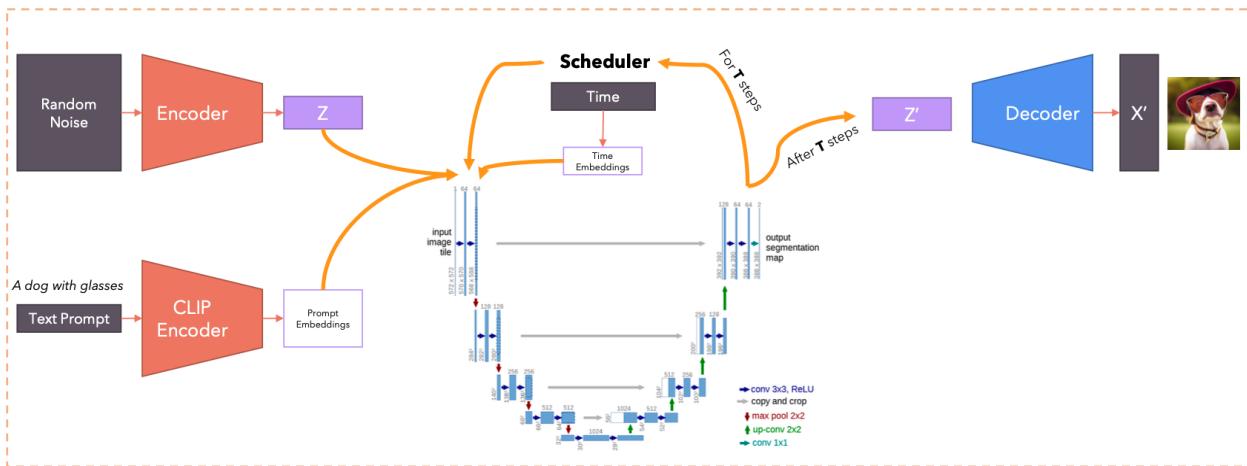


Stable Diffusion: Clearly Explained



High-Resolution Image Synthesis with Latent Diffusion Models

Robin Rombach¹ * Andreas Blattmann¹ * Dominik Lorenz¹ Patrick Esser² Björn Ommer¹
¹Ludwig Maximilian University of Munich & IWR, Heidelberg University, Germany ²Runway ML
<https://github.com/CompVis/latent-diffusion>

Abstract

By decomposing the image formation process into a sequential application of denoising autoencoders, diffusion models (DMs) achieve state-of-the-art synthesis results on image data and beyond. Additionally, their formulation allows for a guiding mechanism to control the image generation process without retraining. However, since these models typically operate directly in pixel space, optimization of powerful DMs often consumes hundreds of GPU

3 Apr 2022

Input	ours ($f = 4$) PSNR: 27.4 R-FID: 0.58	DALL-E ($f = 8$) PSNR: 22.8 R-FID: 32.01	VOGAN ($f = 16$) PSNR: 19.9 R-FID: 4.98

Vinh Dinh Nguyen
PhD in Computer Science

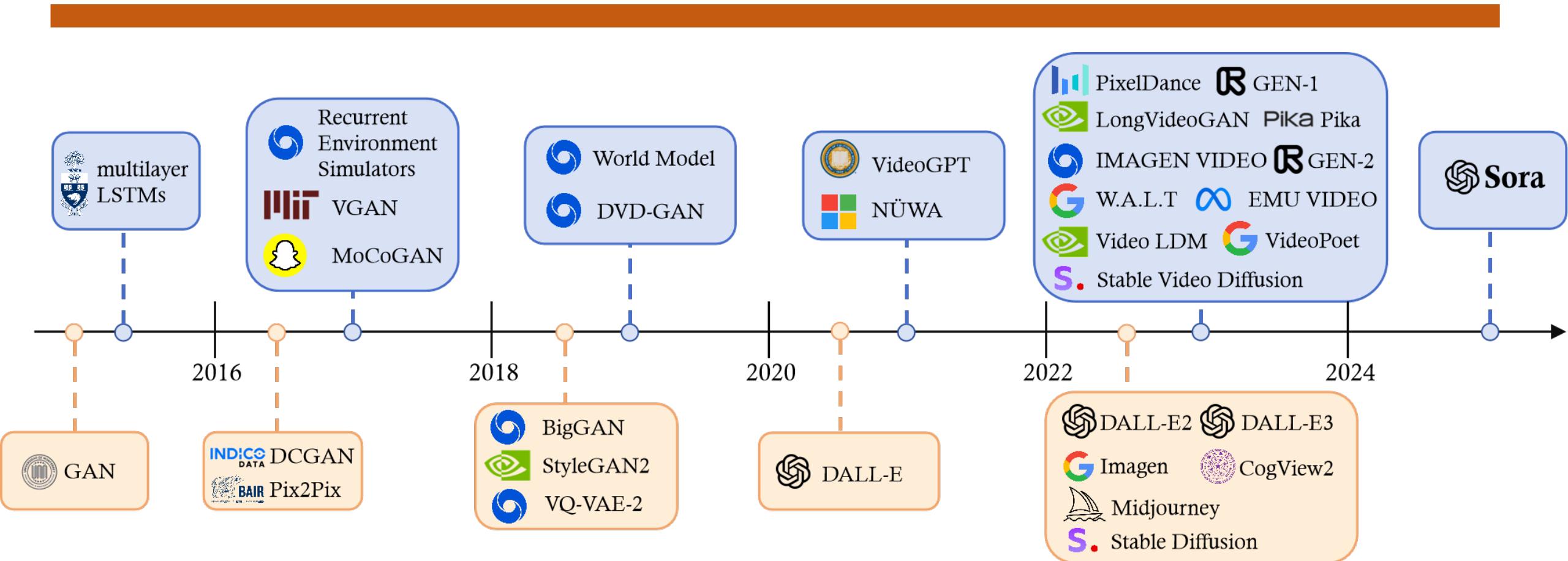
Outline

- **Objective**
- **Introduction to Stable Diffusion**
- **Stable Diffusion: Motivation**
- **Stable Diffusion: Clearly Explained**
- **Stable Diffusion: Demo with API**
- **Summary**

Outline

- **Objective**
- **Introduction to Stable Diffusion**
- **Stable Diffusion: Motivation**
- **Stable Diffusion: Clearly Explained**
- **Stable Diffusion: Demo with API**
- **Summary**

History of Generative AI in Vision Domain



Objective

High-Resolution Image Synthesis with Latent Diffusion Models

Robin Rombach¹ * Andreas Blattmann¹ * Dominik Lorenz¹ Patrick Esser¹ Björn Ommer¹

¹Ludwig Maximilian University of Munich & IWR, Heidelberg University, Germany

<https://github.com/CompVis/latent-diffusion>

Abstract

By decomposing the image formation process into a sequential application of denoising autoencoders, diffusion models (DMs) achieve state-of-the-art synthesis results on image data and beyond. Additionally, their formulation allows for a guiding mechanism to control the image generation process without retraining. However, since these models typically operate directly in pixel space, optimization of powerful DMs often consumes hundreds of GPU days and inference is expensive due to sequential evaluations. To enable DM training on limited computational resources while retaining their quality and flexibility, we apply them in the latent space of powerful pretrained autoencoders. In contrast to previous work, training diffusion models on such a representation allows for the first time to reach a near-optimal point between complexity reduction and detail preservation, greatly boosting visual fidelity. By introducing cross-attention layers into the model architecture, we turn diffusion models into powerful and flexible generators for general conditioning inputs such as text or bounding boxes and high-resolution synthesis becomes possible in a convolutional manner. Our latent diffusion models (LDMs) achieve new state-of-the-art scores for image inpainting and class-conditional image synthesis and highly competitive performance on various tasks, including text-to-image synthesis, unconditional image generation and super-resolution, while significantly reducing computational requirements compared to pixel-based DMs.

1. Introduction



Figure 1. Boosting the upper bound on achievable quality with less aggressive downsampling. Since diffusion models offer excellent inductive biases for spatial data, we do not need the heavy spatial downsampling of related generative models in latent space, but can still greatly reduce the dimensionality of the data via suitable autoencoding models, see Sec. 3. Images are from the DIV2K [1] validation set, evaluated at 512² px. We denote the spatial downsampling factor by f . Reconstruction FIDs [29] and PSNR are calculated on ImageNet-val. [12]; see also Tab. 8.

results in image synthesis [30, 85] and beyond [7, 45, 48, 57], and define the state-of-the-art in class-conditional image synthesis [15, 31] and super-resolution [72]. Moreover, even unconditional DMs can readily be applied to tasks such as inpainting and colorization [85] or stroke-based synthesis [53], in contrast to other types of generative models [19, 46, 69]. Being likelihood-based models, they do not exhibit mode-collapse and training instabilities as GANs and, by heavily exploiting parameter sharing, they can model highly complex distributions of natural images with-

Topics discussed

- 1 • Principle of Diffusion models
- 2 • Model score function of images with UNet model
- 3 • Understanding prompt through contextualized word embedding
- 4 • Let text influence image through cross attention
- 5 • Improve efficiency by adding an autoencoder
- 6 • Latent Diffusion Models in Pytorch

Prerequisites

- 1 • Basics of probability and statistic (multivariate, gaussian, conditional probability, likelihood, Bayes' rule)
- 2 • Basic of Pytorch and Neural Network
- 3 • How the attention mechanism work
- 4 • How convolution layers work

Outline

- **Objective**
- **Introduction to Stable Diffusion**
- **Stable Diffusion: Motivation**
- **Stable Diffusion: Clearly Explained**
- **Stable Diffusion: Demo with API**
- **Summary**

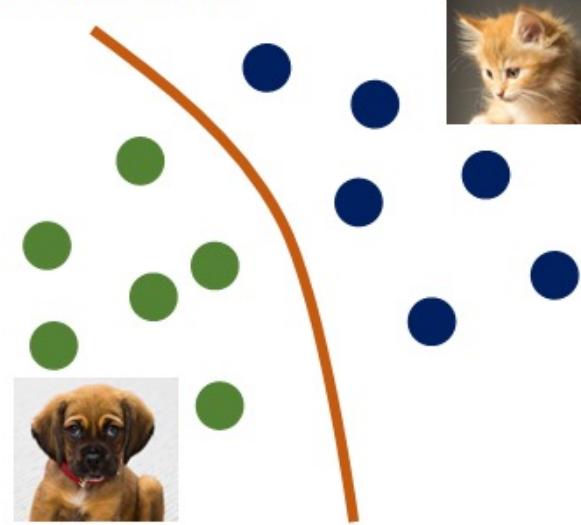
What is a generative model?

A generative model learns a probability distribution of the data set such that we can then sample from the distribution to create new instances of data.

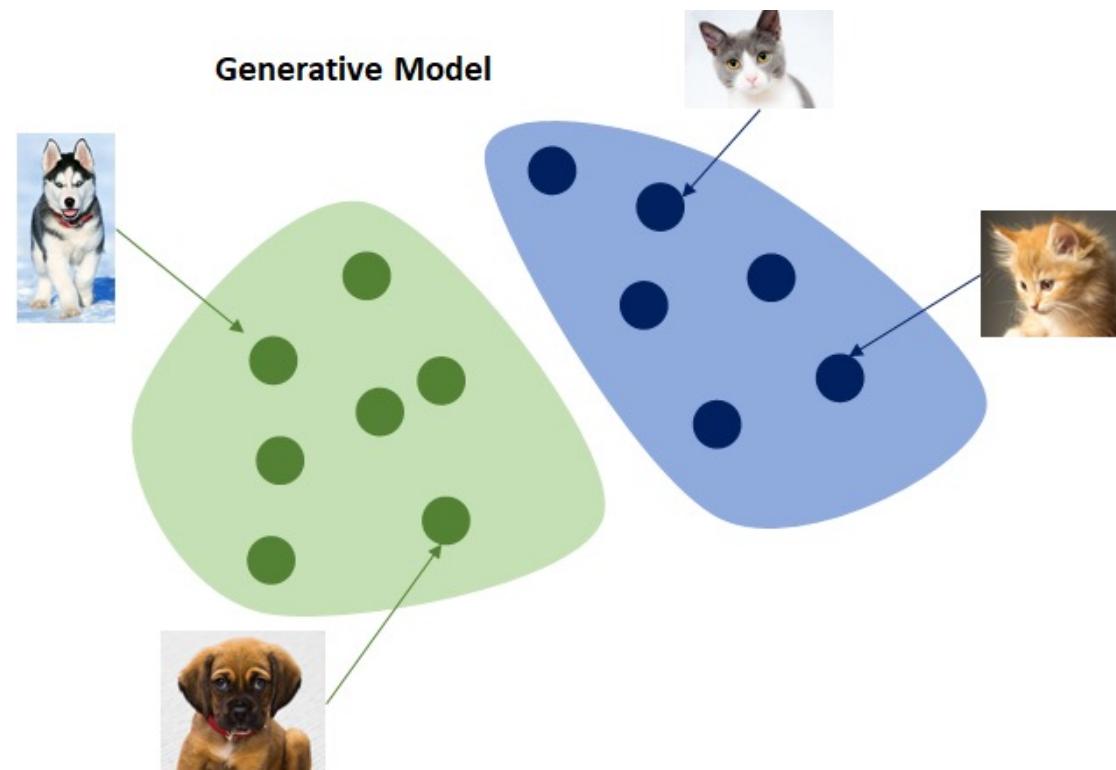
For example, if we have many pictures of cats and we train a generative model on it, we then sample from this distribution to create new images of cats



Discriminant Model



Generative Model

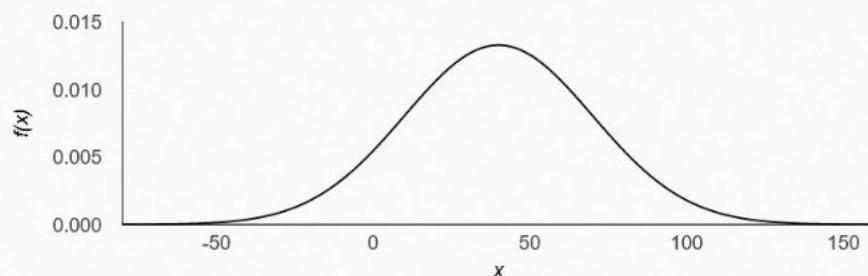


What is a generative model?

Imagine you're a researcher, and you want to generate thousands of fake identities. Each fake identity, is made up of variables, representing the characteristics of a person (Age, Height).

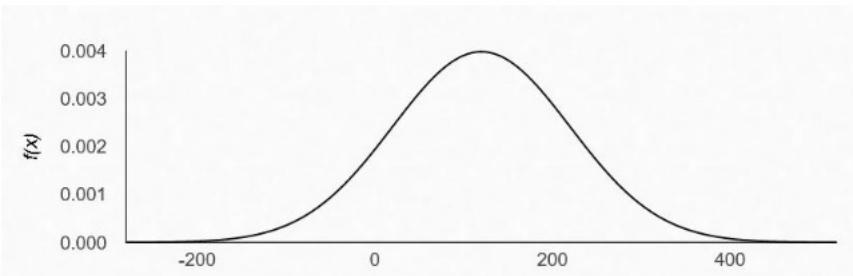


You can ask the Statistics Department of the Government to give you statistics about the age and the height of the population and then sample from these distributions.



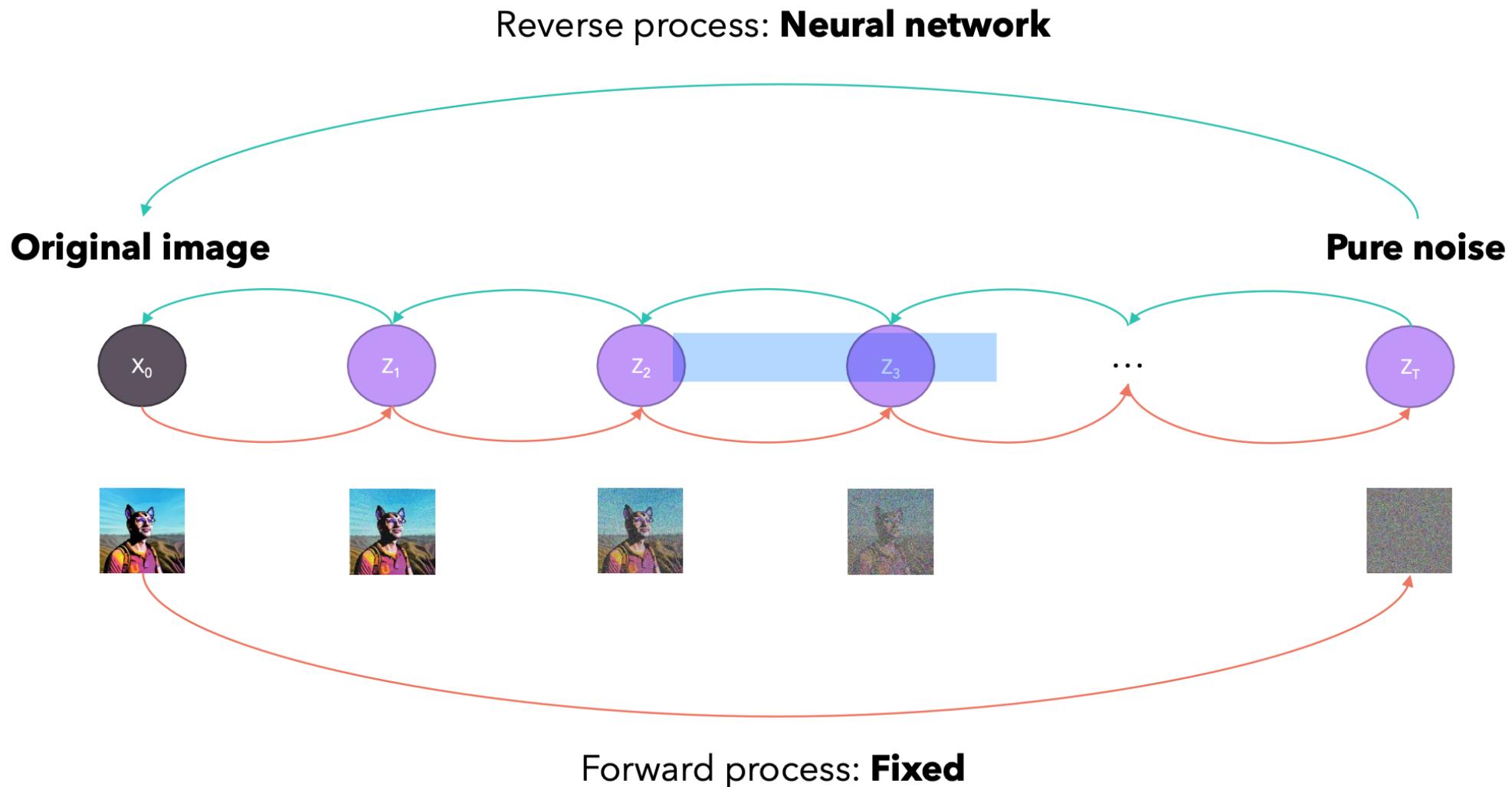
Age: $\mathbf{N}(40, 30^2)$

To generate fake identities that make sense, you need the joint distribution, otherwise you may end up with an unreasonable pair of (Age, Height)



Height: $\mathbf{N}(120, 100^2)$

What is a Diffusion Model?



What is a Diffusion Model?

Reverse process \mathbf{p}

2 Background

Diffusion models [53] are latent variable models of the form $p_\theta(\mathbf{x}_0) := \int p_\theta(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T}$, where $\mathbf{x}_1, \dots, \mathbf{x}_T$ are latents of the same dimensionality as the data $\mathbf{x}_0 \sim q(\mathbf{x}_0)$. The joint distribution $p_\theta(\mathbf{x}_{0:T})$ is called the *reverse process*, and it is defined as a Markov chain with learned Gaussian transitions starting at $p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$:

$$p_\theta(\mathbf{x}_{0:T}) := p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t), \quad p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) := \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)) \quad (1)$$

What distinguishes diffusion models from other types of latent variable models is that the approximate posterior $q(\mathbf{x}_{1:T}|\mathbf{x}_0)$, called the *forward process* or *diffusion process*, is fixed to a Markov chain that gradually adds Gaussian noise to the data according to a variance schedule β_1, \dots, β_T :

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) := \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}), \quad q(\mathbf{x}_t|\mathbf{x}_{t-1}) := \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}) \quad (2)$$

Training is performed by optimizing the usual variational bound on negative log likelihood:

$$\mathbb{E}[-\log p_\theta(\mathbf{x}_0)] \leq \mathbb{E}_q \left[-\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] = \mathbb{E}_q \left[-\log p(\mathbf{x}_T) - \sum_{t \geq 1} \log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] =: L \quad (3)$$

Forward process \mathbf{q}

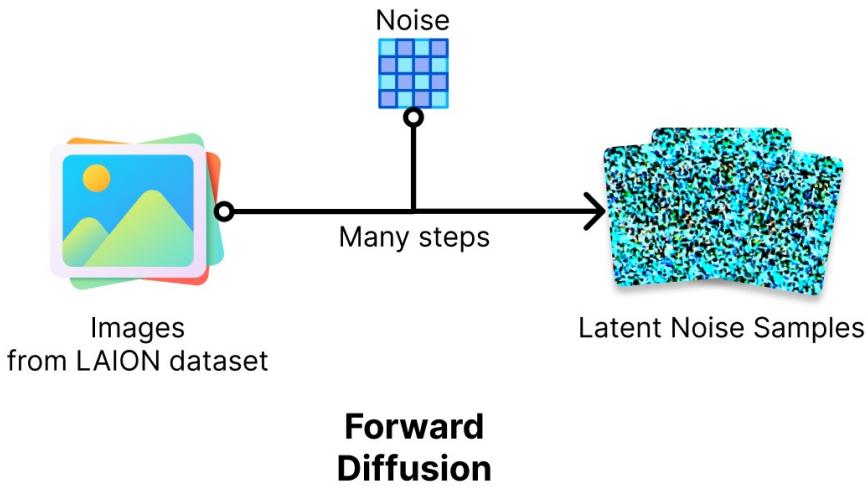
Just like with a VAE, we want to learn the parameters of the latent space

Evidence Lower Bound (ELBO)

The forward process variances β_t can be learned by reparameterization [33] or held constant as hyperparameters, and expressiveness of the reverse process is ensured in part by the choice of Gaussian conditionals in $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$, because both processes have the same functional form when β_t are small [53]. A notable property of the forward process is that it admits sampling \mathbf{x}_t at an arbitrary timestep t in closed form: using the notation $\alpha_t := 1 - \beta_t$ and $\bar{\alpha}_t := \prod_{s=1}^t \alpha_s$, we have

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}) \quad (4)$$

The diffusion model : Forward Process

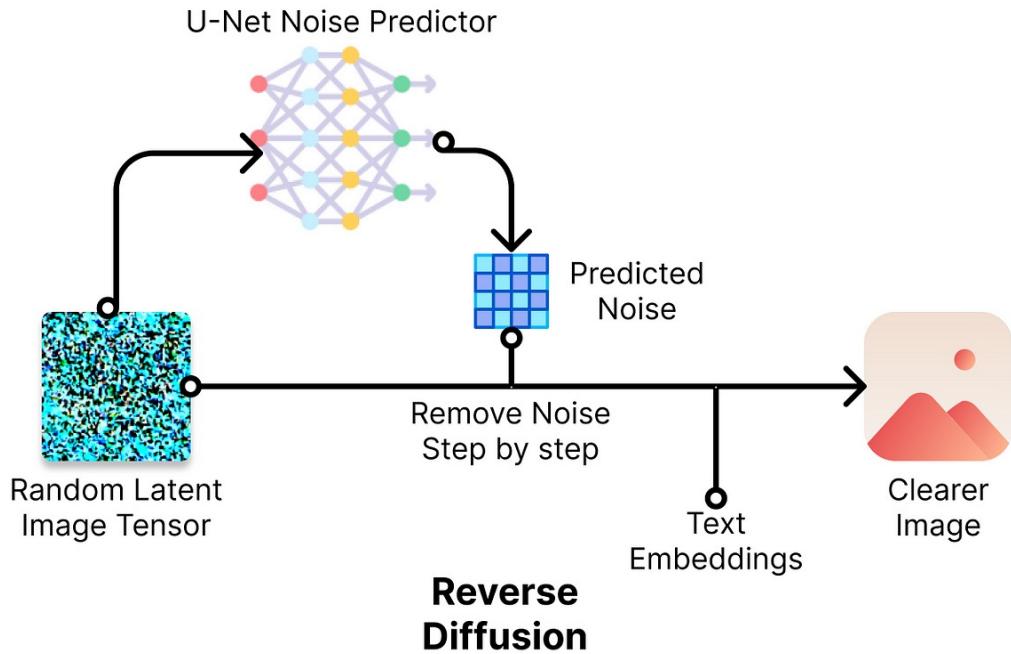


A forward diffusion process to prepare training samples and a reverse diffusion process to generate the images

As the name “diffusion” suggests, this process is just like dropping a drop of ink into the water — the ink droplet gradually diffuses in water, until you can no longer see that it was previously a drop of ink. The noise patterns added to the images are random, just like ink particles diffuse randomly into water particles, but the amount of noise can be controlled



The diffusion model : Reverse Process



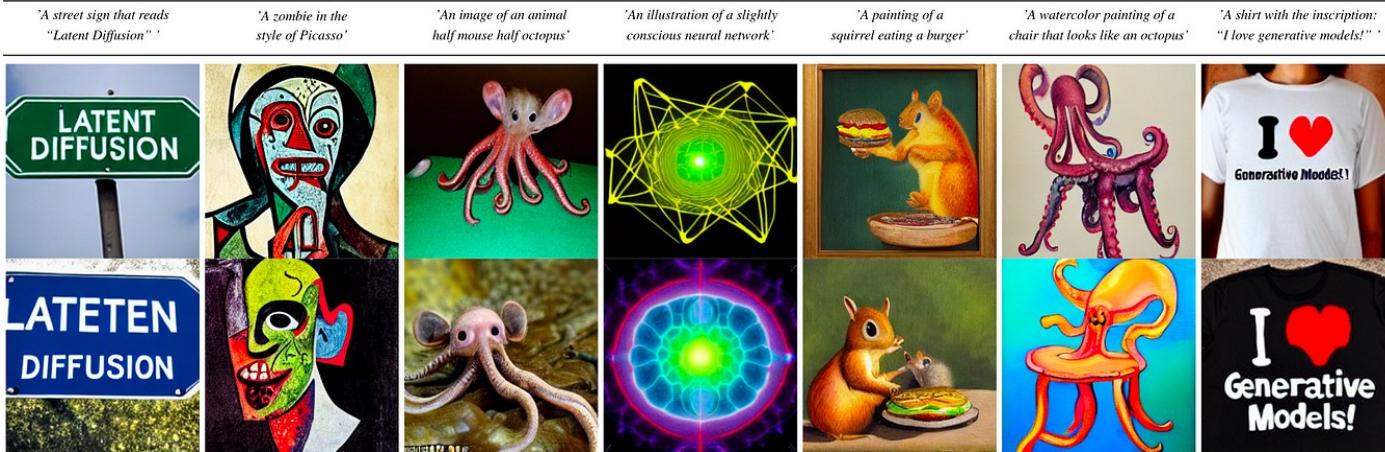
This reverse diffusion process is done gradually through multiple steps to remove the noise. The more denoising steps are done, the clearer the image becomes.



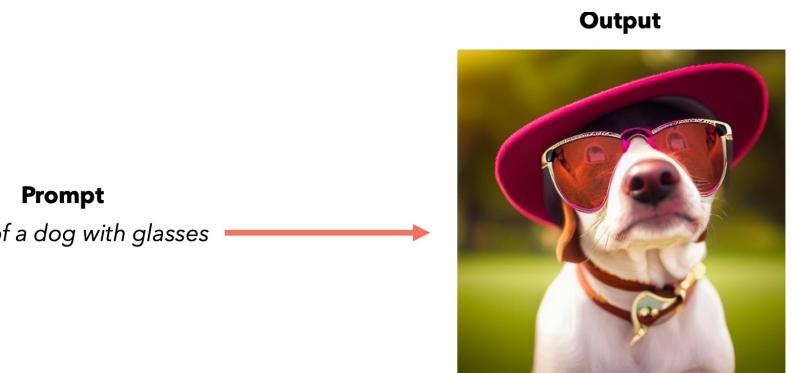
During this reverse diffusion process for getting a clearer image, researchers wanted to control how the image looks through a process called conditioning

In the reverse diffusion process, a noise predictor is trained to predict the noise added to the original image so that the model can remove the predicted noise from the noisy image to get a clearer image (in the latent space). You can think of this process as looking at a partially diffused ink droplet in water and trying to predict the location it dropped earlier.

