

# Advanced Computer Vision

Tanveer Hussain

[htanveer3797@gmail.com](mailto:htanveer3797@gmail.com)



جامعة الملك عبد الله  
لعلوم والتكنولوجيا  
King Abdullah University of  
Science and Technology

KAUST Academy  
King Abdullah University of Science and Technology

Course designed by **Naeem Ullah Khan** ([naeemullah.khan@kaust.edu.sa](mailto:naeemullah.khan@kaust.edu.sa)), updated by Tanveer Hussain

# Generative Adversarial Networks (GANs)

- ▶ Variational Autoencoders are based on maximizing likelihood or approximations

$$p_{\theta}(x) = \underset{z}{\int} p_{\theta}(x|z)p_{\theta}(z)$$

- ▶ What if we give up on explicitly modeling density, and just want ability to sample?
- ▶ **GANs**: don't work with any explicit density function!  
Instead, take game-theoretic approach: learn to generate from training distribution through 2-player game

# Comparing Distributions via Samples

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.,  $P = Q$ ?)



vs.



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

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

# Comparing Distributions via Samples (cont.)

- ▶ Let's consider a test statistic  $T$  for this purpose.
- ▶ Test statistic  $T$  compares  $S_1$  and  $S_2$  e.g., difference in means, variances of the two sets of samples.
- ▶ **Key observation:** Test statistic is likelihood-free since it does not involve the densities  $P$  or  $Q$  (only samples)

- ▶ Let  $S_1 = D = \{x \sim p_{data}\}$  and  $S_2 = \{x \sim p_\theta\}$
- ▶ **Idea:** Train the generative model to minimize a two-sample test objective between  $S_1$  and  $S_2$  i.e., a generator.
- ▶ Okay, now how do we get a two-sample test objective?
- ▶ **Another Idea:** Train another neural network to discriminate between the two samples i.e., a discriminator.
- ▶ And voila, we have a GAN. All that's left is now the training method.

# GANs (cont.)

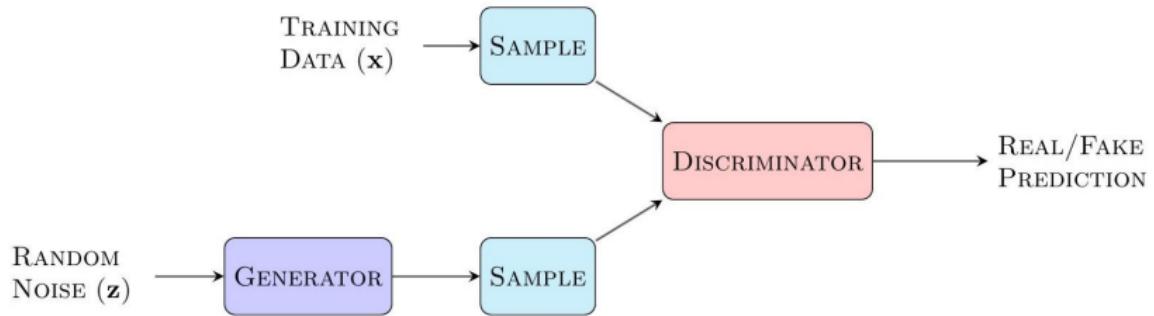
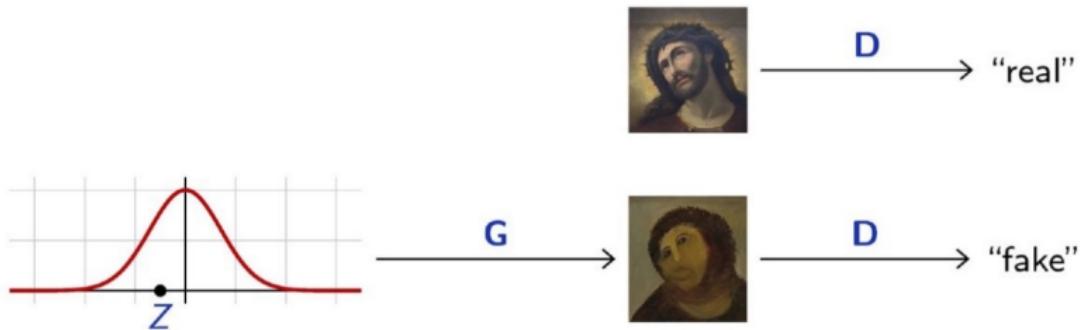


Figure 16: A GAN model diagram<sup>1</sup>

- ▶ Generative Adversarial Networks were introduced by Ian Goodfellow et al. (2014).
- ▶ The idea behind GANs is to train two networks jointly.
- ▶ A **discriminator D** to classify samples as “real” or “fake”,
- ▶ A **generator G** to map a [simple] fixed distribution to samples that fool **D**.
- ▶ The approach is **adversarial** since the two networks have antagonistic objectives.

# GANs (cont.)



**Figure 17:** Generator transforms a simple distribution to a sample from data. Discriminator discriminates between real and fake samples.

<sup>1</sup><https://github.com/JamesAllingham/LaTeX-TikZ-Diagrams>

# Training GANs

- ▶ Both Discriminator and Generator are trained jointly in a min-max game.
- ▶ Minimax objective function:

$$\min_{\theta_G} \max_{\theta_D} = E_{x \sim p_{data}} \log(D_{\theta_D}(x)) + E_{z \sim p(z)} \log(1 - D_{\theta_D}(G_{\theta_G}(z)))$$

Discriminator output for real data  $x$                               Discriminator output for generated fake data  $G(z)$

- ▶ Discriminator  $\theta_D$  wants to maximise objective such that  $D(x)$  is close to 1 (real) and  $D(G(z))$  is close to 0 (fake).
- ▶ Generator  $\theta_G$  wants to minimise objective such that  $D(G(z))$  is close to 1 (discriminator is fooled into thinking generated  $G(z)$  is real).

# Training GANs (cont.)

**for** number of training iterations **do**

**for**  $k$  steps **do**

- Sample minibatch of  $m$  noise samples  $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$  from noise prior  $p_g(\mathbf{z})$ .
- Sample minibatch of  $m$  examples  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  from data generating distribution  $p_{\text{data}}(\mathbf{x})$ .
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[ \log D_{\theta_d}(\mathbf{x}^{(i)}) + \log(1 - D_{\theta_d}(G_{\theta_g}(\mathbf{z}^{(i)}))) \right]$$

**end for**

- Sample minibatch of  $m$  noise samples  $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$  from noise prior  $p_g(\mathbf{z})$ .
- Update the generator by ascending its stochastic gradient (improved objective):

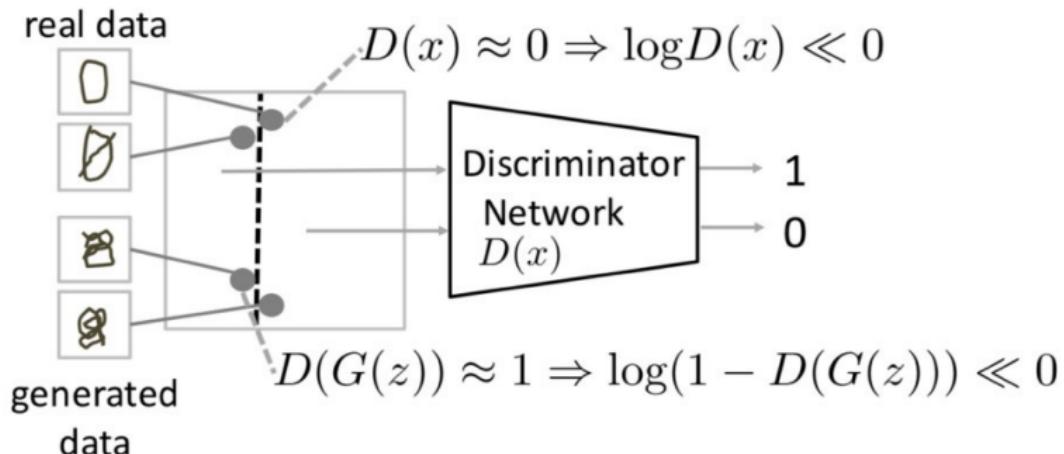
$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log(D_{\theta_d}(G_{\theta_g}(\mathbf{z}^{(i)})))$$

**end for**

# Understanding the Objective function

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

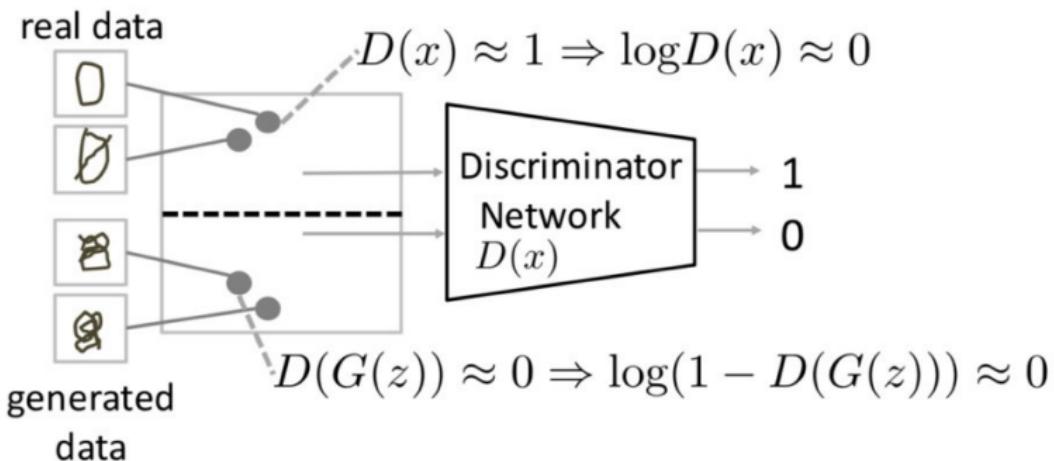
$D(x)$  should be 1       $D(G(z))$  should be 0



