

# Stable Diffusion

Give it a text prompt. It will return  
an image matching the text.

# Contents

- ❑ Image Generation
- ❑ The old method - GANs
- ❑ The new method - Diffusion
- ❑ Diffusion Process
- ❑ Variational AutoEncoder
- ❑ Latent Representation
- ❑ Conditioning



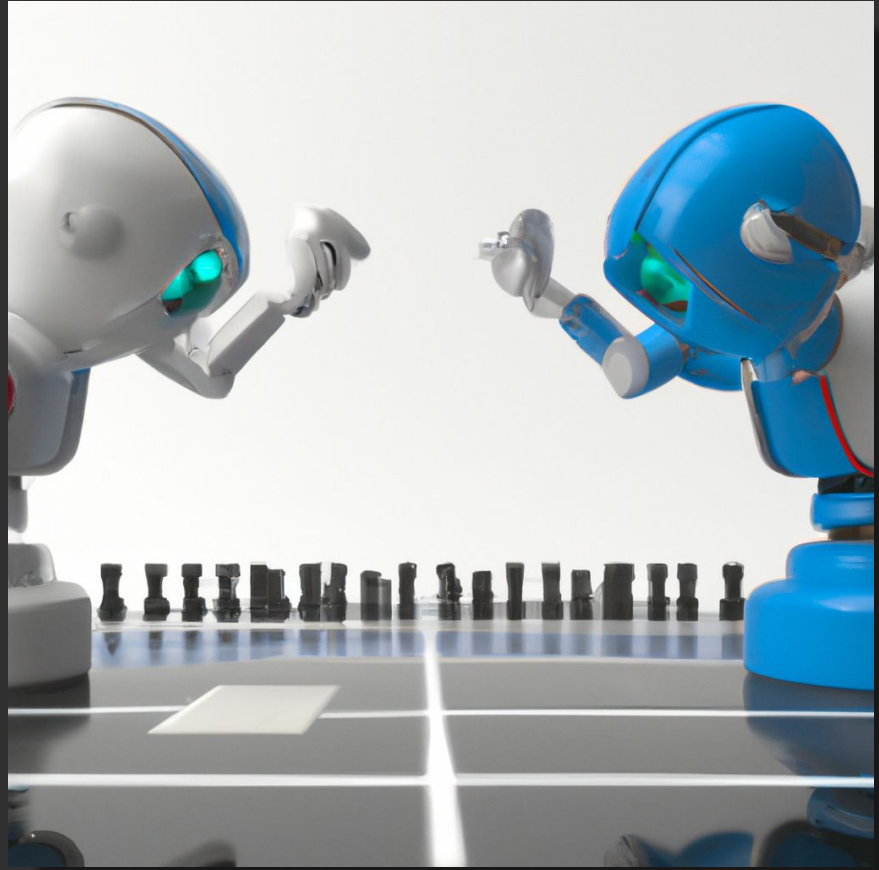
# Introduction

- ❑ Image generation is the creation of artificially generated images that look as realistic as real images.
- ❑ These images can be created by Generation Adversarial Networks(GAN) or with Variational Autoencoders, and more recently, Vector Quantized Variational Autoencoders (VQ-VAE), which create a discrete latent representation and create more variety of images and is easier to train compared to GANs.



