

## Module 9 - Exercise

# Text to Image Generation Using Stable Diffusion Model

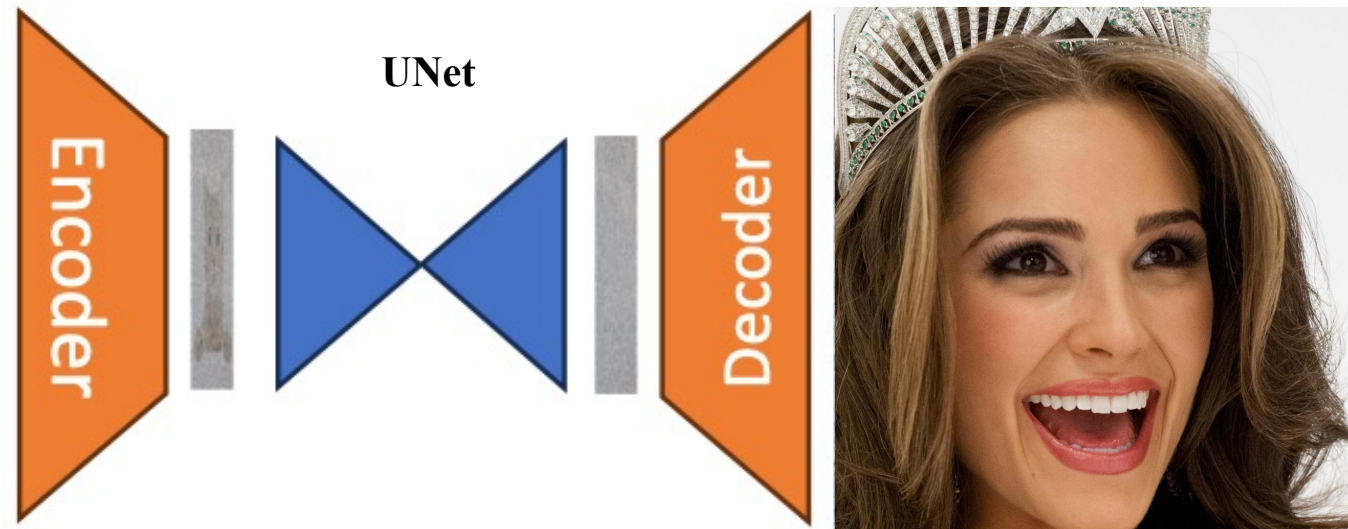
AI VIET NAM  
Nguyen Quoc Thai

# Objectives



## Text-to-Image using Stable Diffusion Model

The person has high cheekbones, and pointy nose. She is wearing lipstick.





# Outline

- **Introduction**
- **Stable Diffusion Model**
- **Text-to-Image Generation using SDM**

# Introduction



## Text to Image Generation

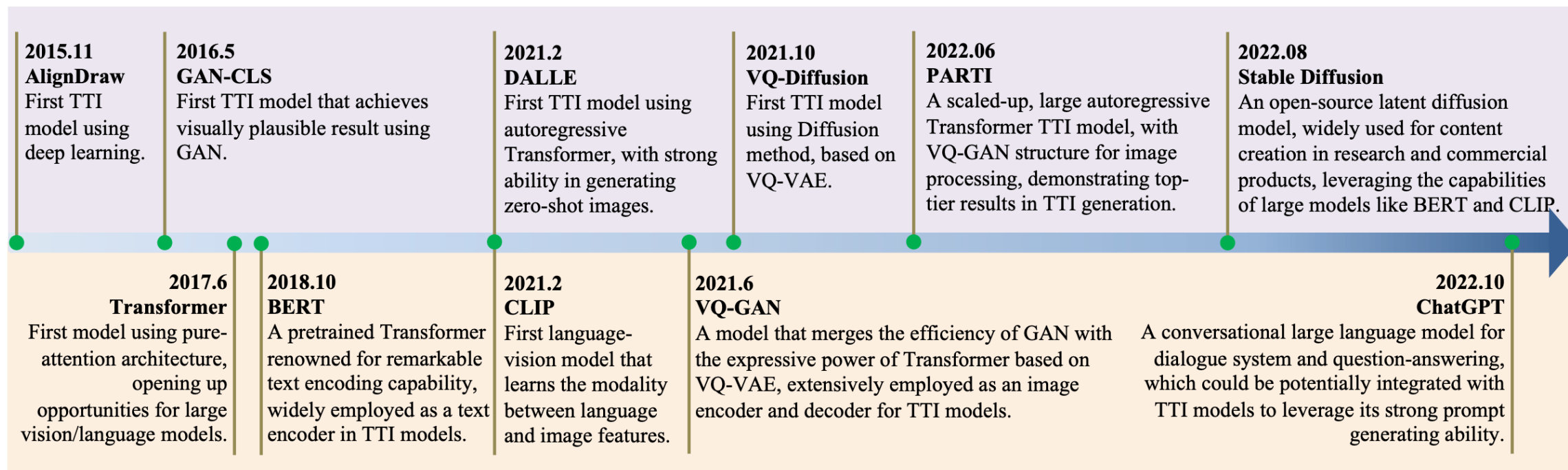
- The person has high cheekbones, and pointy nose. She is wearing lipstick.



# Introduction



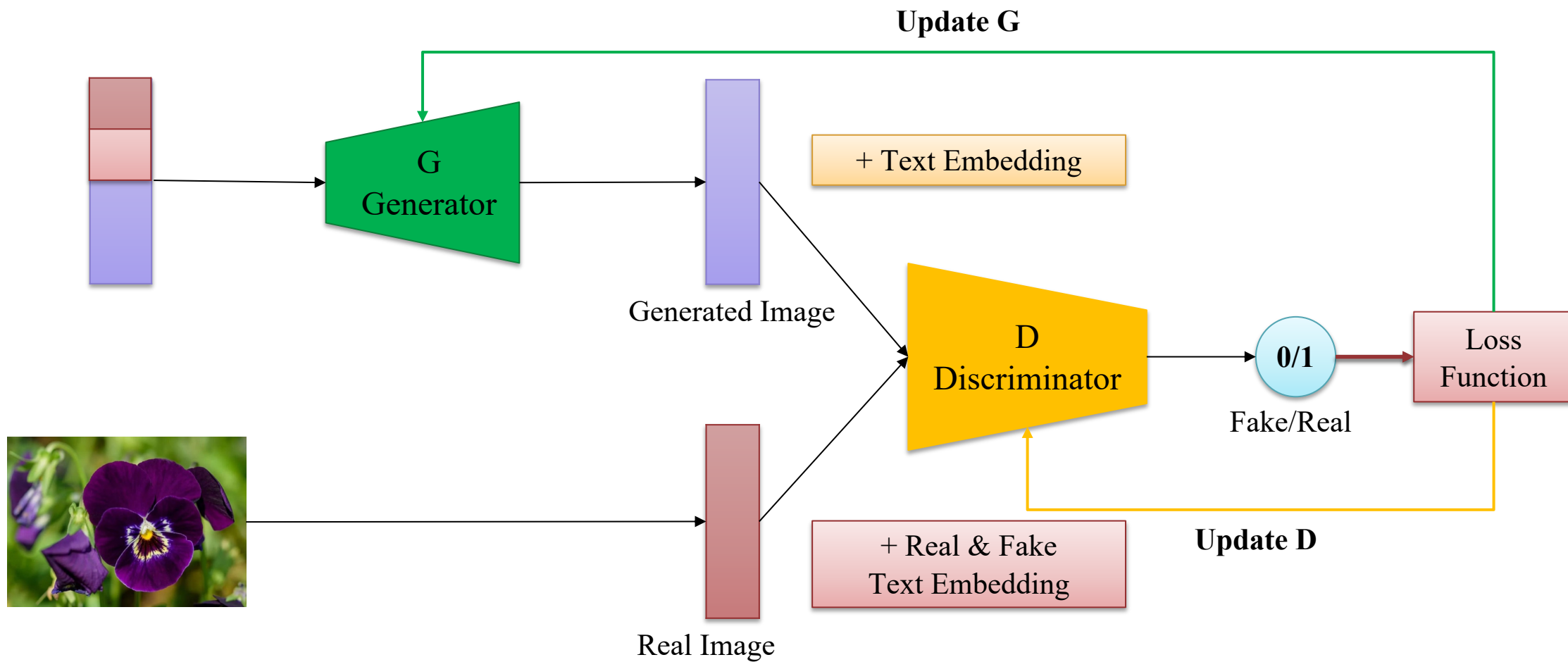
## The milestones of text-to-image models and large models



