

# Stable Diffusion fine tuning approaches

Dominik Lewy

## LLMs Landscape

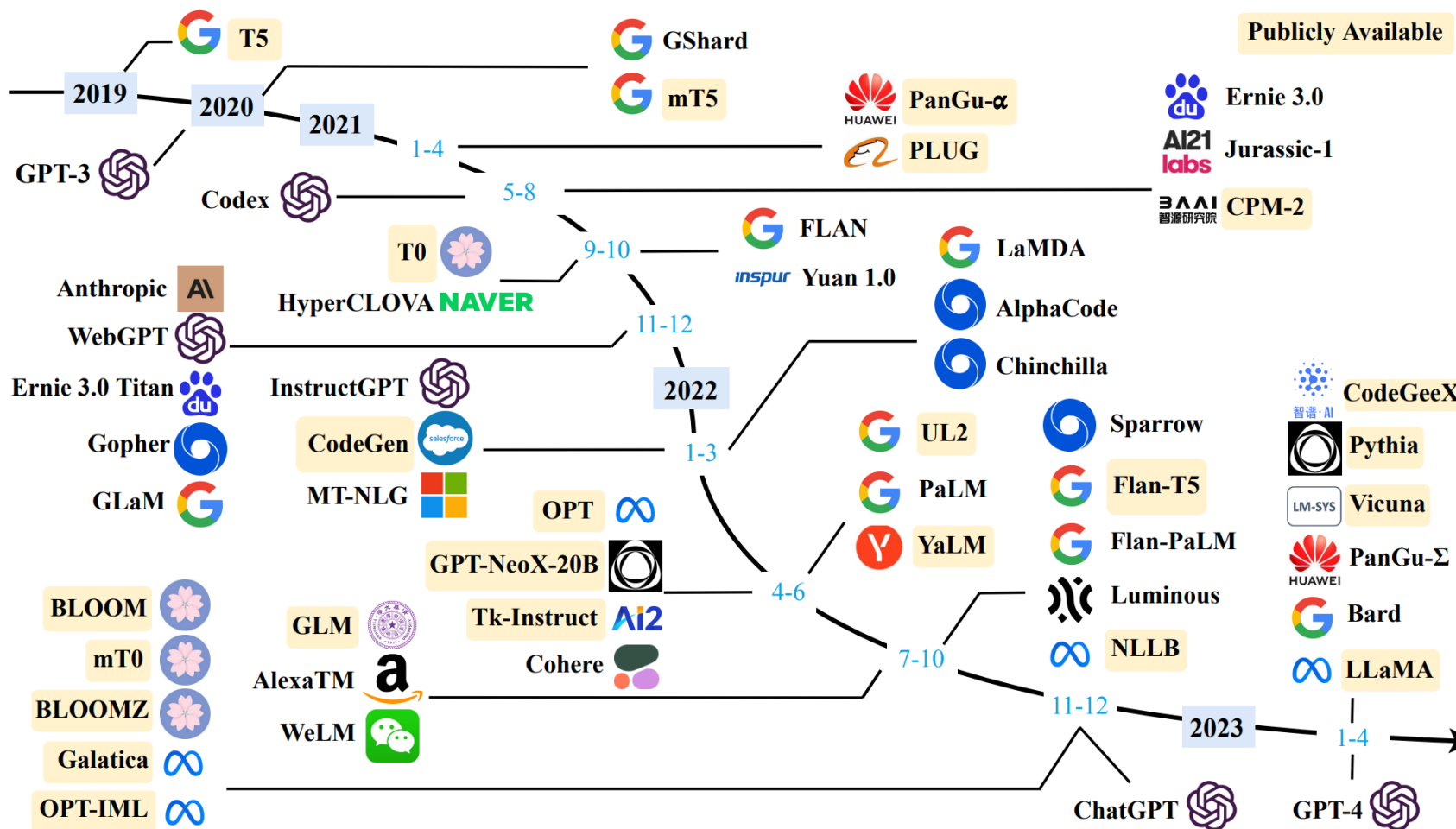


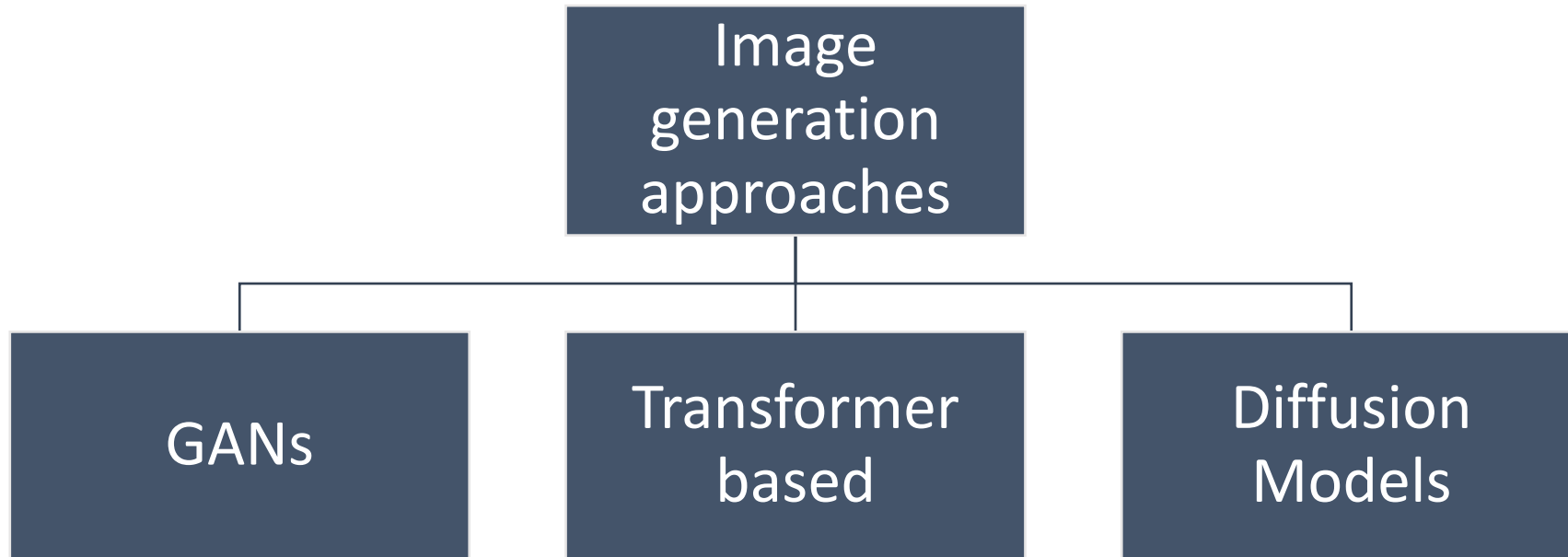
Fig. 1. A timeline of existing large language models (having a size larger than 10B) in recent years.

What will be tomorrow?

	PRE - 2020	2020	2022	2023?	2025?	2030?
TEXT	Spam detection Translation Basic Q&A	Basic copy writing First drafts	Longer form Second drafts	Vertical fine tuning gets good (scientific papers, etc)	Final drafts better than the human average	Final drafts better than professional writers
CODE	1-line auto-complete	Multi-line generation	Longer form Better accuracy	More languages More verticals	Text to product (draft)	Text to product (final), better than full-time developers
IMAGES			Art Logos Photography	Mock-ups (product design, architecture, etc.)	Final drafts (product design, architecture, etc.)	Final drafts better than professional artists, designers, photographers)
VIDEO / 3D / GAMING			First attempts at 3D/video models	Basic / first draft videos and 3D files	Second drafts	AI Roblox Video games and movies are personalized dreams

