

Stable diffusion class models: development and application

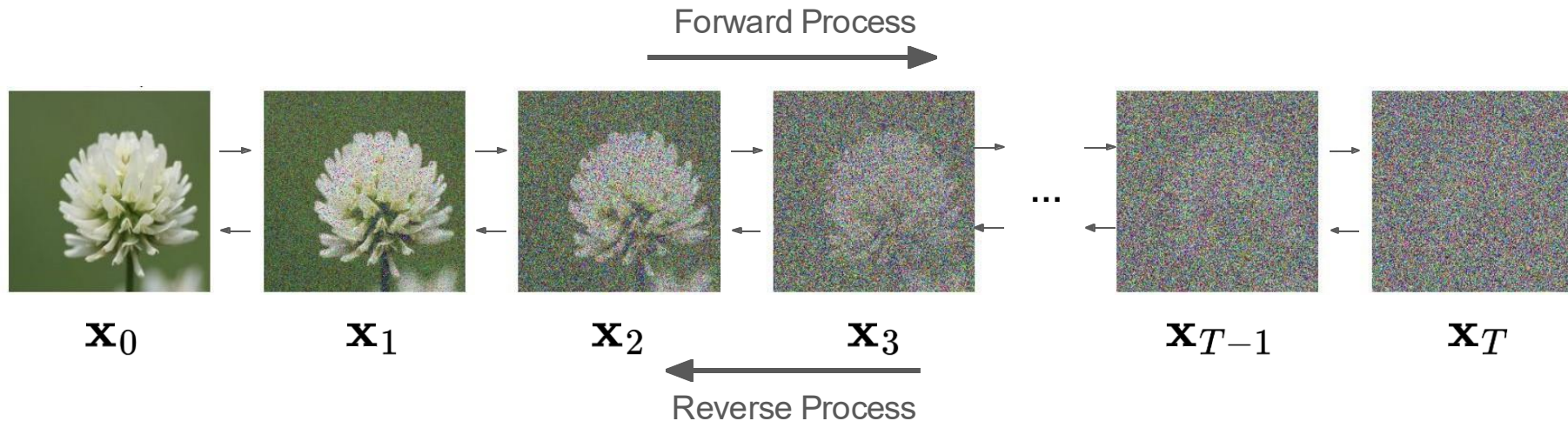
Blagodarniy Artyom



Denoising Diffusion Models

Denoising diffusion models consist of two processes:

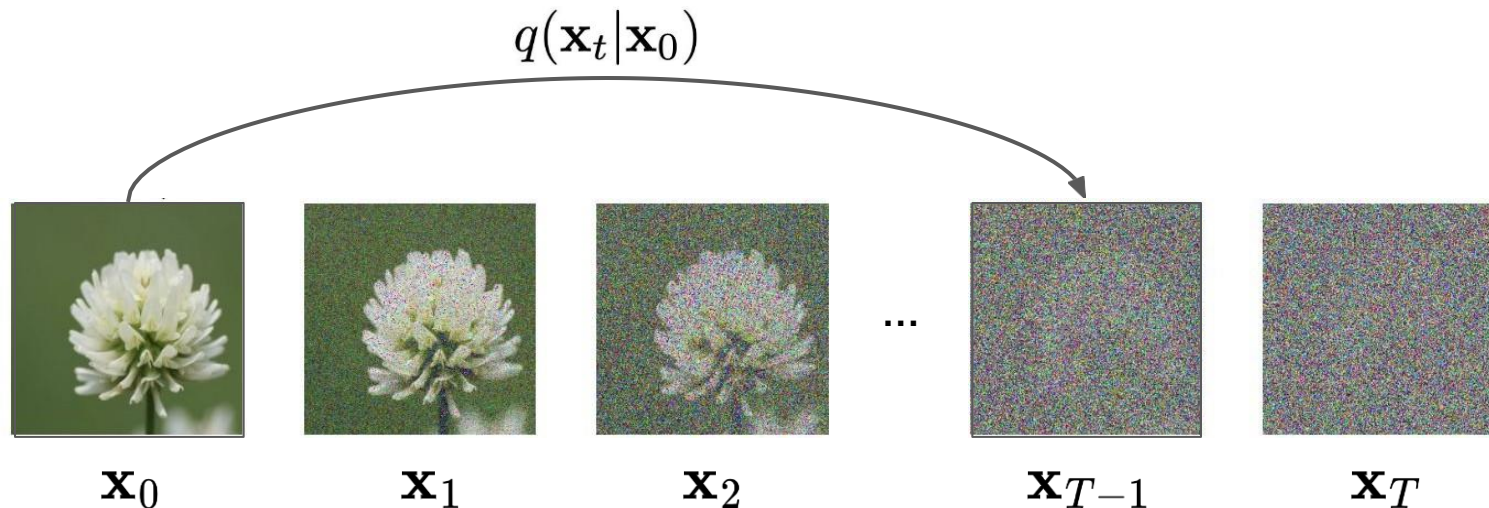
- Forward diffusion process that gradually adds noise to input
- Reverse denoising process that learns to generate data by denoising



Details: Forward Process

Can sample \mathbf{x}_t in closed-form as $q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$

$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1}) \quad \mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \bar{\alpha}_t \in (0, 1)$$



Aside: Noise Schedules

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \bar{\alpha}_t \in (0, 1)$$

- Define the noise schedule in terms of $\bar{\alpha}_t \in (0, 1)$
 - Some monotonically decreasing function from 1 to 0
- Cosine Noise schedule:

