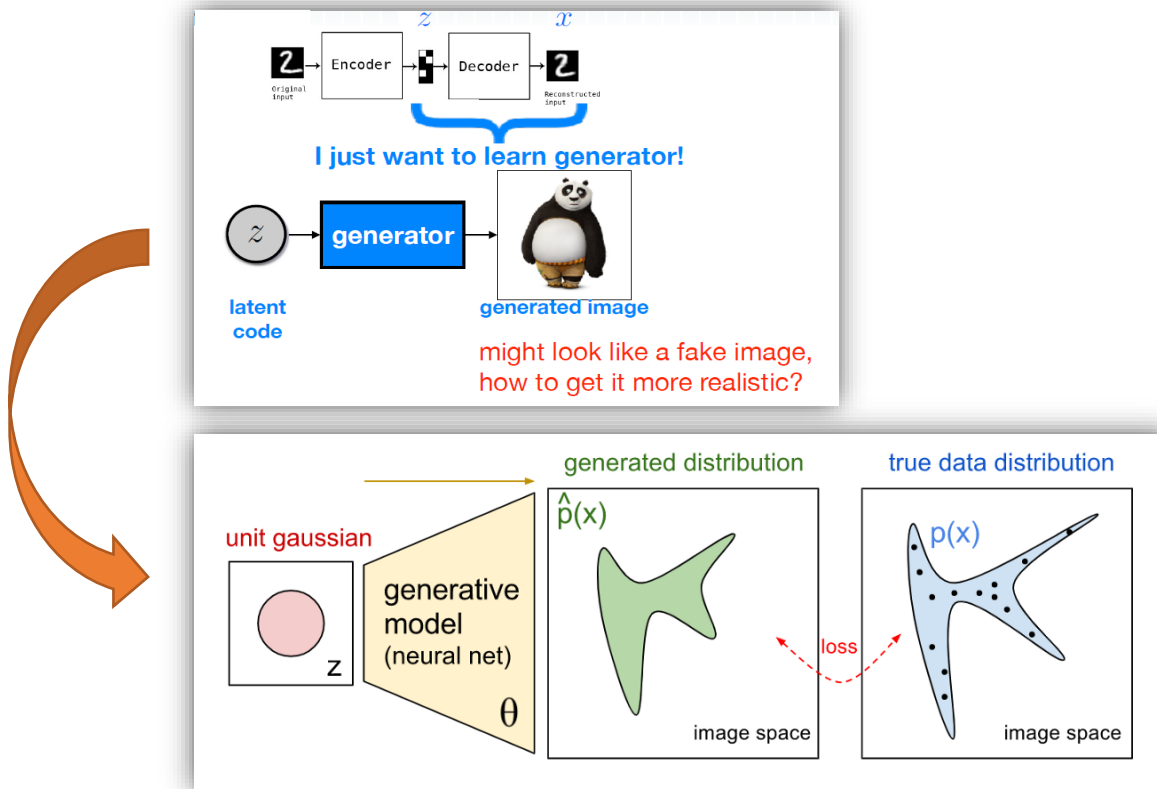
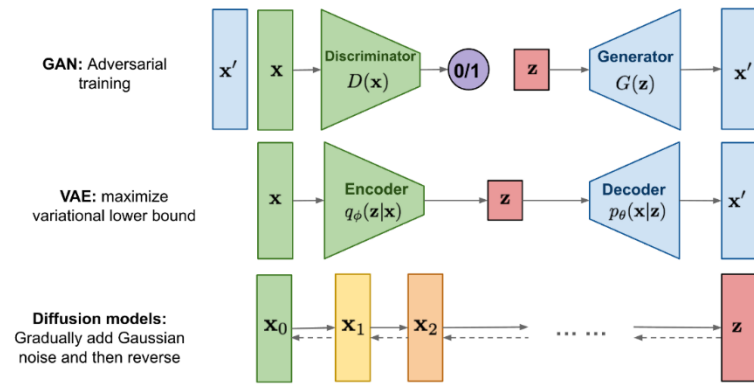
# From VAE to Generative Adversarial Networks (GAN)

Now is a "distribution", we can assume it to be a distribution easy to sample from, e.g. Gaussian



# From VAE to GAN

- Remarks
  - We only need the decoder/generator in practice.
  - We prefer fast generation.
  - How do we know if the output images are sufficiently good?



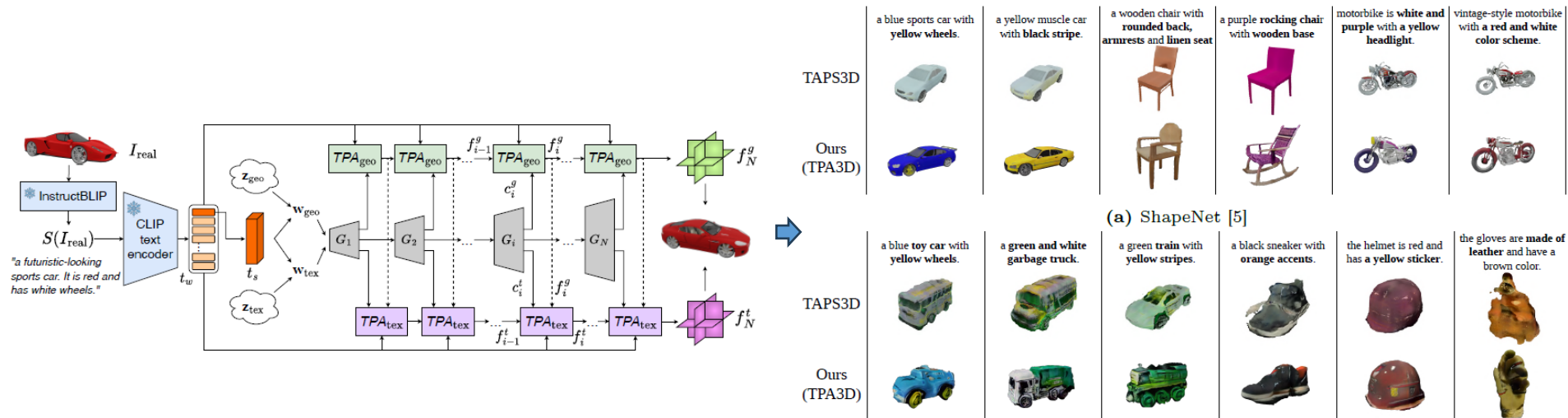
# GAN is NOT an outdated DL technology

- Remarks

- We only need the decoder/generator in practice.
- We prefer fast generation.
- How do we know if the output images are sufficiently good?

- Example

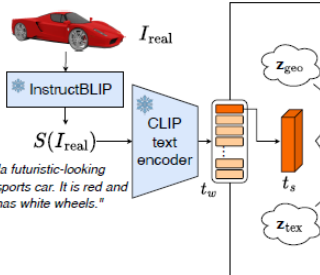
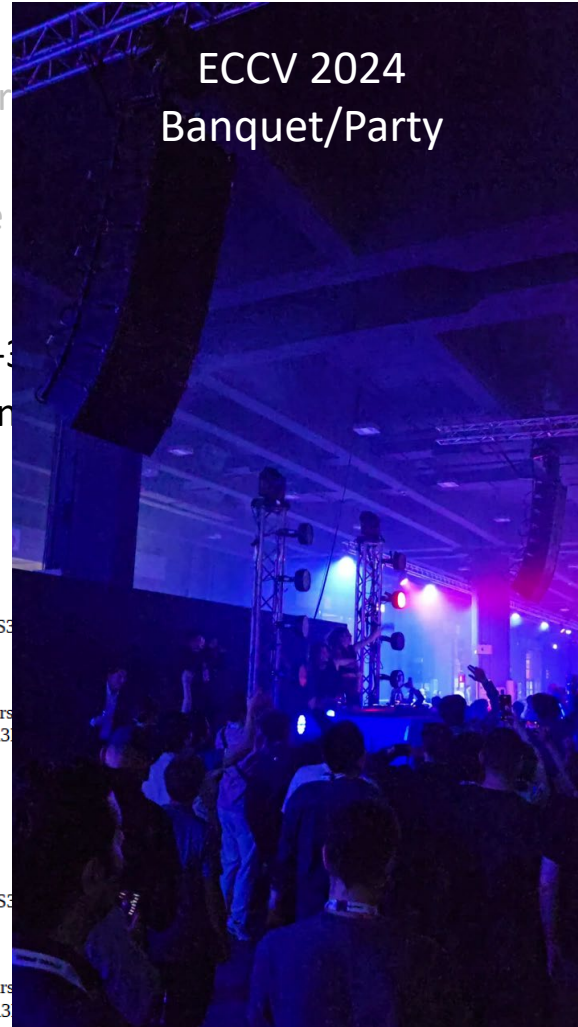
- TPA3D: Triplane Attention for Fast Text-to-3D Generation, [ECCV 2024](#)
- Bin-Shih Wu, Hong-En Chen, Shen-Yu Huang, and Y.-C. Frank Wang





# GAN is NOT an outdated DL technology

- Remarkable performance in prompt-to-image generation
- We demonstrate that GANs can generate high-quality images
- We show that GANs can generate images that are more diverse than those generated by text-to-image models
- How to use GANs for image generation
- Examples of GAN-generated images
- TPA3
- Bir

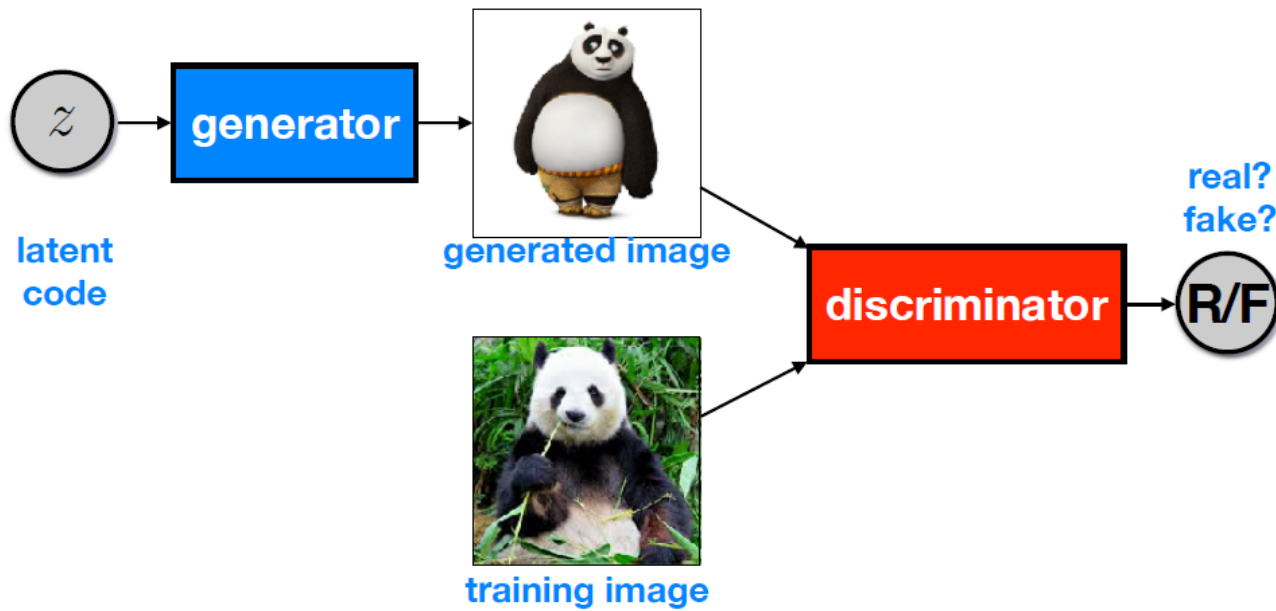


motorbike is white and purple with a yellow headlight.	vintage-style motorbike with a red and white color scheme.
the helmet is red and has a yellow sticker.	the gloves are made of leather and have a brown color.



# Generative Adversarial Network

- Idea
  - **Generator** to convert a vector  $z$  (sampled from  $P_z$ ) into fake data  $x$  (from  $P_G$ ), while we need  $P_G = P_{\text{data}}$
  - **Discriminator** classifies data as real or fake (1/0)
  - How? Impose an **adversarial loss** on the observed data distribution!