Introduction to Stable Diffusion



Stable Diffusion is a text-to-image deep learning model, based on diffusion models.
Introduced in 2022, developed by the CompVis Group at LMU Munich. <https://github.com/Stability-AI/stablediffusion>



High-Resolution Image Synthesis with Latent Diffusion Models

Robin Rombach¹ *

Andreas Blattmann¹ *

Dominik Lorenz¹

Patrick Esser[✉]

Björn Ommer¹

¹Ludwig Maximilian University of Munich & IWR, Heidelberg University, Germany

[✉]Runway ML

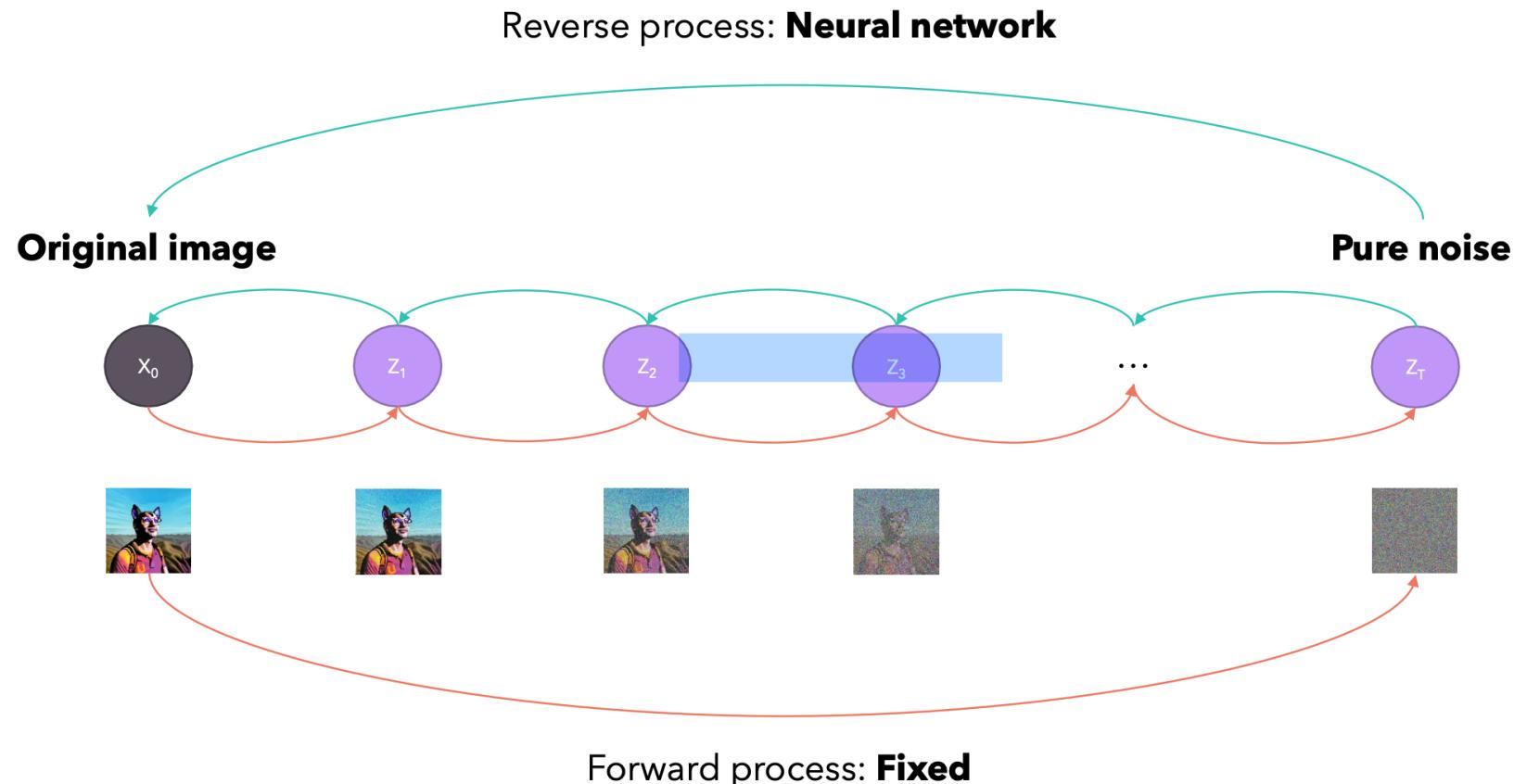
<https://github.com/CompVis/latent-diffusion>

Outline

- **Objective**
- **Introduction to Stable Diffusion**
- **Stable Diffusion: Motivation**
- **Stable Diffusion: Clearly Explained**
- **Stable Diffusion: Demo with API**
- **Summary**

Diffusion Model: Limitations

Since the latent variables have the same dimension (size of the vector) as the original data, if we want to perform many steps to denoise an image, that would result in a lot of steps through the Unet, which can be very slow if the matrix representing our data/latent is large. **What if we could “compress” our data before running it through the forward/reverse process (UNet)?**



Stable Diffusion: Motivation

This is emphasized when dealing with large image sizes and a high number of diffusion steps (T). The time required for denoising the image from Gaussian noise during sampling can become prohibitively long. To address this issue, a group of researchers proposed a novel approach called **Stable Diffusion**, originally known as Latent Diffusion Model (LDM).

High-Resolution Image Synthesis with Latent Diffusion Models

Robin Rombach¹ * Andreas Blattmann¹ * Dominik Lorenz¹ Patrick Esser¹⁸ Björn Ommer¹
¹Ludwig Maximilian University of Munich & IWR, Heidelberg University, Germany ¹⁸Runway ML
<https://github.com/CompVis/latent-diffusion>

Abstract

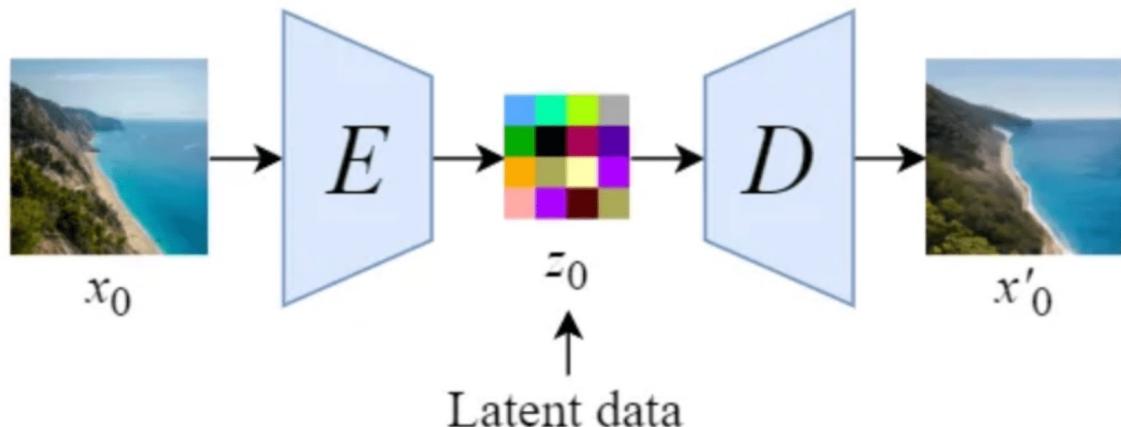
By decomposing the image formation process into a sequential application of denoising autoencoders, diffusion models (DMs) achieve state-of-the-art synthesis results on image data and beyond. Additionally, their formulation allows for a guiding mechanism to control the image generation process without retraining. However, since these models typically operate directly in pixel space, optimization of powerful DMs often consumes hundreds of GPU days and inference is expensive due to sequential evaluations. To enable DM training on limited computational resources while retaining their quality and flexibility, we apply them in the latent space of powerful pretrained autoencoders. In contrast to previous work, training diffusion models on such a representation allows for the first time to reach a near-optimal point between complexity reduction and detail preservation, greatly boosting visual fidelity. By introducing cross-attention layers into the model architecture, we turn diffusion models into powerful and flexible generators for general conditioning inputs such as text or bounding boxes and high-resolution synthesis becomes possible in a convolutional manner. Our latent diffusion models (LDMs) achieve new state-of-the-art scores for image inpainting and class-conditional image synthesis and highly competitive performance on various tasks, including text-to-image synthesis, unconditional image generation and super-resolution, while significantly reducing computational requirements compared to pixel-based DMs.

1. Introduction



Figure 1. Boosting the upper bound on achievable quality with less aggressive downampling. Since diffusion models offer excellent inductive biases for spatial data, we do not need the heavy spatial downampling of related generative models in latent space, but can still greatly reduce the dimensionality of the data via suitable autoencoding models, see Sec. 3. Images are from the DIV2K [1] validation set, evaluated at 512² px. We denote the spatial down-sampling factor by f . Reconstruction FIDs [29] and PSNR are calculated on ImageNet-val. [12]; see also Tab. 8.

results in image synthesis [30, 85] and beyond [7, 45, 48, 57], and define the state-of-the-art in class-conditional image synthesis [15, 31] and super-resolution [72]. Moreover, even unconditional DMs can readily be applied to tasks such as inpainting and colorization [85] or stroke-based synthesis [53], in contrast to other types of generative models [19, 46, 69]. Being likelihood-based models, they do not exhibit mode-collapse and training instabilities as GANs and, by heavily exploiting parameter sharing, they can model highly complex distributions of natural images with-



Stable Diffusion: Motivation

This is emphasized when dealing with large image sizes and a high number of diffusion steps (T). The time required for denoising the image from Gaussian noise during sampling can become prohibitively long. To address this issue, a group of researchers proposed a novel approach called **Stable Diffusion**, originally known as Latent Diffusion Model (LDM).

High-Resolution Image Synthesis with Latent Diffusion Models

Robin Rombach¹ * Andreas Blattmann¹ * Dominik Lorenz¹ Patrick Esser¹⁸ Björn Ommer¹
¹Ludwig Maximilian University of Munich & IWR, Heidelberg University, Germany ¹⁸Runway ML
<https://github.com/CompVis/latent-diffusion>

Abstract

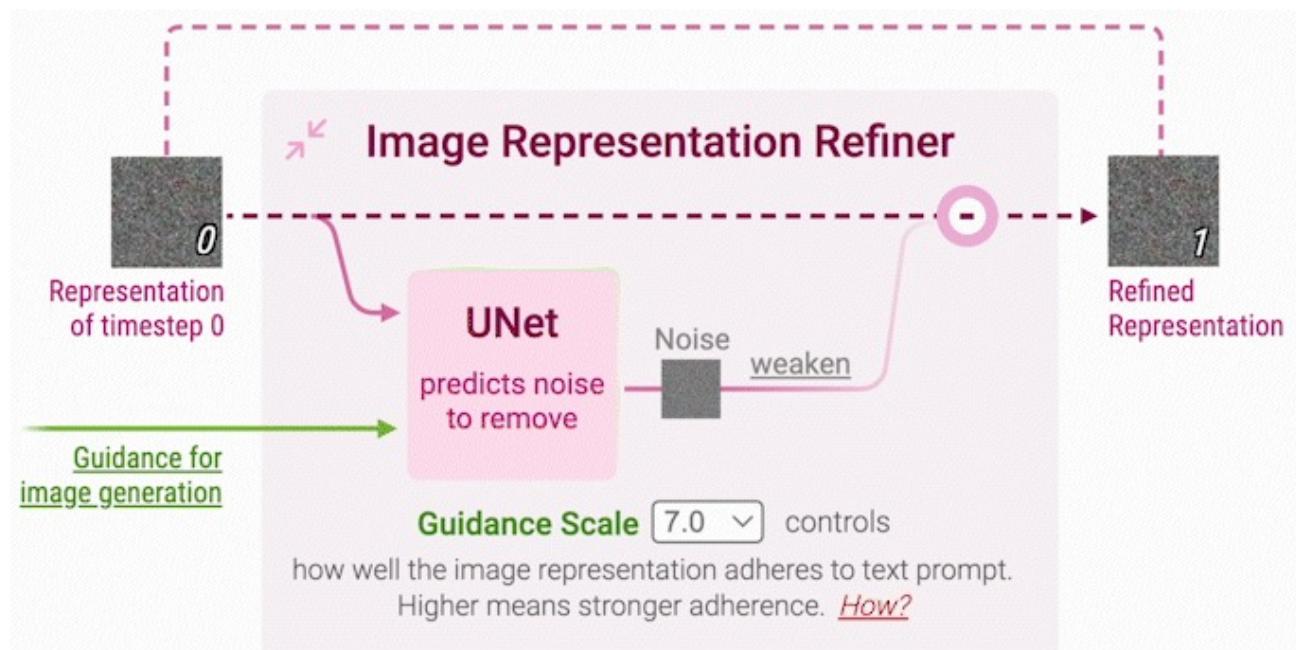
By decomposing the image formation process into a sequential application of denoising autoencoders, diffusion models (DMs) achieve state-of-the-art synthesis results on image data and beyond. Additionally, their formulation allows for a guiding mechanism to control the image generation process without retraining. However, since these models typically operate directly in pixel space, optimization of powerful DMs often consumes hundreds of GPU days and inference is expensive due to sequential evaluations. To enable DM training on limited computational resources while retaining their quality and flexibility, we apply them in the latent space of powerful pretrained autoencoders. In contrast to previous work, training diffusion models on such a representation allows for the first time to reach a near-optimal point between complexity reduction and detail preservation, greatly boosting visual fidelity. By introducing cross-attention layers into the model architecture, we turn diffusion models into powerful and flexible generators for general conditioning inputs such as text or bounding boxes and high-resolution synthesis becomes possible in a convolutional manner. Our latent diffusion models (LDMs) achieve new state-of-the-art scores for image inpainting and class-conditional image synthesis and highly competitive performance on various tasks, including text-to-image synthesis, unconditional image generation and super-resolution, while significantly reducing computational requirements compared to pixel-based DMs.

1. Introduction

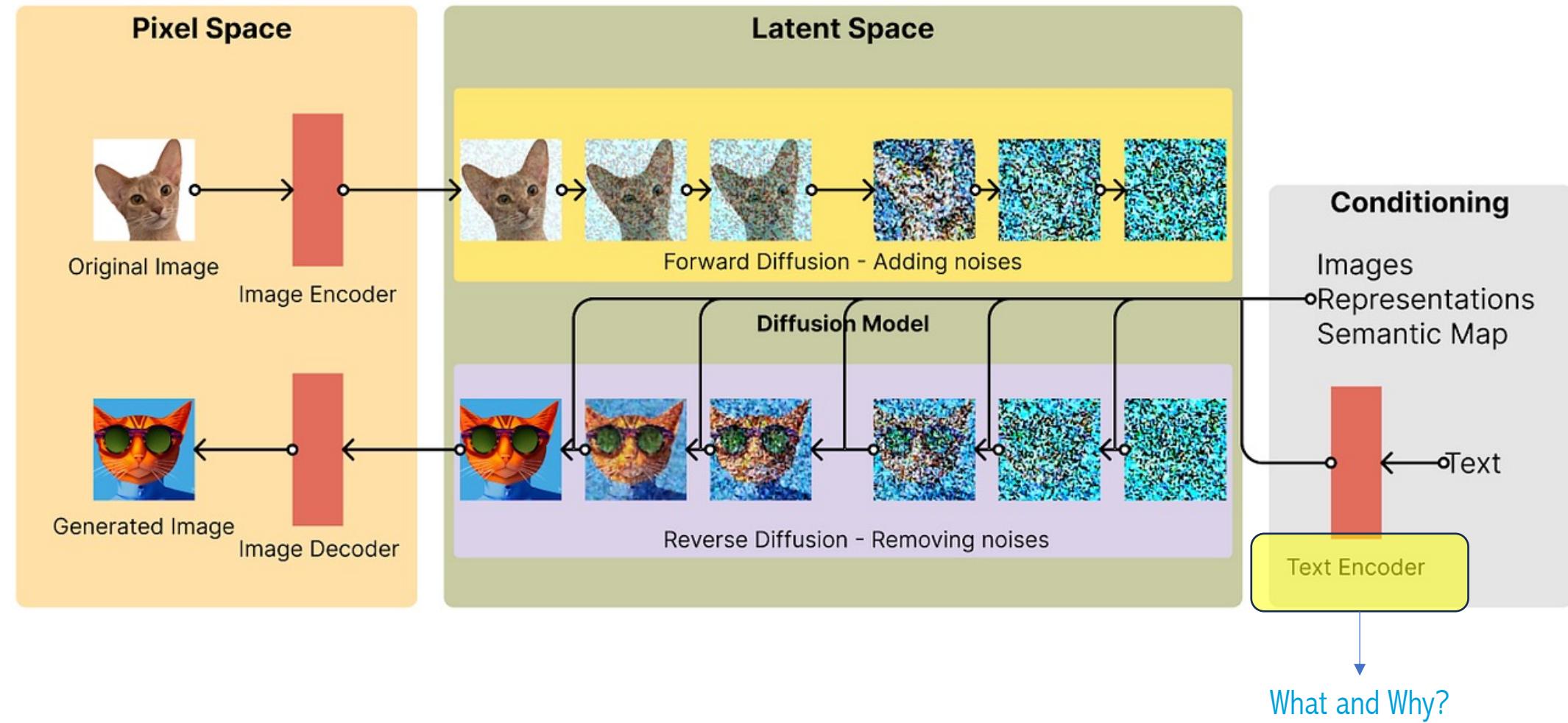


Figure 1. Boosting the upper bound on achievable quality with less aggressive downsampling. Since diffusion models offer excellent inductive biases for spatial data, we do not need the heavy spatial downsampling of related generative models in latent space, but can still greatly reduce the dimensionality of the data via suitable autoencoding models, see Sec. 3. Images are from the DIV2K [1] validation set, evaluated at 512^2 px. We denote the spatial down-sampling factor by f . Reconstruction FIDs [29] and PSNR are calculated on ImageNet-val. [12]; see also Tab. 8.

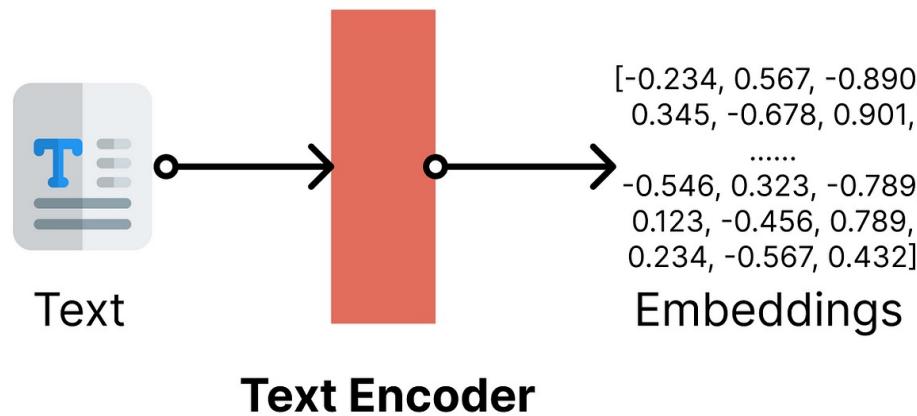
results in image synthesis [30, 85] and beyond [7, 45, 48, 57], and define the state-of-the-art in class-conditional image synthesis [15, 31] and super-resolution [72]. Moreover, even unconditional DMs can readily be applied to tasks such as inpainting and colorization [85] or stroke-based synthesis [53], in contrast to other types of generative models [19, 46, 69]. Being likelihood-based models, they do not exhibit mode-collapse and training instabilities as GANs and, by heavily exploiting parameter sharing, they can model highly complex distributions of natural images with-



Stable Diffusion: Motivation



Text encoder: text input to embeddings



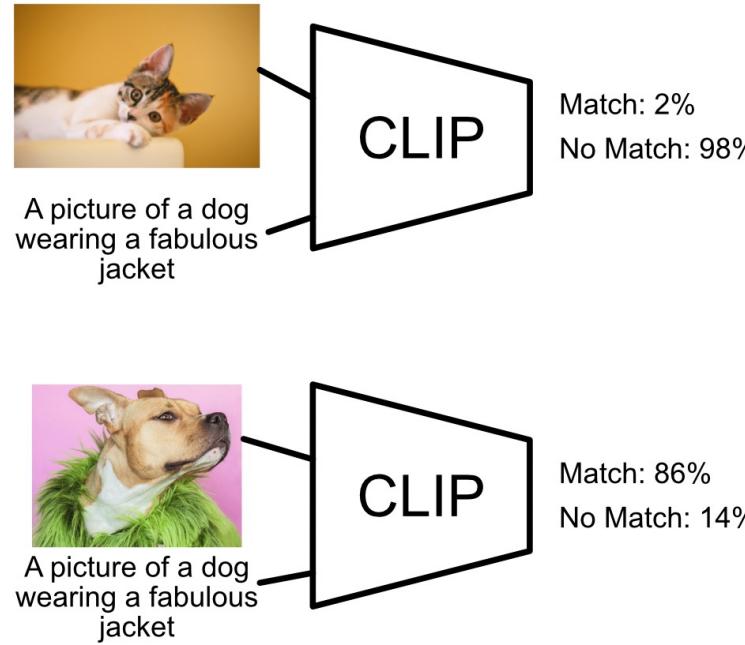
These numbers are not just random numbers — they're called text embeddings, which are high-dimensional vectors that can capture the semantic meaning of the texts (i.e. the relationship between words and their context).

Imagine you want to have a foreign artist draw a painting for you but you don't speak their language, you'll probably use Google Translate or a human translator to translate what you want to them. It's the same with image-generation models — machine learning models don't understand text directly, so they need a text encoder to translate your text instructions into numbers that they can understand

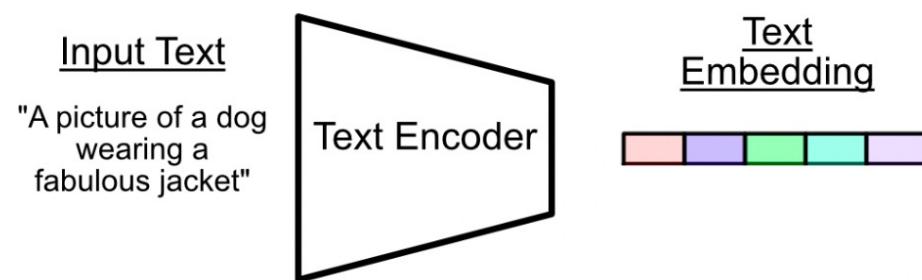
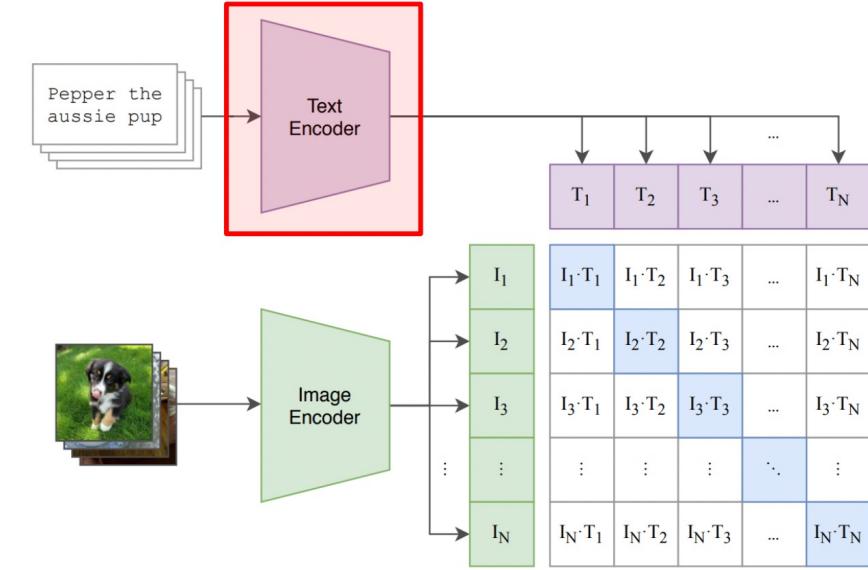
Stable diffusion v1 uses CLIP, which is a GPT-based model by OpenAI. (The original stable diffusion research paper uses BERT, another transformer model; Stable diffusion V2 uses OpenClip, a larger version of CLIP)

→ What and Why?

CLIP (Contrastive Language–Image Pre-training)



ClipText for text encoding.
Input: text.
Output: 77 token embeddings vectors, each in 768 dimensions.



Stable Diffusion: Motivation

Stable Diffusion is a **latent diffusion model**, in which we **don't learn the distribution $p(x)$ of our data set of images**, but rather, **the distribution of a latent representation** of our data by using a **Variational Autoencoder**.

► What and Why?

This allows us to reduce the computation we need to perform the steps needed to generate a sample, because each data will not be represented by a **512x512 image**, but its latent representation, **which is 64x64**

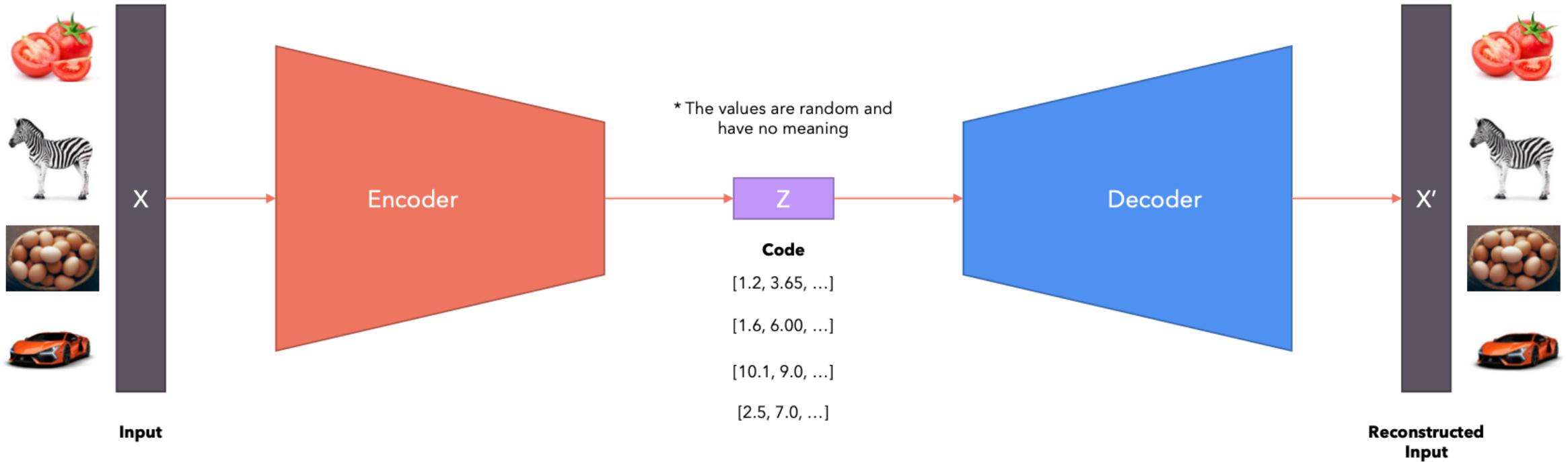
High-Resolution Image Synthesis with Latent Diffusion Models

Robin Rombach¹ * Andreas Blattmann¹ * Dominik Lorenz¹ Patrick Esser[✉] Björn Ommer¹

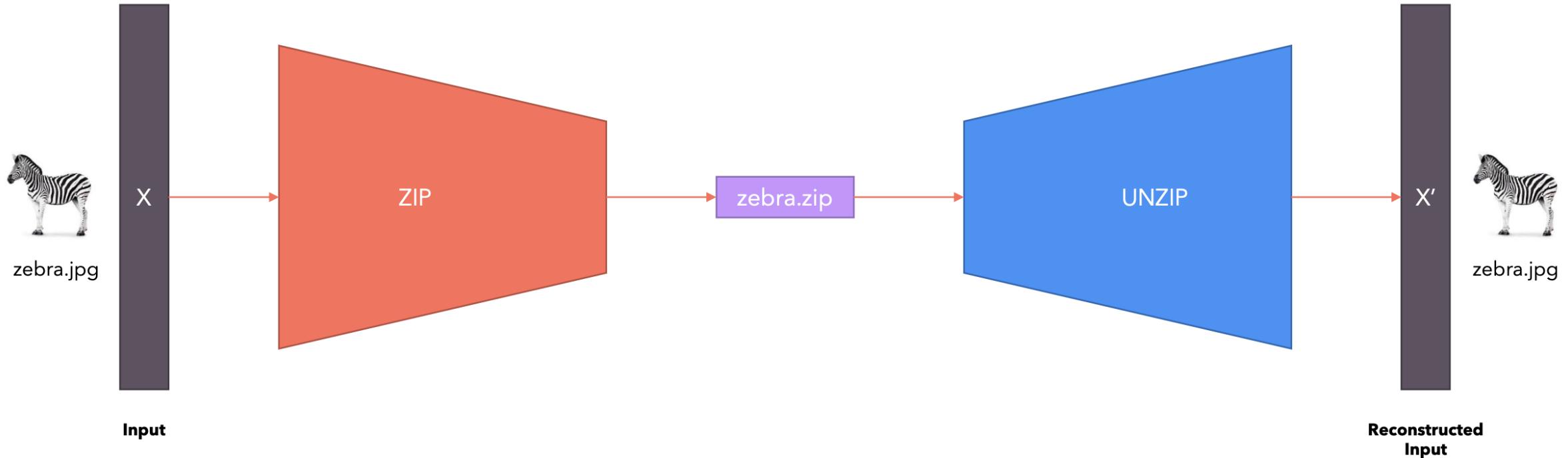
¹Ludwig Maximilian University of Munich & IWR, Heidelberg University, Germany [✉]Runway ML

