# Deep Learning (1470)

## Randall Balestriero

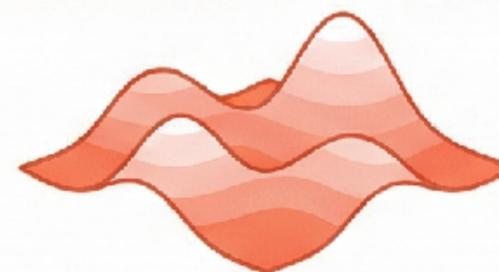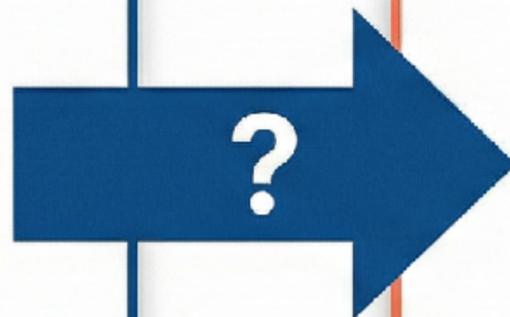**Class 17: Diffusion Models**

# Recap!

# THE GENERATIVE MODELING CHALLENGE

## The Problem

**What We Have**
Training samples from unknown distribution

?

**What We Want**
Learn **p(x)** to generate NEW samples

## Two Approaches

### Implicit Models

noise **z** → Generator **G** → image **x**

- Learn transformation directly
- Example: GANs

✅ **Pros**
Fast sampling

❌ **Cons**
Unstable training, mode collapse

### Explicit Models

image **x** → Model → probability **p(x)**

- Learn the distribution
- Examples: VAEs, Diffusion
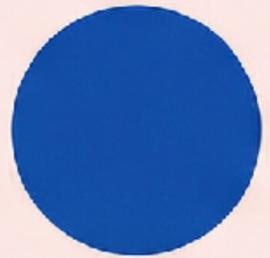
✅ **Pros**
Stable training, likelihood

❌ **Cons**
Can be slow

**We'll focus on this!**

# THE **KEY INSIGHT:** Many Small Steps

## VAE Approach

$z \sim N(0,I)$

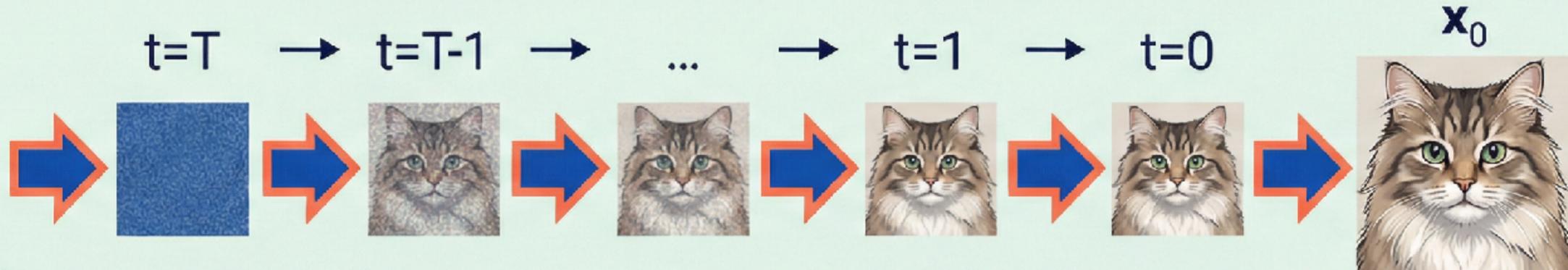One big jump!