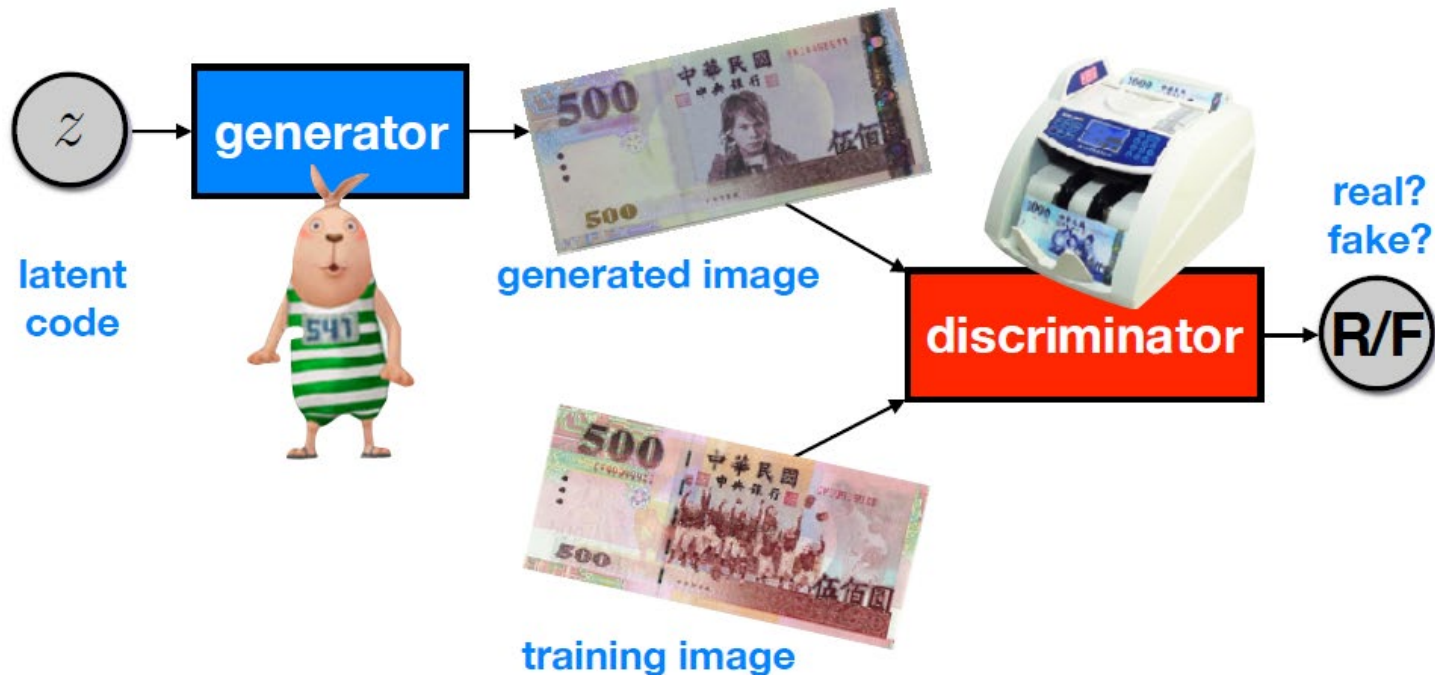


# Generative Adversarial Network (cont'd)

- Key idea:
  - Impose *adversarial loss* on data distribution
  - Let's see a practical example...

generator: try to generate more realistic images to cheat discriminator

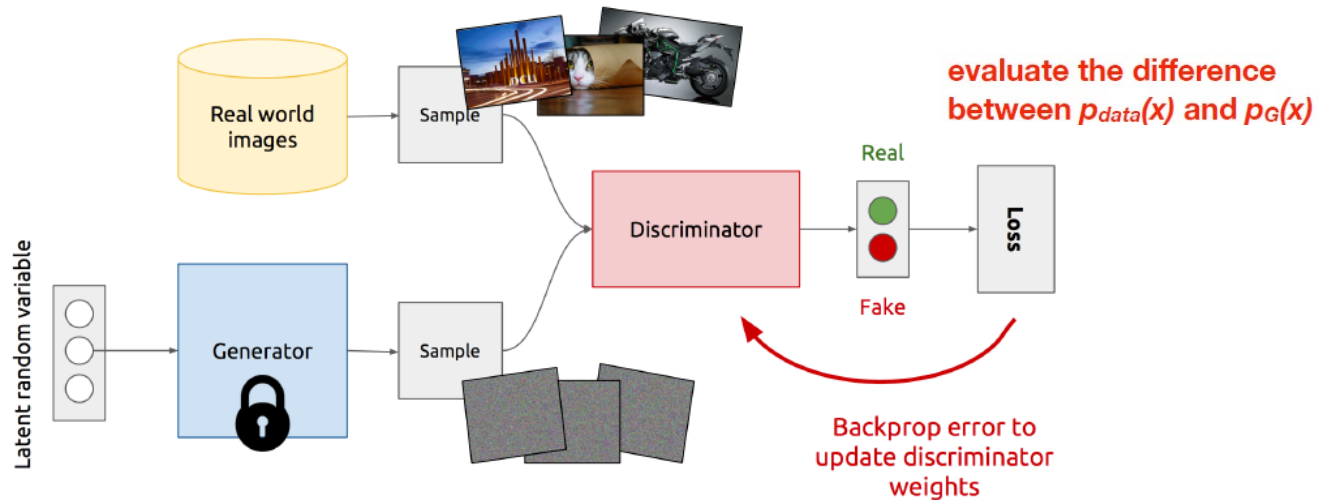
discriminator: try to distinguish whether the image is generated or real



# GAN (cont'd)

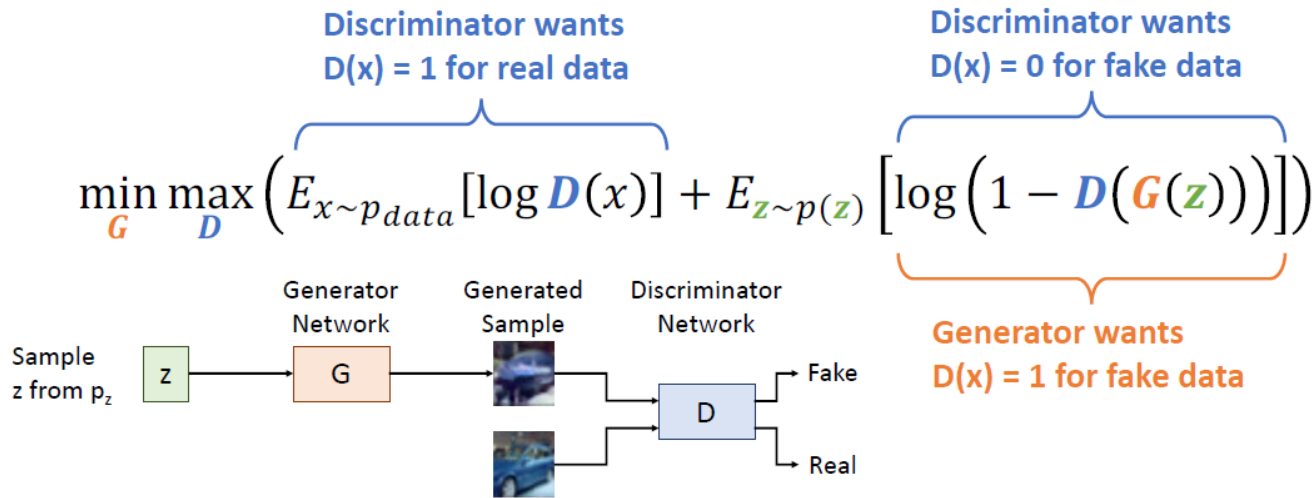
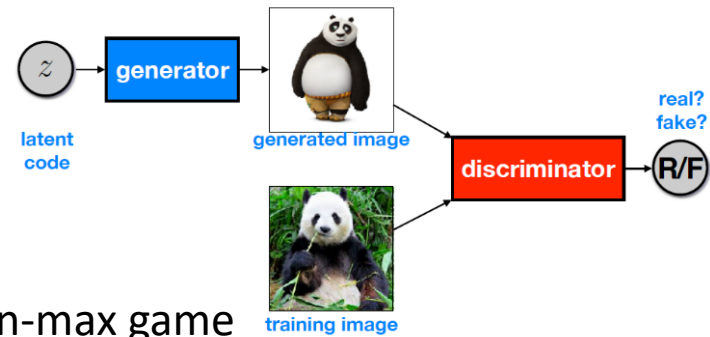
- Remarks
  - A function maps **normal distribution**  $N(\mathbf{0}, I)$  to  $P_{data}$
  - How good we are in mapping  $P_g$  to  $P_{data}$ ?
    - Train & ask the discriminator!
  - Conduct a two-player min-max game (see next slide for more details)

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



# Training Objective of GAN

- Jointly train generator G and discriminator D with a min-max game



- Train G & D with alternating gradient updates

$$\min_G \max_D V(G, D)$$

For t in 1, ... T:

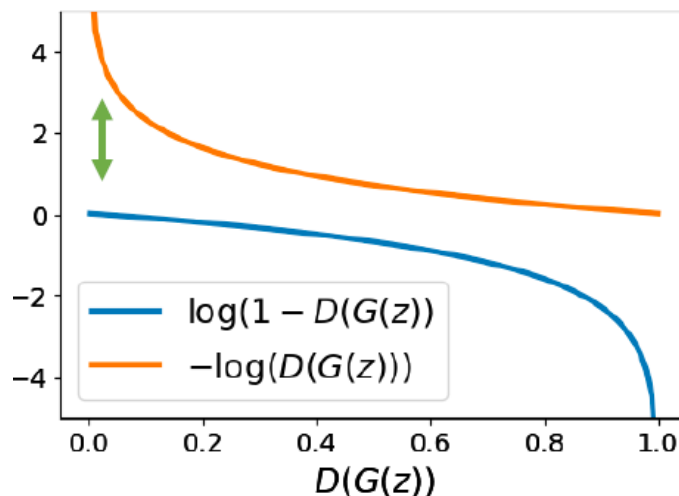
- (Update **D**)  $D = D + \alpha_D \frac{\partial V}{\partial D}$
- (Update **G**)  $G = G - \alpha_G \frac{\partial V}{\partial G}$

# Training Objective of GAN (*optional trick*)

- Potential Problem

$$\min_G \max_D \left( E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p(z)} \left[ \log \left( 1 - D(G(z)) \right) \right] \right)$$

- At start of training, G is not OK yet (obviously); D easily tells apart real/fake data (i.e.,  $D(G(z))$  close to 0).
- **Possible Solution:**
  - Instead of training G to minimize  $\log(1-D(G(z)))$  in the beginning, we train G to minimize  $-\log(D(G(z)))$ .
  - With strong gradients from G, we start the training of the above min-max game.