# GANs

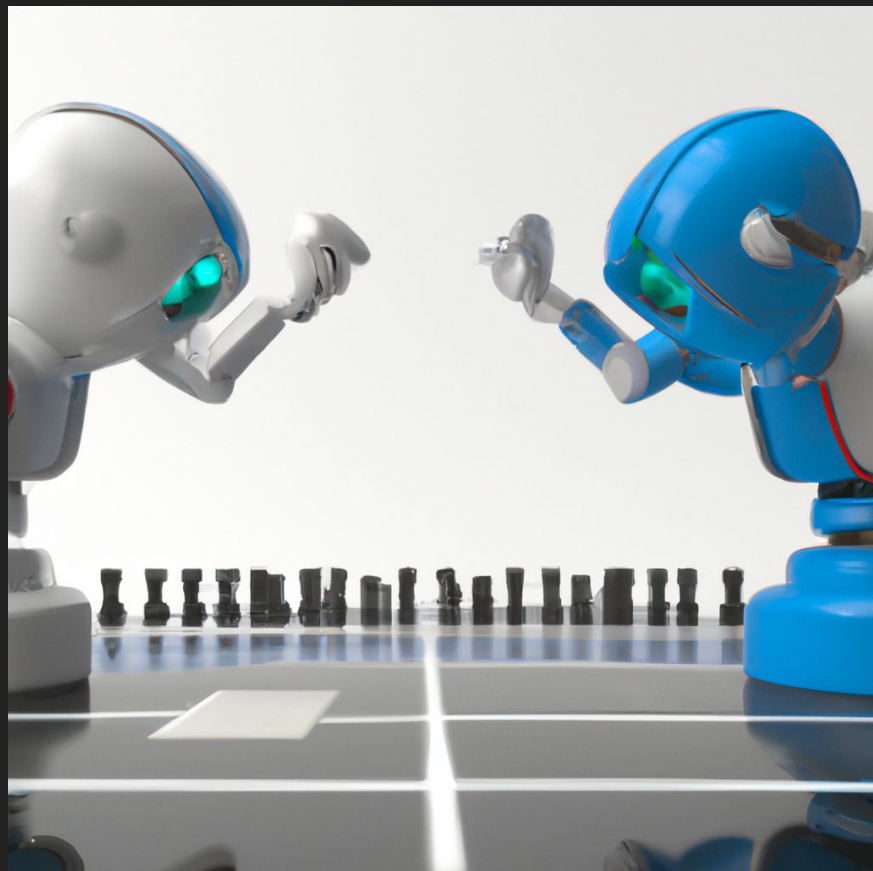
Generative Adversarial  
Networks



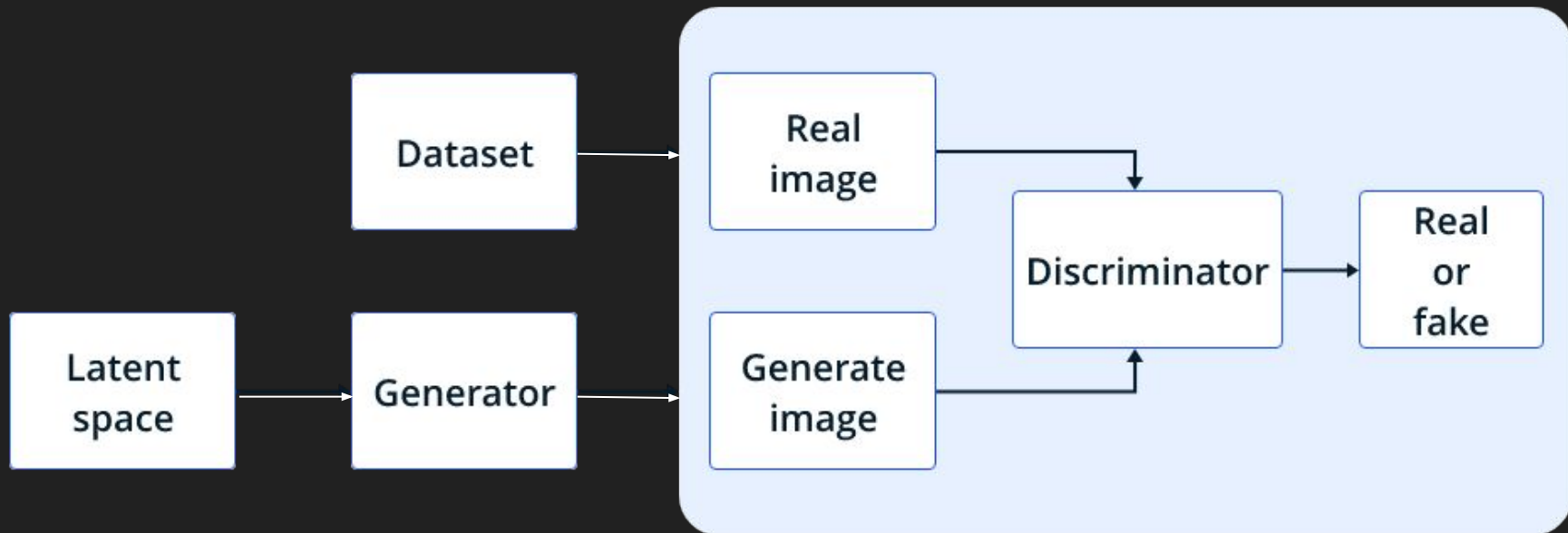


Dataset of labelled hand-drawn numbers

1  $\longleftrightarrow$  7



# GAN Architecture



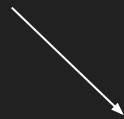


1 2 3  
4 5 6  
7 8 9

[Real Dataset]



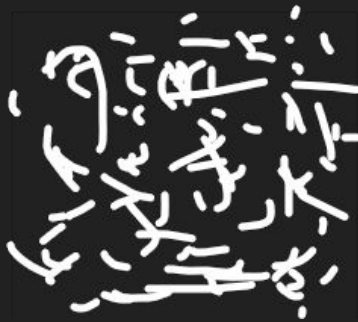
8



Discriminator

1  
[REAL]

0  
[FAKE]



[Noise]

→ Generator →



[Generator Output]

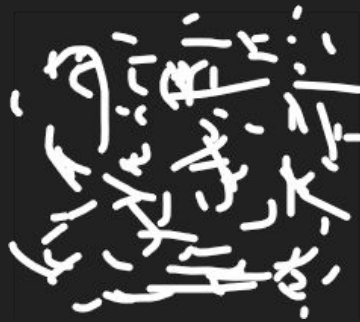
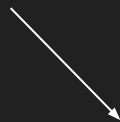


1 2 3  
4 5 6  
7 8 9

[Real Dataset]



8



[Noise]

→ Generator →

5

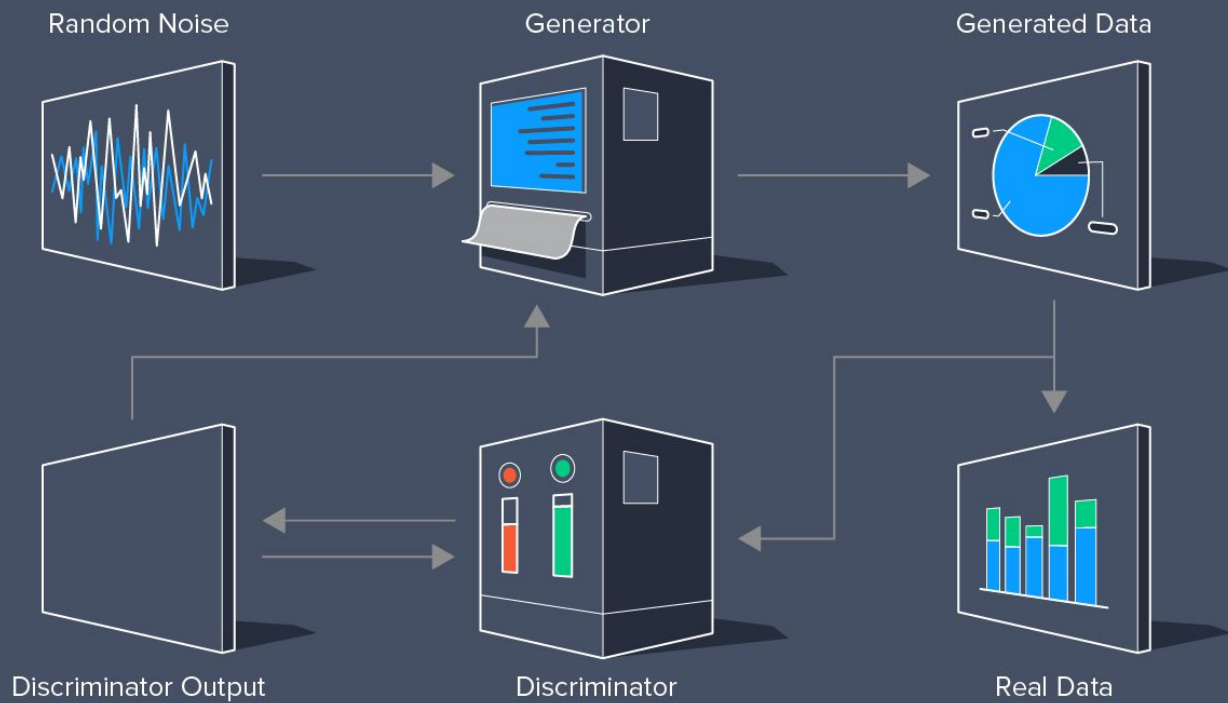
[Generator Output]



Discriminator

1  
[REAL]

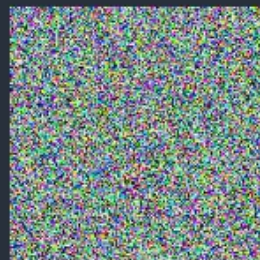
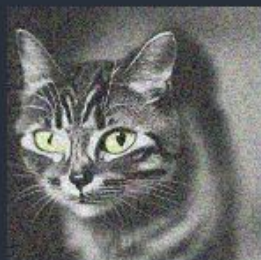
0.9  
[FAKE]



# How Diffusion Models Work ?

Forward and Reverse Diffusion

# Forward diffusion



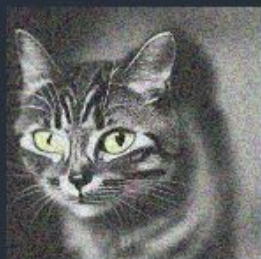
Step

1

2

3

4



Noise 1

Noise 2

Noise 3

Noise 4

Training examples are created by generating **noise** and adding an **amount** of it to the images in the training dataset (forward diffusion)

1

Pick an image



2

Generate some random **noise**

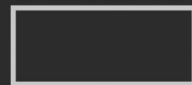


Noise sample 1

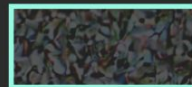
3

Pick an amount of **noise**

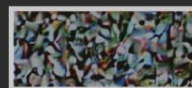
0



1



2



3









4

Add **noise** to the image in that **amount**



# DATASET

# MODEL

INPUT		OUTPUT / LABEL
Noise Amount	Noisy Image	
3		
14		
7		
42		
2		
21		

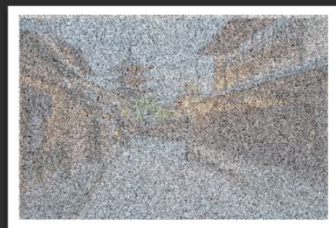
Noise  
Predictor  
(UNet)