$$\bar{\alpha}_t = \cos(.5\pi t/T)^2$$

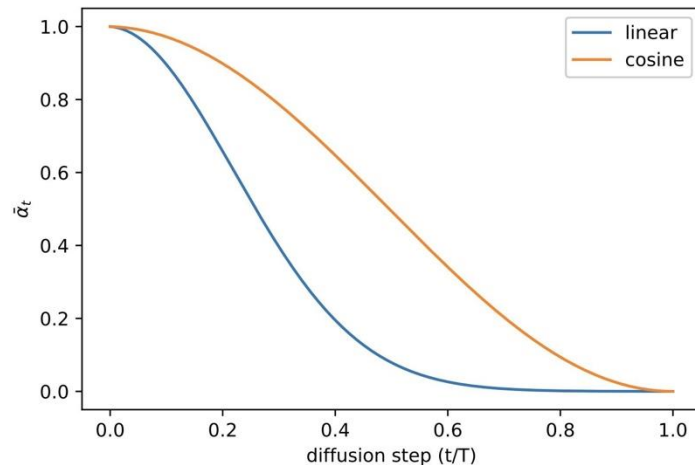
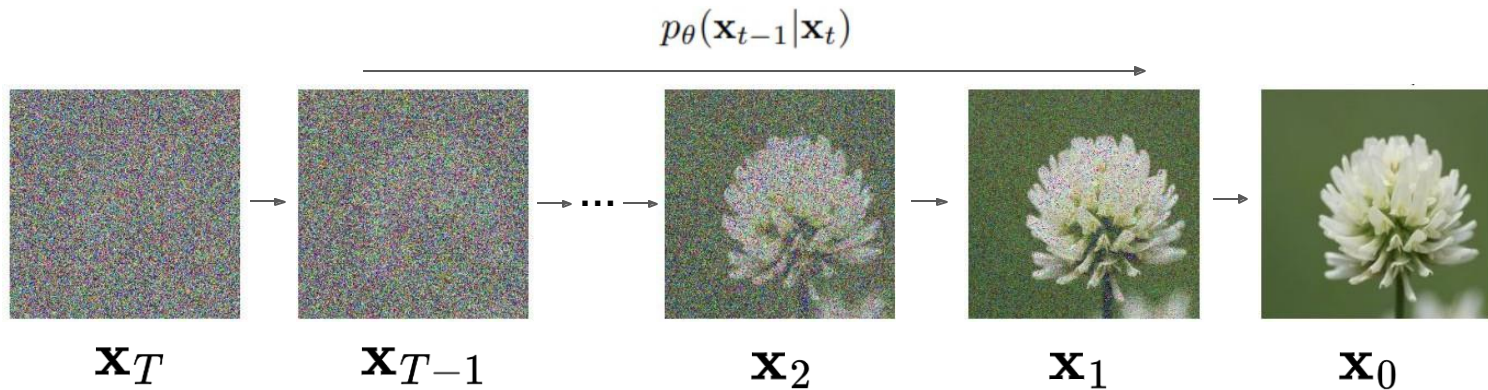


Figure 5. $\bar{\alpha}_t$ throughout diffusion in the linear schedule and our proposed cosine schedule.

Key Idea

We introduce a generative model to approximate the reverse process:

$$\begin{aligned} p(\mathbf{x}_T) &= \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I}) \\ p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) &= \mathcal{N}(\mathbf{x}_{t-1}; \mu_{\theta}(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I}) \end{aligned} \quad \rightarrow \quad p_{\theta}(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$$



Training Objective

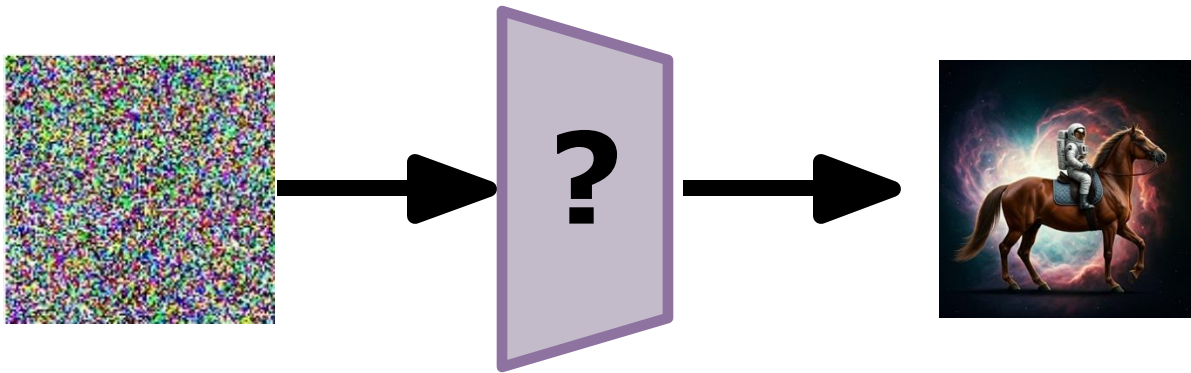
- Bound the likelihood with the ELBO
 - Exactly like VAEs

$$\begin{aligned}\log p(\mathbf{x}) &\geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \\ &= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{prior matching term}} - \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}\end{aligned}$$

Diffusion Models

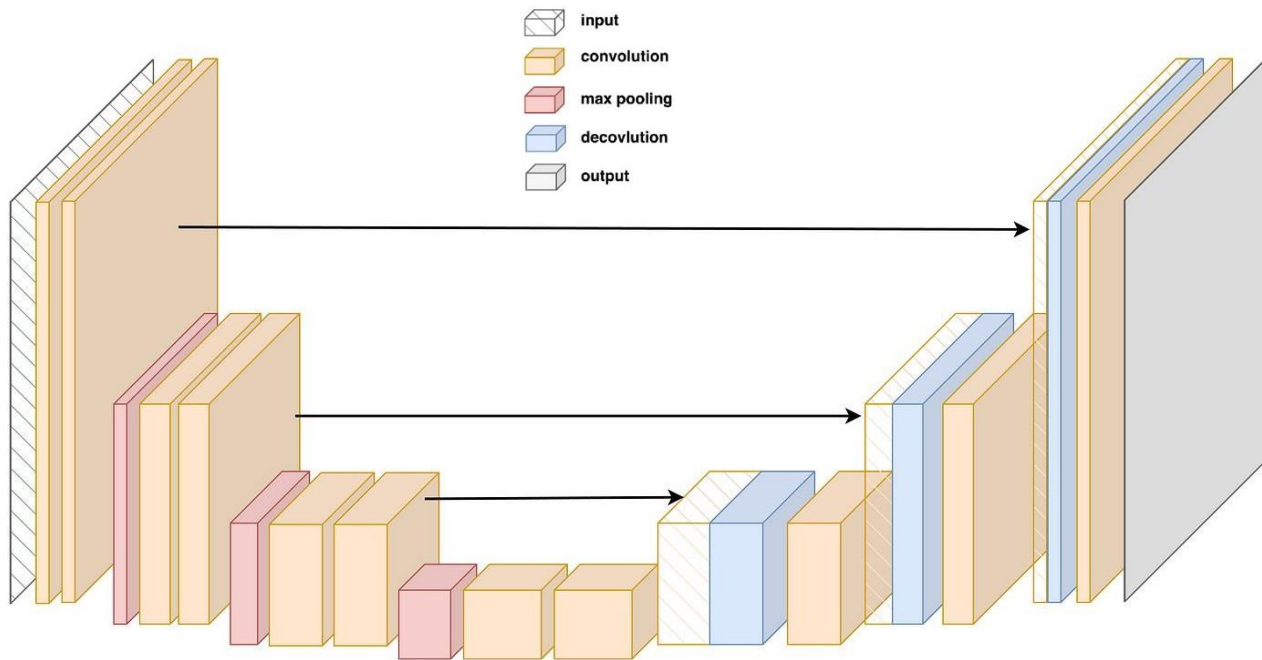
[Sohl-Dickstein et al. 2015]

