# Diffusion

Presented By: Massa Baali

Carnegie Mellon University

# Discriminative vs. Generative Models

- **Discriminative models learn to discriminate**

  - Determine the class given the input

    - Compute P(y|x)

- **Generative models can generate**

  - Produce more instances like the training data
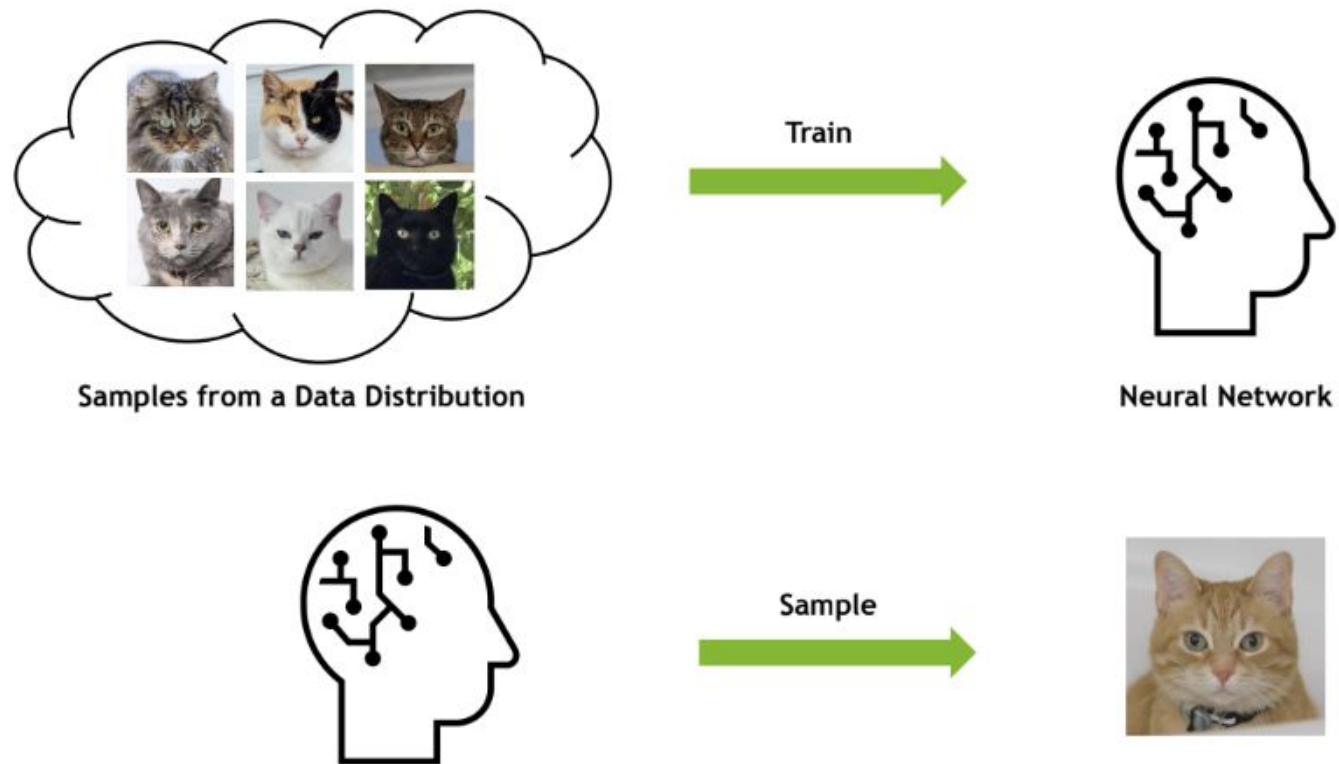
    - Compute and/or draw from P(x,y)

Carnegie Mellon University

# Discriminative vs. Generative Models

**Given a distribution of inputs X and labels Y**

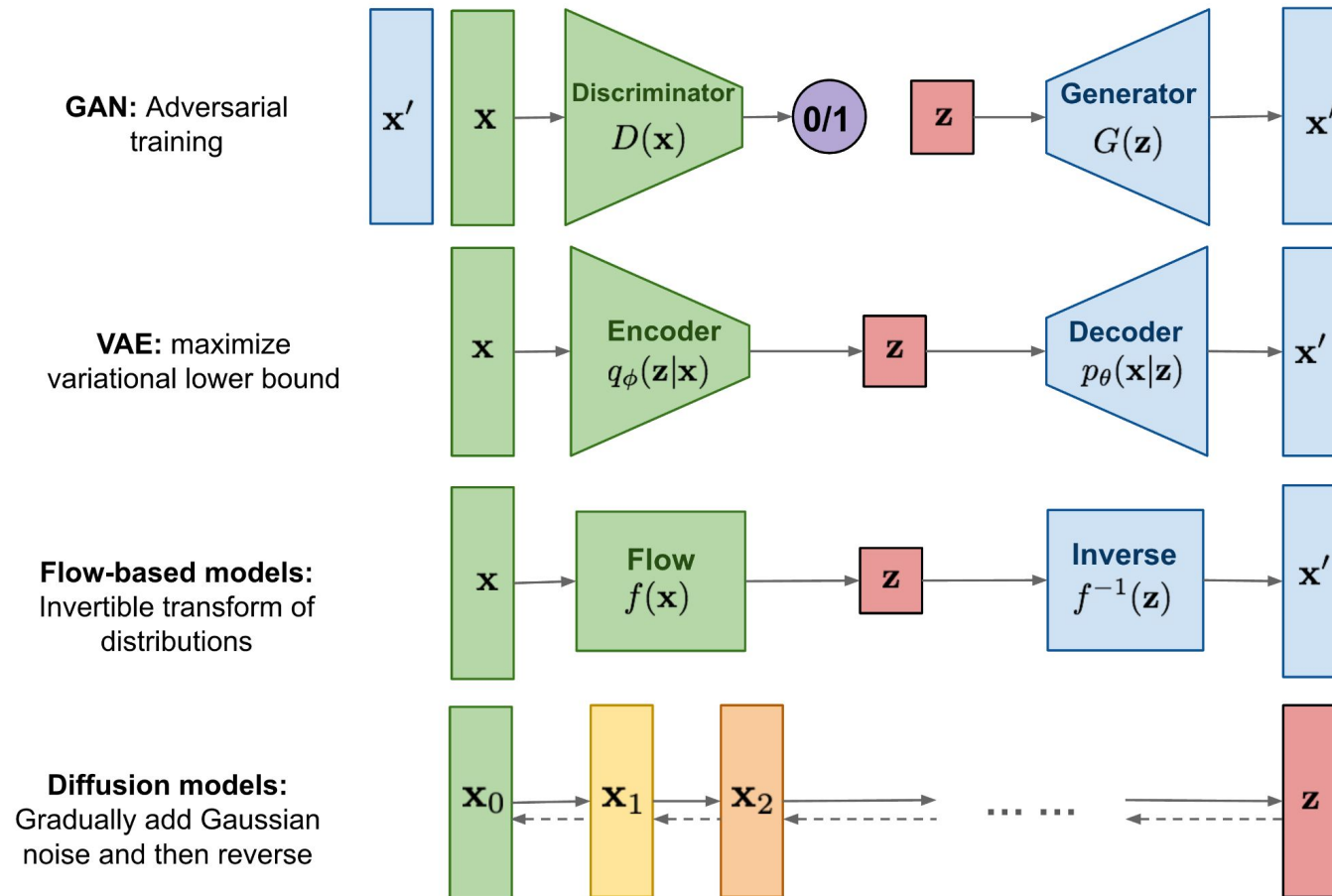| Discriminative models | Generative models |
|---|---|
| Discriminative models learn conditional distribution P(Y \| X) | Generative models learn the joint distribution P(Y, X) |
| Learns decision boundary between classes. | Learns actual probability distribution of data. |
| Limited scope. Can only be used for classification tasks. | Can do both generative and discriminative tasks. |