<https://github.com/CompVis/latent-diffusion>

What is Auto-encoder



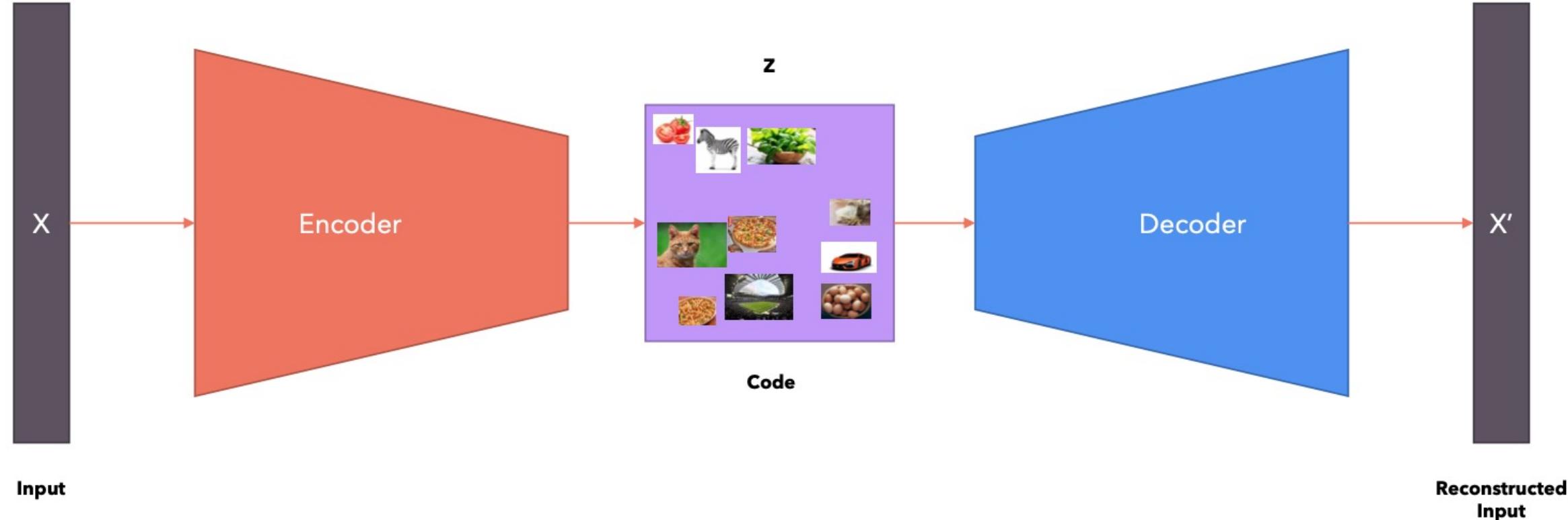
What is Auto-encoder



Analogy with file compression

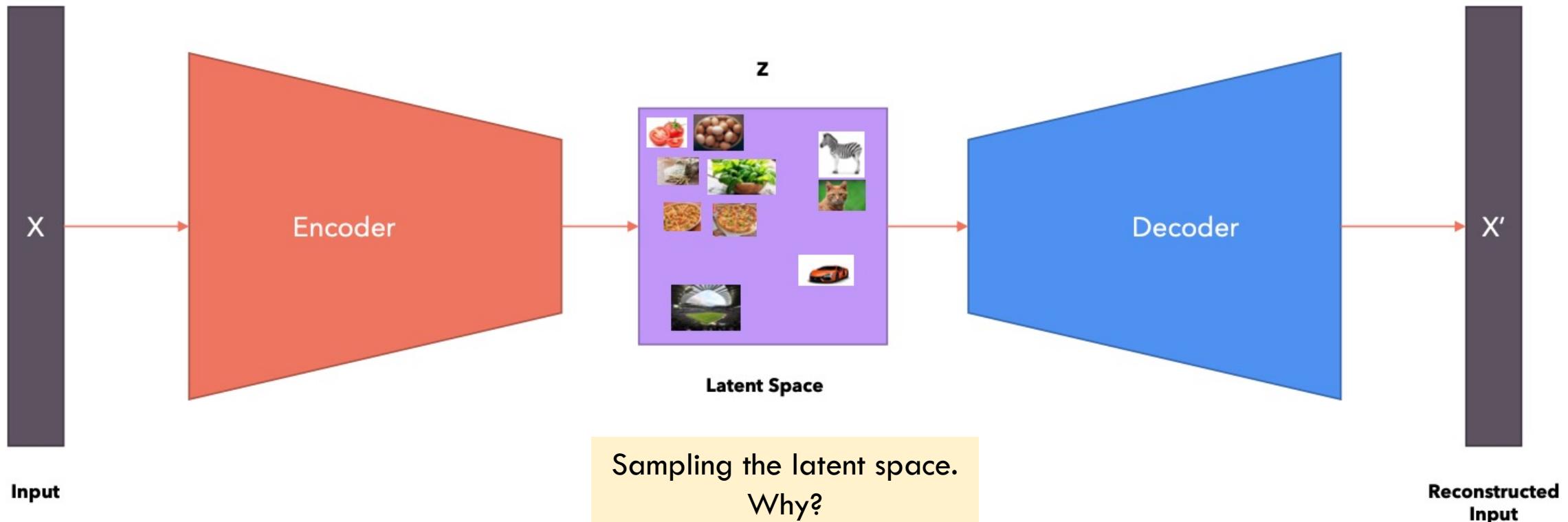


What's the problem with Autoencoders?



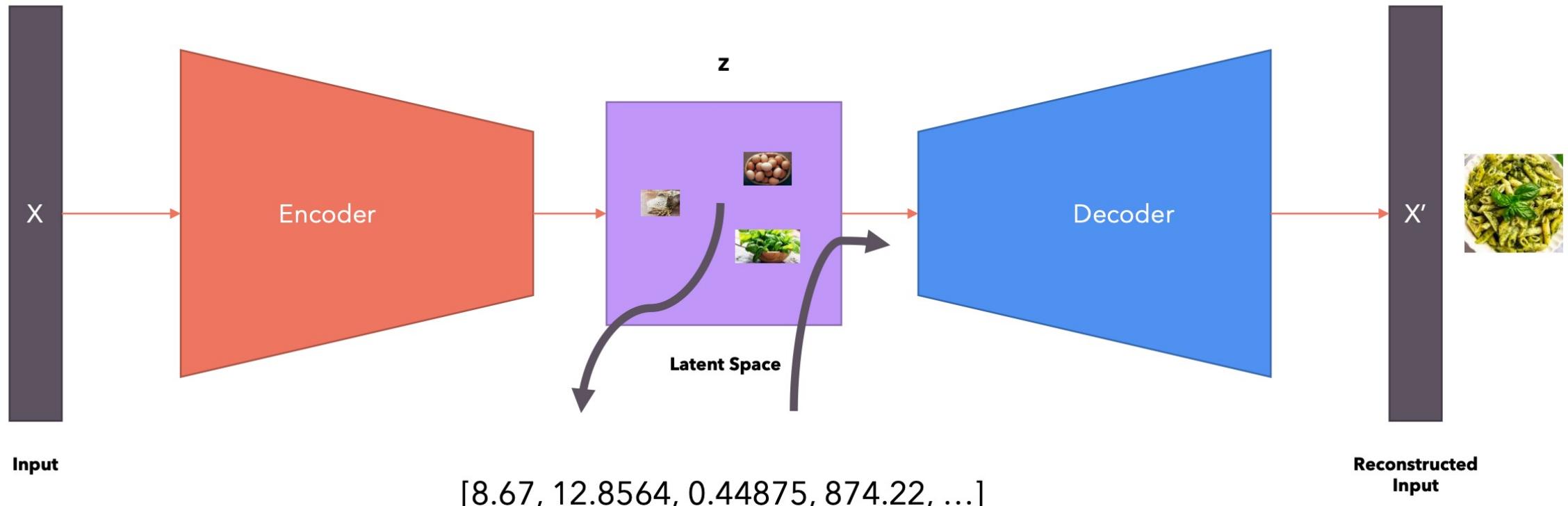
The code learned by the model makes no sense. That is, the model can just assign any vector to the inputs without the numbers in the vector representing any pattern. The model **doesn't capture any semantic relationship between the data**.

Introducing the Variational Autoencoder



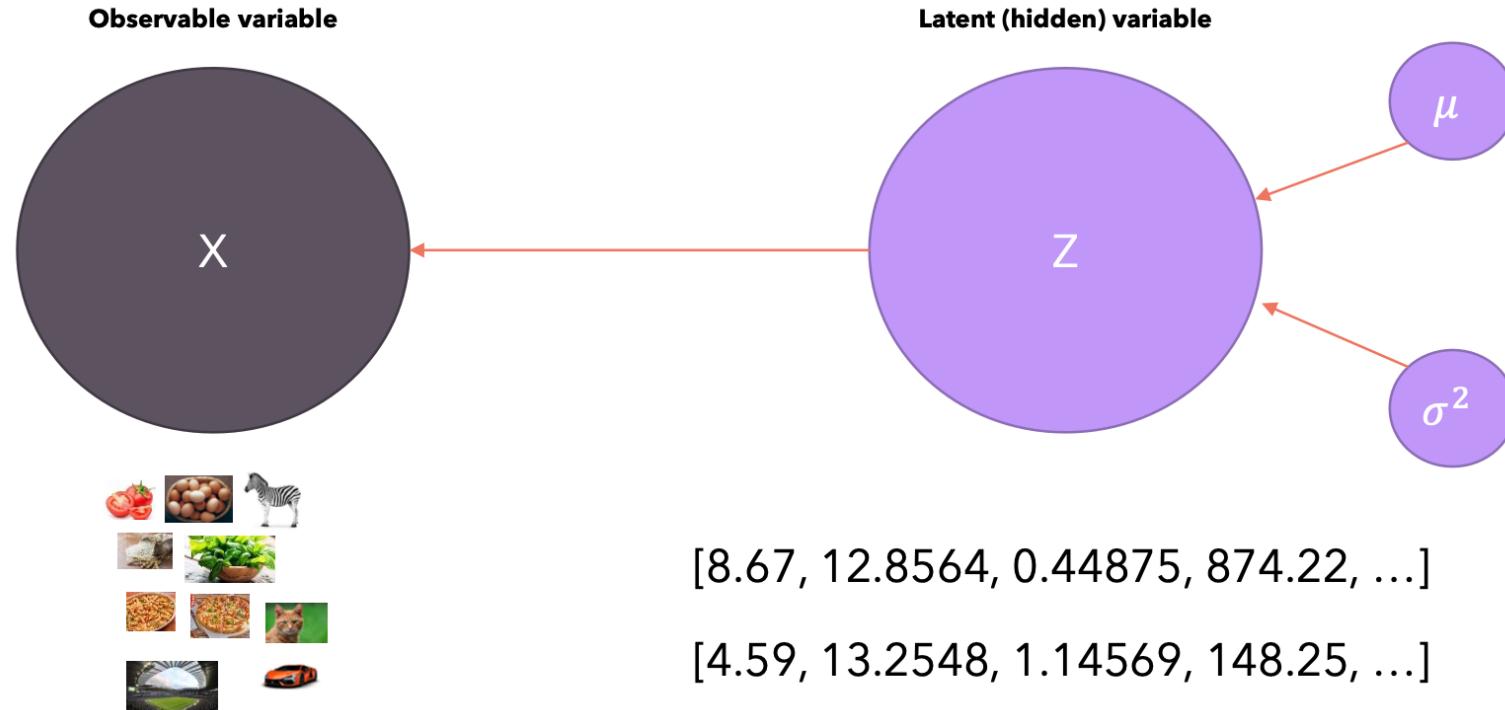
The variational autoencoder, instead of learning a code, learns a “latent space”. The latent space represents the parameters of a (multivariate) distribution.

Sampling the latent space



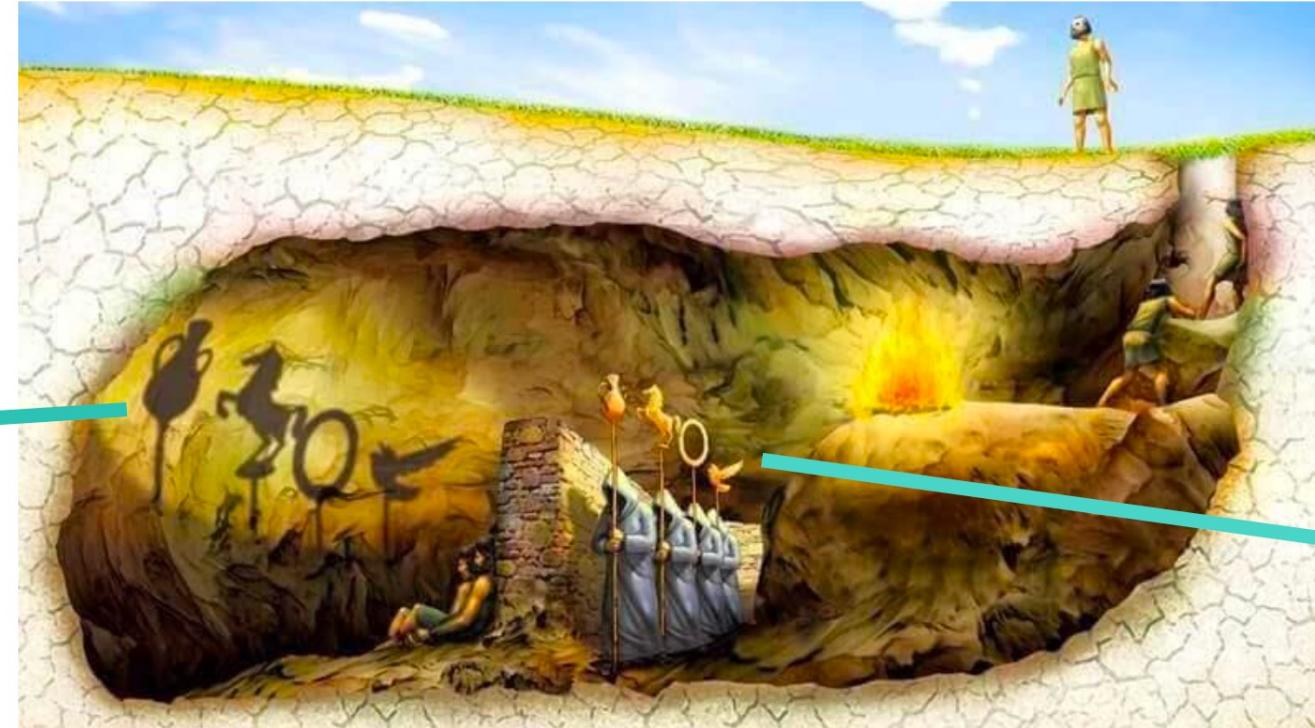
Just like when you use Python to generate a random number between 1 and 100, you're sampling from a uniform (pseudo)random distribution between 1 and 100. In the same way, we can sample from the latent space in order to generate a random vector, give it to the decoder and generate new data.

Why is it called latent space?



Why is it called latent space?

Why is it called latent space?



Observable variable



Latent (hidden) variable

[8.67, 12.8564, 0.44875, 874.22, ...]

[4.59, 13.2548, 1.14569, 148.25, ...]

[1.74, 32.3476, 5.18469, 358.14, ...]



Plato's allegory of the cave

Basic Math Concepts

Expectation of a random variable

$$E_x[f(x)] = \int xf(x)dx$$

Chain rule of probability

$$P(x, y) = P(x|y)P(y)$$

Bayes' Theorem

$$P(x | y) = \frac{P(y|x)P(x)}{P(y)}$$



Math is fun

Kullback-Leibler Divergence

$$D_{KL}(P||Q) = \int p(x) \log\left(\frac{p(x)}{q(x)}\right) dx$$

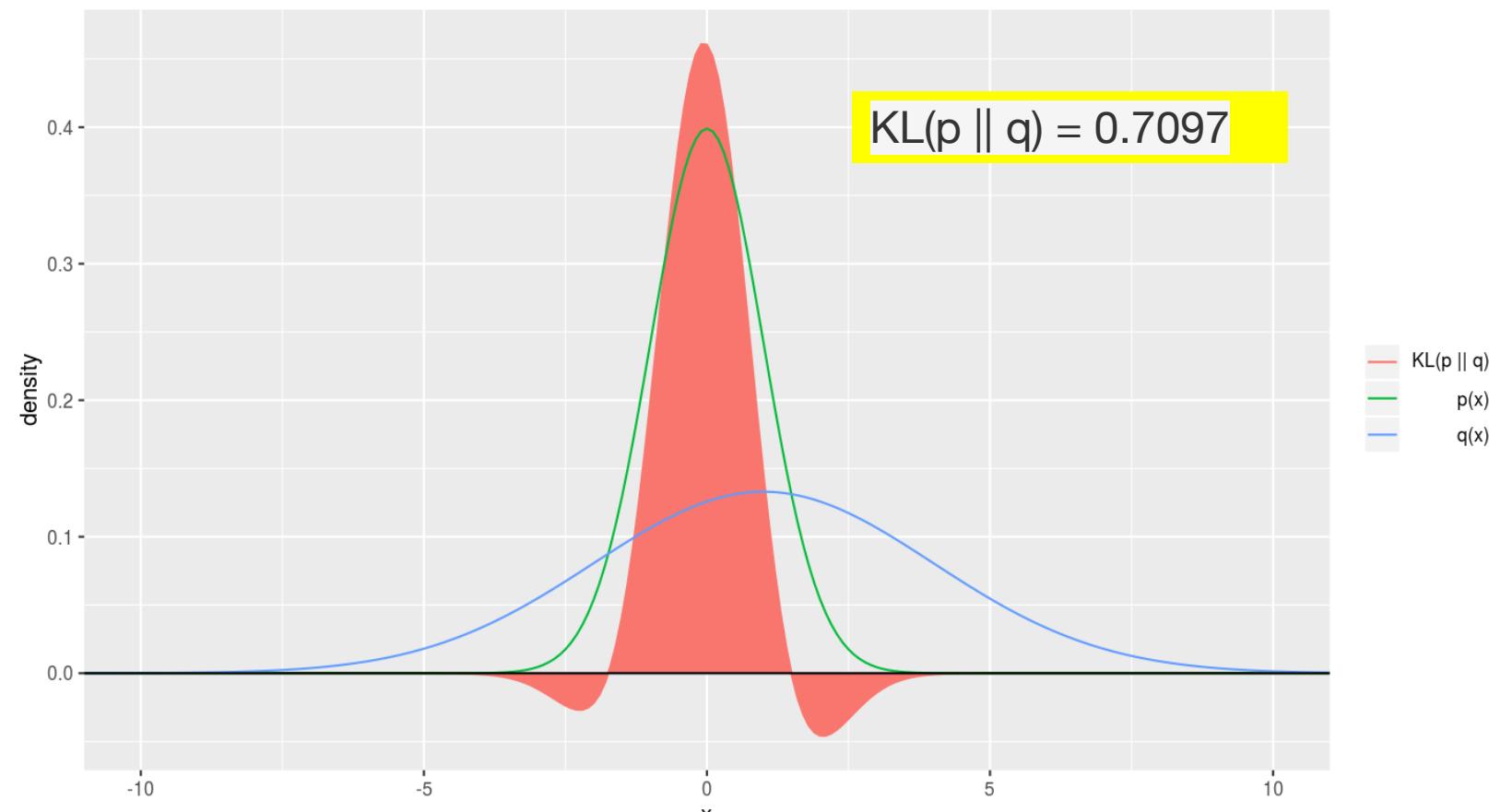
p(x): mean = 0 & std = 1

q(x): mean = 1 & std = 3

Properties:

- Not symmetric.
- Always ≥ 0
- It is equal to 0 if and only if $P = Q$

KL allows you to measure the distance between two probability distributions



Kullback-Leibler Divergence

$$D_{KL}(P||Q) = \int p(x) \log\left(\frac{p(x)}{q(x)}\right) dx$$

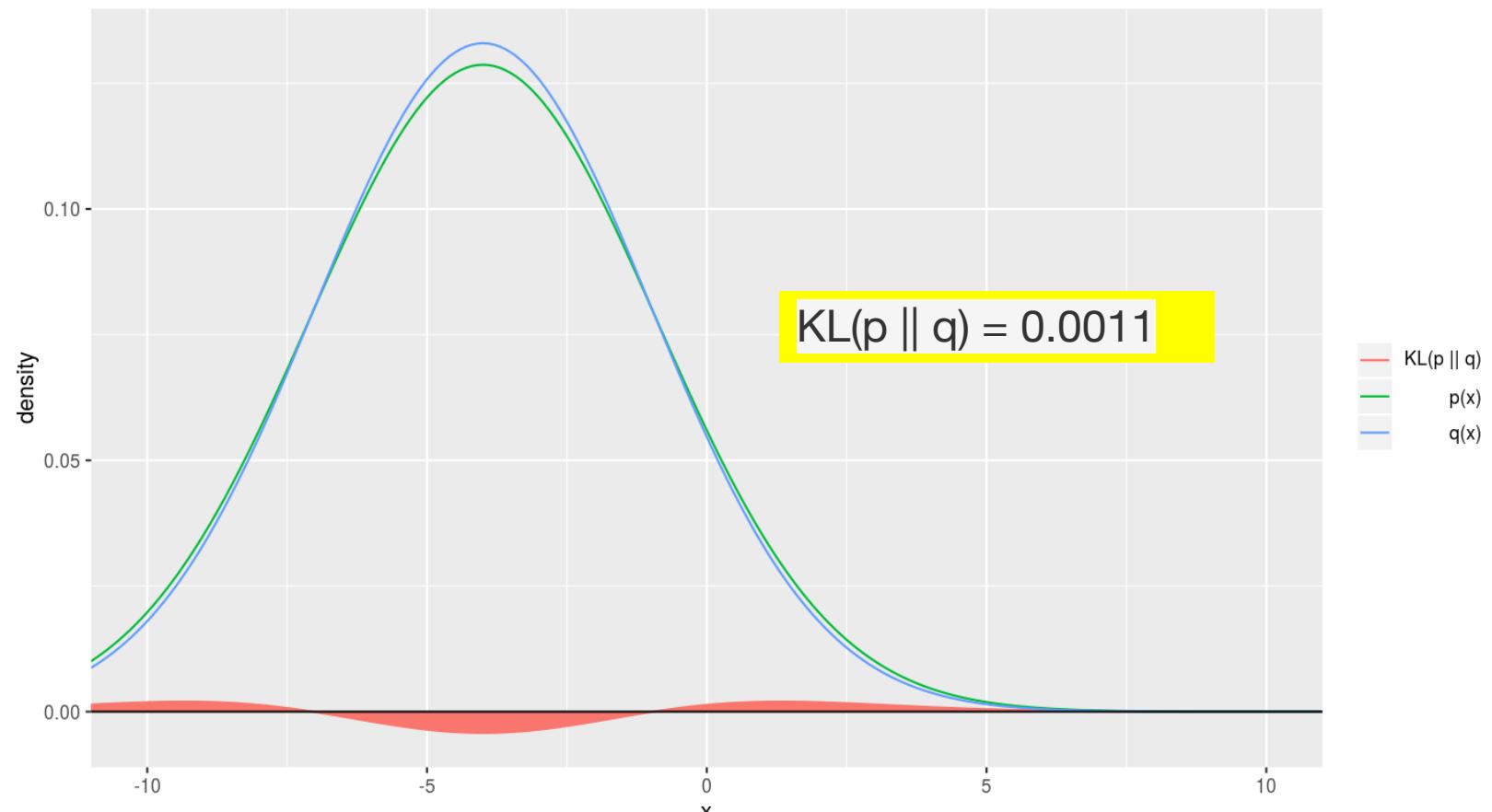
p(x): mean = -4 & std = 3.1

q(x): mean = -4 & std = 3

Properties:

- Not symmetric.
- Always ≥ 0
- It is equal to 0 if and only if $P = Q$

KL allows you to measure the distance between two probability distributions



Overview of the Model

We can define the likelihood of our data as the marginalization over the joint probability with respect to the latent variable

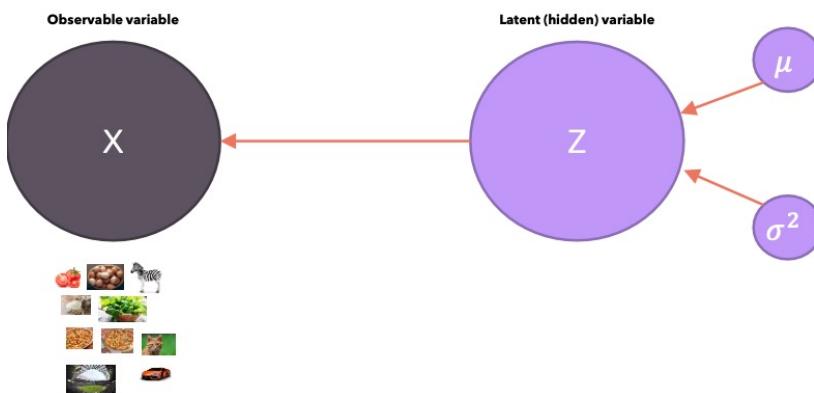
$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z}$$

Is intractable because we would need to evaluate this integral over all latent variables Z.

we can use the Chain rule of probability

$$p(\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{z} | \mathbf{x})}$$

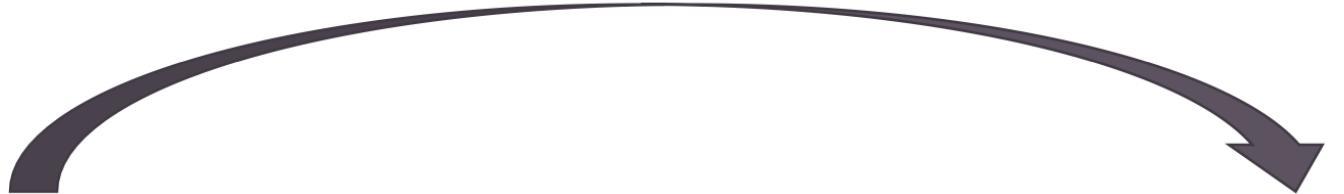
We don't have a ground truth $p(\mathbf{z} | \mathbf{x})$... which is also what we're trying to find!



Intractable problem = a problem that can be solved in theory (e.g. given large but finite resources, especially time), but for which in practice any solution takes too many resources to be useful, is known as an intractable problem.



A chicken and egg problem



In order to have a tractable $p(\mathbf{x})$ we need a tractable $p(\mathbf{z}|\mathbf{x})$

$$p(\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{z}|\mathbf{x})}$$



$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})}$$



In order to have a tractable $p(\mathbf{z}|\mathbf{x})$ we need a tractable $p(\mathbf{x})$

Question: When you can not find what you want?
Answer: We try to approximate it

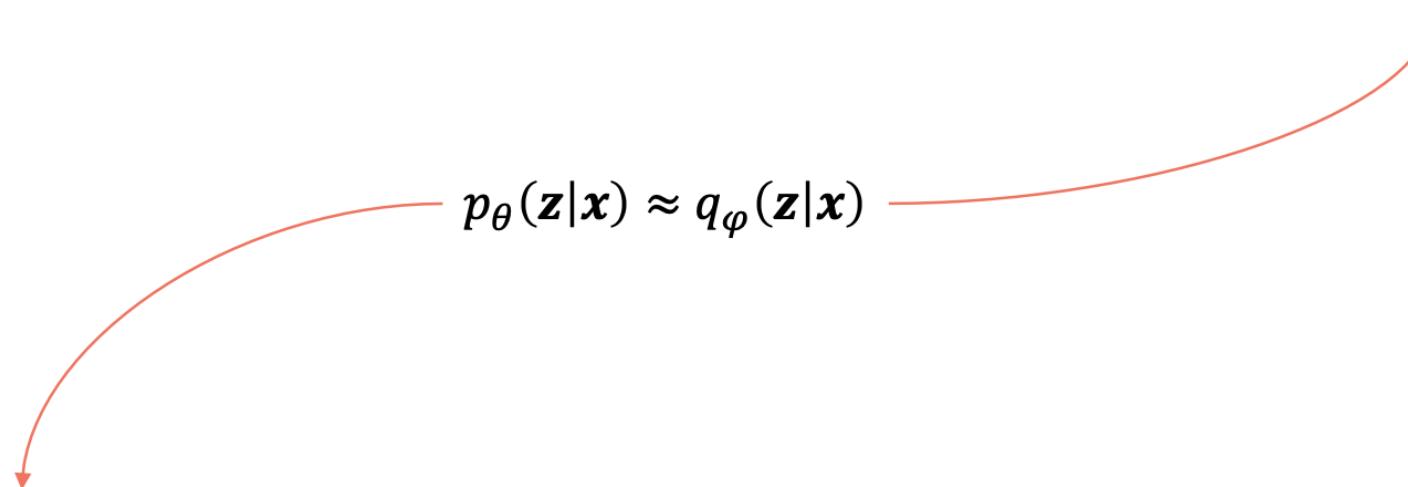


Any Solution?