- Draw a sample of Gaussian noise
- Transform it to obtain a natural image



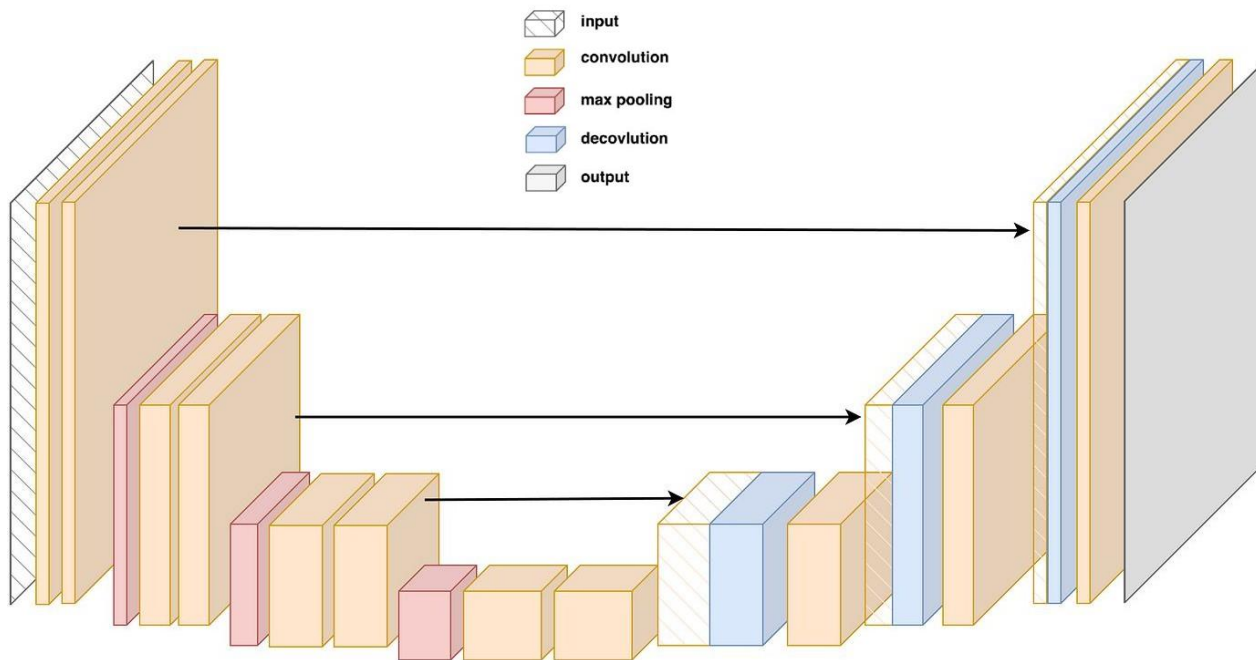
The U-Net

[Ronneberger et al. 2015]



The U-Net

[Ronneberger et al. 2015]