# Introduction



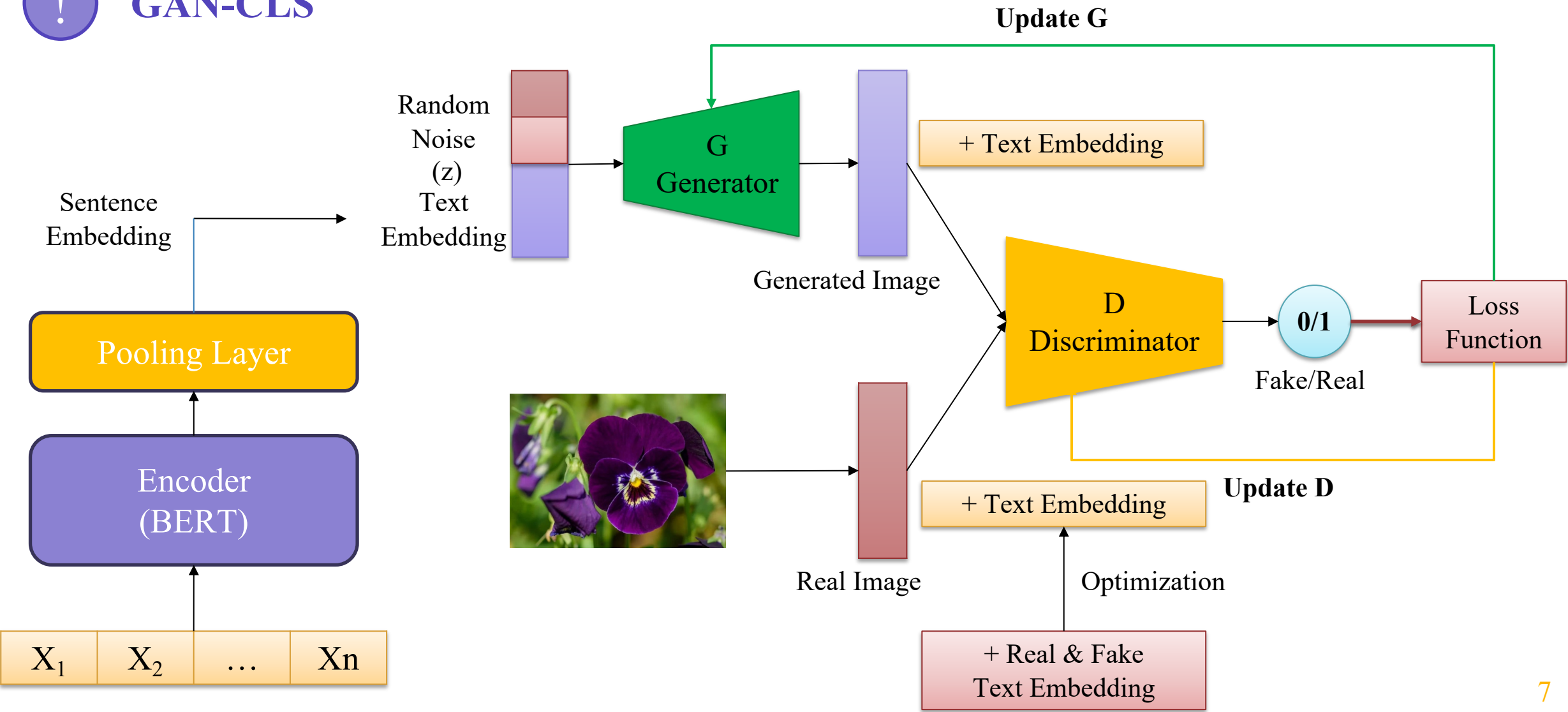
## GAN-CLS



# Introduction



## GAN-CLS

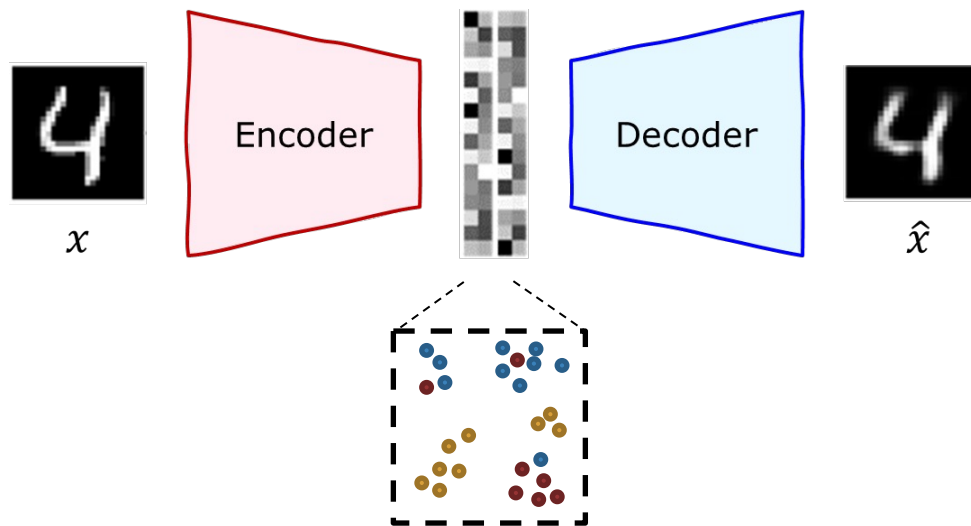


# Introduction

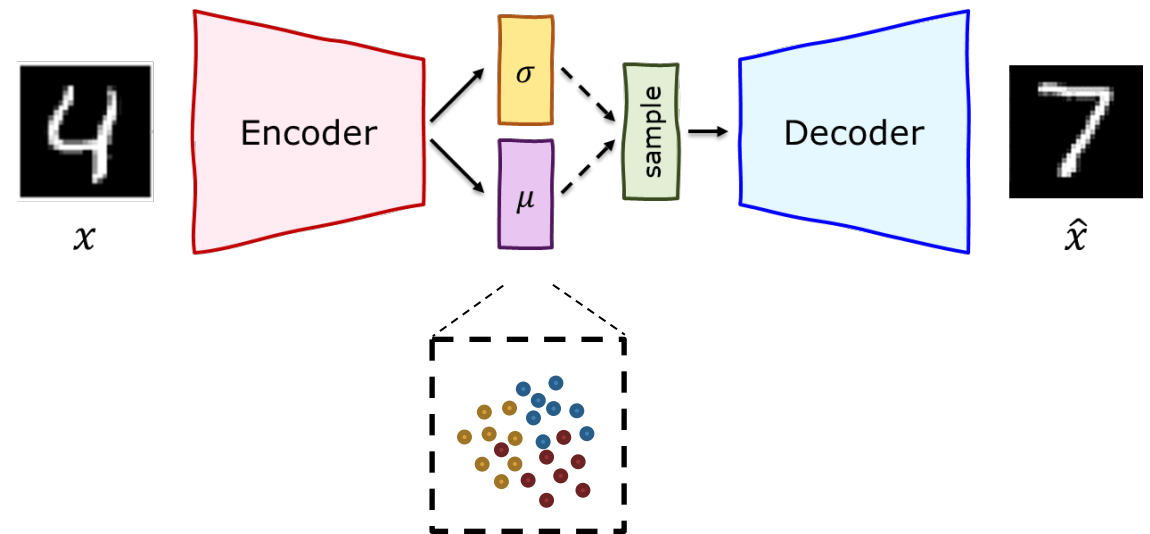


## Variational AutoEncoder

Autoencoder



Variational Autoencoder





# Introduction

## ! Vector Quantized Variational AutoEncoder (VQ-VAE)

- VAE learns a discrete latent representation, instead continuous
- Latents do not necessarily need to be continuous vectors
- It just needs to be some form of numerical representation of the data

**ENCODER**



Image to discrete codes ↓

56	73	67	23	81	19
----	----	----	----	----	----

**DECODER**

56	73	67	23	81	19
----	----	----	----	----	----



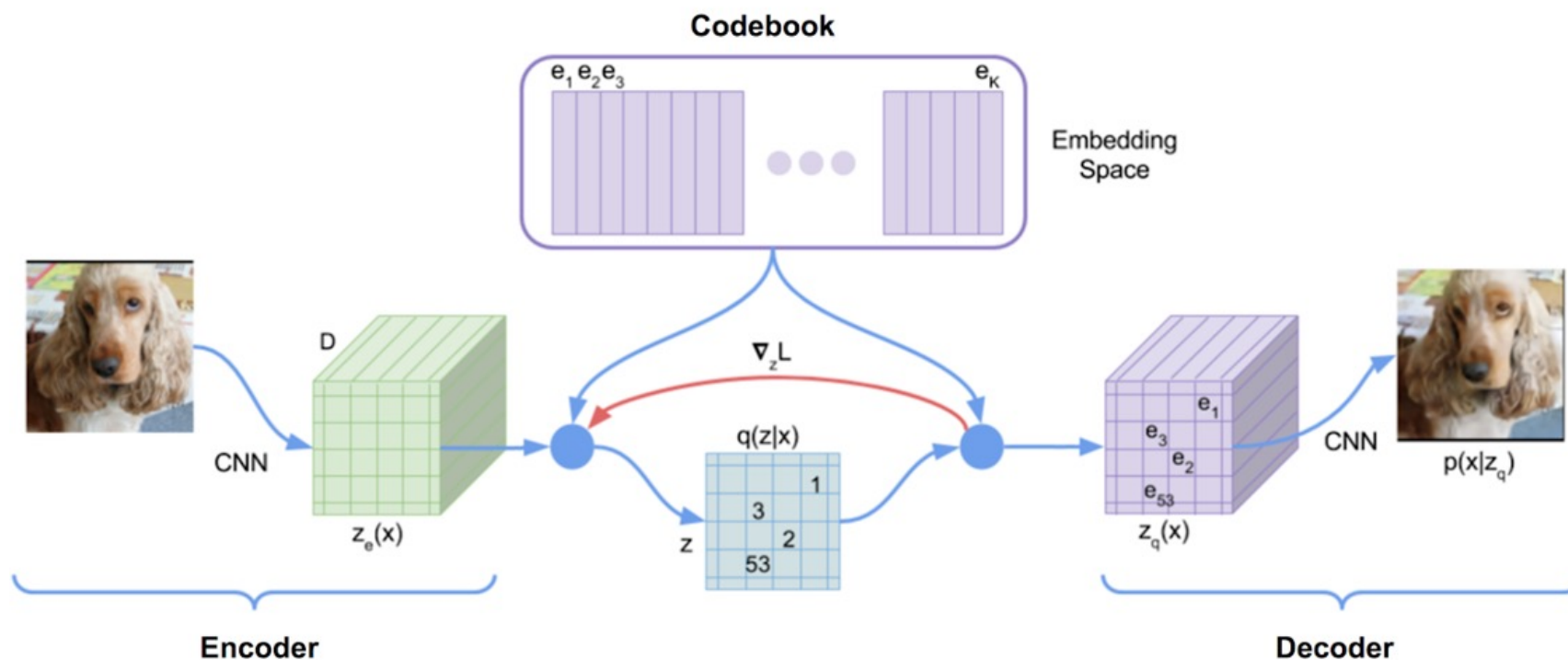
Discrete codes to image



# Introduction

## ! Visual Vocabulary

- Introduces a Discrete Latent Codebook to store a finite set of possible latent vectors
- Describe an image as a sequence of symbols (language tokens)