$x$

**Hard!** Decoder must do everything at once

**Divide and Conquer:** Break one hard problem into many easy problems!
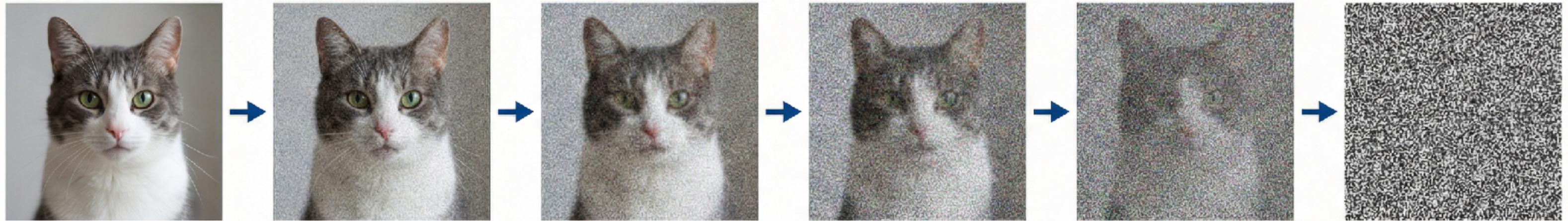
## Diffusion Approach

$x_t \sim N(0,I)$

$t=T \rightarrow t=T-1 \rightarrow ... \rightarrow t=1 \rightarrow t=0$

$x_0$

**Easy!** Each step only removes a little noise

Diffusion = Hierarchical VAE with T latent variables

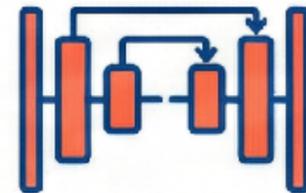# Any Idea?

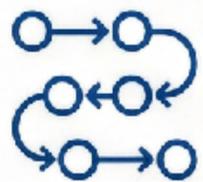# DIFFUSION MODELS
## From Noise to Data

Data $x_0$

Noise $x_t$

**Forward Process** $q(x_t|x_{t-1})$
Add noise (easy, no learning)

**Reverse Process** $p_\theta(x_{t-1}|x_t)$
Remove noise (learned)

U-Net neural network
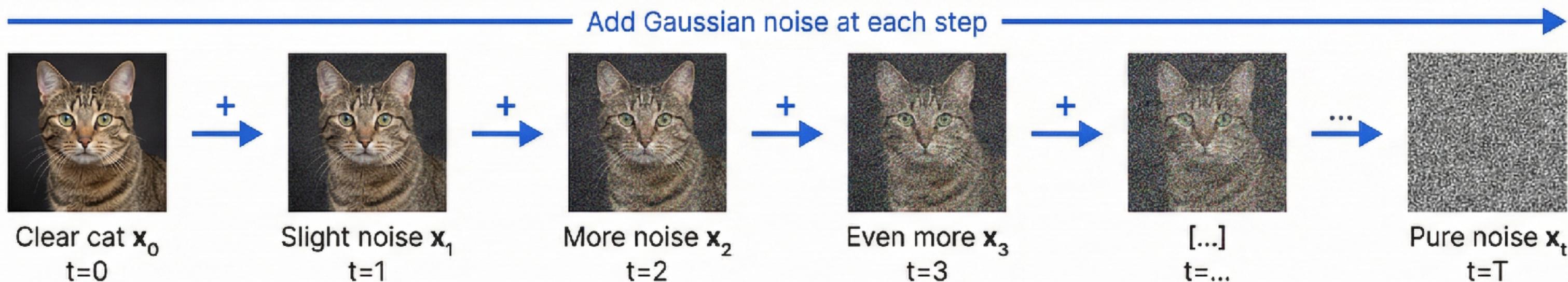
**Iterative**
- Many small steps

**Learnable**
- Neural network denoiser

**Powerful**
- State-of-the-art generation

# THE **FORWARD PROCESS:** Adding Noise

Add Gaussian noise at each step →

| | | | | | |
|---|---|---|---|---|---|
| Clear cat $\mathbf{x}_0$ | Slight noise $\mathbf{x}_1$ | More noise $\mathbf{x}_2$ | Even more $\mathbf{x}_3$ | [...] | Pure noise $\mathbf{x}_t$ |
| t=0 | t=1 | t=2 | t=3 | t=... | t=T |