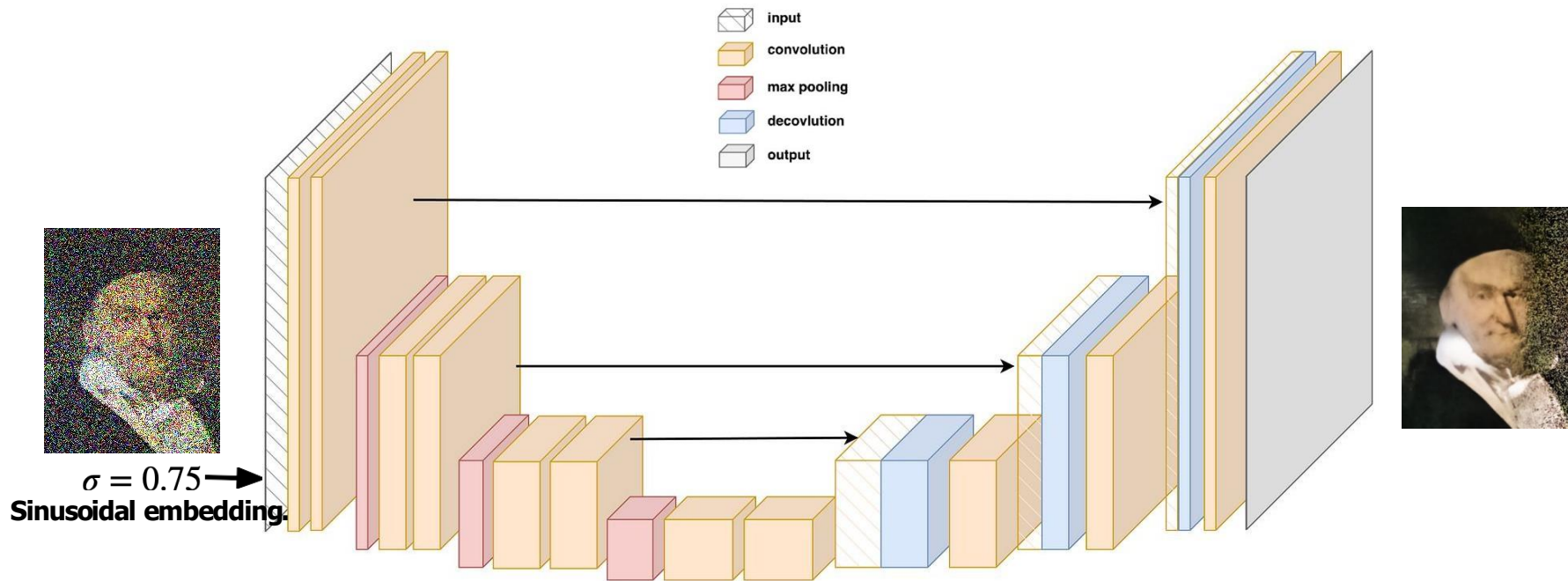


$\sigma = 0.75$



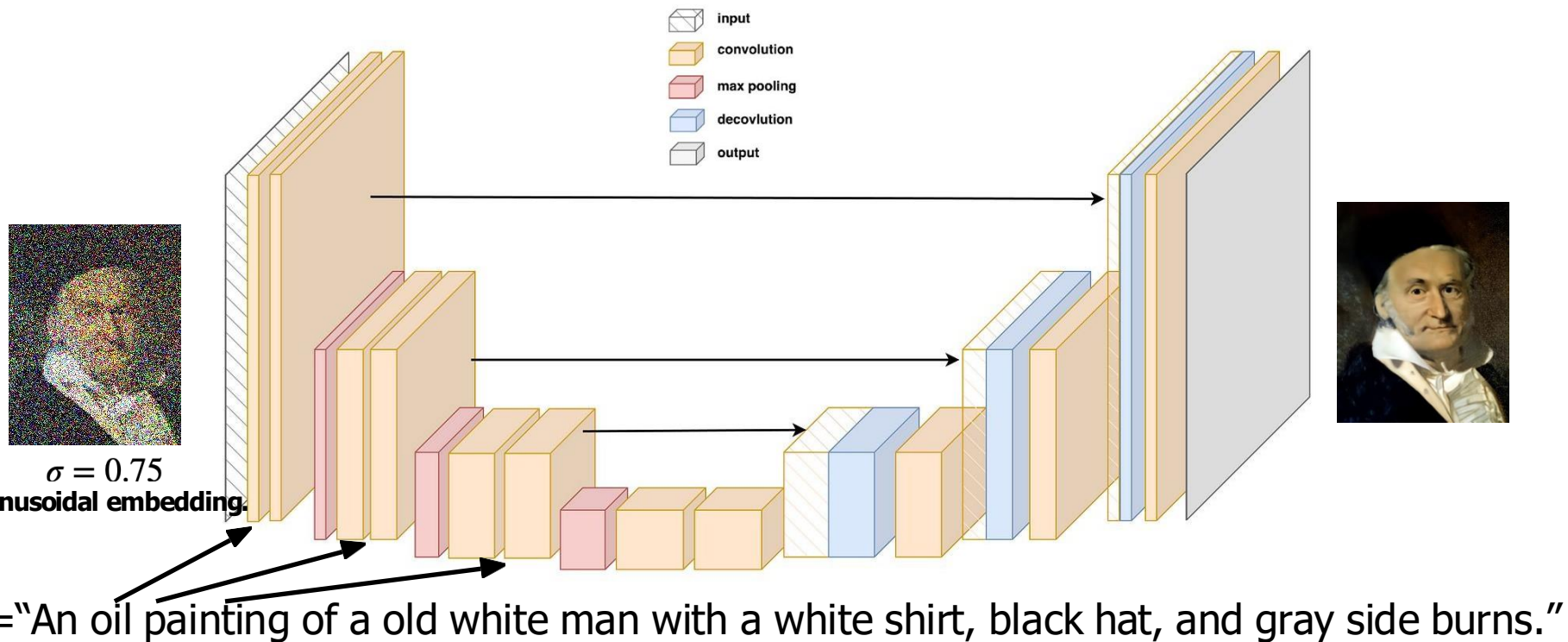
The U-Net

[Ronneberger et al. 2015]



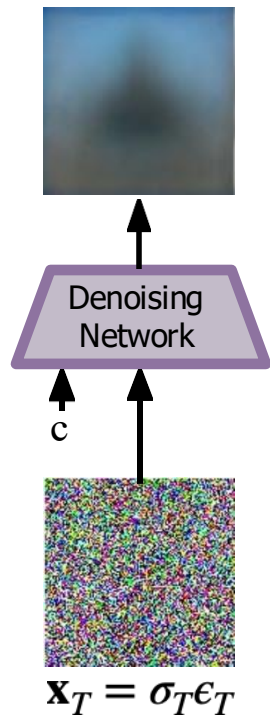
The U-Net

[Ronneberger et al. 2015]



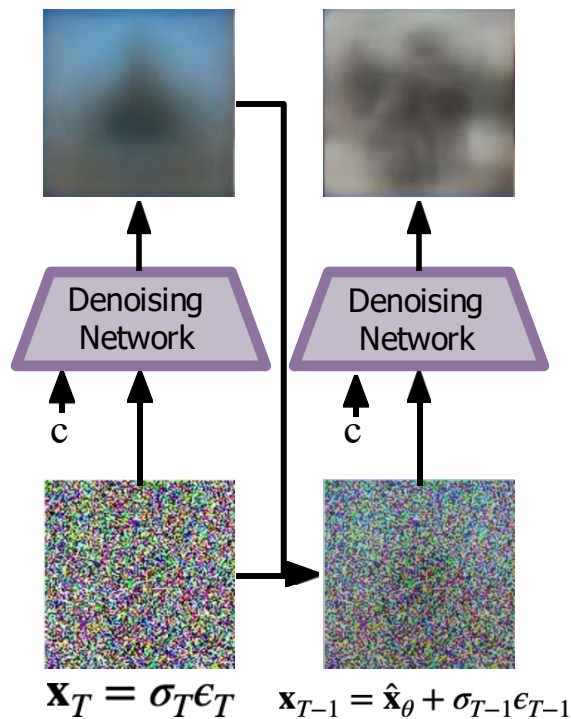
c="An astronaut riding a horse"

Diffusion Sampling



c="An astronaut riding a horse"