Large model availability: ● First attempts ● Almost there ● Ready for prime time

Historical outlook



# Transformer based

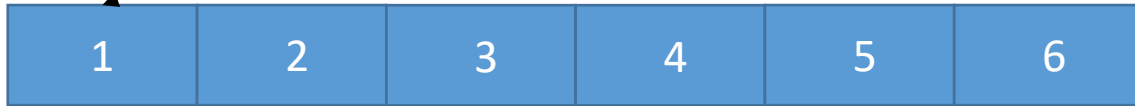
Encoding for images



**DALL-E in practice**

Input: Sentence A

Text BPE vocabulary of size  
**16384**



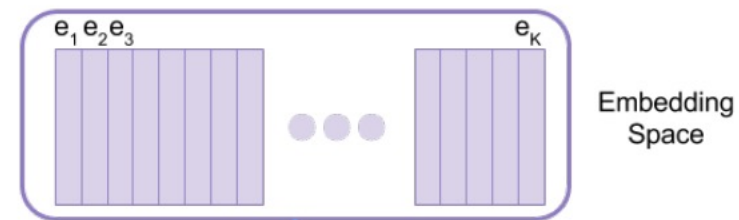
Text token

**DALL-E in practice**

Input: Sentence A

Text BPE vocabulary of size  
**16384**

Image learned vocabulary of size  
**8192**



Text token

Image token

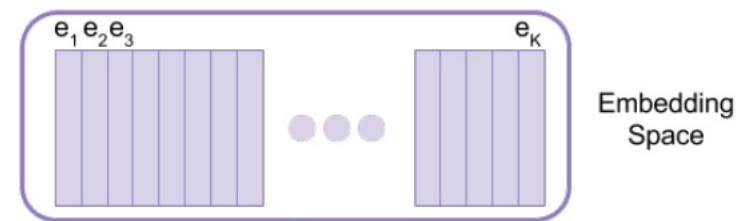


## DALL-E in practice

Input: Sentence A

Text BPE vocabulary of size  
**16384**

Image learned vocabulary of size  
**8192**



Text token are connected to all input tokens via attention

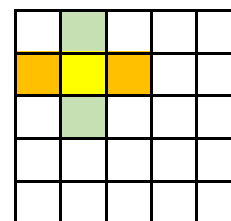
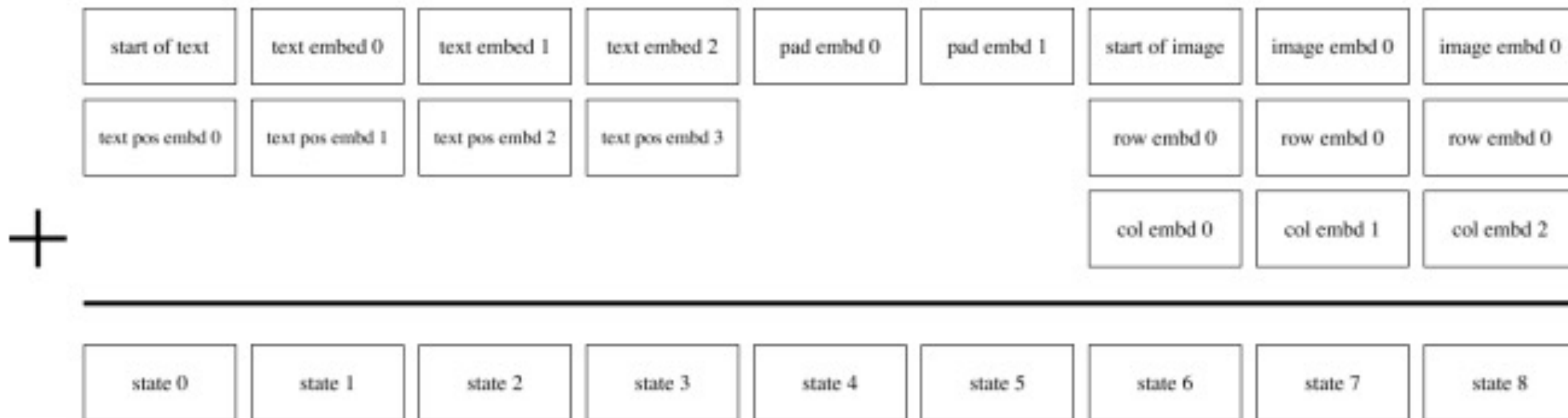