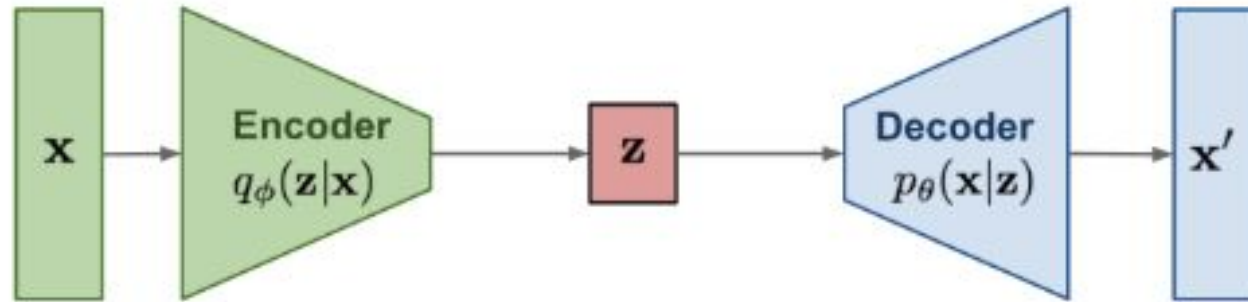Harder problem, requires a deeper understanding of the distribution than discriminative models.

Carnegie Mellon University

# Deep Generative learning

**Learning to generate data**



**Slide credit to: https://cvpr2022-tutorial-diffusion-models.github.io/**

Carnegie Mellon University

# Generative Models



**GAN:** Adversarial training

**VAE:** maximize variational lower bound

**Flow-based models:** Invertible transform of distributions

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

https://lilianweng.github.io/posts/2021-07-11-diffusion-models/
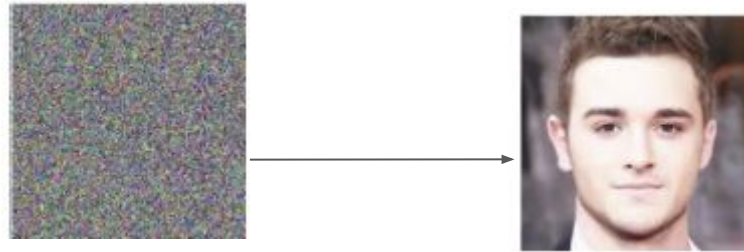
**Carnegie Mellon University**

# VAE (Recap)



A generative model that learns to create new samples by:

- Compressing samples into a small code (*latent space*)

- Reconstructing samples from these codes

https://lilianweng.github.io/posts/2021-07-11-diffusion-models/
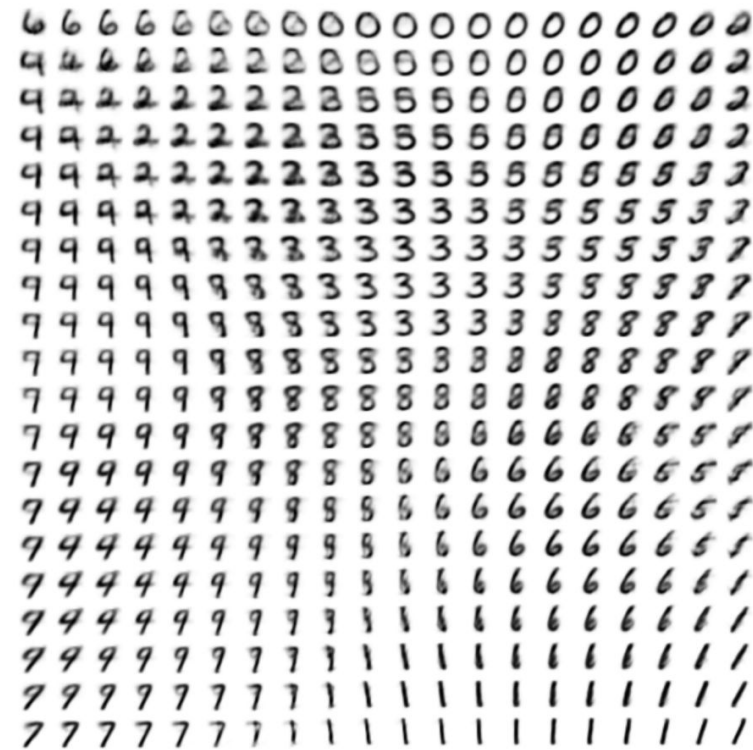
**Carnegie Mellon University**

# Limitations of VAEs



The decoder must transform a standard Gaussian *all the way* to the target distribution in **one step**.

- Too large a gap to bridge in one step
- The decoder has to do ALL the work at once

Carnegie Mellon University
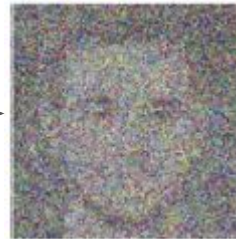
# Limitations of VAEs



Result: Blurry, Low-Quality Images

Kingma et al. Auto-Encoding Variational Bayes

**Carnegie Mellon University**

# From Single-Step to Multi-Step Generation



VAE

Diffusion

**Carnegie Mellon University**

# Hierarchicals VAEs

- Stack multiple VAEs with intermediate latent variables
- Each level: $x_0 \rightarrow x_1 \rightarrow x_2 \rightarrow \ldots \rightarrow x_t$
- Each step is a small, learnable transformation

**Diffusion Models: A Special Case**

- Can be viewed as a hierarchical VAE with specific choices:
  - All latent variables have the same dimension as the original image
  - Forward process (encoder) is fixed (just adds Gaussian noise)
  - Only learn the reverse process (decoder/denoising)
  - Many steps (e.g., T = 1000)

**Carnegie Mellon University**

# Poll 1

Why do VAEs produce blurry images?

- A) Not enough training data
- B) Bad optimizer choice
- C) One-step jump from noise to image is too hard
- D) Network architecture is too small

**Carnegie Mellon University**

# Poll 1

Why do VAEs produce blurry images?

- A) Not enough training data
- B) Bad optimizer choice
- C) One-step jump from noise to image is too hard
- D) Network architecture is too small

**Carnegie Mellon University**

# Outline

- Definition of Diffusion

- Importance of Diffusion

- Explaining the process of diffusion

- Denoising Diffusion Implicit Models (DDIM)

- Diffusion Models from Stochastic Differential Equations and Score
  Matching Perspective

- Classifier-free Guidance (CFG)

- Performance Metrics: FID, IS, Precision and Recall

- Applications of Diffusion Models
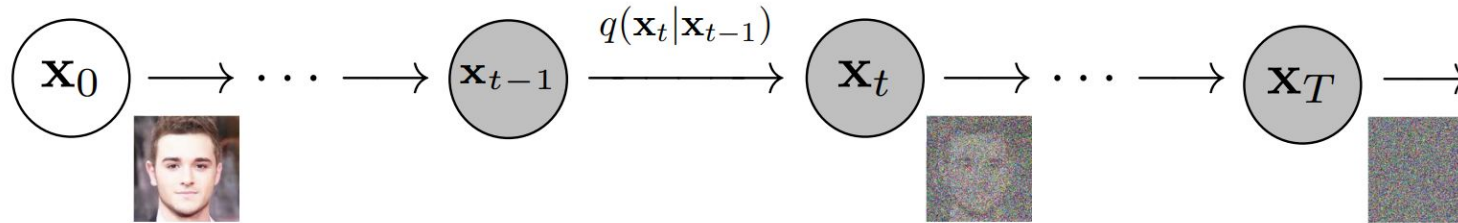
**Carnegie Mellon University**

# Definition



Adapted from: Ho, J., Jain, A. and Abbeel, P., 2020. Denoising
diffusion probabilistic models. Advances in neural information
processing systems, 33, pp.6840-6851.

