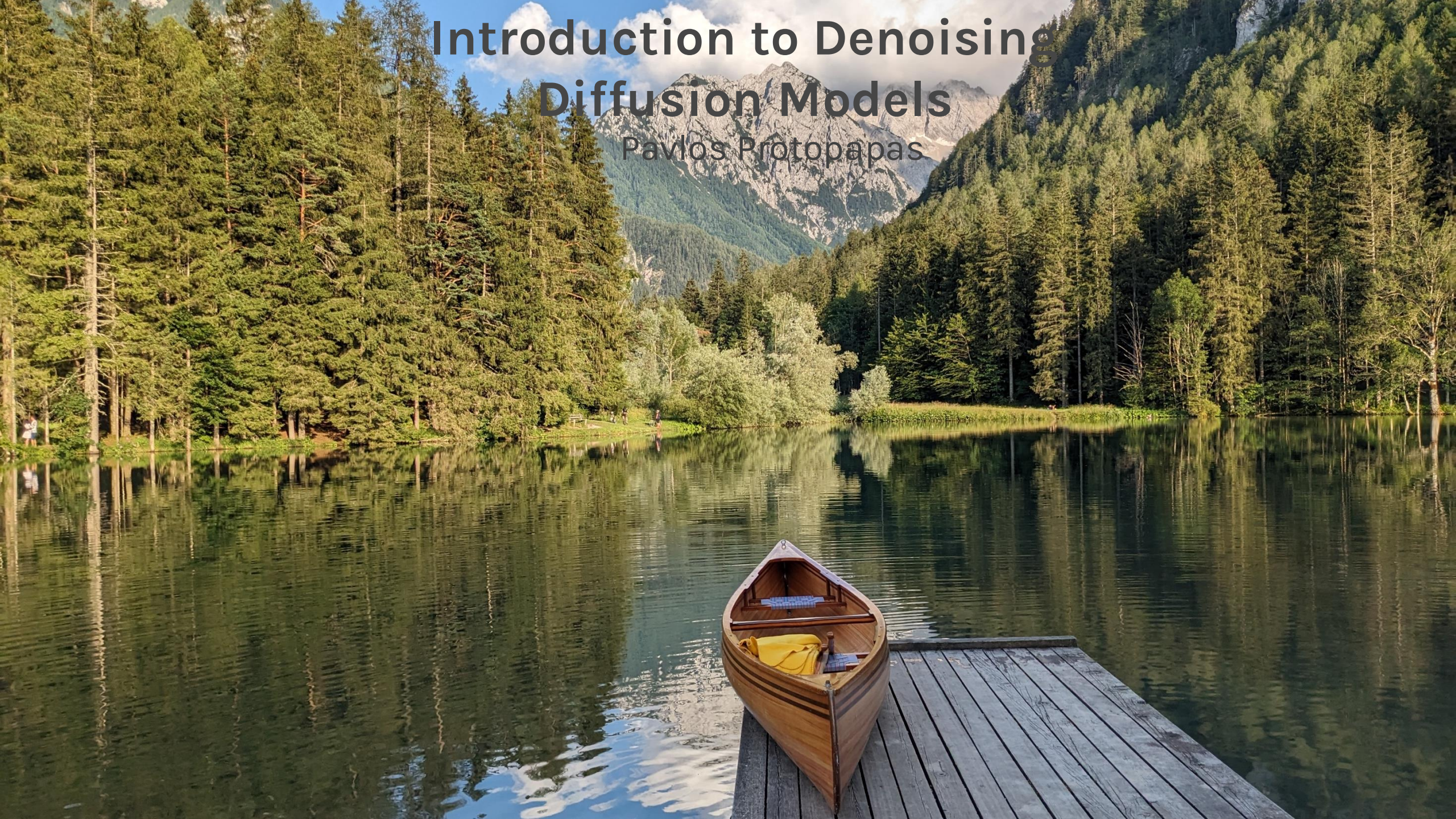


Introduction to Denoising Diffusion Models

Pavlos Protopapas



Outline

- Motivation for Diffusion Models
- Introduction to Denoising Diffusion Models
- Training a diffusion model
- Stable Diffusion
- DALL-E 2

Outline

- **Motivation for Diffusion Models**
- Introduction to Denoising Diffusion Models
- Training a diffusion model
- Stable Diffusion
- DALL-E 2

Motivation for Diffusion Models

Don't say you love the anime

"Professional photograph of bears in sports gear in a triathlon in Kyoto"



"Oriental Painting of Sun Tzu playing a game of Warcraft II on his desktop computer"



"Ukiyo-e painting of a cat hacker wearing VR headsets, on a postage stamp"



"Oriental painting of a dragon programming on a laptop in the Song dynasty, cyberpunk"



If you haven't read the manga

Deep Unsupervised Learning using Noisy Stochastic Thermodynamics

2.2. Reverse Trajectory

The generative distribution will be trained to describe the same trajectory, but in reverse,

$$p(x^{(1)}) = p(x^{(T)}) \quad (4)$$

$$p(x^{(1:T-1)}) = p(x^{(T)}) \prod_{t=1}^{T-1} p(x^{(t-1)}|x^{(t)}). \quad (5)$$

For both Gaussian and binomial diffusion, for continuous diffusion (limit of small step size), the reversal of the diffusion process has the identical functional form as the forward process (Fukunaga, 1999). Since $q(x^{(t)}|x^{(t-1)})$ is a Gaussian (binomial) distribution, and if Δt is small, then $q(x^{(t-1)}|x^{(t)})$ will also be a Gaussian (binomial) distribution. The longer the trajectory the smaller the diffusion rate β can be made.

During learning only the mean and covariance for a Gaussian diffusion kernel, or the bit flip probability for a binomial kernel, need be estimated. As shown in Table App. 1, $q(x^{(t)}|x^{(t-1)})$ and $q(x^{(t-1)}|x^{(t)})$ are functions of the mean and covariance of the reverse Markov transition for a Gaussian, and $f_1(x^{(t)}|x^{(t-1)})$ is a function providing the bit flip probability for a binomial distribution. The computational cost of running this algorithm is the cost of these functions, times the number of time-steps. For all results in this paper, multi-layer perceptrons are used to define these functions. A wide range of regression or function fitting techniques would be applicable however, including support vector methods.

2.3. Model Probability

The probability the generative model assigns to the data is

$$p(x^{(1:T)}) = \int dx^{(2:T-1)} p(x^{(1)}, x^{(T)}). \quad (6)$$

This can be evaluated rapidly by averaging over samples from the forward trajectory $q(x^{(1:T-1)}|x^{(T)})$. For instance, and if the forward and reverse distribution over trajectories can be made identical (see Section 2.2). If they are identical then only a single sample from $q(x^{(1:T-1)}|x^{(T)})$ is required to exactly evaluate the above integral, as can be seen by calculation. This corresponds to the case of a quasi-static process in statistical physics (Gymer & Ford, 2013; Jarzynski, 2011).

2.4. Training

Training involves maximizing the model log likelihood,

$$L = \int dx^{(1:T)} q(x^{(1:T)}) \log p(x^{(1:T)}) \quad (7)$$

$$= \int dx^{(1:T)} q(x^{(1:T)}) \log \left[\frac{\int dx^{(2:T-1)} q(x^{(1:T-1)}|x^{(T)}) p(x^{(1)}, x^{(T)})}{p(x^{(1:T)}) \prod_{t=1}^{T-1} q(x^{(t)}|x^{(t-1)})} \right] \quad (8)$$

which has a lower bound provided by Jensen's inequality,

$$L \geq \int dx^{(1:T)} q(x^{(1:T)}) \log \left[\frac{p(x^{(1)}, x^{(T)})}{\prod_{t=1}^{T-1} q(x^{(t)}|x^{(t-1)})} \right] \quad (9)$$

As described in Appendix A, for our diffusion trajectories this reduces to,

$$L \geq N \quad (10)$$

$$N = \sum_{t=2}^T \int dx^{(2:T-1)} q(x^{(2:T-1)}) \mathcal{D}_{KL} \left(q(x^{(2-1)}|x^{(2)}, x^{(T)}) \parallel p(x^{(1-1)}|x^{(1)}) \right)$$