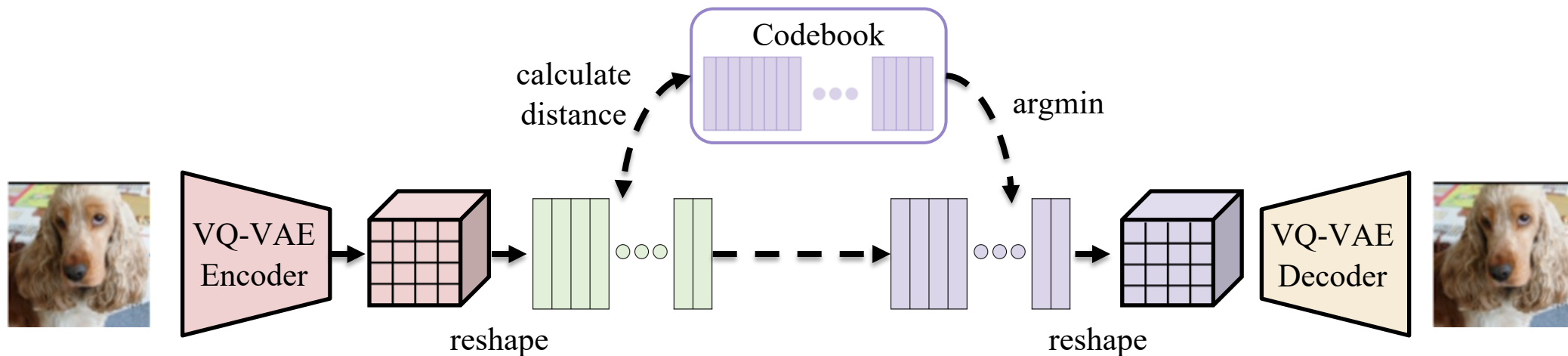


# Introduction

## ! Quantization Layer

### During Training

- Output of encoder is compared to all the vectors in the codebook
- The codebook vector closet in Euclidean distance is fed to decoder

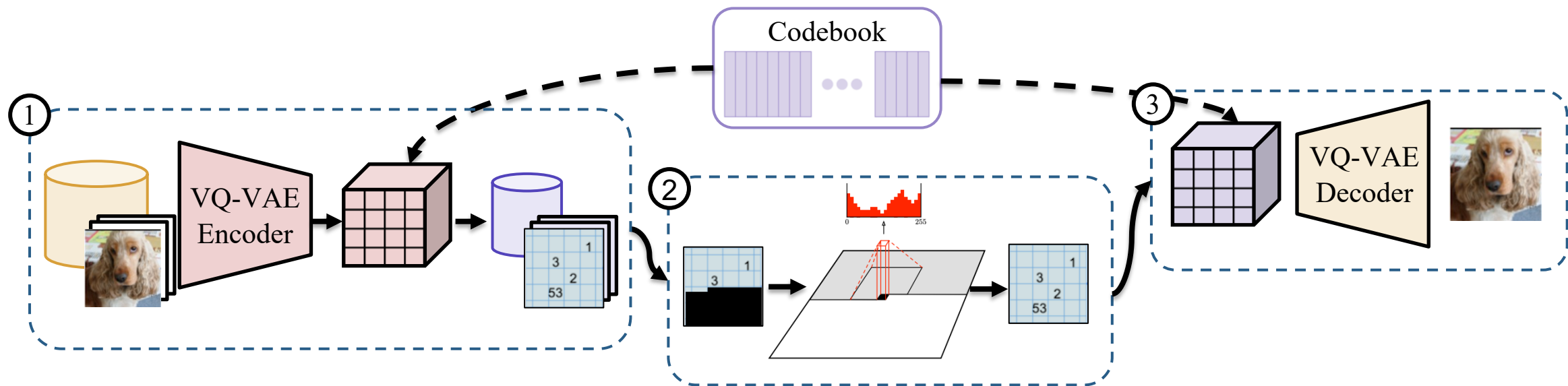


# Introduction

## ! Generate Image from Codebook

Train a PixelCNN as prior on the discretized 32x32 latent space

- Use VQ-VAE Encoder to extract latent space (codebook indicates) from dataset
- Train PixelCNN to auto-regressively complete the latent codebook
- Use VQ-VAE Decoder to generate image from the completed latent codebook



# Introduction



## DALL-E

Text-to-Image Generator model using a transformer that autoregressively models the text and image tokens as a single stream of data

- Uses Discrete VAE
- Switch PixelCNN with a 12-billion parameter GPT-3
- Trained on 250 million image-text pairs