An **approximate** posterior.
Parametrized by φ .

$$p_\theta(\mathbf{z}|\mathbf{x}) \approx q_\varphi(\mathbf{z}|\mathbf{x})$$

Our **true** posterior (that we can't evaluate due to its intractability)
Parametrized by θ .



$$\log p_{\theta}(\mathbf{x}) = \log p_{\theta}(\mathbf{x})$$

$$= \log p_{\theta}(\mathbf{x}) \int q_{\varphi}(\mathbf{z}|\mathbf{x}) d\mathbf{z}$$

Multiply by 1

This is the integral over the domain of a probability distribution which is always equal to 1

$$= \int \log p_{\theta}(\mathbf{x}) q_{\varphi}(\mathbf{z}|\mathbf{x}) d\mathbf{z}$$

Bring inside the integral

$$= E_{q_{\varphi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x})]$$

Definition of expectation

$$E_x[f(x)] = \int xf(x)dx$$

$$= E_{q_{\varphi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x})} \right]$$

Apply the equation $p_{\theta}(\mathbf{x}) = \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x})}$

$$= E_{q_{\varphi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z}) q_{\varphi}(\mathbf{z}|\mathbf{x})}{p_{\theta}(\mathbf{z}|\mathbf{x}) q_{\varphi}(\mathbf{z}|\mathbf{x})} \right]$$

Multiply by 1

$$= E_{q_{\varphi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\varphi}(\mathbf{z}|\mathbf{x})} \right] + E_{q_{\varphi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{q_{\varphi}(\mathbf{z}|\mathbf{x})}{p_{\theta}(\mathbf{z}|\mathbf{x})} \right]$$

Split the expectation

$$= E_{q_{\varphi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\varphi}(\mathbf{z}|\mathbf{x})} \right] + D_{KL} \left(q_{\varphi}(\mathbf{z}|\mathbf{x}) \| p_{\theta}(\mathbf{z}|\mathbf{x}) \right)$$

Definition of KL divergence

$$\underbrace{\quad}_{\geq 0}$$

Math is fun

$$\text{Total Compensation} = \text{Base Salary} + \underbrace{\text{Bonus}}_{\geq 0}$$



employee

$$\text{Total Compensation} \geq \text{Base Salary}$$

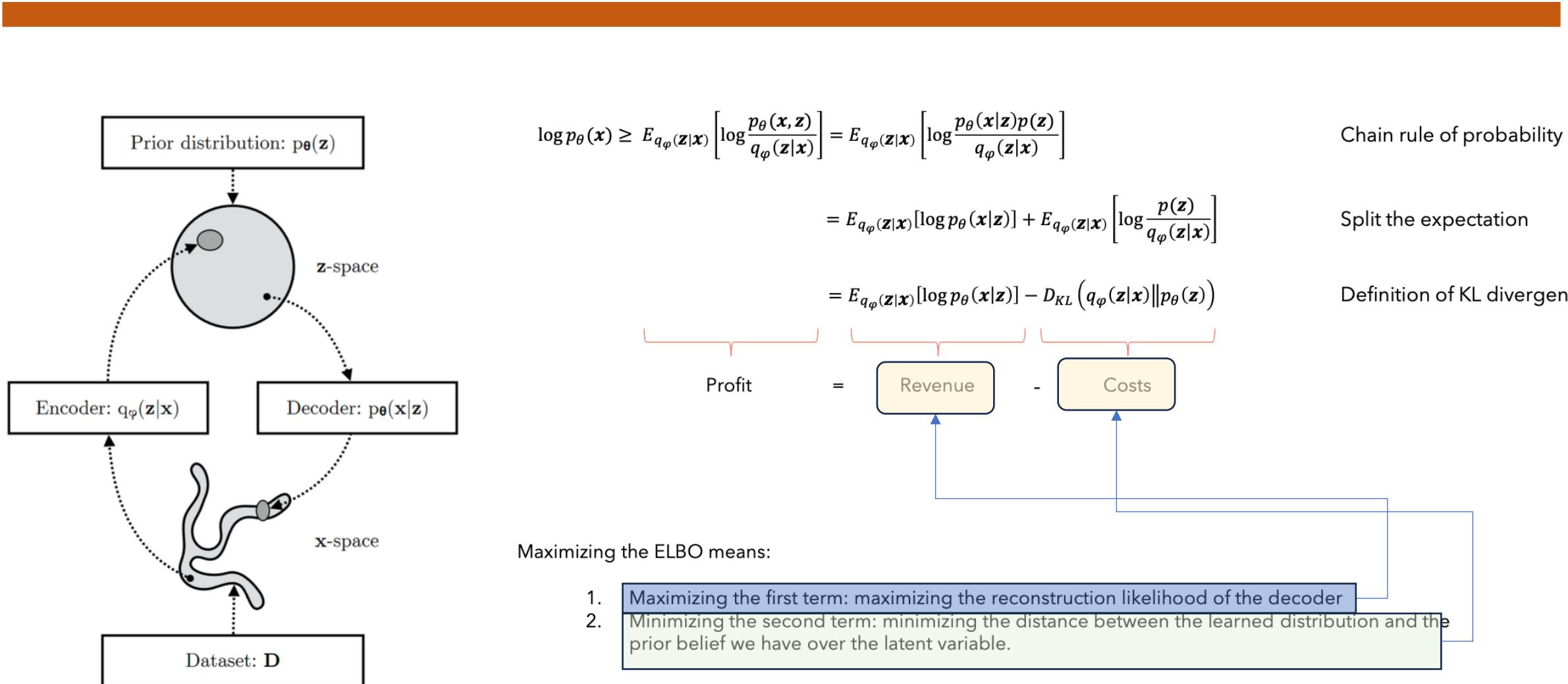
$$\log p_\theta(\mathbf{x}) = E_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] + D_{KL} \left(q_\phi(\mathbf{z}|\mathbf{x}) \| p_\theta(\mathbf{z}|\mathbf{x}) \right)$$

ELBO = Evidence Lower Bound **ELBO** ≥ 0



$$\log p_\theta(\mathbf{x}) \geq E_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right]$$

Math is fun



Maximizing the ELBO



When we have a **function** we want to **maximize**, we usually take the **gradient** and adjust the weights of the model so that they move **along the gradient direction**.

When we have a function we want to **minimize**, we usually take the **gradient**, and adjust the weights of the model so that they move **against the gradient direction**.

Stochastic Gradient Descent

When used to minimize the above function, a standard (or "batch") **gradient descent** method would perform the following iterations:

$$w := w - \eta \nabla Q(w) = w - \frac{\eta}{n} \sum_{i=1}^n \nabla Q_i(w),$$

where η is a step size (sometimes called the **learning rate** in machine learning).

$$L(\theta, \varphi, \mathbf{x}) = E_{q_\varphi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\varphi(\mathbf{z}|\mathbf{x})} \right] = E_{q_\varphi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - D_{KL} (q_\varphi(\mathbf{z}|\mathbf{x}) \| p_\theta(\mathbf{z}))$$

Maximizing the ELBO



ELBO

$$\mathcal{L}(\theta, \phi; \mathbf{x}^{(i)}) = -D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}^{(i)})||p_\theta(\mathbf{z})) + \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_\theta(\mathbf{x}^{(i)}|\mathbf{z})] \quad (3)$$

We want to differentiate and optimize the lower bound $\mathcal{L}(\theta, \phi; \mathbf{x}^{(i)})$ w.r.t. both the variational parameters ϕ and generative parameters θ . However, the gradient of the lower bound w.r.t. ϕ is a bit problematic. The usual (naïve) Monte Carlo gradient estimator for this type of problem is: $\nabla_\phi \mathbb{E}_{q_\phi(\mathbf{z})} [f(\mathbf{z})] = \mathbb{E}_{q_\phi(\mathbf{z})} [f(\mathbf{z}) \nabla_{q_\phi(\mathbf{z})} \log q_\phi(\mathbf{z})] \simeq \frac{1}{L} \sum_{l=1}^L f(\mathbf{z}) \nabla_{q_\phi(\mathbf{z}^{(l)})} \log q_\phi(\mathbf{z}^{(l)})$ where $\mathbf{z}^{(l)} \sim q_\phi(\mathbf{z}|\mathbf{x}^{(i)})$. This gradient estimator exhibits very high variance (see e.g. [BIP12]) and is impractical for our purposes.

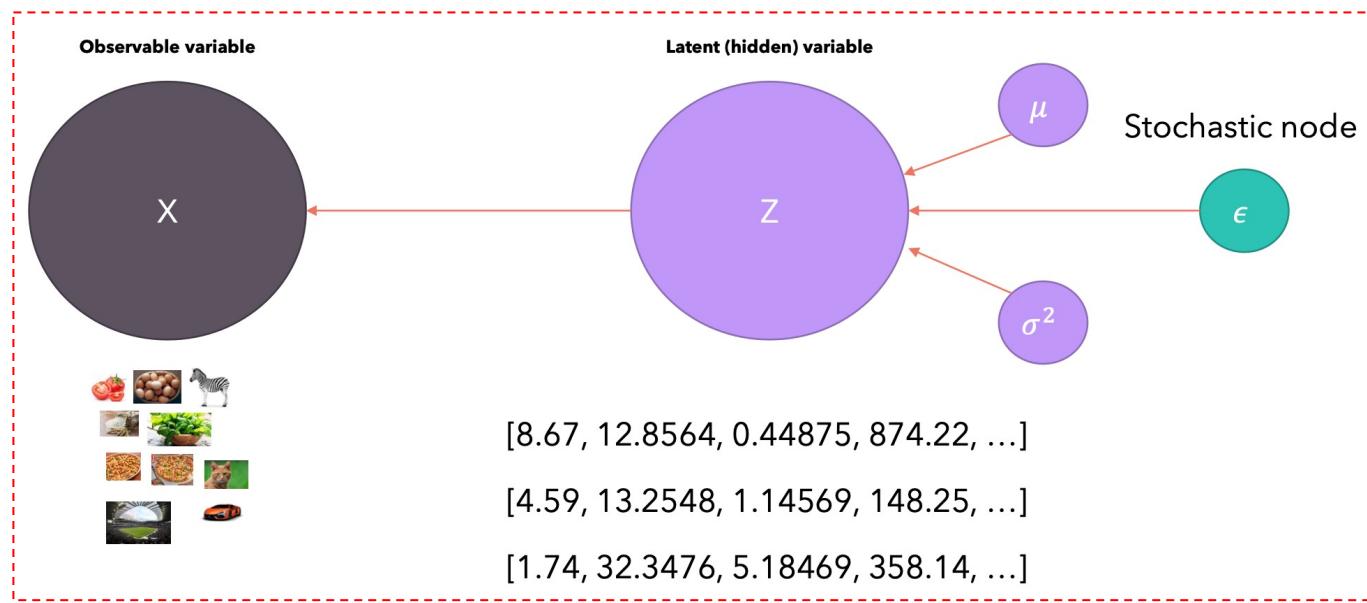
SCORE estimator

Kingma, D.P. and Welling, M., 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.

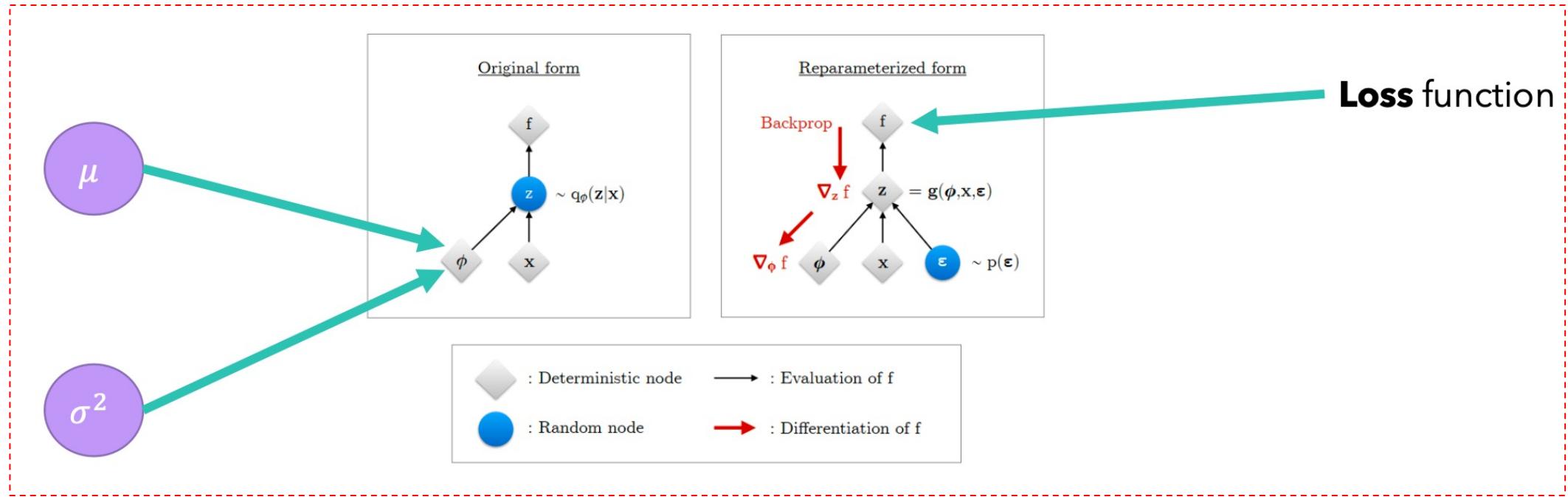
This estimator is **unbiased**, meaning that even if at every step it may not be equal to the true expectation, on average it will converge to it, but as it is stochastic, it also has a variance and it happens to be high for practical use. Plus, we can't run backpropagation through it!

We need a new estimator

The reparameterization trick



The reparameterization trick



A new estimator



$$L(\theta, \varphi, \mathbf{x}) = E_{q_\varphi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\varphi(\mathbf{z}|\mathbf{x})} \right] = E_{p(\epsilon)} \left[\log \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\varphi(\mathbf{z}|\mathbf{x})} \right]$$

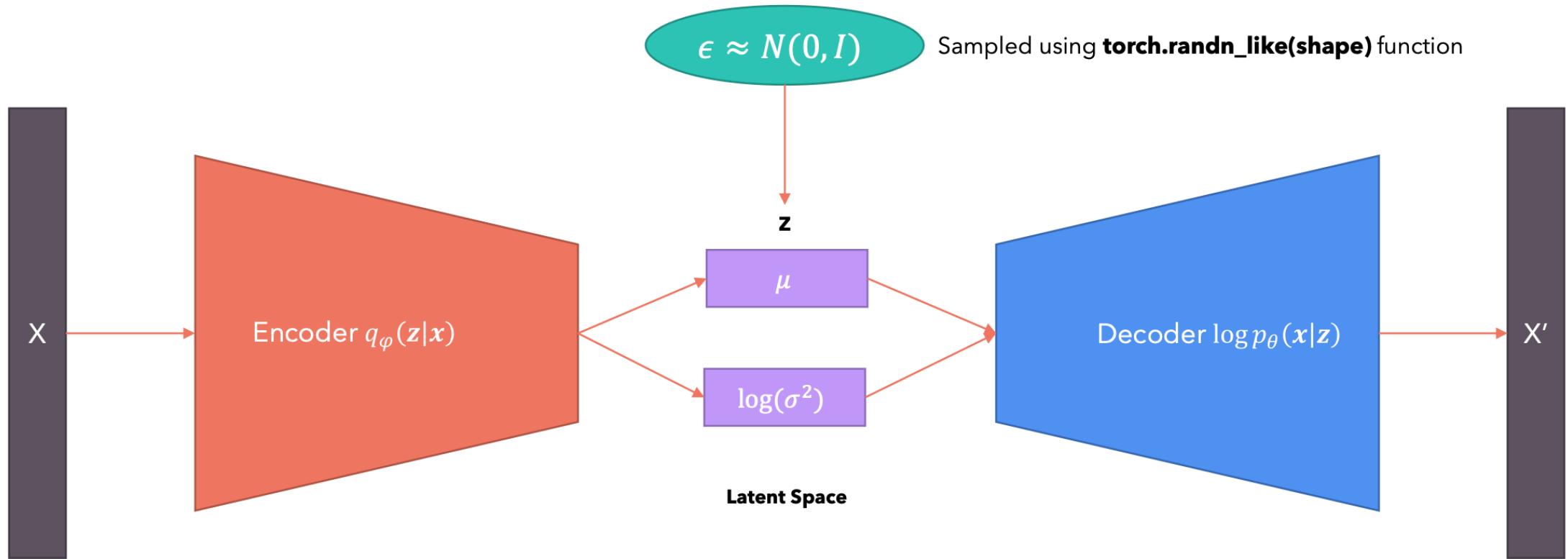
$$\epsilon \approx p(\epsilon)$$

ELBO

$$\mathbf{z} = g(\varphi, \mathbf{x}, \epsilon)$$

$$E_{p(\epsilon)} \left[\log \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\varphi(\mathbf{z}|\mathbf{x})} \right] \cong \tilde{L}(\theta, \varphi, \mathbf{x}) = \log \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\varphi(\mathbf{z}|\mathbf{x})}$$

VAE: Example Network



Input

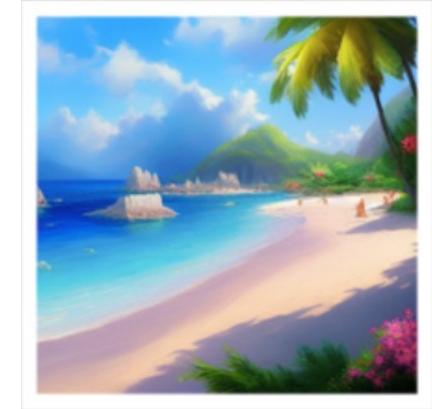
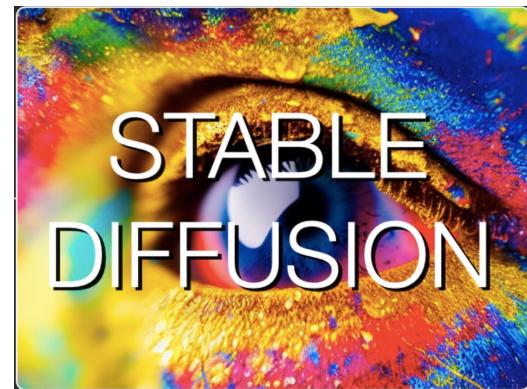


We prefer learning $\log(\sigma^2)$ because it can be negative, so the model doesn't need to be forced to produce only positive values for it.

Reconstructed Input

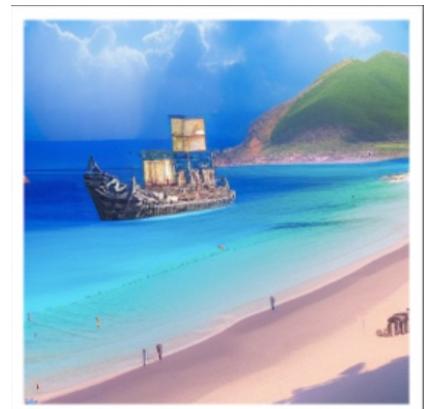
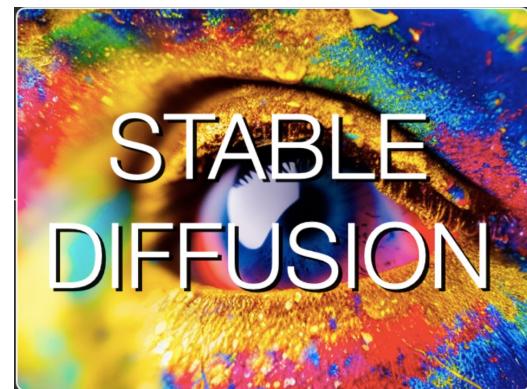
Stable Diffusion: Example

Paradise
Cosmic Beach
Input



Output

Pirate Ship



Output

Input

Outline

- **Objective**
- **Introduction to Stable Diffusion**
- **Stable Diffusion: Motivation**
- **Stable Diffusion: Clearly Explained**
- **Stable Diffusion: Demo with API**
- **Summary**

The Components of Stable Diffusion

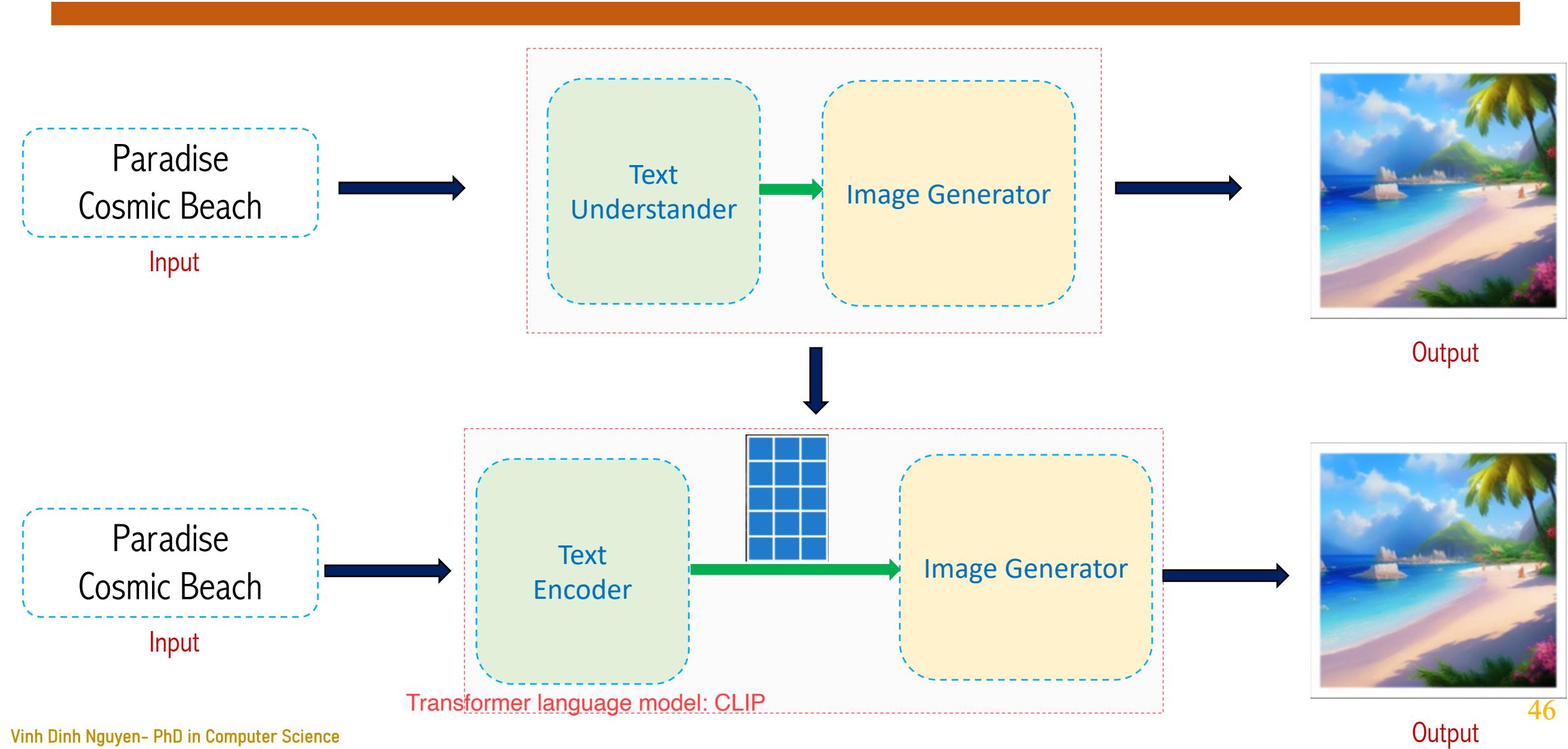
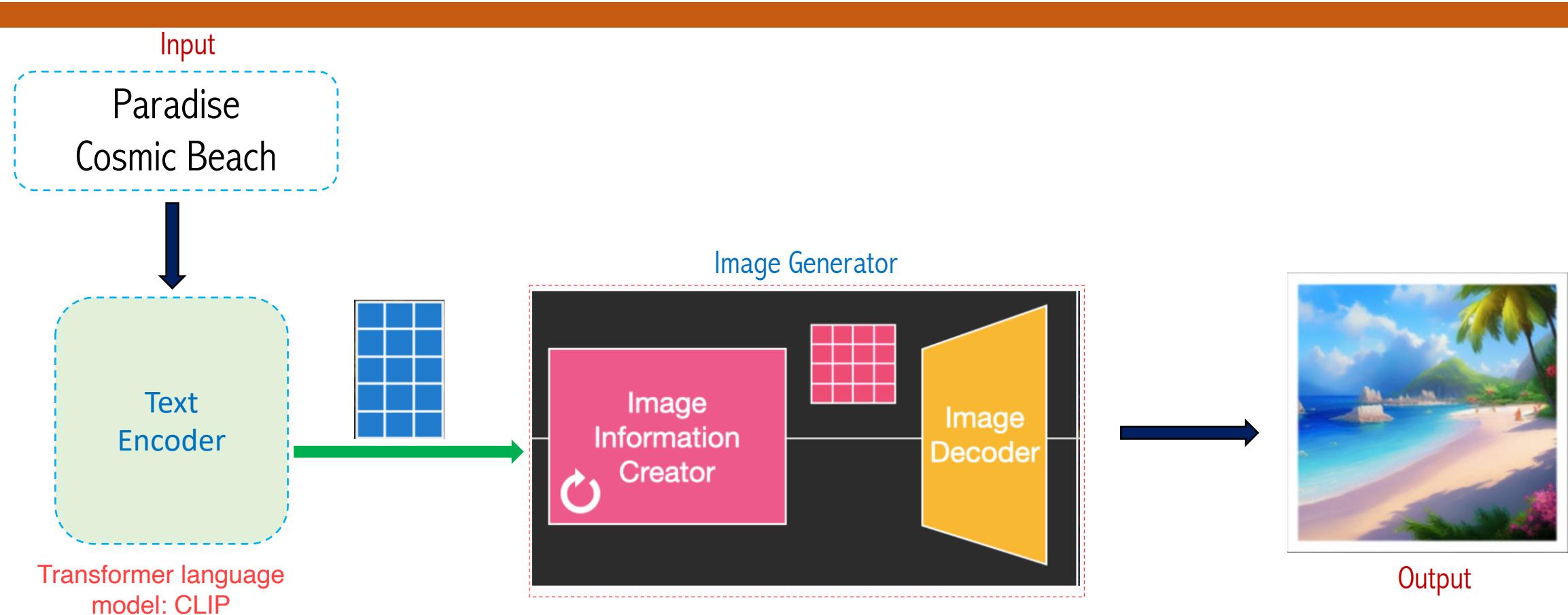
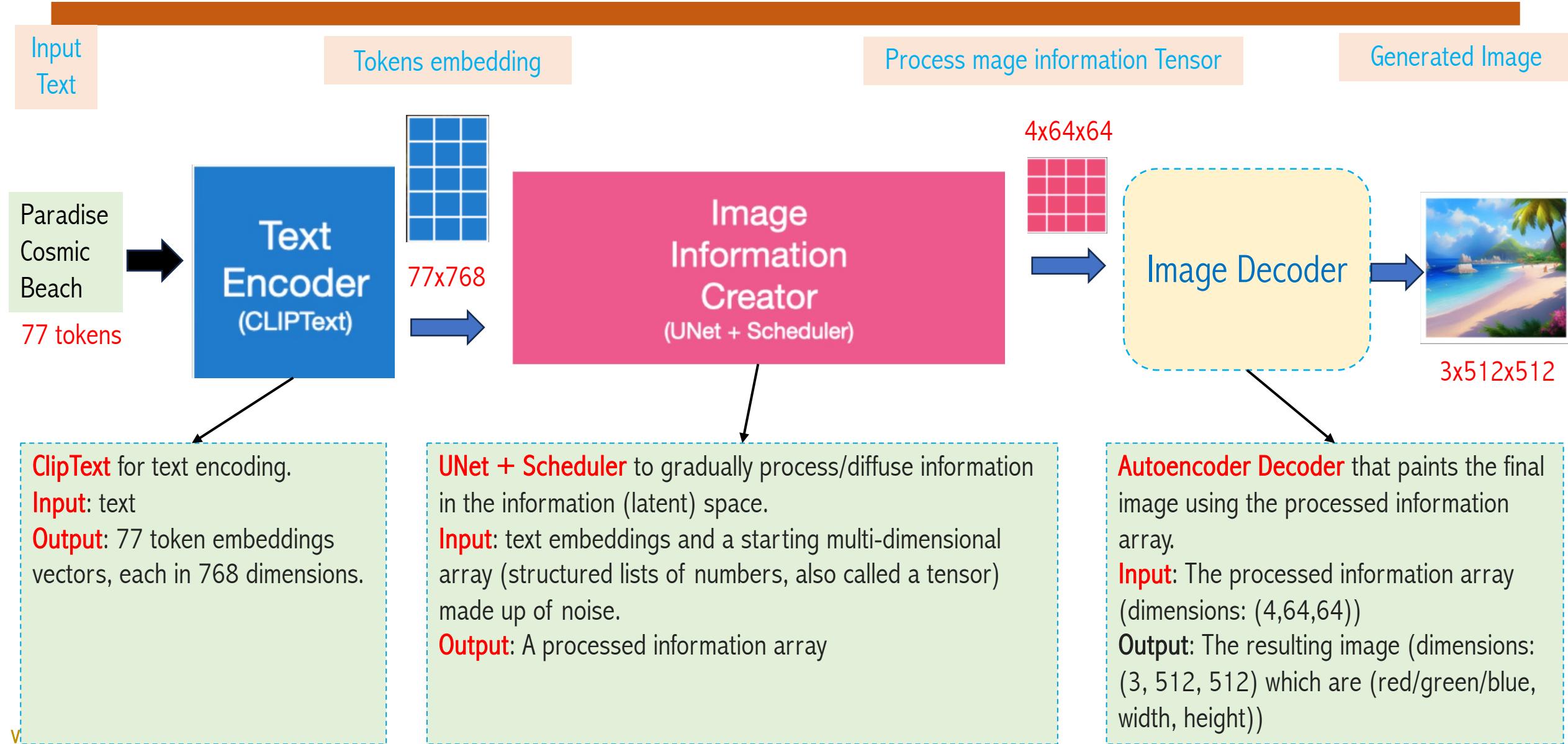