## Step-by-Step
One Step

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}\left(\sqrt{1-\beta_t}\cdot\mathbf{x}_{t-1}, \beta_t\mathbf{I}\right)$$
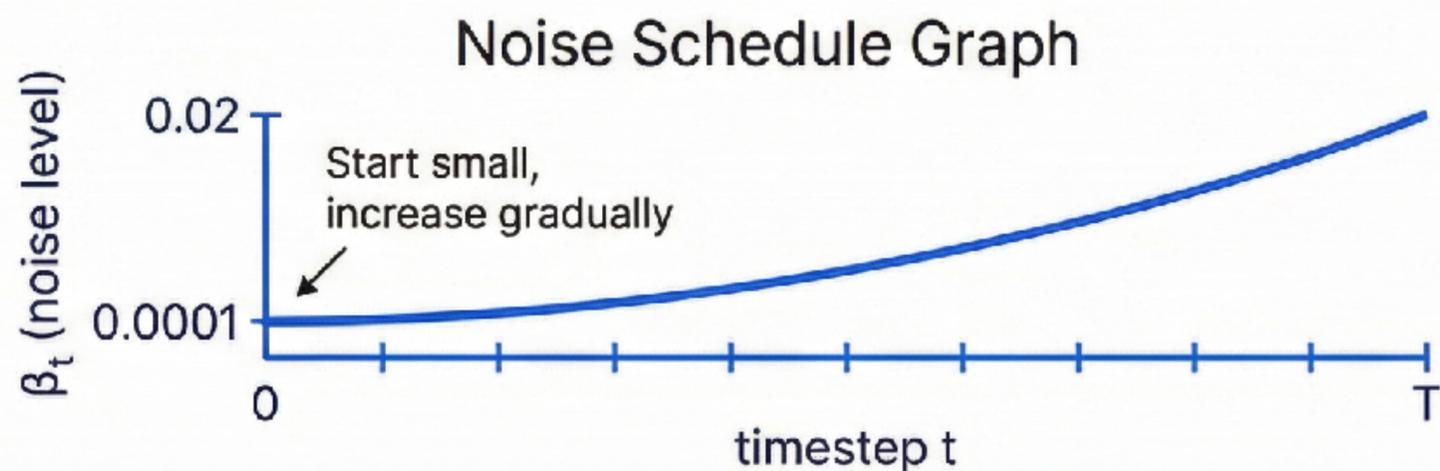
→

