

Deep Generative Learning

Valeriya Strizhkova

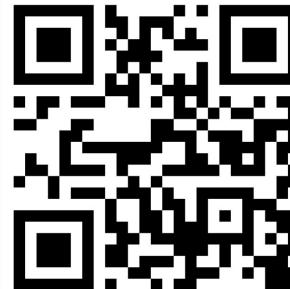
About myself

Valeriya Strizhkova

PhD @ INRIA, CEO @ Facila

Research interests:

- computer vision
- facial expression recognition
- human behavior understanding
- multi-modal fusion



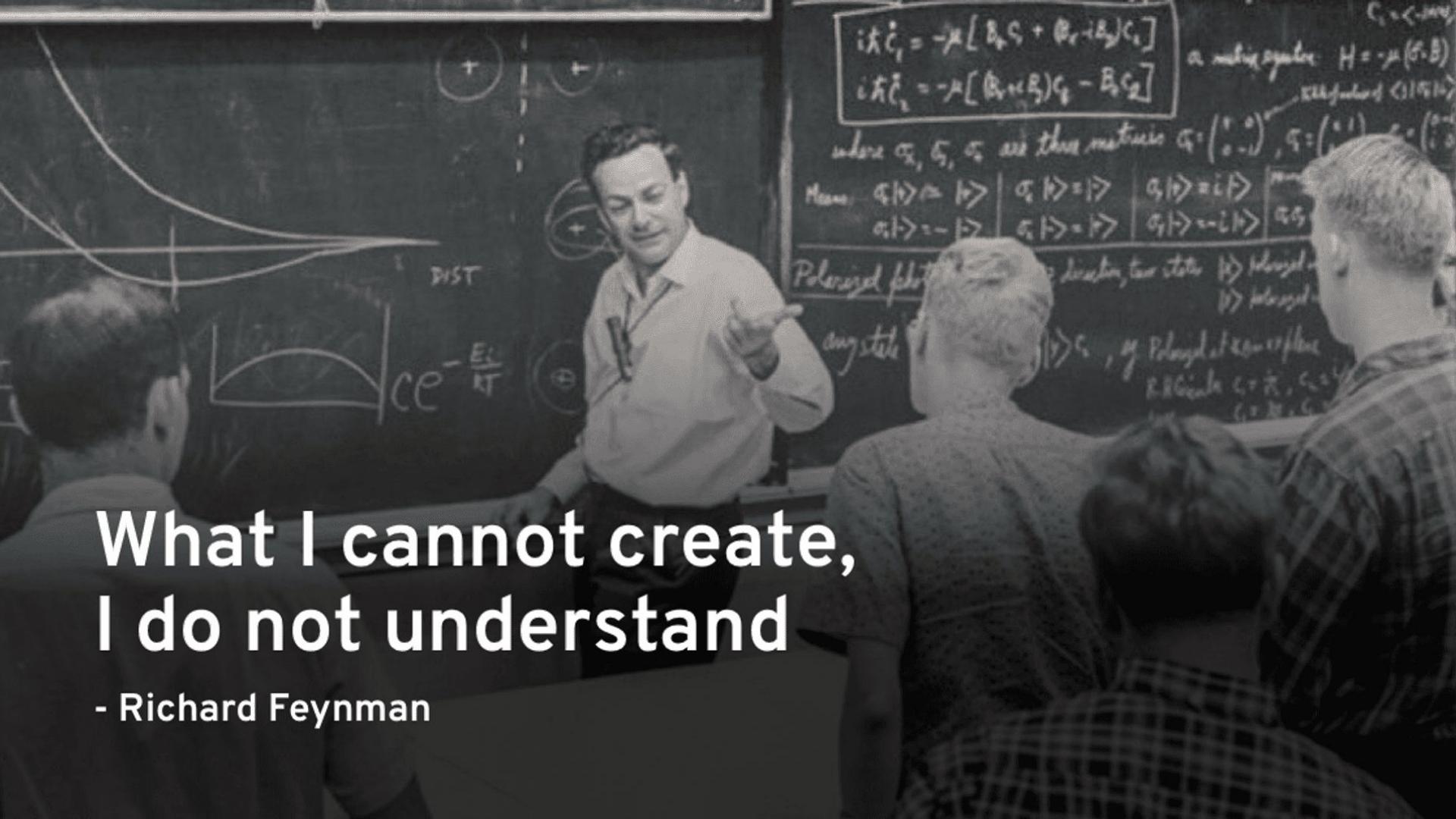
Outline

- Introduction
 - What is a generative model?
 - Types of deep generative models
- Image Generation
 - Autoregressive models
 - Generative Adversarial Networks
 - Variational Autoencoders
 - Diffusion Models
- Evaluation
 - Inception Score
 - Frechet Inception Distance