# Optimality of GAN

- Why the min-max game as objective a good idea?

$$\min_G \max_D \left( E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p(z)} [\log (1 - D(G(z)))] \right)$$

$$= \min_G \max_D \left( E_{x \sim p_{data}} [\log D(x)] + E_{x \sim p_G} [\log (1 - D(x))] \right)$$

$$= \min_G \int_X \max_D (p_{data}(x) \log D(x) + p_G(x) \log (1 - D(x))) dx$$

$$f(y) = a \log y + b \log(1 - y) \quad \left. \vphantom{f(y)} \right\} f'(y) = 0 \iff y = \frac{a}{a+b} \text{ (local max)}$$

$$f'(y) = \frac{a}{y} - \frac{b}{1-y} \quad \text{Optimal Discriminator: } D_G^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_G(x)}$$

# Optimality of GAN

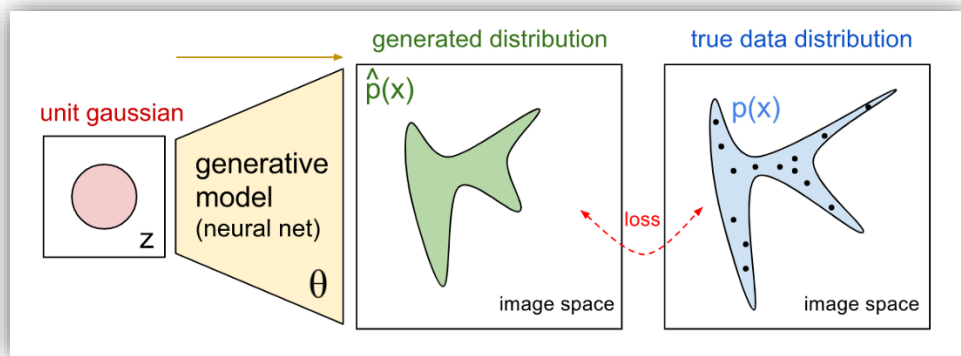
- Why the min-max game as objective a good idea? (cont'd)

$$\begin{aligned} & \min_G \max_D \left( E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p(z)} [\log (1 - D(G(z)))] \right) \\ \Rightarrow & \min_G \int_X \left( p_{data}(x) \log \frac{p_{data}(x)}{p_{data}(x) + p_G(x)} + p_G(x) \log \frac{p_G(x)}{p_{data}(x) + p_G(x)} \right) dx \\ & = \min_G \left( E_{x \sim p_{data}} \left[ \log \frac{2}{2} \frac{p_{data}(x)}{p_{data}(x) + p_G(x)} \right] + E_{x \sim p_G} \left[ \log \frac{2}{2} \frac{p_G(x)}{p_{data}(x) + p_G(x)} \right] \right) \\ & = \min_G \left( E_{x \sim p_{data}} \left[ \log \frac{2 * p_{data}(x)}{p_{data}(x) + p_G(x)} \right] + E_{x \sim p_G} \left[ \log \frac{2 * p_G(x)}{p_{data}(x) + p_G(x)} \right] - \log 4 \right) \end{aligned}$$

# Optimality of GAN

- Why the min-max game as objective a good idea? (cont'd)

$$\begin{aligned} & \min_G \max_D \left( E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p(z)} [\log (1 - D(G(z)))] \right) \\ &= \min_G \left( E_{x \sim p_{data}} \left[ \log \frac{2 * p_{data}(x)}{p_{data}(x) + p_G(x)} \right] + E_{x \sim p_G} \left[ \log \frac{2 * p_G(x)}{p_{data}(x) + p_G(x)} \right] - \log 4 \right) \\ &= \min_G \left( KL \left( p_{data}, \frac{p_{data} + p_G}{2} \right) + KL \left( p_G, \frac{p_{data} + p_G}{2} \right) - \log 4 \right) \end{aligned}$$



**Kullback-Leibler Divergence:**

$$KL(p, q) = E_{x \sim p} \left[ \log \frac{p(x)}{q(x)} \right]$$



# Optimality of GAN

- Why the min-max game as objective a good idea? (cont'd)

$$\begin{aligned} & \min_G \max_D \left( E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p(z)} [\log (1 - D(G(z)))] \right) \\ &= \min_G \left( E_{x \sim p_{data}} \left[ \log \frac{2 * p_{data}(x)}{p_{data}(x) + p_G(x)} \right] + E_{x \sim p_G} \left[ \log \frac{2 * p_G(x)}{p_{data}(x) + p_G(x)} \right] - \log 4 \right) \\ &= \min_G \left( KL \left( p_{data}, \frac{p_{data} + p_G}{2} \right) + KL \left( p_G, \frac{p_{data} + p_G}{2} \right) - \log 4 \right) \\ &= \min_G (2 * JSD(p_{data}, p_G) - \log 4) \end{aligned}$$

JSD is always nonnegative, and zero only when the two distributions are equal!  
Thus  $p_{data} = p_G$  is the global min, QED

**Jensen-Shannon Divergence:**

$$JSD(p, q) = \frac{1}{2} KL \left( p, \frac{p + q}{2} \right) + \frac{1}{2} KL \left( q, \frac{p + q}{2} \right)$$

## Remarks on Optimality of GAN

$$\min_G \max_D \left( E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p(z)} [\log (1 - D(G(z)))] \right)$$
$$= \min_G (2 * JSD(p_{data}, p_G) - \log 4)$$

- Summary

- The global min of the minmax game happens when

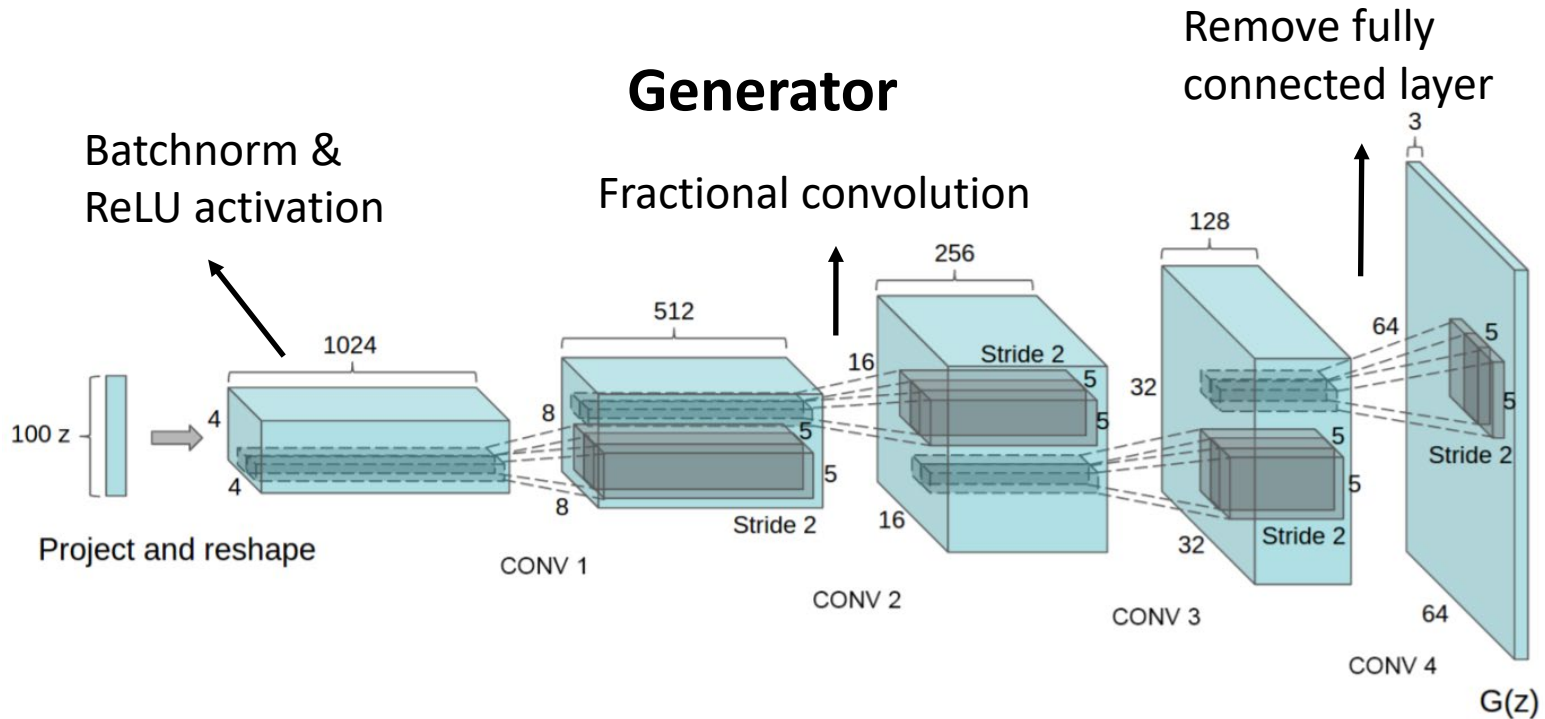
1.  $D_G^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_G(x)}$  (Optimal discriminator for any G) ➡
2.  $p_G(x) = p_{data}(x)$  (Optimal generator for optimal D) ➡

- **Caution!**

- G and D are learned models (i.e., DNNs) with fixed architectures.  
We don't know whether we can actually represent the **optimal D & G**.
- Optimality of GAN does not tell anything about **convergence** to the optimal D/G.

# Deep Convolutional GAN (DC-GAN)

- Remarks
  - ICLR 2016
  - A CNN+GAN architecture
  - Empirically make training of GAN more stable



# Deep Convolutional GAN (DC-GAN)

- Example Results



Collected face dataset

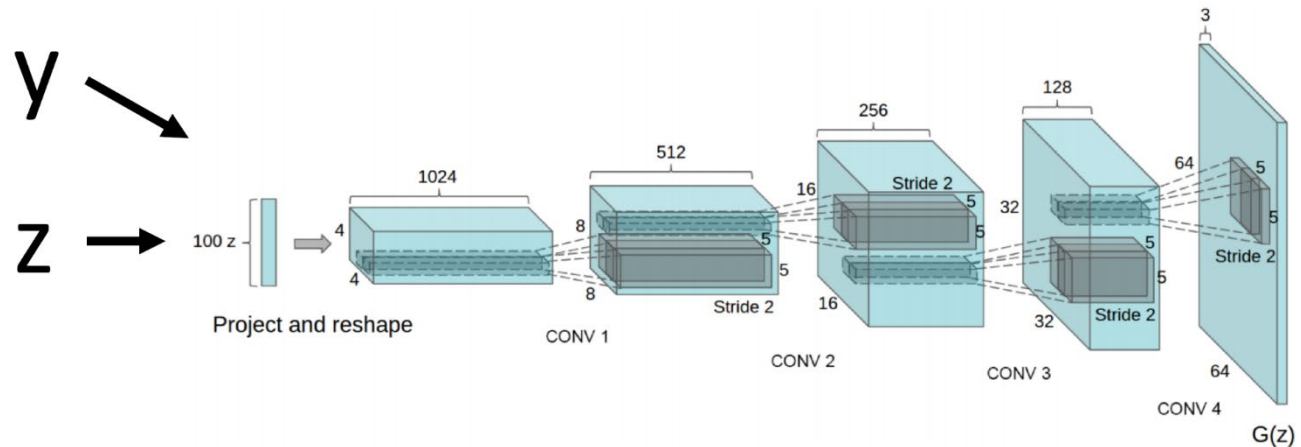


LSUN dataset

# Conditional GANs

- Remarks

- ICLR 2016
- Conditional generative model  $p(x|y)$  instead of  $p(x)$
- Both G and D take the label  $y$  as an additional input... i.e., a **conditional discriminator** is deployed...Why?
- **Why not just use D as designed in the standard GAN?**





# Conditional GANs

- Example Results

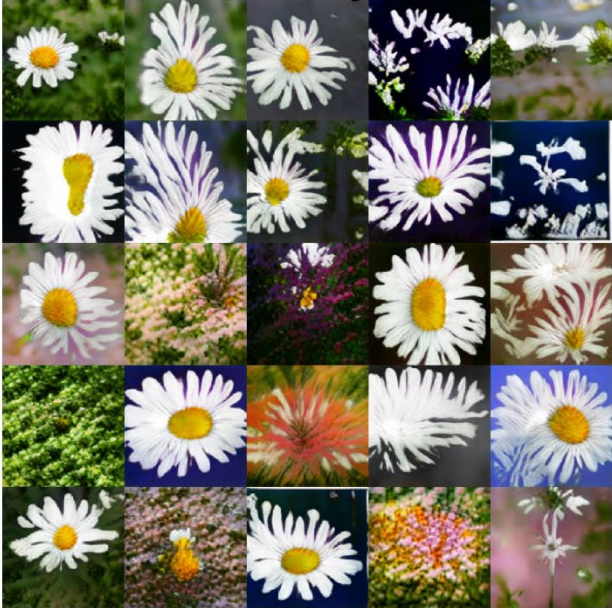
Welsh springer spaniel



Fire truck

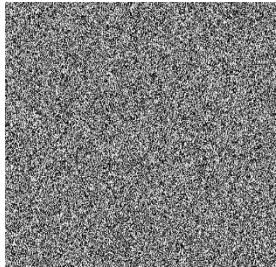


Daisy

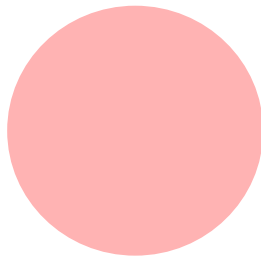


# Problems in Training GANs: Vanishing Gradients

- What Might Go Wrong?
  - GAN training is often unstable.
  - In other words, training might not converge properly.
  - The discriminator which we prefer is...