	VP [42]	VE [42]	DDPM [35] + DDIM [40]	Others
Sampling (Section 3)	ODE solver	Euler	Euler	2 nd order Runge-Kutta
Time steps	$t \in [N]$	$t \in [N]$	$t \in [N]$	$t \in [N]$
Schedule	$\sigma(t) = \sqrt{t/N}$	$\sigma(t) = \sqrt{t/N}$	$\sigma(t) = \sqrt{t/N}$	$\sigma(t) = \sqrt{t/N}$
Scaling	$\sigma(t) = 1/\sqrt{t/N}$	$\sigma(t) = 1/\sqrt{t/N}$	$\sigma(t) = 1/\sqrt{t/N}$	$\sigma(t) = 1/\sqrt{t/N}$
Network and preconditioning (Section 5)	DDIM++	DDIM++	DDIM++	(any)
Architecture of P_θ	DDIM++	DDIM++	DDIM++	(any)
Skip-connection	$c_{sk}(x) = 1$	$c_{sk}(x) = 1$	$c_{sk}(x) = 1$	$c_{sk}(x) = 1$
Output scaling	$c_{out}(x) = \sigma$	$c_{out}(x) = \sigma$	$c_{out}(x) = \sigma$	$c_{out}(x) = \sigma$
Input scaling	$c_{in}(x) = 1/\sqrt{\sigma^2 + 1}$	$c_{in}(x) = 1/\sqrt{\sigma^2 + 1}$	$c_{in}(x) = 1/\sqrt{\sigma^2 + 1}$	$c_{in}(x) = 1/\sqrt{\sigma^2 + 1}$
Noise cond.	$c_{cond}(x) = [d-1]\sigma^{-1}(\sigma)$	$c_{cond}(x) = [d-1]\sigma^{-1}(\sigma)$	$c_{cond}(x) = [d-1]\sigma^{-1}(\sigma)$	$c_{cond}(x) = [d-1]\sigma^{-1}(\sigma)$
Training (Section 5)	Noise distribution	Noise distribution	Noise distribution	Noise distribution
Loss weighting	$N(x)$	$N(x)$	$N(x)$	$N(x)$
Parameters	$\beta_0 = 10^{-4}, \beta_{max} = 0.1$ $\sigma_0 = 10^{-4}, \sigma_{max} = 0.1$ $N = 1000$	$\beta_0 = 10^{-4}, \beta_{max} = 0.1$ $\sigma_0 = 10^{-4}, \sigma_{max} = 0.1$ $N = 1000$	$\beta_0 = 10^{-4}, \beta_{max} = 0.1$ $\sigma_0 = 10^{-4}, \sigma_{max} = 0.1$ $N = 1000$	$\beta_0 = 10^{-4}, \beta_{max} = 0.1$ $\sigma_0 = 10^{-4}, \sigma_{max} = 0.1$ $N = 1000$

[†] DDPM also employs a second loss term L_{loss} . From tests, $\beta_0 = 8$ yielded better FID than $\beta_0 = 1$ used by DDPM.

Motivation for Diffusion Models

Denoising diffusion models have shown impressive sample quality and diversity for several [image generation](#) tasks.



Motivation for Diffusion Models

Denoising diffusion models have shown impressive results for [text-to-image](#) tasks.

[DALL-E 2](#)



“A Univ.AI student about to write a paper on neural differential equations painted in a pointilistic style”

[Imagen](#)



“A robot couple fine dining with the Eiffel Tower in the background”

Motivation for Diffusion Models

Denoising diffusion models have shown impressive results for [image editing](#) and [inpainting](#) tasks.



“A man with red hair”

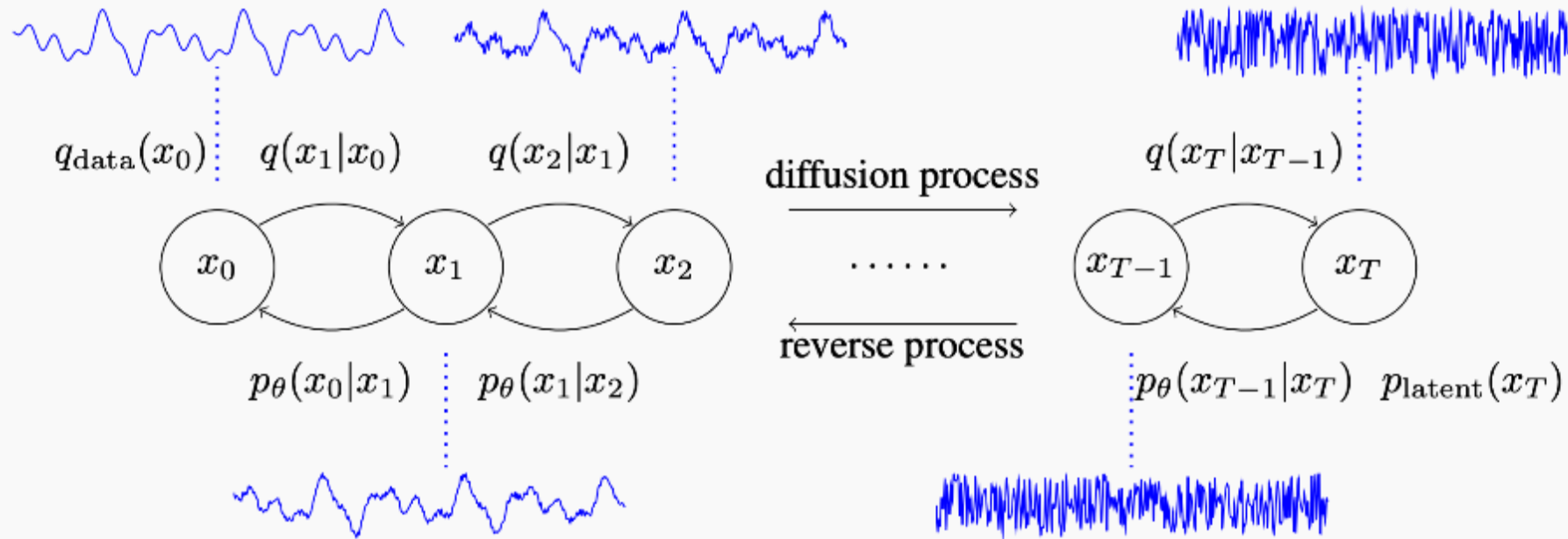


“Zebras roaming the field”

[GLIDE](#)

Motivation for Diffusion Models

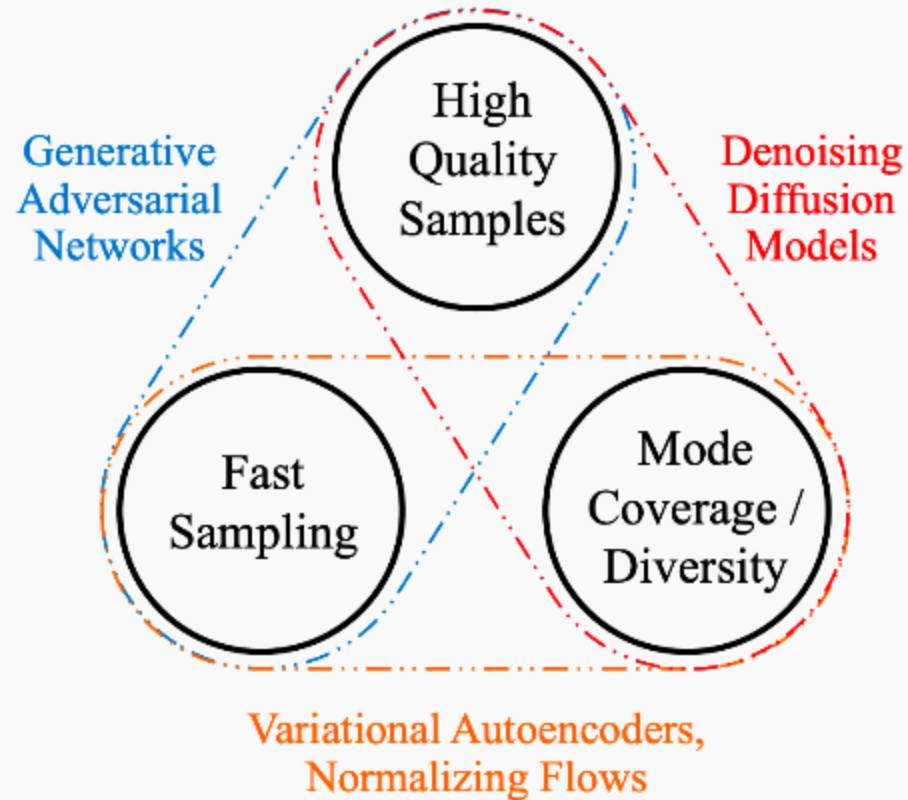
Denoising diffusion models have also been effective in [non-visual domains](#).



[Diffwave: A versatile diffusion model for audio synthesis](#)

Motivation for Diffusion Models

But are denoising diffusion models all you need?