Introduction

What I cannot create, I do not understand

- Richard Feynman



Generative AI in Image Applications

Art & Design



Content Generation



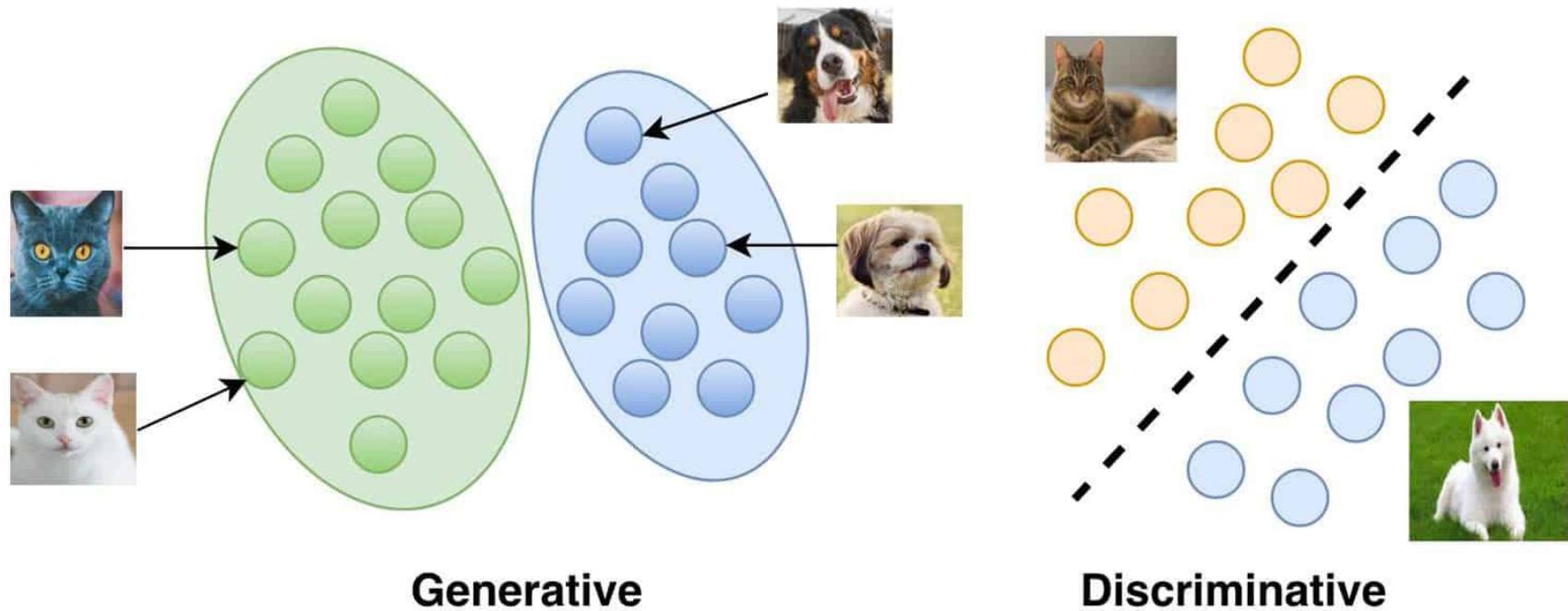
Representation Learning



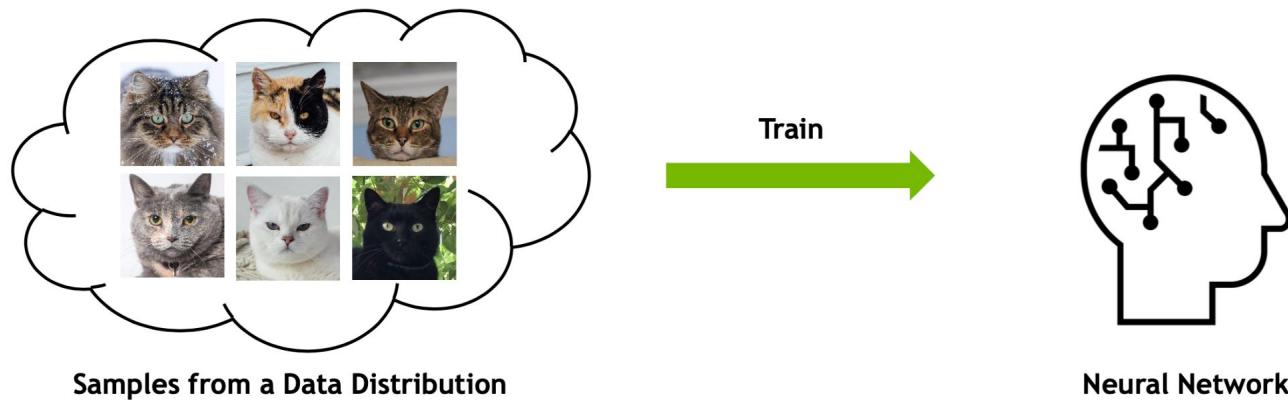
Entertainment



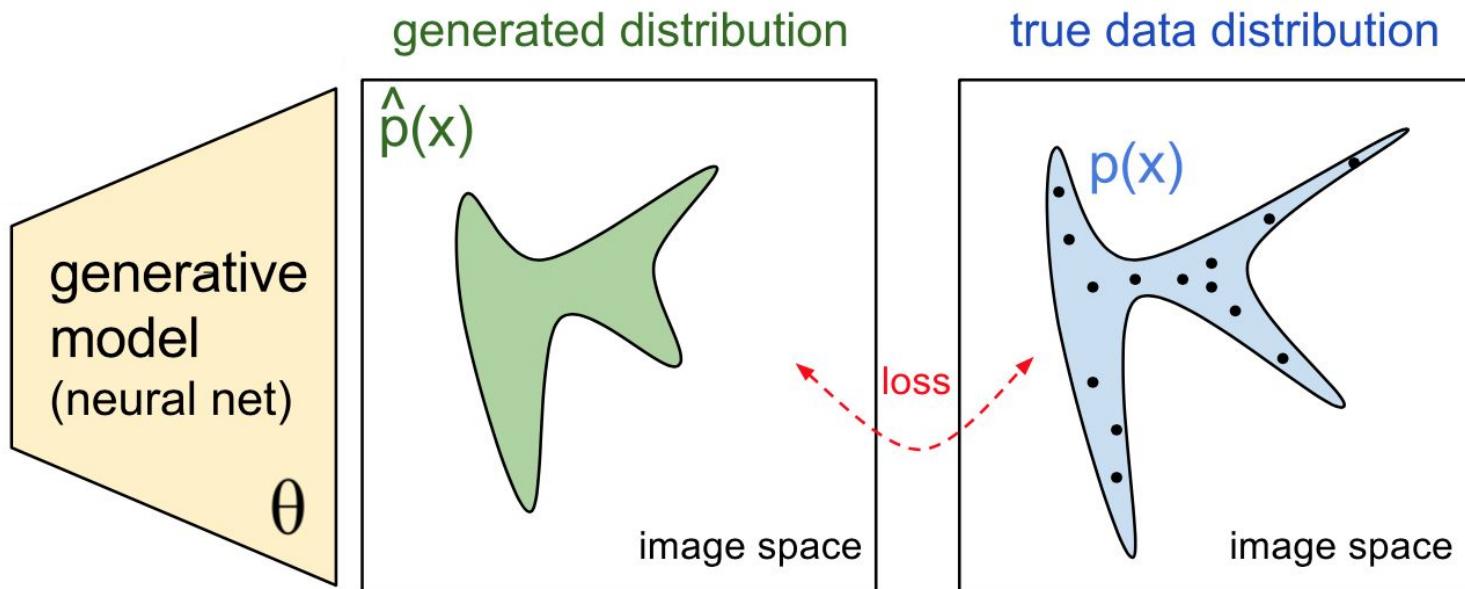
What is a Generative Model?



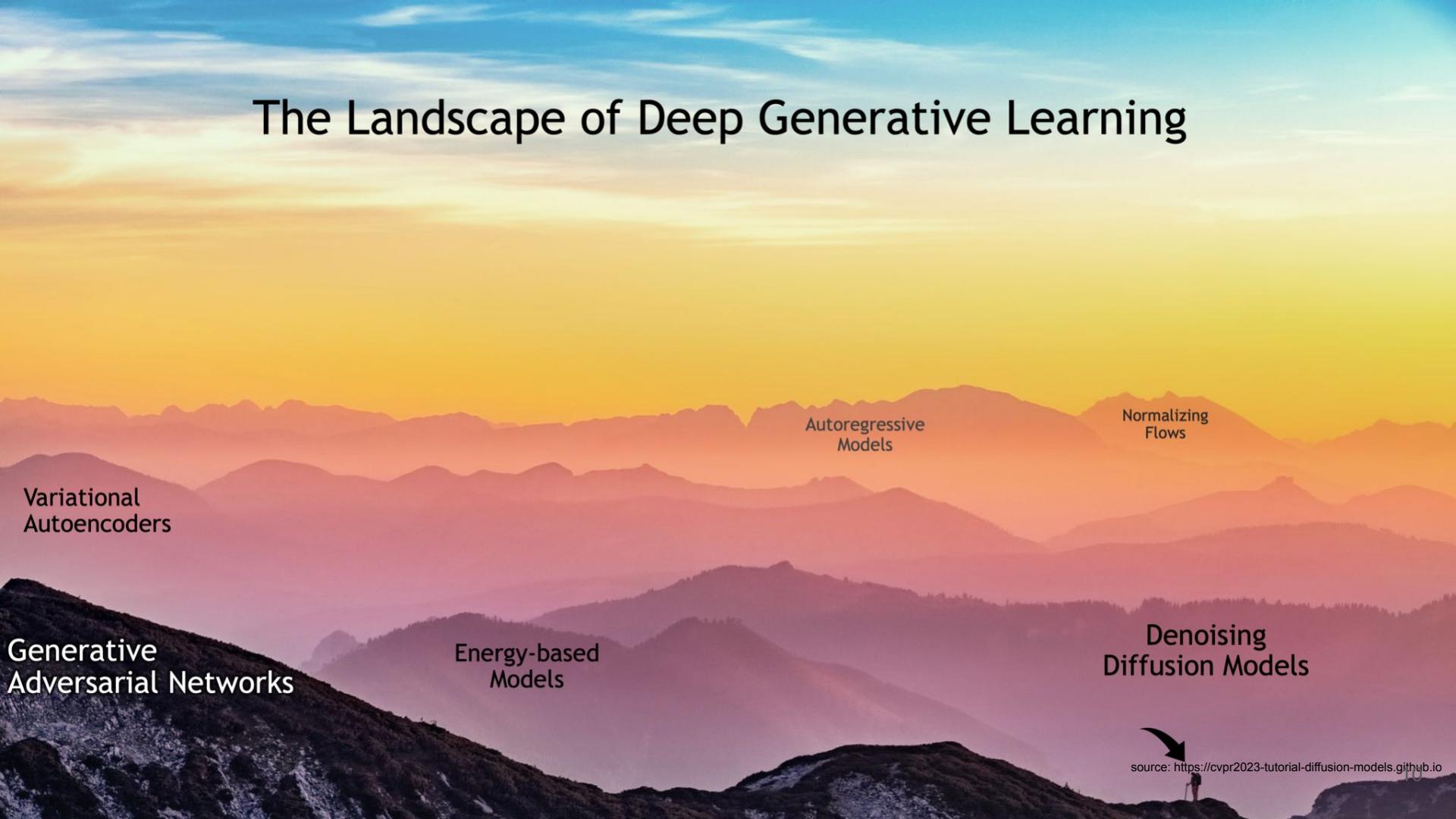
What is a Generative Model?



What is a Generative Model?



The Landscape of Deep Generative Learning



Variational
Autoencoders

Generative
Adversarial Networks

Energy-based
Models

Autoregressive
Models

Normalizing
Flows

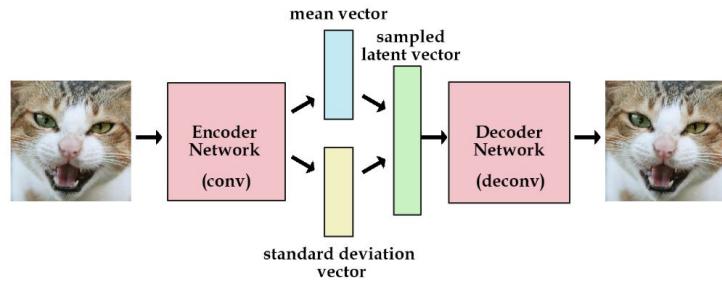
Denoising
Diffusion Models



source: <https://cvpr2023-tutorial-diffusion-models.github.io>

Types of generative models

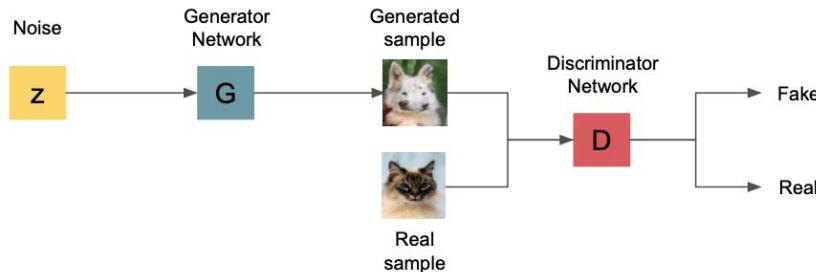
Variational Autoencoder:



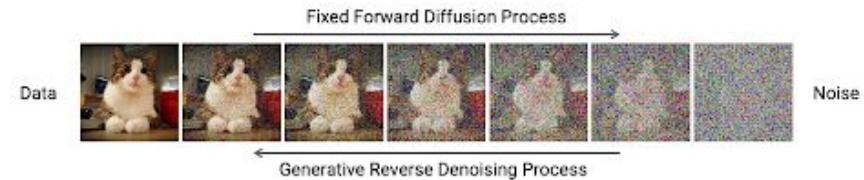
Autoregressive model:



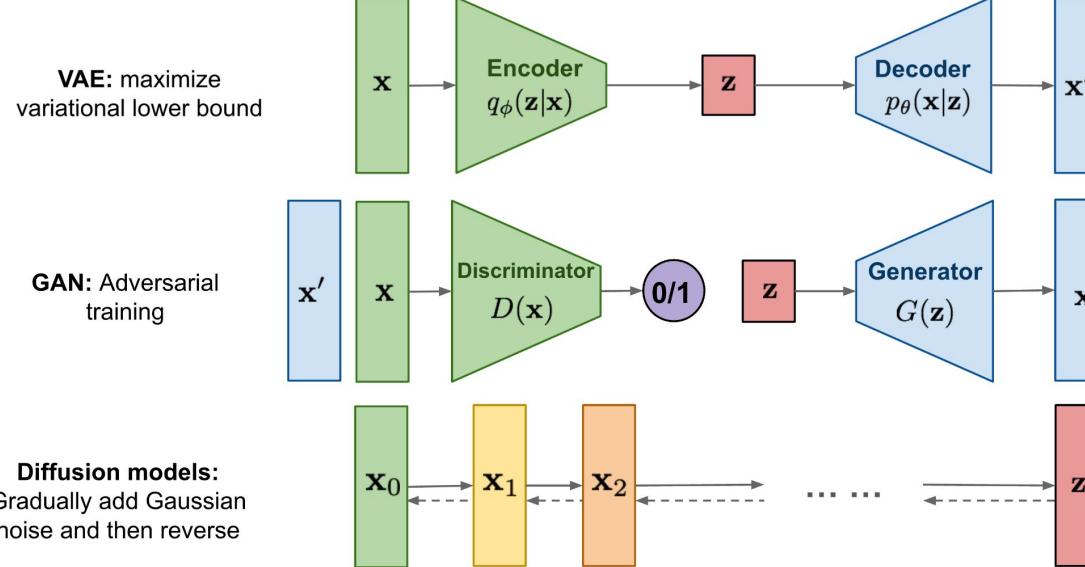
Generative Adversarial Network:



Diffusion model:



Types of generative models



Generative learning trilemma

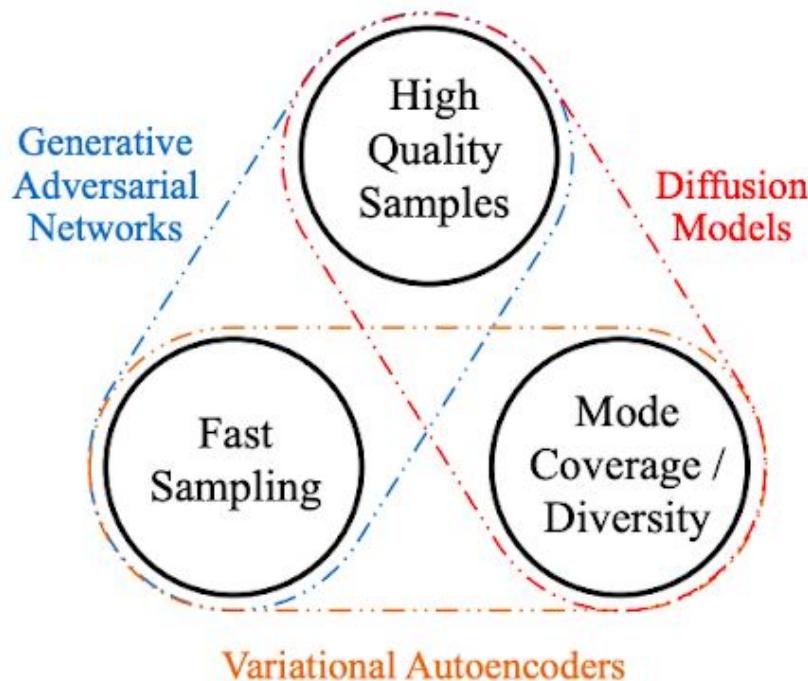
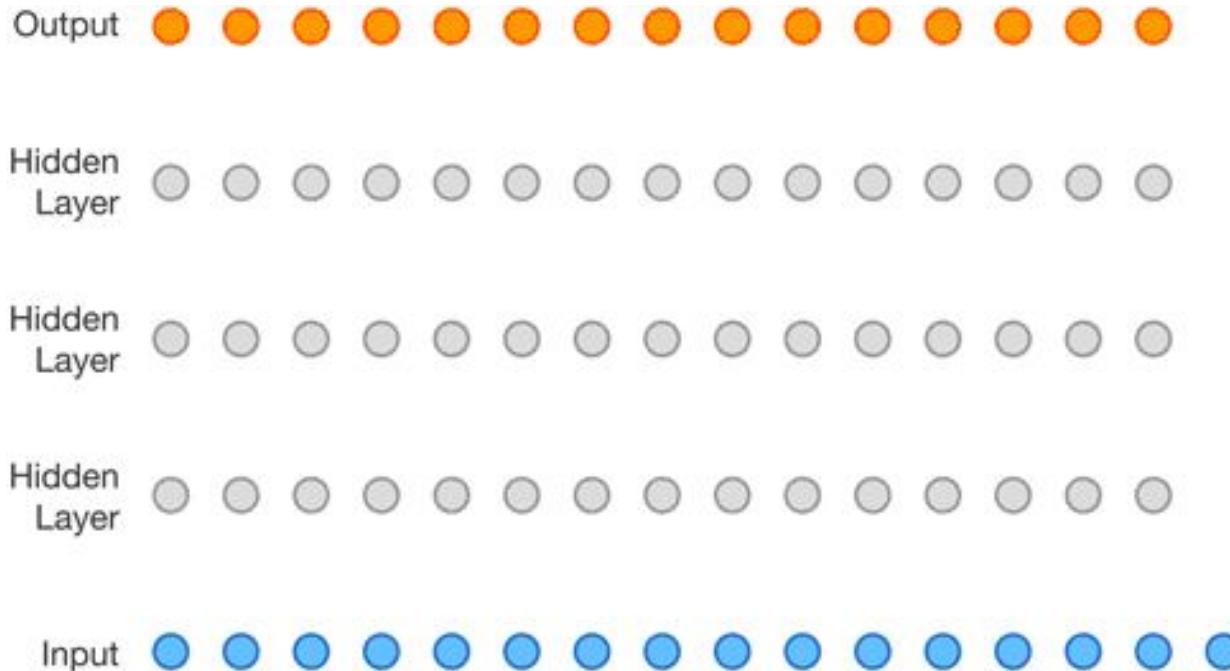
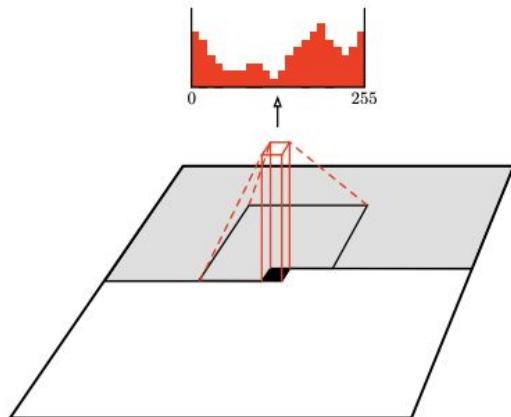


Image Generation

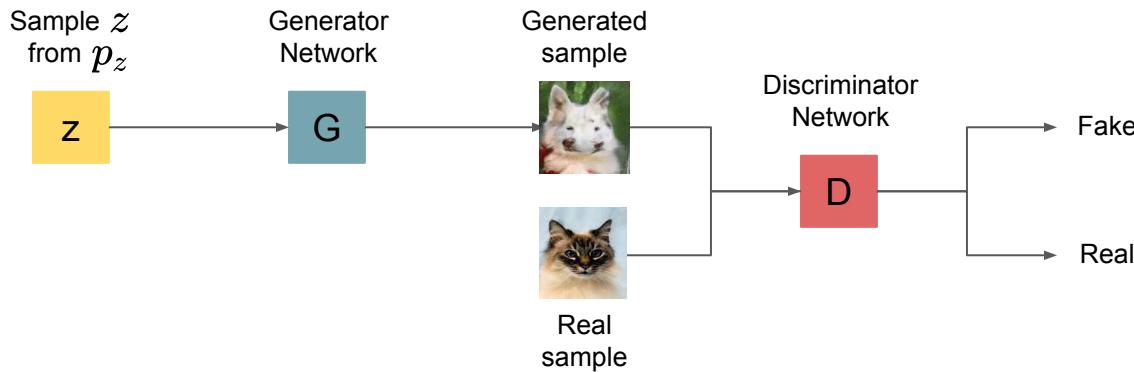
Autoregressive Models



Autoregressive Models: PixelCNN



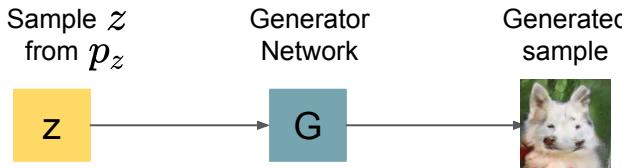
Generative Adversarial Networks (GAN)



$$\min_{\mathbf{G}} \max_{\mathbf{D}} (E_{x \sim p_{data}} [\log \mathbf{D}(x)] + E_{\mathbf{z} \sim p(\mathbf{z})} [\log(1 - \mathbf{D}(\mathbf{G}(\mathbf{z})))])$$

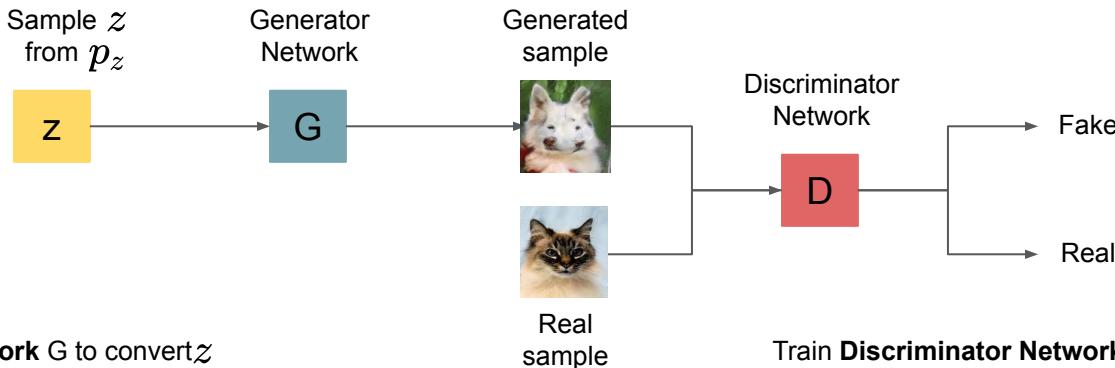
Generative Adversarial Networks

- Setup: Assume we have data x_i drawn from distribution $p_{data}(x)$. Want to sample from p_{data} .
- Idea: Introduce a latent variable z with simple prior $p(z)$.
- Sample $z \sim p(z)$ and pass to a Generator Network $x = G(z)$
- Then x is a sample from the Generator distribution p_G . Want $p_G = p_{data}$