Image Generator



Stable Diffusion: Model



Stable Diffusion: Model

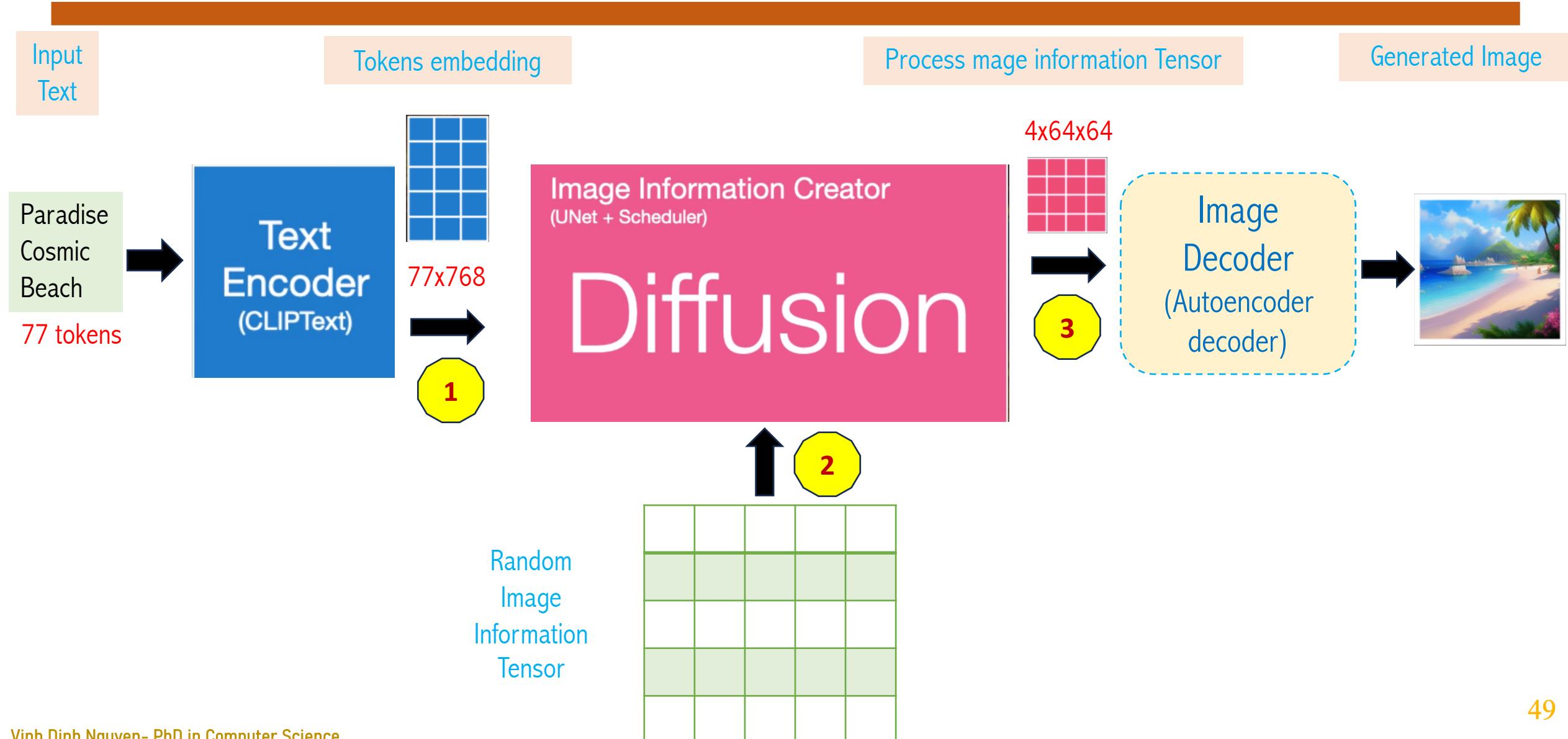
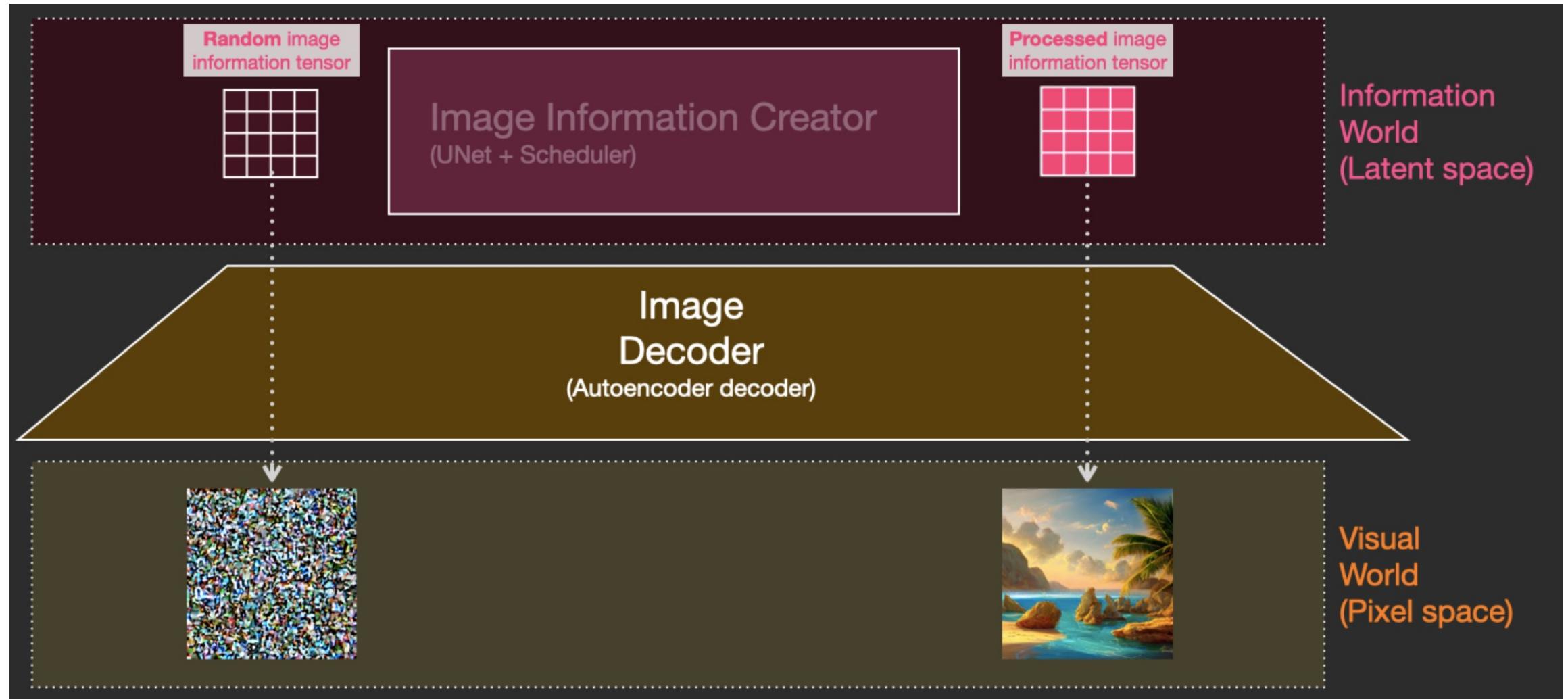
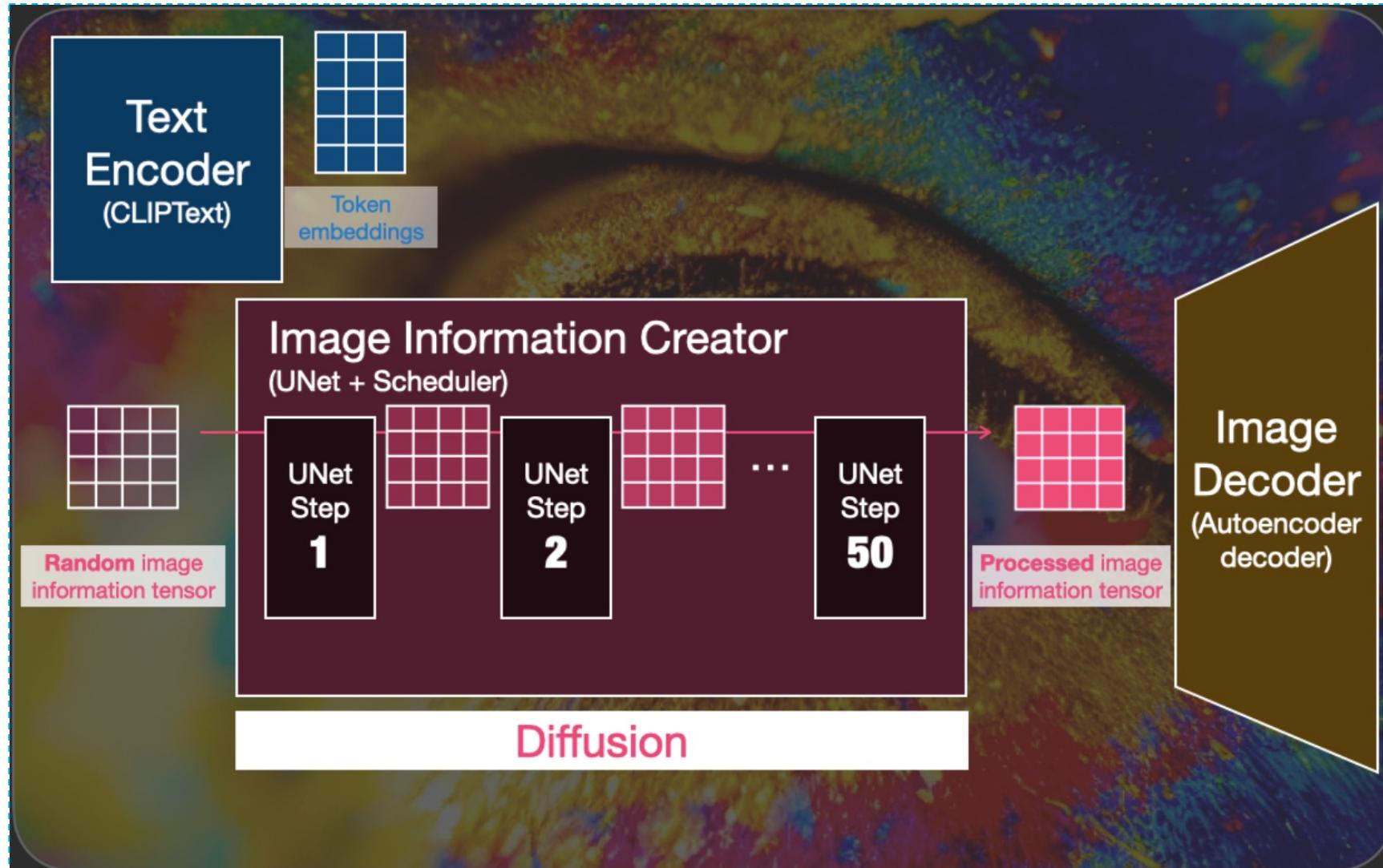


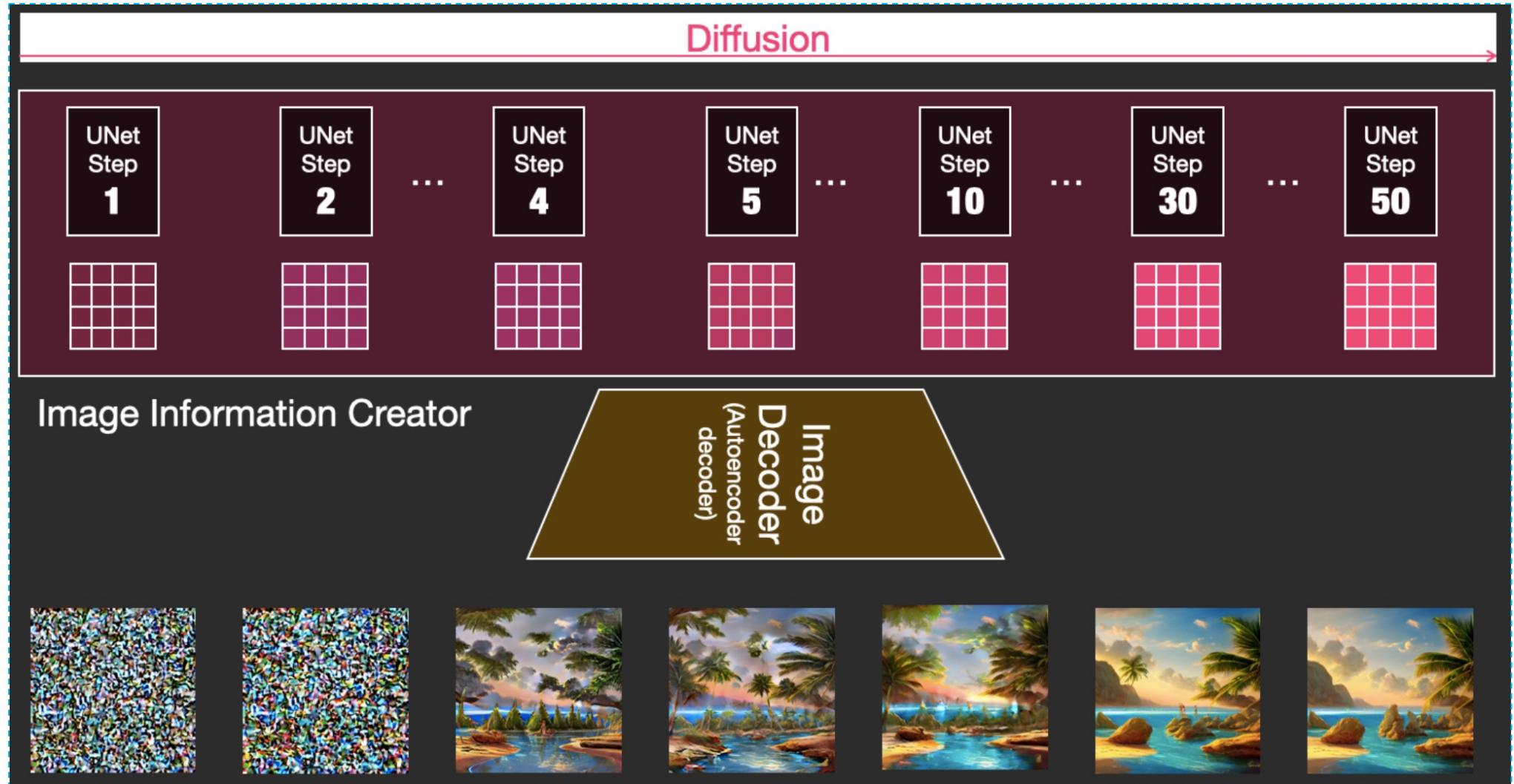
Image Decoder: Visualization



Diffusion Process: Where



Laten Feature: Visualization



Forward Diffusion: Image Generation

Pickup an image

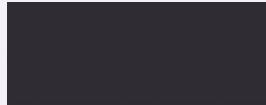


Generate some random noise



Pickup an amount of noise

0



1



2



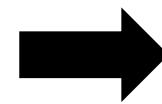
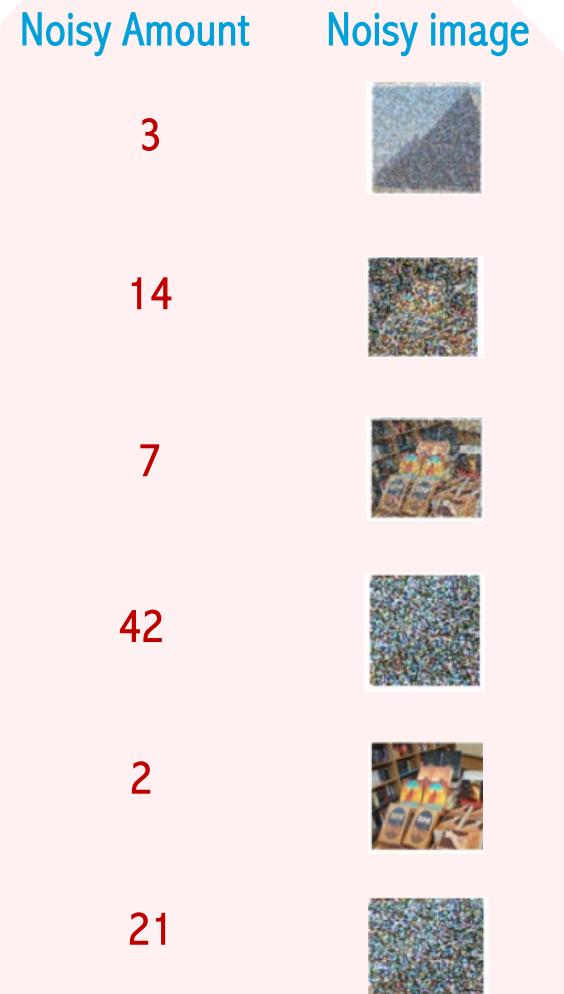
3



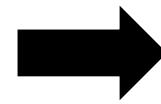
Add noise to the image in that amount



Diffusion Training: Noise Predictor



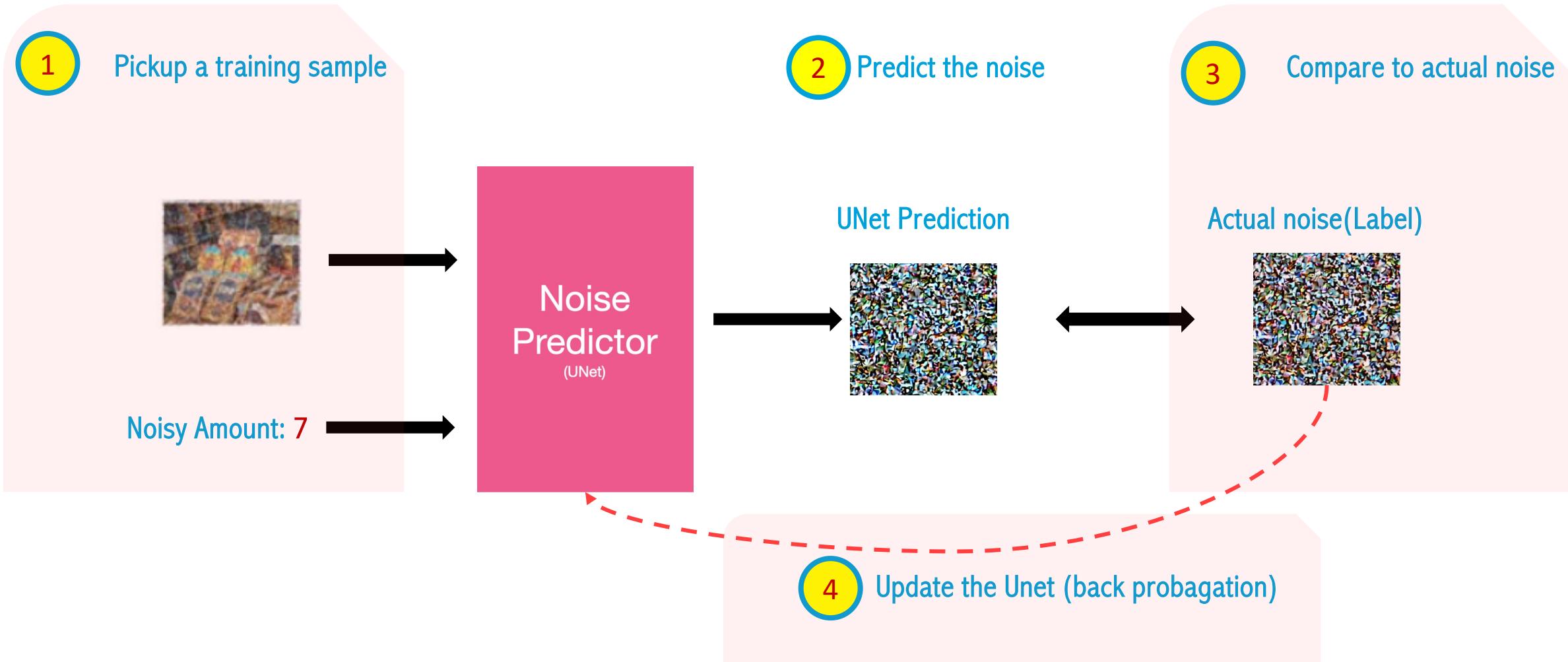
Noise
Predictor
(UNet)



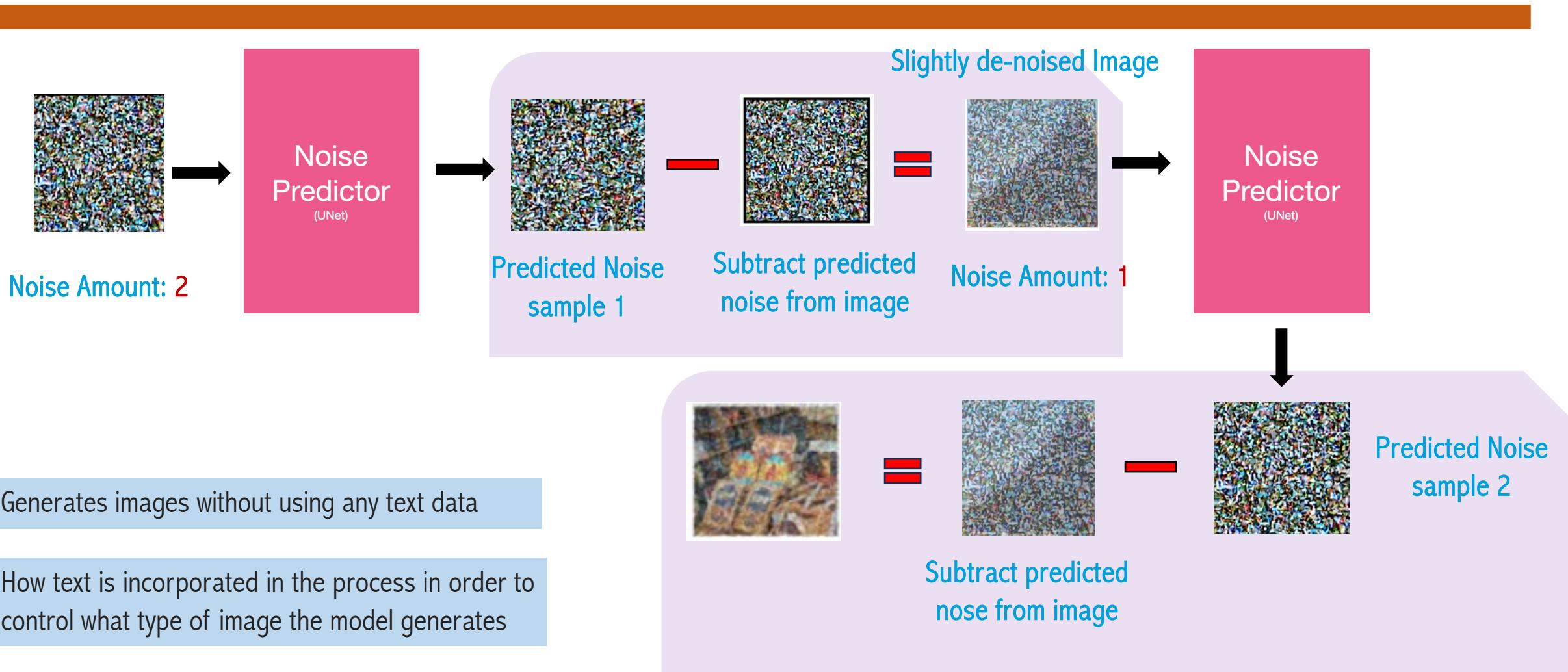
Noise sample



Diffusion Training: Noise Predictor



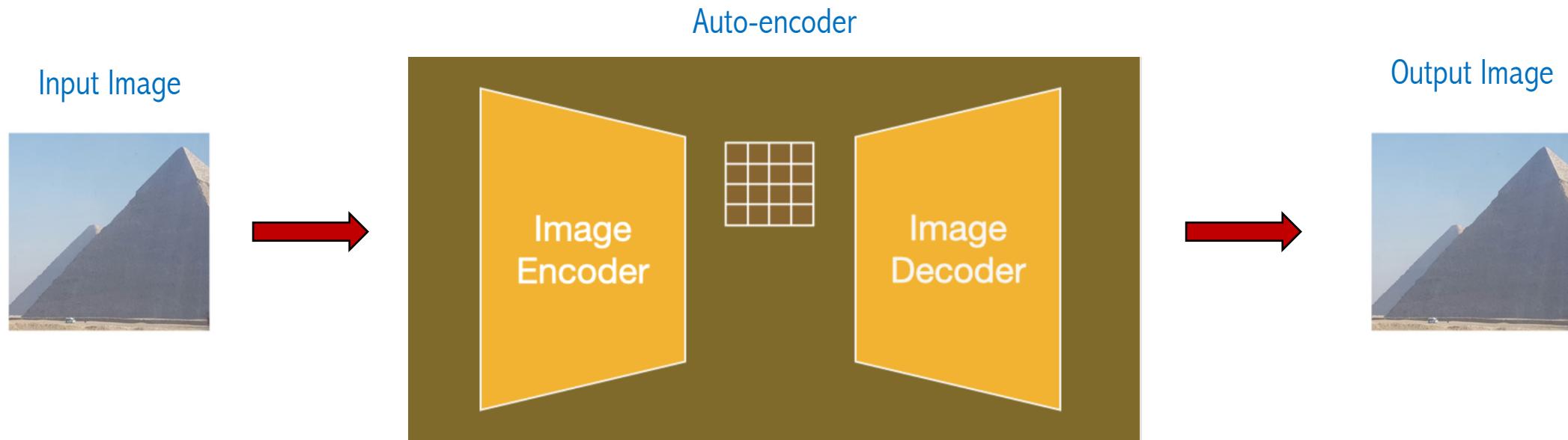
Reverse Diffusion (Denoising)



Speed Boost

- Speed Boost: Diffusion on Compressed (Latent) Data Instead of the Pixel Image

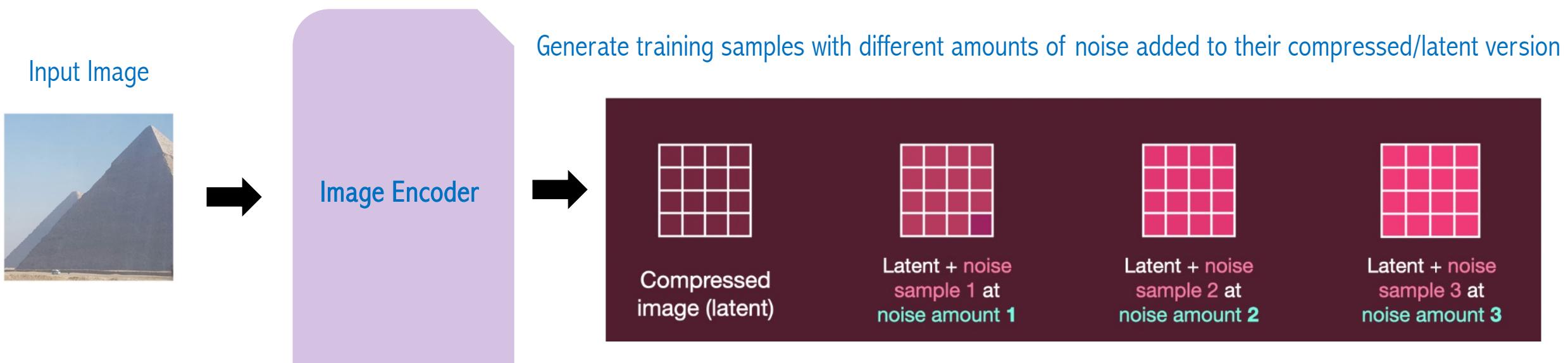
The Stable Diffusion paper runs the diffusion process not on the pixel images themselves, but on a compressed version of the image