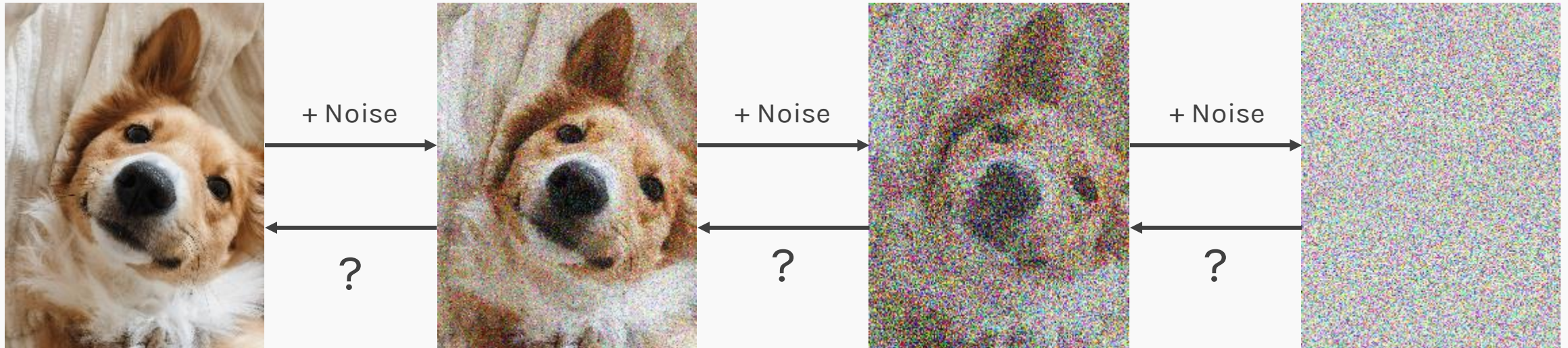
Generative learning trilemma
[[Source](#)]

Outline

- Motivation for Diffusion Models
- **Introduction to Denoising Diffusion Models**
- Training a diffusion model
- Stable Diffusion
- DALL-E 2

Denoising Diffusion Models

The general idea of diffusion models is quite simple.

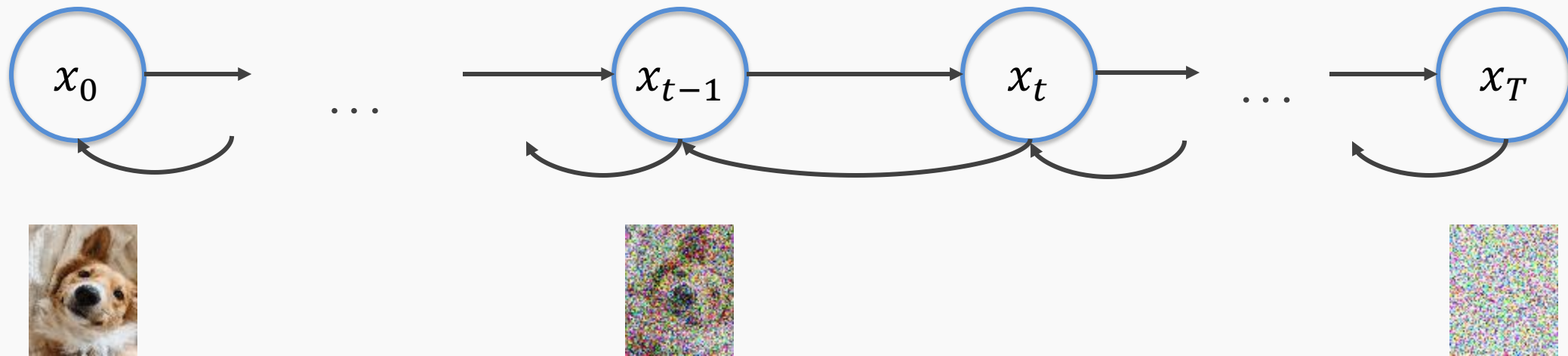


Diffusion models help undoing this process: Starting with noise, it can gradually move toward a coherent image.

Denoising Diffusion Models

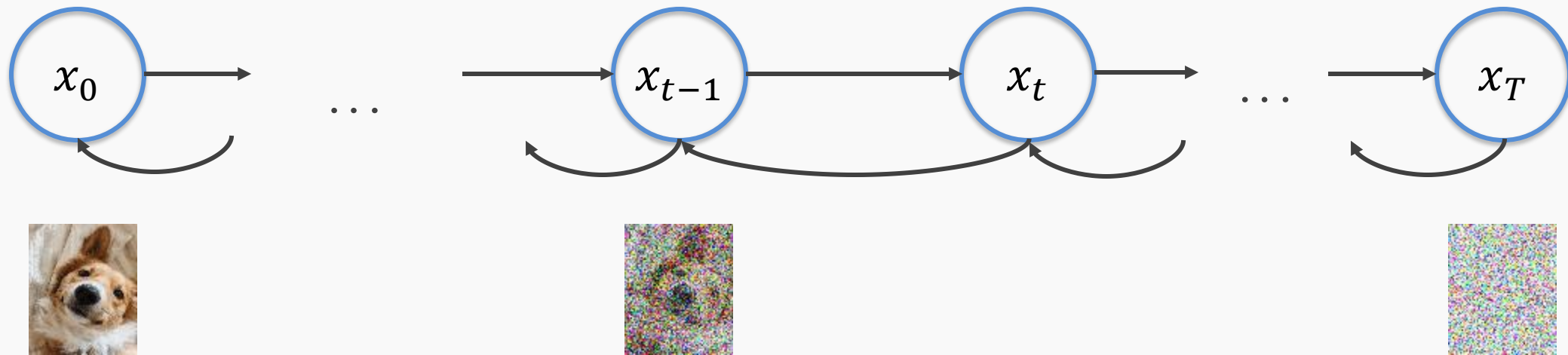
A Markov chain means that each step only depends on the previous one.

In practice, they are formulated using a Markov chain of T steps.



- They take the input image and gradually add Gaussian noise to it through a series of small steps.
- A neural network is trained to recover the original data by reversing the noising process. By being able to model the reverse process, we can generate new data.

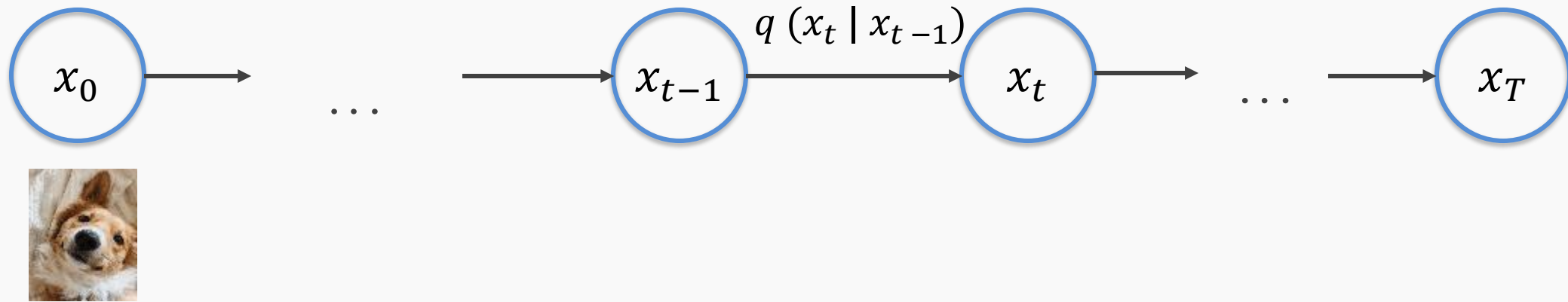
Denoising Diffusion Models



Therefore, there are two processes in a diffusion model:

1. Forward process
2. Reverse process

Denoising Diffusion Models – Forward Process



Given a data point x_0 sampled from the real data distribution $q(x)$.

At each forward step of the Markov chain, we add a small amount of Gaussian noise with variance β_t to x_{t-1} , producing a new latent variable x_t with distribution $q(x_t | x_{t-1})$.

This diffusion process can be formulated as follows:

$$q(x_t | x_{t-1}) = N(\mu_t = \sqrt{1 - \beta_t} x_{t-1}, \Sigma_t = \beta_t I)$$

As $T \rightarrow \infty$, you eventually end up with an isotropic gaussian (i.e. pure random noise $\sim N(0, I)$).

Denoising Diffusion Models – Forward Process

In general:

$$q(x_{1:T} | x_0) = \prod_{t=1}^T q(x_t | x_{t-1})$$

where:

$$q(x_t | x_{t-1}) = N(\mu_t = \sqrt{1 - \beta_t} x_{t-1}, \Sigma_t = \beta_t I)$$

But how would you find $q(x_{1000} | x_0)$?

You would have to go from x_0 till x_{999} step-by-step until you can compute $q(x_{1000} | x_{999})$!

Denoising Diffusion Models – Forward Process

By reparametrizing such that:

$$\alpha_t = 1 - \beta_t$$
$$\overline{\alpha}_t = \prod_{i=1}^t \alpha_i$$

We get:

$$q(x_t | x_0) = N(\mu_t = \sqrt{\overline{\alpha}_t} x_0, \Sigma_t = (1 - \overline{\alpha}_t) I)$$

Using the above equation, we can now generate any time step t directly from x_0 instead of iteratively going through x_{t-1} , x_{t-2} ... and so on.