Generative Adversarial Networks

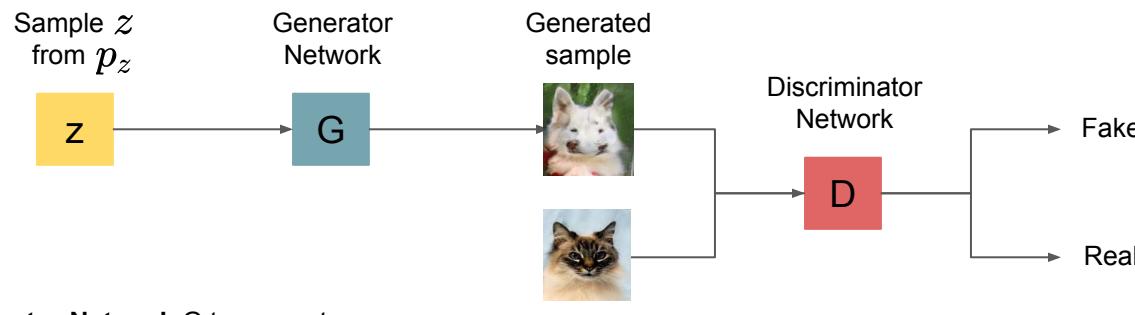
- Setup: Assume we have data x_i drawn from distribution $p_{data}(x)$. Want to sample from p_{data} .
- Idea: Introduce a latent variable z with simple prior $p(z)$.
- Sample $z \sim p(z)$ and pass to a Generator Network $x = G(z)$
- Then x is a sample from the Generator distribution p_G . Want $p_G = p_{data}$



Generative Adversarial Networks: Training Objective

Jointly train generator G and discriminator D with a **minimax game**

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



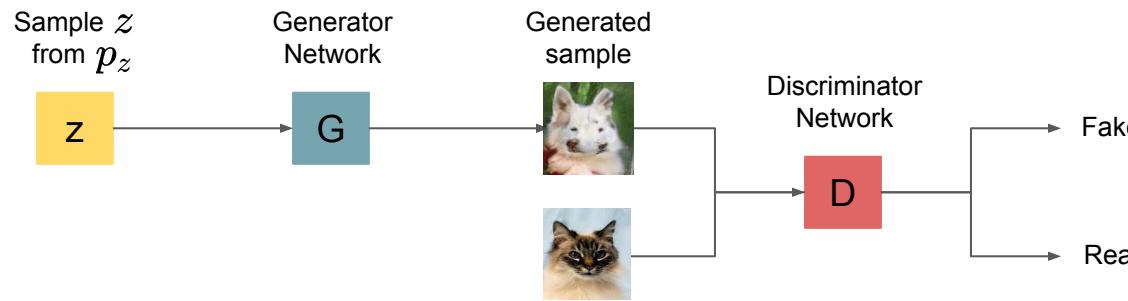
Train **Generator Network** G to convert z into fake data x sampled from p_G by fooling the Discriminator D

Train **Discriminator Network** D to classify data as real or fake (1/0)

Generative Adversarial Networks: Training Objective

Jointly train generator G and discriminator D with a **minimax game**

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



Train **Generator Network** G to convert z into fake data x sampled from p_G by fooling the Discriminator D

Train **Discriminator Network** D to classify data as real or fake (1/0)

Generative Adversarial Networks: Training Objective

Jointly train generator G and discriminator D with a **minimax game**

$$\min_{\mathbf{G}} \max_{\mathbf{D}} (E_{x \sim p_{data}} [\log \mathbf{D}(x)] + E_{\mathbf{z} \sim p(\mathbf{z})} [\log(1 - \mathbf{D}(\mathbf{G}(\mathbf{z})))])$$

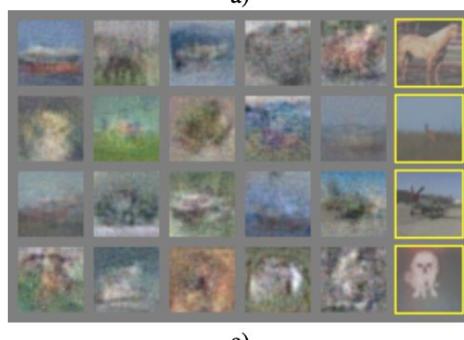
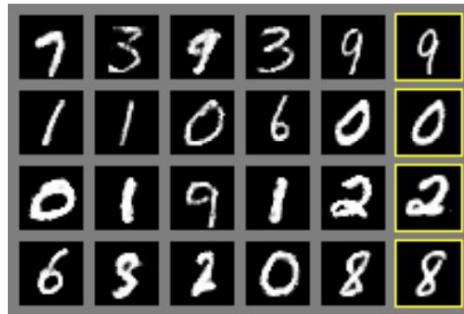
$$= \min_{\mathbf{G}} \max_{\mathbf{D}} \mathbf{V}(\mathbf{G}, \mathbf{D})$$

Train G and D using alternating gradient updates:

$$1. \text{ Update } \mathbf{D} = \mathbf{D} + \alpha_{\mathbf{D}} \frac{\delta \mathbf{V}}{\delta \mathbf{D}}$$

$$2. \text{ Update } \mathbf{G} = \mathbf{G} - \alpha_{\mathbf{G}} \frac{\delta \mathbf{V}}{\delta \mathbf{G}}$$

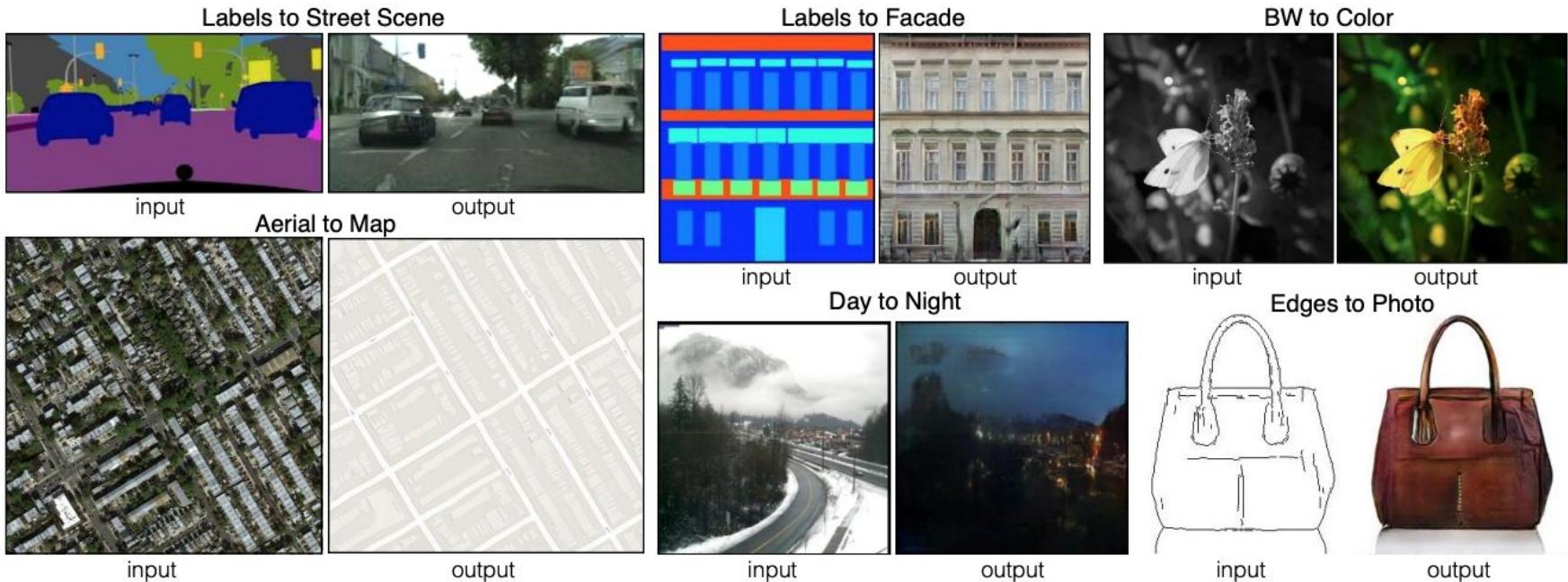
Generative Adversarial Networks: first results