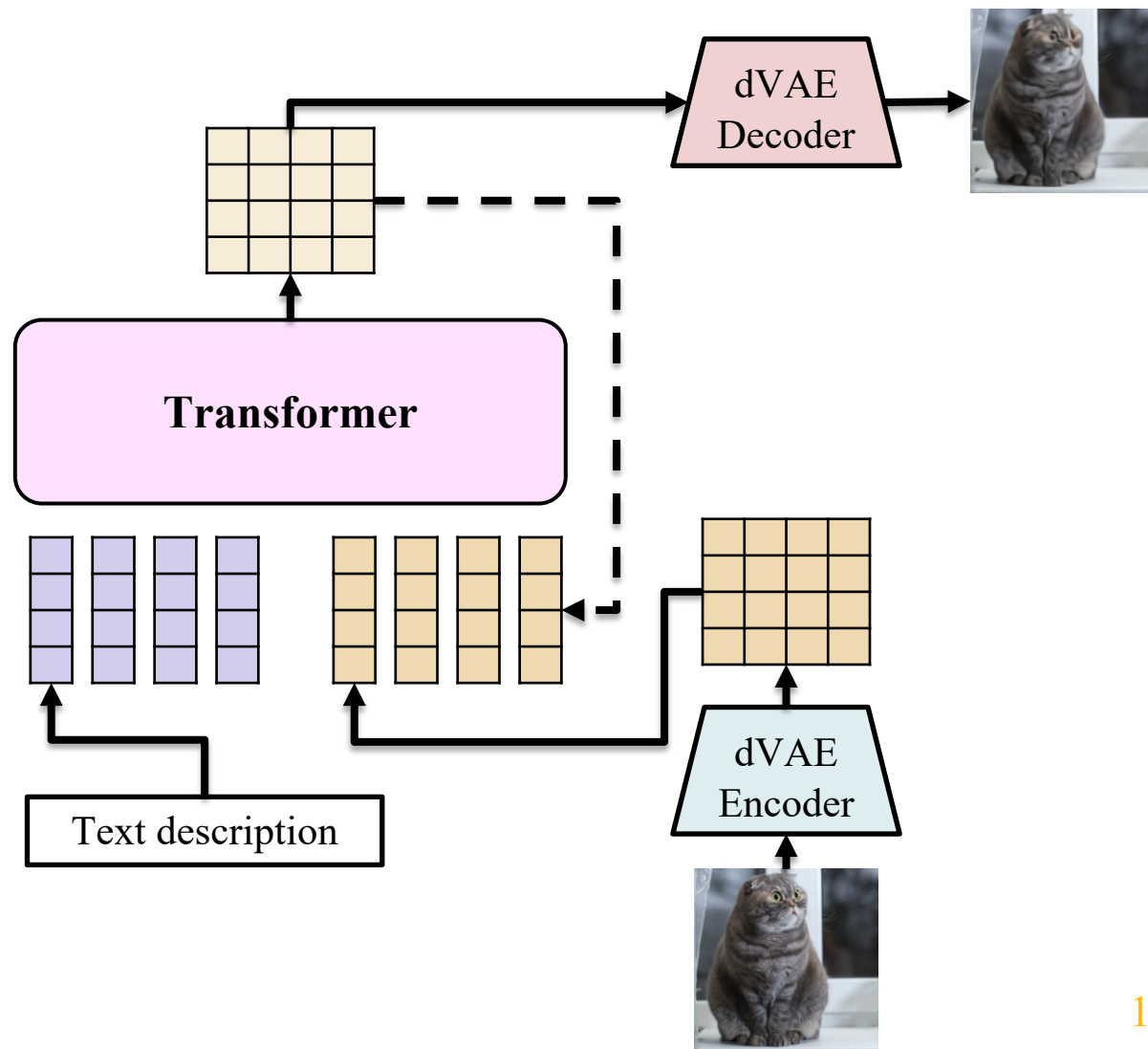


An armchair in the shape of an avocado

# Introduction

## ! DALL-E Parts

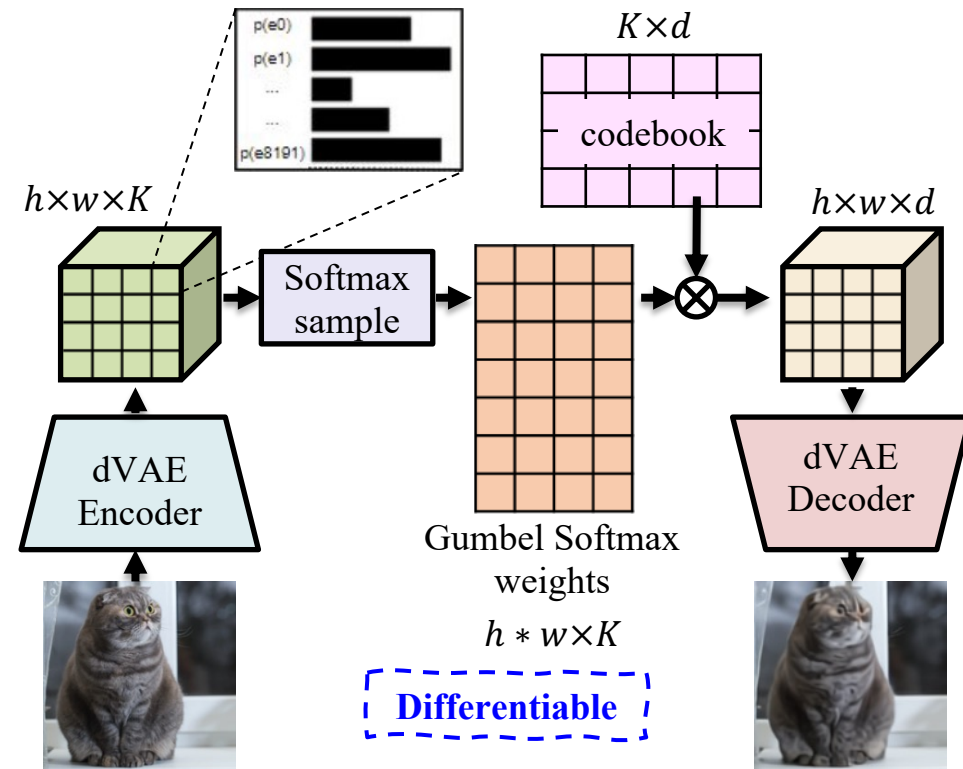
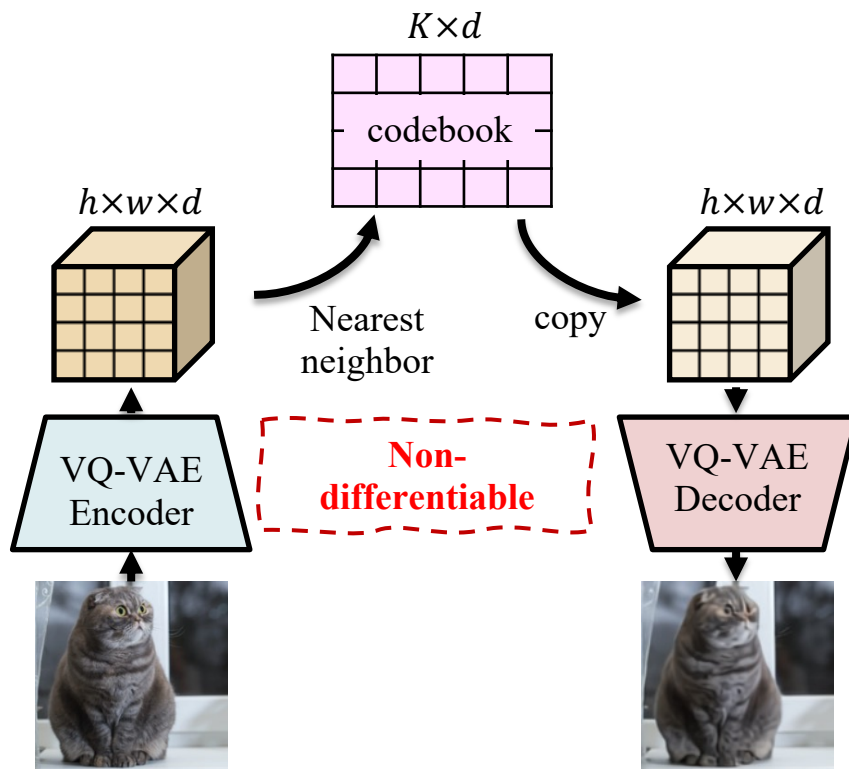
- Discrete VAE encoder and decoder
  - Inspired by VQ-VAE-2
  - Compress 256x256 RGB images into a 32x32 grid of image tokens
  - With 8192 possible codebook tokens
- Transformer Decoder
  - Concatenate text tokens with image tokens into single array
  - Train to predict text image token from the preceding tokens



# Introduction

## ! Gumbel Softmax Relaxation

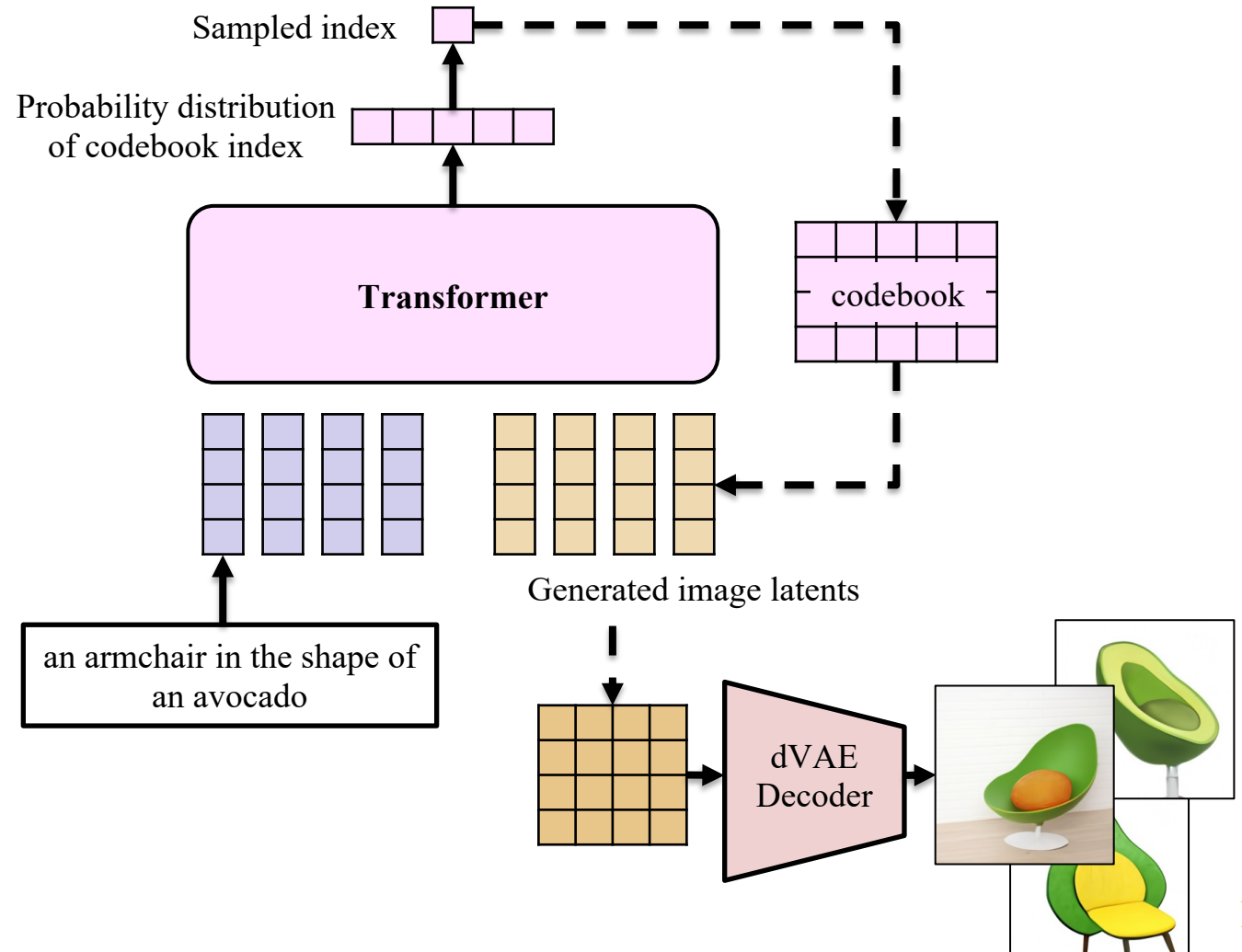
- Outputs a distribution over codebook vectors for each latent code instead of mapping deterministically to a single codebook vector.



# Introduction

## ! Dall-E Inference

- Not directly predict (choose) the next latent index, but predict the distribution
- Then sample the index from that distribution