Denoising Diffusion Models – Forward Process

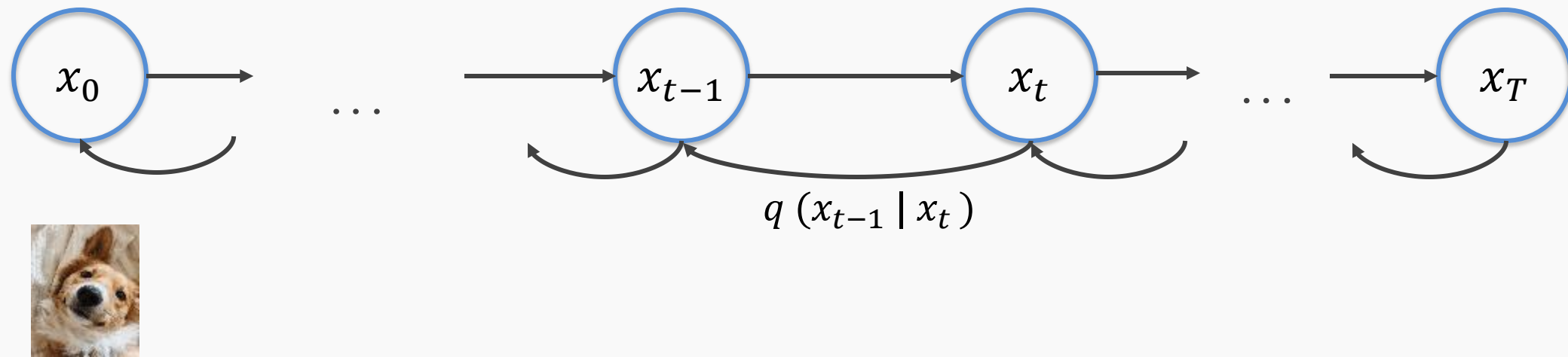
- The variance parameter β_t can be fixed to a constant or chosen as a schedule over the T timesteps.
- The original DDPM authors utilized a linear schedule, however the Improving DDPM paper showed that employing a cosine schedule works even better.



Linear Scheduler (top), Cosine scheduler (bottom)

[Source](#)

Denoising Diffusion Models – Reverse Process



- The goal of a diffusion model is to learn a neural network to reverse the denoising process by iteratively undoing the forward process.
- If we learn the reverse distribution $q(x_{t-1} | x_t)$, we can then:
 1. Sample x_T from $N(0, I)$
 2. Run the reverse process iteratively and acquire a sample from $q(x_0)$, generating a novel data point from the original data distribution.

Denoising Diffusion Models – Reverse Process

But finding $q(x_{t-1} | x_t)$ is difficult as it depends on the entire data distribution and is therefore intractable.

Turns out that for small enough β_t , $q(x_{t-1} | x_t)$ will also be Gaussian.

We can then approximate $q(x_{t-1} | x_t)$ with a gaussian distribution $p_\theta(x_{t-1} | x_t)$ which is parameterized by some θ :

$$p_\theta(x_{t-1} | x_t) = N(\mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$$

and

$$p_\theta(x_{0:T}) = p_\theta(x_T) \prod_{t=1}^T p_\theta(x_{t-1} | x_t)$$

Denoising Diffusion Models – Reverse Process

How do you train a model to learn these parameters θ ? What is the loss function we need to optimize?

We can use the maximum likelihood estimate by optimizing the negative log-likelihood of the training data.

Using variational inference and the evidence lower bound, we get our loss to be:

$$L = \mathbb{E}_{q(x_{0:T})} \left[\log \frac{q(x_{1:T}|x_0)}{p_\theta(x_{0:T})} \right]$$

On further simplification, we get:

$$L = \mathbb{E}[D_{KL}(q(x_T|x_0) \parallel p_\theta(x_T)) + \sum_{t=2}^T D_{KL}(q(x_{t-1}|x_t, x_0) \parallel p_\theta(x_{t-1}|x_t)) - \log p_\theta(x_0|x_1)]$$

Denoising Diffusion Models – Reverse Process

$$L = \underbrace{\mathbb{E}[D_{KL}(q(x_T|x_0) \parallel p(x_T))]}_{L_T} + \sum_{t=2}^T \underbrace{D_{KL}(q(x_{t-1}|x_t, x_0) \parallel p_{\theta}(x_{t-1}|x_t))}_{L_{t-1}} \underbrace{- \log p_{\theta}(x_0|x_1)}_{L_0}$$

Therefore,

$$L = L_T + L_{T-1} + L_{T-2} + \dots + L_0$$

where

$$L_T = D_{KL}(q(x_T|x_0) \parallel p(x_T))$$

$$L_t = D_{KL}(q(x_t|x_{t+1}, x_0) \parallel p_{\theta}(x_t|x_{t+1})) \text{ for } 1 \leq t \leq T-1$$

$$L_0 = -\log p_{\theta}(x_0|x_1)$$

Denoising Diffusion Models – Reverse Process

$$L = \underbrace{\mathbb{E}[D_{KL}(q(x_T|x_0) \parallel p(x_T))]}_{L_T} + \sum_{t=2}^T \underbrace{D_{KL}(q(x_{t-1}|x_t, x_0) \parallel p_{\theta}(x_{t-1}|x_t))}_{L_{t-1}} - \underbrace{\log p_{\theta}(x_0|x_1)}_{L_0}$$

Here, L_T is constant and does not vary with the parameter we are trying to optimize over θ .

Thus, we can ignore this term to get:

$$L = \sum_{t=2}^T D_{KL}(q(x_{t-1}|x_t, x_0) \parallel p_{\theta}(x_{t-1}|x_t)) - \log p_{\theta}(x_0|x_1)$$

Denoising Diffusion Models – Reverse Process

Looking only at term $L_t = D_{KL}(q(x_t|x_{t+1}, x_0) \parallel p_\theta(x_t|x_{t+1}))$ for $1 \leq t \leq T - 1$


$$q(x_t|x_{t+1}, x_0) = N(\tilde{\mu}_t(x_t, x_0), \tilde{\beta}_t I)$$

Where we reparametrize the normal distribution:

$$\tilde{\mu}_t = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right)$$

$$\tilde{\beta}_t = \beta_t \left(\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \right)$$

where

$$\alpha_t = 1 - \beta_t, \quad \bar{\alpha}_t = \prod_{i=1}^t \alpha_i$$

$\epsilon_t \sim N(0, I)$, β_t is the variance scheduler

$$p_\theta(x_{t-1} | x_t) = N(\mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$$

where

$\mu_\theta(x_t, t)$ is learnt by the neural network

$\Sigma_\theta(x_t, t)$ is fixed as $\tilde{\beta}_t$ or β_t

Denoising Diffusion Models – Reverse Process

$$\tilde{\mu}_t = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right) \quad \left| \quad \mu_\theta(x_t, t) \text{ is learnt by the neural network} \right.$$

Our goal is to now reduced to ensuring that $\mu_\theta(x_t, t) \approx \tilde{\mu}_t$.

We can do this by minimizing the MSE between the two:

$$\begin{aligned} L_t &= \frac{1}{2\sigma^2} \| \tilde{\mu}_t - \mu_\theta(x_t, t) \|^2 \\ &= \frac{1}{2\sigma^2} \left\| \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right) - \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t, t) \right) \right\|^2 \end{aligned}$$

Denoising Diffusion Models – Reverse Process

$$= \frac{1}{2\sigma^2} \left\| \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right) - \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(x_t, t) \right) \right\|^2$$
$$L_t = \frac{\beta_t^2}{2\sigma^2 \alpha_t (1 - \bar{\alpha}_t)} \left\| \epsilon - \epsilon_{\theta}(x_t, t) \right\|^2$$

The authors have found that dropping the constant term have led to better sample quality, so:

$$L_t = \left\| \epsilon - \epsilon_{\theta}(x_t, t) \right\|^2$$

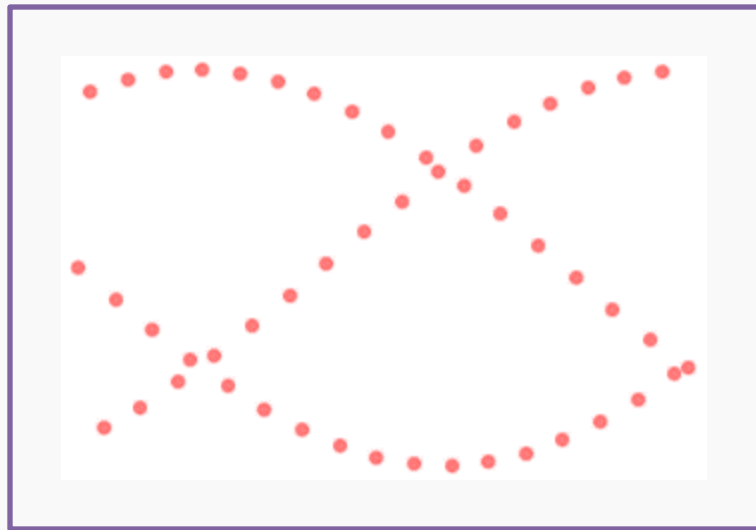
$$\text{where } x_t = \sqrt{\alpha_t} x_0 + \sqrt{1 - \alpha_t} \epsilon$$