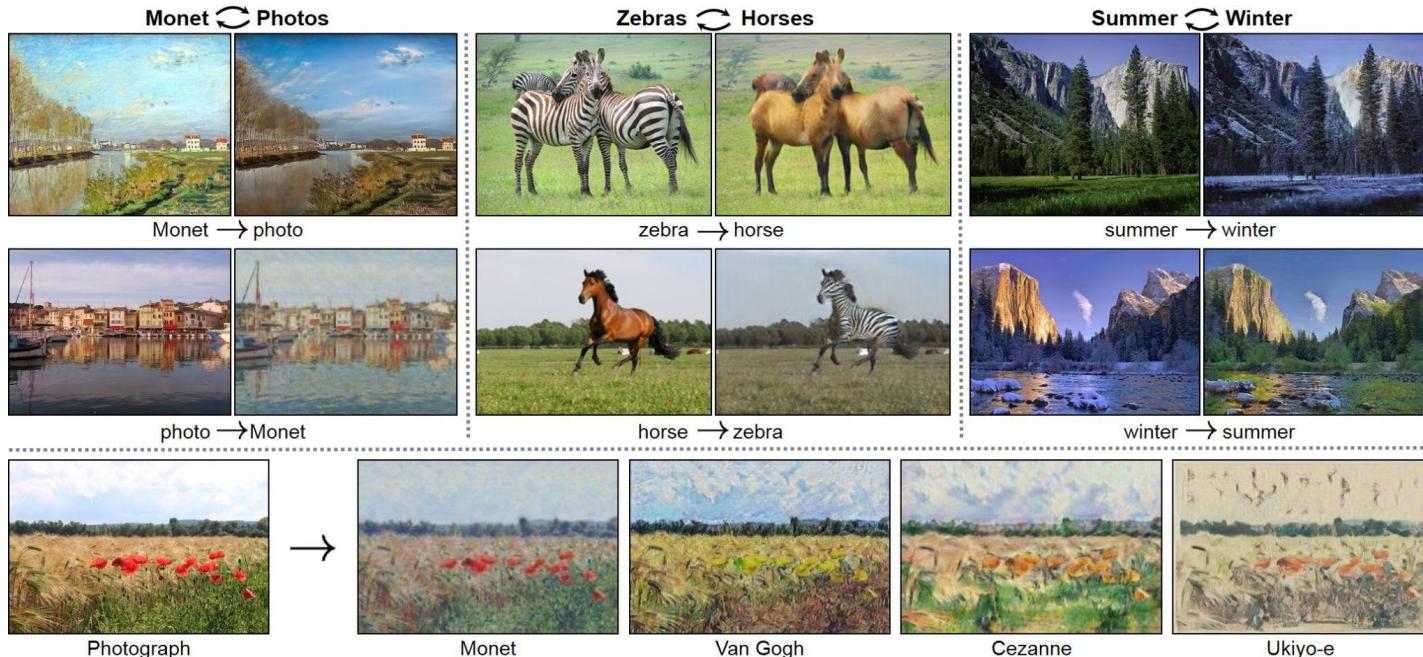
StyleGAN



Image-to-Image Translation: Pix2Pix



Unpaired Image-to-Image Translation: CycleGAN



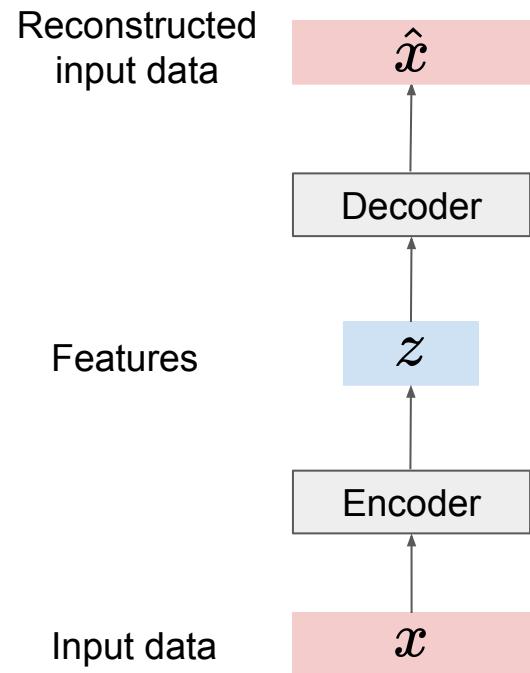
Unpaired Image-to-Image Translation: CycleGAN



Autoencoders (non-variational)

Unsupervised method for learning latent features from data without any labels.

$$L = \|\hat{x} - x\|_2^2$$

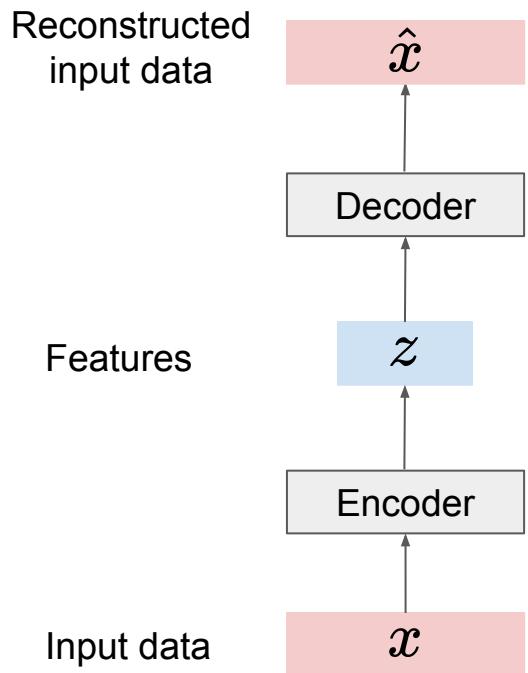


Autoencoders (non-variational)

Unsupervised method for learning latent features from data without any labels.

Features need to be **lower dimensional** than the data.

$$L = \|\hat{x} - x\|_2^2$$



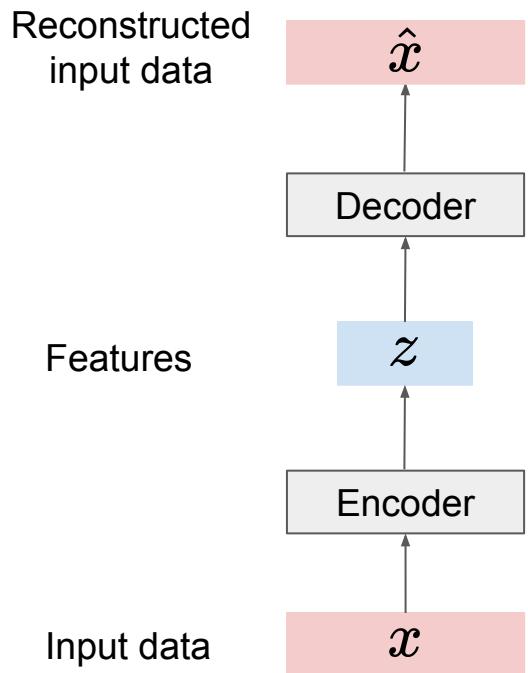
Autoencoders (non-variational)

Unsupervised method for learning latent features from data without any labels.

Features need to be **lower dimensional** than the data.

Limitation: no way to produce any new content

$$L = \|\hat{x} - x\|_2^2$$



Variational Autoencoders (VAE)

Add a probabilistic constraint between the encoder and decoder.

Variational Autoencoders

Add a probabilistic constraint between the encoder and decoder.

VAE is an autoencoder that learns **latent features** from data and enables **generative process**.

Variational Autoencoders

Add a probabilistic constraint between the encoder and decoder

VAE is an autoencoder that learns **latent features** from data and enables **generative process**.

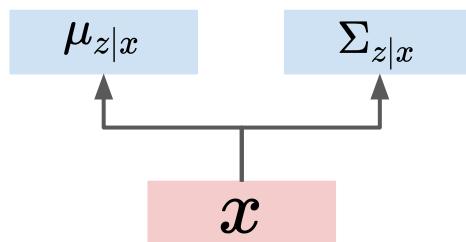
Instead of encoding an input as a single point, VAE encodes it as a distribution over the latent space.

Variational Autoencoders

Encoder network inputs data x and outputs distribution over latent codes z

Encoder Network

$$q_\phi(z|x) = N(\mu_{z|x}, \Sigma_{z|x})$$



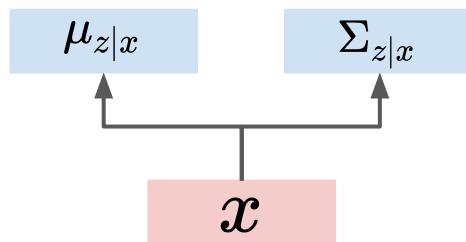
Variational Autoencoders

Encoder network inputs data x and outputs distribution over latent codes z

Decoder network inputs latent code z and outputs distribution over data x

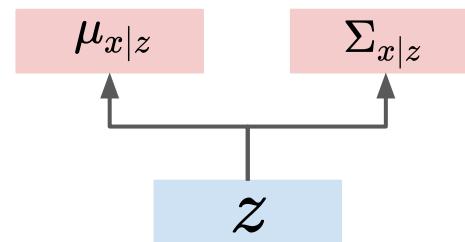
Encoder Network

$$q_\phi(z|x) = N(\mu_{z|x}, \Sigma_{z|x})$$



Decoder Network

$$p_\theta(x|z) = N(\mu_{x|z}, \Sigma_{x|z})$$



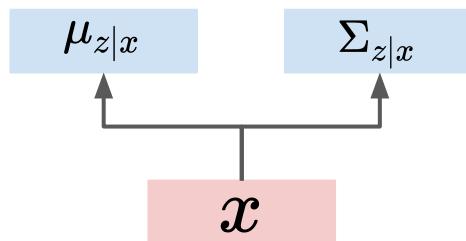
Variational Autoencoders

Jointly train **encoder** q and **decoder** p to maximize the **variational lower bound** on the data likelihood

$$\log p_\theta(x) \geq E_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - KL(q_\phi(z|x), p(z))$$

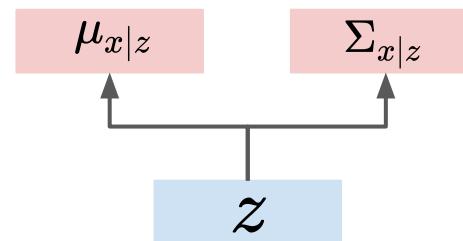
Encoder Network

$$q_\phi(z|x) = N(\mu_{z|x}, \Sigma_{z|x})$$



Decoder Network

$$p_\theta(x|z) = N(\mu_{x|z}, \Sigma_{x|z})$$

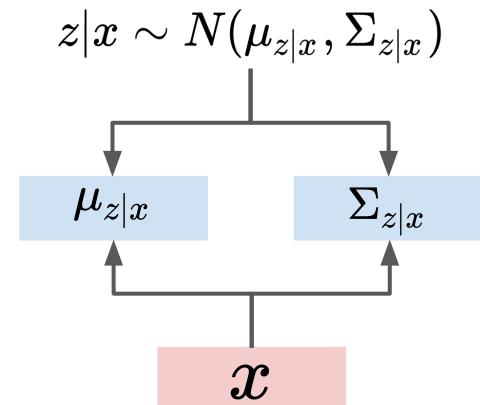


Variational Autoencoders

Train by maximize the **variational lower bound**.

$$E_{z \sim q_\phi(z|x)} [\log p_\Theta(x|z)] - KL(q_\phi(z|x), p(z))$$

1. The input is **encoded** as distribution over the latent space



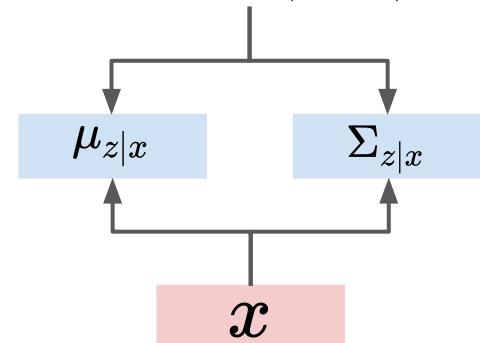
Variational Autoencoders

Train by maximize the **variational lower bound**.

$$E_{z \sim q_\phi(z|x)} [\log p_\Theta(x|z)] - KL(q_\phi(z|x), p(z))$$

1. The input is **encoded** as distribution over the latent space
2. **Encoder output should match prior $p(z)$**

$$z|x \sim N(\mu_{z|x}, \Sigma_{z|x})$$



Variational Autoencoders

Train by maximize the **variational lower bound**.