# UNet training step

1

Pick a training example from the training dataset



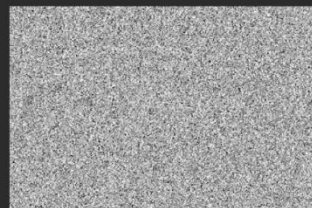
Noise  
amount:  
3

2

Predict the noise

UNet

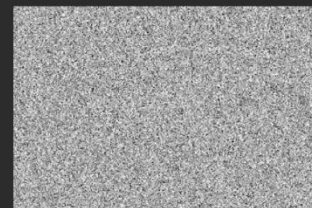
Unet Prediction



3

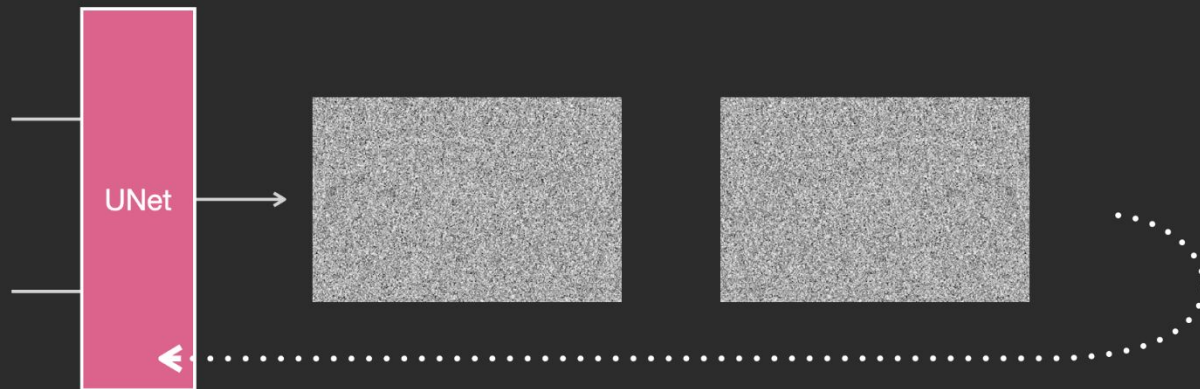
Compare to actual noise (calculate loss)

Actual noise (Label)

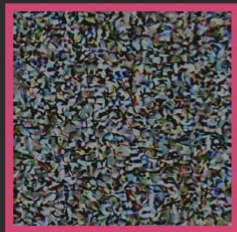


4

Update model  
(backprop)



Predicted  
noise sample



Noise amount:  
3



Trained  
Noise  
Predictor  
(UNet)