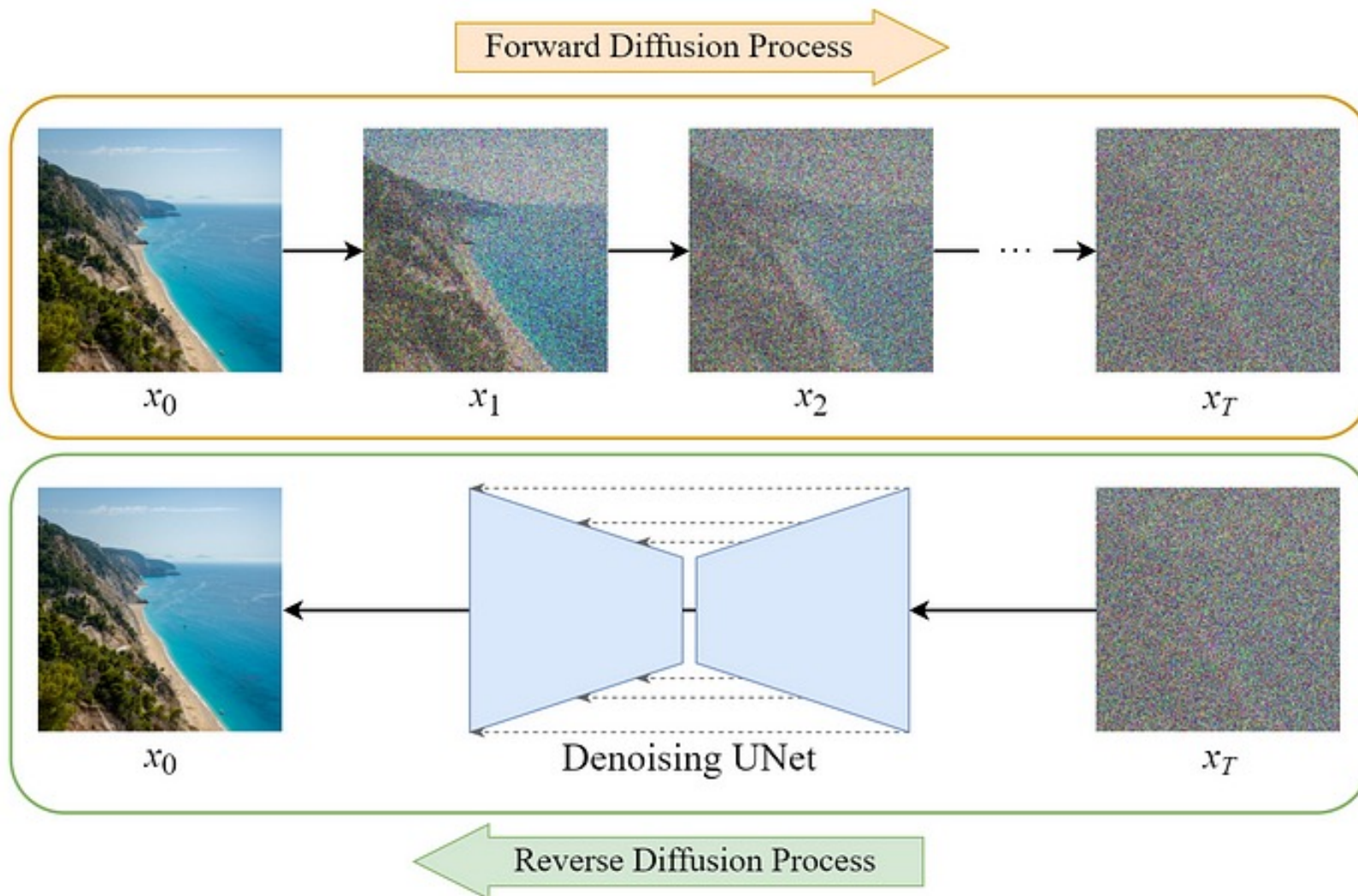




# Introduction



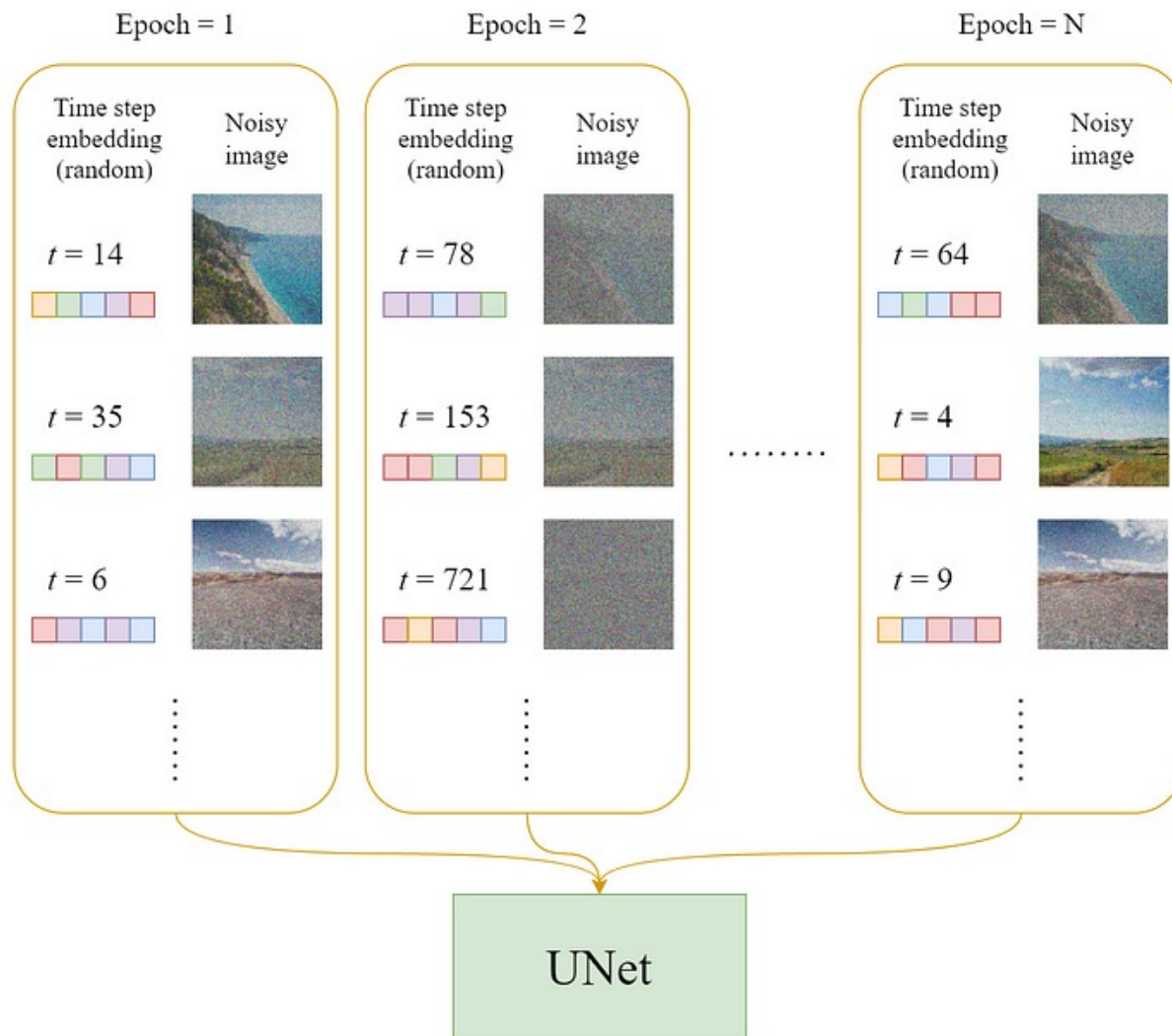
## Diffusion Model



# Introduction



## Diffusion Model - Training

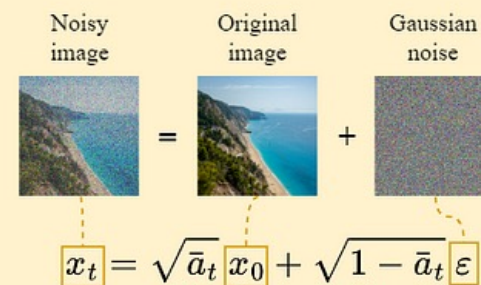


For each training step:

1. Randomly select a time step & encode it



2. Add noise to image



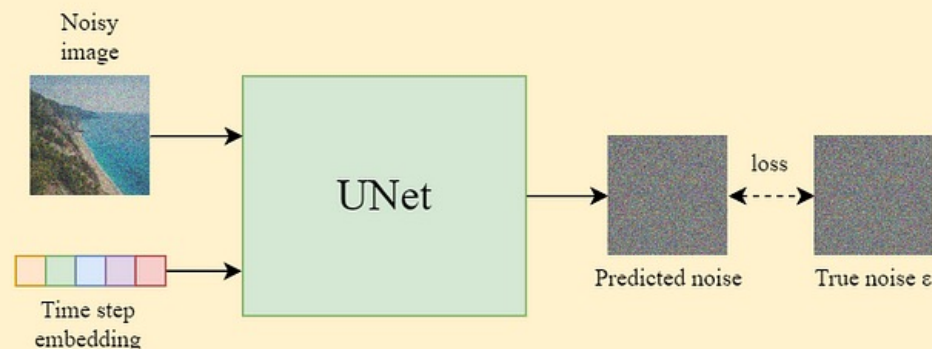
Adjust the amount of noise according to the time step  $t$

$$\epsilon \sim \mathcal{N}(0, 1)$$

$$\alpha_t = 1 - \beta_t$$

$$\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$$

3. Train the UNet



# Introduction



## Diffusion Model - Inference

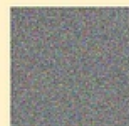
Reverse Diffusion / Denoising / Sampling

1. Sample a Gaussian noise

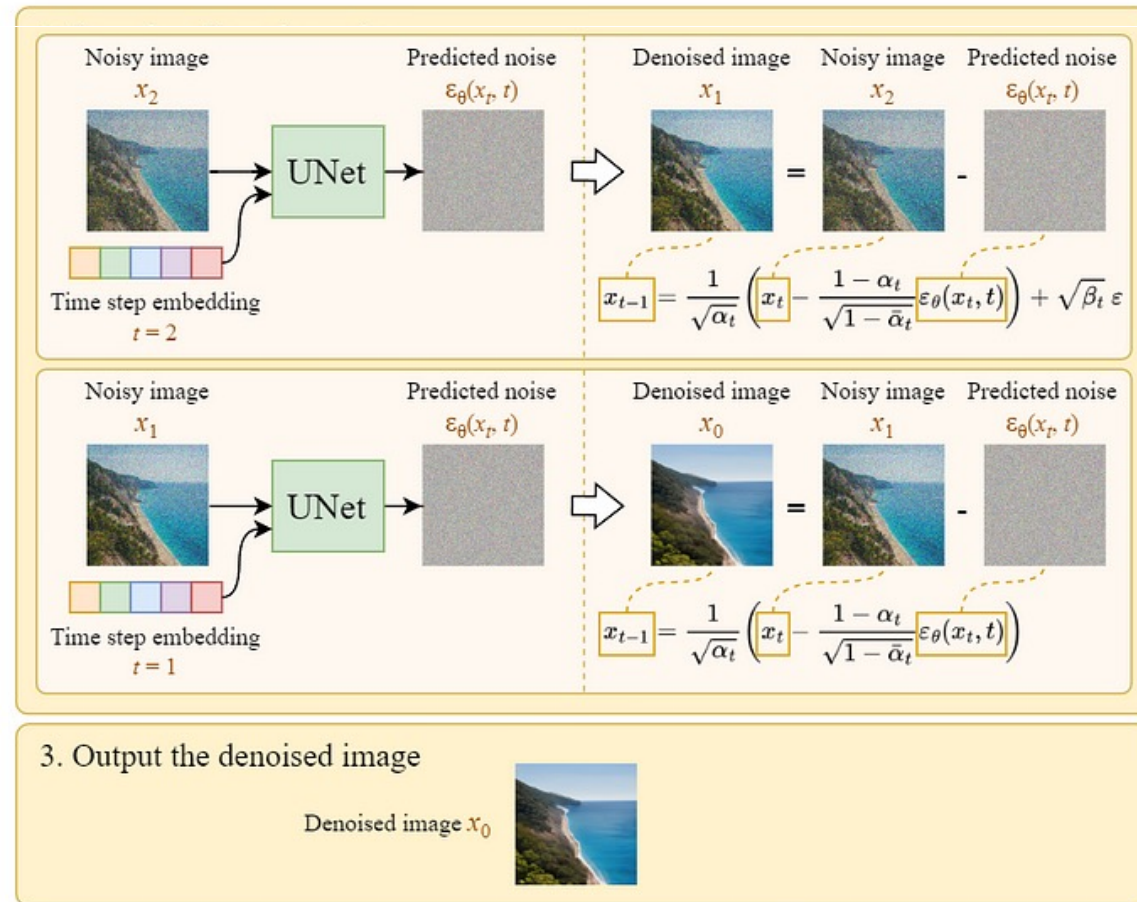
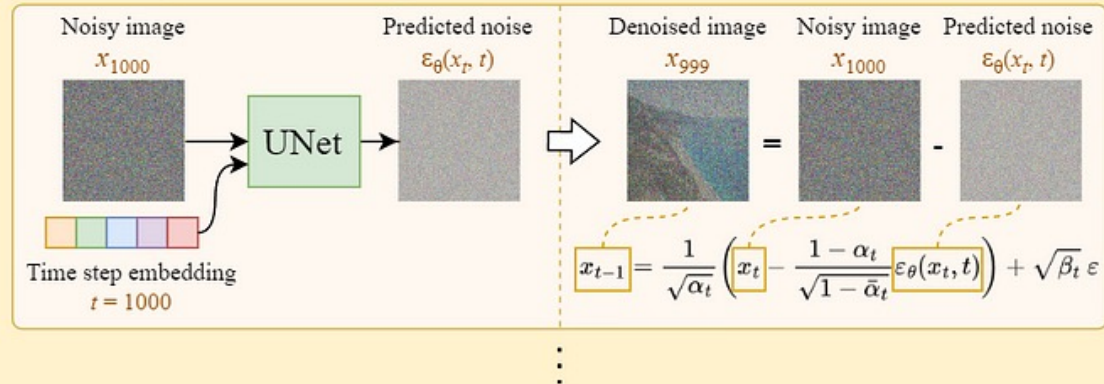
$$x_T \sim N(0, I)$$