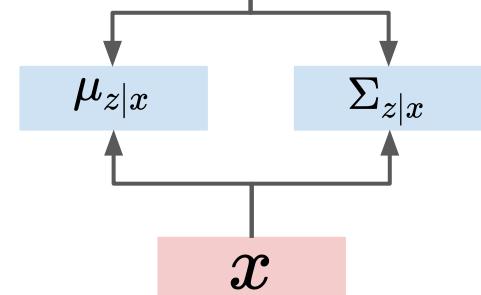
$$E_{z \sim q_\phi(z|x)} [\log p_\Theta(x|z)] - KL(q_\phi(z|x), p(z))$$

1. The input is **encoded** as distribution over the latent space
2. **Encoder output should match prior $p(z)$**
3. A point from the latent space is sampled from that distribution

z

Sample from z

$$z|x \sim N(\mu_{z|x}, \Sigma_{z|x})$$

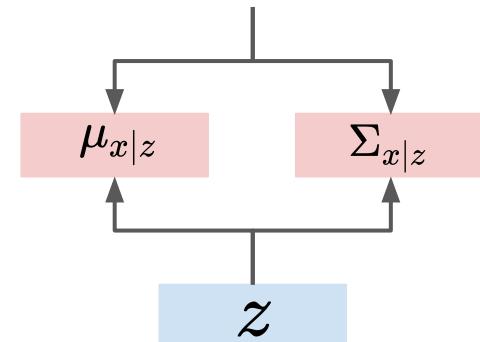


Variational Autoencoders $x|z \sim N(\mu_{x|z}, \Sigma_{x|z})$

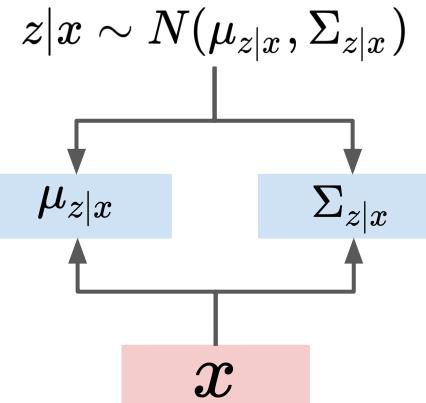
Train by maximize the **variational lower bound**.

$$E_{z \sim q_\phi(z|x)} [\log p_\Theta(x|z)] - KL(q_\phi(z|x), p(z))$$

1. The input is **encoded** as distribution over the latent space
2. **Encoder output should match prior $p(z)$**
3. A point from the latent space is sampled from that distribution
4. The sampled point is **decoded**



Sample from z

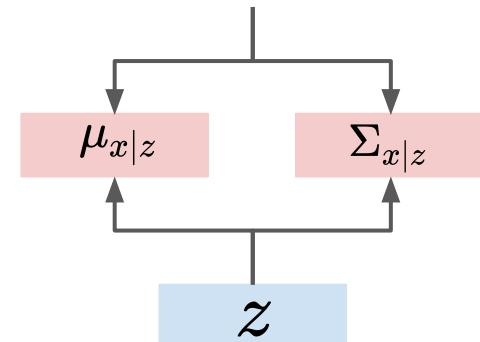


Variational Autoencoders $x|z \sim N(\mu_{x|z}, \Sigma_{x|z})$

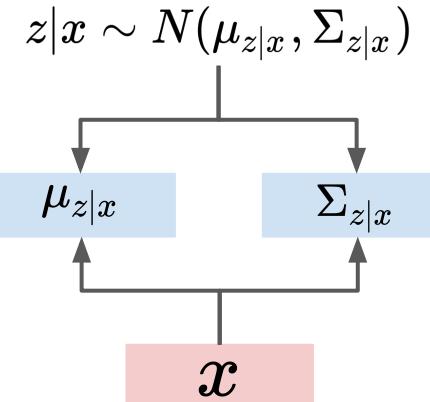
Train by maximize the **variational lower bound**.

$$E_{z \sim q_\phi(z|x)} [\log p_\Theta(x|z)] - KL(q_\phi(z|x), p(z))$$

1. The input is **encoded** as distribution over the latent space
2. **Encoder output should match prior $p(z)$**
3. A point from the latent space is sampled from that distribution
4. The sampled point is **decoded**
5. **The reconstruction error is computed**



Sample from z



Variational Autoencoders: Results

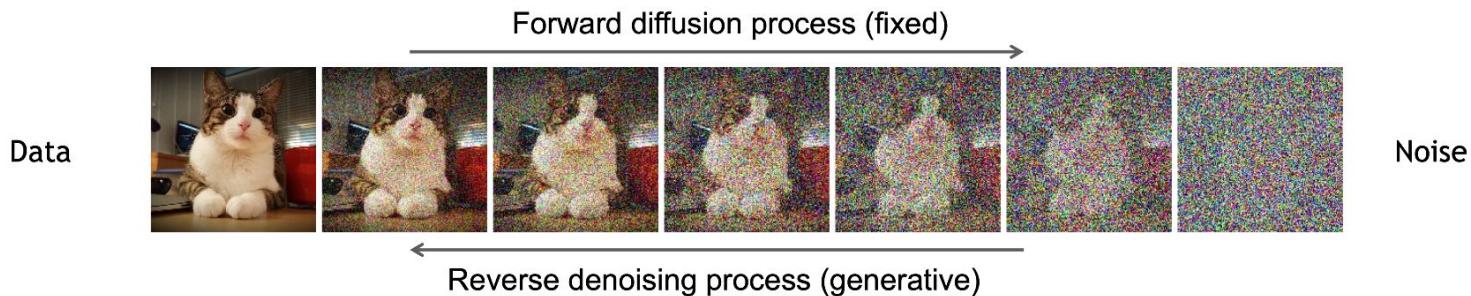


Denoising Diffusion Models

Learning to generate by denoising

Denoising diffusion models consist of two processes:

- Forward diffusion process that gradually adds noise to input
- Reverse denoising process that learns to generate data by denoising

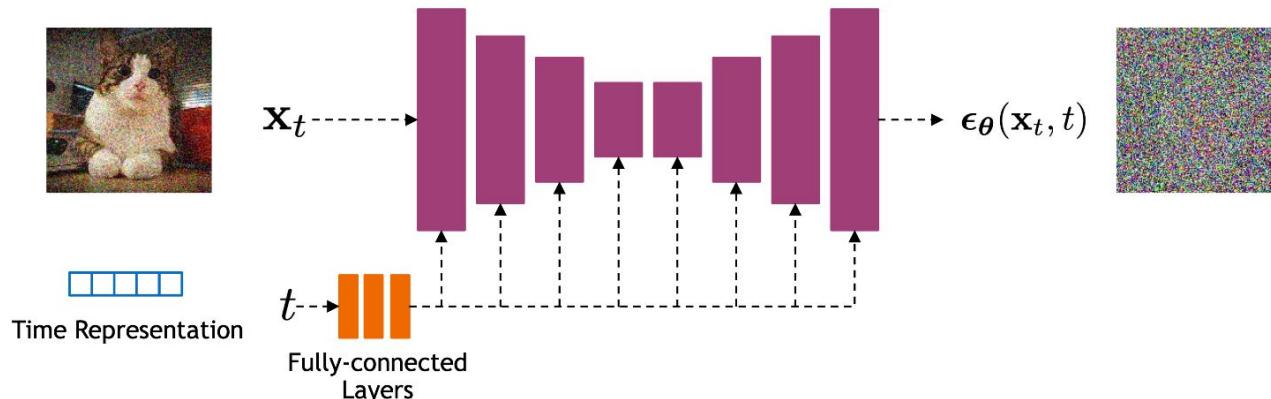


[Sohl-Dickstein et al., Deep Unsupervised Learning using Nonequilibrium Thermodynamics, ICML 2015](#)
[Ho et al., Denoising Diffusion Probabilistic Models, NeurIPS 2020](#)
[Song et al., Score-Based Generative Modeling through Stochastic Differential Equations, ICLR 2021](#)