0



0.2



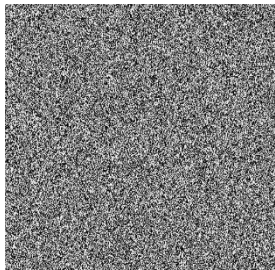
0.6



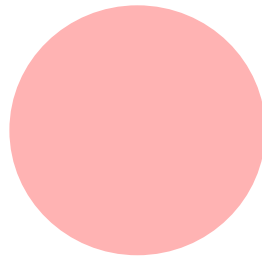
1

# Problems in Training GANs: Vanishing Gradients (cont'd)

- What Might Go Wrong?
  - GAN training is often unstable.
  - In other words, training might not converge properly.
  - The discriminator we trained might be as follows.  
In other words, no gradient to guide the generator to output proper images.



0



0



0



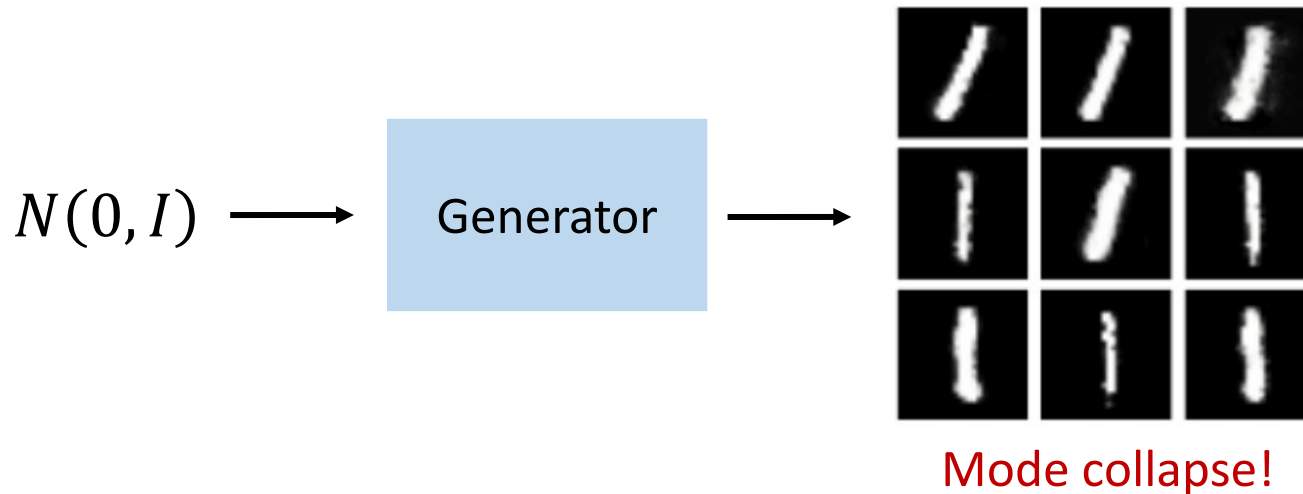
1

- This is known as the problem of *vanishing gradients*.



# Problems in Training GANs: Mode Collapse

- Remarks
  - The generator only outputs a limited number of image variants regardless of the inputs.



# Problems in Training GANs: Mode Collapse (cont'd)

- Remarks
  - The generator only outputs a limited number of image variants regardless of the inputs.

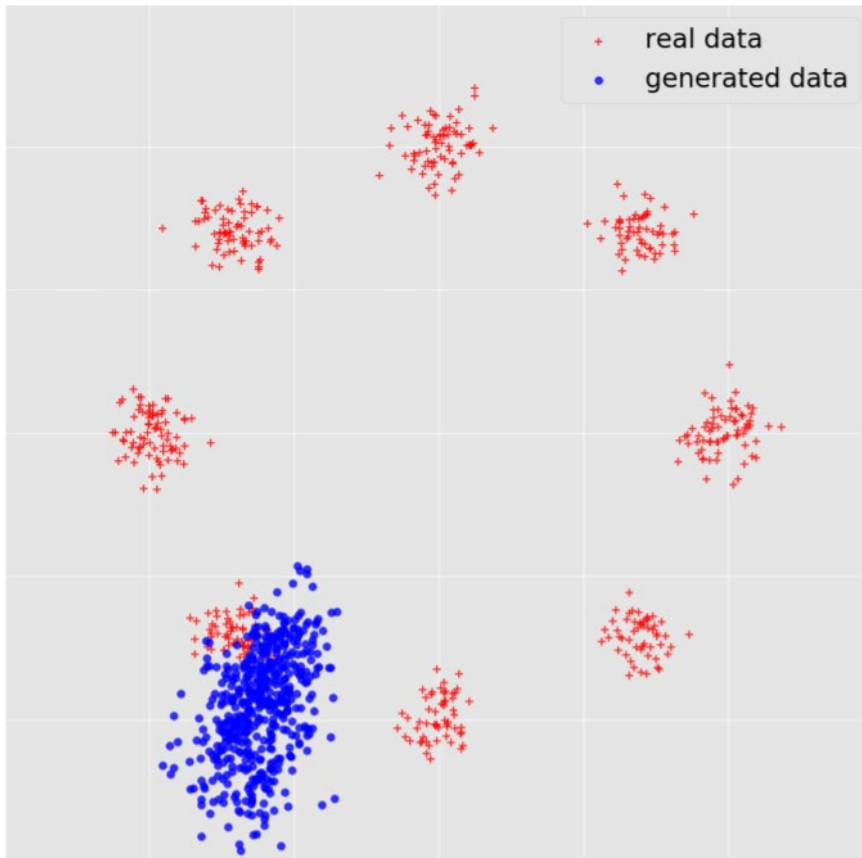
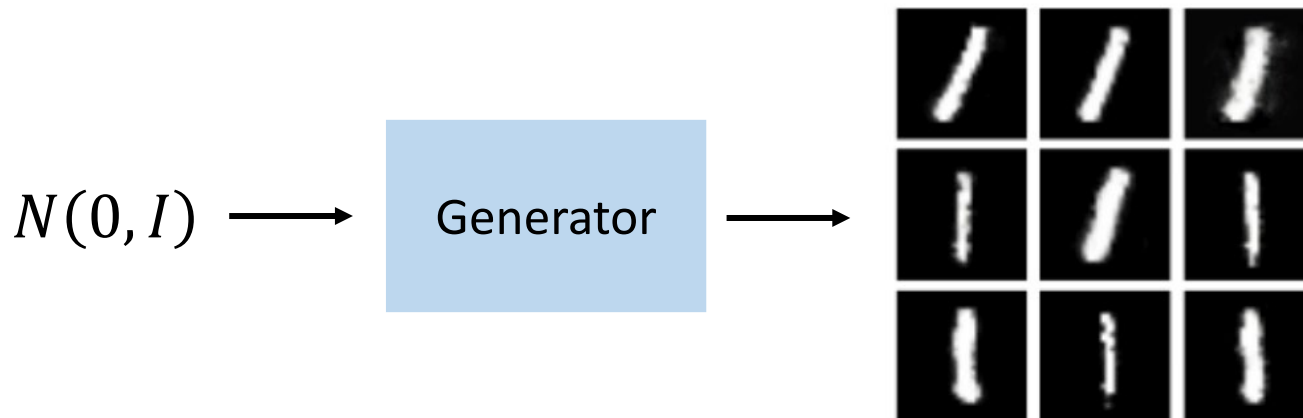


Photo credit:  
<https://openreview.net/pdf?id=rkmu5b0a->

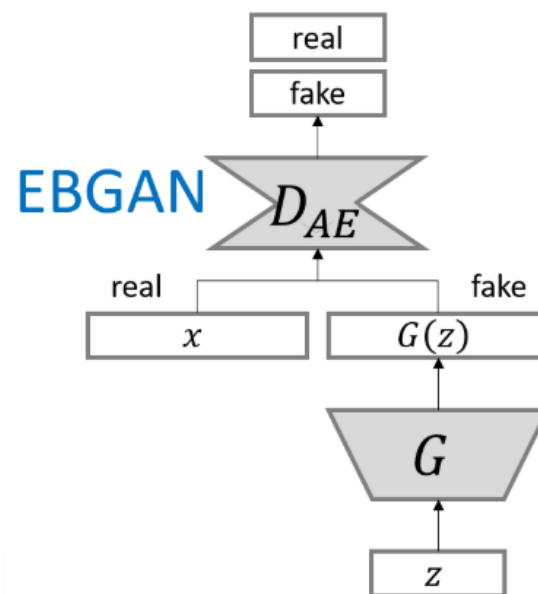
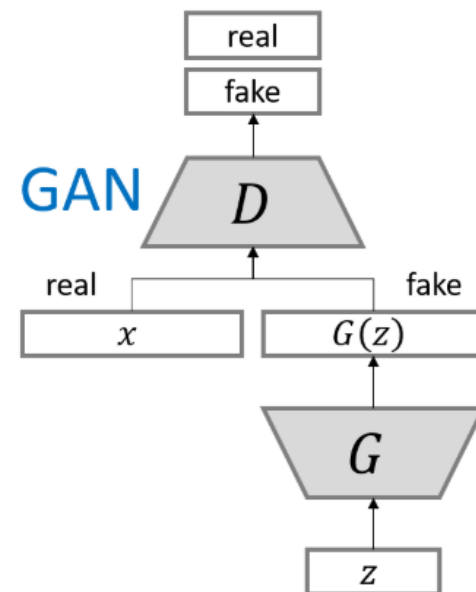
# Problems in Training GANs: Mode Collapse (cont'd)

- Why Mode Collapse Happens?
  - The objective of GANs assesses the **image authenticity**, not **diversity**.
  - Imbalance training between generator/discriminator (exploding/vanishing gradients)



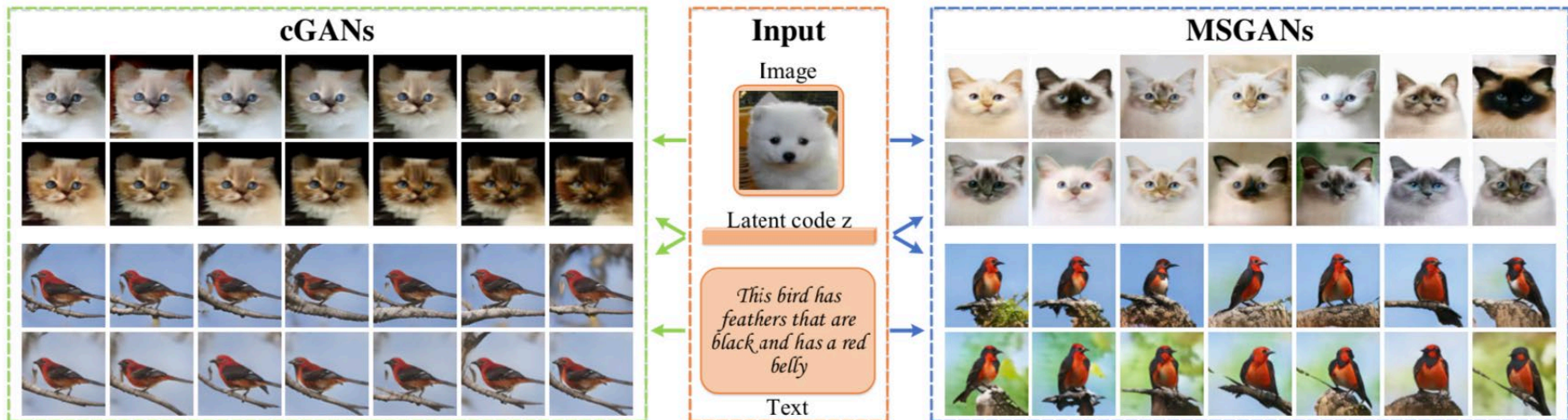
# Energy-Based GAN

- Energy Function
  - Converting input data into scalar outputs, viewed as energy values
  - Desired configuration is expected to output low energy values & vice versa.
- Energy Function as Discriminator
  - Use of autoencoder; can be pre-trained!
  - Reconstruction loss outputs a range of values instead of binary logistic loss.
  - Empirically better convergence