Outline

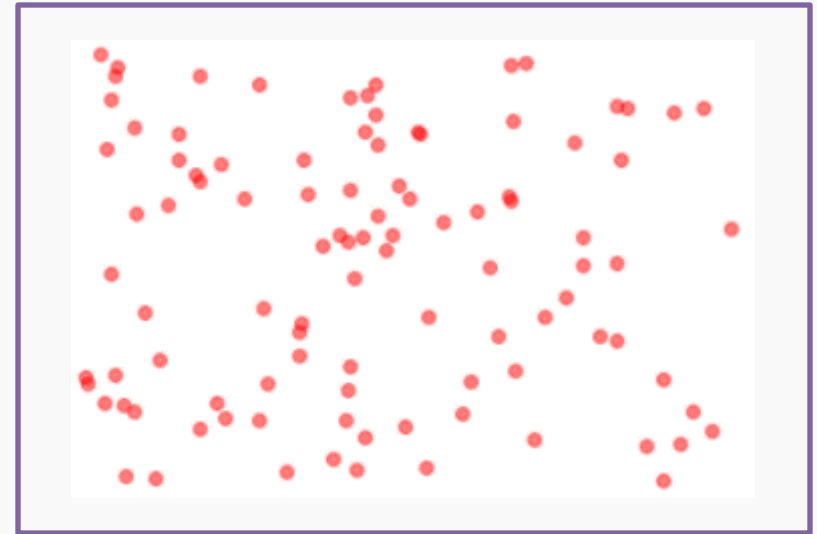
- Motivation for Diffusion Models
- Introduction to Denoising Diffusion Models
- **Training a diffusion model**
- Stable Diffusion
- DALL-E 2

Training a diffusion model

Abstractly, our goal is to learn noise from structured data, which is nontrivial.



Structured data



Noise

Training a diffusion model

Now that we have found our loss function:

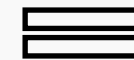
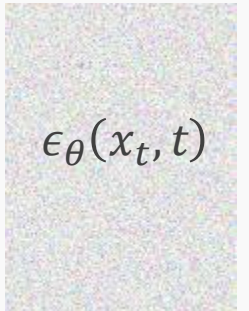
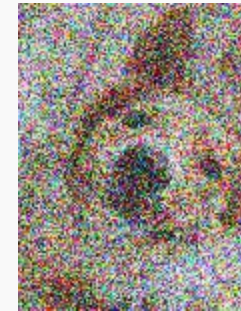
$$L_t = \| \epsilon - \epsilon_{\theta}(x_t, t) \|^2$$

$$\text{where } x_t = \sqrt{\alpha_t} x_0 + \sqrt{1 - \alpha_t} \epsilon$$

We can use a neural network to take in x_t and predict a noise $\epsilon_{\theta}(x_t, t)$ that describes the noise that was added to x_{t-1} to get x_t .

Therefore, by passing x_t through our U-Net architecture, we get a noise that we can subtract from x_t to get x_{t-1} . By doing this iteratively, we can get back x_0 .

Given x_t , the model must predict the noise $\epsilon_{\theta}(x_t, t)$ such that
$$x_t - \epsilon_{\theta}(x_t, t) = x_{t-1}$$



Algorithm 1 Training

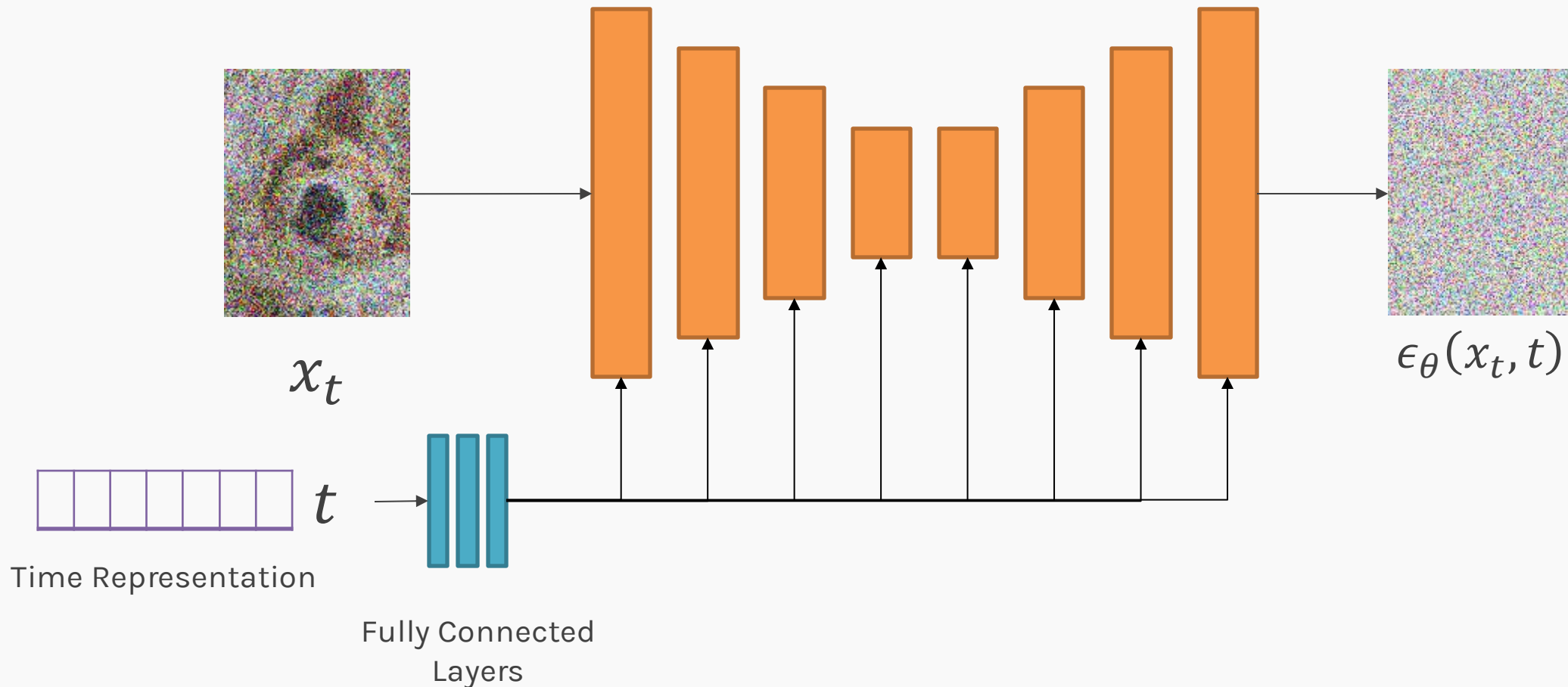
- 1: **repeat**
 - 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
 - 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
 - 4: $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 5: Take gradient descent step on
$$\nabla_{\theta} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t}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

Algorithm 2 Sampling

```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$   
2: for  $t = T, \dots, 1$  do  
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  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  $\mathbf{x}_0$ 
```

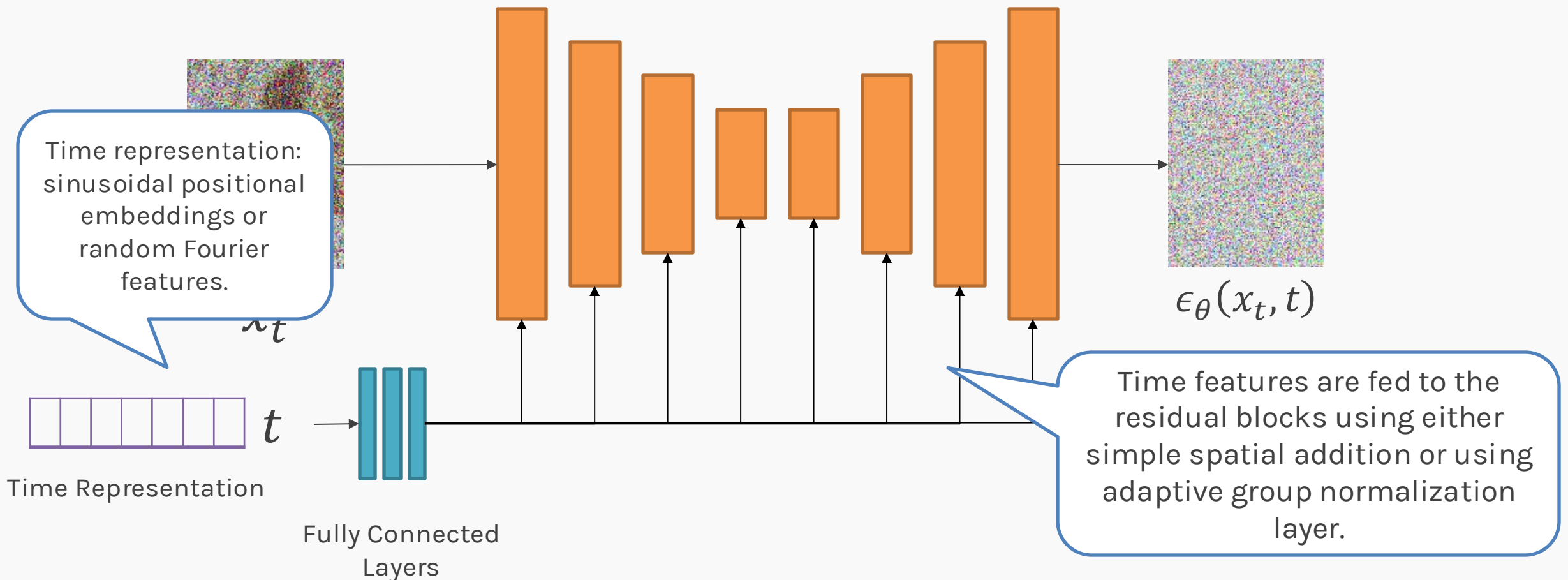
Training a diffusion model

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



Training a diffusion model

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



Outline

- Motivation for Diffusion Models
- Introduction to Denoising Diffusion Models
- Training a diffusion model
- **Stable Diffusion**
- DALL-E 2