# Reverse Diffusion (Denoising)

## Step 1

Slightly  
de-noised  
image



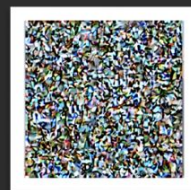
=



-

Subtract  
predicted noise  
from image

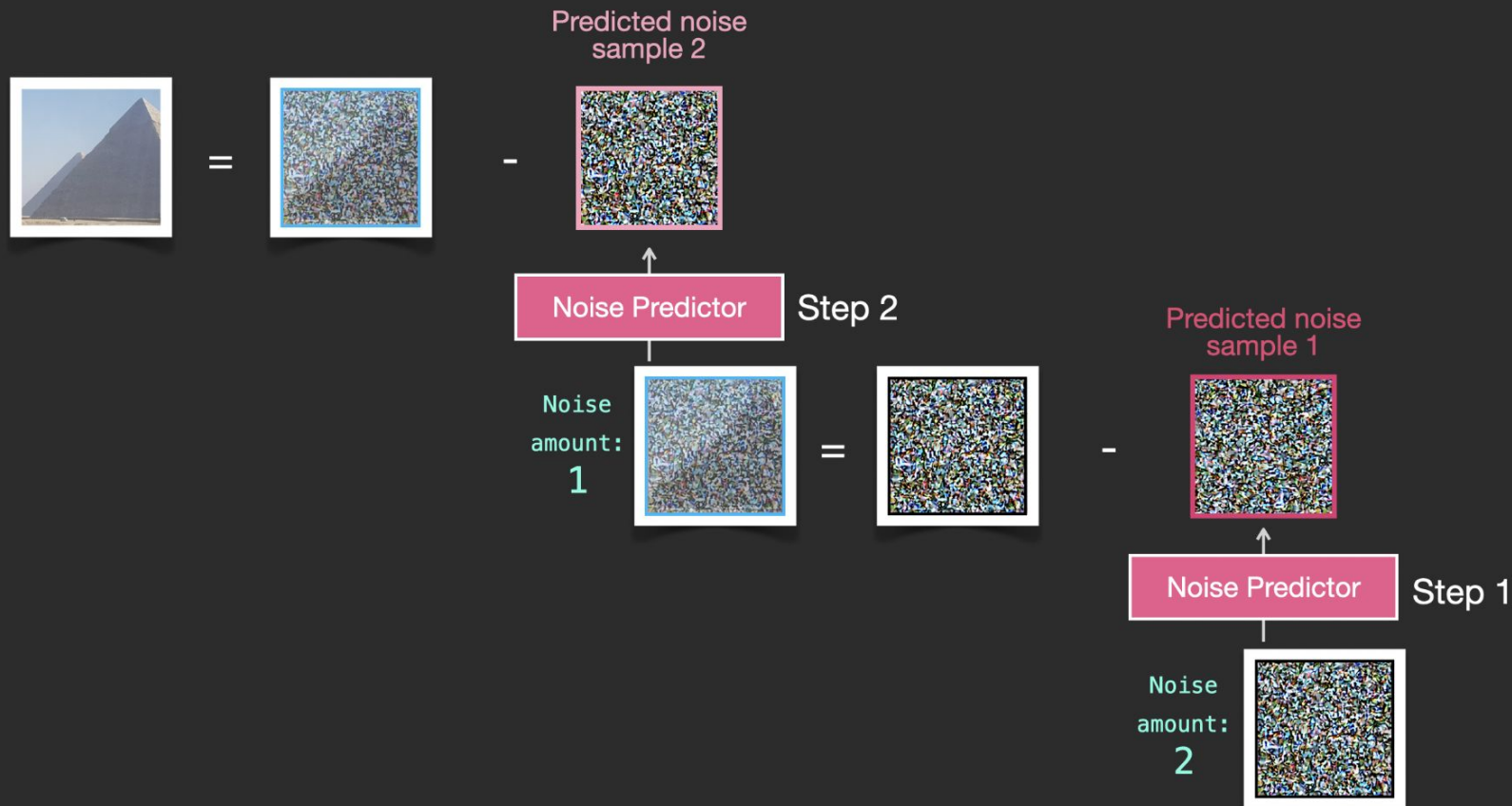
Predicted  
noise sample



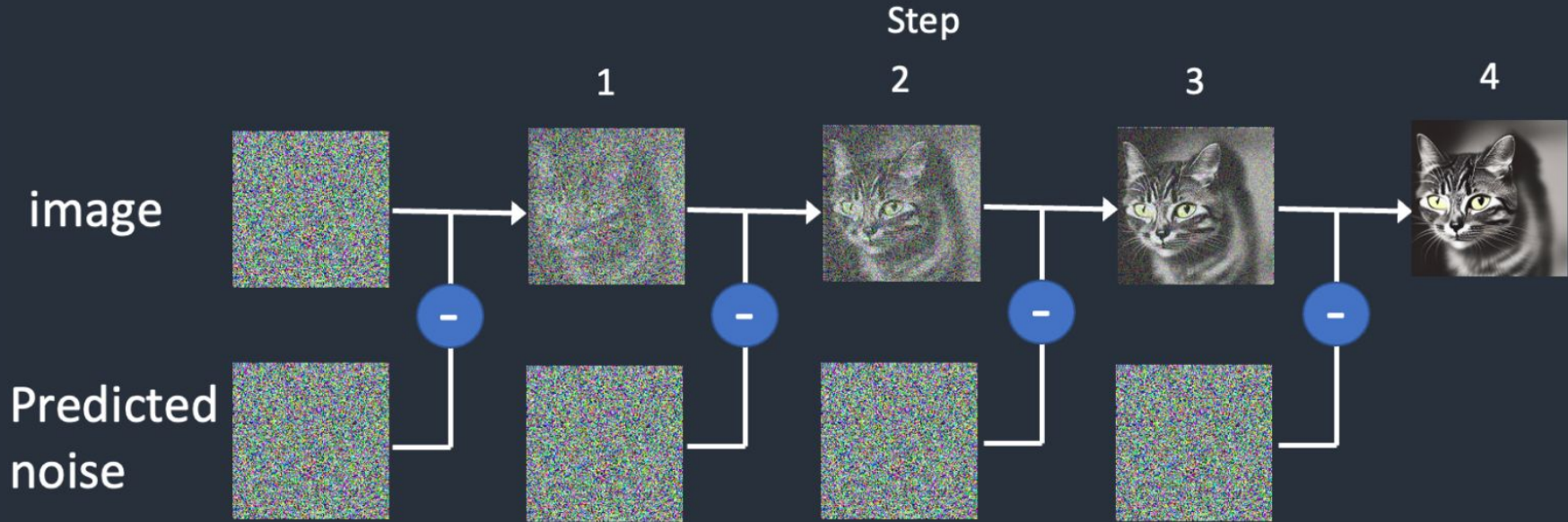
Noise amount:  
3

Trained  
Noise  
Predictor  
(UNet)

# Image Generation by Reverse Diffusion (Denoising)



# Reverse Diffusion



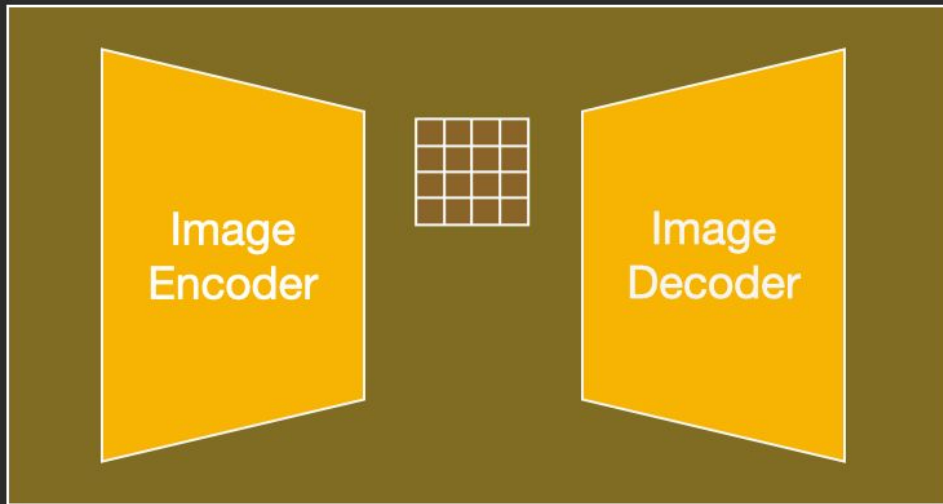
Reverse diffusion works by subtracting the predicted noise from the image successively.