Shrink signal, add noise

## Direct Jump
Shortcut

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

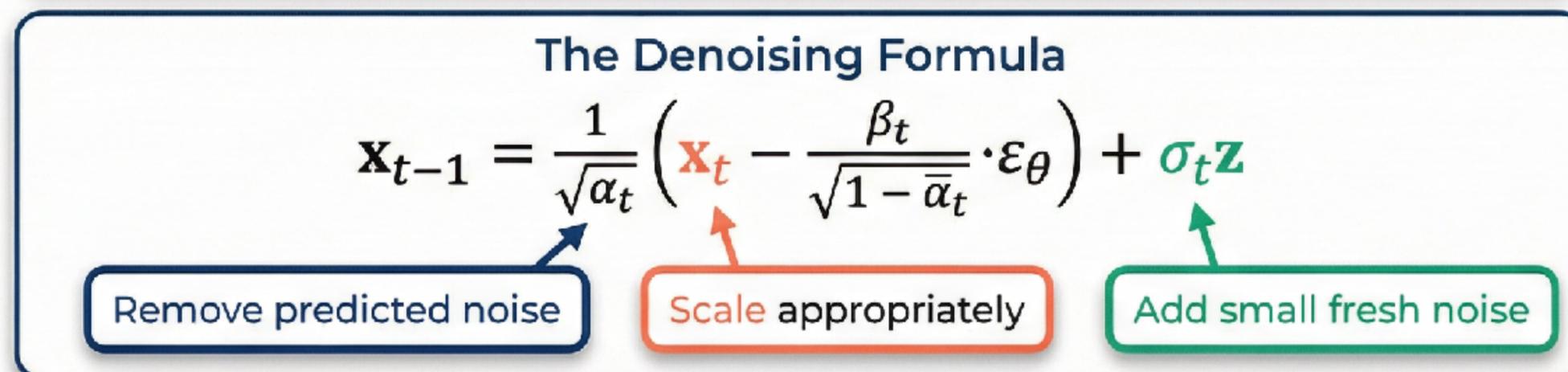$\mathbf{x}_0$ ⟶ $\mathbf{x}_t$

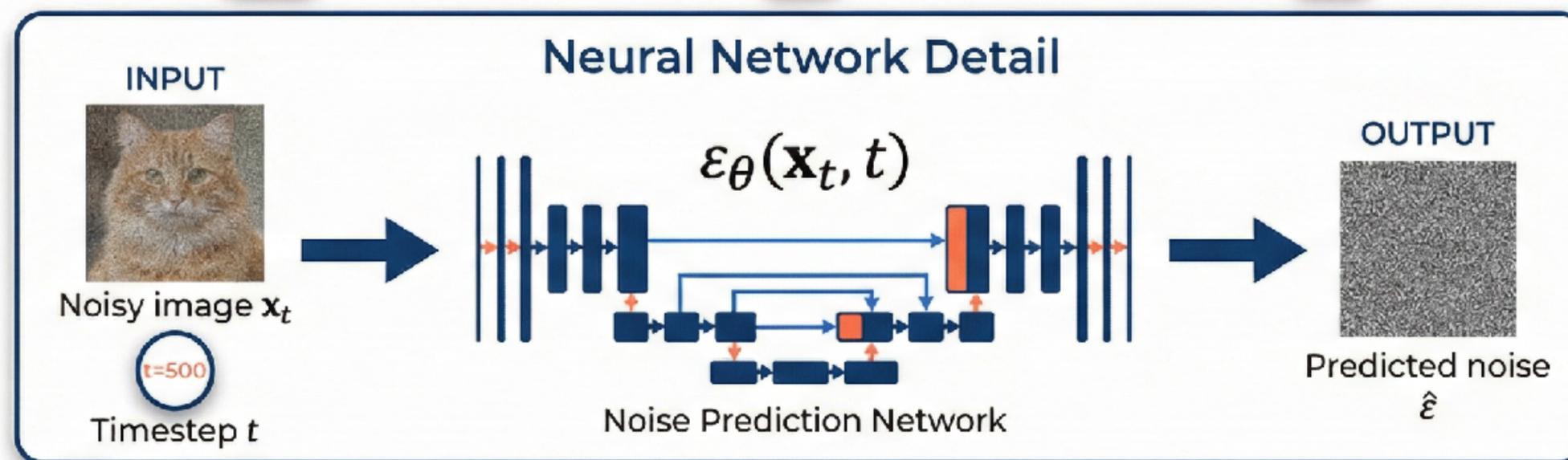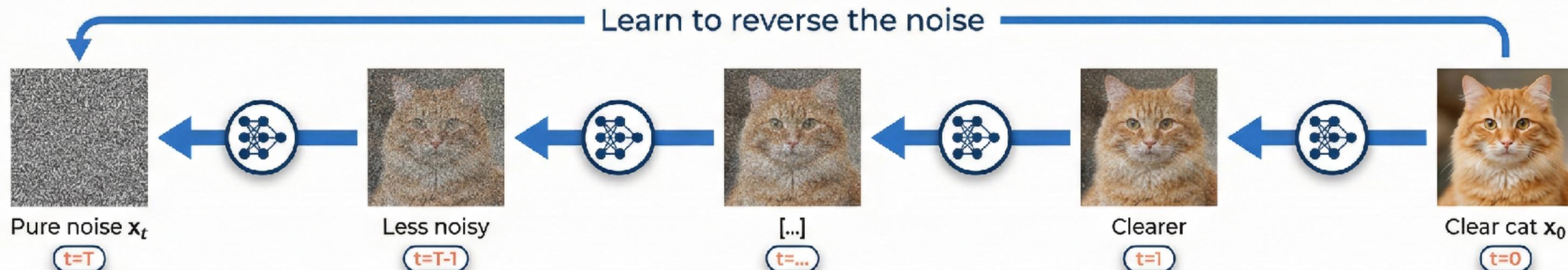Jump to any timestep directly!