Image token in addition to seeing text tokens have column and row attention

Text token

Image token

## DALL-E in practice



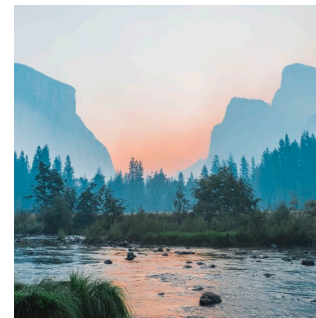
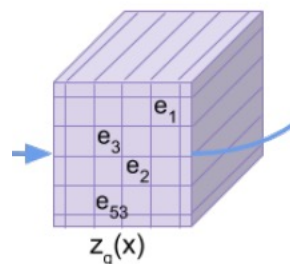
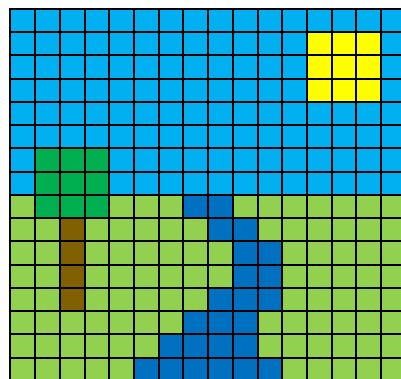
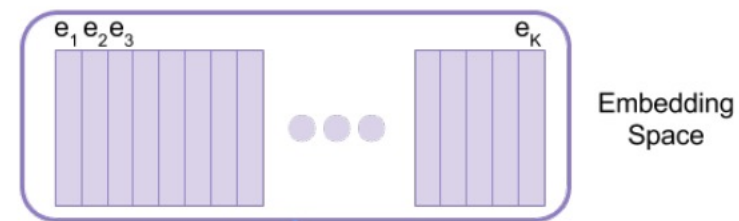
*Figure 10.* Illustration of the embedding scheme for a hypothetical version of our transformer with a maximum text length of 6 tokens. Each box denotes a vector of size  $d_{\text{model}} = 3968$ . In this illustration, the caption has a length of 4 tokens, so 2 padding tokens are used (as described in Section 2.2). Each image vocabulary embedding is summed with a row and column embedding.