# Departure to Latent Space

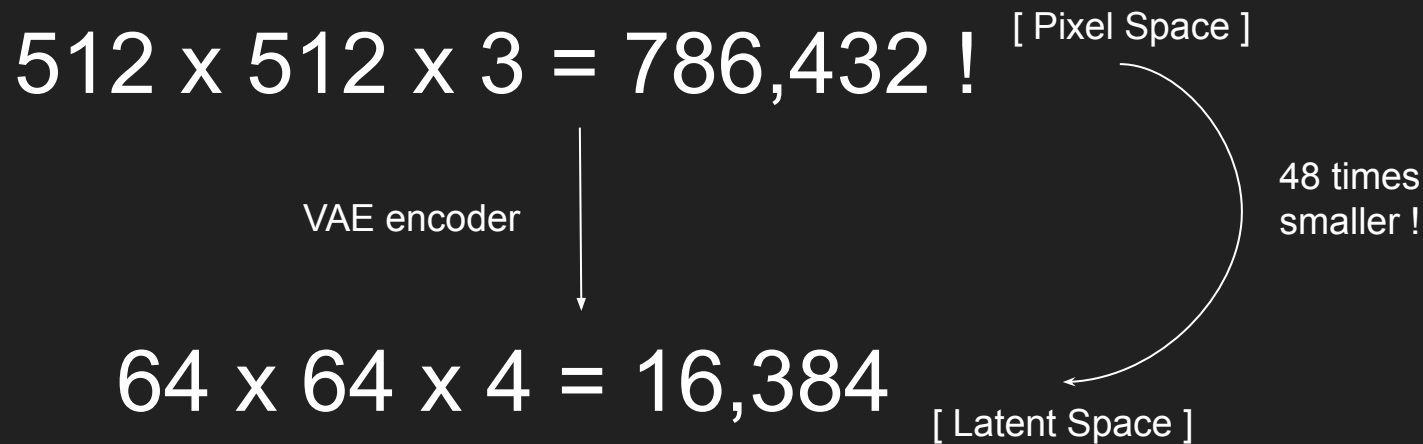
Original  
image



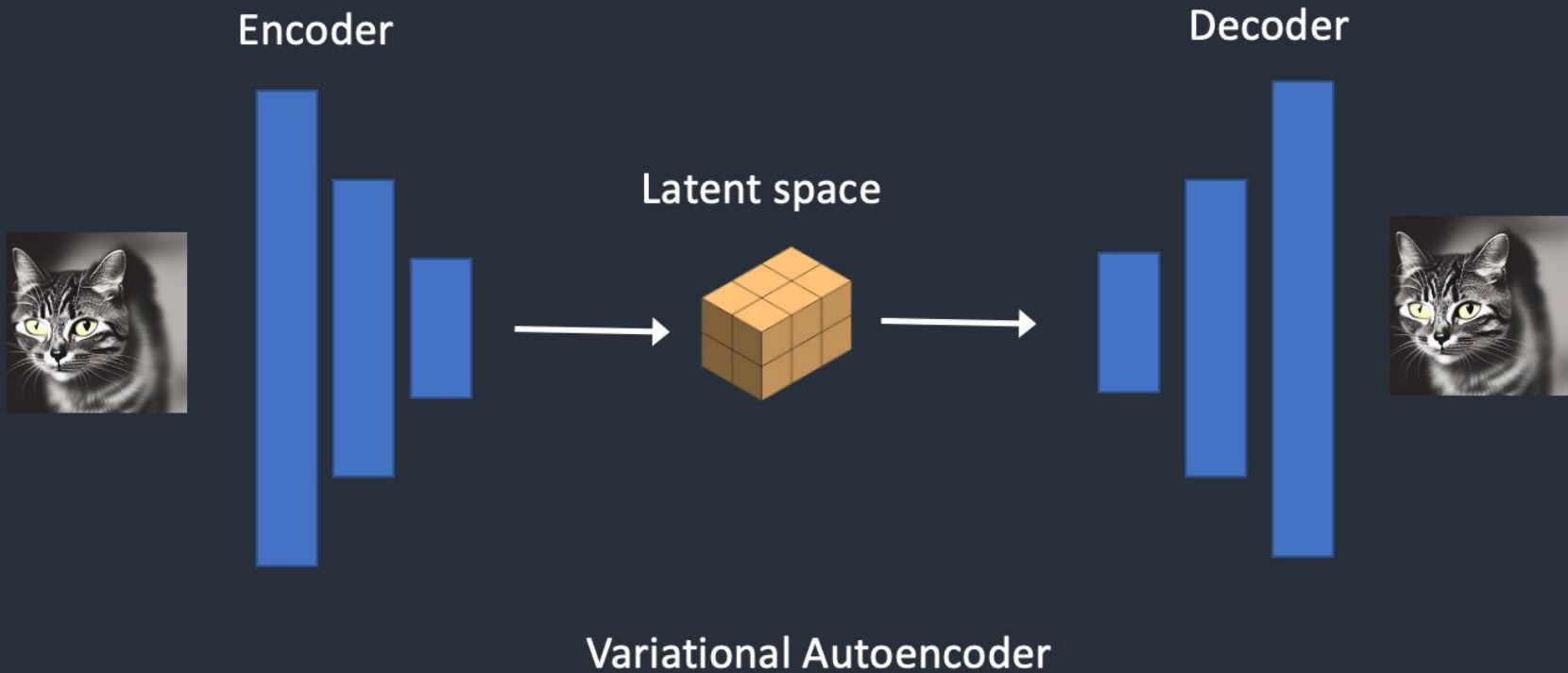
Generated  
image



# Latent Space



# Variational AutoEncoder (VAE)





# Latent Space

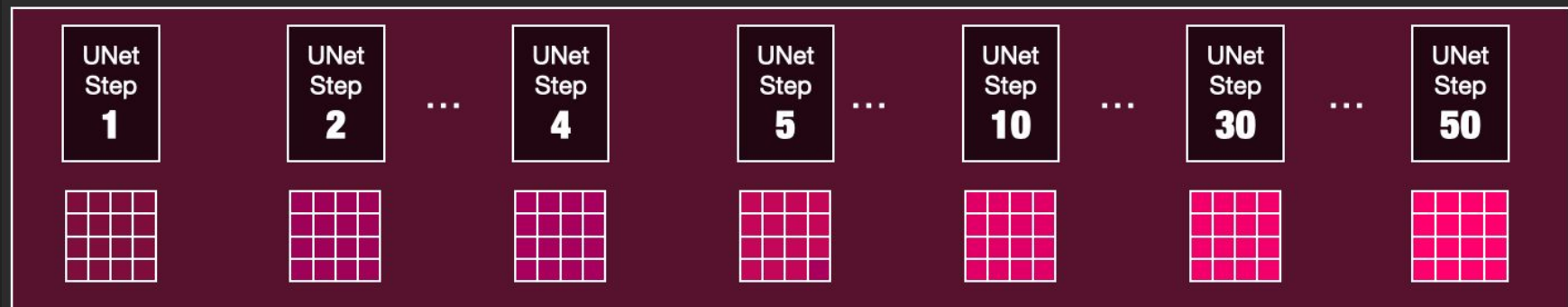


Image Information Creator

Image  
Decoder  
(Autoencoder  
decoder)





Random image  
information tensor



Image Information Creator  
(UNet + Scheduler)

Processed image  
information tensor



Information  
World  
(Latent space)

Image  
Decoder  
(Autoencoder decoder)



Visual  
World  
(Pixel space)

# Quick Recap

Original image



Image  
Encoder

Generate training examples with different amounts  
of noise added to their compressed/latent version



Compressed  
image (latent)



Latent + noise  
sample 1 at  
noise amount 1



Latent + noise  
sample 2 at  
noise amount 2



Latent + noise  
sample 3 at  
noise amount 3



Image  
Decoder

Image Generation by Reverse Diffusion (Denoising)



Processed  
Image  
Information

UNet  
Step  
**50**



UNet  
Step  
**2**



UNet  
Step  
**1**



Complete  
noise

Image Information Creator

Generated image

OK but when does the text  
prompt come into the picture ?

[ Pun Intended ]

# Conditioning

## Stable Diffusion


	paradise
	cosmic
	beach

Text  
Understander  
(Encoder)

Image  
Generator



# Stable Diffusion



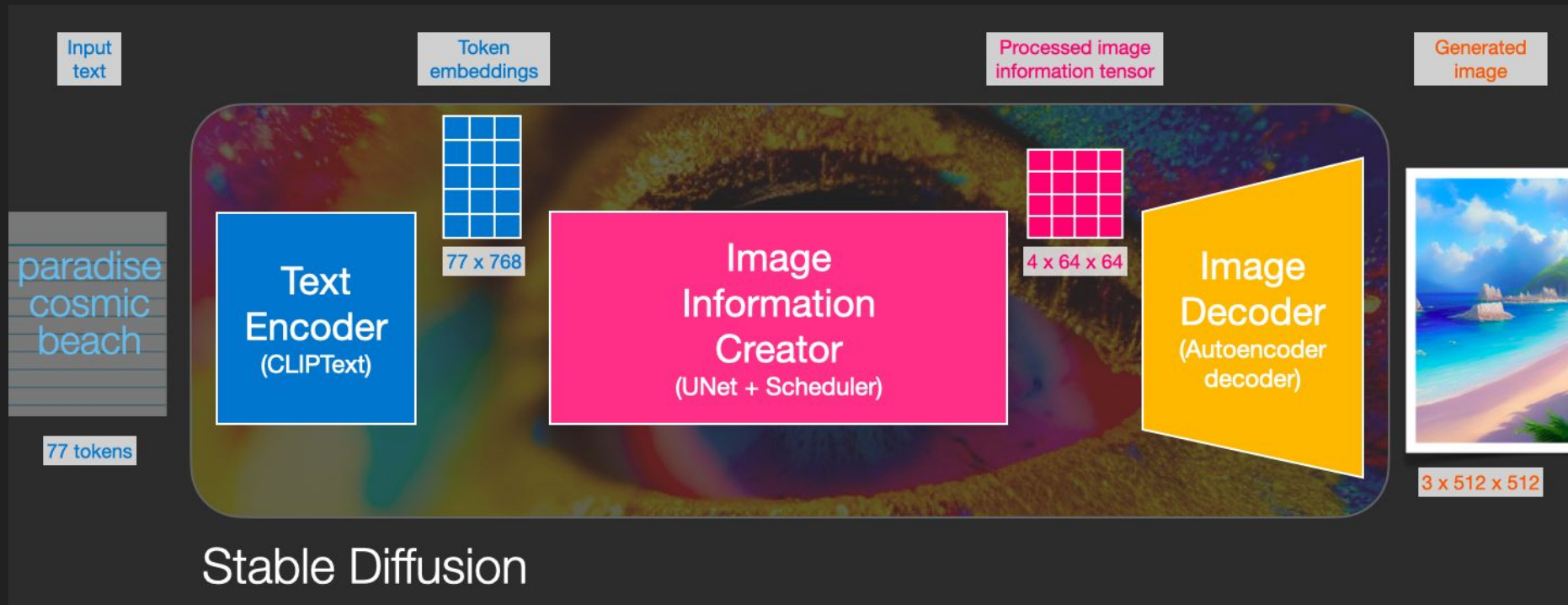
The diagram illustrates the architecture of Stable Diffusion, which is composed of three main components arranged horizontally. The background of the diagram is a colorful, abstract image. The first component on the left is a blue square labeled 'Text Encoder (CLIPText)'. The middle component is a pink rectangle labeled 'Image Information Creator (UNet + Scheduler)'. The third component on the right is a yellow trapezoid labeled 'Image Decoder (Autoencoder decoder)'.

**Text  
Encoder**  
(CLIPText)

**Image  
Information  
Creator**  
(UNet + Scheduler)

**Image  
Decoder**  
(Autoencoder  
decoder)





# Stable Diffusion

paradise  
cosmic  
beach

77 tokens

Text  
Encoder  
(CLIPText)



Token  
embeddings



Random image  
information tensor

Image Information Creator  
(UNet + Scheduler)

UNet  
Step  
**1**



UNet  
Step  
**2**



...