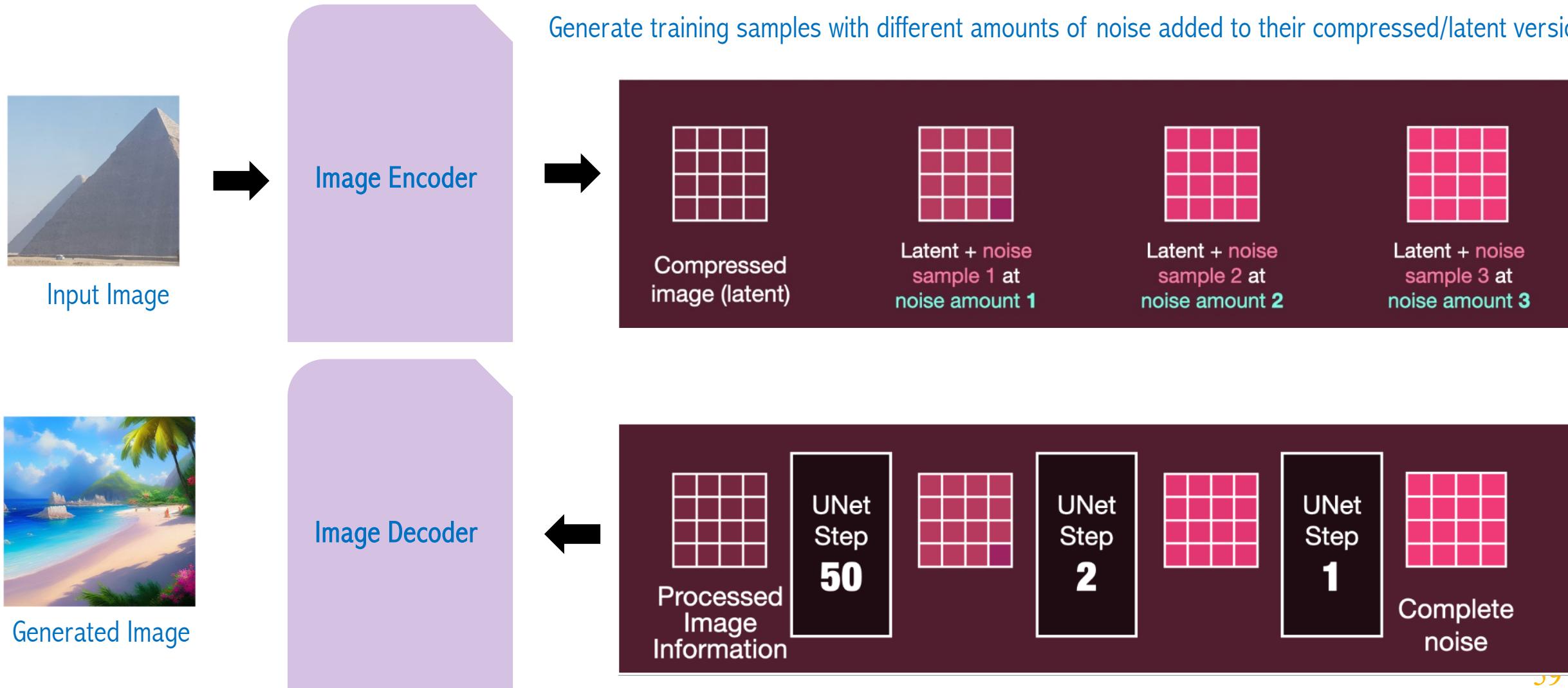
Forwarding Process

- Speed Boost: Diffusion on Compressed (Latent) Data Instead of the Pixel Image

The Stable Diffusion paper runs the diffusion process not on the pixel images themselves, but on a compressed version of the image

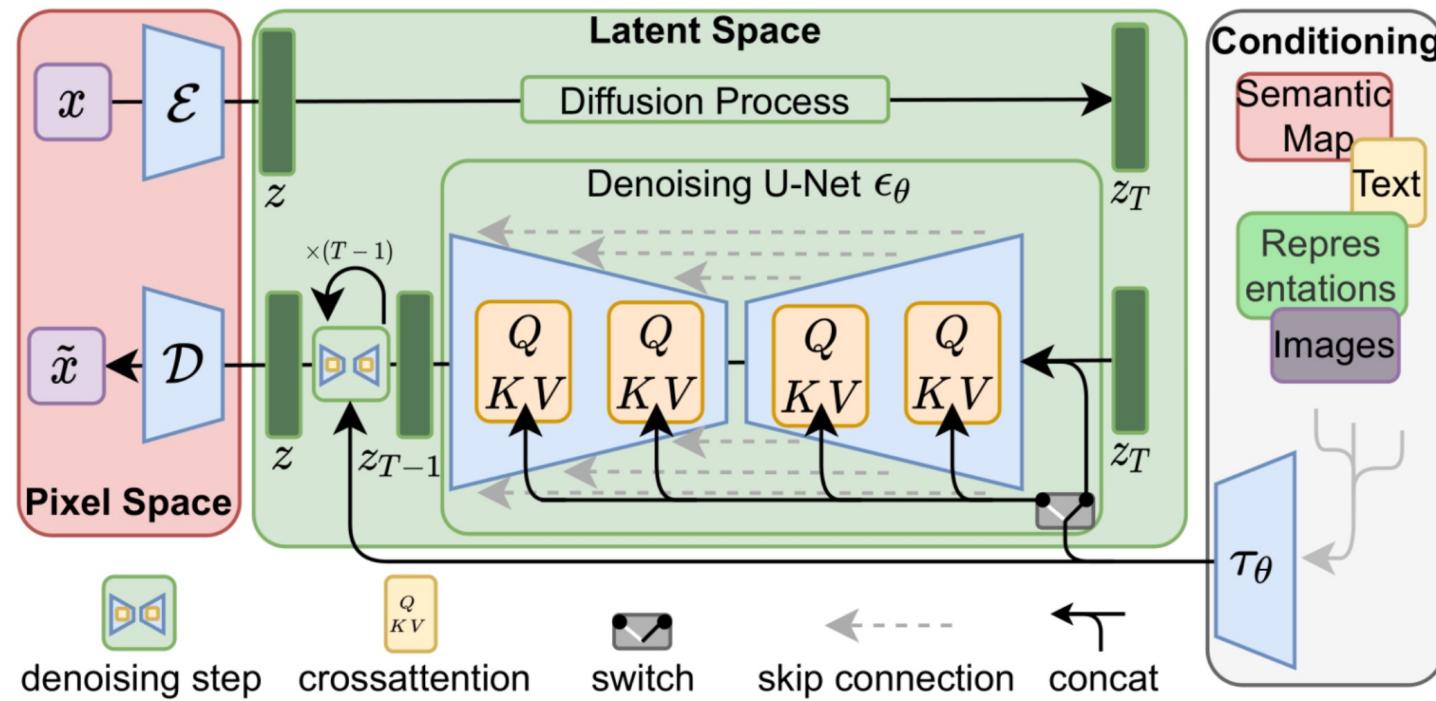


Reverse Process



LDM/Stable Diffusion

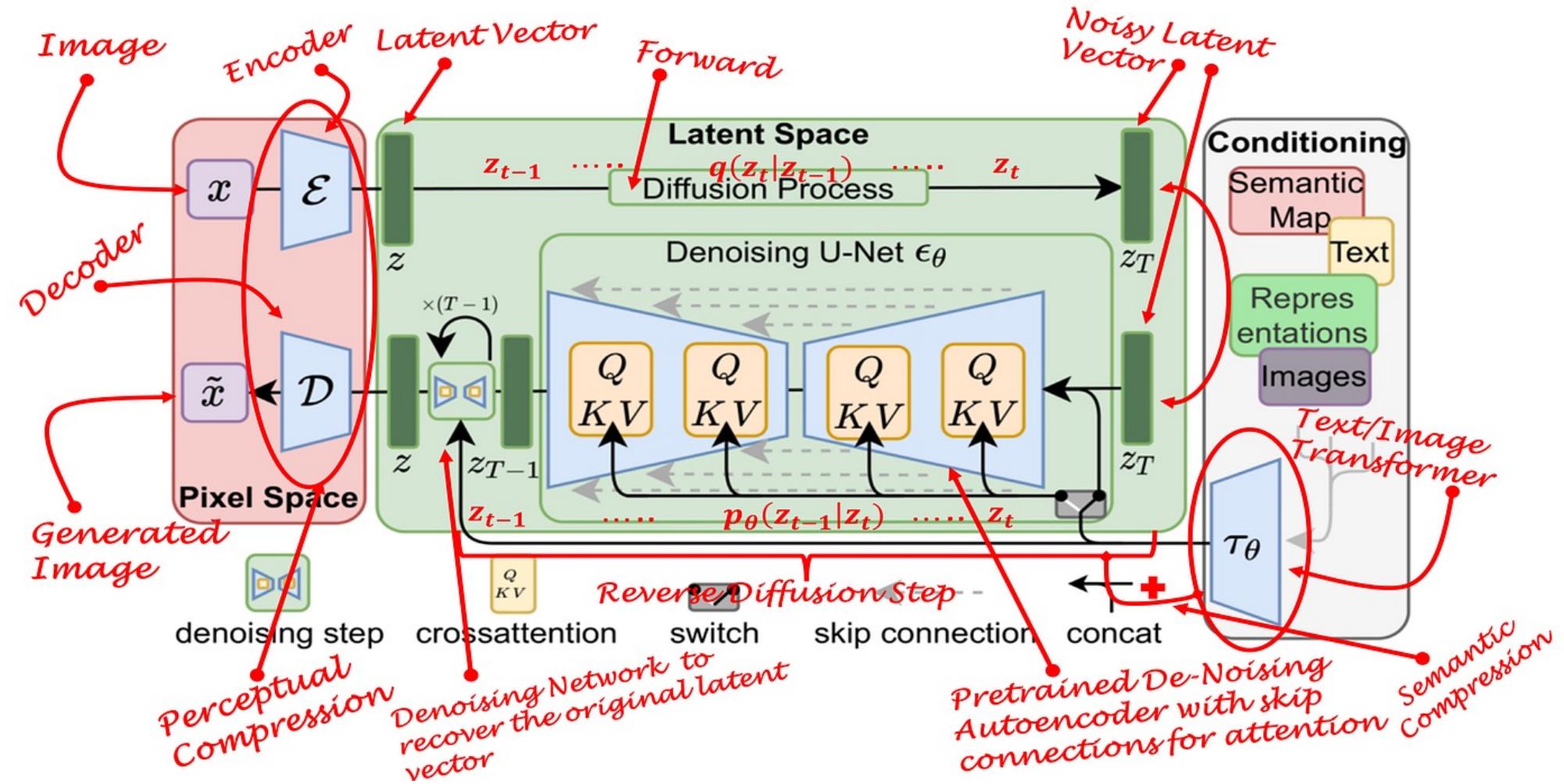
The text prompts describing what image the model should generate



$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}}\right) \cdot \mathbf{V}$$

where $\mathbf{Q} = \mathbf{W}_Q^{(i)} \cdot \varphi_i(\mathbf{z}_i)$, $\mathbf{K} = \mathbf{W}_K^{(i)} \cdot \tau_\theta(y)$, $\mathbf{V} = \mathbf{W}_V^{(i)} \cdot \tau_\theta(y)$
and $\mathbf{W}_Q^{(i)} \in \mathbb{R}^{d \times d_e^i}$, $\mathbf{W}_K^{(i)}, \mathbf{W}_V^{(i)} \in \mathbb{R}^{d \times d_\tau}$, $\varphi_i(\mathbf{z}_i) \in \mathbb{R}^{N \times d_e^i}$, $\tau_\theta(y) \in \mathbb{R}^{M \times d_\tau}$

LDM/Stable Diffusion



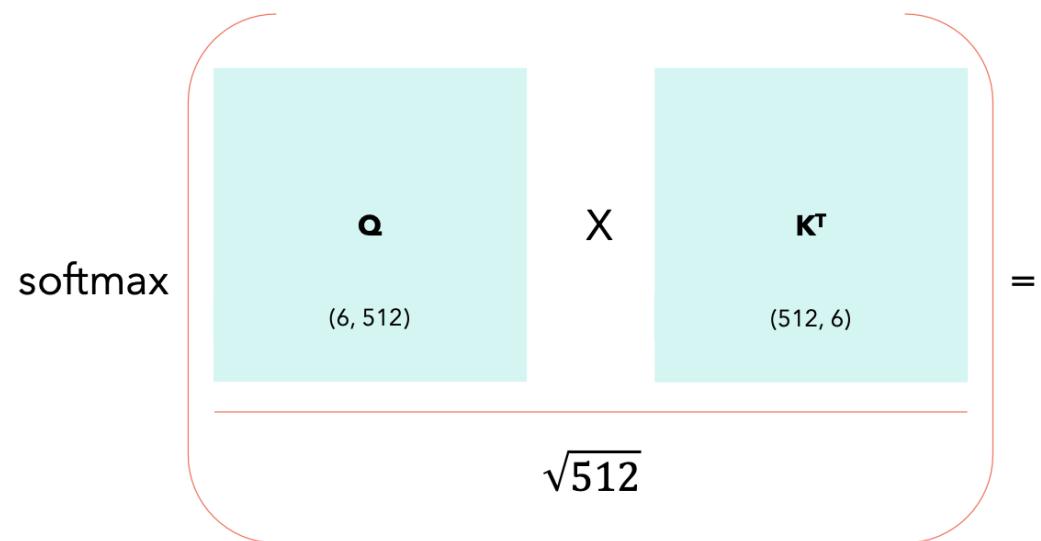
What is Self-Attention?

Self-Attention allows the model to relate words to each other.

In this simple case we consider the sequence length **seq** = 6 and $d_{\text{model}} = d_k = 512$.

The matrices **Q**, **K** and **V** are just the input sentence.

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



* for simplicity I considered only one head, which makes $d_{\text{model}} = d_k$.

	YOUR	CAT	IS	A	LOVELY	CAT	Σ
YOUR	0.268	0.119	0.134	0.148	0.179	0.152	1
CAT	0.124	0.278	0.201	0.128	0.154	0.115	1
IS	0.147	0.132	0.262	0.097	0.218	0.145	1
A	0.210	0.128	0.206	0.212	0.119	0.125	1
LOVELY	0.146	0.158	0.152	0.143	0.227	0.174	1
CAT	0.195	0.114	0.203	0.103	0.157	0.229	1

* all values are random.

(6, 6)

What is Self-Attention?

	YOUR	CAT	IS	A	LOVELY	CAT
YOUR	0.268	0.119	0.134	0.148	0.179	0.152
CAT	0.124	0.278	0.201	0.128	0.154	0.115
IS	0.147	0.132	0.262	0.097	0.218	0.145
A	0.210	0.128	0.206	0.212	0.119	0.125
LOVELY	0.146	0.158	0.152	0.143	0.227	0.174
CAT	0.195	0.114	0.203	0.103	0.157	0.229

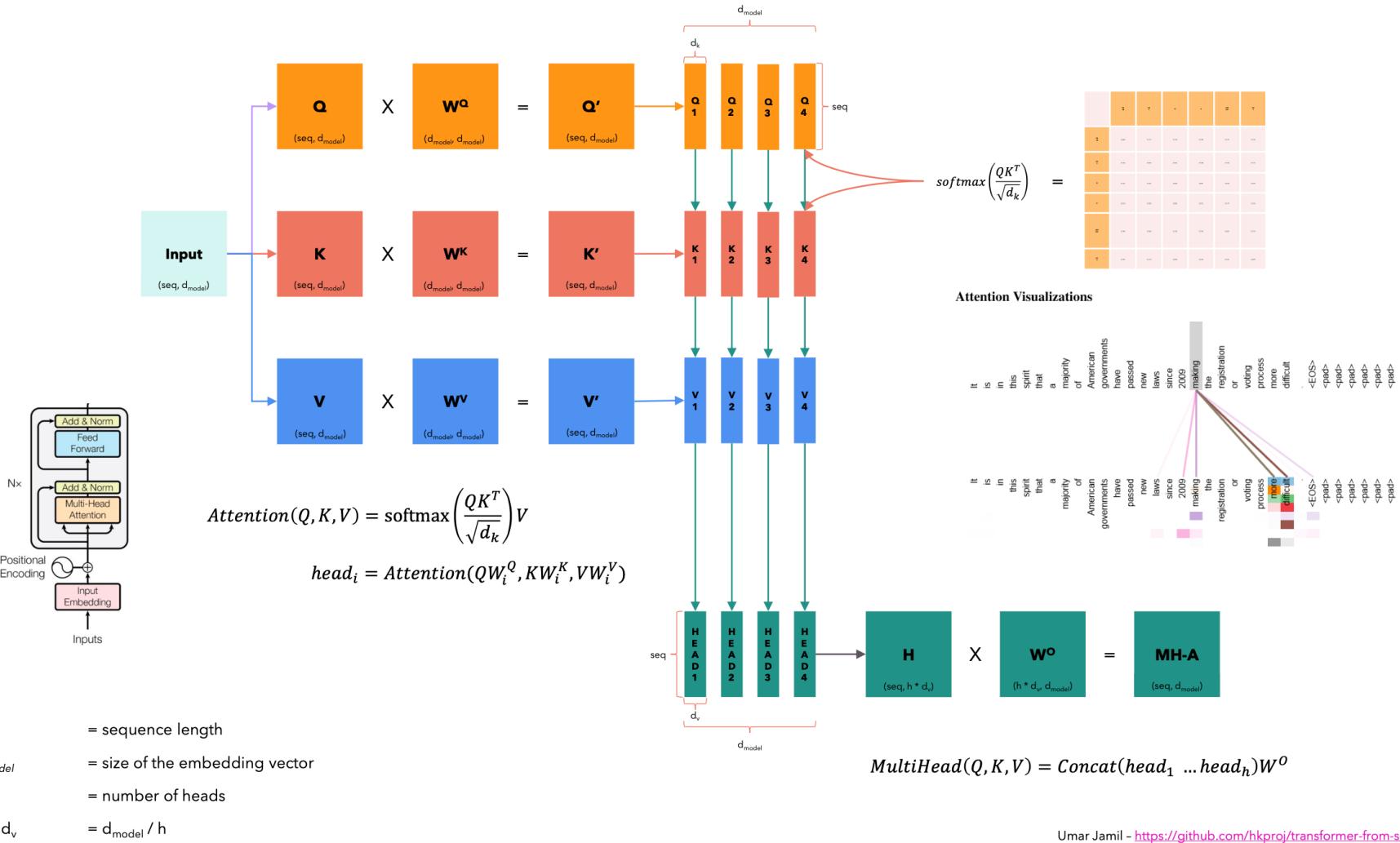
(6, 6)

$$X \quad v \quad = \quad \text{Attention} \\ (6, 512) \quad (6, 512)$$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

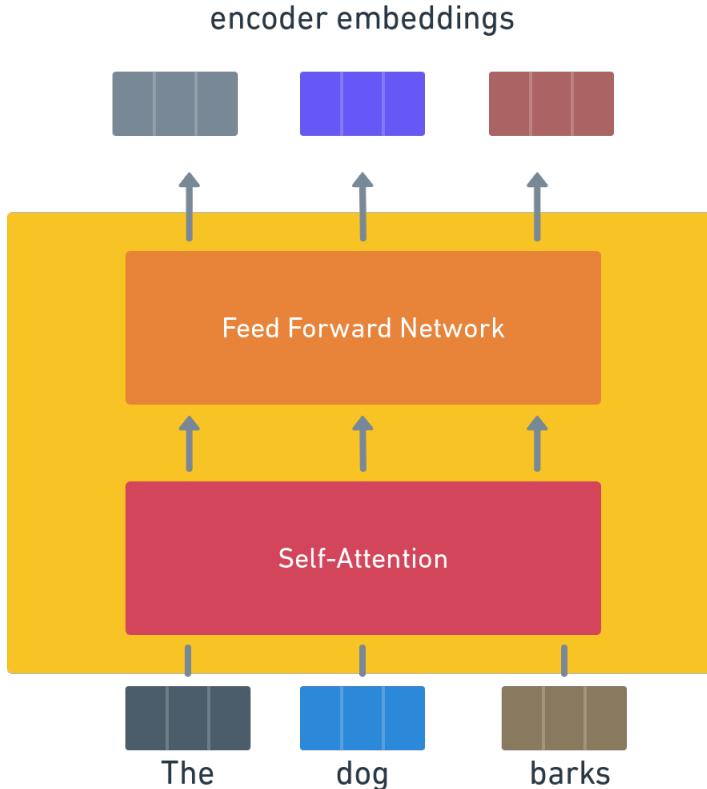
Each row in this matrix captures not only the meaning (given by the embedding) or the position in the sentence (represented by the positional encodings) but also each word's interaction with other words.

What is Self-Attention?

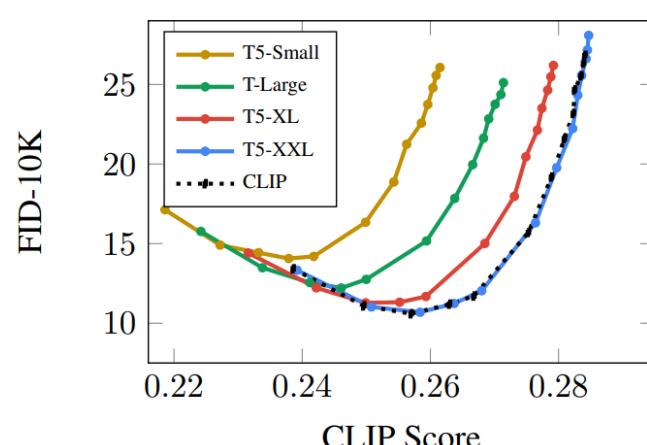


Umar Jamil - <https://github.com/hkproj/transformer-from-scratch-notes>

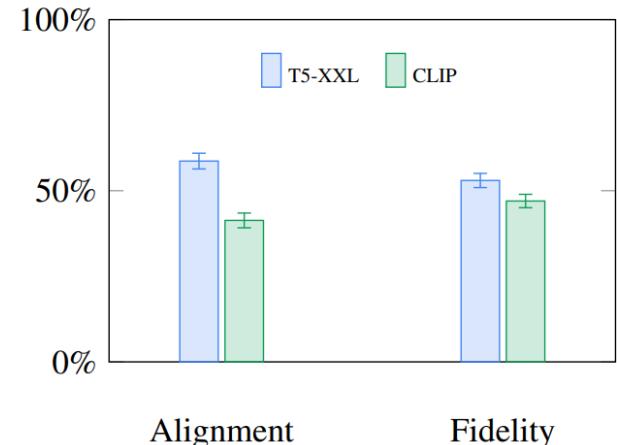
The Text Encoder: A Transformer Language Model



Takes the text prompt and produces token embeddings



(a) Pareto curves comparing various text encoders.



(b) Comparing T5-XXL and CLIP on DrawBench.

Larger/better language models have a significant effect on the quality of image generation models

How CLIP is trained

Image



Caption

Photo pour Japanese pagoda and old house in Kyoto at twilight - image libre de droit



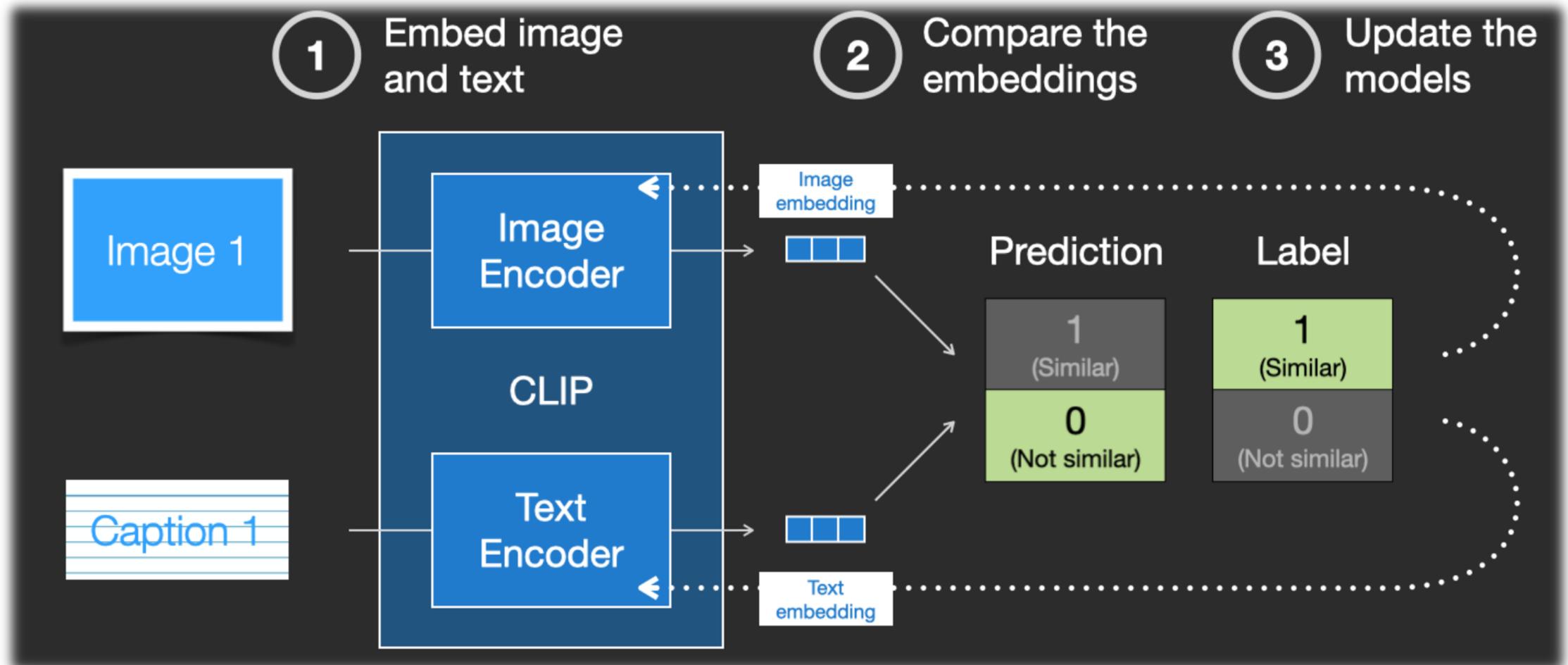
Soaring by Peter Eades



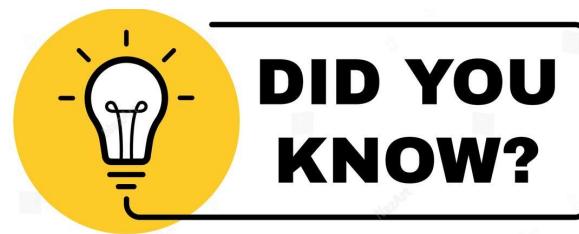
far cry 4 concept art is the reason why it's a beautiful game vg247.
Black Bedroom Furniture Sets. Home Design Ideas

The early Stable Diffusion models just plugged in the pre-trained ClipText model released by OpenAI

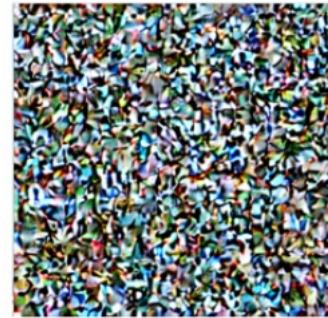
How CLIP is trained



Feeding Text Information Into The Image Generation Process

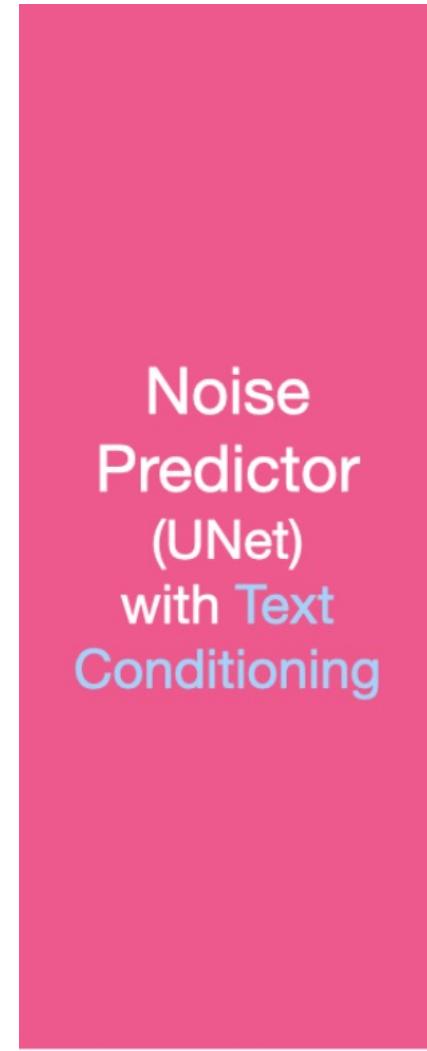
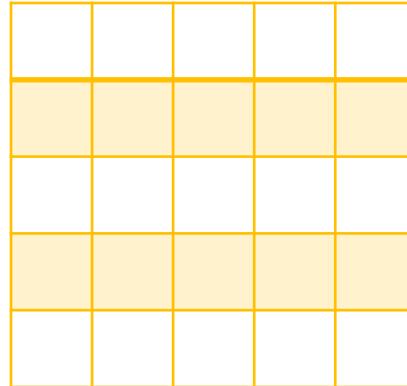