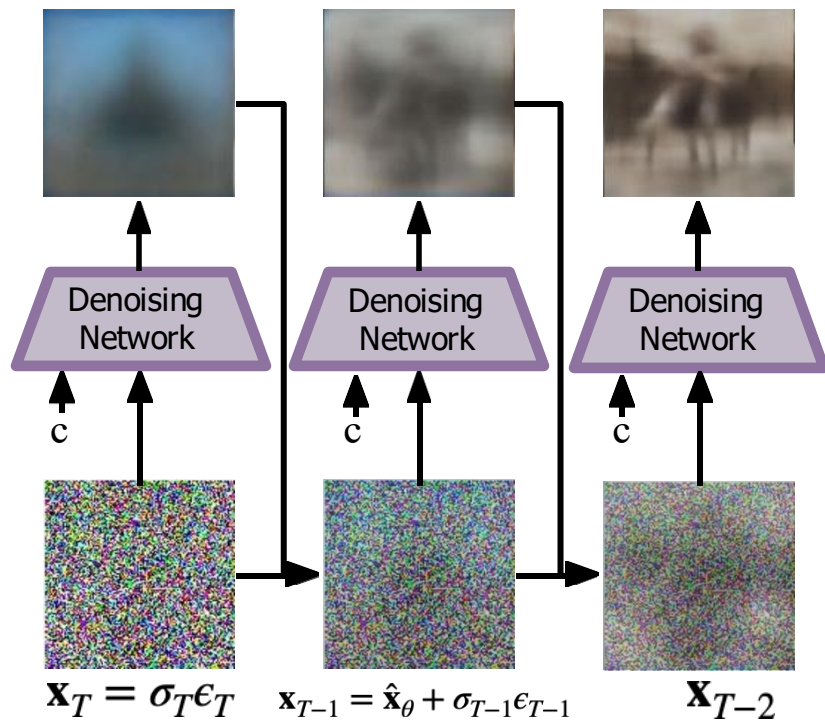
Diffusion Sampling



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

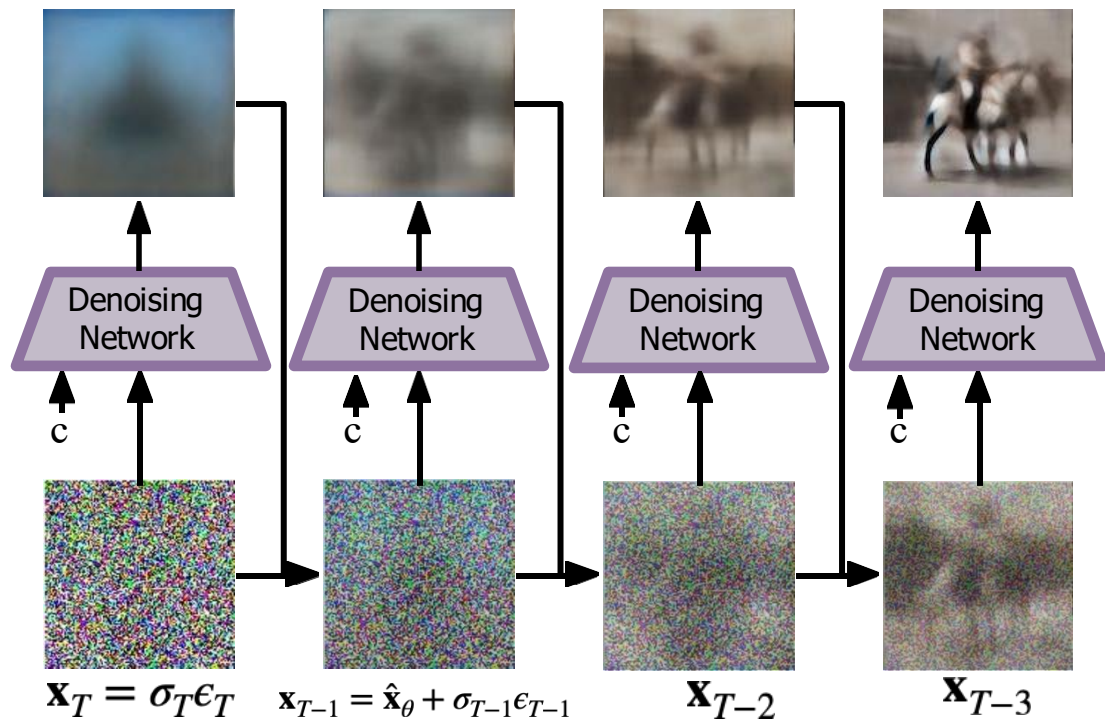
c="An astronaut riding a horse"

Diffusion Sampling



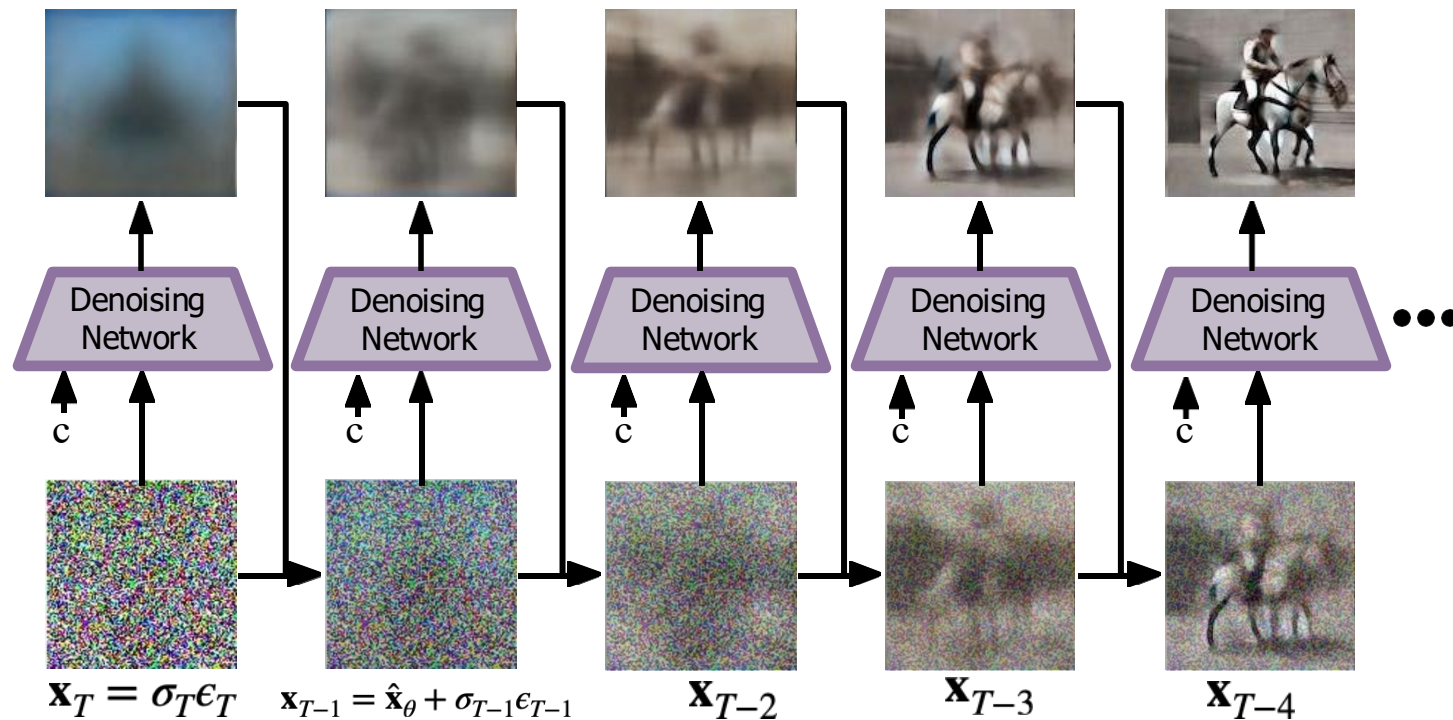
c ="An astronaut riding a horse"