E.g.  $T = 1000$

$$x_{1000} \sim N(0, I)$$



2. Iteratively denoise the image

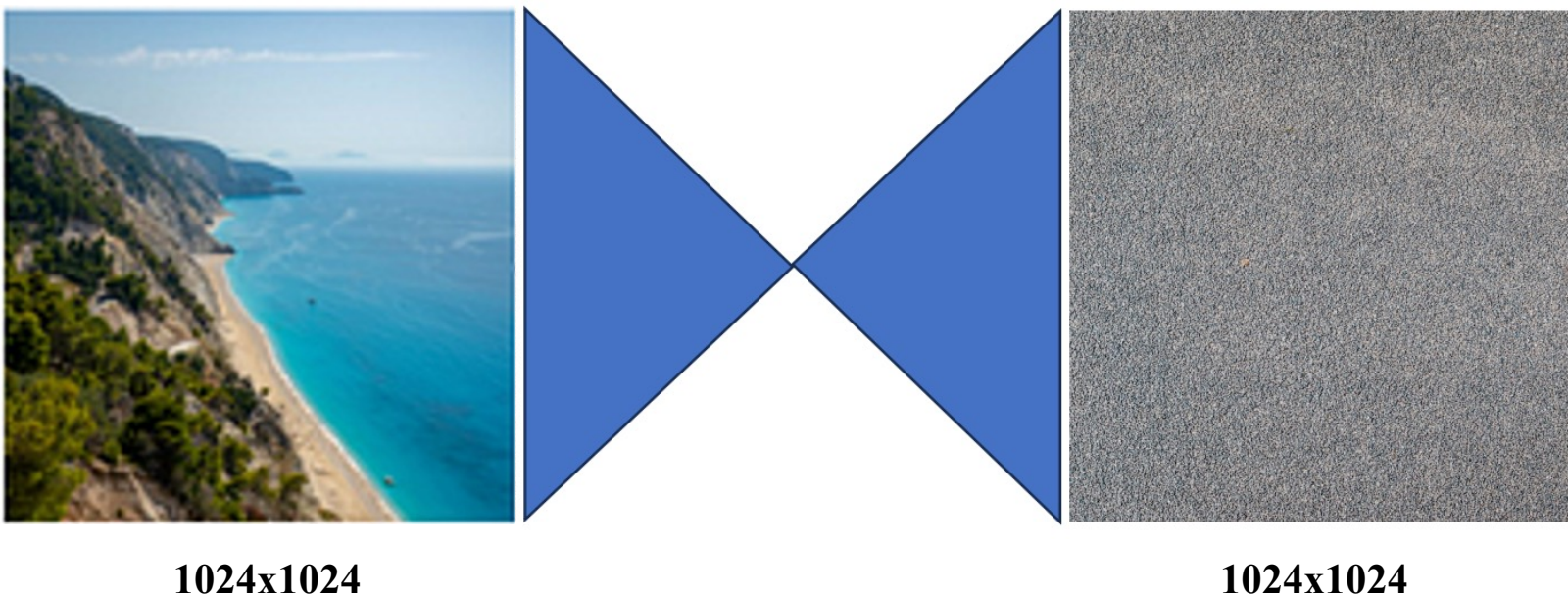




# Introduction

## ! Diffusion Model - Problem

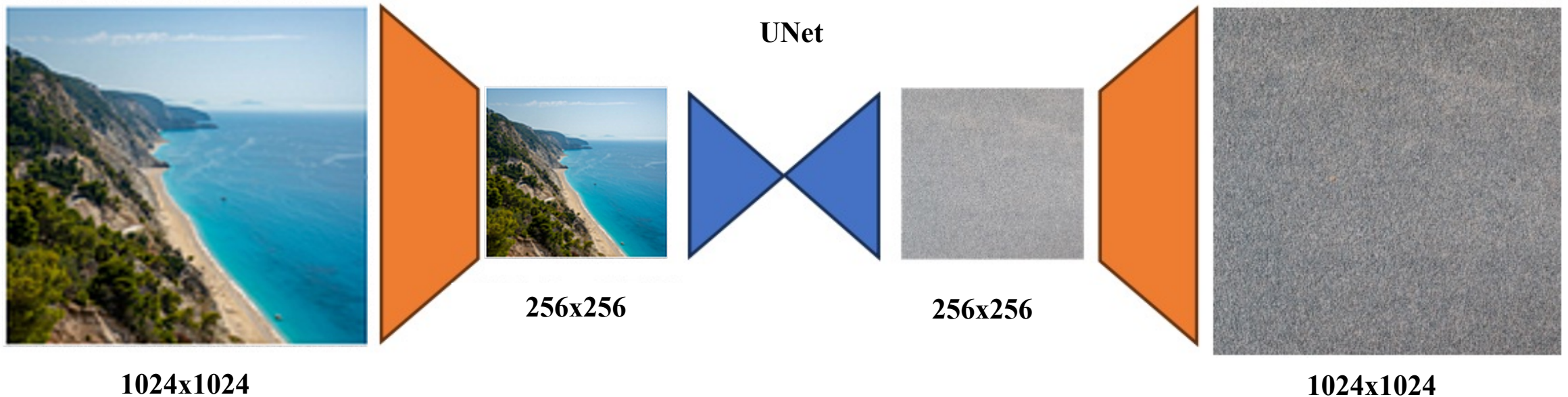
- Operating in the input space is very computationally expensive



# Introduction



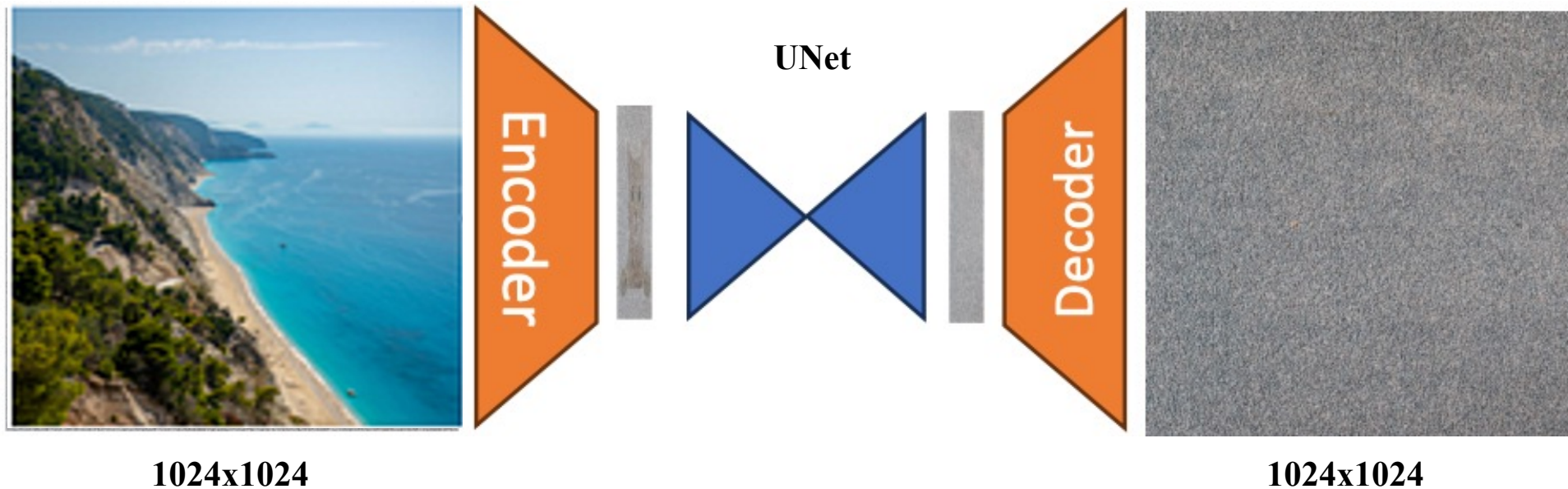
## Diffusion Model – Generate Low-Resolution + Upsample