Carnegie Mellon University

# Importance of Diffusion Models

- **High-Quality Outputs:** Exceptional generation of images, audio, and more with high fidelity.
- **Stable Training:** More consistent training outcomes compared to other models like GANs.
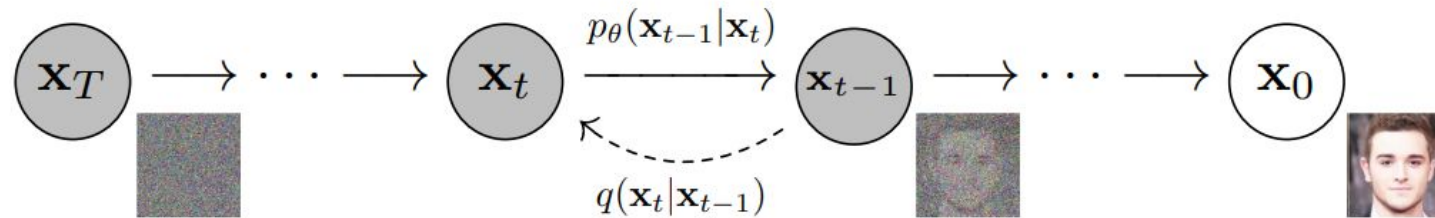
**Carnegie Mellon University**

# Forward Process



- **Gradual Noise Addition:** Over t time steps, noise is progressively added to an image.
- **Markov Chain Model:** The process, denoted by q takes the form of Markov Chain where the distribution at a particular time steps only depends on the sample from the immediate previous step.

**Carnegie Mellon University**

# Reverse Process



- **Denoise:** The model learns to gradually remove noise from the data at each step.

- **Reconstruction:** Aims to reconstruct the original image or data from its noisy version.

- The main goal of training a diffusion model is learning the reverse process specifically training pθ(xt−1|xt).

**Carnegie Mellon University**

# How do we add noise?

- **Gaussian Noise:** The noise added in diffusion models follows a Gaussian (normal) distribution. This means at each diffusion step, we inject some random noise that has the familiar "bell curve" distribution.
- **Standard Normal Sampling:** We sample the noise from a standard normal distribution $\mathcal{N}(0, \boldsymbol{I})$, which has a mean of 0 and a variance of 1.

**Carnegie Mellon University**

# What Does the Model Actually Learn?

**The Task:**

- Input: Noisy image $x_t$, timestep t
- Output: Predict the noise $\varepsilon$ that was added

**Why This Works:**

- If we know the noise, we can subtract it
- Subtracting noise = denoising = moving toward clean data



**DIFFUSION MODEL TRAINING LOOP**

1. Take Real Image $x_0$ — CLEAN IMAGE — $x_0 \sim q(x_0)$
2. Add Known Noise $\varepsilon$ at time $t$ — $x_t = \sqrt{\alpha}\,x_0 + \sqrt{(1-\bar{\alpha}_c)}\,\varepsilon$
3. Train Network to Predict $\varepsilon$ — NETWORK → $\varepsilon$ — Loss $= \|\varepsilon - \dot{\varepsilon}\|^2$ — Backpropagate

**Carnegie Mellon University**

# Training Objective

- **Minimize Error:** The goal is to minimize the Mean Squared Error (MSE) between the model's predicted noise and the actual noise added during the forward process.
- **High-Quality Reconstruction:** Effective training enables the model to accurately reconstruct original data from noise, enhancing its capability to generate high-quality samples.

$$L_{\text{simple}}(\theta) := \mathbb{E}_{t,\mathbf{x}_0,\boldsymbol{\epsilon}} \left[ \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta \left( \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t}\boldsymbol{\epsilon}, t \right) \right\|^2 \right]$$

**Carnegie Mellon University**

# Training

**Algorithm 1** Training

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

Carnegie Mellon University

# Sampling

---

**Algorithm 2** Sampling

1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
2: **for** $t = T, \ldots, 1$ **do**
3:   $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = \mathbf{0}$
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** $\mathbf{x}_0$

---

# Poll 2

What does the neural network actually predict during training?

- A) The clean image $x_0$
- B) The noisy image $x_t$
- C) The noise $\varepsilon$ that was added
- D) The next timestep $t+1$

**Carnegie Mellon University**

# Poll 2

What does the neural network actually predict during training?

- A) The clean image $x_0$
- B) The noisy image $x_t$
- C) The noise $\varepsilon$ that was added
- D) The next timestep t+1

**Carnegie Mellon University**

# Outline

- Definition of Diffusion

- Importance of Diffusion

- Explaining the process of diffusion

- Denoising Diffusion Implicit Models (DDIM)

- Diffusion Models from Stochastic Differential Equations and Score Matching Perspective

- Classifier-free Guidance (CFG)

- Performance Metrics: FID, IS, Precision and Recall
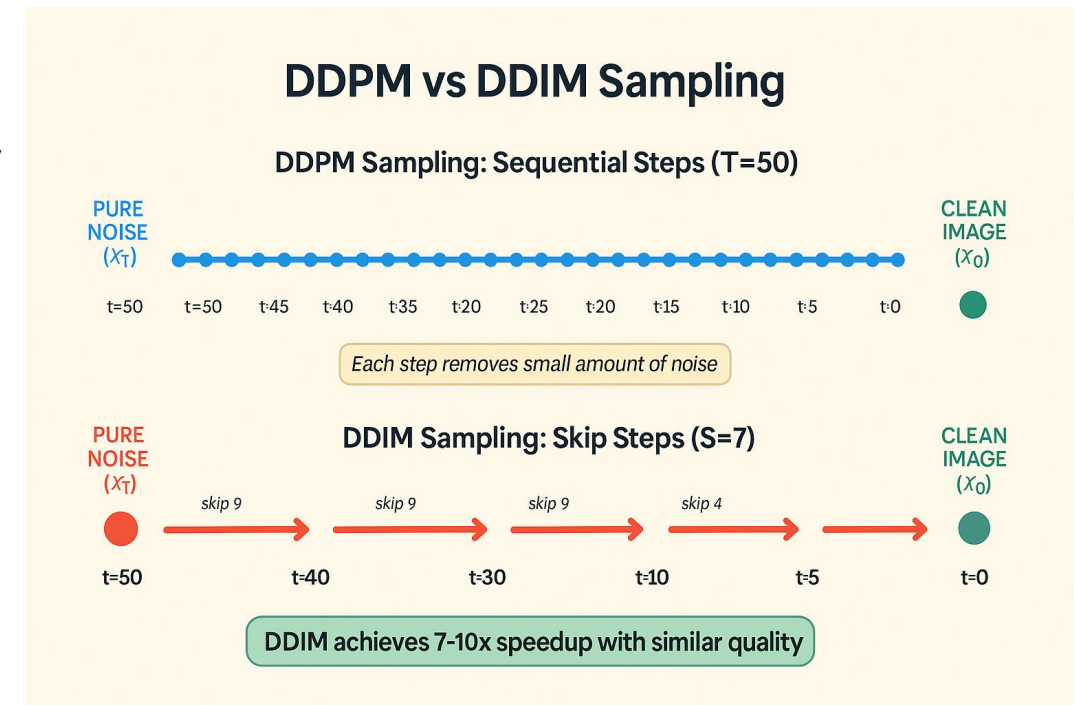
- Applications of Diffusion Models