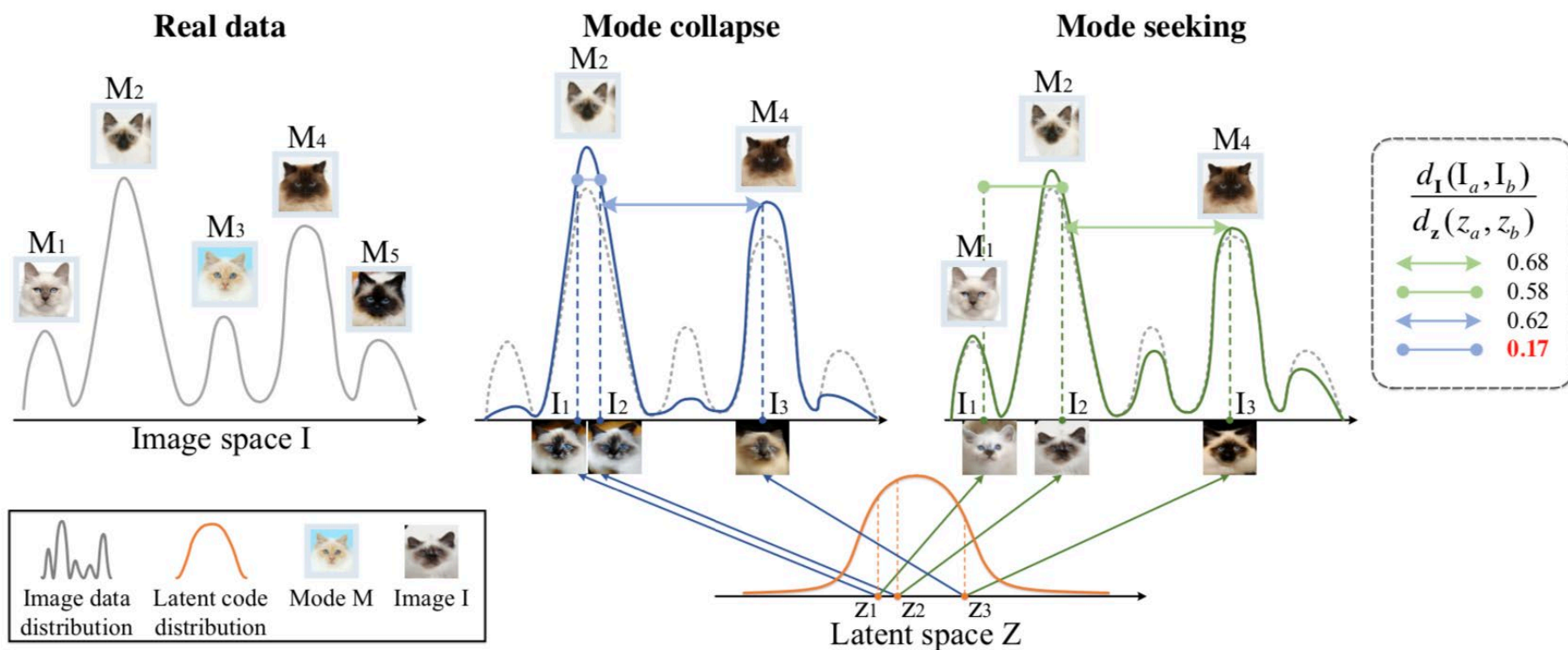
# MSGAN

- To address the **mode collapse** issue by **conditional GANs**
- Mode Seeking Generative Adversarial Networks for Diverse Image Synthesis
- With the goal of producing **diverse** image outputs.



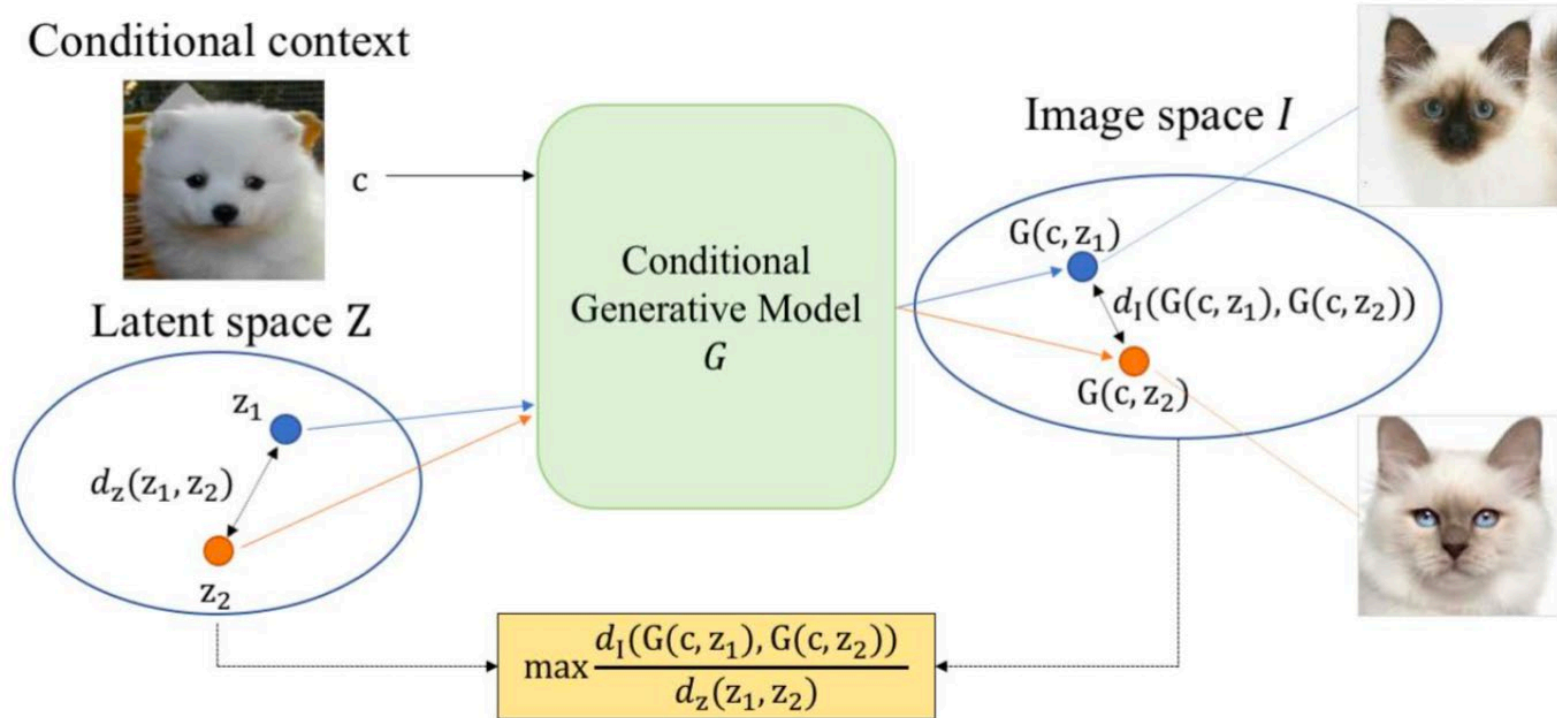
# MSGAN (cont'd)

- Motivation (for unconditional GAN)



# MSGAN (cont'd)

- Proposed Regularization (for conditional GAN)





# MSGAN (cont'd)

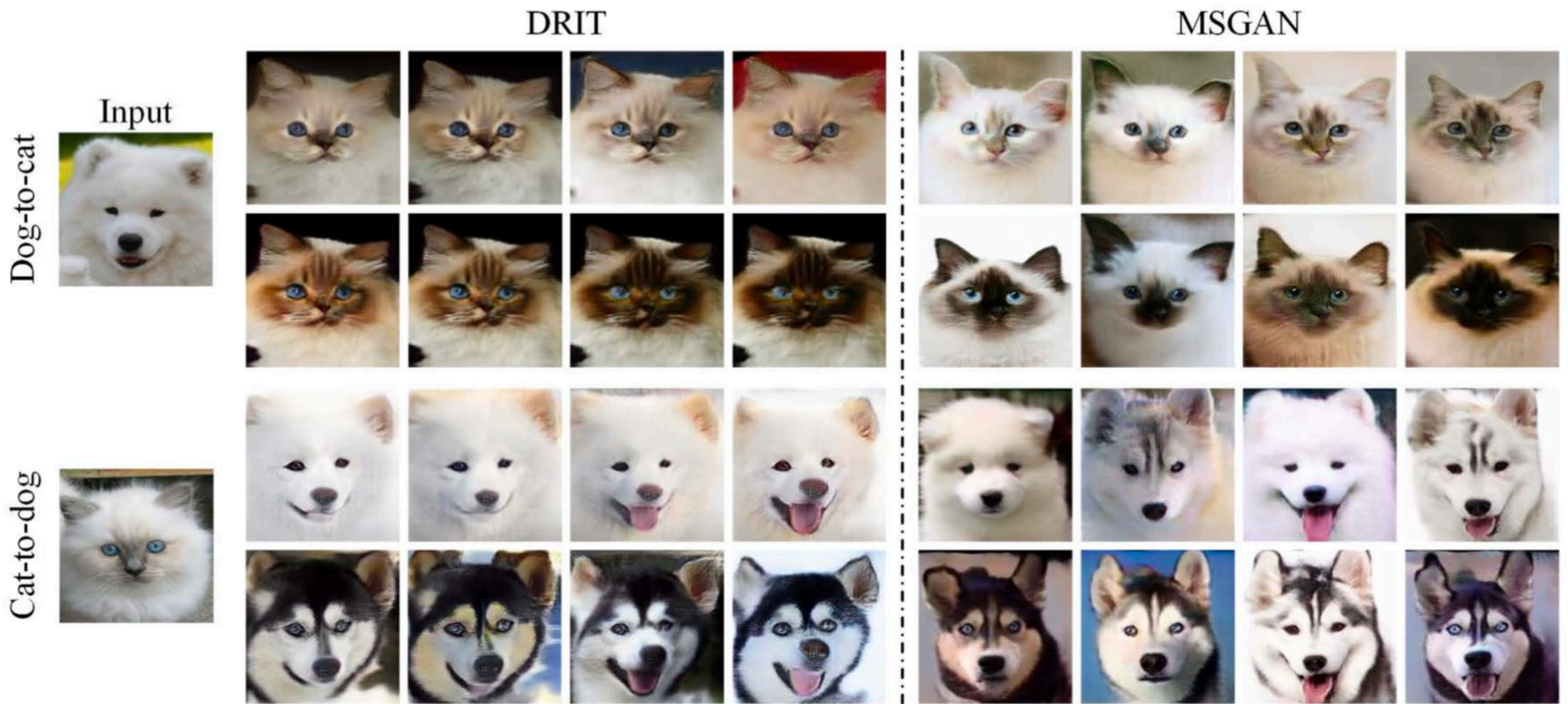
- Qualitative results
  - Conditioned on [paired images](#)