Examples of recent diffusion models

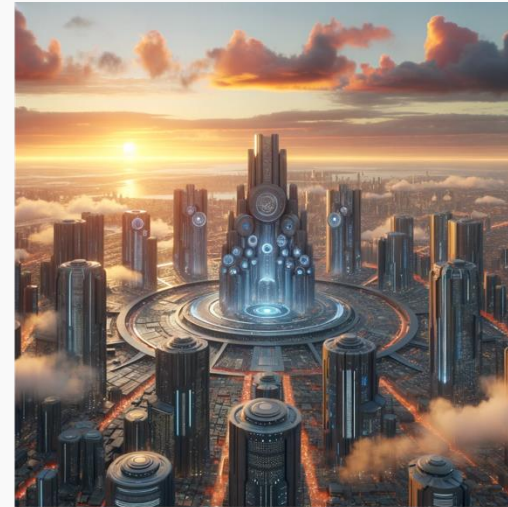
There are many recent text-to-image or image generation diffusion models:

- DALL-E 2
- DALL-E 3
- GLIDE
- Imagen
- Imagen 2
- MidJourney
- Stable Diffusion 1.x
- Stable Diffusion 2.x
- Adobe Firefly
- And many more!

"A futuristic cityscape at sunset."



Imagen 2



DALL-E 3



Adobe Firefly

In today's class, we will be talking about Stable Diffusion and DALL-E 2, once you understand this, you can understand other diffusion models!

Latent Diffusion Models

- Training models in pixel space is very computationally expensive as images are high dimensional data.
- Latent diffusion models tackle this issue. Instead of applying the diffusion process directly on a high-dimensional input, they project the input into a smaller latent space and apply the diffusion there.

Latent Diffusion Models

There are two stages of latent diffusion models:

1. Train perceptual compression models to learn a semantically equivalent latent space by removing high-level details. The loss includes reconstruction, adversarial, and regularization terms.
2. Perform diffusion process in the latent space for efficiency and focusing on relevant semantic data bits.

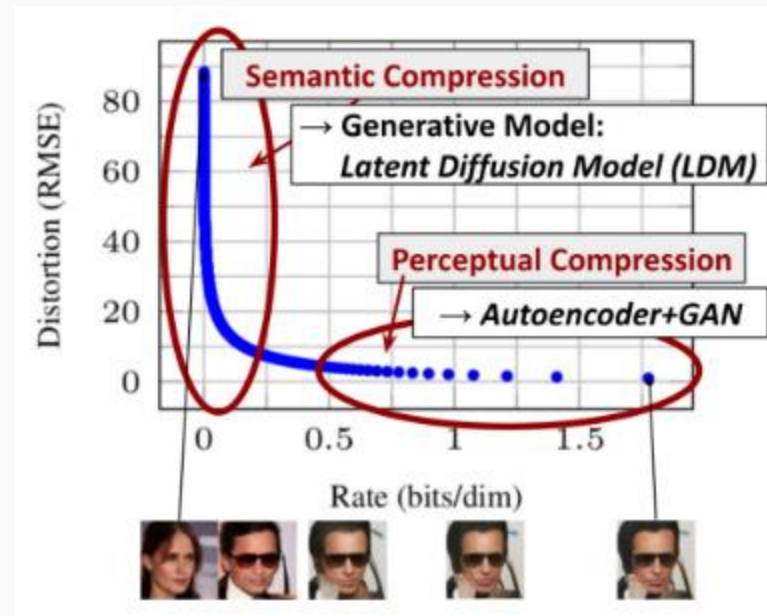


Image from Rombach et al. in
*High-resolution image synthesis
with latent diffusion models*

Stable Diffusion ([R Rombach et al. 2021](#))

- Stable Diffusion (SD) is a family of open-source conditional **latent diffusion** models, able to generate high-quality images from various conditions (such as textual descriptions).



Stable Diffusion

- The model was trained to generate realistic images from textual descriptions, using a subset of LAION-5B, which is an open dataset of (image, caption) pairs.


Backend url:

Index:


Search: 🔍 📷 ⬇️

[Clip retrieval](#) works by converting the text query to a CLIP embedding, then using that embedding to query a knn index of clip image embeddings


Display captions ☒
Display full captions ☐
Display similarities ☐
Safe mode ☒
Hide duplicate urls ☒
Hide (near) duplicate images ☒
Search over
Search with multilingual clip ☐




french cat




french cat




How to tell if your feline is french. He wears a b...




イケメン猫モデル「トキ・ナンタケツト」がかっこいい - NAVER まとめ




Hilarious pics of funny cats! funnycatsgif.com




Hipster cat




網友挑戰「加幾筆畫出最創意貓咪圖片」，笑到岔氣之後我也手...



cat in a suit Georgian sells tomatoes





French Bread Cat Loaf Metal Print

Stable Diffusion – The different parts

A **Variational Autoencoder** (VAE) that is used to both map an image to a lower dimensional latent representation (using the Encoder), and to go from the latent representation to an image (using the Decoder).

A **UNet** which is used to predict the noise added to an input image in each time step.

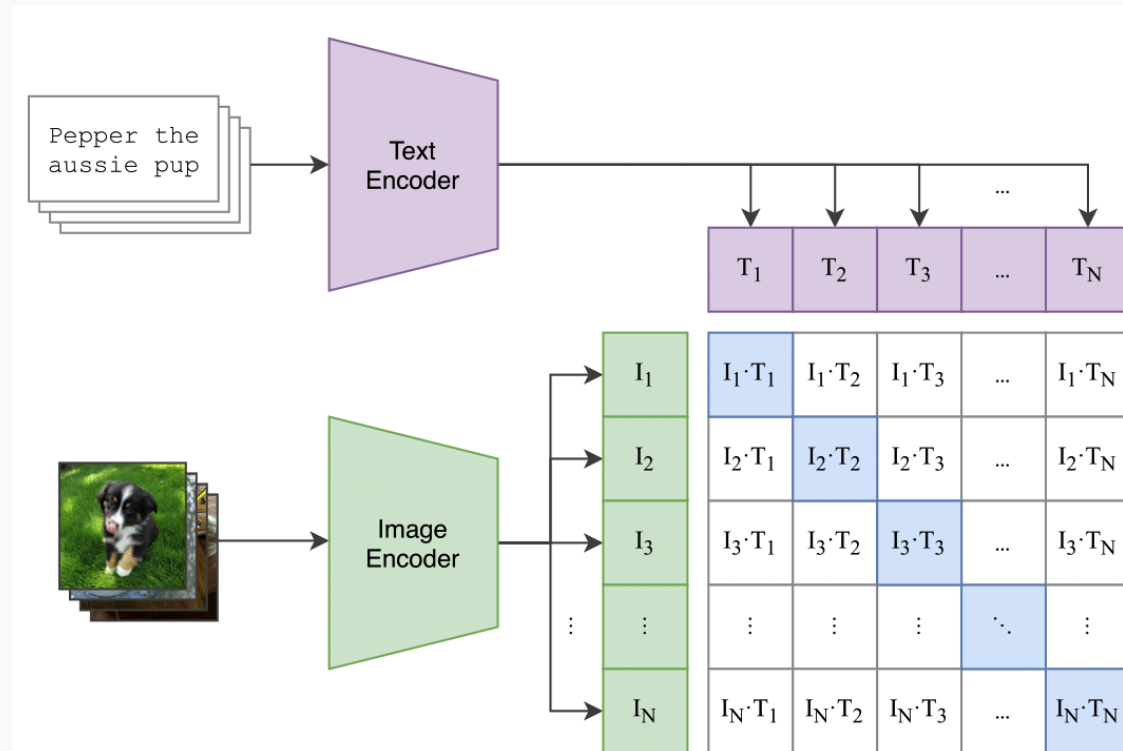
The **CLIP** model, particularly the text encoder, to incorporate textual information which are rich in image as well as conditional information for the model.

SD – Variational Autoencoder

- The main ideas behind latent diffusion models is to perform the diffusion process in a latent space, instead of directly using pixel space.
- In SD, an image $x_0 \in \mathbb{R}^{512 \times 512 \times 3}$ is first mapped using the **Encoder** to a low dimensional representation $z_0 \in \mathbb{R}^{64 \times 64 \times 4}$, where the diffusion process happens, obtaining a noisy latent z_T .
- Now the goal is to denoise z_T and map it back to an image using the **Decoder**.
- $512 \times 512 \times 3 \approx 768\text{k}$ and $64 \times 64 \times 4 \approx 16\text{k}$, this is a reduction of **48 times** in the dimension of the input!