**Carnegie Mellon University**

# The Speed Problem

**DDPM Challenge:**

- Must run model T=1000 times sequentially
- Each step removes tiny amount of noise
- Generation takes minutes per image

**DDIM Solution:**

- Skip timesteps intelligently
- Use 10-50 steps instead of 1000
- 10-100× faster with similar quality



DDPM vs DDIM Sampling

DDPM Sampling: Sequential Steps (T=50)

PURE NOISE ($x_T$)

CLEAN IMAGE ($x_0$)

t=50  t=50  t=45  t=40  t=35  t=20  t=25  t=20  t=15  t=10  t=5  t=0

*Each step removes small amount of noise*

DDIM Sampling: Skip Steps (S=7)

PURE NOISE ($x_T$)

CLEAN IMAGE ($x_0$)

skip 9   skip 9   skip 9   skip 4

t=50   t=40   t=30   t=10   t=5   t=0

DDIM achieves 7-10x speedup with similar quality

**Carnegie Mellon University**

# Denoising Diffusion Implicit Models (DDIM)

**Limitations of DDPM Inference**

**Sequential Denoising:** Must process each timestep in reverse to remove Gaussian noise. This happens due to the markov chain structure of the reverse process.

**Algorithm 2** Sampling

1: $x_T \sim \mathcal{N}(0, I)$         $\triangleright$*Initial isotropic gaussian noise sampling*
2: **for** t = T, ..., 1 **do**
3:     $z \sim \mathcal{N}(0, I)$ if $t > 1$ else $z = 0$     $\triangleright$*Sample random noise (if not last step)*
4:     $\tilde{\epsilon} = \epsilon_\theta(x_t, t)$     $\triangleright$*Estimated noise in current noisy data*
5:     $\tilde{x}_0 = \frac{1}{\sqrt{\bar{\alpha}_t}}(x_t - \sqrt{1 - \bar{\alpha}_t}\tilde{\epsilon})$     $\triangleright$*Estimated $x_0$ from estimated noise*
6:     $\tilde{\mu} = \mu_t(x_t, \tilde{x}_0) \left(= \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\epsilon_\theta(x_t, t)\right)\right)$     $\triangleright$*Mean for previous step sampling*
7:     $x_{t-1} = \tilde{\mu} + \sigma_t z$     $\triangleright$*Previous step sampling*
8: **end for**
9: **return** $x_0$

**Carnegie Mellon University**

# DDIM Overview

**Random Sampling:** In DDPM the inference step involves stochastic (random) sampling to reverse the diffusion process.

$$\tilde{x}_{t-1} = \mu_\theta(\tilde{x}_t, t) + \sigma_t(z_t - \epsilon_\theta(\tilde{x}_t, t))$$

noise schedule

estimated noise

denoising function

**Carnegie Mellon University**

# DDIM Overview

**Deterministic Sampling:** DDIM uses a non-Markovian process allowing for fewer timesteps.It modifies the inference step by modifying the reverse process, making the process deterministic.

$$\boldsymbol{x}_{t-1} = \sqrt{\alpha_{t-1}} \underbrace{\left( \frac{\boldsymbol{x}_t - \sqrt{1-\alpha_t}\,\epsilon_\theta^{(t)}(\boldsymbol{x}_t)}{\sqrt{\alpha_t}} \right)}_{\text{`` predicted } \boldsymbol{x}_0\text{''}} + \underbrace{\sqrt{1-\alpha_{t-1}-\sigma_t^2} \cdot \epsilon_\theta^{(t)}(\boldsymbol{x}_t)}_{\text{``direction pointing to } \boldsymbol{x}_t\text{''}} + \underbrace{\sigma_t \epsilon_t}_{\text{random noise}}$$

**Carnegie Mellon University**

# Poll 3

DDIM is 10-50× faster than DDPM because:

- A) It uses a smaller neural network
- B) It skips timesteps intelligently
- C) It trains faster
- D) It uses a better GPU

**Carnegie Mellon University**

# Poll 3

DDIM is 10-50× faster than DDPM because:

- A) It uses a smaller neural network
- B) It skips timesteps intelligently
- C) It trains faster
- D) It uses a better GPU

**Carnegie Mellon University**

# Outline

- Definition of Diffusion

- Importance of Diffusion

- Explaining the process of diffusion

- Denoising Diffusion Implicit Models (DDIM)

- Diffusion Models from Stochastic Differential Equations and Score Matching Perspective

- Classifier-free Guidance (CFG)

- Performance Metrics: FID, IS, Precision and Recall

- Applications of Diffusion Models

**Carnegie Mellon University**

# Why Move to Continuous Time?

**Discrete View (DDPM):**

- Fixed timesteps: t = 1, 2, 3, ..., 1000
- Fixed noise schedule: $\beta_1$, $\beta_2$, ..., $\beta_{1000}$
- Limited flexibility

**Continuous View (SDE):**

- Continuous time: t $\in$ [0, T]
- Smooth noise function: $\beta(t)$
- Can discretize flexibly at sampling time

# From ODEs to SDEs - Modeling Diffusion

- Diffusion Models (Discrete View)
  - Forward: Add noise in T discrete steps (e.g., T=1000)
  - Backward: Remove noise in T discrete steps

**Problem:**

- Can't flexibly adjust **steps** during generation
- Limited flexibility in how we add/remove noise
- Can't use powerful continuous-time solvers

**Carnegie Mellon University**

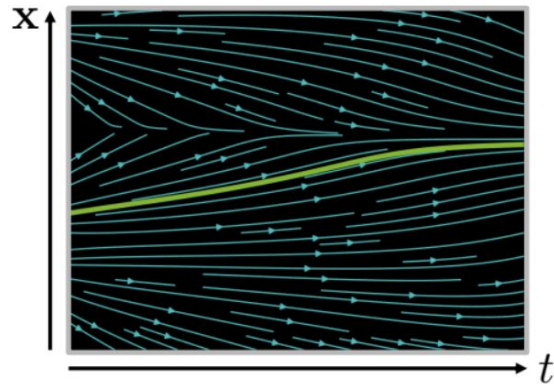# From ODEs to SDEs - Modeling Diffusion

**Why SDEs for Diffusion Models?**

- Deterministic Sampling via Probability Flow ODEs
  - Derive equivalent ODE from any diffusion SDE (no randomness in generation)
  - Enable exact reconstruction, smooth interpolation, and image editing

- Unified Framework
  - DDPM, score-based, and flow models as special cases of one SDE
  - Transfer insights: train with one method, sample with another

**Carnegie Mellon University**

# Ordinary Differential Equations

Ordinary Differential Equation (ODE):

$$\frac{d\mathbf{x}}{dt} = \mathbf{f}(\mathbf{x}, t) \quad \text{or} \quad d\mathbf{x} = \mathbf{f}(\mathbf{x}, t)dt$$



Analytical Solution: $\quad \mathbf{x}(t) = \mathbf{x}(0) + \int_0^t \mathbf{f}(\mathbf{x}, \tau)d\tau$
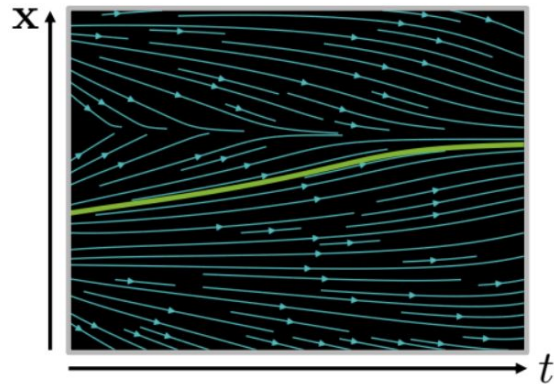
Iterative Numerical Solution: $\quad \mathbf{x}(t + \Delta t) \approx \mathbf{x}(t) + \mathbf{f}(\mathbf{x}(t), t)\Delta t$