# Understanding the Objective function (cont.)

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

$D(x)$  should be 1       $D(G(z))$  should be 0

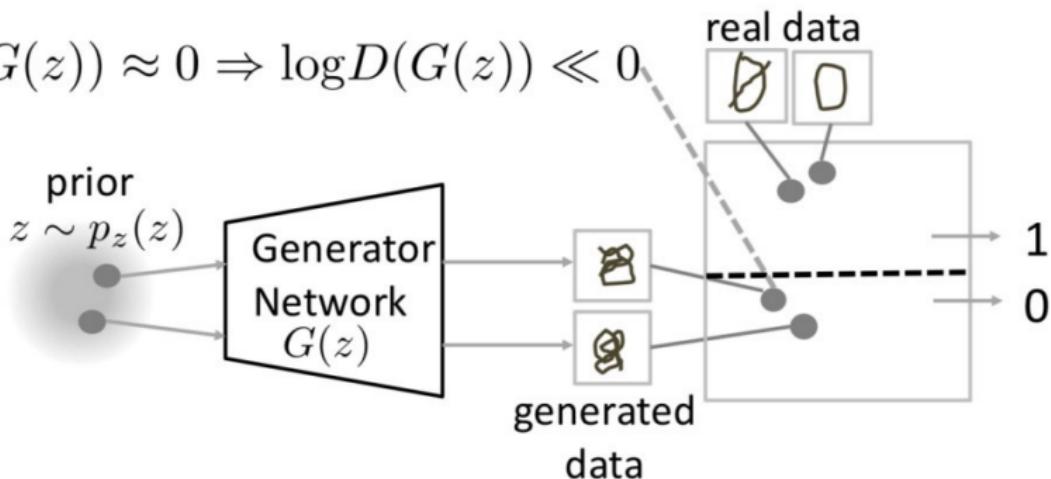


# Understanding the Objective function (cont.)

$$\max_G (\mathbb{E}_{z \sim p_z(z)} [\log(D(G(z)))] )$$

$D(G(z))$  should be 1

$$D(G(z)) \approx 0 \Rightarrow \log D(G(z)) \ll 0$$

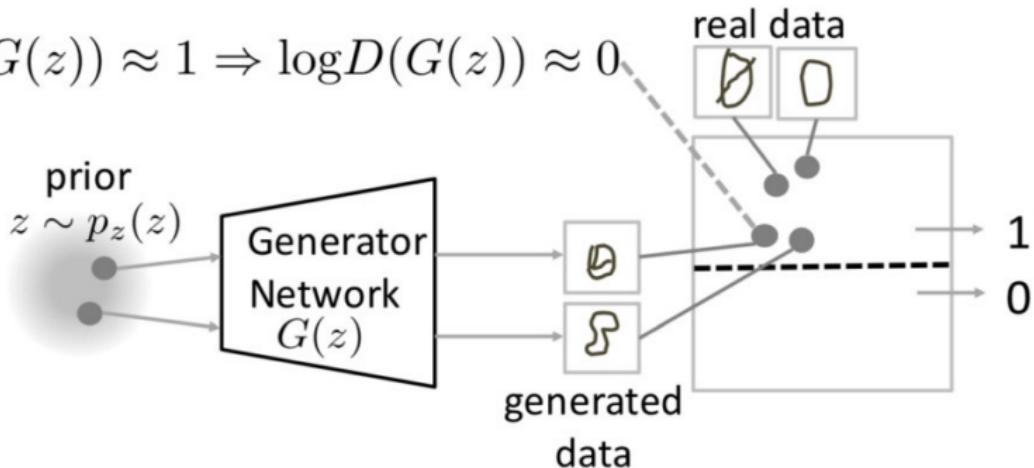


# Understanding the Objective function (cont.)

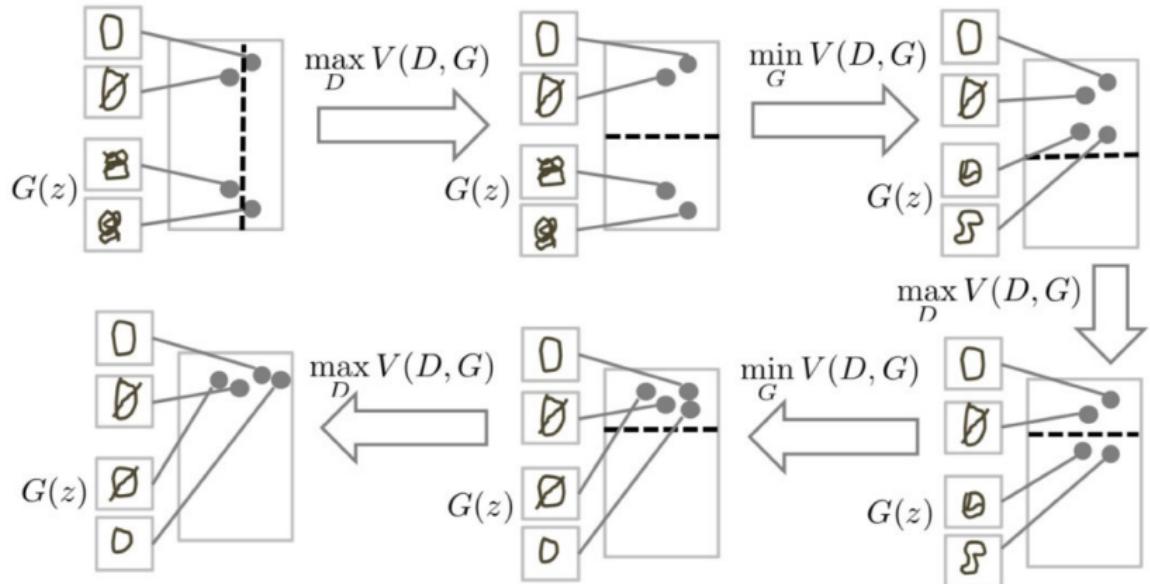
$$\max_G (\mathbb{E}_{z \sim p_z(z)} [\log(D(G(z)))])$$

$D(G(z))$  should be 1

$$D(G(z)) \approx 1 \Rightarrow \log D(G(z)) \approx 0$$



# Understanding the Objective function (cont.)



<sup>1</sup><https://www.slideshare.net/cmarkohchang/generative-adversarial-networks>

# GANs - Interactive Demo

<https://poloclub.github.io/ganlab/>

generated\_images



Figure 18: GAN generated samples for MNIST digits dataset

# GANs - Results (cont.)

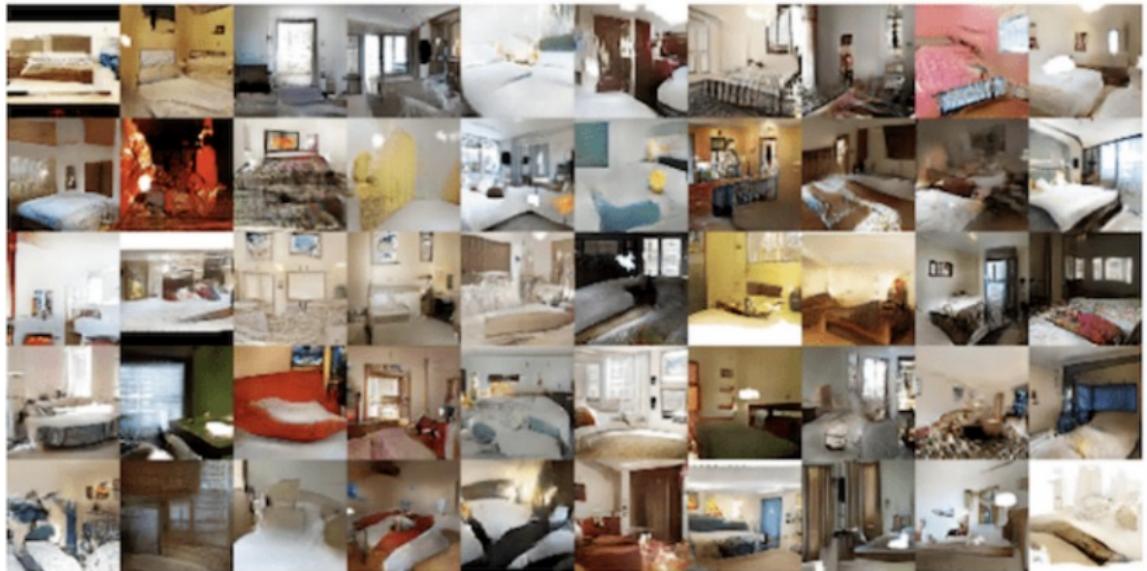


Figure 19: GAN generated samples for bedroom images

# GANs - Results (cont.)



Figure 20: GAN generated samples for anime character faces

## ► Vanishing Gradients

- When the discriminator is perfect, we are guaranteed with  $D(x) = 1, \forall x \in p_{data}$  and  $D(x) = 0, \forall x \in p_G$ .
- Loss function falls to zero and we end up with no gradient to update.
- **Solution:** Perform gradient ascent on generator i.e., different objective.

$$\max_{\theta_G} E_{x \sim p(z)} \log(D_{\theta_D}(G_{\theta_G}(z)))$$

# Problems with GANs (cont.)

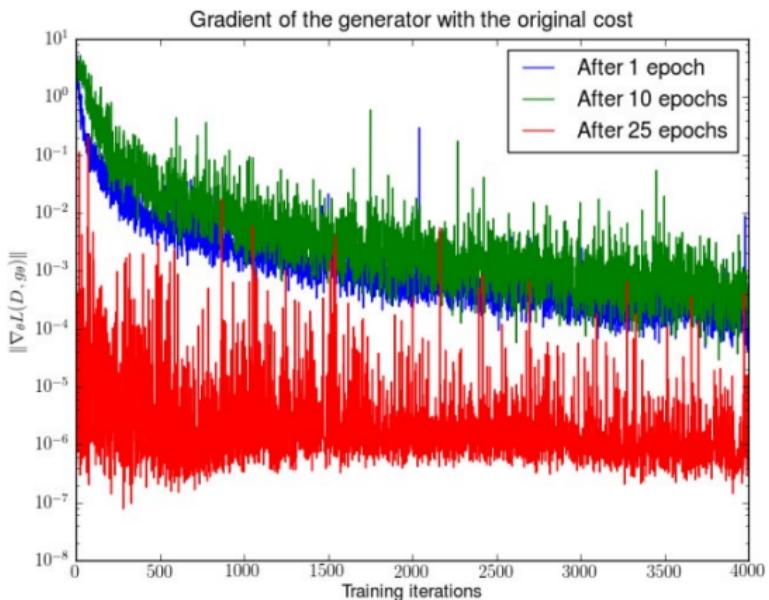


Figure 21: Vanishing gradient in GANs (Image source: [Arjovsky and Bottou, 2017](#))

# Problems with GANs (cont.)

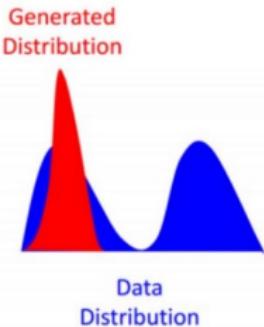
## ► Hard to achieve Nash equilibrium

- Two models are trained simultaneously to find a Nash equilibrium to a two player non-cooperative game. However, each model updates its cost independently with no respect to another player in the game. Updating the gradient of both models concurrently cannot guarantee a convergence
- How to Train a GAN? Tips and tricks to make GANs work by Soumith Chintala  
<https://github.com/soumith/ganhacks>

# Problems with GANs (cont.)

## ► Mode collapse

- Remember the generator's goal is to trick the discriminator  $D$  into thinking that the outputs by  $G$  are real, if generator  $G$  produces one sample which is realistic to the original data, then the discriminator finds it hard to distinguish between them.
- Fixes to mode collapse are mostly empirically driven: alternative architectures, alternative GAN loss, adding regularization terms, etc.



# Problems with GANs (cont.)



Figure 22: GAN mode collapse on MNIST digits dataset

- ▶ After the invention of GANs, there has been done a lot of research around them.
- ▶ A named list of GANs can be found [here](#).
- ▶ We will discuss the following variants here:
  - Wasserstein GAN
  - Conditional GAN
  - Cycle GAN

- ▶ Wasserstein GAN uses wasserstein distance instead of crossentropy loss.
- ▶ Wasserstein distance that has a smoother gradient everywhere.