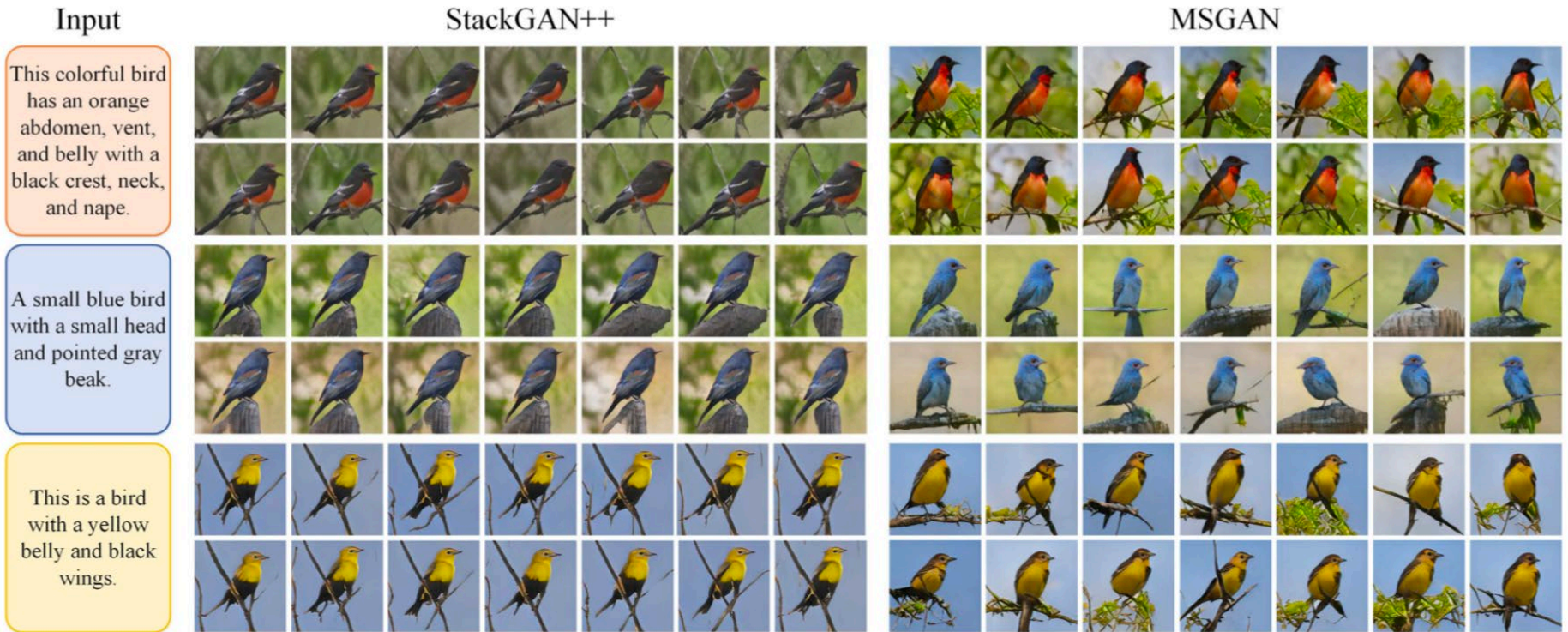
# MSGAN (cont'd)

- Qualitative results
  - Conditioned on **unpaired images**



# MSGAN (cont'd)

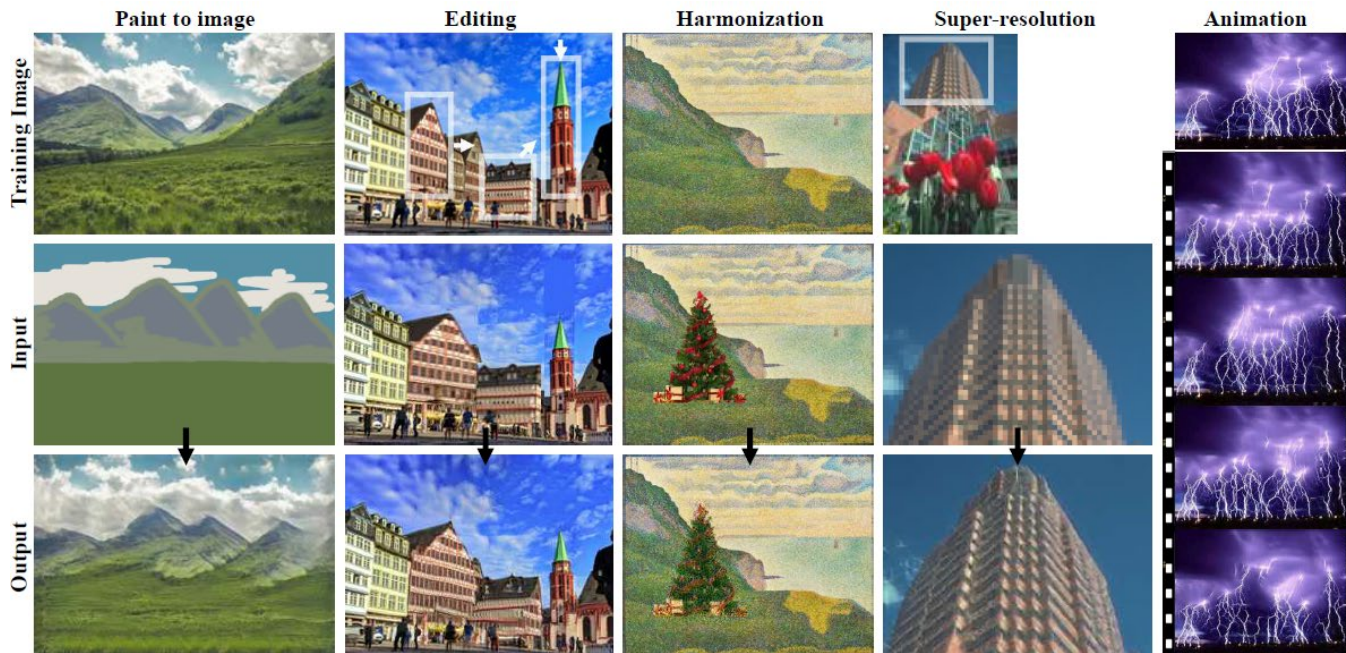
- Qualitative results
  - Conditioned on **text** (will talk about Vision & Language later this semester)





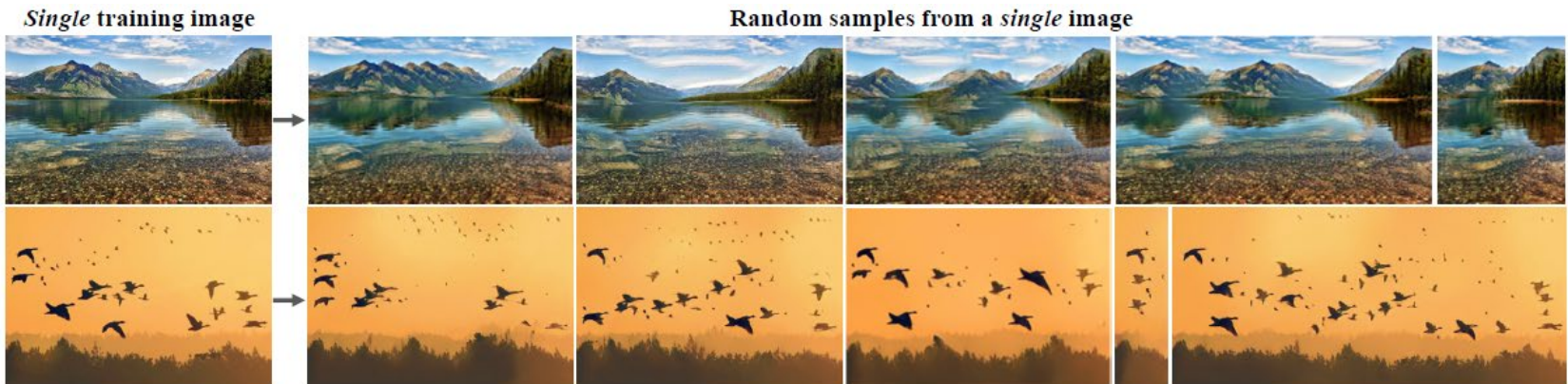
# SinGAN (if time permits): Learning a Generative Model from a Single Natural Image

- ICCV 2019 Best Paper Award
- Remarks:
  - Learning from a **single image**
  - Handle **multiple image manipulation tasks**
    - Super-resolution, style conversion, harmonization, image editing, et.



# SinGAN: Learning a Generative Model from a Single Natural Image

- Goal
  - Output images with arbitrary sizes and aspect ratios (via fully conv models) by changing dimensions of noise and the input size

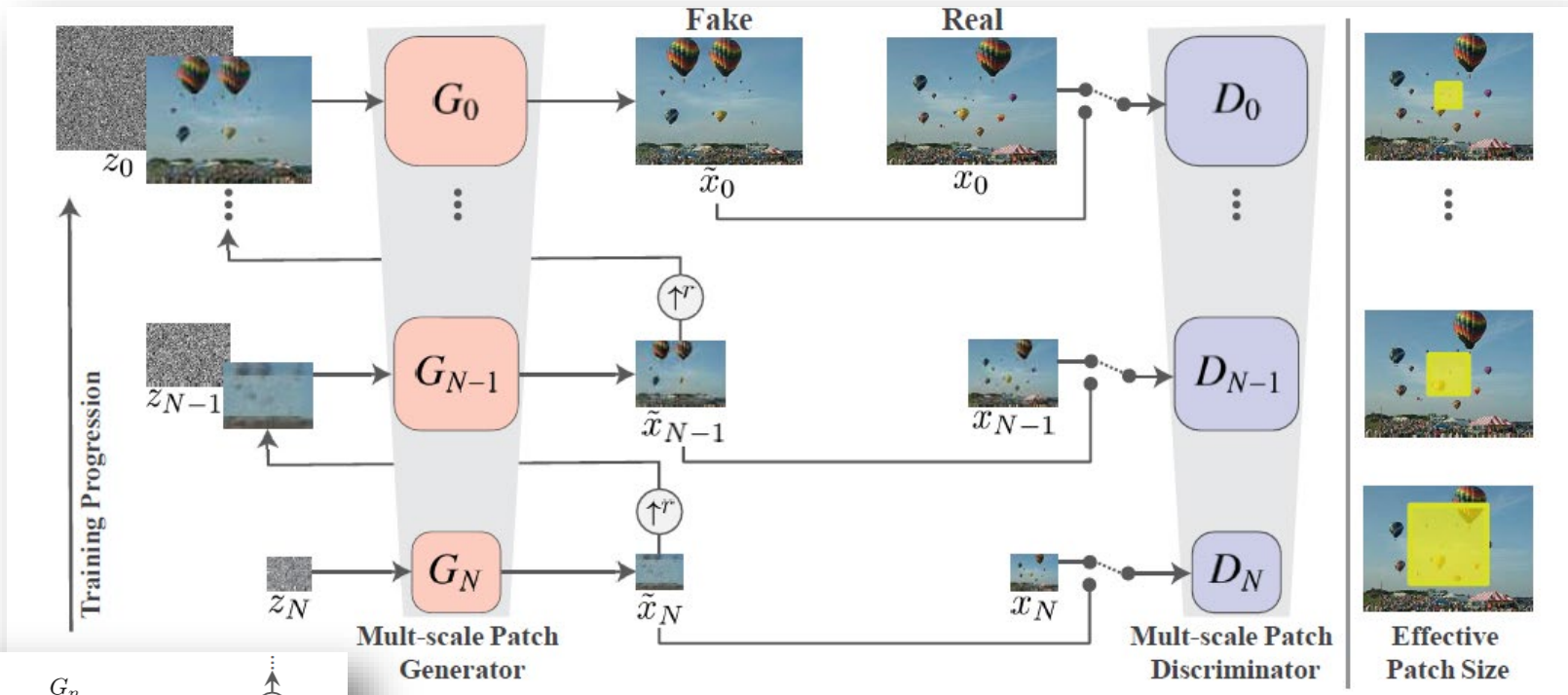


# SinGAN: Learning a Generative Model from a Single Natural Image

- Framework

$$\min_{G_n} \max_{D_n} \mathcal{L}_{\text{adv}}(G_n, D_n) + \alpha \mathcal{L}_{\text{rec}}(G_n)$$