**Slide credit to: https://cvpr2022-tutorial-diffusion-models.github.io/**

**Carnegie Mellon University**

# Stochastic Differential Equations



**Ordinary Differential Equation (ODE):**

$$\frac{\mathrm{d}\mathbf{x}}{\mathrm{d}t} = \mathbf{f}(\mathbf{x}, t) \quad \text{or} \quad \mathrm{d}\mathbf{x} = \mathbf{f}(\mathbf{x}, t)\mathrm{d}t$$

Analytical Solution:
$$\mathbf{x}(t) = \mathbf{x}(0) + \int_0^t \mathbf{f}(\mathbf{x}, \tau)\mathrm{d}\tau$$
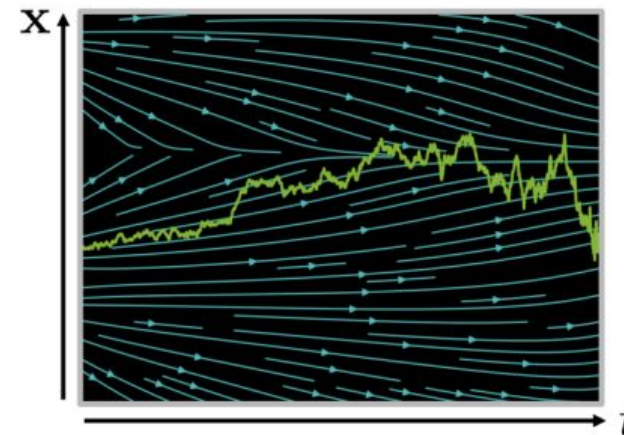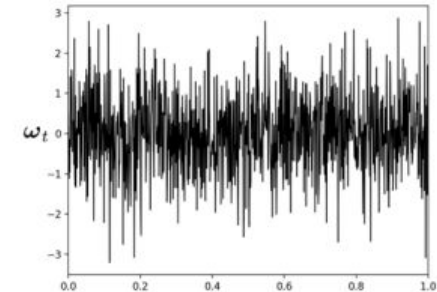
Iterative Numerical Solution:
$$\mathbf{x}(t + \Delta t) \approx \mathbf{x}(t) + \mathbf{f}(\mathbf{x}(t), t)\Delta t$$

**Stochastic Differential Equation (SDE):**

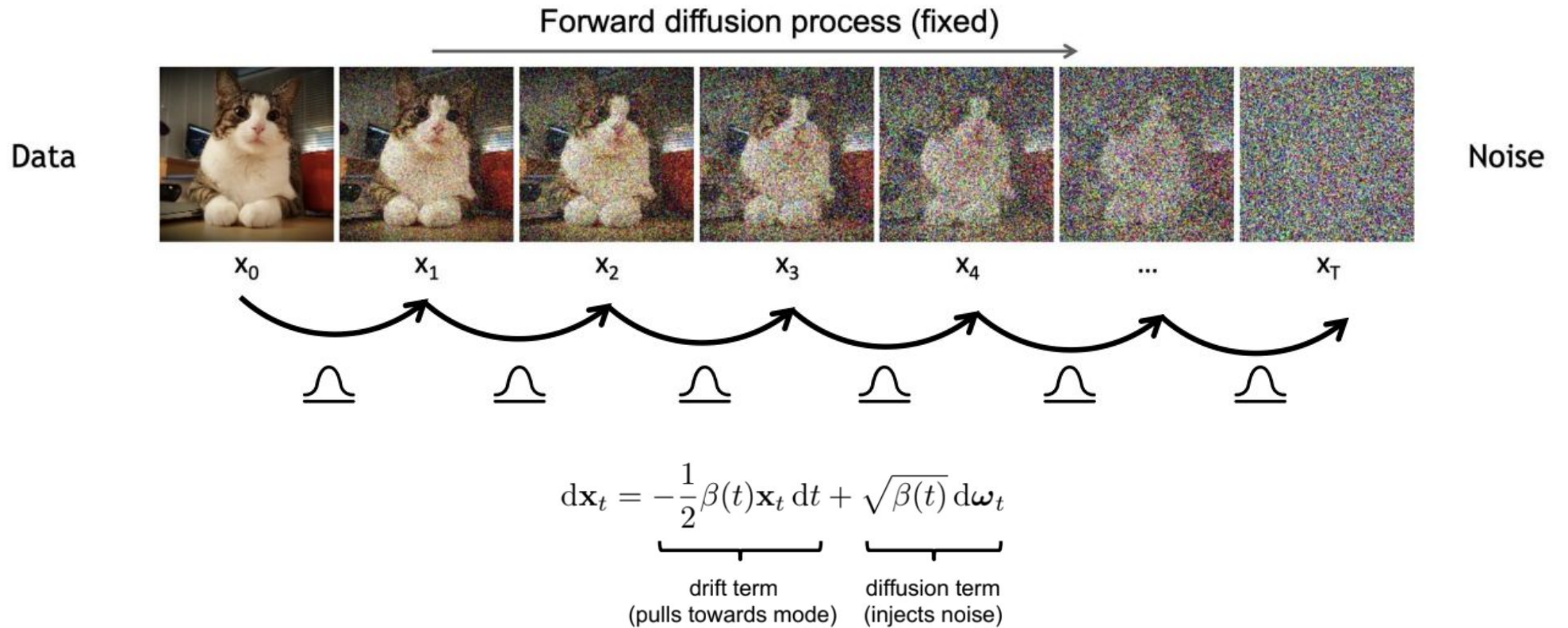$$\frac{\mathrm{d}\mathbf{x}}{\mathrm{d}t} = \underbrace{\mathbf{f}(\mathbf{x}, t)}_{\text{drift coefficient}} + \underbrace{\sigma(\mathbf{x}, t)}_{\text{diffusion coefficient}}\boldsymbol{\omega}_t \longleftarrow \text{Wiener Process (Gaussian White Noise)}$$

$$\left( \mathrm{d}\mathbf{x} = \mathbf{f}(\mathbf{x}, t)\mathrm{d}t + \sigma(\mathbf{x}, t)\mathrm{d}\boldsymbol{\omega}_t \right)$$
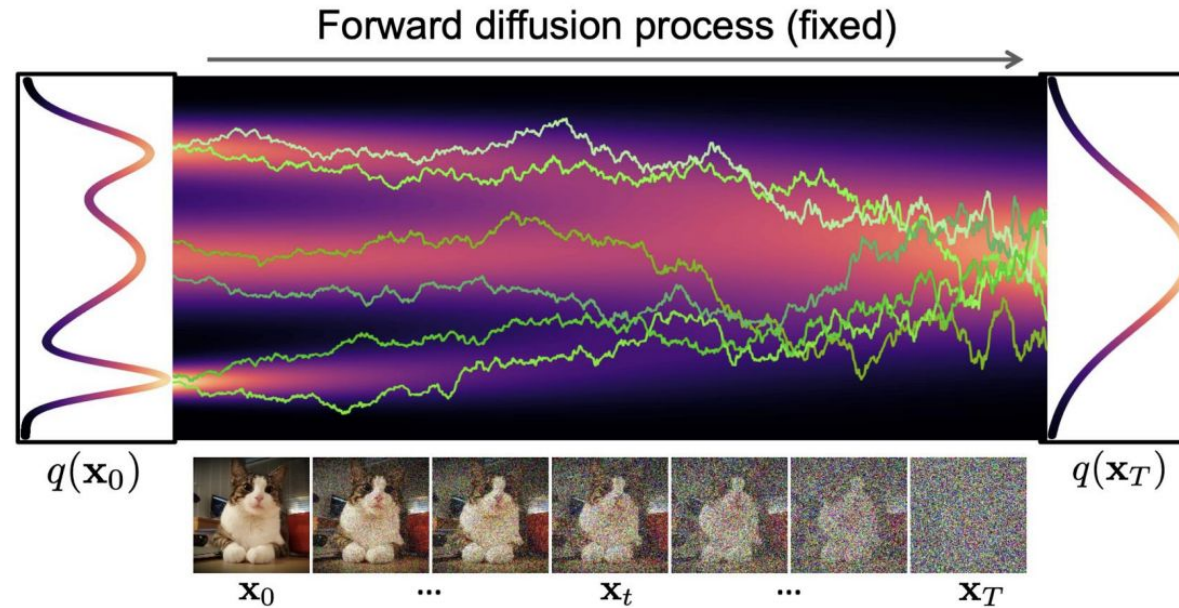
$$\mathbf{x}(t + \Delta t) \approx \mathbf{x}(t) + \mathbf{f}(\mathbf{x}(t), t)\Delta t + \sigma(\mathbf{x}(t), t)\sqrt{\Delta t}\mathcal{N}(\mathbf{0}, \mathbf{I})$$