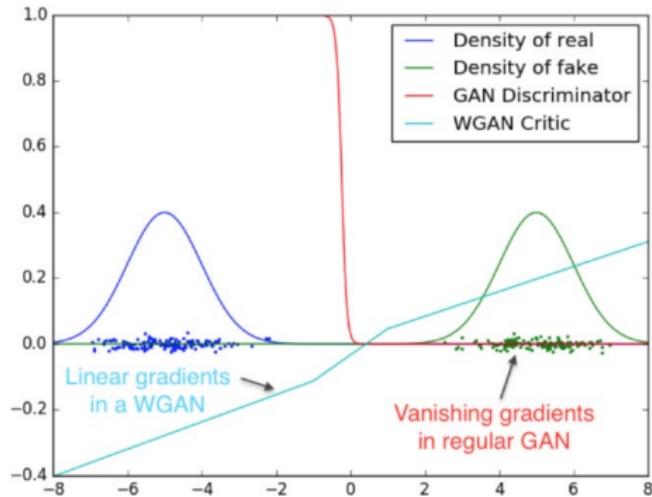
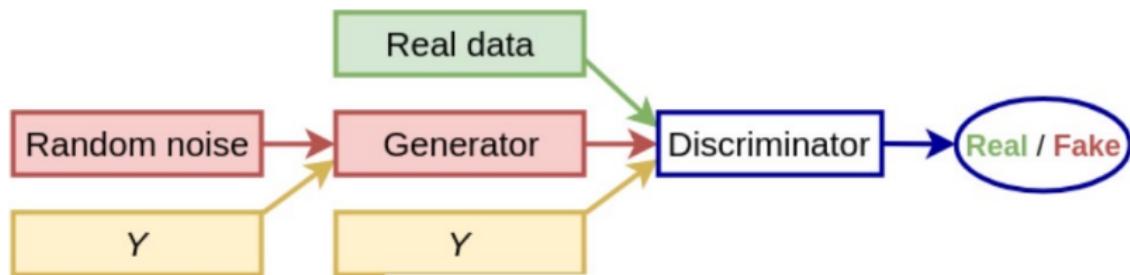


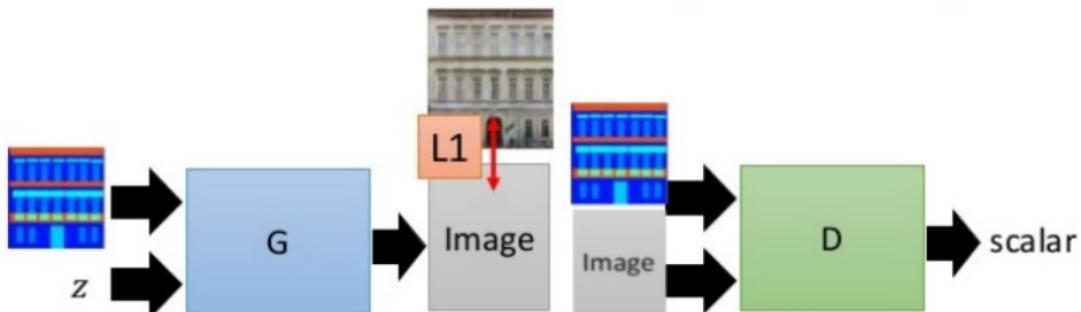


Figure 23: WGAN generation results on bedroom images

- ▶ Add a condition in addition to a random noise in the generator input. Similar case with discriminator.



# Conditional GAN - Image to Image



Testing:



Figure 24: Using L1 loss in addition in GAN loss can help in image to image translation

# Conditional GAN - Image to Image (cont.)

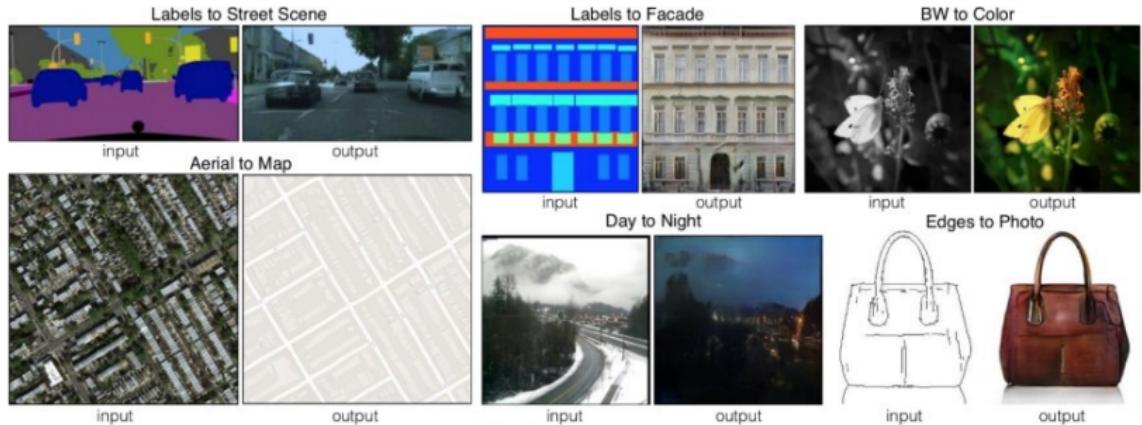
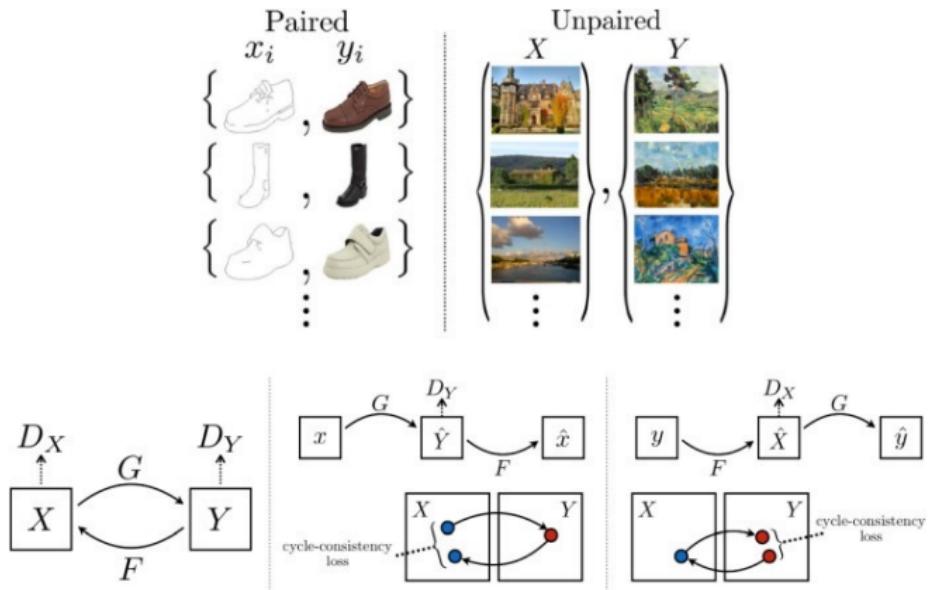


Figure 25: Image to Image translation with conditional GANs

# Cycle GANs

- ▶ Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks, Efros, ICCV 2017



# Cycle GANs (cont.)

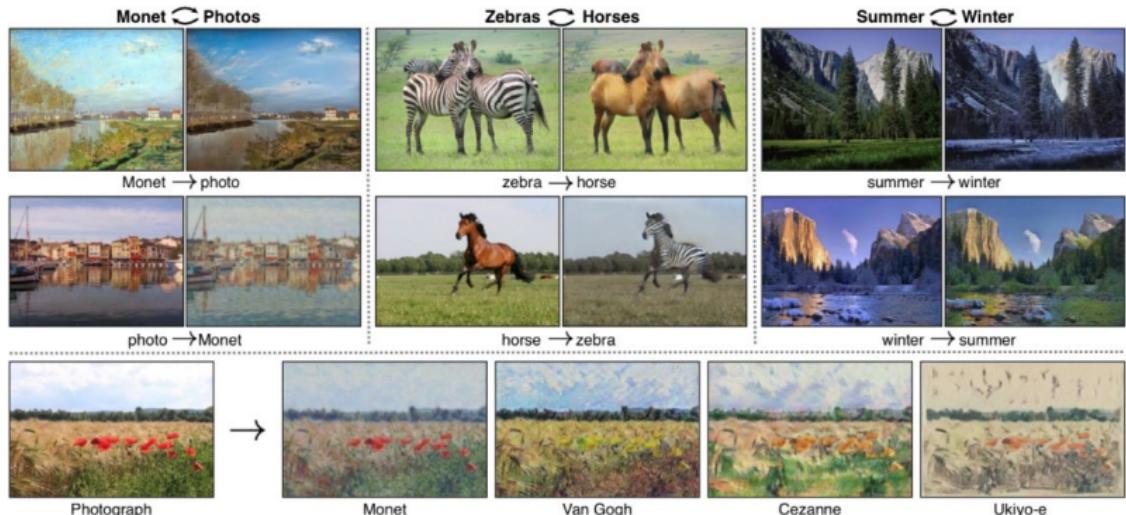


Figure 26: Image to Image translation with Cycle GANs

# GANs - Latent space Interpolation

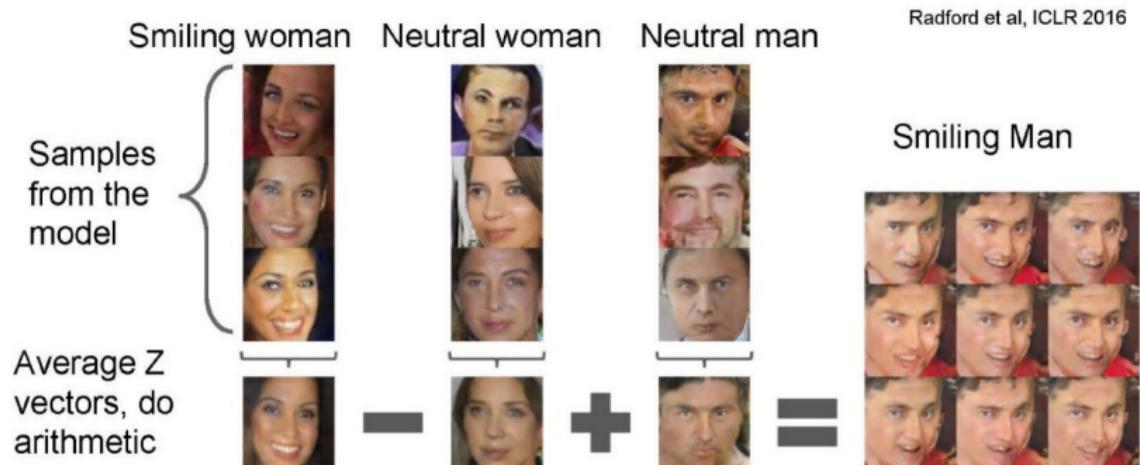


Figure 27: Latent space interpolation with GANs

# GANs - Latent space Interpolation (cont.)

Latent space interpolation with StyleGAN  
(Demo by Xander Steenburge)

<https://colab.research.google.com/drive/1mH70YxGNlnEaSOn0J8Lsgkl-QOvslb3MscrollTo=uEhxBvAR-7y3>

# Some more GAN results



Figure 28: Image Inpainting.