SD – CLIP

- To be able to generate images conditioned on text, SD uses CLIP, a model jointly trained to learn image and text representations at the same time.
- The idea is that the text embeddings are rich in visual content as well and can help guide the image generation towards the given textual description.



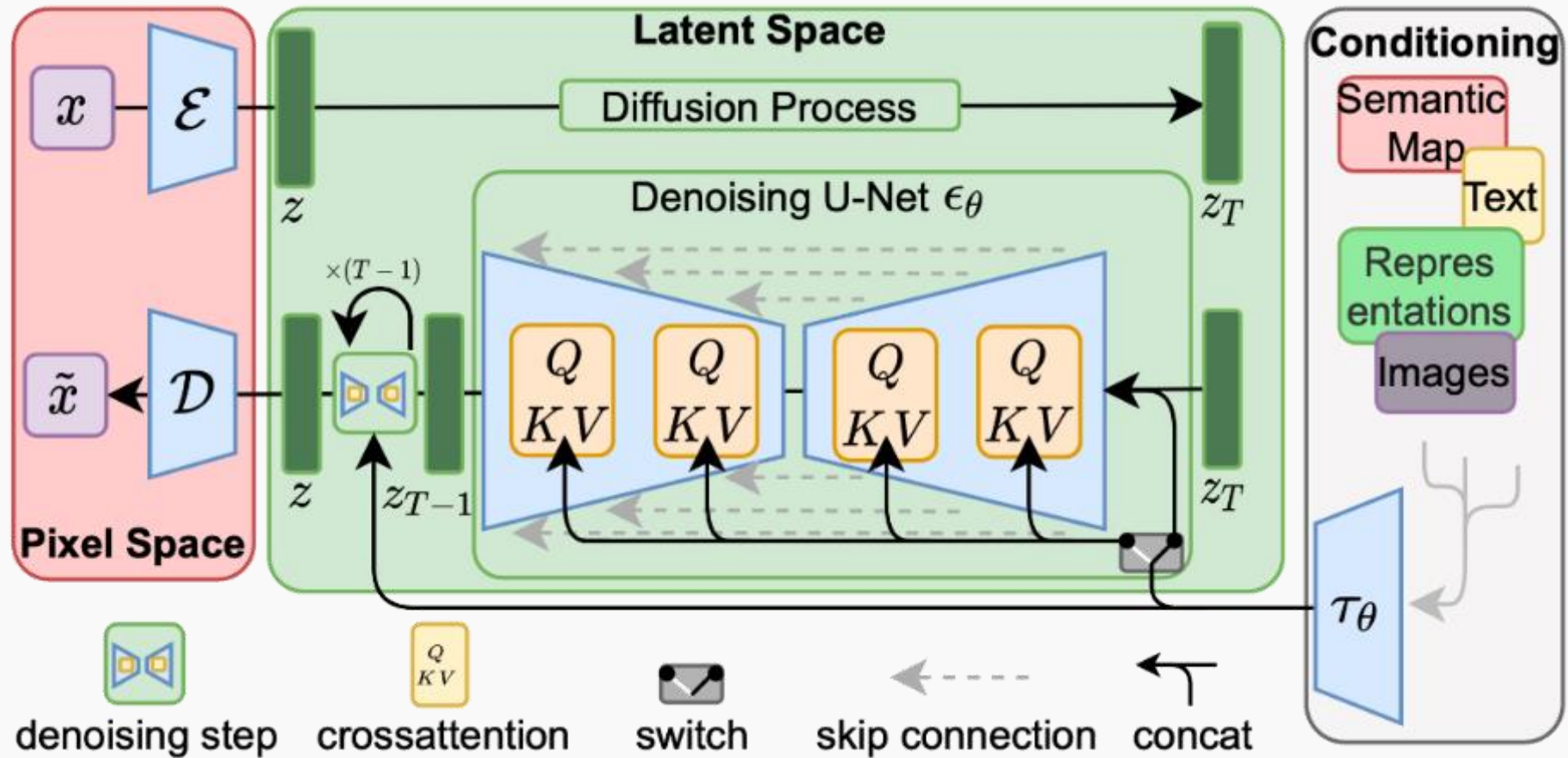
SD – CLIP

- CLIP was trained on a dataset of 400 million (image, text) pairs, and is composed of an image transformer and a text transformer.
- During training if we have image text pairs $\{(x_i, y_i)\}_{i=1}^N$, we first use the transformers to obtain representations $\{(I_i, T_i)\}_{i=1}^N$.
- We then maximize the similarity between the embeddings of (image, text) pairs that are present in the dataset, and minimize the similarities of those which are not

SD – UNet

- Finally, we have the UNet backbone, which it is the main part of the diffusion process.
- Given a noisy latent z_t , a time step t and the CLIP text embedding $\tau_\theta(y)$ of a prompt y , the UNet predicts the noise $\epsilon_\theta(z_t, t, \tau_\theta(y))$ which was added to go from z_{t-1} to z_t .
- The UNet uses cross-attention, where the queries refer to intermediate representations of the UNet, and the key and values to the conditional information $\tau_\theta(y)$

Stable Diffusion – Putting everything together



Outline

- Motivation for Diffusion Models
- Introduction to Denoising Diffusion Models
- Training a diffusion model
- Latent Diffusion Models
- **DALL-E 2**

DALL-E 2 ([Ramesh et al. 2022](#))



a shiba inu wearing a beret and black turtleneck



a close up of a handpalm with leaves growing from it

DALL-E 2

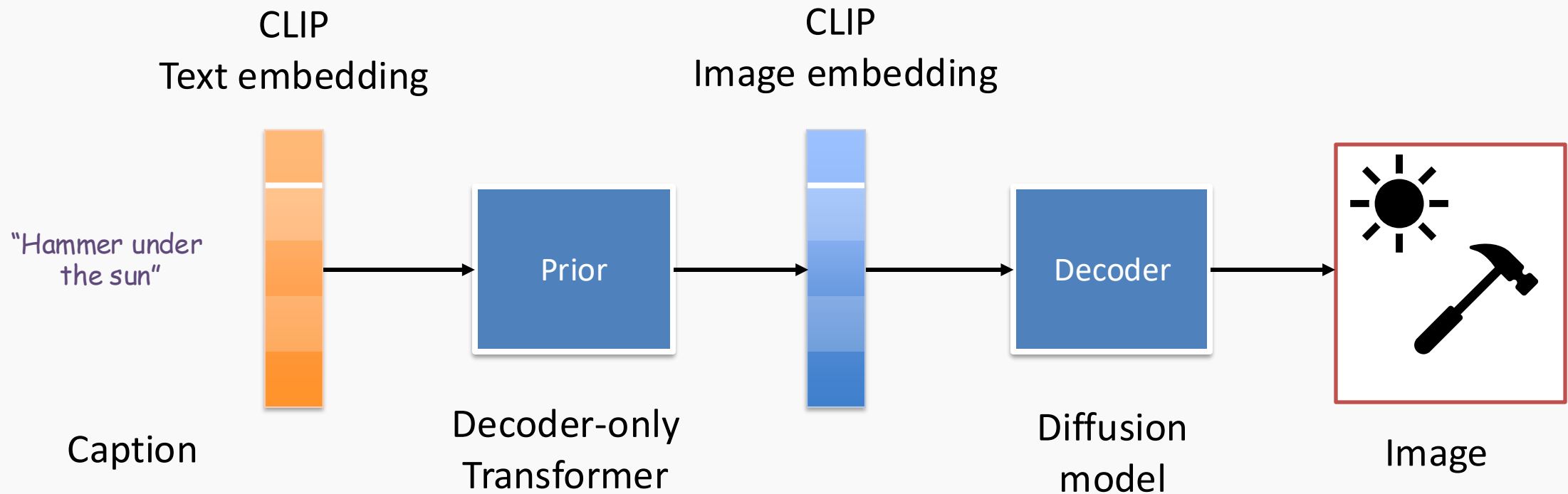
From a high-level, DALL-E 2's works in three steps:

1. A text prompt is input into a **text encoder** that is trained to map the prompt y to a text encoding z_t .
2. A diffusion model $p(z_i | y)$ called the **prior** maps the text encoding to a corresponding **image encoding** z_i that captures the semantic information of the prompt contained in the text encoding.
3. An **image decoder** $p(x | z_i, y)$ stochastically generates an image which is a visual manifestation of this semantic information.
4. The full conditional generative model is defined by:

$$p(x | y) = p(x | z_i, y)p(z_i | y)$$

DALL-E 2

- CLIP is used to map a caption y to a text embedding z_t which is fed to the prior to generate the image embedding z_i .
- The authors used both an autoregressive prior and a diffusion prior (decoder only transformer), but found the diffusion prior worked better.