**Slide credit to: https://cvpr2022-tutorial-diffusion-models.github.io/**

**Carnegie Mellon University**

# Forward Diffusion Process as SDEs



Forward diffusion process (fixed)

Data → Noise

$x_0$, $x_1$, $x_2$, $x_3$, $x_4$, ..., $x_T$

$$\mathrm{d}\mathbf{x}_t = -\frac{1}{2}\beta(t)\mathbf{x}_t\,\mathrm{d}t + \sqrt{\beta(t)}\,\mathrm{d}\boldsymbol{\omega}_t$$

drift term (pulls towards mode)

diffusion term (injects noise)

**Slide credit to: https://cvpr2022-tutorial-diffusion-models.github.io/**

**Carnegie Mellon University**

# Forward Diffusion Process as SDEs



Forward diffusion process (fixed)

$q(\mathbf{x}_0)$     $q(\mathbf{x}_T)$

$\mathbf{x}_0$   ...   $\mathbf{x}_t$   ...   $\mathbf{x}_T$

$$\mathrm{d}\mathbf{x}_t = -\frac{1}{2}\beta(t)\mathbf{x}_t\,\mathrm{d}t + \sqrt{\beta(t)}\,\mathrm{d}\boldsymbol{\omega}_t$$

drift term (pulls towards mode)     diffusion term (injects noise)

**Carnegie Mellon University**

# Forward Diffusion Process as SDEs



Figure credit to: https://yang-song.net/blog/2021/score/

Carnegie Mellon University

# Generative Reverse SDEs



Forward diffusion process (fixed)

$q(\mathbf{x}_0)$  $\mathbf{x}_0$  ...  $\mathbf{x}_t$  ...  $\mathbf{x}_T$  $q(\mathbf{x}_T)$

drift term  diffusion term

$$\mathrm{d}\mathbf{x}_t = \left[ -\frac{1}{2}\beta(t)\mathbf{x}_t - \beta(t)\nabla_{\mathbf{x}_t}\log q_t(\mathbf{x}_t) \right]\mathrm{d}t + \sqrt{\beta(t)}\,\mathrm{d}\bar{\boldsymbol{\omega}}_t$$

"Score Function"

Carnegie Mellon University

# Generative Reverse SDEs



**Figure credit to: https://yang-song.net/blog/2021/score/**

Carnegie Mellon University

# The Score Function: Why It Matters

**Traditional Approach:**

Learn p(x) directly → requires normalizing constant Z → Z is intractable for complex data

**Score-Based Approach:**

Learn $\nabla_x \log p(x)$ instead → no Z needed! → Z cancels in gradient

**Carnegie Mellon University**

# The Score Function: Why It Matters

**Connection to Diffusion:**

Predicting noise $\varepsilon \Leftrightarrow$ Estimating score $\nabla_x \log p(x)$

Same objective, different interpretation

Score tells us direction toward high probability

Carnegie Mellon University

# Score Matching

- General form of probability density function

$$p_\theta(\mathbf{x}) = \frac{e^{-f_\theta(\mathbf{x})}}{Z_\theta}$$

- Maximizing the log-likelihood requires us to know $Z_\theta$
  - Often intractable

- Instead, we can model the score function

$$\nabla_{\mathbf{x}} \log p(\mathbf{x})$$

**Carnegie Mellon University**

# Denoising Score Matching

Forward SDE (data → noise)

$$\mathbf{x}(0) \qquad \mathrm{d}\mathbf{x} = \mathbf{f}(\mathbf{x}, t)\mathrm{d}t + g(t)\mathrm{d}\mathbf{w} \qquad \mathbf{x}(T)$$

**score function**

$$\mathbf{x}(0) \quad \longleftarrow \quad \mathrm{d}\mathbf{x} = \left[\mathbf{f}(\mathbf{x}, t) - g^2(t)\boxed{\nabla_{\mathbf{x}} \log p_t(\mathbf{x})}\right] \mathrm{d}t + g(t)\mathrm{d}\bar{\mathbf{w}} \quad \longleftarrow \quad \mathbf{x}(T)$$

Reverse SDE (noise → data)

**Figure credit to: https://yang-song.net/blog/2021/score/**
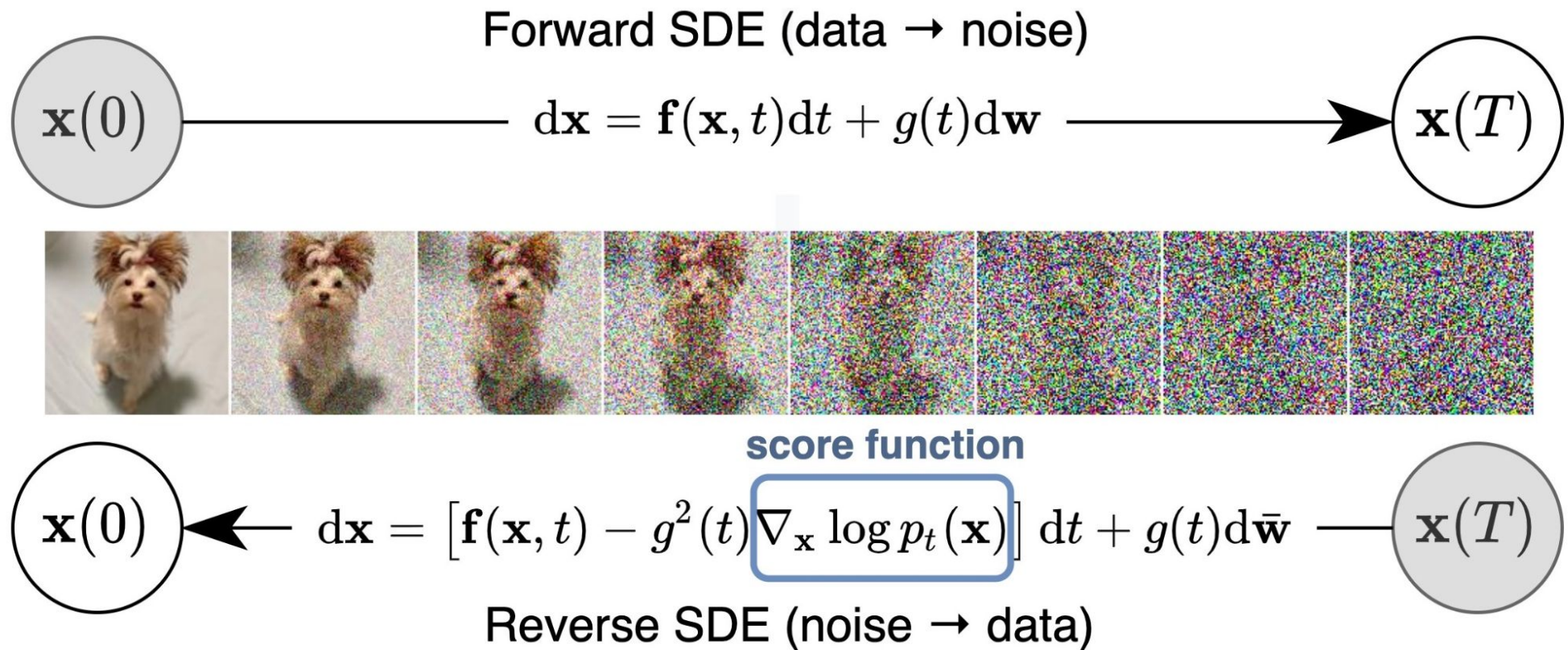
Carnegie Mellon University