Diffusion Sampling



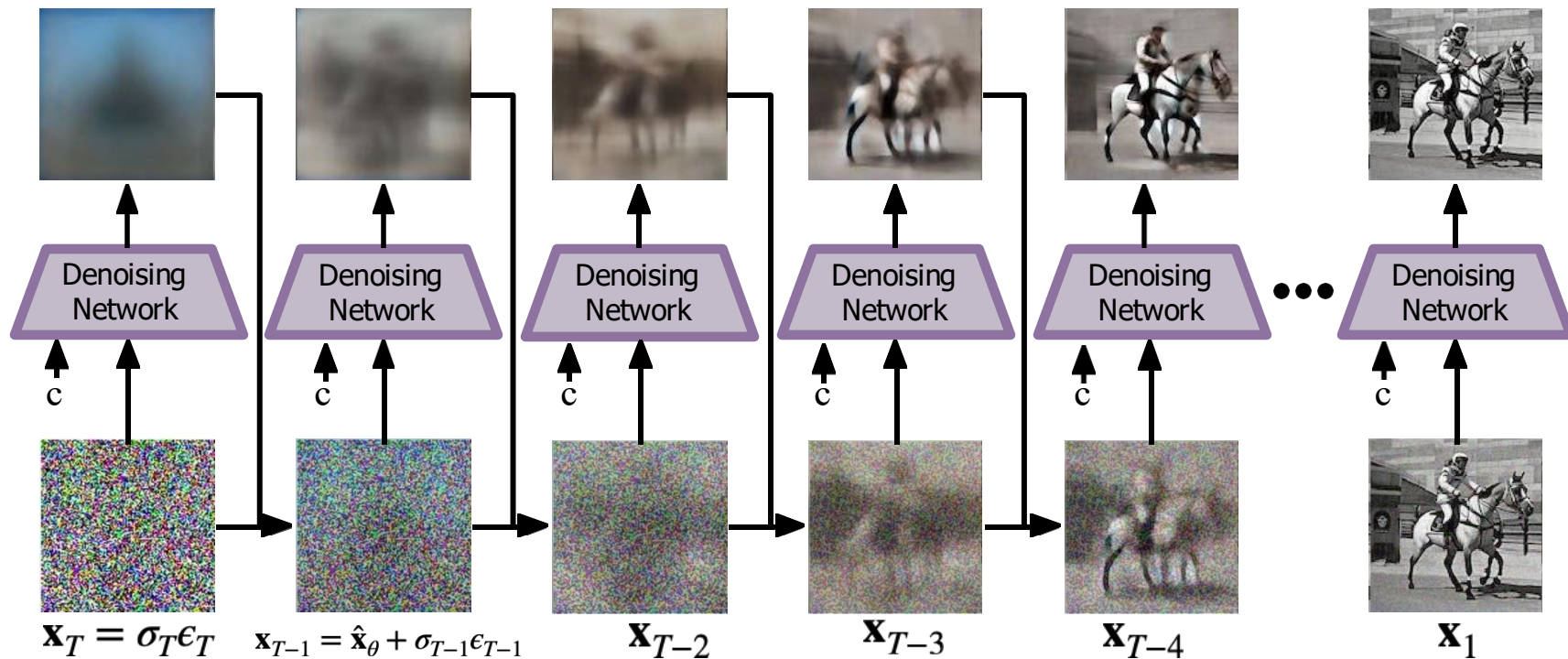
c ="An astronaut riding a horse"

Diffusion Sampling



c ="An astronaut riding a horse"

Diffusion Sampling



Diffusion Models



An astronaut riding a horse

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Training Algorithm

Repeat until convergence

1. $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ \leftarrow Sample original image from image distribution
2. $t \sim U\{1, 2, \dots, T\}$ \leftarrow Sample random time step uniformly
3. $\epsilon \sim \mathcal{N}(0, 1)$ \leftarrow Sample Gaussian noise
4. Optimizer step on $L(\theta) = \mathbb{E}_{t, \mathbf{x}_0, \epsilon} [\|\epsilon - \epsilon_\theta(\mathbf{x}_t, t)\|^2]$
 \leftarrow Model predicts noise applied at time step t and calculate loss

Inference Sampling Algorithm

$\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ \leftarrow Sample pure Gaussian noise