UNet  
Step  
**50**



Processed image  
information tensor

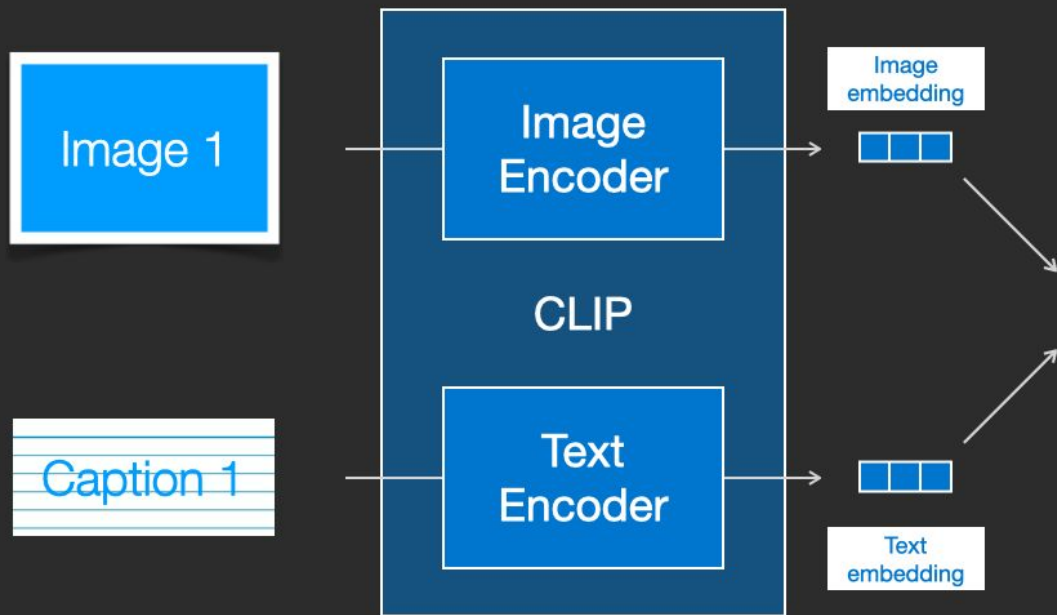
Image  
Decoder  
(Autoencoder  
decoder)

Diffusion

Generated  
image

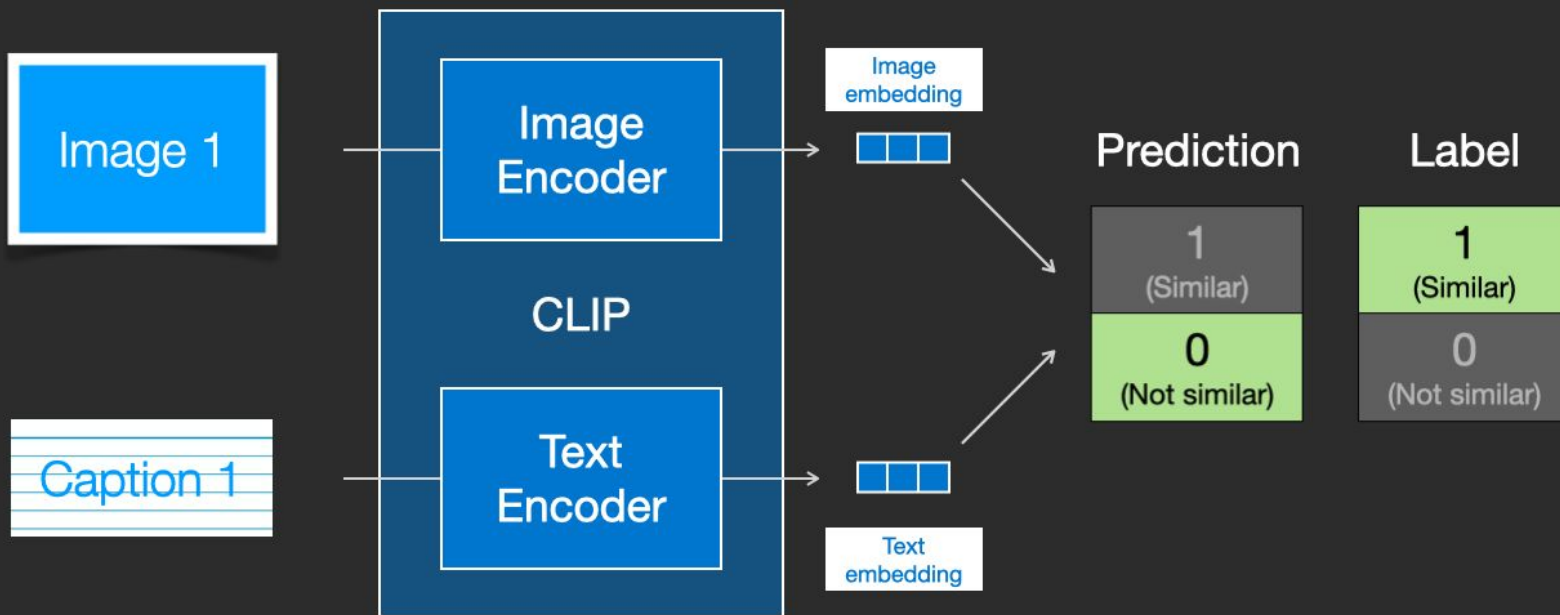


**1** Embed image and text



1 Embed image and text

2 Compare the embeddings



1

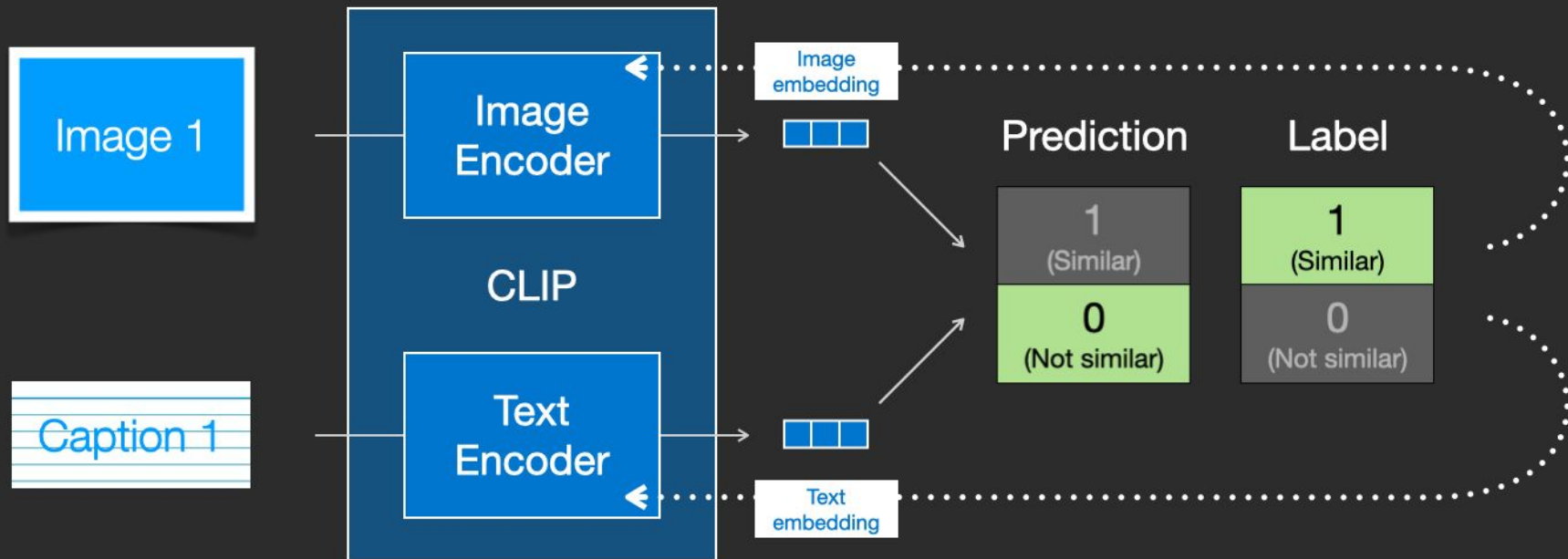
Embed image  
and text

2

Compare the  
embeddings

3

Update the  
models





Noise amount:  
3

Prompt text  
information  
(token embeddings)





