<https://www.nvidia.com/en-us/research/ai-demos/>

# Some more GAN results (cont.)

this small bird has a pink breast and crown, and black primaries and secondaries.



this white and yellow flower have thin white petals and a round yellow stamen



Figure 29: Text to Image Synthesis with GANs

## Some more GAN results (cont.)

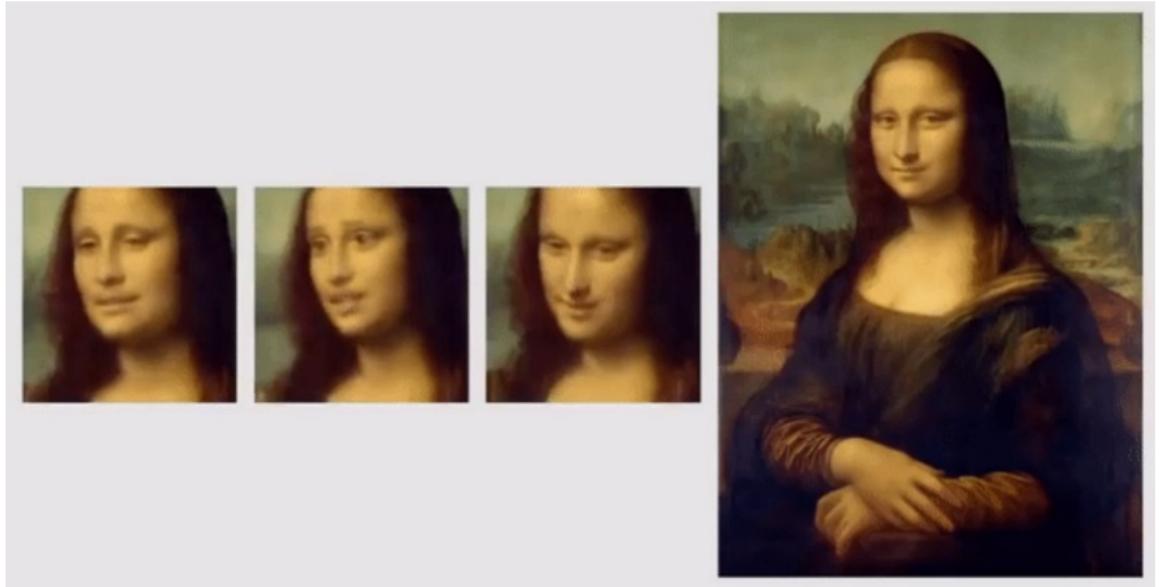
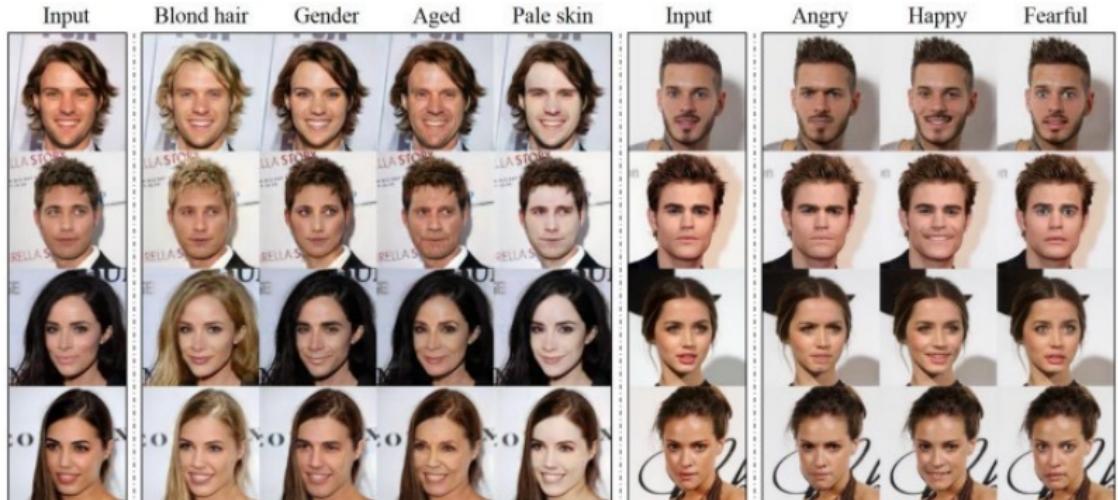


Figure 30: Living Portraits with GANs

# Some more GAN results (cont.)



**Figure 31:** Image to Image translation in multiple domains with StyleGAN (Choi et al.)

# Denoising Diffusion Probabilistic Models

- ▶ Denoising diffusion models, now also known as score-based generative models, have recently emerged as a powerful class of generative models. They demonstrate astonishing results in high-fidelity image generation, often even outperforming GANs.

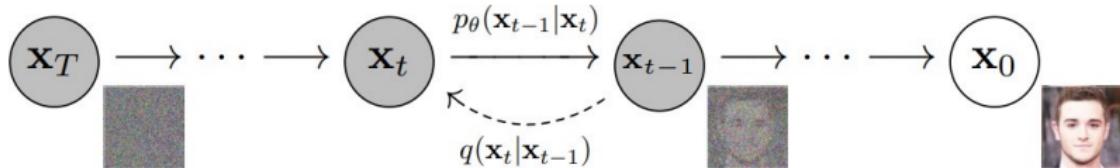


Figure 32: Diffusion Models Beat GANs on Image Synthesis [Dharwal & Nichol, OpenAI, 2021](#)

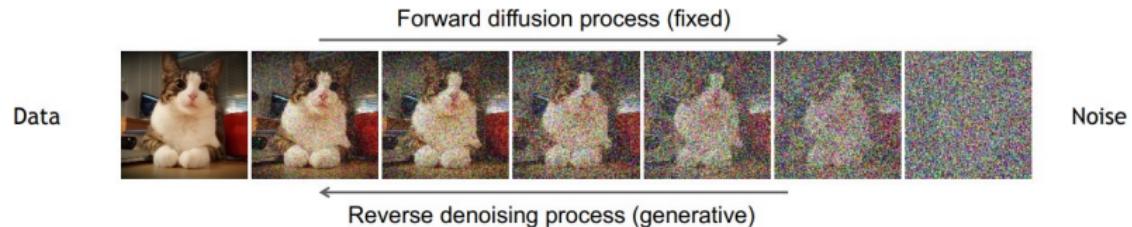
# Denoising Diffusion Probabilistic Models (cont.)

Denoising diffusion models consist of two processes:

- ▶ A fixed (or predefined) forward diffusion process  $q$  of our choosing, that gradually adds Gaussian noise to an image, until you end up with pure noise
- ▶ a learned reverse denoising diffusion process  $p_\theta$ , where a neural network is trained to gradually denoise an image starting from pure noise, until you end up with an actual image.



# Denoising Diffusion Probabilistic Models (cont.)



# Forward Diffusion Process

- ▶ Iteratively add Gaussian noise to the image.

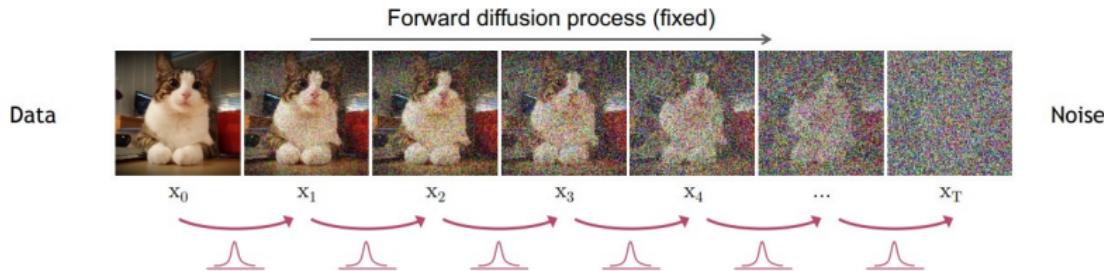
$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = N(\mathbf{x}_t; \sqrt{\frac{1 - \beta_t}{\beta_t}} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \longrightarrow q(\mathbf{x}_{1:T} | \mathbf{x}_0) = \prod_{t=1}^{T-1} q(\mathbf{x}_t | \mathbf{x}_{t-1})$$

where  $\beta$  is a known variance schedule such that

$0 < \beta_1 < \beta_2 < \dots < \beta_T < 1$ . Usually  $\beta$  is fixed but it can also be learned. This variance schedule can be linear, quadratic, cosine, etc.

- ▶ Basically, each new (slightly noisier) image at time step t is drawn from a conditional Gaussian distribution with  $\mu_t = \sqrt{1 - \beta_t} \mathbf{x}_{t-1}$  and  $\sigma_t^2 = \beta_t$ , which we can do by sampling  $\epsilon \sim N(\mathbf{0}, \mathbf{I})$  and then setting  $\mathbf{x}_t = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \beta_t \epsilon$ .

# Forward Diffusion Process (cont.)



► What if want to go directly from  $q_0$  to  $q_t$ ?

► Let  $\bar{\alpha}_t = \prod_{s=1}^t (1 - \beta_s)$ , then

$$q(\mathbf{x}_t | \mathbf{x}_0) = N(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$$

► This allows us, during training, to optimize random terms of the loss function  $L$  (or in other words, to randomly sample  $t$  during training and optimize  $L_t$ ).

# Reverse Denoising process

- ▶ In the reverse denoising process, we denoise Gaussian noise to generate an image from it.
- ▶ Basically, sample  $p(\mathbf{x}_T = N(\mathbf{x}_T; 0, 1)$
- ▶ Then,  $\mathbf{x}_{t-1} \sim p(\mathbf{x}_{t-1} | \mathbf{x}_t)$
- ▶ However, we don't know  $p(\mathbf{x}_{t-1} | \mathbf{x}_t)$ .
- ▶ So, let's approximate(learn) it with a neural network.
- ▶ If  $\beta_t$  is small enough at each time step we can assume  $p(\mathbf{x}_{t-1} | \mathbf{x}_t)$  to be a Gaussian Distribution. So, we can parameterize the process as

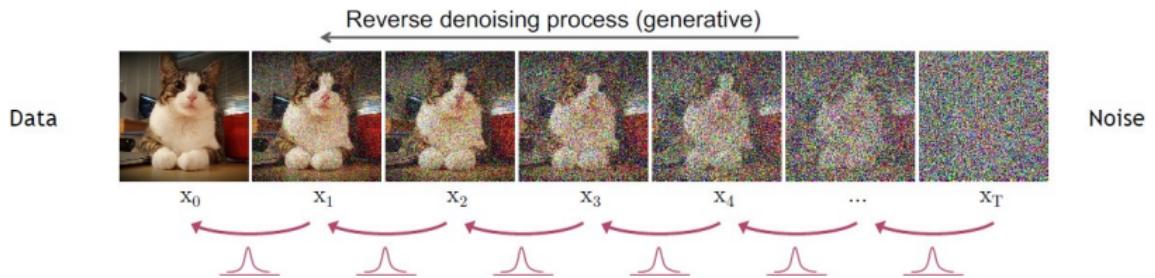
$$p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t) = N(\mathbf{x}_{t-1}; \mu_{\theta}(\mathbf{x}_t, t), \Sigma_{\theta}(\mathbf{x}_t, t))$$

where  $\mu_{\theta}$  and  $\Sigma_{\theta}$  are neural networks. Mean and variance are also conditioned on  $t$ .

# Reverse Denoising process (cont.)

- ▶ But we fix variance for now, so we only need  $\mu_\theta$ . Later variants also learn  $\Sigma_\theta$ .