# Introduction



## Diffusion Models – Generate in Latent Space





# Outline

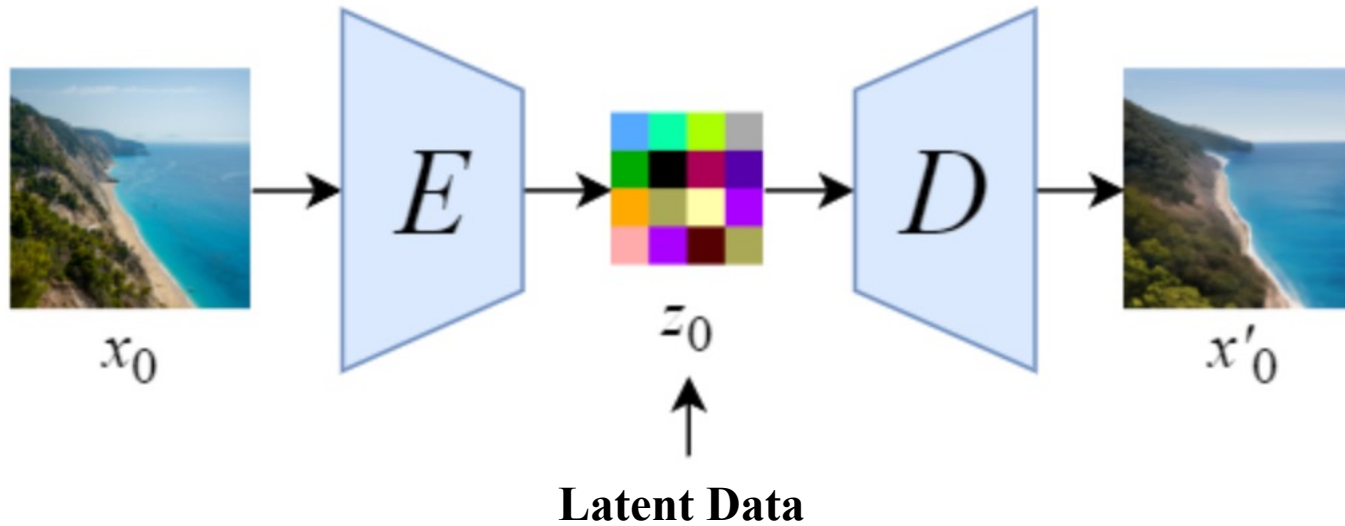
- **Introduction**
- **Stable Diffusion Model**
- **Text-to-Image Generation using SDM**

# Stable Diffusion Model



## Stable Diffusion Model (Latent Diffusion Model)

- The Diffusion process happens in the latent space
- First, train an autoencoder to learn to compress the image data into low-dimensional representation



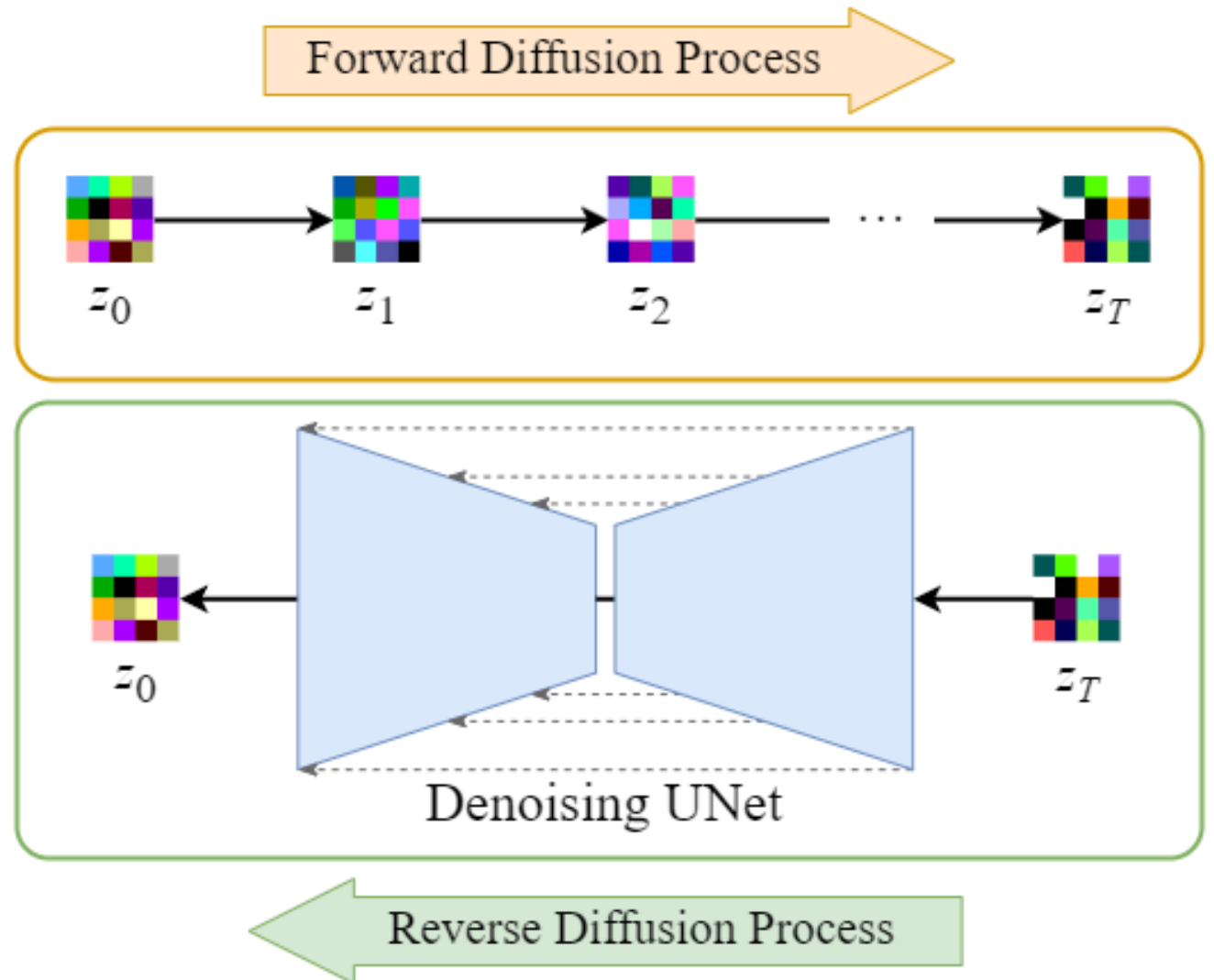


# Stable Diffusion Model



## Stable Diffusion Model (Latent Diffusion Model)

- After encoding the images into latent data, the forward and reverse diffusion processes will be done in the latent space.

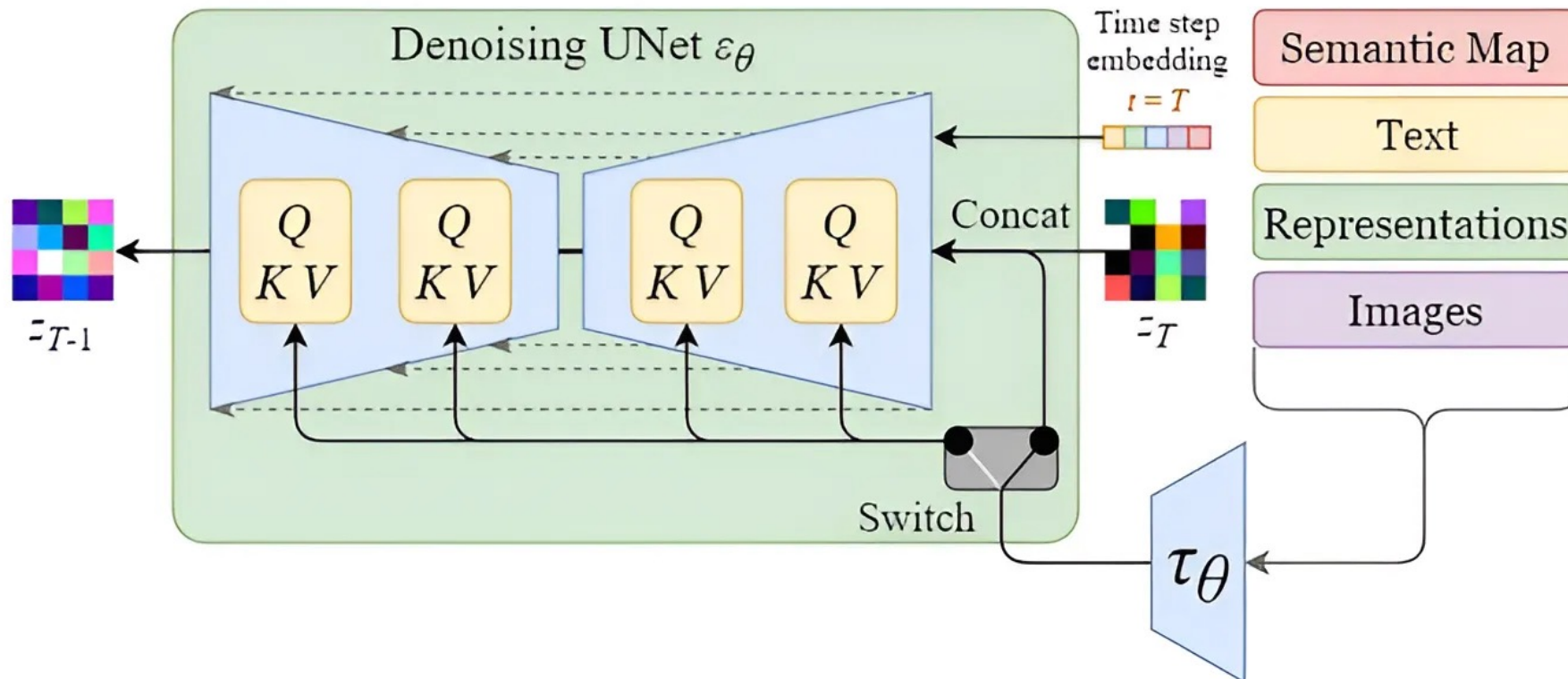


# Stable Diffusion Model



## Conditional Generation

- Condition denoising on text, images, etc,...