DALL-E in practice

Input: Sentence A

Text BPE vocabulary of size  
**16384**

Image learned vocabulary of size  
**8192**



## VQ-VAE

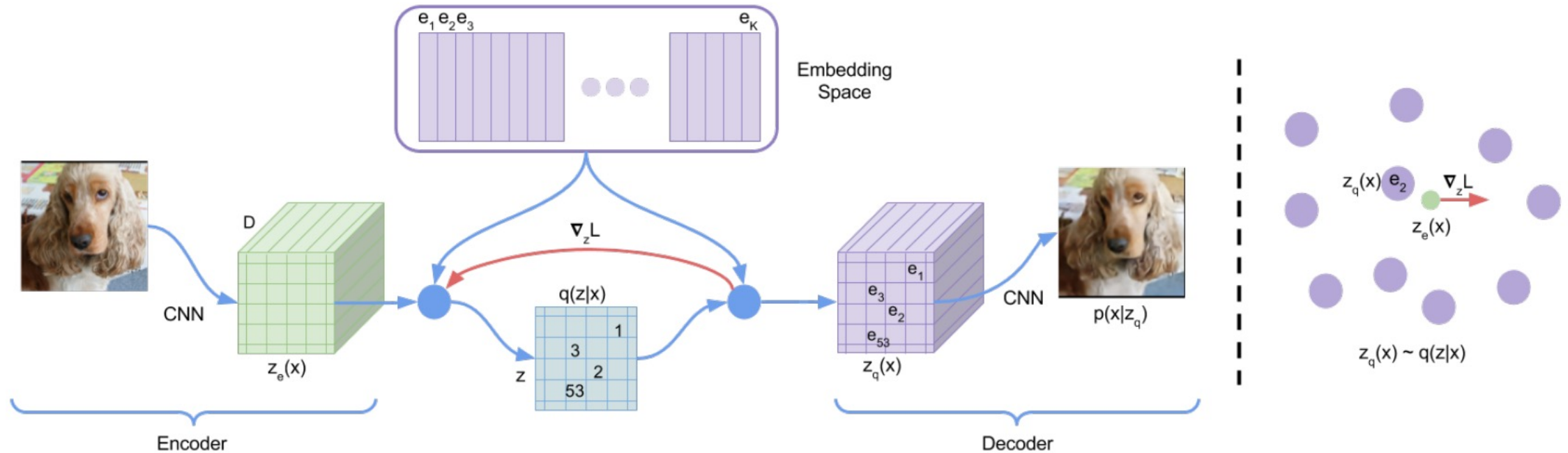
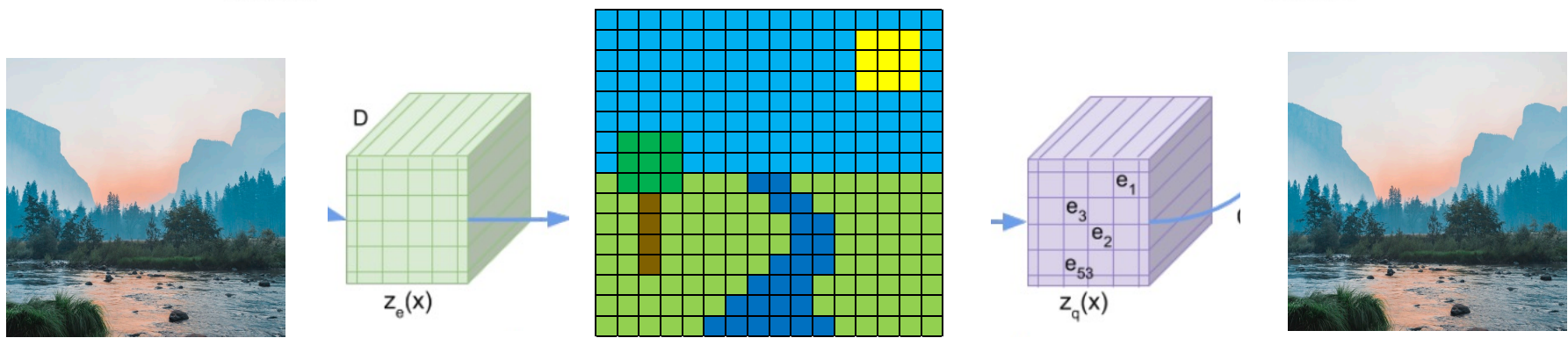
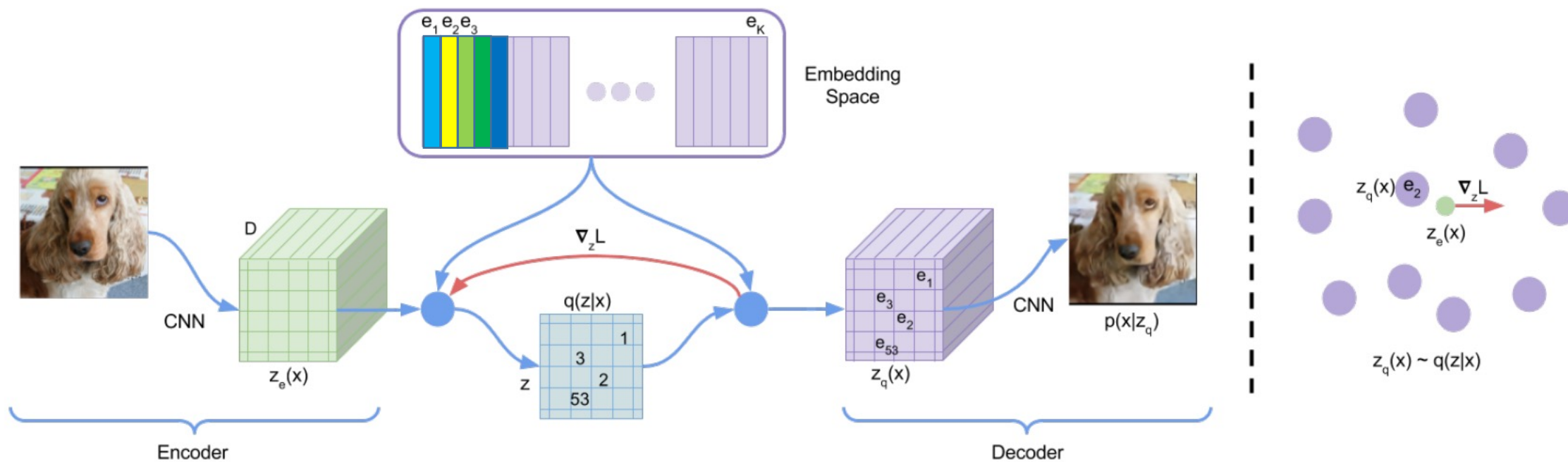