# Weighted Diffusion Objective

Denoising score matching objective with loss weighting

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{x}_0 \sim q_0(\mathbf{x}_0)} \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0},\mathbf{I})} \frac{\lambda(t)}{\sigma_t^2} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\mathbf{x}_t, t) \right\|_2^2$$

Problem: AND Why weighting $\lambda(t)$?

- Different timesteps contribute differently
- Early timesteps (high noise): easy to predict, less important
- Late timesteps (low noise): hard to predict, very important

Common choices:

- $\lambda(t) = 1$: Simple, uniform weighting (DDPM)
- $\lambda(t) = \sigma\_t^2$: Better perceptual quality
- $\lambda(t) = \beta(t)$: Maximum likelihood

**Carnegie Mellon University**

# Poll 4

The score function $\nabla \log p(x)$ points toward:

- A) Random directions
- B) Low probability regions
- C) High probability regions (data)
- D) The origin

**Carnegie Mellon University**

# Poll 4

The score function $\nabla \log p(x)$ points toward:

- A) Random directions
- B) Low probability regions
- C) High probability regions (data)
- D) The origin

**Carnegie Mellon University**

# Outline

- Definition of Diffusion

- Importance of Diffusion

- Explaining the process of diffusion

- Denoising Diffusion Implicit Models (DDIM)

- Diffusion Models from Stochastic Differential Equations and Score Matching Perspective

- Classifier-free Guidance (CFG)

- Performance Metrics: FID, IS, Precision and Recall

- Applications of Diffusion Models

**Carnegie Mellon University**

## Conditional Generation: The Control Problem

**Unconditional Diffusion:**

- Sample: Random noise → Random image from training distribution
- No control over output class/content

**The Need for Control:**

- "Draw a sunset over mountains"
- "Create a portrait in Van Gogh style"

**Carnegie Mellon University**

## Conditional Generation: The Control Problem

1. **Classifier Guidance:** Train separate classifier, use gradients
   - ✗ Requires training extra model
   - ✗ Classifier can be noisy at high $t$
2. **Classifier-Free Guidance (CFG):** Train one model for both
   - ✓ Single model
   - ✓ Better quality

**Carnegie Mellon University**

# Classifier-free Guidance (CFG)

**Limited Class Control:** Previous method is incapable of generating an image for a given class

**Purpose of CFG:** it allows for targeted generation of images by conditioning the model on a specific class label.

**Carnegie Mellon University**

# How Classifier-Free Guidance Works

**Training:**

- Randomly drop condition 10-20% of time
- Model learns both:
  - $\epsilon\_\theta(x\_t, class) \leftarrow$ conditional
  - $\epsilon\_\theta(x\_t, \varnothing) \leftarrow$ unconditional

**Sampling:**

$$\tilde{\epsilon} = \epsilon\_uncond + \gamma(\epsilon\_cond - \epsilon\_uncond)$$

$$\uparrow \qquad \uparrow \quad \uparrow$$

baseline     scale  direction toward class

Guidance scale γ:

- γ = 0: Ignore condition (random)
- γ = 1: Standard conditioning
- γ > 1: Stronger conditioning (more typical of class)

**Carnegie Mellon University**

# Classifier-free Guidance (CFG)

---

**Algorithm 1** Classifier guided diffusion sampling, given a diffusion model $(\mu_\theta(x_t), \Sigma_\theta(x_t))$, classifier $p_\phi(y|x_t)$, and gradient scale $s$.

Input: class label $y$, gradient scale $s$
$x_T \leftarrow$ sample from $\mathcal{N}(0, \mathbf{I})$
**for all** $t$ from $T$ to 1 **do**
    $\mu, \Sigma \leftarrow \mu_\theta(x_t), \Sigma_\theta(x_t)$
    $x_{t-1} \leftarrow$ sample from $\mathcal{N}(\mu + s\Sigma \nabla_{x_t} \log p_\phi(y|x_t), \Sigma)$
**end for**
**return** $x_0$

**Carnegie Mellon University**

# Outline

- Definition of Diffusion

- Importance of Diffusion

- Explaining the process of diffusion

- Denoising Diffusion Implicit Models (DDIM)

- Diffusion Models from Stochastic Differential Equations and Score Matching Perspective

- Classifier-free Guidance (CFG)

- Performance Metrics: FID, IS, Precision and Recall

- Applications of Diffusion Models

**Carnegie Mellon University**

# Performance Metrics

**Fréchet Inception Distance (FID)**

Measures the distance between feature vectors of real and generated images; lower scores indicate better image quality and similarity to real images.

**Inception Score (IS)**

Assesses the diversity and clarity of generated images using a pre-trained model; higher scores denote better image clarity and variety.

**Precision and Recall**

Evaluates the quality and diversity of generated images; high precision indicates realistic images, while high recall shows variety close to the actual dataset.

**Carnegie Mellon University**

# Outline

- Definition of Diffusion

- Importance of Diffusion

- Explaining the process of diffusion

- Denoising Diffusion Implicit Models (DDIM)

- Diffusion Models from Stochastic Differential Equations and Score Matching Perspective

- Classifier-free Guidance (CFG)

- Performance Metrics: FID, IS, Precision and Recall
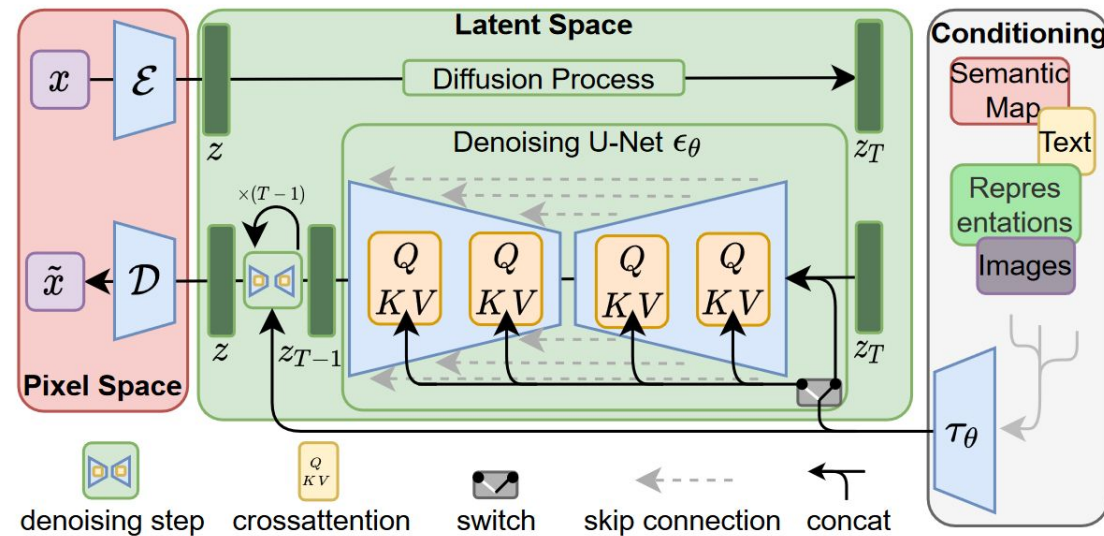
- Applications of Diffusion Models