# Stable Diffusion Model



## Training

- The training objective (loss function) is pretty similar to the one in the pure diffusion model. The only changes are:
  - Input latent data  $z$  instead of the image  $x$
  - Added conditioning input  $\tau_\theta(y)$  to the UNet

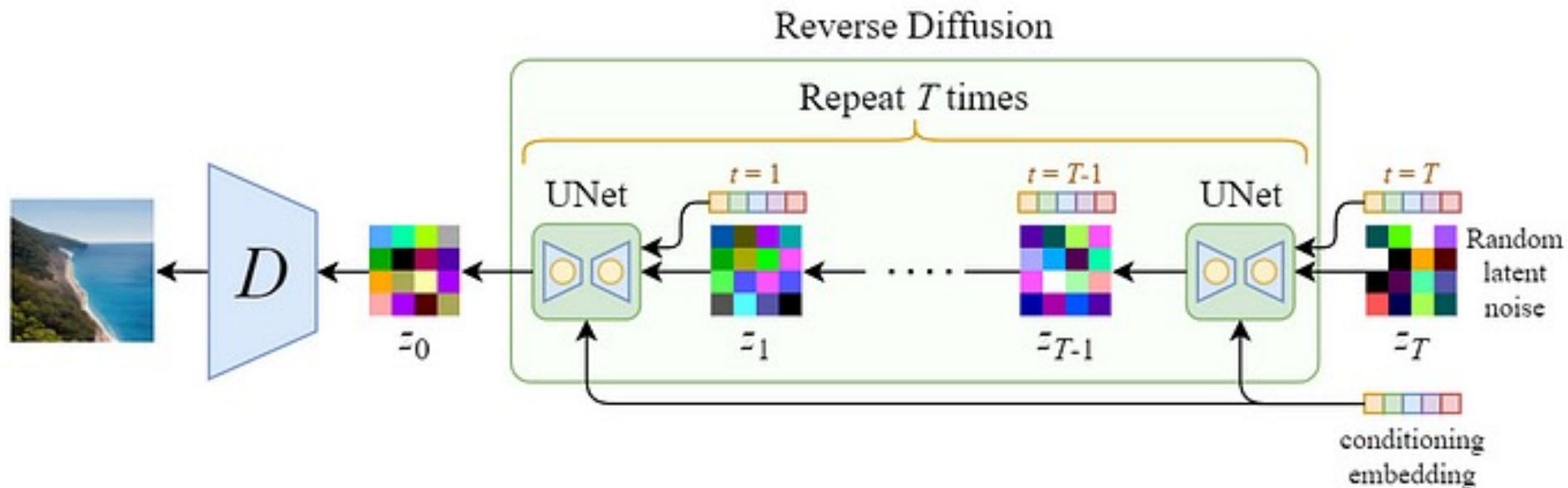
$$\begin{aligned} z_0 &= E(x_0) \\ z_t &= \sqrt{\bar{\alpha}_t} z_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon \\ L_{\text{LDM}} &= \mathbb{E}_{t, z_0, \epsilon, y} \left[ \left\| \epsilon - \epsilon_\theta(z_t, t, \tau_\theta(y)) \right\|^2 \right] \end{aligned}$$

Conditioning

# Stable Diffusion Model

## ! Sampling

- Stable Diffusion sampling process use the latent data
- The size of the latent data is much smaller than the original images, the denoising process will be much faster





# Outline

- **Introduction**
- **Stable Diffusion Model**
- **Text-to-Image Generation using SDM**



# Text-to-Image using SDM



## Celeb-HQ Dataset



# Text-to-Image using SDM



## Celeb-HQ Dataset

- The person has high cheekbones, and pointy nose. She is wearing lipstick.
- The person has high cheekbones, and pointy nose. She is wearing lipstick.
- She is wearing lipstick. She is young, and smiling and has big lips, mouth slightly open, pointy nose, and high cheekbones.
- This attractive woman has high cheekbones, pointy nose, bushy eyebrows, mouth slightly open, wavy hair, arched eyebrows, and bags under eyes.
- ...





# Text-to-Image using SDM



## Celeb-HQ Dataset

- The person has high cheekbones, and pointy nose. She is wearing lipstick.

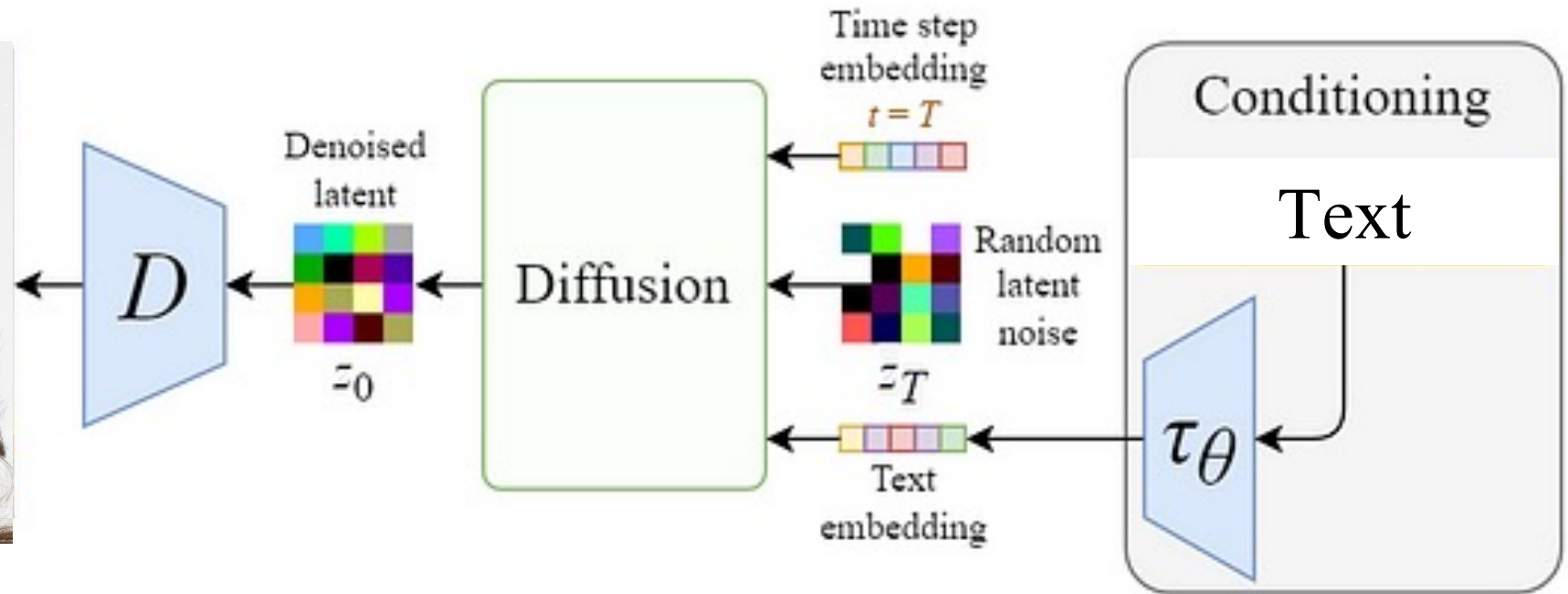




# Text-to-Image using SDM

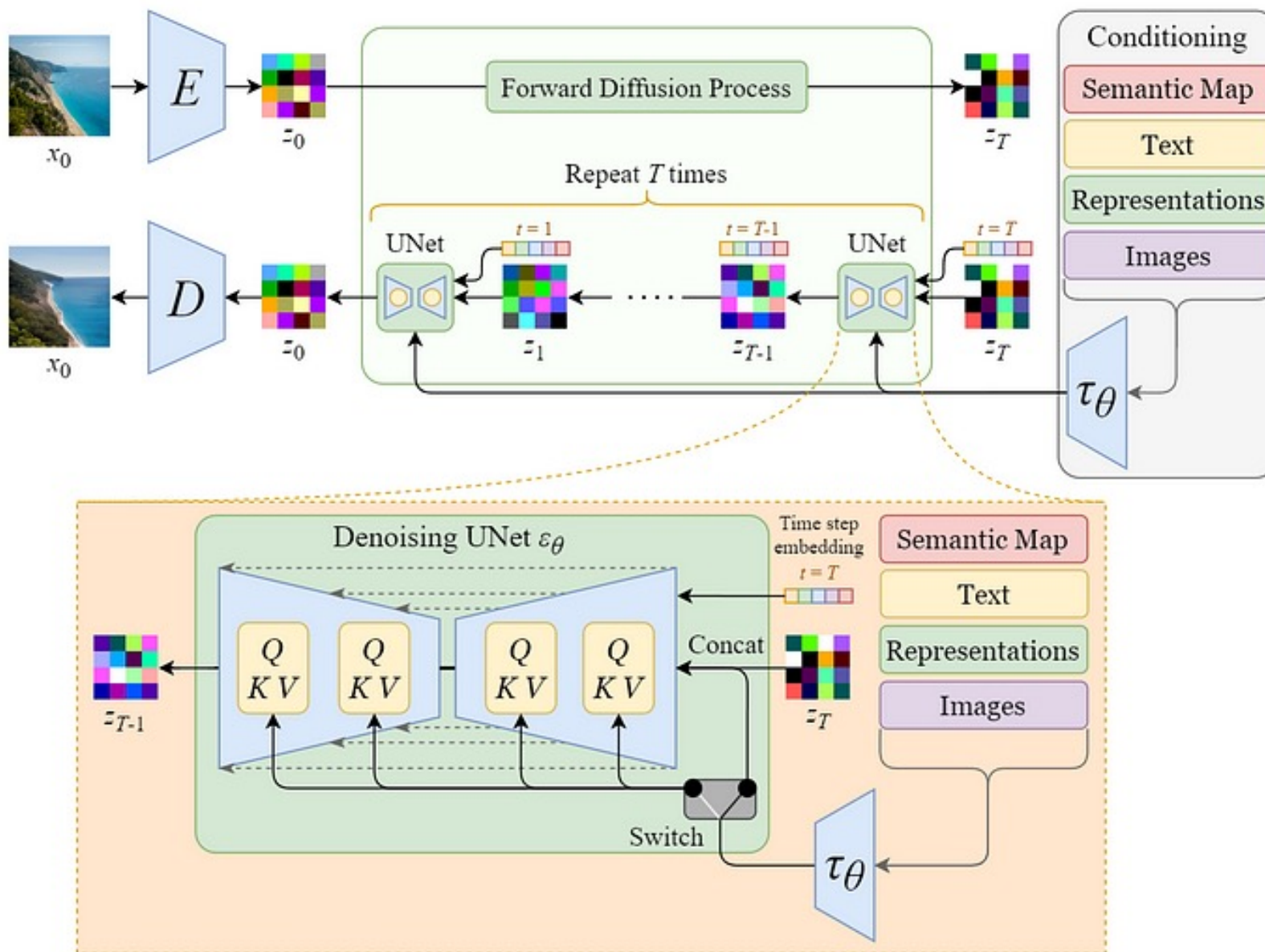


## Stable Diffusion Model – Text Condition



**BERT  
CLIP**

# Summary





AI VIET NAM

@aivietnam.edu.vn

# Thanks!

## Any questions?