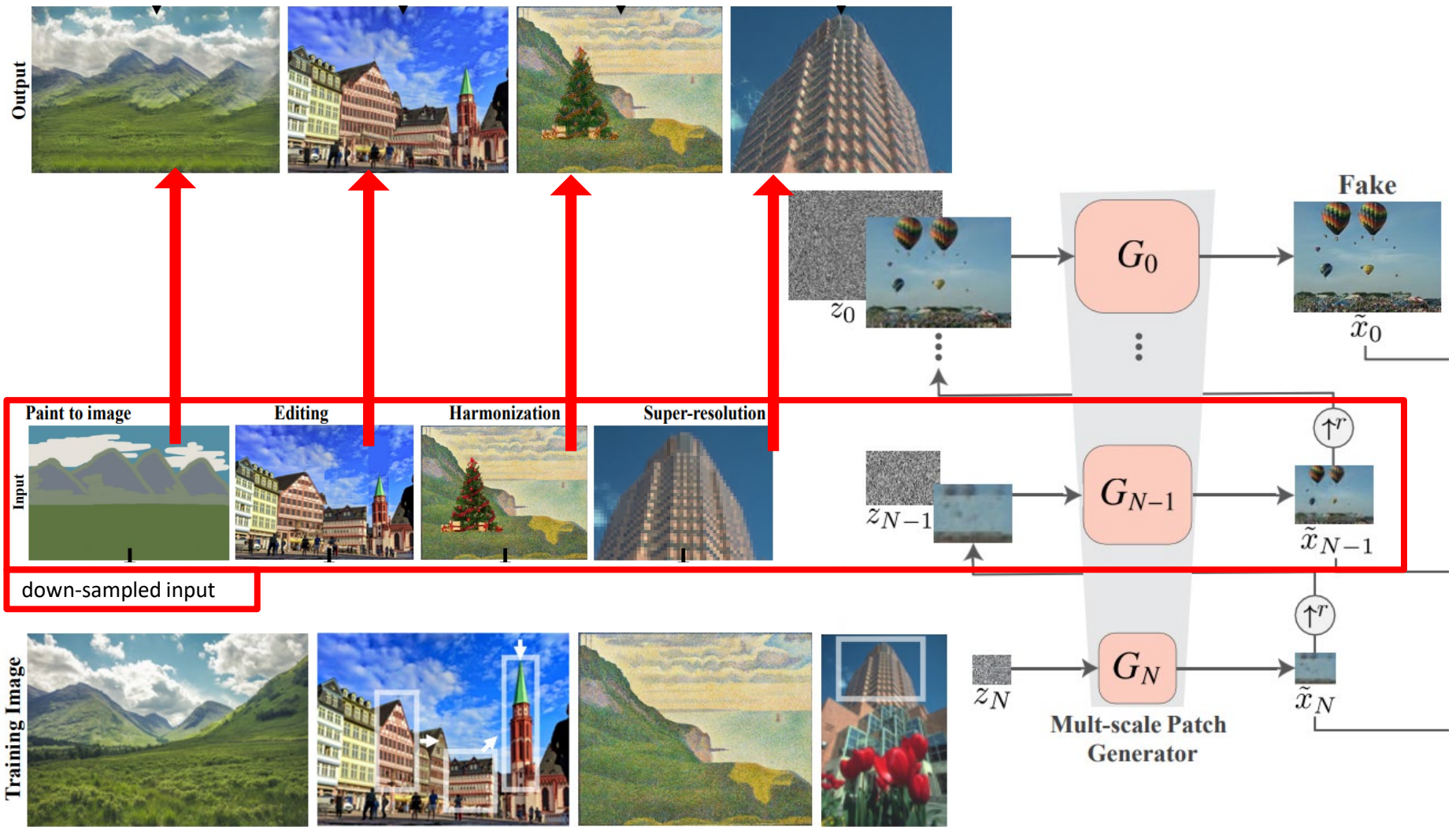
fix kernel (receptive field) size at each scale:  
capture structures of decreasing size as we go up



add **noise** before **Conv**:  
ensure that GAN does not  
disregard the noise

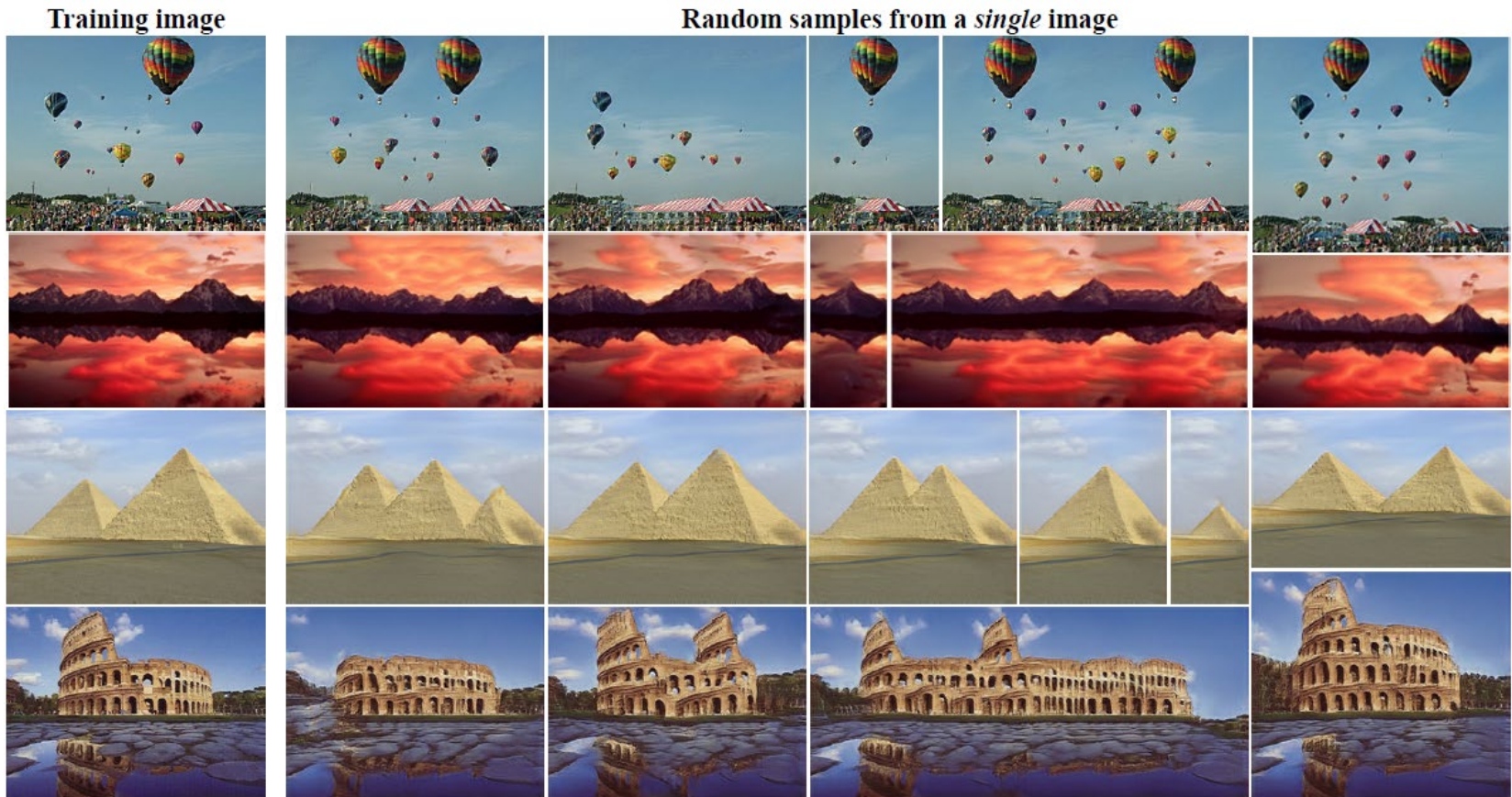


# Inference Stage for SinGAN



# SinGAN: Learning a Generative Model from a Single Natural Image

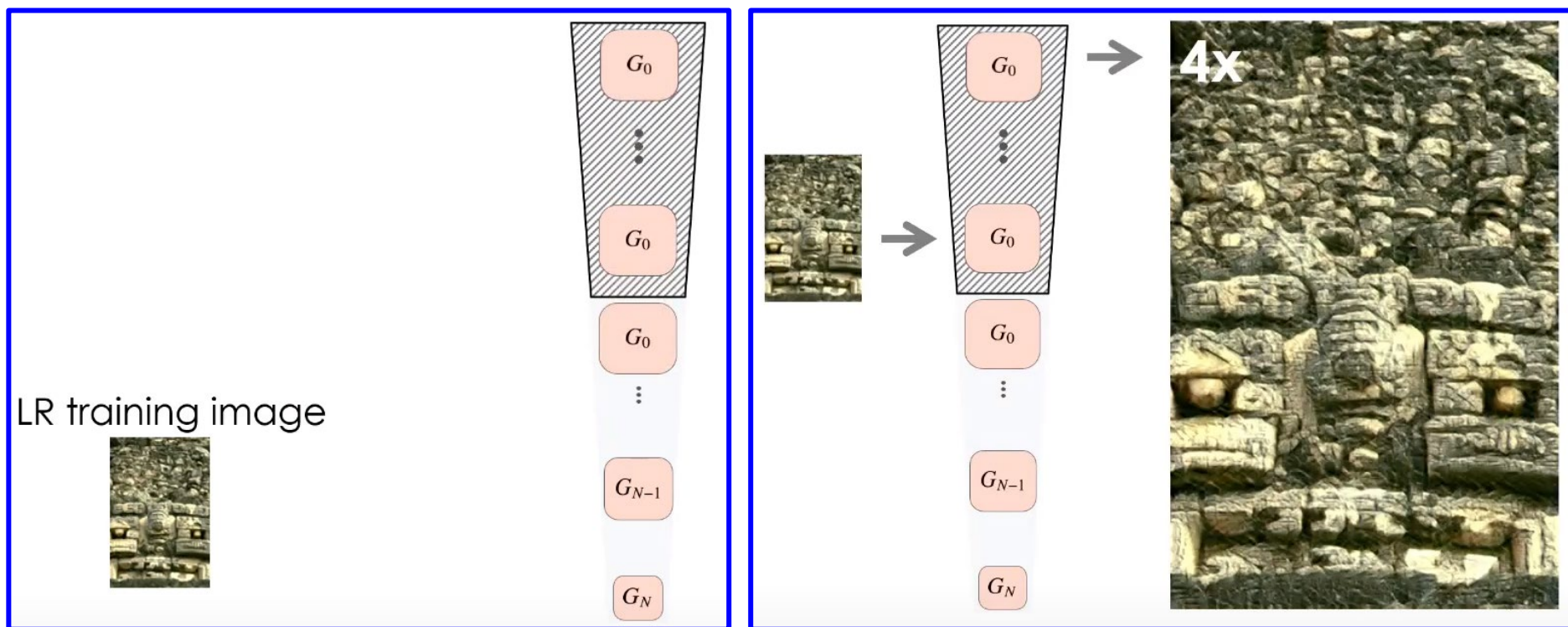
- Random image generation





# SinGAN: Learning a Generative Model from a Single Natural Image

- Super-Resolution





# SinGAN: Learning a Generative Model from a Single Natural Image

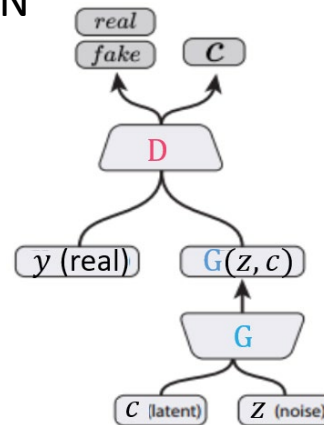
- Editing

Image	Injection scale	Total number of scales
Rock1	$n = 5$	$N = 7$
Rock2	$n = 5$	$N = 7$
Rock3 (also Fig. 12, main text)	$n = 5$	$N = 7$
Tree	$n = 7$	$N = 9$
Mountain	$n = 4$	$N = 8$
Red cliff	$n = 5$	$N = 9$
Hay	$n = 6$	$N = 9$



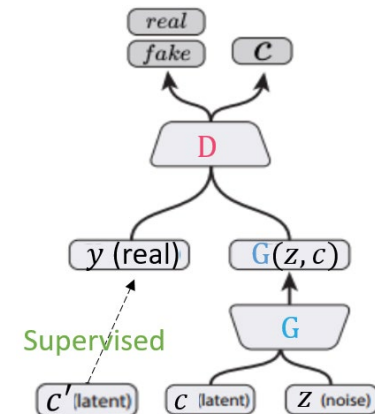
# Representation Disentanglement: Conditional GAN