$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = N(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \beta_t \mathbf{I}) \longrightarrow p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_t) \prod_{t=1}^{T-1} p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)$$



# Learning Denoising Models

- ▶ For training, we can form variational upper bound that is commonly used for training variational autoencoders:

$$E_{q(x_0)}[-\log p_\theta(x_0)] \leq E_{q(x_0)q(x_{1:T} | x_0)} \left( -\log \frac{p_\theta(x_{0:T})}{q(x_{1:T} | x_0)} \right)^1$$

- ▶ It can be shown that the ELBO for this process is a sum of losses at each time step  $t$ ,

$$L = L_0 + L_1 + \dots + L_T$$

# Learning Denoising Models (cont.)

- ▶ Lastly, we can reparametrize the mean to make the neural network learn (predict) the added noise (via a network  $\epsilon_\theta(\mathbf{x}_t, t)$  for noise level  $t$  in the KL terms which constitute the losses. This means that our neural network becomes a noise predictor, rather than a (direct) mean predictor. The mean can be computed as follows:

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

- ▶ The final objective function  $L_t$  then looks as follows (for a random time step  $t$  given  $\epsilon \sim N(\mathbf{0}, \mathbf{I})$ ):

$$\|\epsilon - \epsilon_\theta(\mathbf{x}_t, t)\|_2 = \|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)} \epsilon, t)\|_2.$$

- ▶ For complete derivation, read [here](#)

---

## Algorithm 1 Training

---

1: **repeat**

2:  $\mathbf{x}_0 \sim q(\mathbf{x}_0)$

3:  $t \sim \text{Uniform}(\{1, \dots, T\})$

4:  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

5: Take gradient descent step on

$$\nabla_{\theta} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t)\|^2$$

6: **until** converged

---

---

## Algorithm 2 Sampling

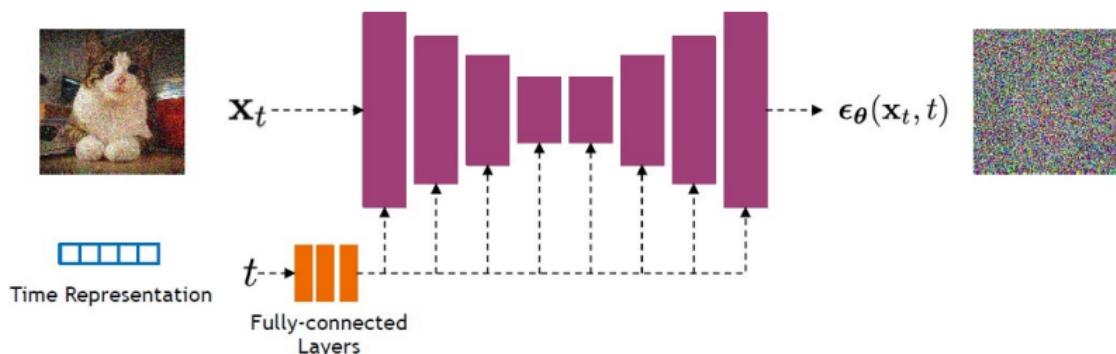
---

```
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:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
```

---

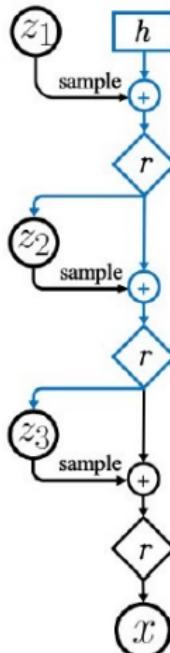
# Network Architectures

- ▶ Diffusion models often use U-Net architectures with ResNet blocks and self-attention layers to represent  $\epsilon_\theta(x_t, t)$
- ▶ Time representation: sinusoidal positional embeddings or random Fourier features.



# Connection to VAEs

- ▶ Diffusion models can be considered as a special form of hierarchical VAEs (one VAE after another).
- ▶ However, in diffusion models:
  - The encoder is fixed
  - The latent variables have the same dimension as the data
  - The denoising model is shared across different timestep
  - The model is trained with some reweighting of the variational bound.



# Problems with Diffusion Models

- ▶ For now, it seems that the major limitation of the Diffusion Models is its notoriously slow sampling procedure which normally requires hundreds to thousands of time discretization steps of the learned diffusion process to reach the desired accuracy.

# Trilemma of generative learning

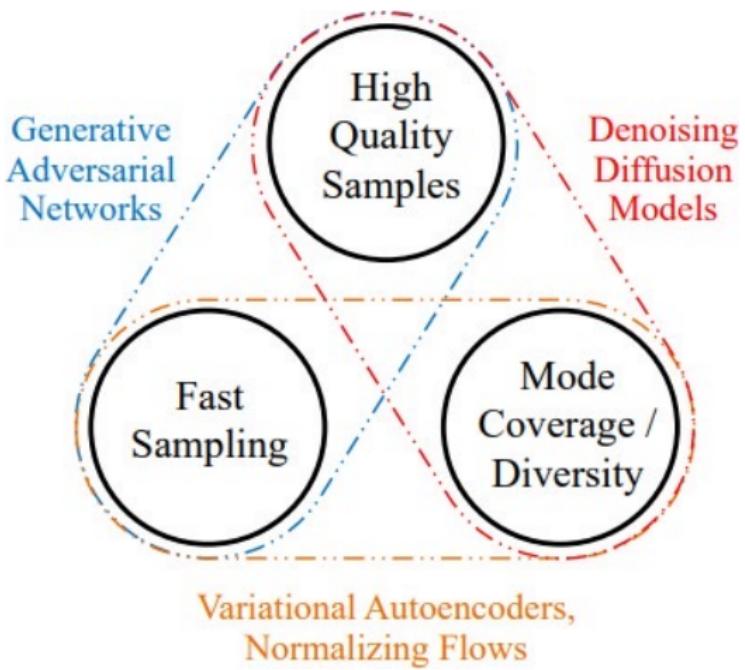


Figure 33: Generative learning Trilemma

<sup>1</sup>[Tackling the Generative Learning Trilemma with Denoising Diffusion GANs](#)

► Denoising Diffusion Implicit Model

- DDIM roughly sketches the final sample then refine it with the reverse process.

► Improved Denoising Diffusion Probabilistic Models

- Train  $\sigma^2$  while training the diffusion model instead of fixing it.

► Score-Based Generative Modeling through Stochastic Differential Equations

- Model the gradient of the log probability density function, a quantity often known as the (Stein) score function.

► Tackling the Generative Learning Trilemma with Denoising Diffusion GANs

- Introduce denoising diffusion generative adversarial networks (denoising diffusion GANs) that model each denoising step using a multimodal conditional GAN.

► Cascaded Diffusion Models for High Fidelity Image Generation

- Cascaded diffusion models to boost sample quality.



“a man wearing a white hat”

Figure 34: Image Inpainting with GLIDE

<sup>1</sup>[GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models](#)





a shiba inu wearing a beret and black turtleneck



a close up of a handpalm with leaves growing from it

Figure 35: Text to image generation with DALL.E 2

<sup>1</sup>[Hierarchical Text-Conditional Image Generation with CLIP Latents](#)

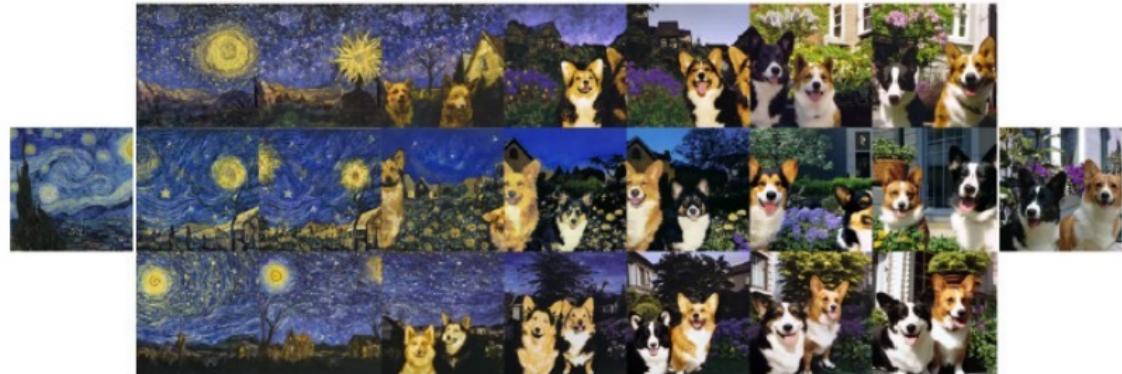


Fix the CLIP embedding  $z$ .

Decode using different decoder latents  $x_T$ .

Figure 36: Image Variations

<sup>1</sup>[Hierarchical Text-Conditional Image Generation with CLIP Latents](#)



Interpolate image CLIP embeddings  $z_i$

Use different  $x_T$  to get different interpolation trajectories.

Figure 37: Image interpolation

<sup>1</sup>[Hierarchical Text-Conditional Image Generation with CLIP Latents](#)



Change the image CLIP embedding towards the difference of the text CLIP embeddings of two prompts.

Decoder latent is kept as a constant.

Figure 38: Text Difference Image interpolation

<sup>1</sup>[Hierarchical Text-Conditional Image Generation with CLIP Latents](#)



A brain riding a rocketship heading towards the moon.

<sup>1</sup>[Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding](#)





A dragon fruit wearing karate belt in the snow.

<sup>1</sup>[Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding](#)





A relaxed garlic with a blindfold reading a newspaper while floating in a pool of tomato soup.

<sup>1</sup>[Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding](#)



# Diffusion Autoencoders

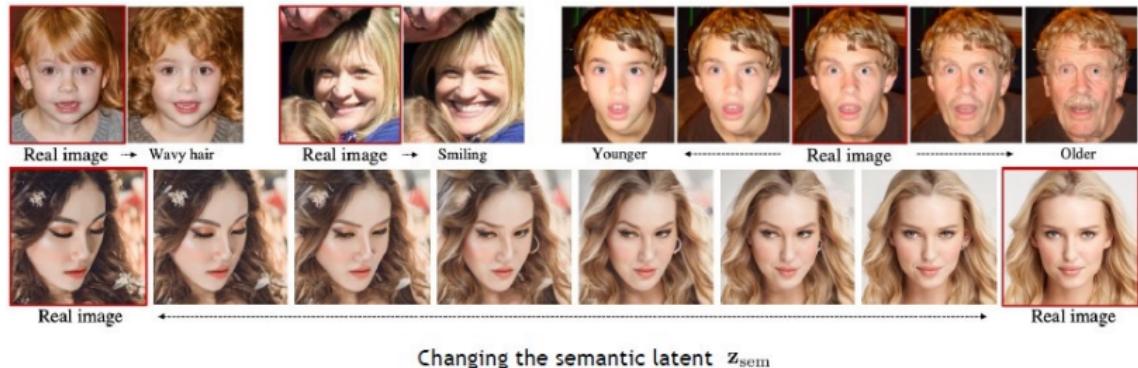


Figure 39: Learning semantic meaningful latent representations in diffusion models

<sup>1</sup>[Diffusion Autoencoders: Toward a Meaningful and Decodable Representation](#)