Implementation Considerations

Network Architectures

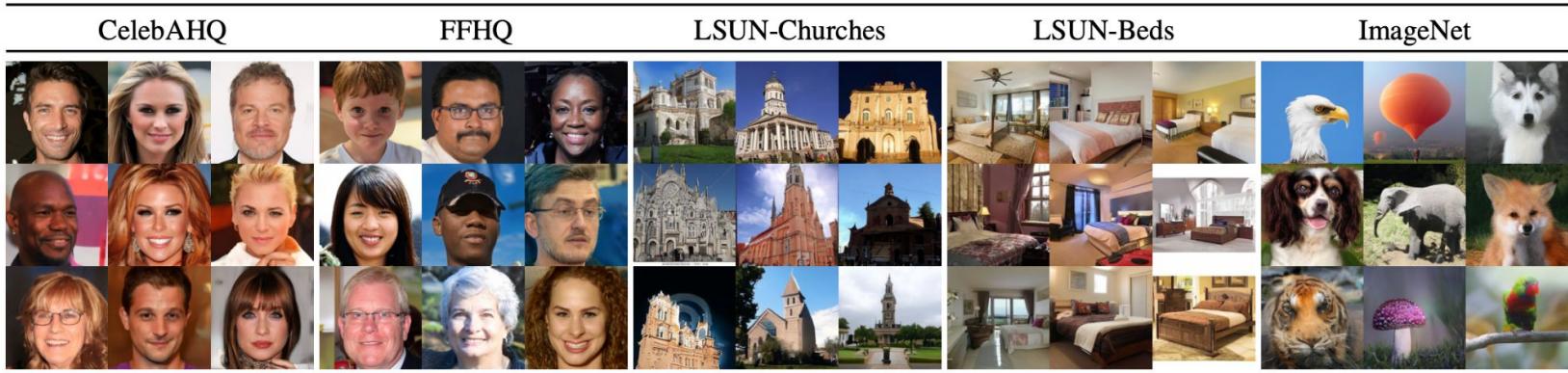
Diffusion models often use U-Net architectures with ResNet blocks and self-attention layers to represent $\epsilon_\theta(\mathbf{x}_t, t)$



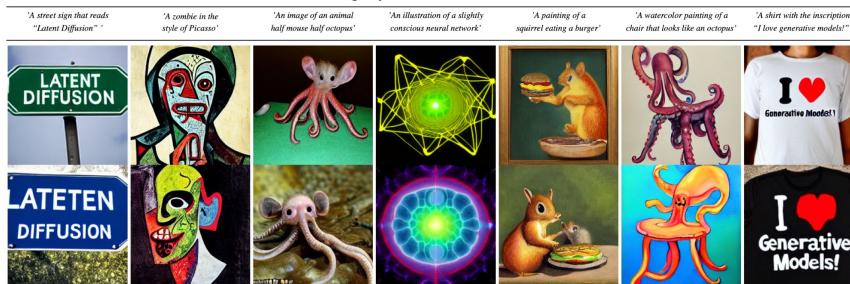
Time representation: sinusoidal positional embeddings or random Fourier features.

Time features are fed to the residual blocks using either simple spatial addition or using adaptive group normalization layers. (see [Dhariwal and Nichol NeurIPS 2021](#))

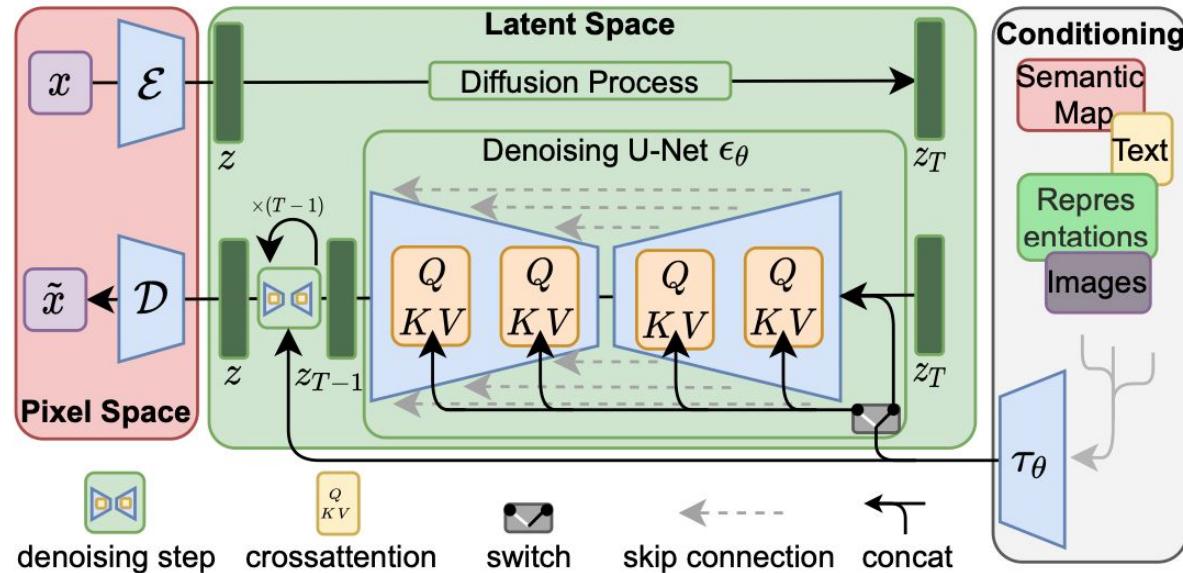
Latent Diffusion Models



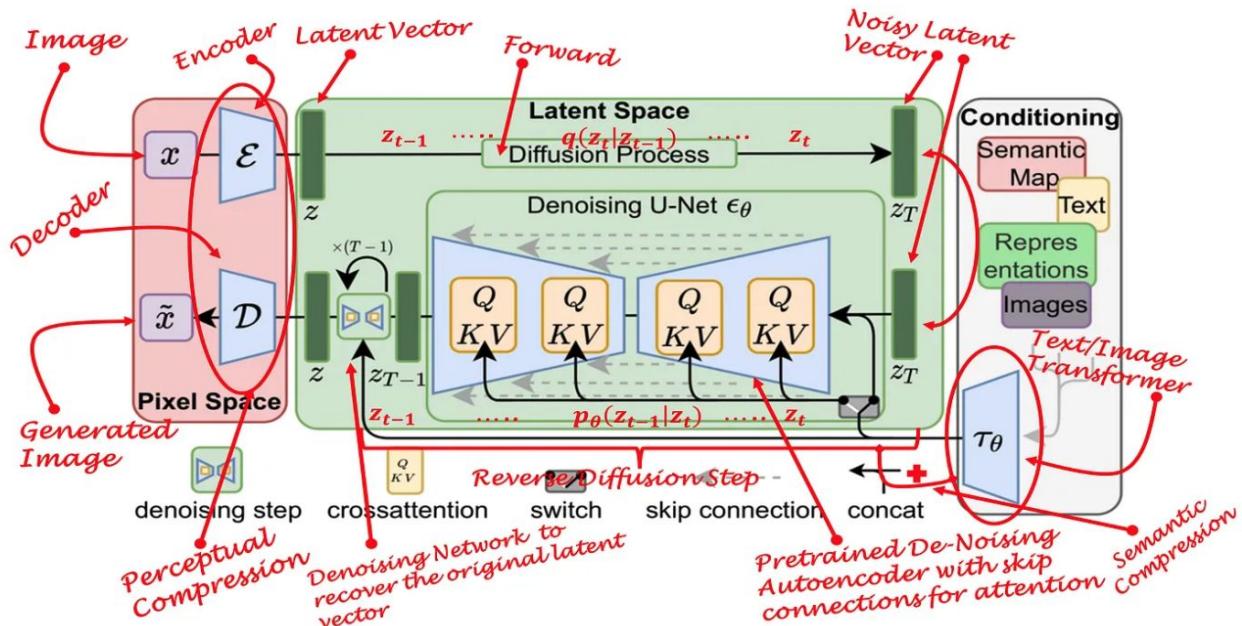
Text-to-Image Synthesis on LAION, 1.45B Model.



Latent Diffusion Models



Latent Diffusion Models



Generative models evaluation

Generative Models Evaluation

Which of these images looks better?



Generative Models Evaluation

Which of these images looks more realistic?



Generative Models Evaluation

Which of these images appears to be more similar to the text prompt?



Prompt: The saying "BE EXCELLENT TO EACH OTHER" written on a red brick wall with a graffiti image of a green alien wearing a tuxedo. A yellow fire hydrant is on a sidewalk in the foreground.

Generative Models Evaluation

- Human-based ratings and preference judgments
- Inception Score (IS) [1]
- Frechet Inception Distance (FID) [2]

[1] Salimans et al. Improved Techniques for Training GANs. NeurIPS 2016

[2] Heusel et al. GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. NeurIPS 2017

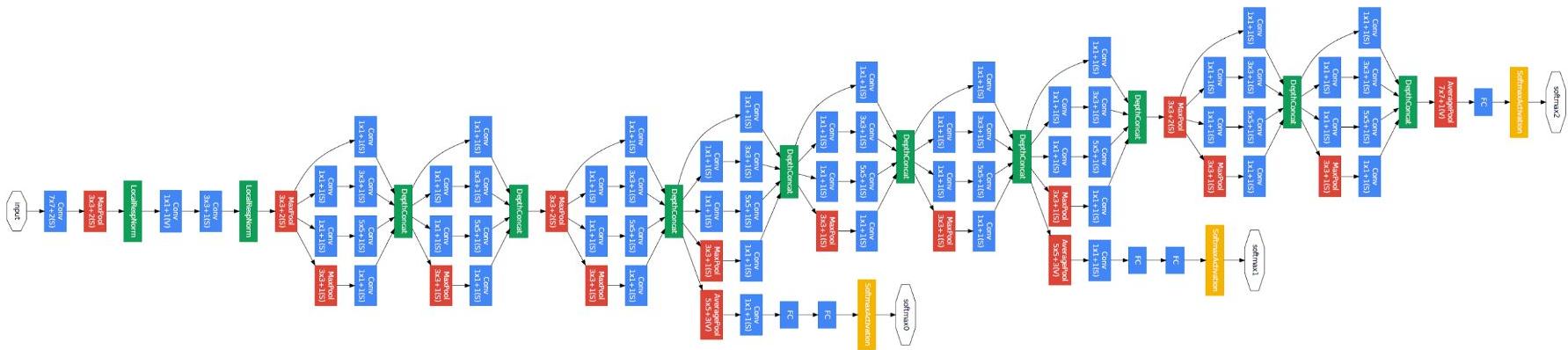
Inception Score (IS)