where $\bar{\alpha}_t = \prod_i (1-\beta_i)$

## Noise Schedule Graph

Start small, increase gradually

$\beta_t$ (noise level): 0.0001 to 0.02 over timestep t (0 to T)

## Key Point

✓ No learning required - just math!

# THE REVERSE PROCESS: Learning to Denoise

Learn to reverse the noise

Pure noise $\mathbf{x}_t$
t=T

Less noisy
t=T-1

[...]
t=...

Clearer
t=1

Clear cat $\mathbf{x}_0$
t=0

## Neural Network Detail

INPUT

Noisy image $\mathbf{x}_t$

t=500

Timestep $t$

$$\varepsilon_\theta(\mathbf{x}_t, t)$$

Noise Prediction Network

OUTPUT

Predicted noise $\hat{\varepsilon}$

## The Denoising Formula

$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}}\left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\cdot\varepsilon_\theta\right) + \sigma_t\mathbf{z}$$

Remove predicted noise

Scale appropriately

Add small fresh noise

# TRAINING: The Simple Denoising Objective

$$\mathcal{L} = \mathbb{E}\Big[\,\|\boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)\|^2\,\Big]$$

Just predict the noise that was added!

Sample $\mathbf{x}_0$ from data

Sample $t \sim \mathrm{Uniform}(1, T)$ and $\varepsilon \sim \mathcal{N}(0, \mathbf{I})$

Create noisy:
$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t \cdot \mathbf{x}_0} + \sqrt{(1 - \bar{\alpha}_t)} \cdot \varepsilon$$

Predict $\hat{\varepsilon} = \varepsilon_\theta(\mathbf{x}_t, t)$, compute MSE

repeat

**Simple**
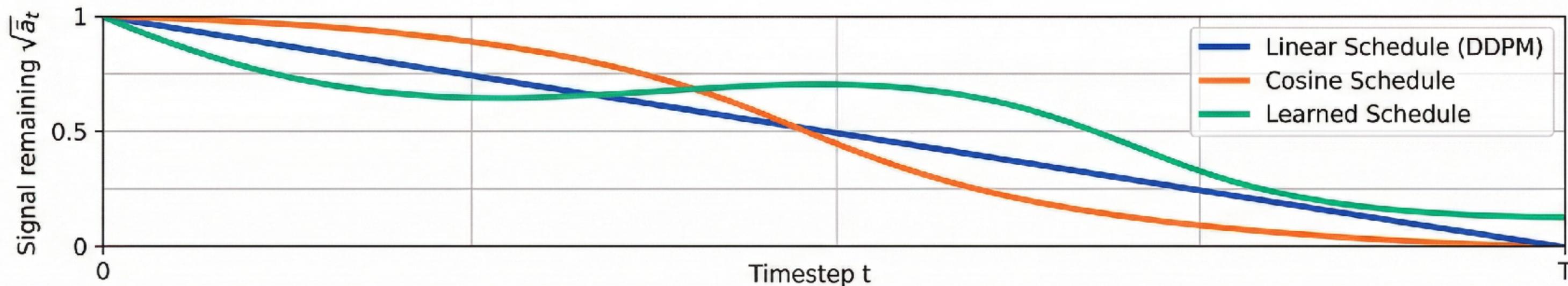Just MSE regression

**Stable**
No adversarial training

**Effective**
State-of-the-art results

# How to schedule alpha?

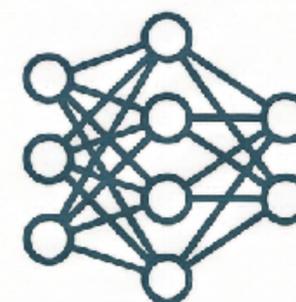# NOISE SCHEDULES: How Fast to Add Noise?