# Super Resolution

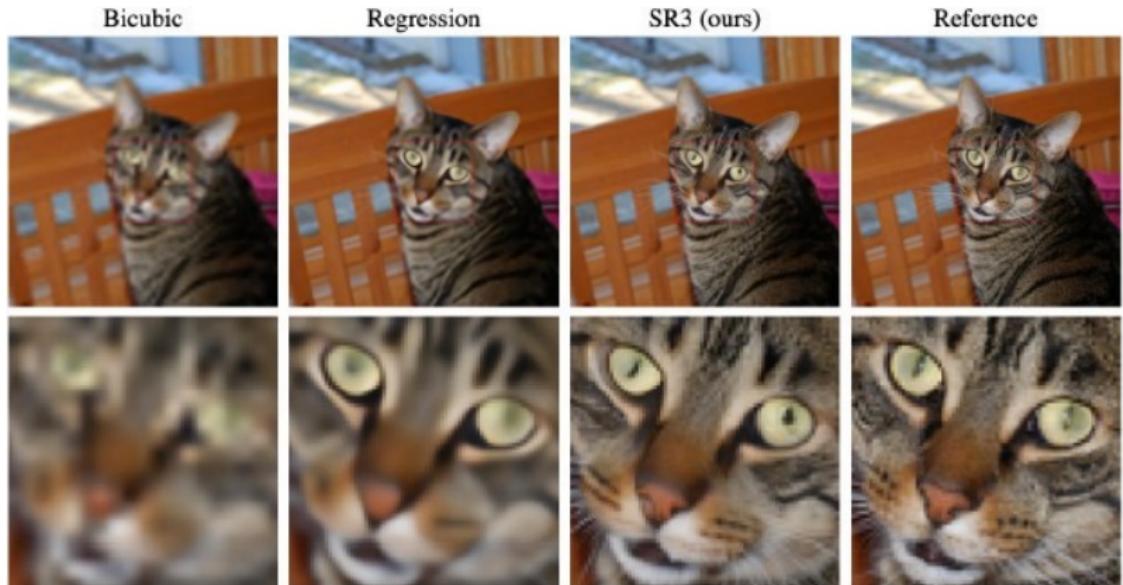
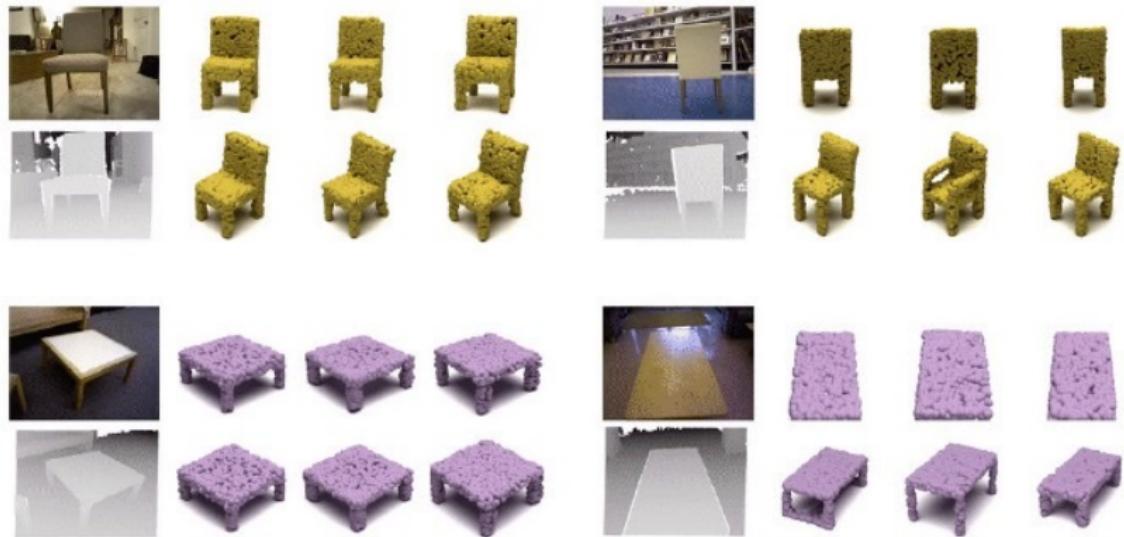


Figure 40: Natural Image Super Resolution  $64 \times 64 \rightarrow 256 \times 256$

<sup>1</sup>[Image Super-Resolution via Iterative Refinement](#)

# 3D Shape Generation



<sup>1</sup>[3D Shape Generation and Completion through Point-Voxel Diffusion](#)

# Try It Yourself!

Text to Image generation with stable diffusion.

<https://huggingface.co/spaces/stabilityai/stable-diffusion>

---

<sup>1</sup>[High-Resolution Image Synthesis with Latent Diffusion Models](#)

# The Generative AI Landscape

## Text

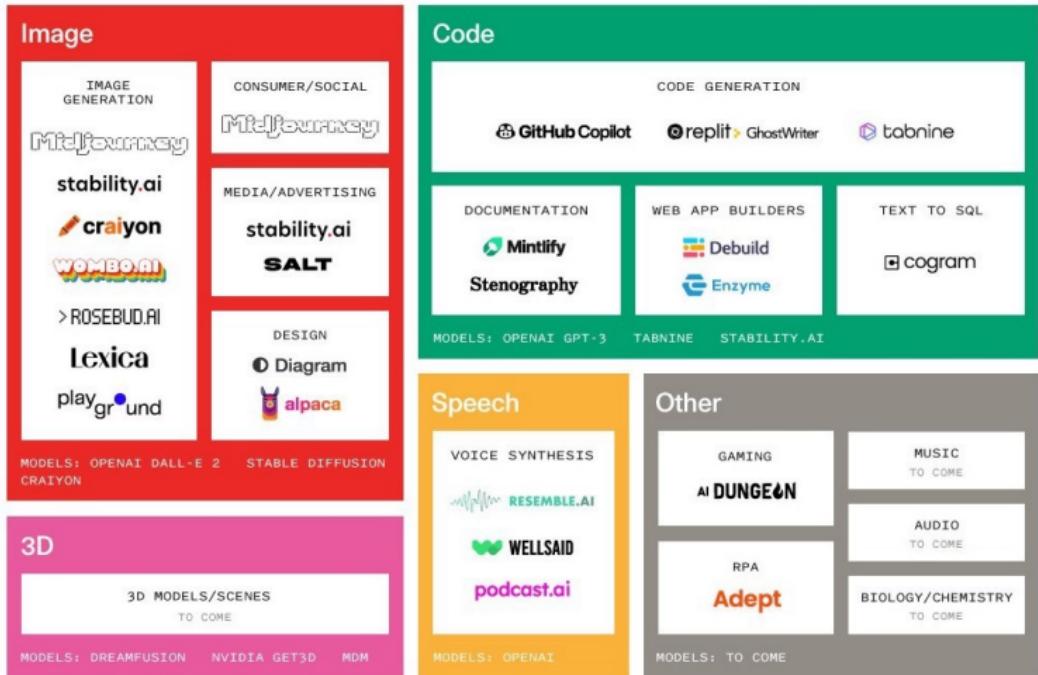


## Video



<sup>1</sup>[Sonya Huang \(@sonyatweetybird\)](#)

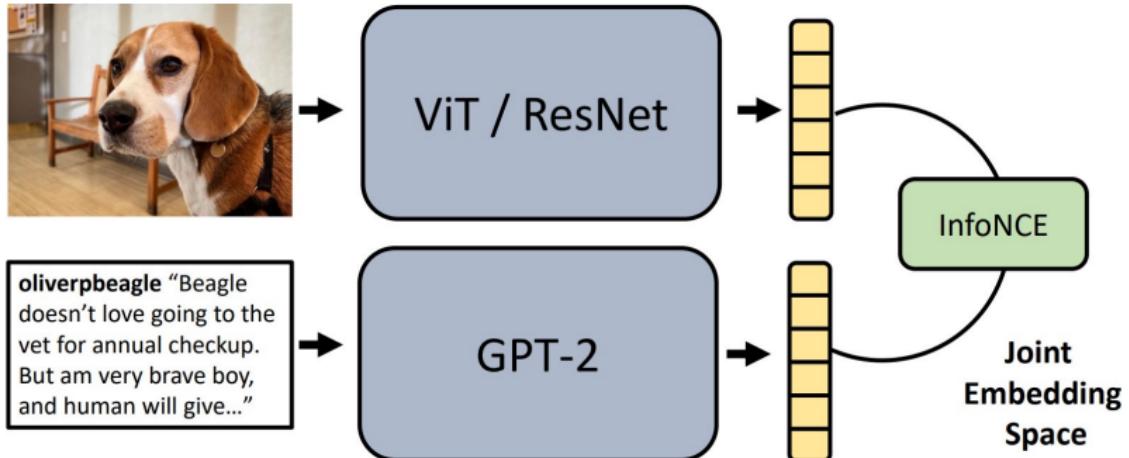
# The Generative AI Landscape



<sup>1</sup>Sonya Huang (@sonyatweetybird)

# Image and Text

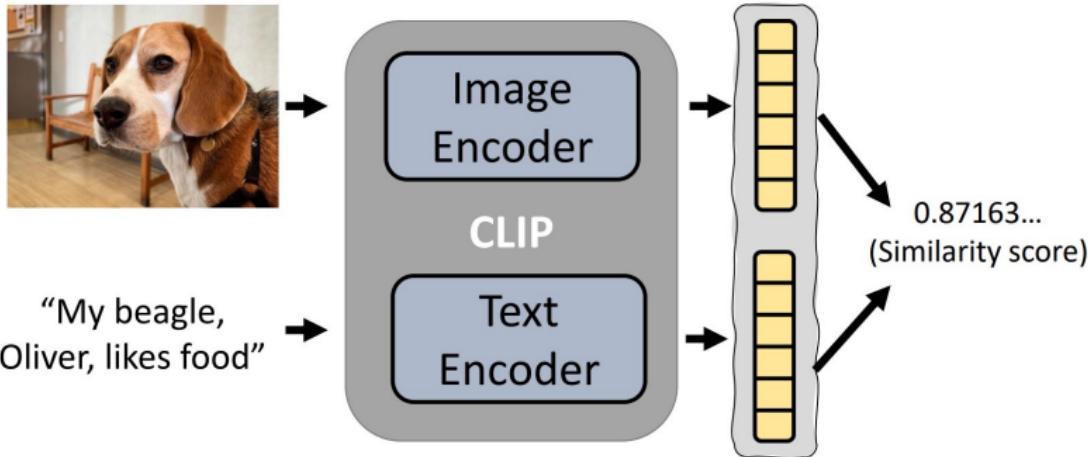
# CLIP – Connecting Images and Text



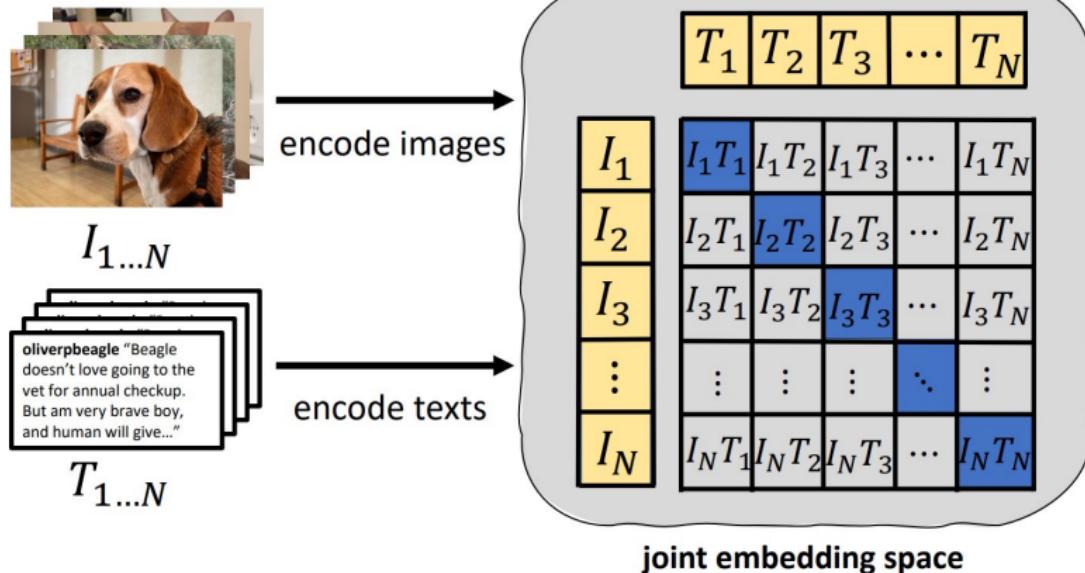
<sup>0</sup>Radford et. al, Learning Transferable Visual Models From Natural Language Supervision, ArXiv'21