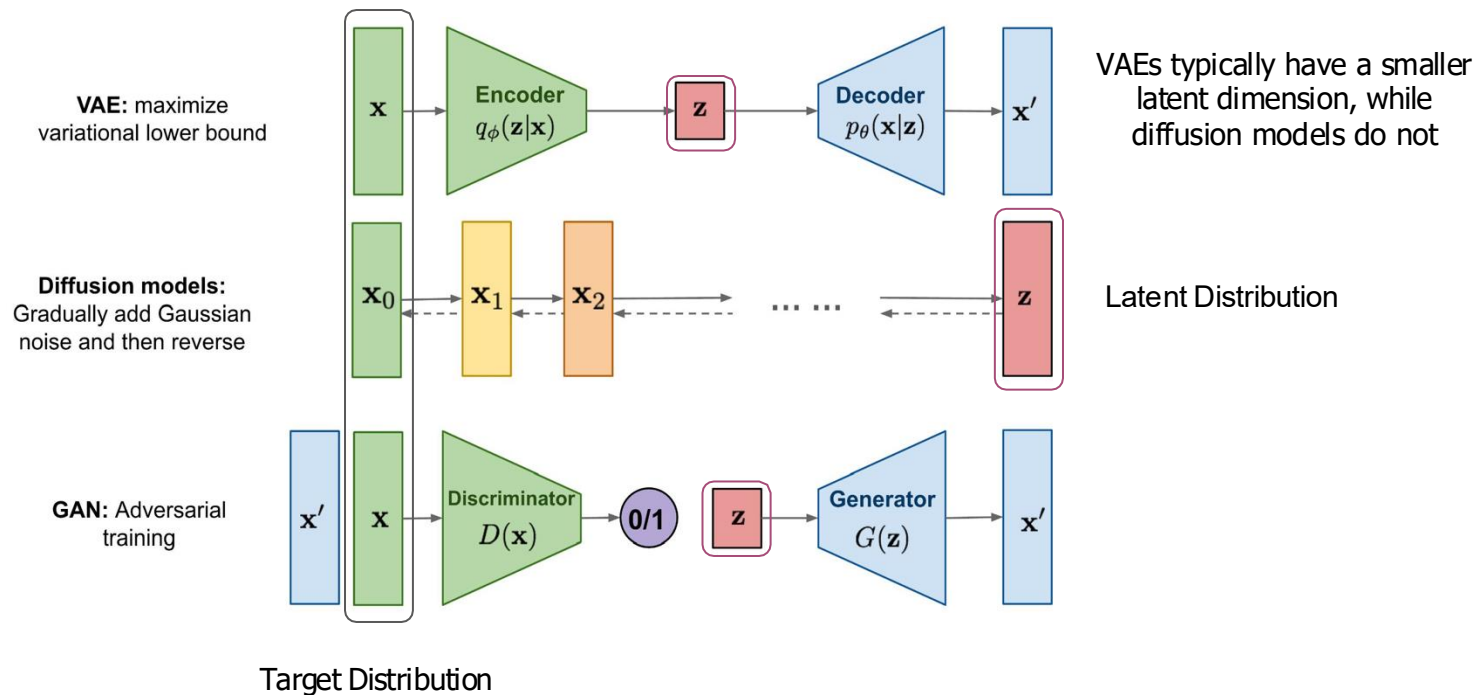
For $t = T, T - 1, \dots, 1$

$\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = \mathbf{0}$ \leftarrow Sample Gaussian noise to apply to image

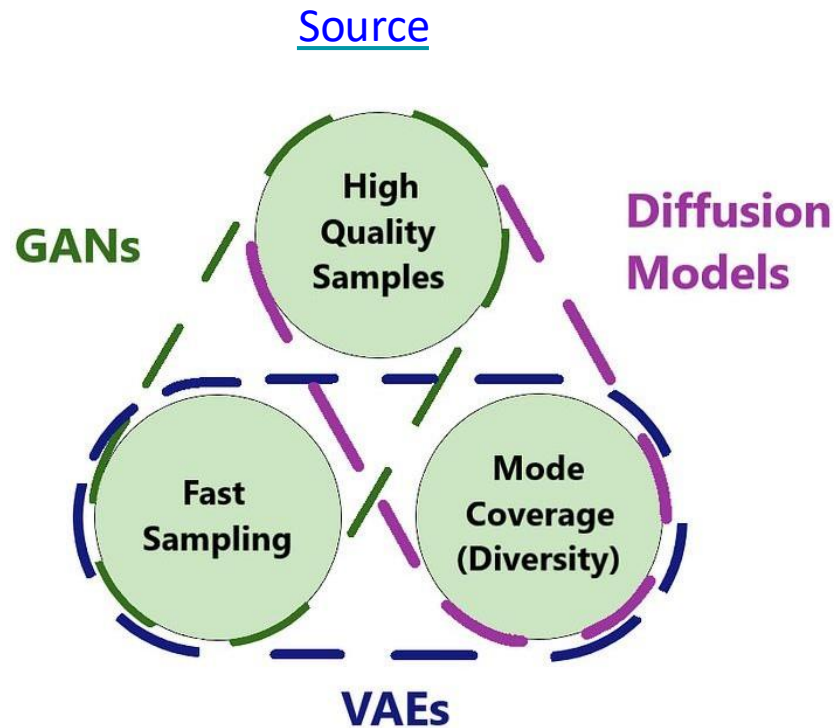
$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ \leftarrow Predict noise applied to image and remove that noise

Return \mathbf{x}_0

Generative Modeling



Diffusion Models vs. VAEs vs. GAN





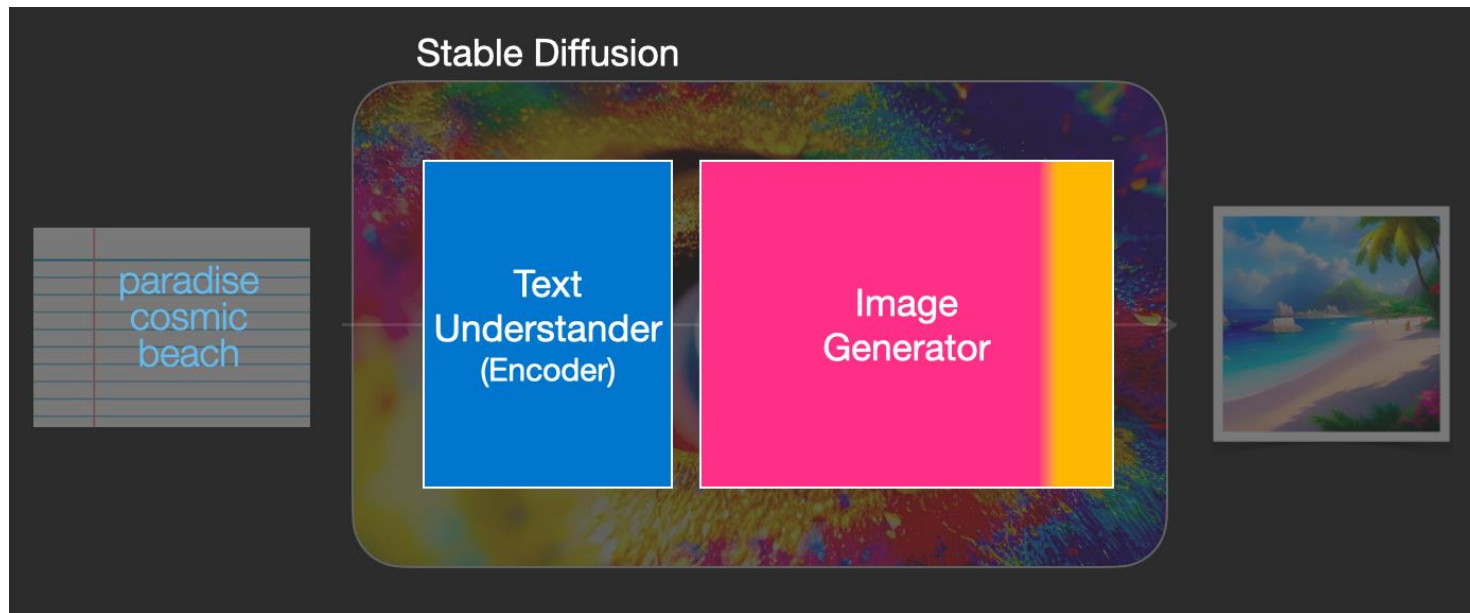
Stable Diffusion

Stable Diffusion is a deep learning, text-to-image model released in 2022 based on diffusion techniques