IS measures:

- the **quality** of the generated images
- their **diversity**

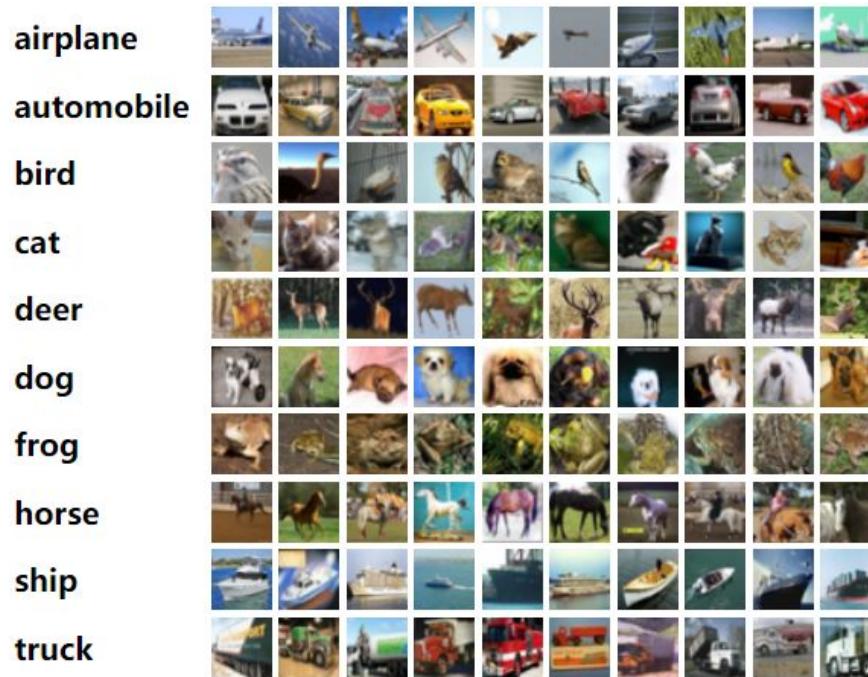
Inception Score

Inception image classifier pre-trained on CIFAR10



Inception Score

CIFAR10 dataset:

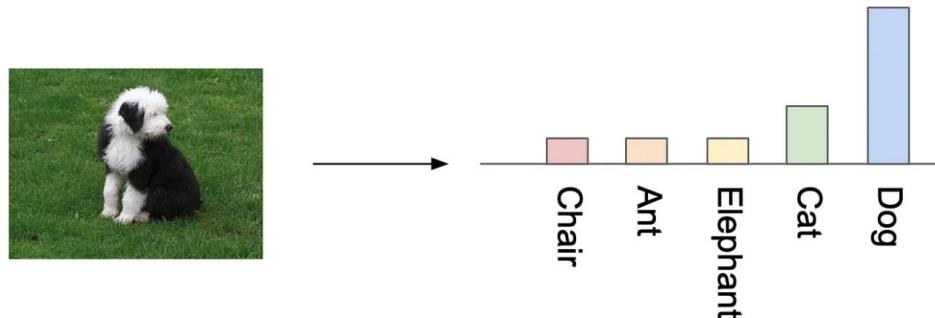


Inception Score

Generated images are passed through a pre-trained image classifier (e.g. Inception network trained on CIFAR-10).

The classifier predicts the label distribution $p(y|x)$, where:

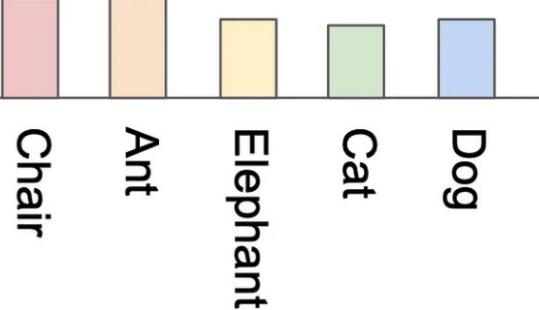
- x : generated image
- y : predicted class label



Inception Score: Low Confidence Example

If the predicted label distribution is **flat or uncertain**, the image is likely **ambiguous or low quality**.

- ✓ Low confidence → classifier spreads probability across many classes
- ✗ This leads to a **low Inception Score**

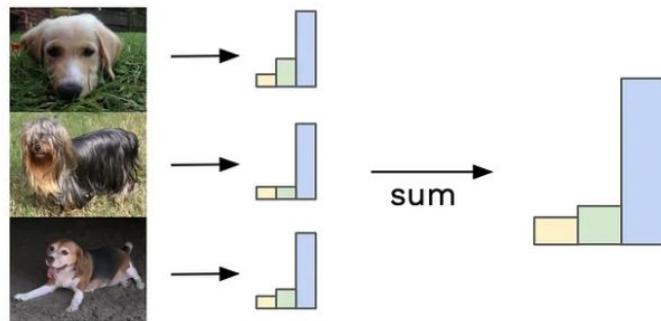


Inception Score: Measuring Diversity with Marginal Distribution

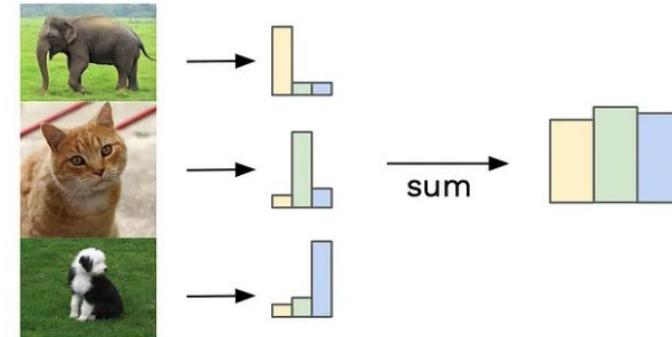
Calculate marginal probability $p(y) = \int_z p(y|x=G(z))dz$

The marginal $p(y)$ shows how diverse the generated samples are.

Similar labels sum to give focussed distribution



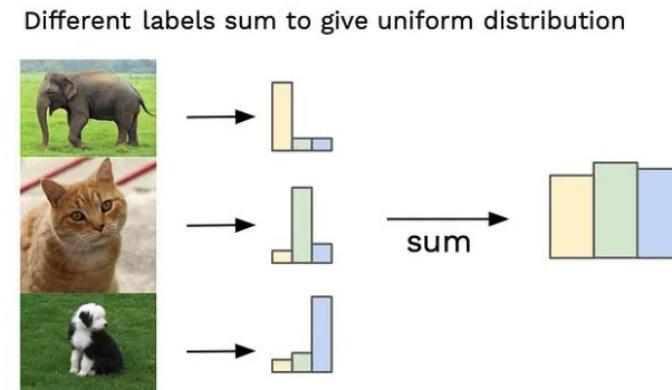
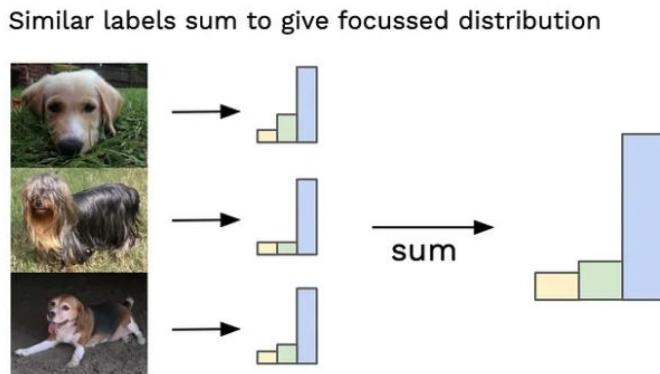
Different labels sum to give uniform distribution



Inception Score: Measuring Diversity with Marginal Distribution

The marginal probability tells us how much variety the generator produces.

- **Low diversity:** similar predictions → focused marginal
- **High diversity:** varied predictions → uniform marginal



Inception Score

- **Quality:** conditional probability $p(y|x)$
- **Diversity:** marginal probability $p(y)$

We want

- the conditional probability $p(y|x)$ to be highly predictable (**low entropy**) i.e. given an image, we should know the object type easily.
- the marginal probability $p(y)$ to be uniform (**high entropy**).

Inception Score

Compute their KL-divergence to combine these two criteria:

$$IS(G) = \exp(E_{x \sim p_g} KL(p(y|x) || p(y)))$$

A **high Inception Score** means:

- Clear, confident predictions → **high image quality**
- Diverse labels across samples → **high diversity**

Frechet Inception Distance (FID)

- Use the **Inception network** to extract features from an intermediate layer

Frechet Inception Distance (FID)

- Use the **Inception network** to extract features from an intermediate layer
- Model data distribution for these features using a multivariate Gaussian distribution with mean μ and covariance Σ

Frechet Inception Distance (FID)

- Use the **Inception network** to extract features from an intermediate layer
- Model data distribution for these features using a multivariate Gaussian distribution with mean μ and covariance Σ
- The FID between the real images x and generated images g :

$$FID(x, g) = \|\mu_x - \mu_g\|_2^2 + Tr(\Sigma_x + \Sigma_g - 2(\Sigma_x \Sigma_g)^{\frac{1}{2}})$$

where Tr sums up all the diagonal elements

Frechet Inception Distance

- **Lower FID values mean better** image quality and diversity
- FID is sensitive to mode collapse, the distance increases when modes are missed
- FID is more robust to noise than IS. If the model only generates one image per class, the distance will be high

A landscape photograph of a sunset or sunrise over a range of mountains. The sky is filled with warm colors, transitioning from deep orange and yellow at the horizon to a lighter, pinkish-purple higher up. The mountains in the foreground are dark, silhouetted against the bright sky. In the center of the image, the words "Thank you" are written in a large, black, sans-serif font.

Thank you