# CLIP – Connecting Images and Text

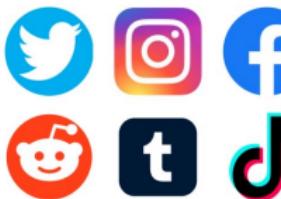
- Contrastive Language Image Pretraining



# CLIP – Connecting Images and Text



# CLIP – Connecting Images and Text

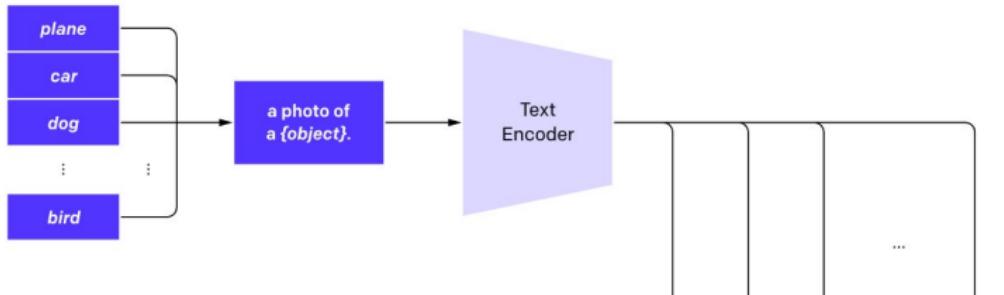


**oliverpbeagle** "Beagle doesn't love going to the vet for annual checkup. But am very brave boy, and human will give me many treats afterwards" 🐶 #oliverpbeagle #vet #beagle #beaglesofinstagram

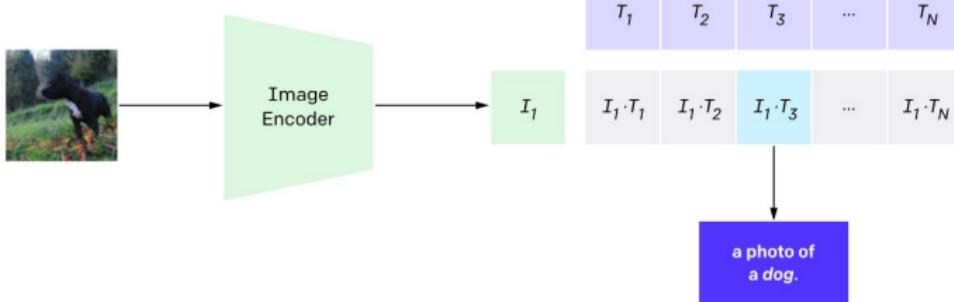
) × 400 Million

# CLIP - Zero Shot Capabilities

## 2. Create dataset classifier from label text



## 3. Use for zero-shot prediction



<sup>0</sup><https://openai.com/research/clip>

# CLIP - Zero Shot Capabilities

## Caltech-101

Kangaroo (99.8%) Ranked 1 out of 102 labels



- a photo of a kangaroo.
- a photo of a gerenuk.
- a photo of a rhea.
- a photo of a ruru.
- a photo of a wild cat.
- a photo of a scorpion.

## Oxford-IIIT Pets

Maine Coon (100.0%) Ranked 1 out of 37 labels



- a photo of a maine coon, a type of pet.
- a photo of a persian, a type of pet.
- a photo of a ragdoll, a type of pet.
- a photo of a birman, a type of pet.
- a photo of a siamese, a type of pet.

## ImageNet-R (Rendition)

Siberian Husky (76.0%) Ranked 1 out of 200 labels



- a photo of a siberian husky.
- a photo of a german shepherd dog.
- a photo of a collie.
- a photo of a border collie.
- a photo of a rottweiler.

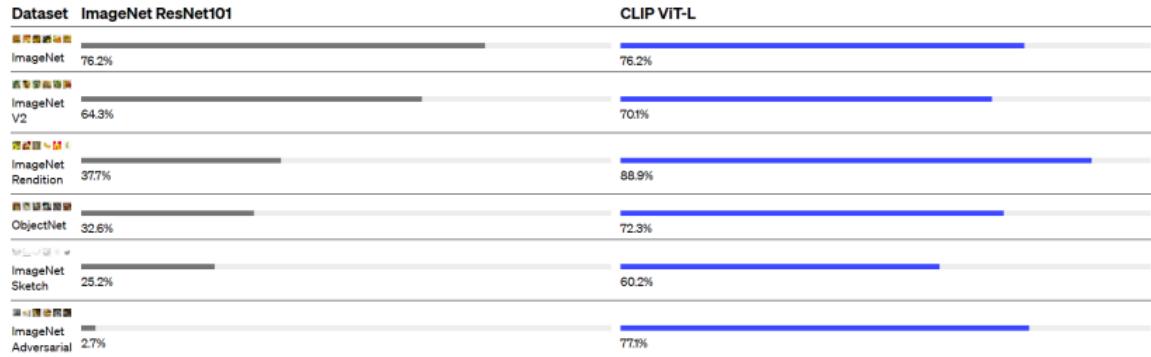
## CIFAR-100

snake (38.0%) Ranked 1 out of 100 labels



- a photo of a snake.
- a photo of a sweet pepper.
- a photo of a flatfish.
- a photo of a turtle.
- a photo of a lizard.

# CLIP - Zero Shot Capabilities

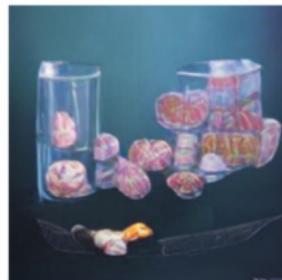


<sup>0</sup><https://openai.com/research/clip>

# Using CLIP for generative tasks

Generation

A beautiful painting of a building in a serene landscape

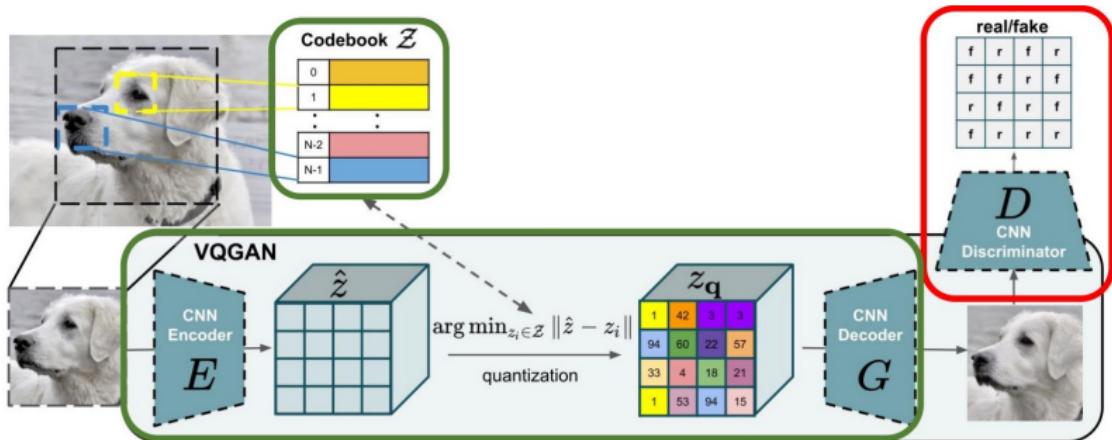


Editing

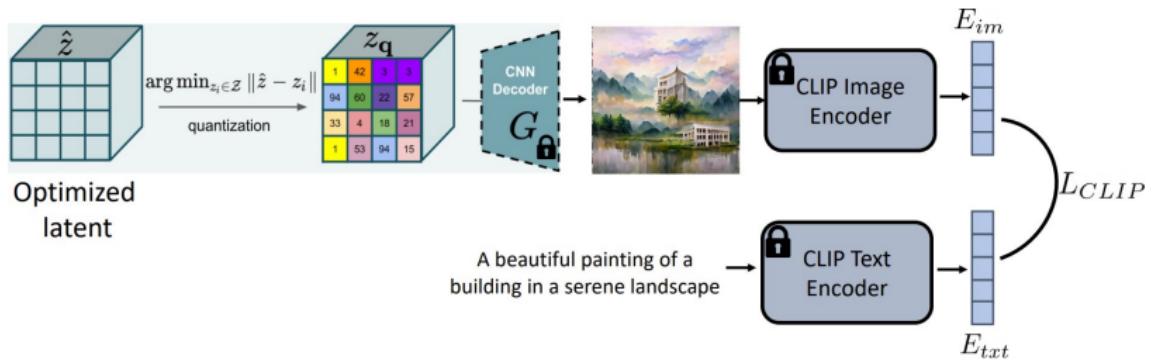


"A cake made of ice"





# VQ-GAN + CLIP

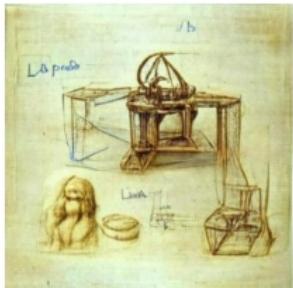


$$L_{CLIP} = 1 - \text{Cos}(E_{im}, E_{txt})$$

---

<sup>0</sup>VQGAN-CLIP: Open Domain Image Generation and Editing with Natural Language Guidance, Crowson et al. ICCV 13 2021

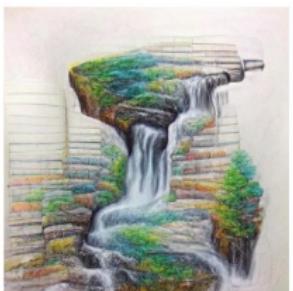
# VQ-GAN + CLIP Results



“A sketch of 3D printer by da Vinci”

“An autogyro flying, artstation”

“A futuristic city in synthwave style”



“A colored pencil drawing of a waterfall”

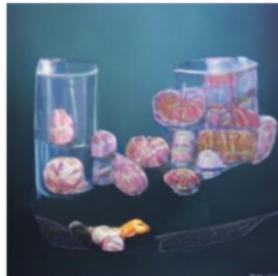
“A painting of a city in a deep valley”

“Baba Yaga’s house, fantasy art”

# Using CLIP for generative tasks

**Generation**

A beautiful  
painting of a  
building in a  
serene landscape

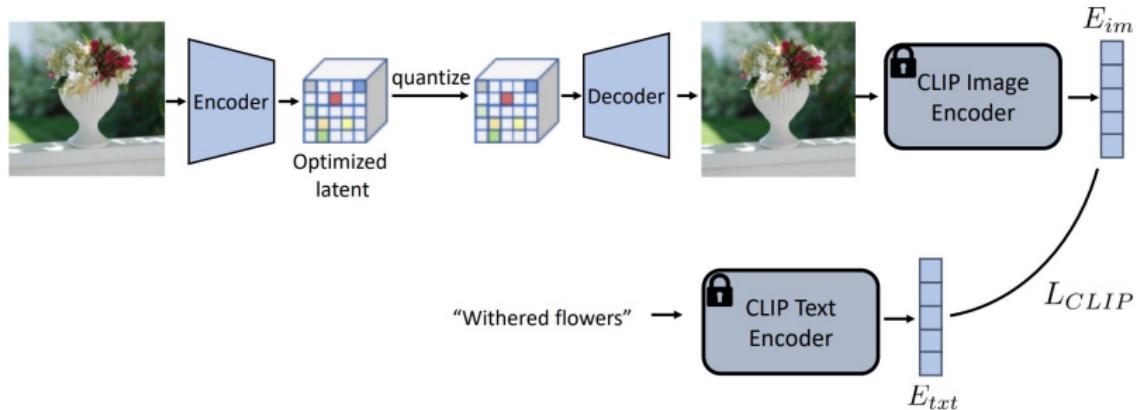


**Editing**

"A cake made of ice"



# VQ-GAN + CLIP Editing



<sup>0</sup>VQGAN-CLIP: Open Domain Image Generation and Editing with Natural Language Guidance, Crowson et al. ICCV 13 2021

# VQ-GAN + CLIP Editing Results

Instruction

“Wooden”

Original



VQGAN-CLIP



“Withered Flowers”

“Focused”



<sup>0</sup>VQGAN-CLIP: Open Domain Image Generation and Editing with Natural Language Guidance, Crowson et al. ICCV 13 2021

# StyleCLIP - Results



Input

"Beyonce"

"A woman  
without  
makeup"

"Elsa from Frozen"



Input

"A man with  
..."