- Goal
  - **Interpretable** deep feature representation
  - Disentangle attribute of interest  $c$  from the derived latent representation  $z$ 
    - Unsupervised: InfoGAN
    - Supervised: AC-GAN



InfoGAN

Chen et al.  
NIPS '16

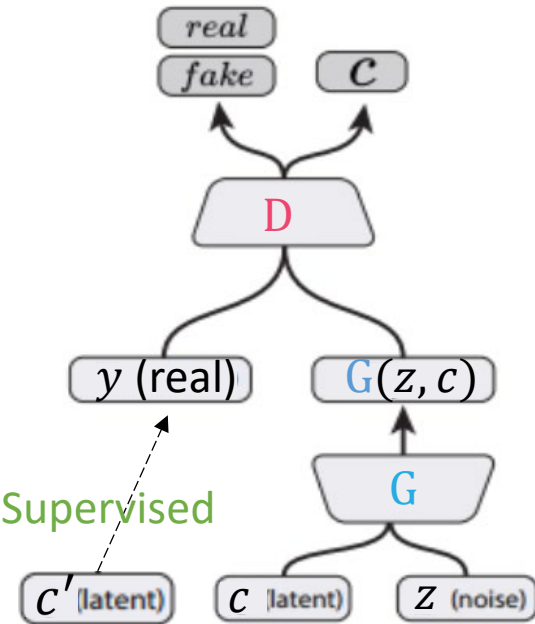


ACGAN

Odena et al.  
ICML '17

# AC-GAN

- Supervised Disentanglement



- Learning**

- Overall objective function

$$G^* = \arg \min_G \max_D \mathcal{L}_{GAN}(G, D) + \mathcal{L}_{cls}(G, D)$$

- Adversarial Loss

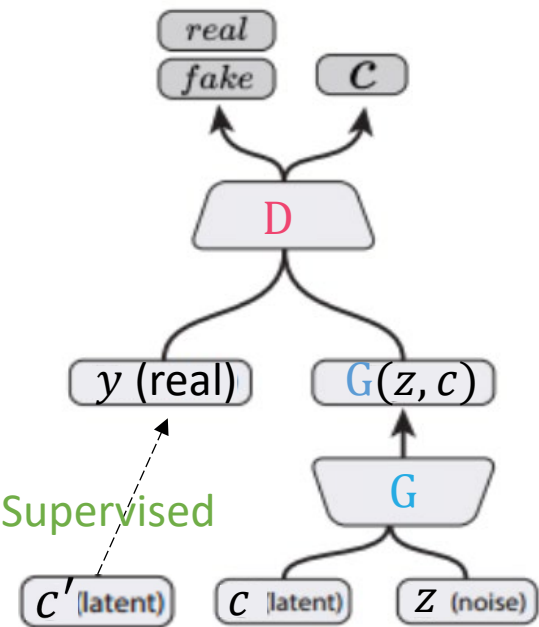
$$\mathcal{L}_{GAN}(G, D) = \mathbb{E}[\log(1 - D(G(z, c)))] + \mathbb{E}[\log D(y)]$$

- Disentanglement loss

$$\mathcal{L}_{cls}(G, D) = \underbrace{\mathbb{E}[-\log D_{cls}(c'|y)]}_{\text{Real data w.r.t. its domain label}} + \underbrace{\mathbb{E}[-\log D_{cls}(c|G(x, c))]}_{\text{Generated data w.r.t. assigned label}}$$

# AC-GAN

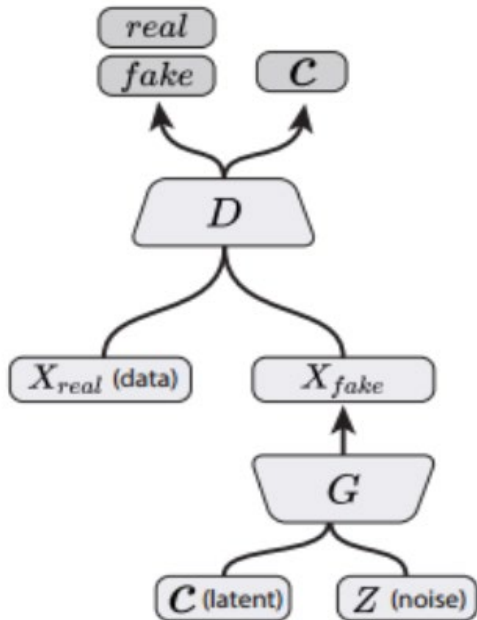
- Supervised Disentanglement



Different  $c$  values

# InfoGAN

- Unsupervised Disentanglement



- Learning

- Overall objective function

$$G^* = \arg \min_G \max_D \mathcal{L}_{GAN}(G, D) + \mathcal{L}_{cls}(G, D)$$

- Adversarial Loss

$$\mathcal{L}_{GAN}(G, D) = \mathbb{E}[\log(1 - D(G(z, c)))] + \mathbb{E}[\log D(y)]$$

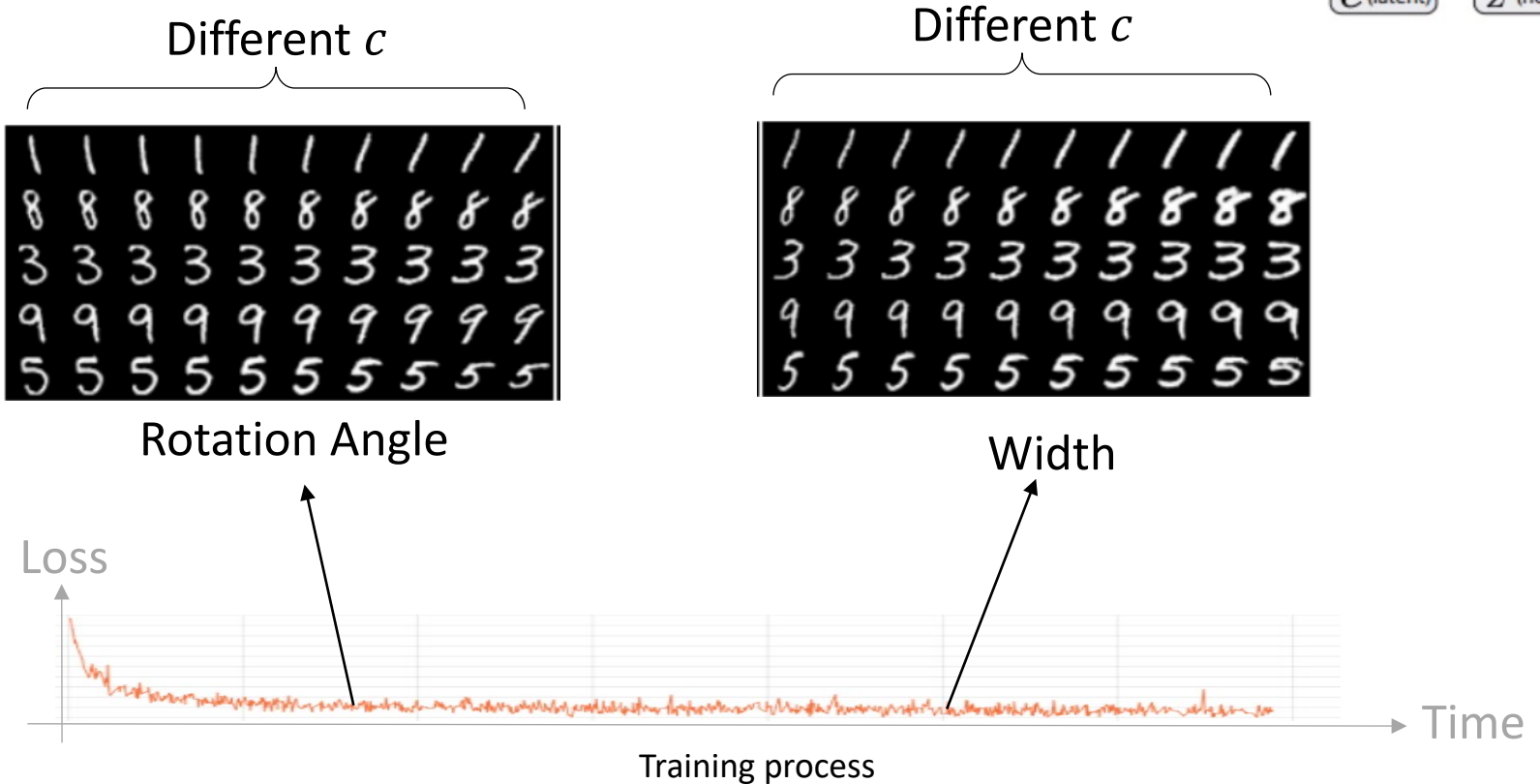
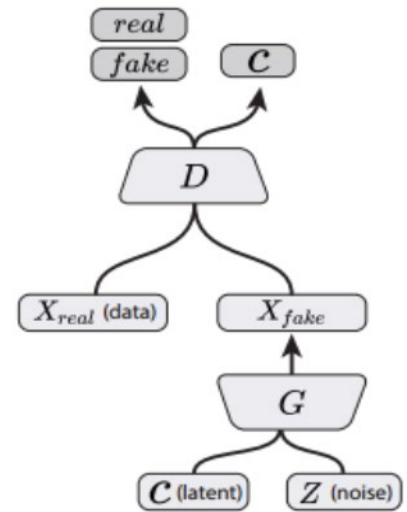
- Disentanglement loss

$$\mathcal{L}_{cls}(G, D) = \mathbb{E}[-\log D_{cls}(c | G(x, c))]$$

Generated data  
w.r.t. assigned label

# InfoGAN

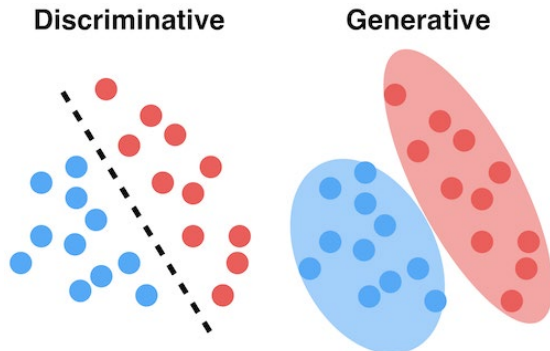
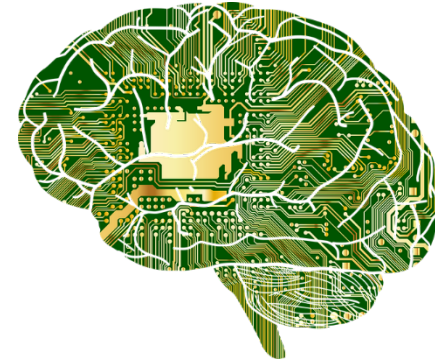
- Unsupervised Disentanglement
  - No guarantee in disentangling particular semantics
  - It can be viewed as...



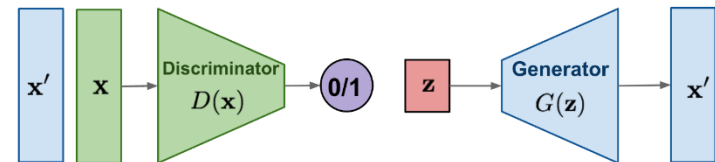


# What We've Covered Today...

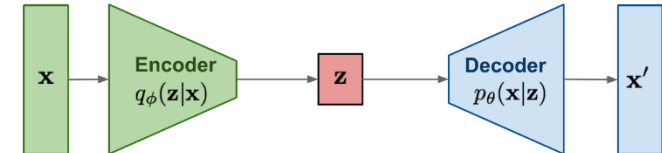
- Generative Models
  - Diffusion Model
    - Conditional Diffusion Model
      - Classifier Guidance
      - Classifier-Free Guidance
      - Text/Image Guidance
    - Personalization via Diffusion Model
  - Generative Adversarial Network
  - **HW #2 is out! (due 10/29)**



GAN: Adversarial training



VAE: maximize variational lower bound



Diffusion models:  
Gradually add Gaussian noise and then reverse