Noise  
Predictor  
(UNet)  
with Text  
Conditioning

Predicted noise  
sample



# DATASET

# MODEL

Step	INPUT		OUTPUT / LABEL
	Image	Text	noise sample
3			
14			
7			
42			
2			
21			

Noise Predictor  
(UNet)  
with Text  
Conditioning

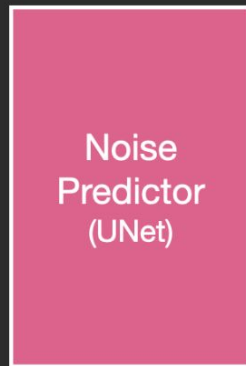
Without Text



Noisy compressed  
image  
(latents)



Noise amount:  
3



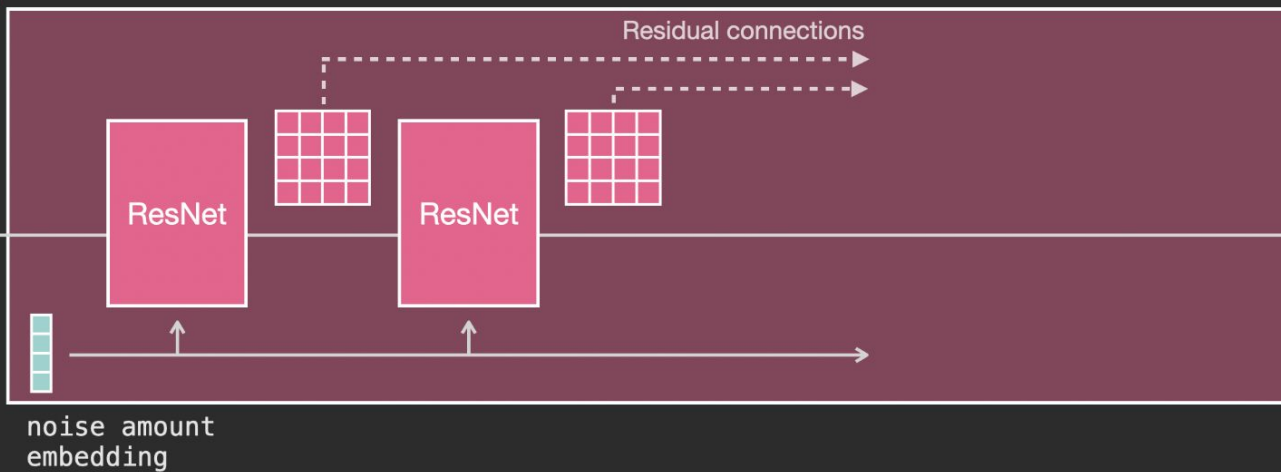
Predicted  
noise sample

Noisy  
compressed  
image  
(latents)



Noise  
amount:  
3

## Noise Predictor (UNet)



Predicted  
noise sample



With Text

Noisy compressed  
image  
(latents)



Text information  
(token embeddings)



Noise amount:  
3

Noise  
Predictor  
(UNet)

Predicted  
noise sample



Noisy  
compressed  
image  
(latents)

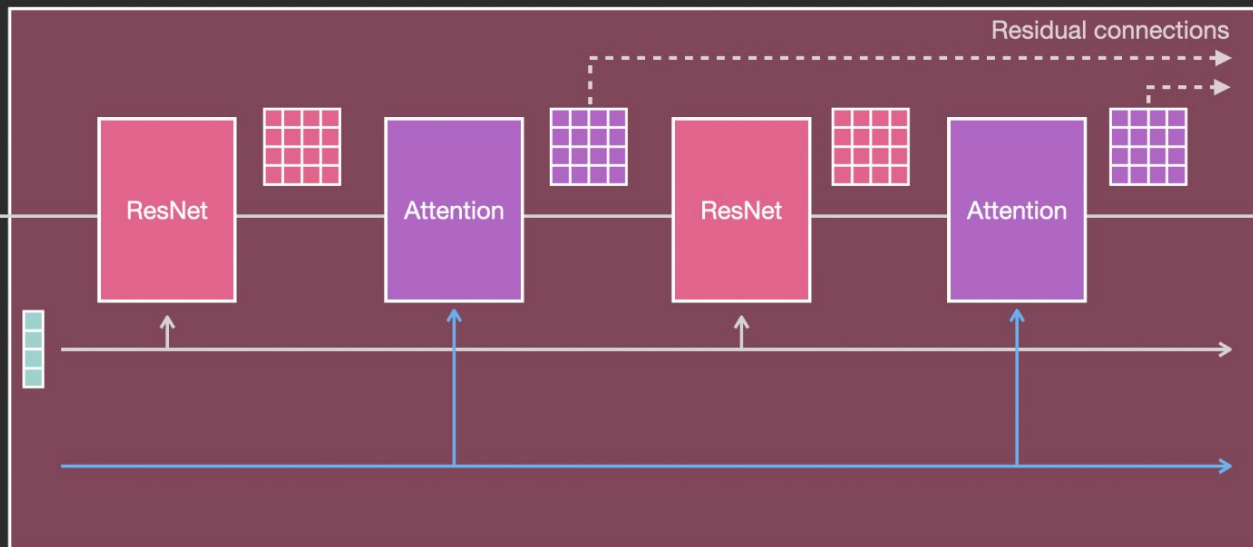


Noise amount:  
3



Text information  
(token embeddings)