**Carnegie Mellon University**

# Latent DDPM

**Why Latent Space?**

- **Efficiency:** Running diffusion models directly in pixel space is computationally expensive. Latent Diffusion Models (LDMs) operate in a compressed latent space, drastically reducing computational cost and complexity.
- **Data Compression:** LDMs capture the essential *features* of high-dimensional data like images in a more compact and structured form, enhancing processing efficiency.
- **Improved Performance:** Working in latent space focuses the model on **relevant aspects** of the *data*, improving both the efficiency and effectiveness of the generation process.

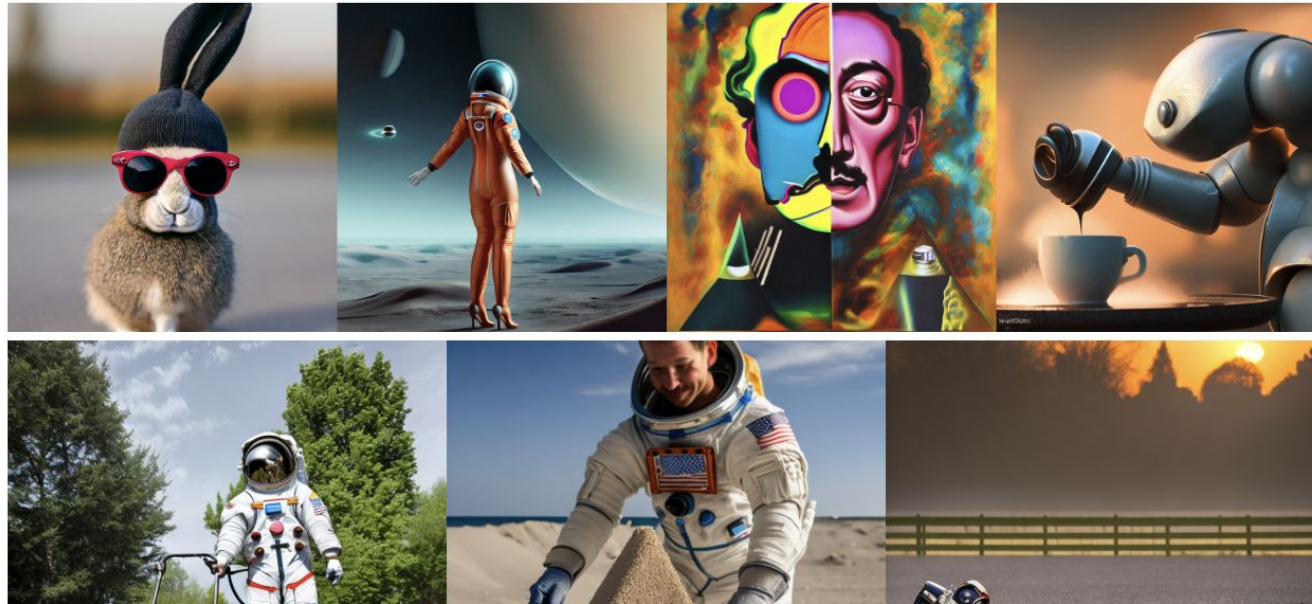**Carnegie Mellon University**

# Latent DDPM

- Map data into latent space using VAE or any other similar approach
1. **Encoder:** Compresses input data into a latent representation.
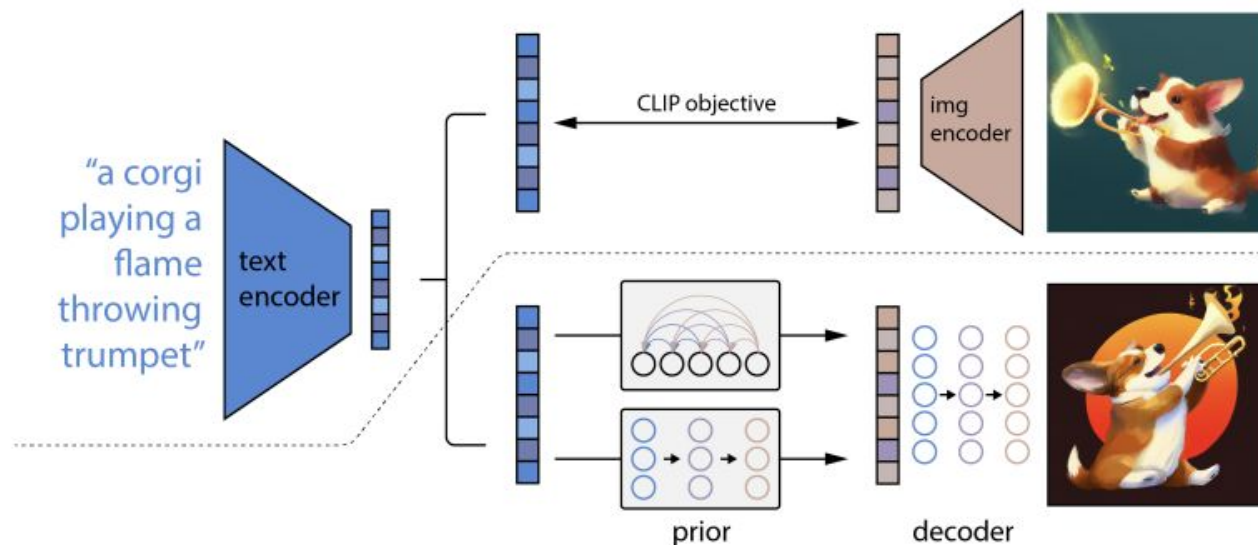2. **Decoder:** Reconstructs the original data from the latent representation.

Carnegie Mellon University

# Stable Diffusion

- Stable Diffusion is a latent text-to-image diffusion model.



Stability AI. https://github.com/Stability-AI/stablediffusion

Carnegie Mellon University

# DALLE

- DALLE is a text-to-image generative model that creates coherent, high-resolution images from natural-language prompts.



Ramesh et al. Hierarchical Text-Conditional Image Generation with CLIP Latents

Carnegie Mellon University

# DiT

- DiT (Diffusion Transformer) is a diffusion model architecture that replaces U-Nets with Transformers to more effectively denoise and generate high-quality images.



Peebles et al. Scalable Diffusion Models with Transformers. 2020.

Carnegie Mellon University

# MAR

- An autoregressive model with diffusion loss



Li et al. Autoregressive Image Generation without Vector Quantization. 2024.

**Carnegie Mellon University**

# Key Takeaways

**1. Many Small Steps Beat One Big Jump**

VAEs try to generate in one step and fail (blurry images). Diffusion uses 1000 tiny steps instead. Each small change is easy to learn. Result: sharp, high-quality outputs.

**2. Forward is Fixed, Reverse is Learning**

Forward process: add noise using a fixed equation (no training needed). Reverse process: train neural network to predict the noise. Training objective: MSE loss.

**3. Three Views, Same Model**

Denoising, score-based, and SDEs are equivalent perspectives. Understanding one helps understand all.

**Carnegie Mellon University**

Thank you :-)

**Carnegie Mellon University**