Noise predictor with Text Input

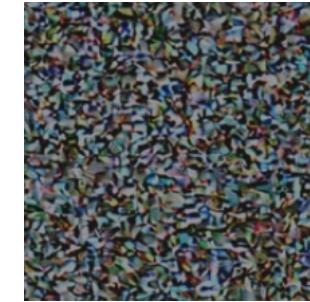


Noise Amount: 2

Prompt text information (token embeddings)



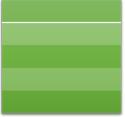
Predicted Noise Sample



Diffusion Training: Noise Predictor

Noisy Amount Noisy image Text Promt

3



14



7



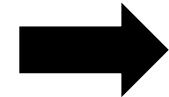
42



2



21



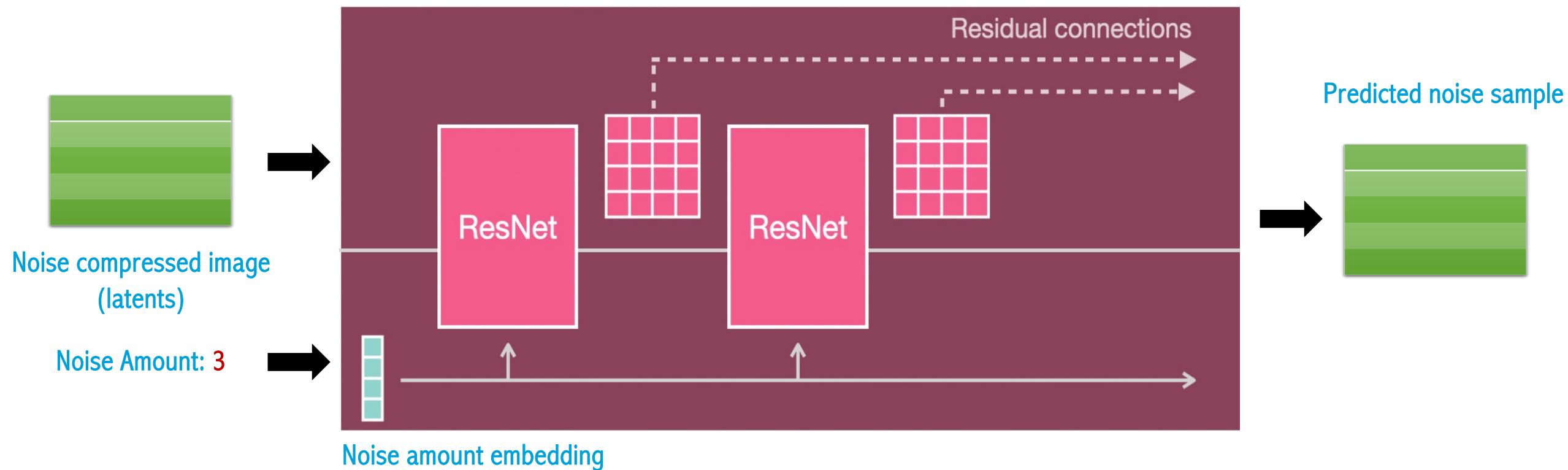
Noise
Predictor
(UNet)



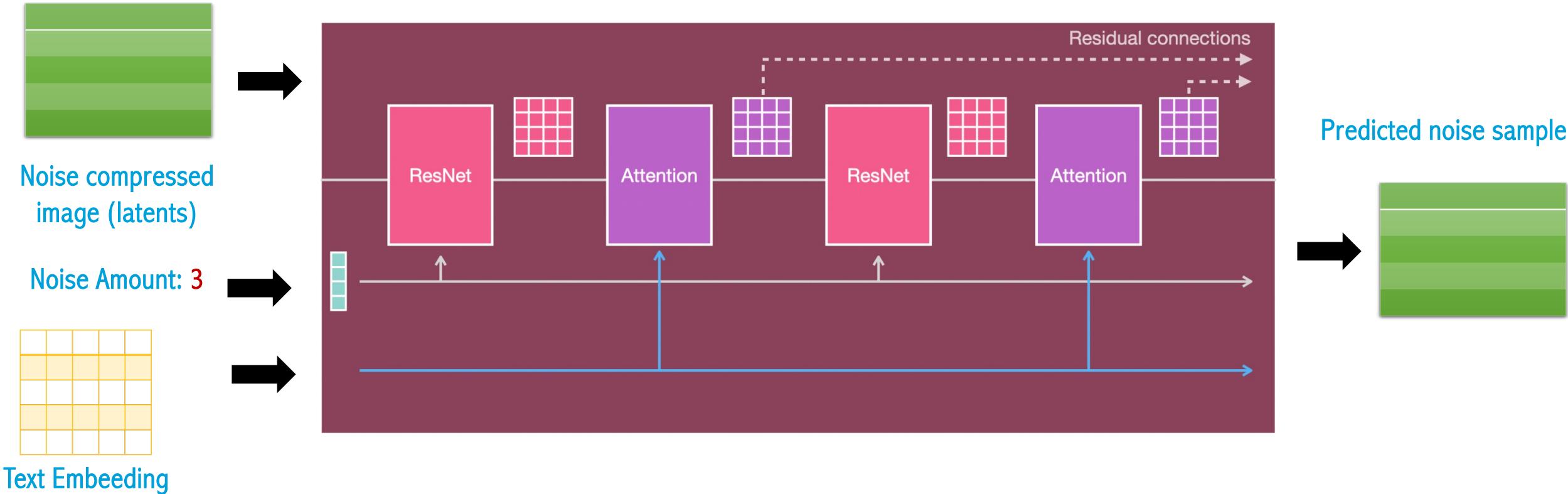
Noise sample



Unet Noise predictor (without text)



Unet Noise predictor (with text)



Latent Diffusion Loss: Discussion

$$L_{DM} = \mathbb{E}_{x, \epsilon \sim \mathcal{N}(0,1), t} \left[\|\epsilon - \epsilon_\theta(x_t, t)\|_2^2 \right]$$

$$L_{LDM} := \mathbb{E}_{\mathcal{E}(x), \epsilon \sim \mathcal{N}(0,1), t} \left[\|\epsilon - \epsilon_\theta(z_t, t)\|_2^2 \right]$$

Latent Vector

Latent State of the
Denoising U-Net at
time 't'

Domain Specific Encoder
(i.e. Transformer) on a
condition 'y'

Conditional Loss Function

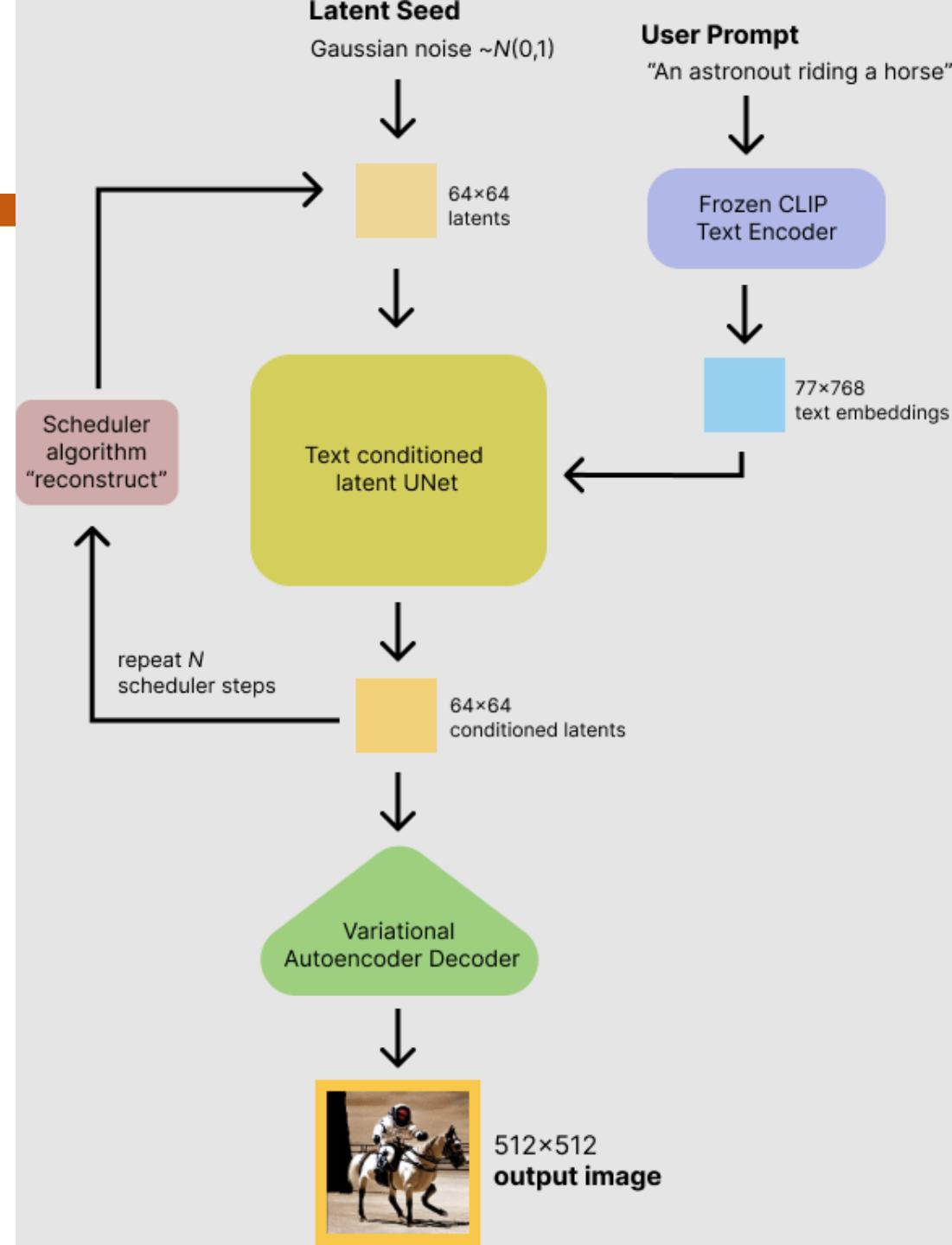
$$L_{LDM} := \mathbb{E}_{\mathcal{E}(x), y, \epsilon \sim \mathcal{N}(0,1), t} \left[\|\epsilon - \epsilon_\theta(z_t, t, \tau_\theta(y))\|_2^2 \right]$$

QUIZ TIME

Stable Diffusion Inference Process

There are mainly three main components in latent diffusion:

1. An autoencoder (VAE).
2. A U-Net.
3. A text-encoder, e.g. CLIP's Text Encoder.



VAE: Implementation

During latent diffusion training, the encoder converts a 512*512*3 image into a low dimensional latent representation of image of size say 64*64*4 for the forward diffusion process

```
[2] from torchvision import transforms as tfms
from diffusers import AutoencoderKL

# Load the autoencoder model which will be used to decode the latents into image space.
vae = AutoencoderKL.from_pretrained("CompVis/stable-diffusion-v1-4", subfolder="vae")

# To the GPU we go!
vae = vae.to("cpu")

# Convert PIL image to latents

def pil_to_latent(input_im):
    # Single image -> single latent in a batch (so size 1, 4, 64, 64)
    with torch.no_grad():
        latent = vae.encode(tfms.ToTensor()(input_im).unsqueeze(0).to(torch_device)*2-1) # Note scaling
    return 0.18215 * latent.latent_dist.sample()
```

Unet: Implementation

The Unet that takes in the noisy latents (x) and predicts the noise. We use a conditional model that also takes in the timestep (t) and our text embedding as guidance.

```
from diffusers import UNet2DConditionModel

# The UNet model for generating the latents.
unet = UNet2DConditionModel.from_pretrained("CompVis/stable-diffusion-v1-4", subfolder="unet")

# To the GPU
unet = unet.to(torch_device);
noise_pred = unet(latents, t, encoder_hidden_states=text_embeddings)["sample"]
```

The Text-encoder: Implementation

The text-encoder transforms the input prompt into an embedding space that goes as input to the U-Net

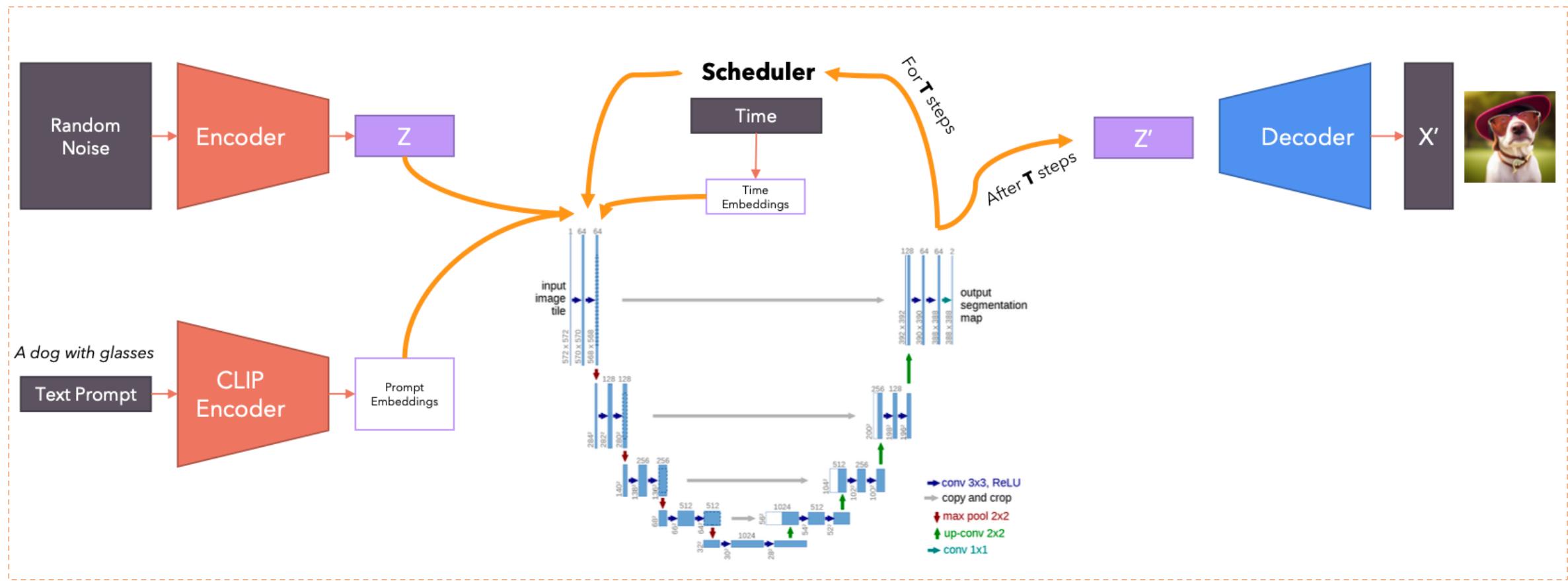
```
from transformers import CLIPTextModel, CLIPTokenizer

# Load the tokenizer and text encoder to tokenize and encode the text.
tokenizer = CLIPTokenizer.from_pretrained("openai/clip-vit-large-patch14")
text_encoder = CLIPTextModel.from_pretrained("openai/clip-vit-large-patch14")

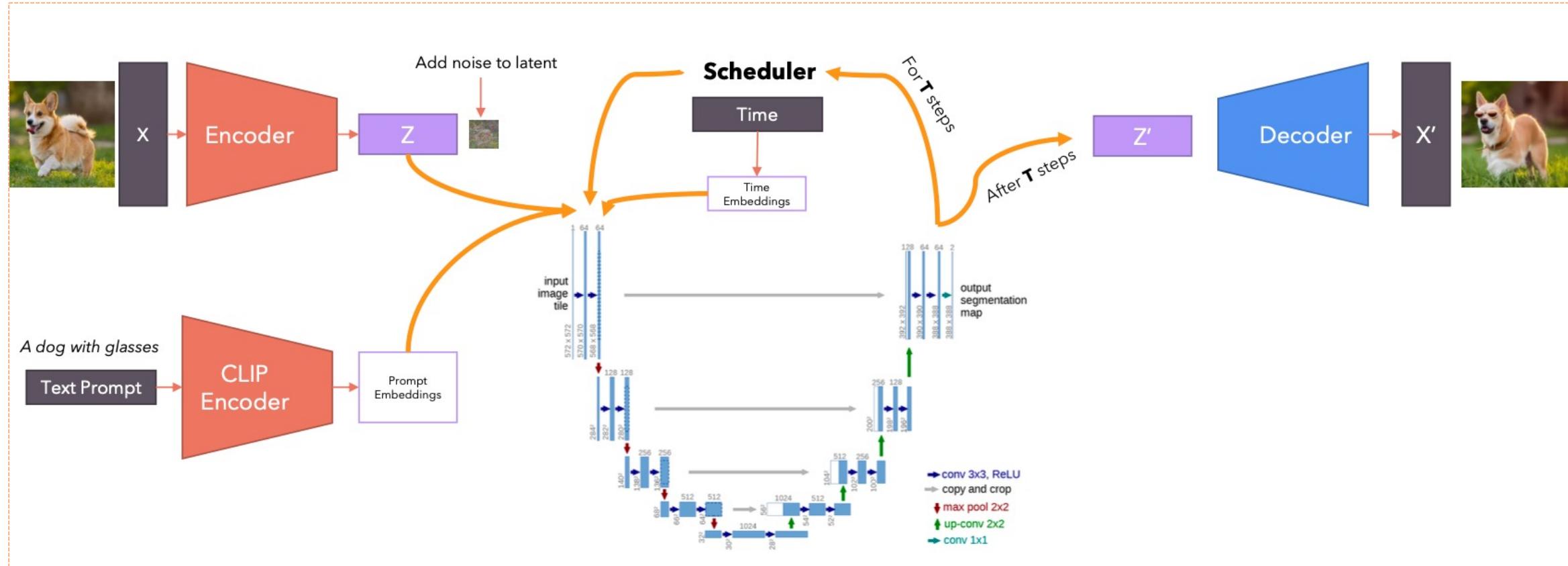
# To the GPU
text_encoder = text_encoder.to(torch_device)

prompt = 'An astronaut riding a horse'
# Turn the text into a sequence of tokens:
text_input = tokenizer(prompt, padding="max_length", max_length=tokenizer.model_max_length,
input_ids = text_input.input_ids.to(torch_device)
# Get output embeddings from tokens
output_embeddings = text_encoder(text_input.input_ids.to(torch_device))[0]
print('Shape:', output_embeddings.shape)
```

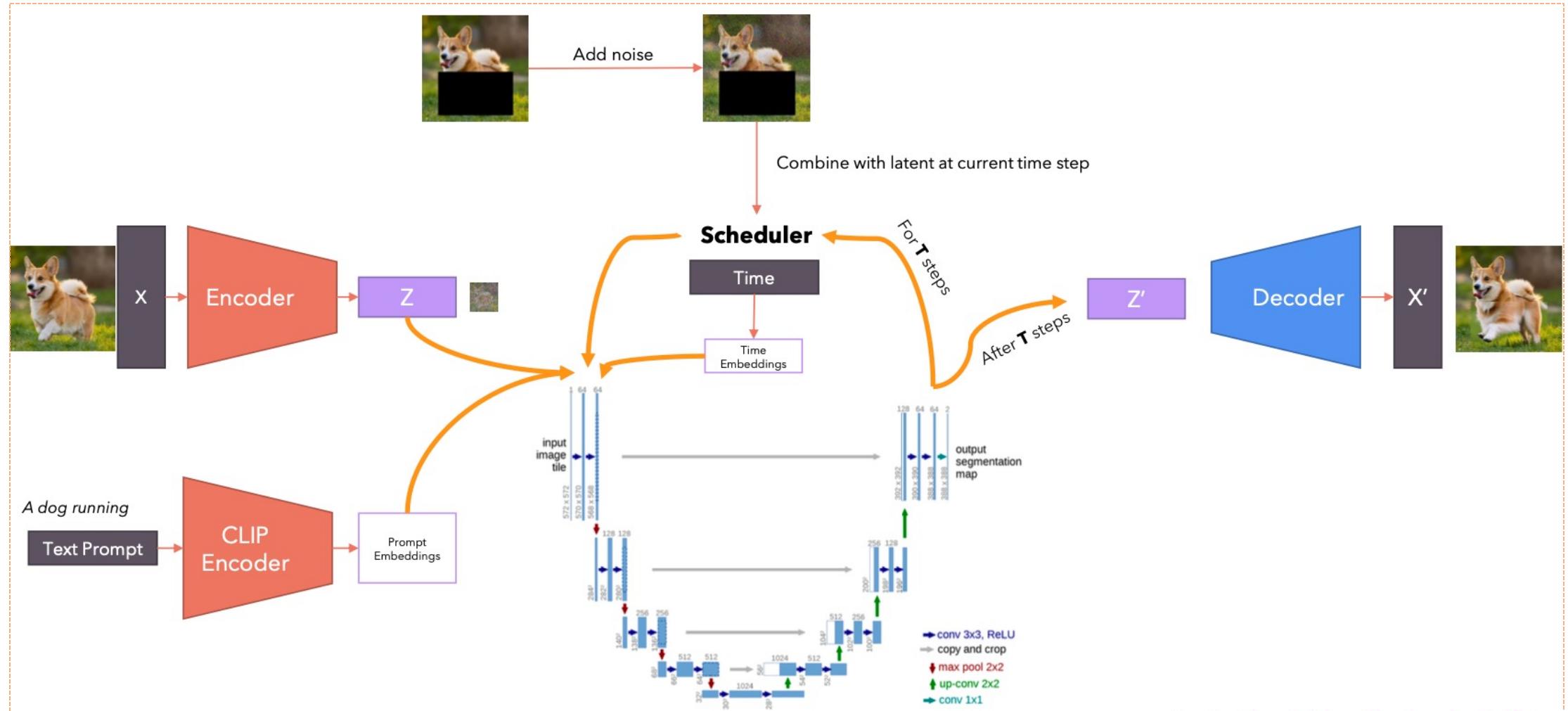
Achitecture (Text-To-Image)



Achitecture (Image-To-Image)



Achitecture (In-painting)



Outline

- **Objective**
- **Introduction to Stable Diffusion**
- **Stable Diffusion: Motivation**
- **Stable Diffusion: Clearly Explained**
- **Stable Diffusion: Demo with API**
- **Summary**

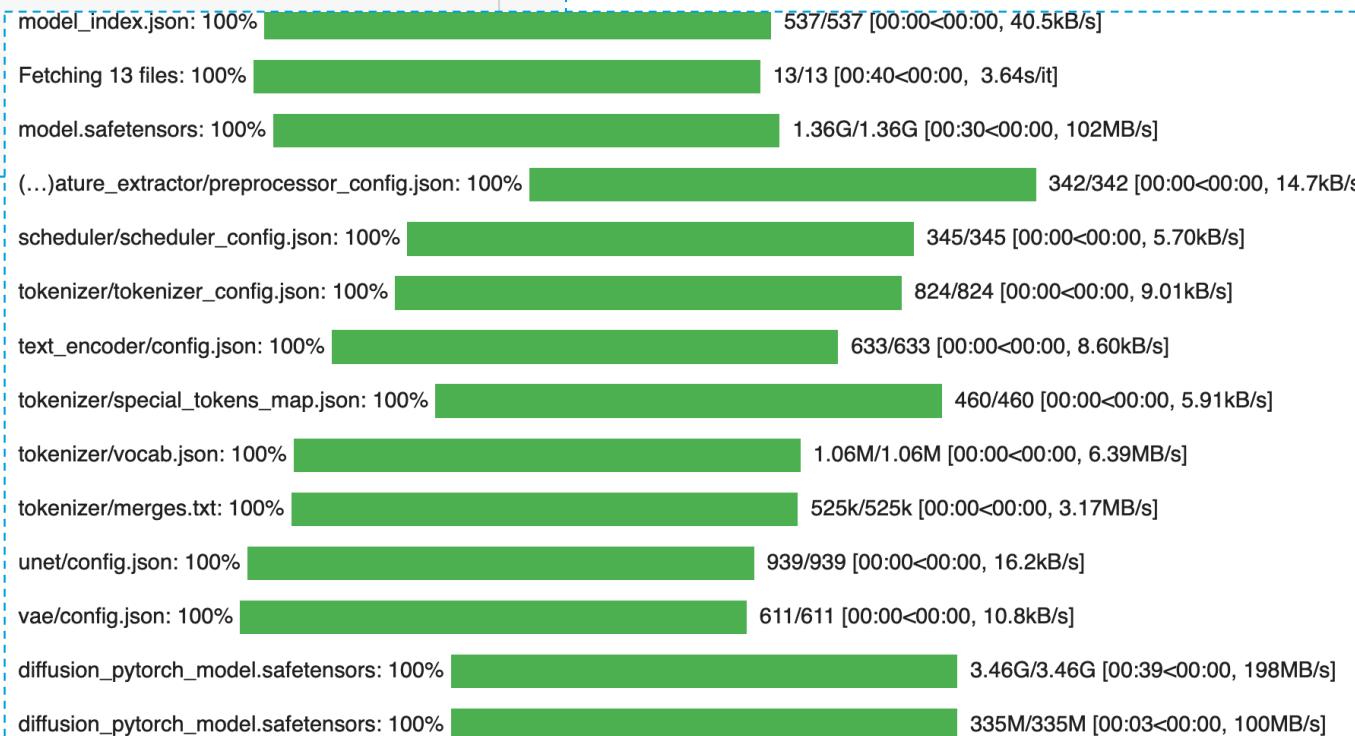
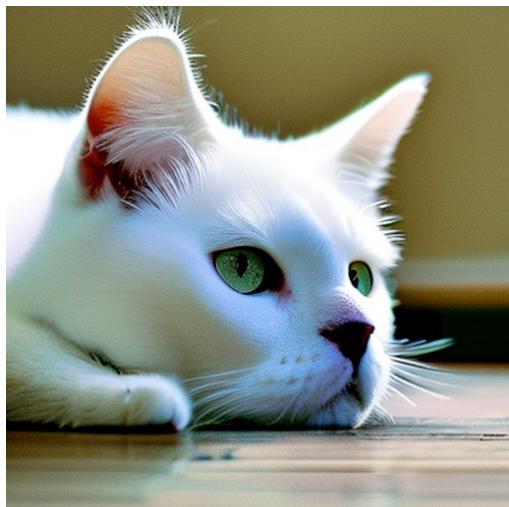
Stable Diffusion with API

```
from diffusers import StableDiffusionPipeline, DPMSolverMultistepScheduler

model_id = "stabilityai/stable-diffusion-2-1"

# Use the DPMSolverMultistepScheduler (DPM-Solver++) scheduler here instead
pipe = StableDiffusionPipeline.from_pretrained(model_id, torch_dtype=torch.float16)
pipe.scheduler = DPMSolverMultistepScheduler.from_config(pipe.scheduler.config)
pipe = pipe.to("cuda")

prompt = "a photo of a white cat"
image = pipe(prompt, height=512, width=512).images[0]
image
```



Stable Diffusion with API

```
▶ prompt = "a lovely cat running in the desert in Van Gogh style, trending art."  
image = pipe(prompt).images[0] # image here is in [PIL format](https://pillow.readthedocs.io/en/stable/handbook/image-file-formats.html#pil-image-objects)  
  
# Now to display an image you can do either save it such as:  
image.save(f"lovely_cat.png")  
image
```



100%

50/50 [00:17<00:00, 2.72it/s]



Stable Diffusion with API

```
▶ generator = torch.Generator("cuda").manual_seed(1024)

prompt = "a sleeping cat enjoying the sunshine."
image = pipe(prompt, generator=generator).images[0] # image here is in [P

# Now to display an image you can do either save it such as:
image.save(f"lovely_cat_sun.png")
image

→ 100% [50/50 [00:17<00:00, 2.69it/s]
```



Stable Diffusion with API

```
▶ prompt = "a sleeping cat enjoying the sunshine."  
image = pipe(prompt, generator=generator,  
            negative_prompt="tree and leaves").images[0] #  
  
# Now to display an image you can do either save it such as:  
image.save(f"lovely_cat_sun_no_trees.png")  
image  
→ 100% 50/50 [00:15<00:00, 3.14it/s]
```