## Noise Predictor with Text Conditioning (UNet with attention)



Predicted  
noise sample



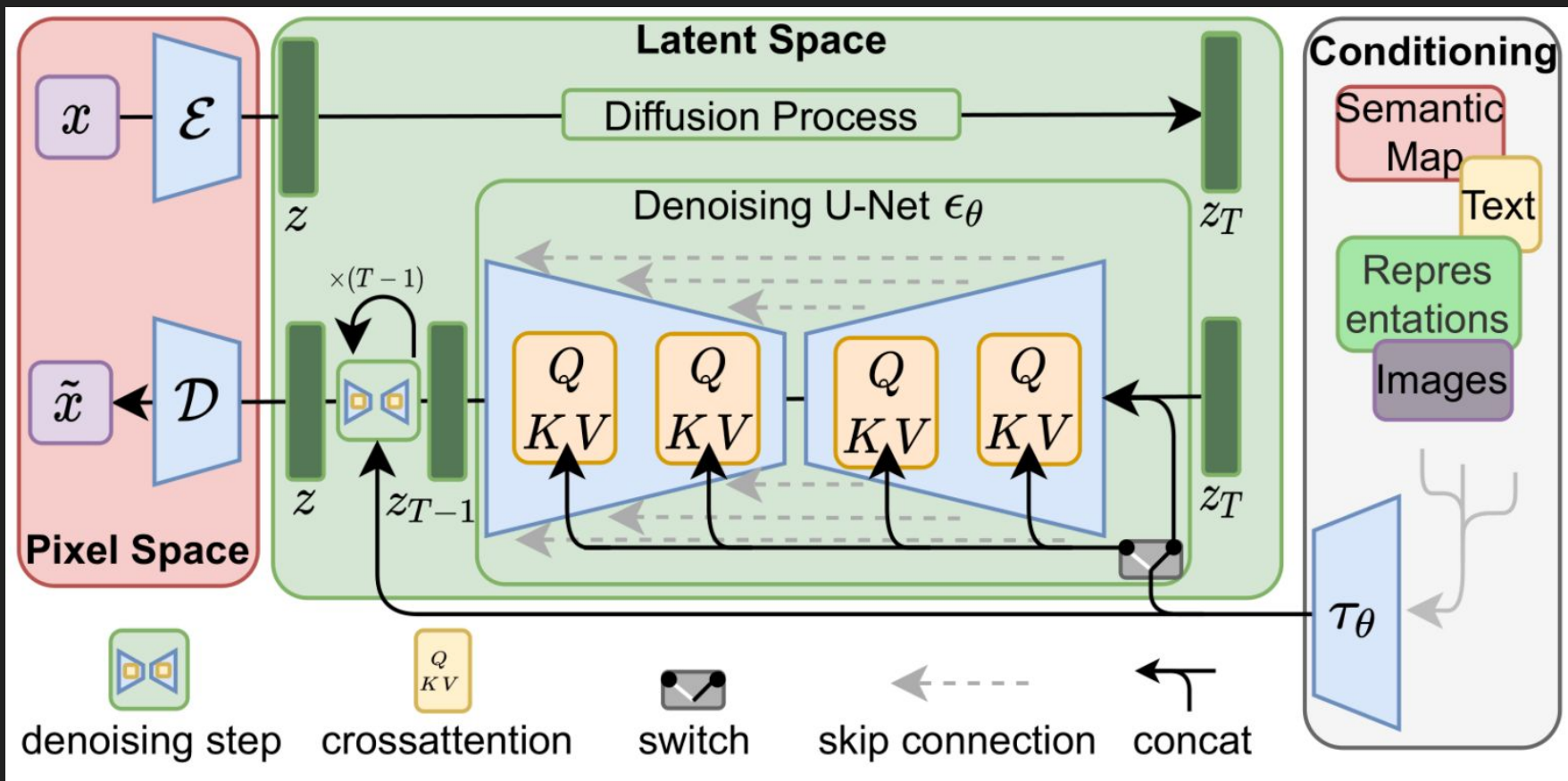
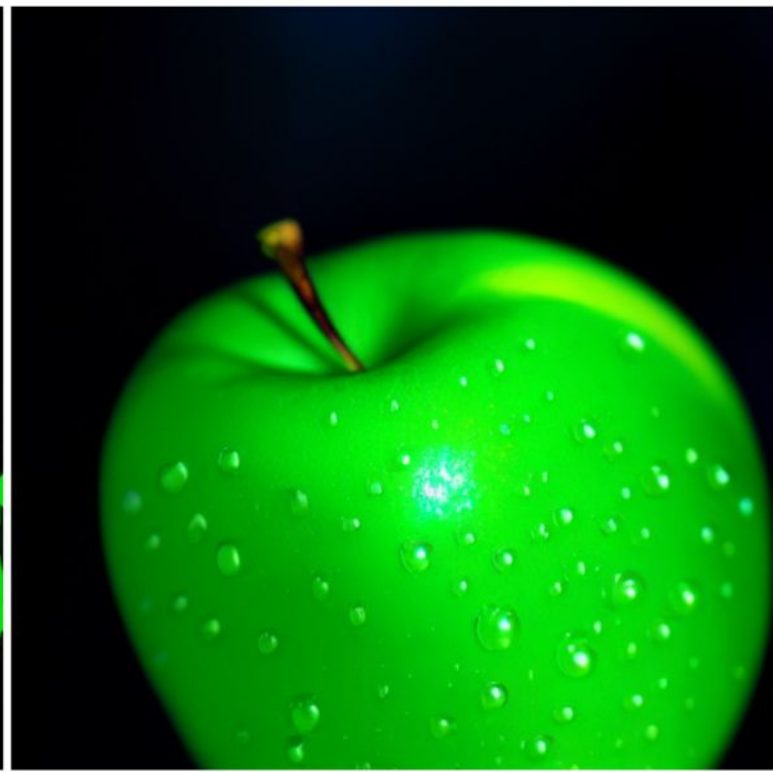
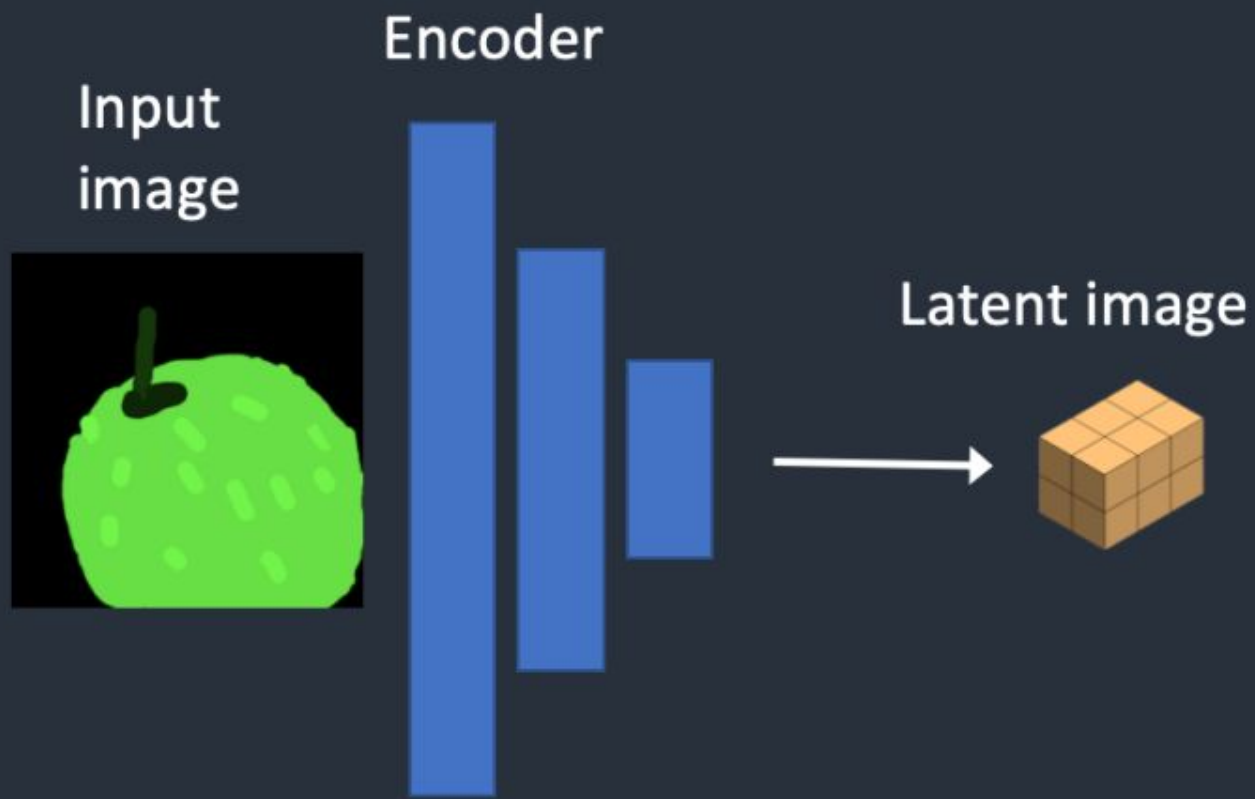


Image-to-Image







# REFERENCES

- ❏ <https://huggingface.co/blog/annotated-diffusion>
- ❏ <https://jalammar.github.io/illustrated-stable-diffusion/>
- ❏ <https://stable-diffusion-art.com/how-stable-diffusion-work/>