"A blonde man"

"Donald Trump"

---

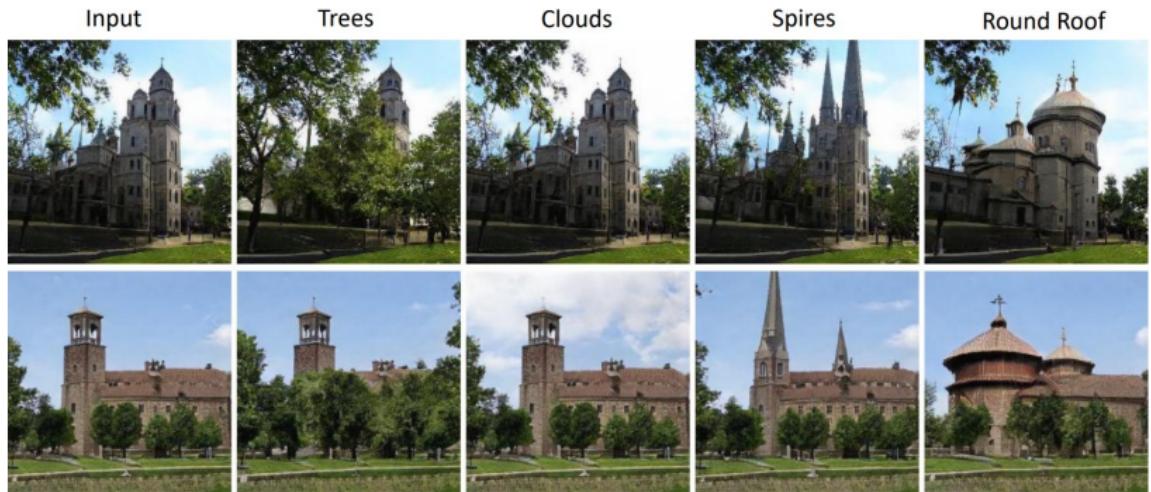
<sup>0</sup>StyleCLIP: Text-Driven Manipulation of StyleGAN Imagery, Patashnik and Wu et al. ICCV 2021

# StyleCLIP - Results



<sup>0</sup>StyleCLIP: Text-Driven Manipulation of StyleGAN Imagery, Patashnik and Wu et al. ICCV 2021

# StyleCLIP - Results



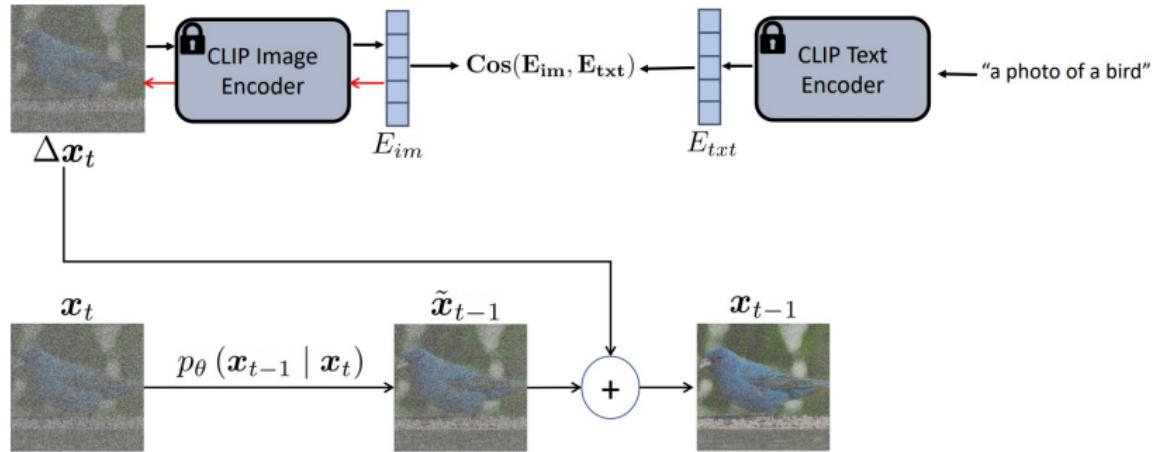
<sup>0</sup>StyleCLIP: Text-Driven Manipulation of StyleGAN Imagery, Patashnik and Wu et al. ICCV 2021

# Text2Live - Results



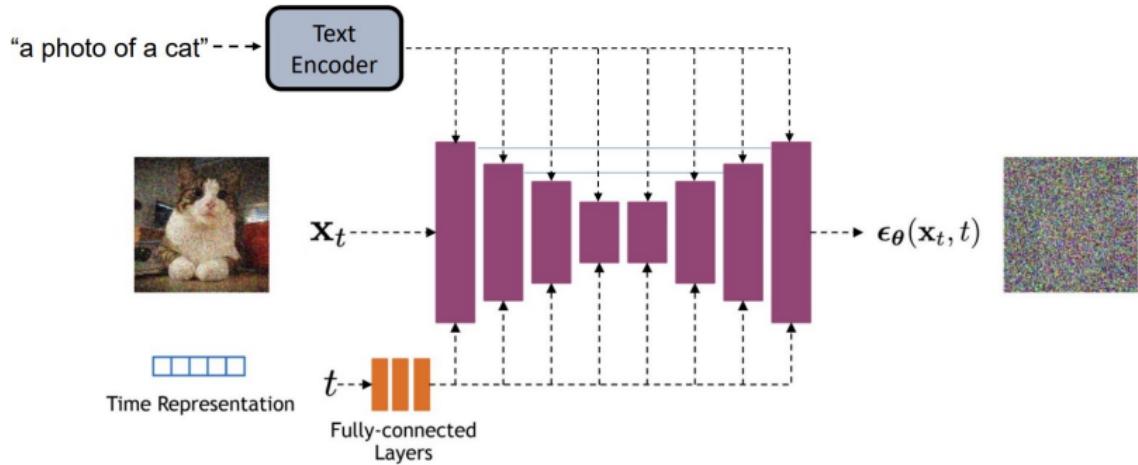
<sup>0</sup>Text2LIVE: Text-Driven Layered Image and Video Editing, Bar-Tal Ofri-Amar and Fridman et al. ECCV 2022

# Conditional image generation with CLIP Using Diffusion Models

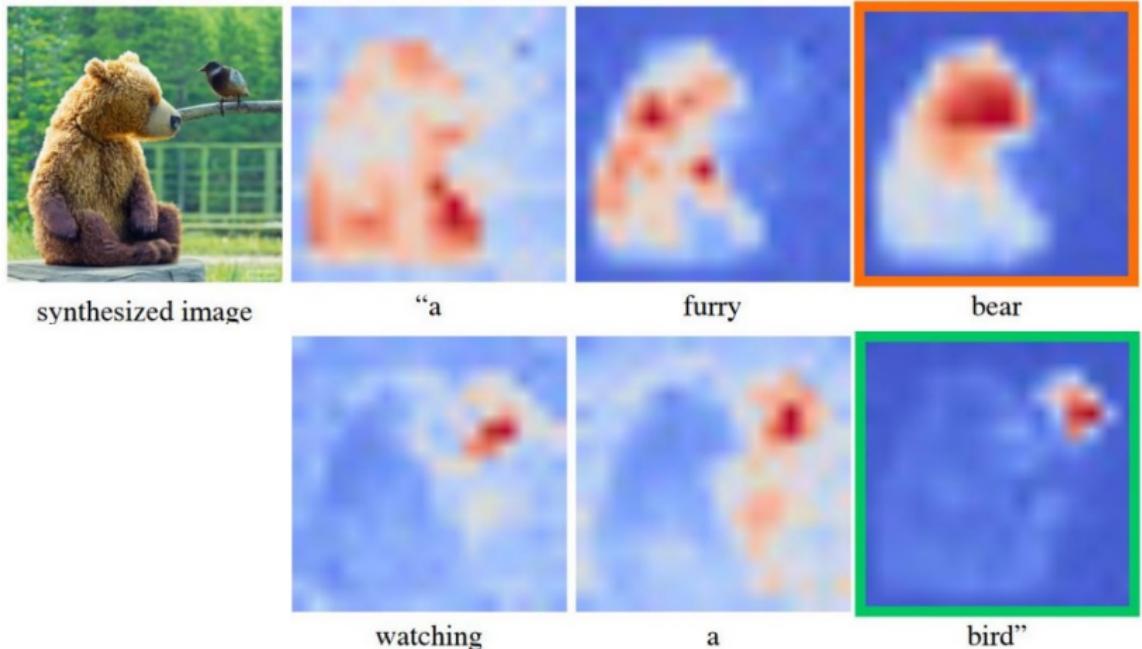


<sup>0</sup>Diffusion Models Beat GANs on Image Synthesis, Dhariwal and Nichol et al.  
NeurIPS 2021

# Conditional image generation with CLIP Using Diffusion Models



# Text conditioning in Diffusion Models



Average attention maps across all timestamps

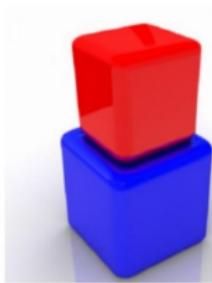
# Results



“a hedgehog using a calculator”



“a corgi wearing a red bowtie and purple party hat”



“a red cube on top of a blue cube”



“a high-quality oil painting of a psychedelic hamster dragon”

<sup>0</sup>GLIDE: Towards Photorealistic Image Generation and Editing with Text Nichol et al., 2021, GLIDE -Guided Diffusion Models, Nichol et al. [2021](#)

## Reference Slides

- ▶ Fei-Fei Li "Generative Deep Learning" CS231
- ▶ Hao Dong "Deep Generative Models"
- ▶ Hung-Yi Lee "Machine Learning"
- ▶ Francois Fleuret "Deep Learning" EE559
- ▶ Murtaza Taj "Deep Learning" CS437
- ▶ Kreis, Gao & Vahdat [CVPR 2022 Tutorial](#)
- ▶ Rogge & Rasul [The Annotated Diffusion Model](#)