**Linear**

Uniform noise addition
DDPM (2020)

**Cosine**

Preserves structure longer
Improved DDPM (2021)

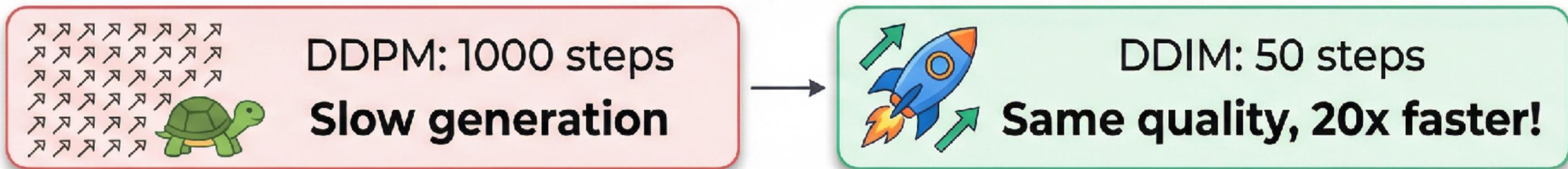**Learned**

Optimized during training
VDM (2021)

💡 Cosine schedule keeps images recognizable longer → better training signal!

# Is Generation Fast?

# Is Generation Fast Now?

# LATENT DIFFUSION: Stable Diffusion (Rombach et al., 2022)

## The Problem

**LEFT**

512×512×3

**Expensive!**

Pixel Space: 786,432 dimensions

**RIGHT**

64×64×4

**48x smaller!**

Latent Space: 16,384 dimensions

## Architecture Diagram

[Image 512×512] → **[ENCODER E]** → [Latent 64×64] → **[U-NET Diffusion]** → [Latent 64×64] → **[DECODER D]** → [Output 512×512]

Pretrained VAE (frozen)