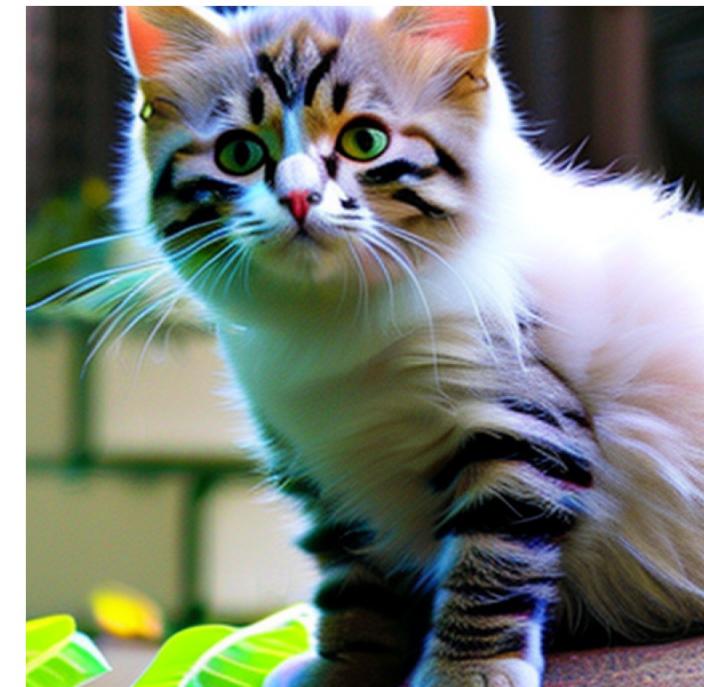


Write a simple text2img sampling function

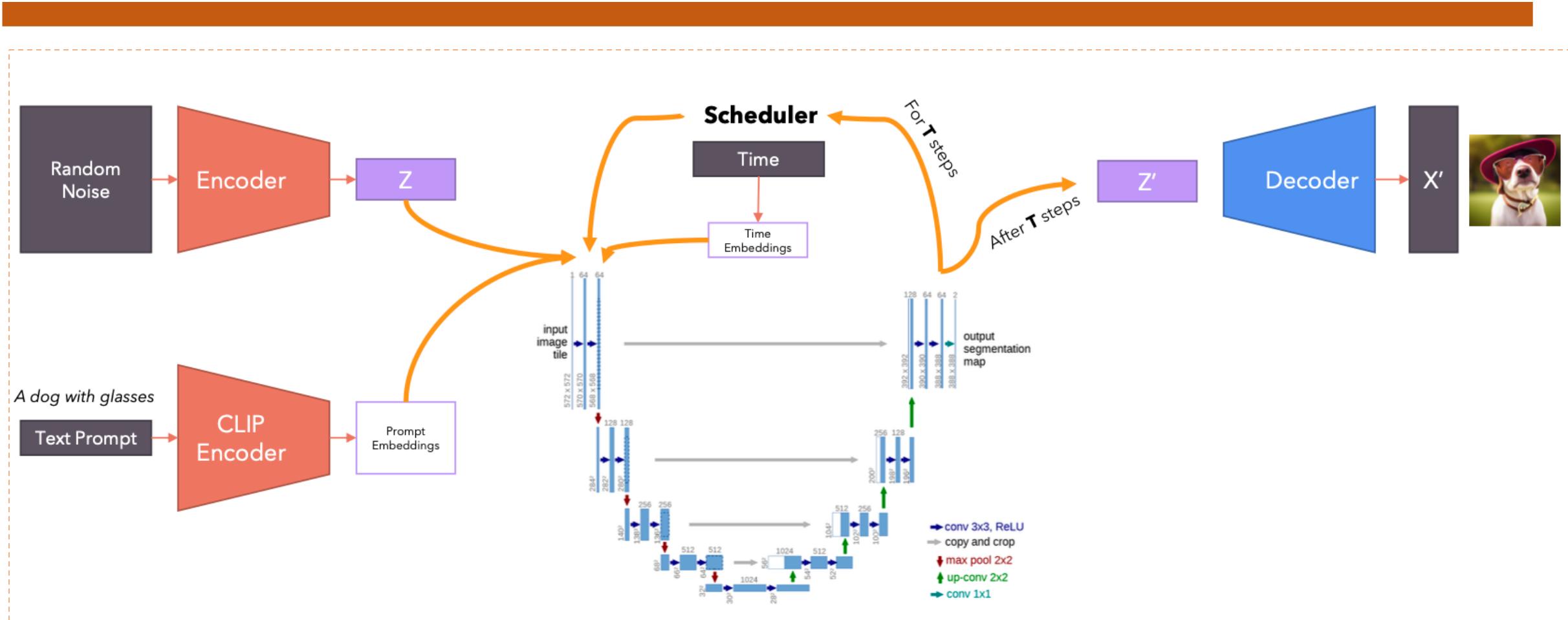
```
[23] @torch.no_grad()
def generate_simplified(
    prompt = ["a lovely cat"],
    negative_prompt = [""],
    num_inference_steps = 50,
    guidance_scale = 7.5):
    # do_classifier_free_guidance
    batch_size = 1
    height, width = 512, 512
    generator = None
    # here `guidance_scale` is defined analog to the guidance weight `w` of equation (2)
    # of the Imagen paper: https://arxiv.org/pdf/2205.11487.pdf . `guidance_scale = 1`
    # corresponds to doing no classifier free guidance.

    # get prompt text embeddings
    text_inputs = pipe.tokenizer(
        prompt,
        padding="max_length",
        max_length=pipe.tokenizer.model_max_length,
        return_tensors="pt",
    )
```

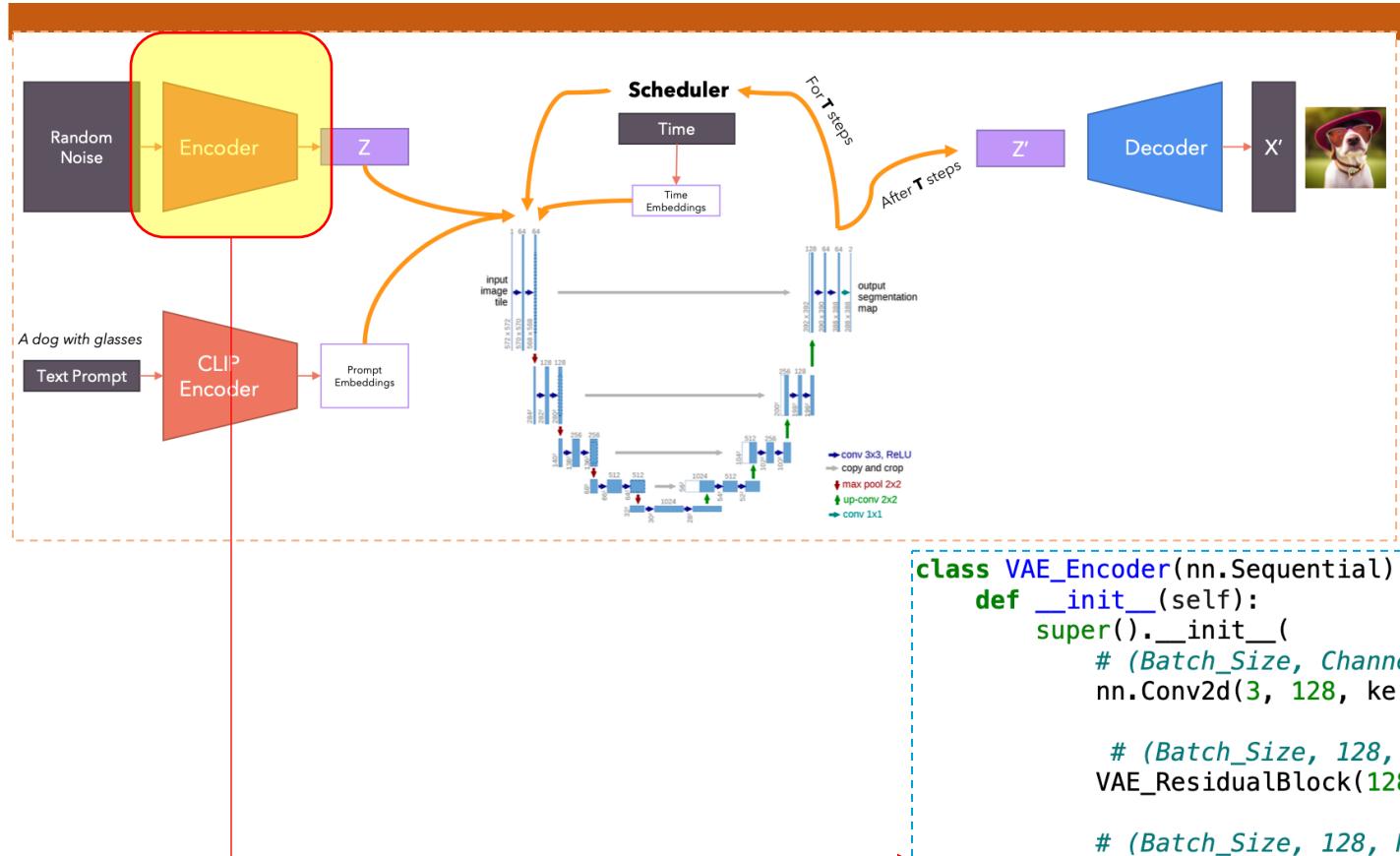
```
image = generate_simplified(
    prompt = ["a lovely cat"],
    negative_prompt = ["Sunshine"], )
plt_show_image(image[0])
```



Stable Diffusion: Implementation



Encoder Implementation



Input: 512x512
Output: 64x64

```
class VAE_Encoder(nn.Sequential):
    def __init__(self):
        super().__init__(
            # (Batch_Size, Channel, Height, Width) -> (Batch_Size, 128, Height, Width)
            nn.Conv2d(3, 128, kernel_size=3, padding=1),

            # (Batch_Size, 128, Height, Width) -> (Batch_Size, 128, Height, Width)
            VAE_ResidualBlock(128, 128),

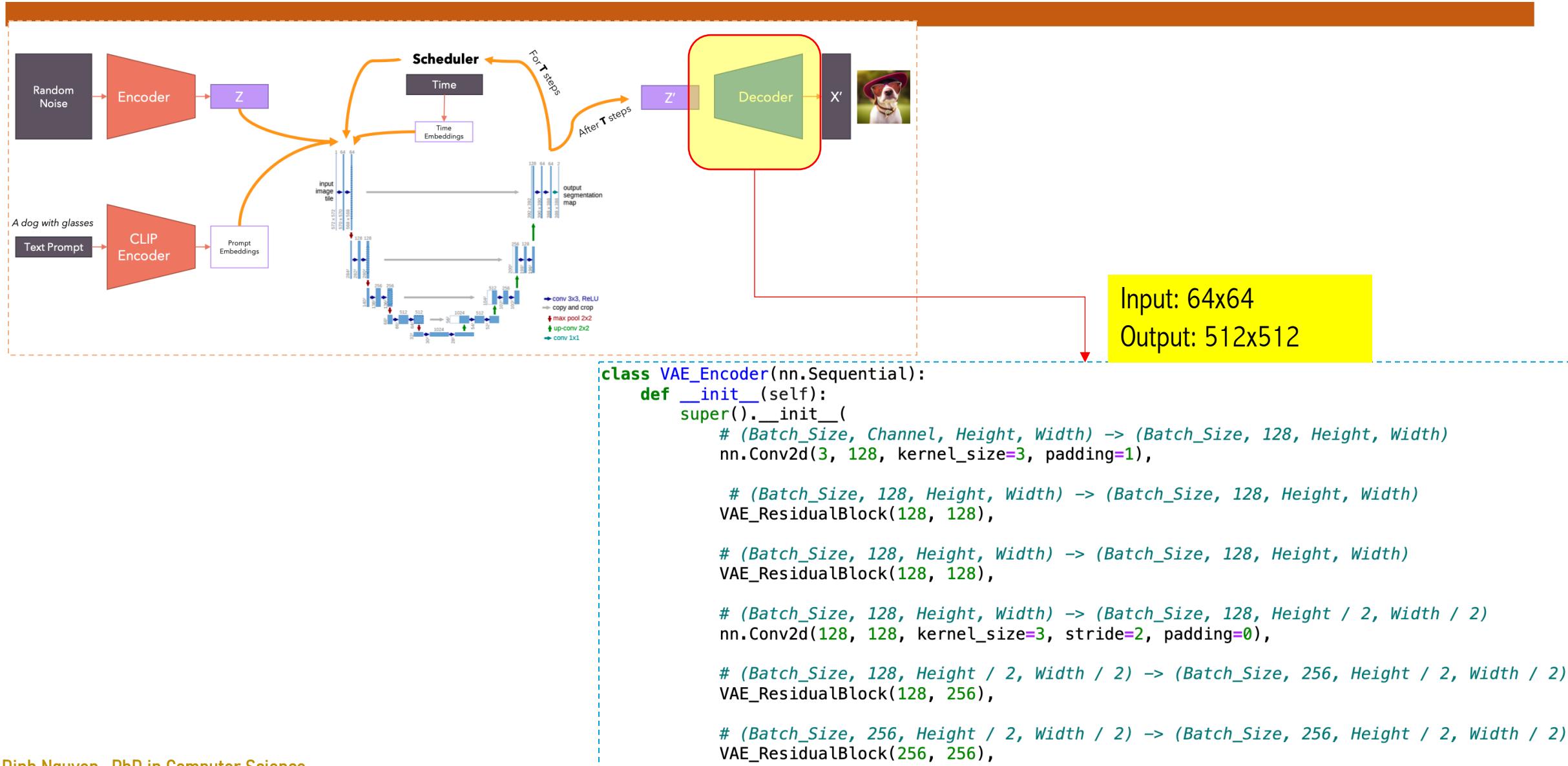
            # (Batch_Size, 128, Height, Width) -> (Batch_Size, 128, Height, Width)
            VAE_ResidualBlock(128, 128),

            # (Batch_Size, 128, Height, Width) -> (Batch_Size, 128, Height / 2, Width / 2)
            nn.Conv2d(128, 128, kernel_size=3, stride=2, padding=0),

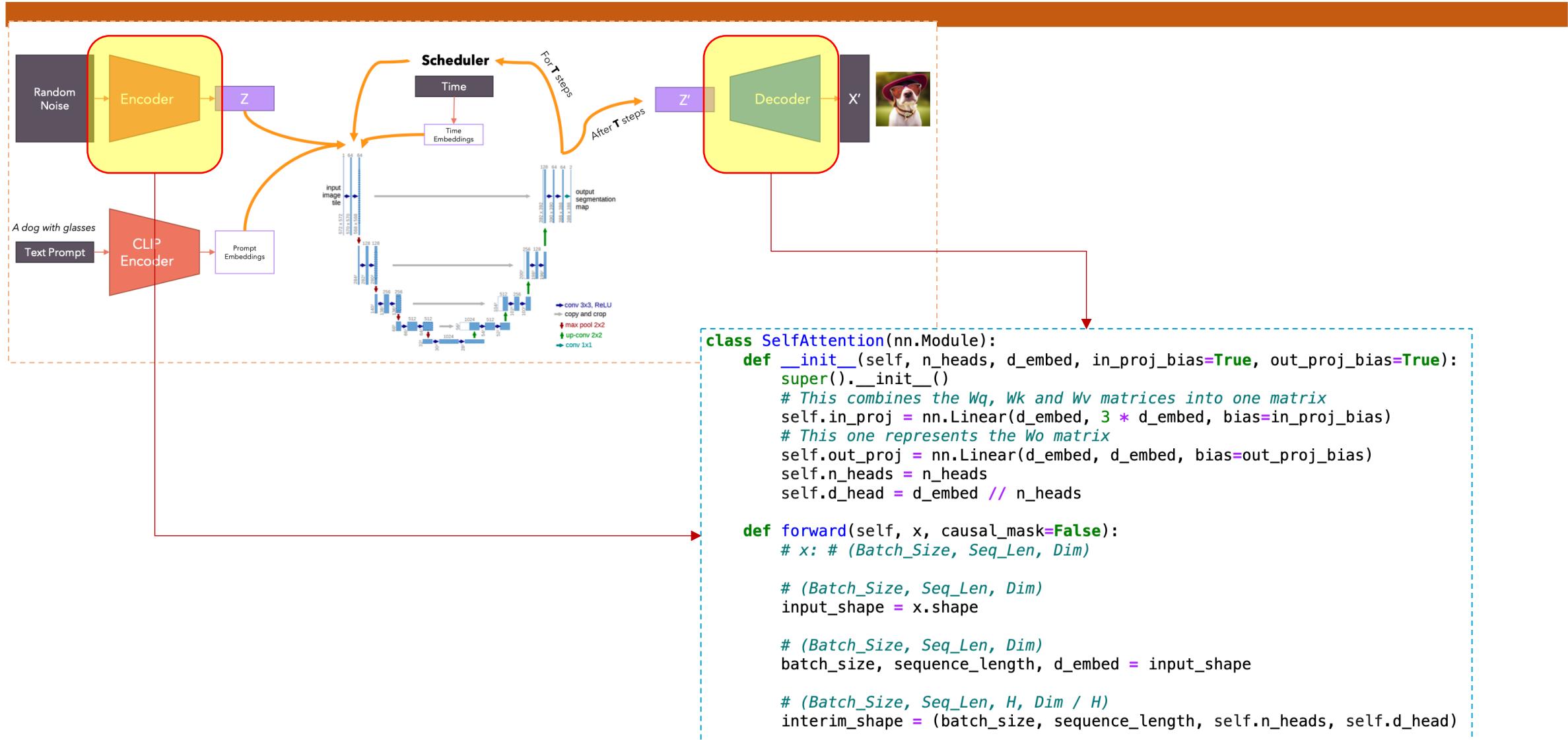
            # (Batch_Size, 128, Height / 2, Width / 2) -> (Batch_Size, 256, Height / 2, Width / 2)
            VAE_ResidualBlock(128, 256),

            # (Batch_Size, 256, Height / 2, Width / 2) -> (Batch_Size, 256, Height / 2, Width / 2)
            VAE_ResidualBlock(256, 256),
```

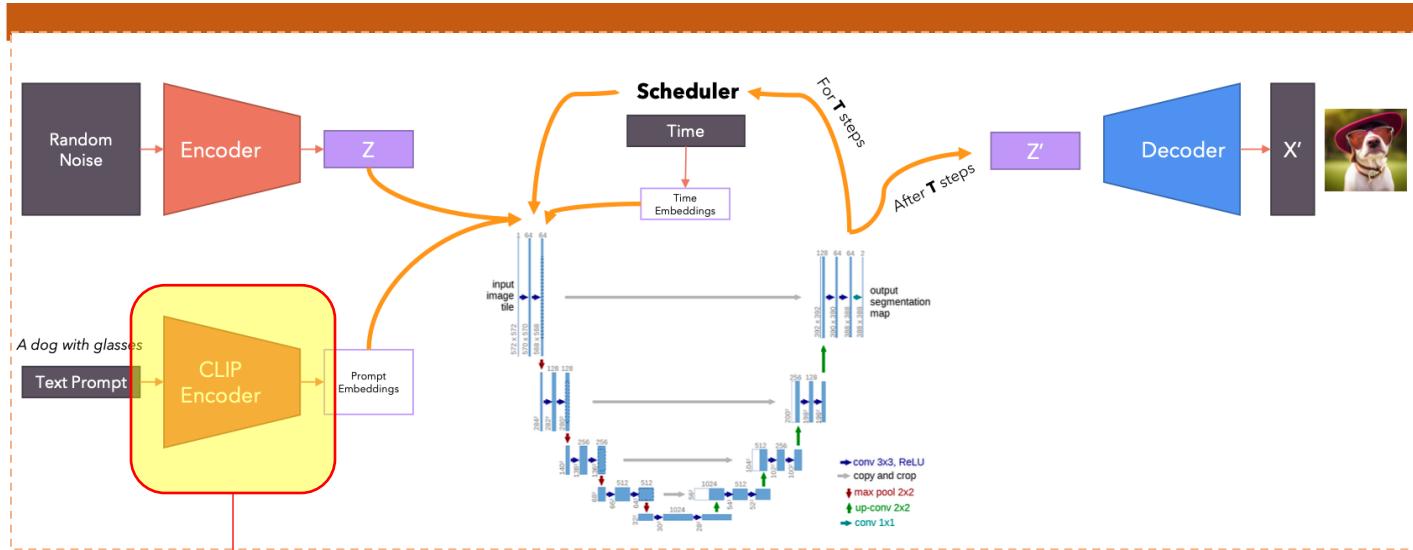
Encoder Implementation



Self-attention Implementation



CLIP Implementation



```

class CLIPEmbedding(nn.Module):
    def __init__(self, n_vocab: int, n_embd: int, n_token: int):
        super().__init__()

        self.token_embedding = nn.Embedding(n_vocab, n_embd)
        # A learnable weight matrix encodes the position information for each token
        self.position_embedding = nn.Parameter(torch.zeros((n_token, n_embd)))

    def forward(self, tokens):
        # (Batch_Size, Seq_Len) -> (Batch_Size, Seq_Len, Dim)
        x = self.token_embedding(tokens)
        # (Batch_Size, Seq_Len) -> (Batch_Size, Seq_Len, Dim)
        x += self.position_embedding

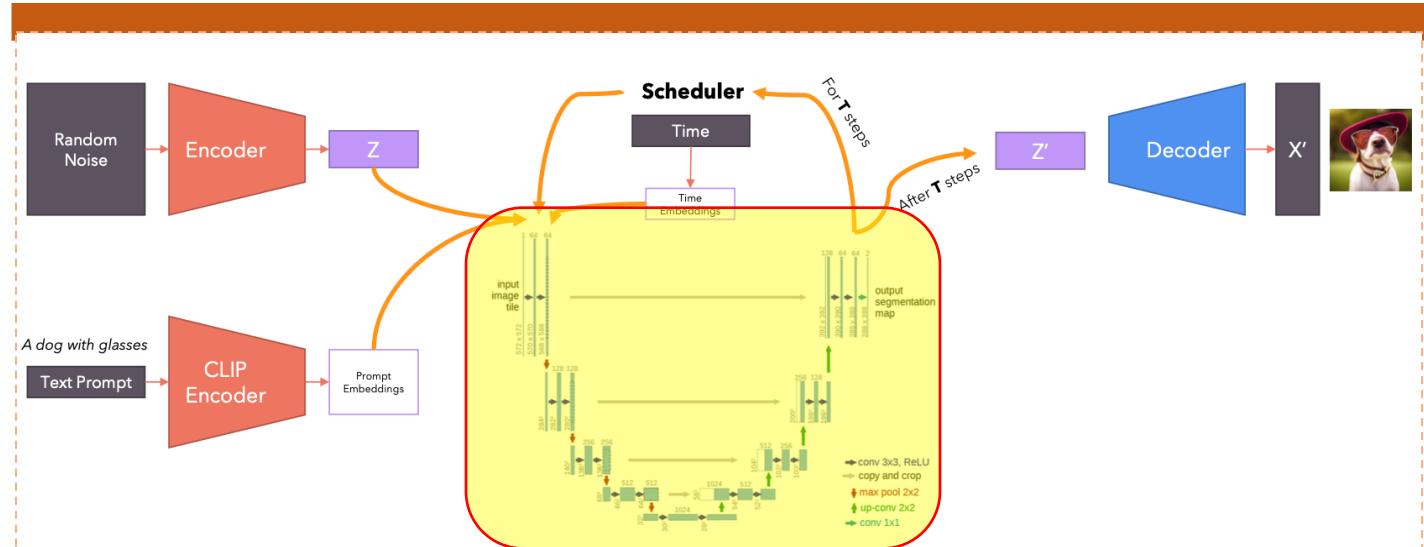
        return x

class CLIPLayer(nn.Module):
    def __init__(self, n_head: int, n_embd: int):
        super().__init__()

        # Pre-attention norm
        self.layernorm_1 = nn.LayerNorm(n_embd)
        # Self attention
        self.attention = SelfAttention(n_head, n_embd)
        # Pre-FNN norm
        self.layernorm_2 = nn.LayerNorm(n_embd)
        # Feedforward layer

```

Unet Implementation



```

class UNET_AttentionBlock(nn.Module):
    def __init__(self, n_head: int, n_embd: int, d_context=768):
        super().__init__()
        channels = n_head * n_embd

        self.groupnorm = nn.GroupNorm(32, channels, eps=1e-6)
        self.conv_input = nn.Conv2d(channels, channels, kernel_size=1, padding=0)

        self.layernorm_1 = nn.LayerNorm(channels)
        self.attention_1 = SelfAttention(n_head, channels, in_proj_bias=False)
        self.layernorm_2 = nn.LayerNorm(channels)
        self.attention_2 = CrossAttention(n_head, channels, d_context, in_proj_bias=False)
        self.layernorm_3 = nn.LayerNorm(channels)
        self.linear_gelu_1 = nn.Linear(channels, 4 * channels * 2)
        self.linear_gelu_2 = nn.Linear(4 * channels, channels)

        self.conv_output = nn.Conv2d(channels, channels, kernel_size=1, padding=0)

    def forward(self, x, context):
        # x: (Batch_Size, Features, Height, Width)
        # context: (Batch_Size, Seq_Len, Dim)

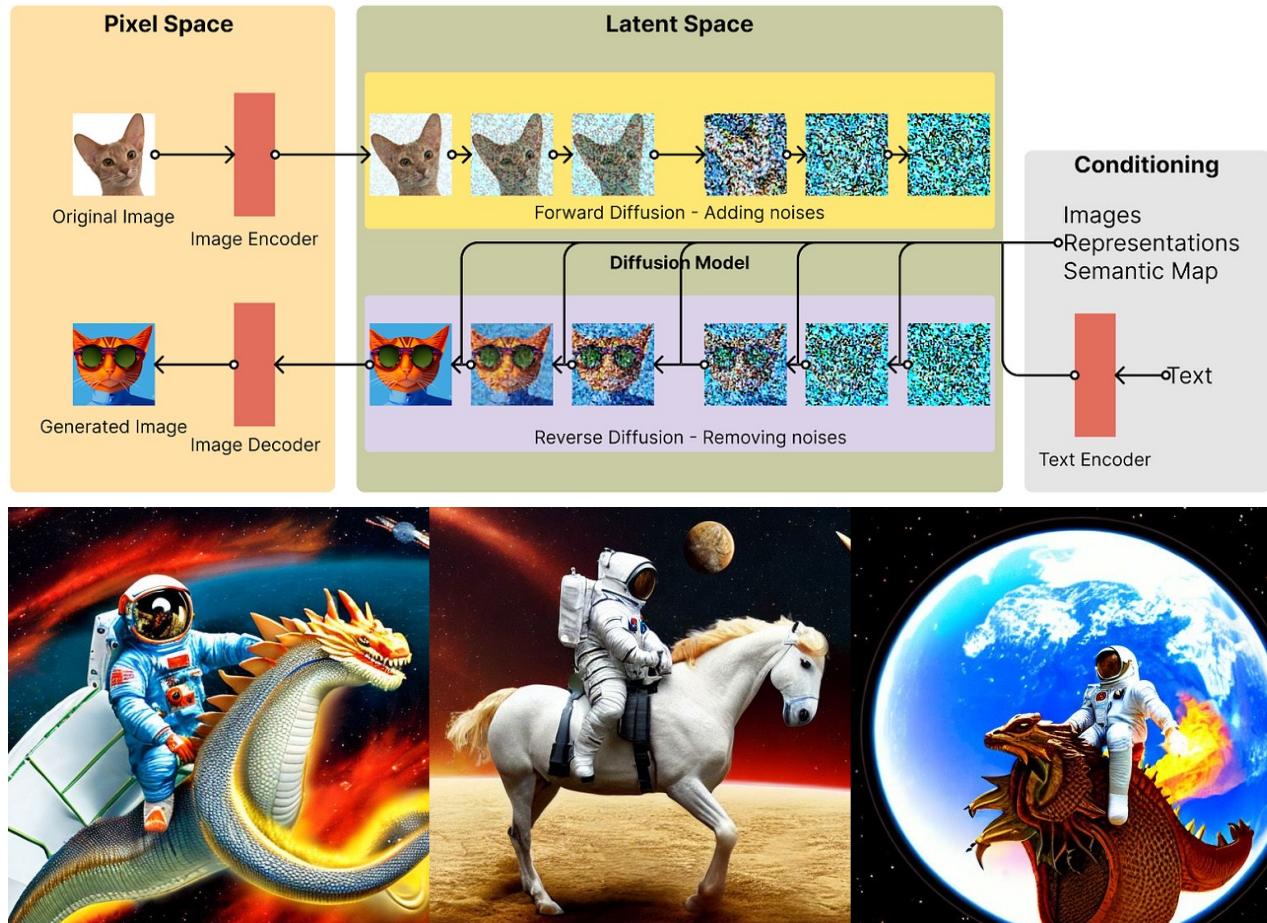
        residue_long = x

```

Outline

- **Objective**
- **Introduction to Stable Diffusion**
- **Stable Diffusion: Motivation**
- **Stable Diffusion: Clearly Explained**
- **Stable Diffusion: Demo with API**
- **Summary**

Summary



- 1 • Principle of Diffusion models.
- 2 • Model score function of images with UNet model
- 3 • Understanding prompt through contextualized word embedding
- 4 • Let text influence image through cross attention
- 5 • Improve efficiency by adding an autoencoder