DALL-E 2 Prior

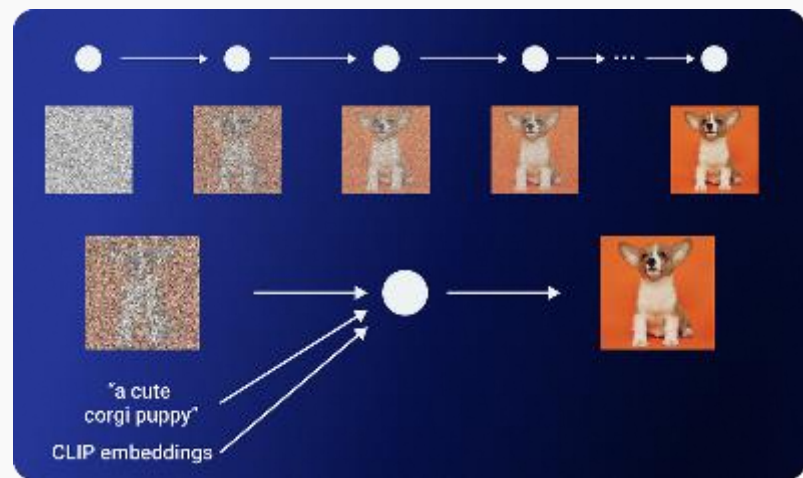
- The use of a prior is not strictly necessary. For example, the model could condition only on the caption. However, the authors found that the prior increased image diversity.
- It is a decoder only Transformer of width 2048 and 24 blocks, inputting:
 - Encoded text
 - CLIP text embedding
 - Embedding for the diffusion timestep
 - Noised CLIP image embedding
 - A final embedding whose output from the Transformer is used to predict the unnoised CLIP image embedding

DALL-E 2 Prior

- The idea is that we use a transformer to roughly translate and decode the text embedding into the image embedding.
- Unlike in typical stable diffusion, the authors use a loss function to predict the unnoised image embedding directly as opposed to predicting the ϵ noise.
- At sampling time, they generate two candidate image embeddings, and take the one that is more aligned with the text embedding (maximum dot product between the two embeddings, which are vectors).

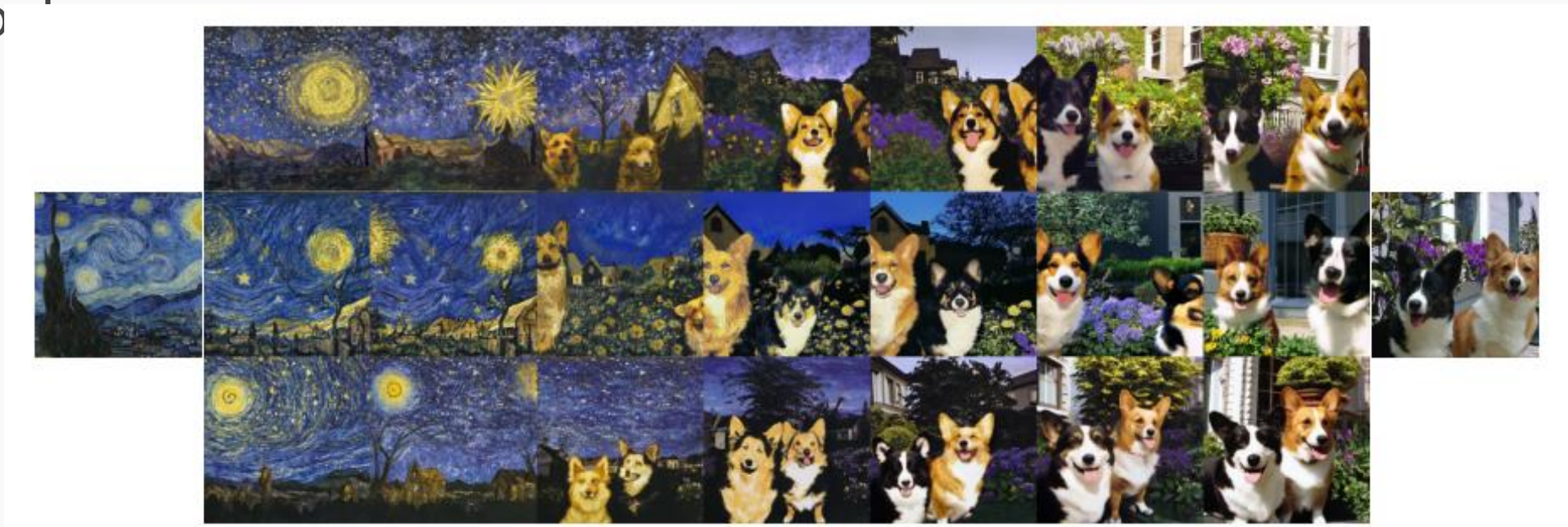
DALL-E 2 Decoder

- The decoder is a diffusion model. The goal of the decoder is to turn the image embeddings into actual images.
- They also optionally condition on the text captions in classifier-free guidance 50% of the time to reduce the variance of the outputs, and randomly setting the CLIP image embeddings to 0 about 10% of the time.

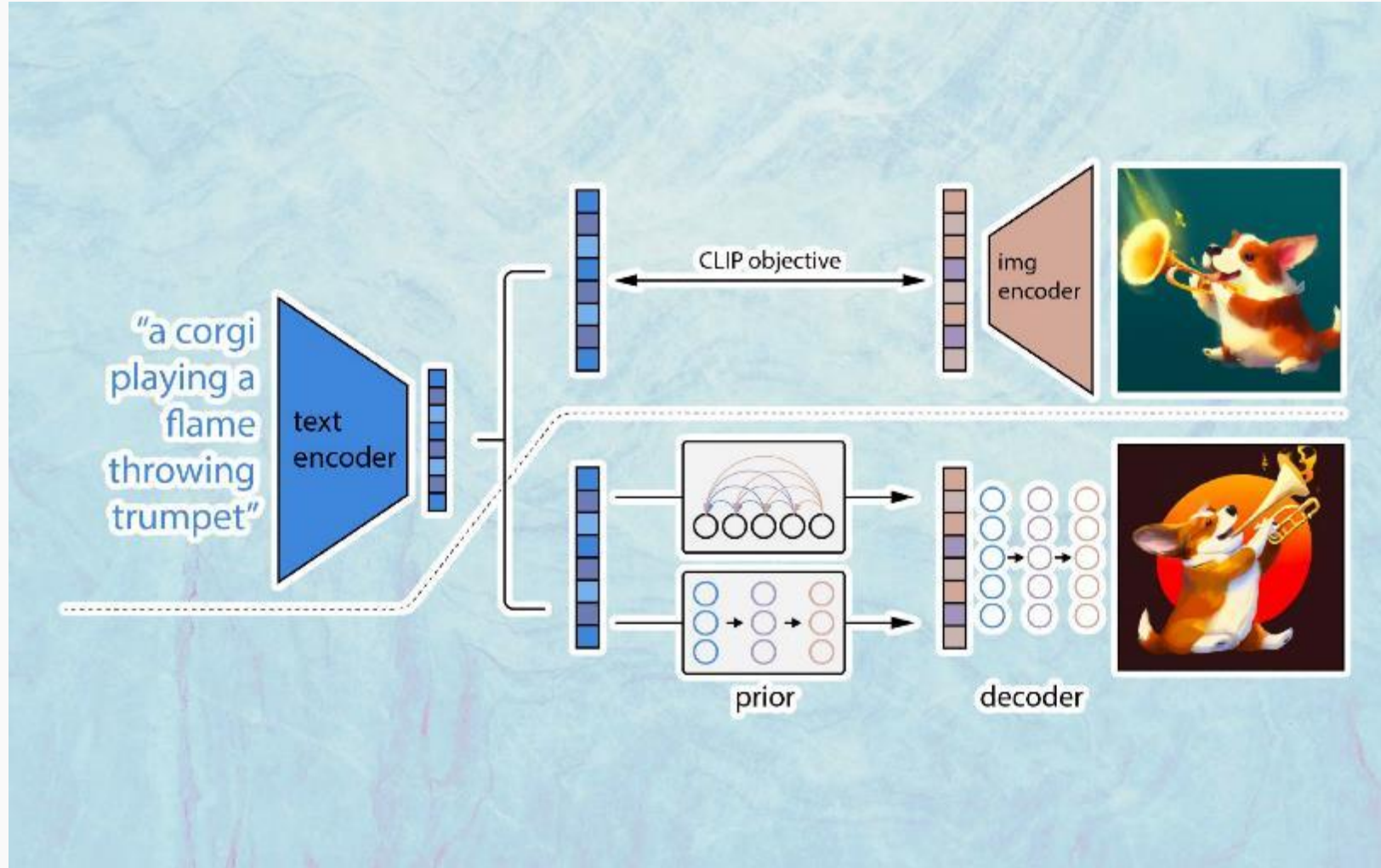


DALL-E 2 Decoder

- The decoder model allows users to create relevant images for any image embedding. As an example, the authors get two images and created image embeddings using CLIP.
- They then interpolated between these two embeddings and were able to generate images for these interpolated images as seen below.



DALL-E 2





That's all Folks!