**Diffusion happens HERE**

Pretrained VAE (frozen)
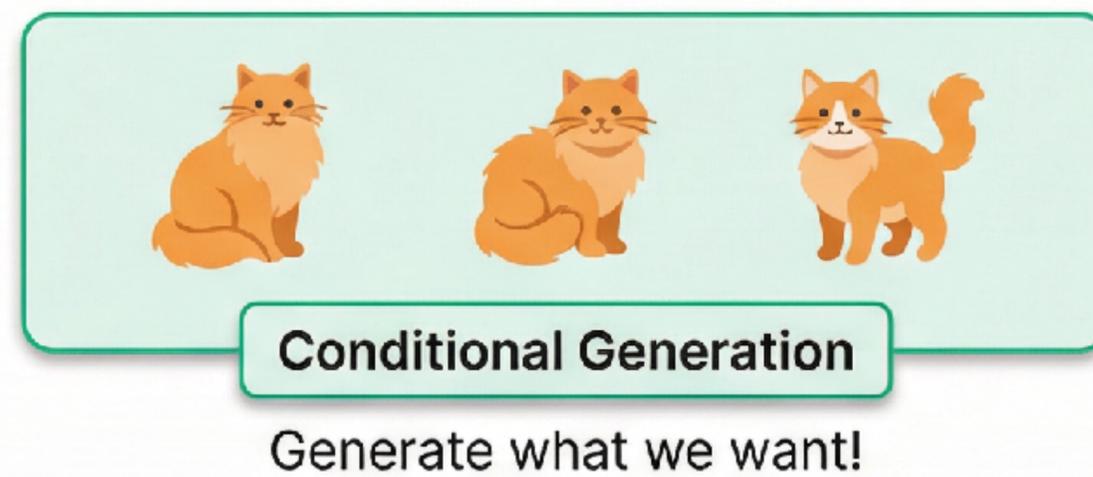
**Efficient**
Train on 1 GPU
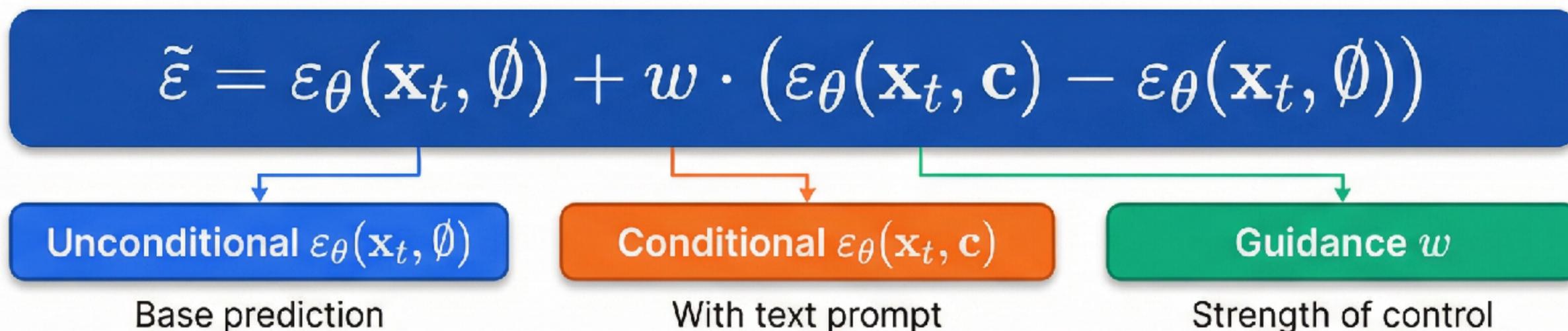
**High Quality**
Perceptually same

**Flexible**
Easy conditioning

This architecture powers Stable Diffusion!

# CLASSIFIER-FREE GUIDANCE: Controlling Generation

**TOP SECTION**
**The Problem**



**Unconditional Generation**

No control over output



**Conditional Generation**

Generate what we want!

**MIDDLE SECTION**
**The Method**

$$\tilde{\varepsilon} = \varepsilon_\theta(\mathbf{x}_t, \emptyset) + w \cdot (\varepsilon_\theta(\mathbf{x}_t, \mathbf{c}) - \varepsilon_\theta(\mathbf{x}_t, \emptyset))$$

**Unconditional** $\varepsilon_\theta(\mathbf{x}_t, \emptyset)$

Base prediction

**Conditional** $\varepsilon_\theta(\mathbf{x}_t, \mathbf{c})$

With text prompt

**Guidance** $w$

Strength of control

**BOTTOM SECTION**
**Guidance Scale**

Too weak 💩          Just right 🌟          Too strong ⚠️

w=1          w=7          w=15

**Key insight** | Train with random prompt dropout → one model does both!

# OPEN QUESTIONS: Food for Thought

**Why does denoising work?**
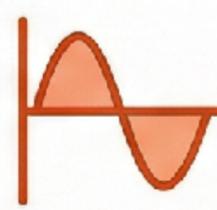What structure is the network learning about natural images?

**How fast can we go?**
$1000 \rightarrow 50 \rightarrow 4 \rightarrow 1$ step?
What's the limit?

**What's the best latent space?**
Pixels? VAE? Something new?

**Discrete or continuous time?**
Steps vs SDEs - which is optimal?

**Beyond images?**
Video, 3D, audio, molecules, robotics...

See you on Wednesday!