The logo for Stability.ai, featuring the word "stability" in a white, lowercase, sans-serif font, followed by a red dot and the letters ".ai" in the same font, all set against a solid purple rectangular background.

stability.ai

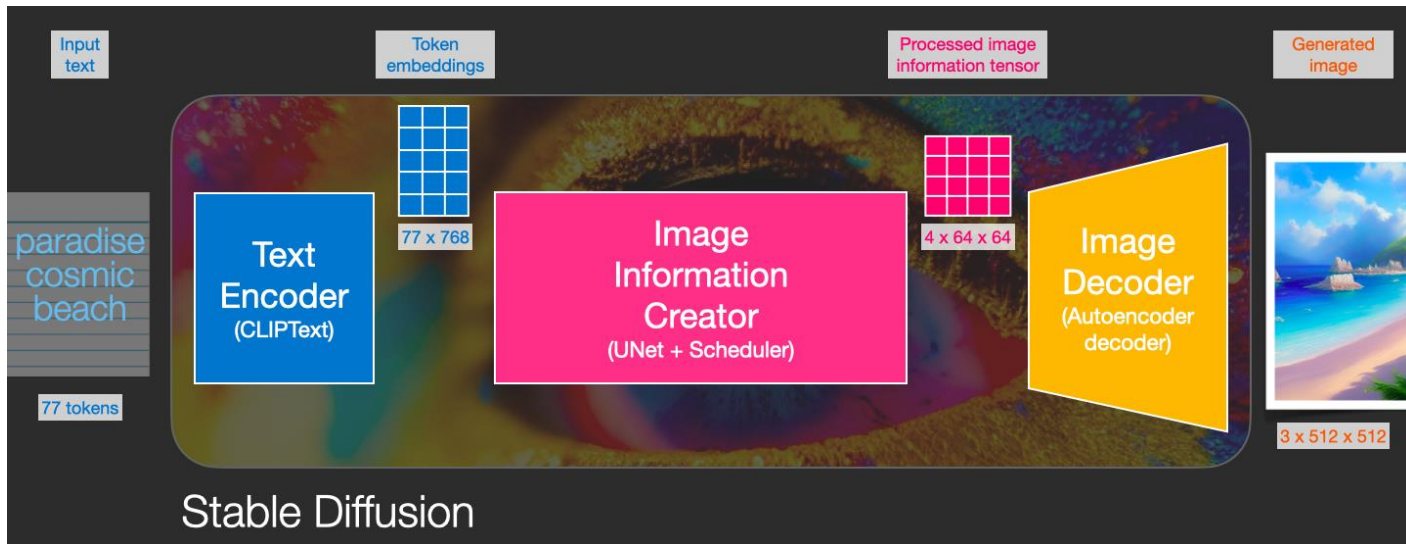
Text encoder is a special Transformer language model
(CLIP model)



ClipText for text encoding.

UNet + Scheduler to gradually process/diffuse information in the information (latent) space.

Autoencoder Decoder that paints the final image using the processed information array.



Stable Diffusion

paradise
cosmic
beach

77 tokens

Text
Encoder
(CLIPText)



Token
embeddings

1

2



Random image
information tensor

Image Information Creator
(UNet + Scheduler)

Diffusion

3



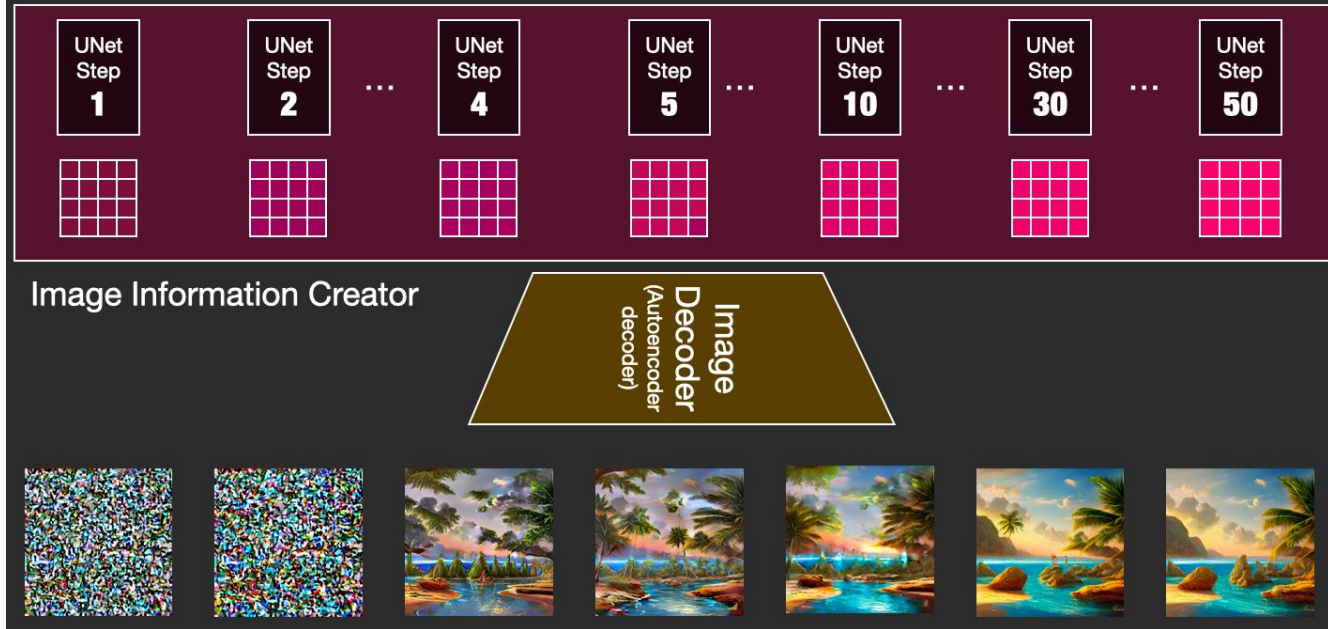
Processed image
information tensor

Image
Decoder
(Autoencoder
decoder)

Generated
image



Diffusion



Application Stable Diffusion

Sample input: "messi as a real madrid player"



Application Stable Diffusion



Prompt: 3D animation of a small, round, fluffy creature with big, expressive eyes explores a vibrant, enchanted forest.



Prompt: A cat waking up its sleeping owner demanding breakfast.