Figure 1: Left: A figure describing the VQ-VAE. Right: Visualisation of the embedding space. The output of the encoder  $z(x)$  is mapped to the nearest point  $e_2$ . The gradient  $\nabla_z L$  (in red) will push the encoder to change its output, which could alter the configuration in the next forward pass.

VQ-VAE



# Diffusion Models

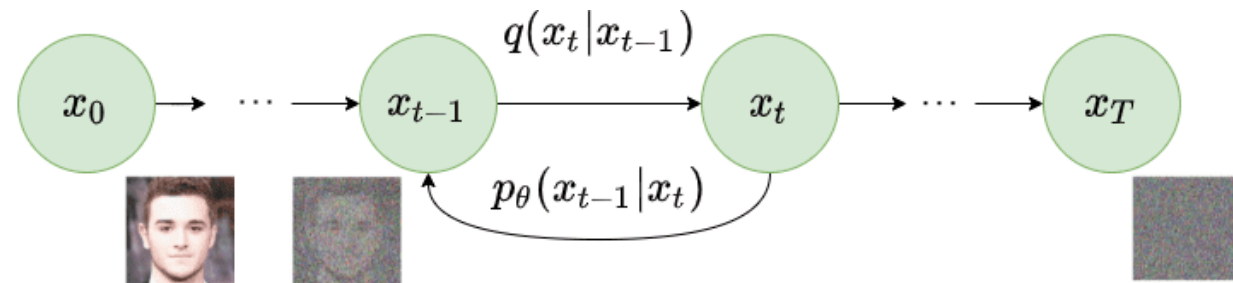
## Diffusion Models

### Models:

- Diffusion models have emerged as a new [state-of-the-art family of deep generative models](#).
- They broke the long domination of [Generative Adversary Networks \(GANs\)](#) and [Transformers](#).
- Over the past year 2022, numerous architectures have been developed:
  - Glide [1],
  - Dalle2 [2],
  - Imagen [3]
  - StableDiffusion [4]
    - v 1.0
    - v 2.0, 2.1

### Principles:

- Diffusion models aim to decompose the image generation process (sampling) in many small “denoising” steps.
- In forward process, diffusion models take the input image and gradually add Gaussian noise to it through a series of T steps (Markov chain).
- Neural network is then trained to recover the original data by reversing the noising process.



# Stable Diffusion



Stable Diffusion – What are the key components?

### ClipText

- Text encoder used to translate prompt into information guiding the process of noise removal
- Input: text.
- Output: 77 token embeddings vectors, each in 768 dimensions.

### VAE

- Image Encoder/Decoder – runs once at the end of the process to create final image in pixel space
- Input / Output: The processed information array (dimensions: (4,64,64))
- Output / Input: The resulting image (dimensions: (3, 512, 512) which are (red/green/blue, width, height))

### UNet

- Image information creator – runs iteratively for N steps to produce the result. It works in the latent space (image information space)
- Input: text embeddings and a starting multi-dimensional array (structured lists of numbers, also called a *tensor*) made up of noise.
- Output: A processed information array

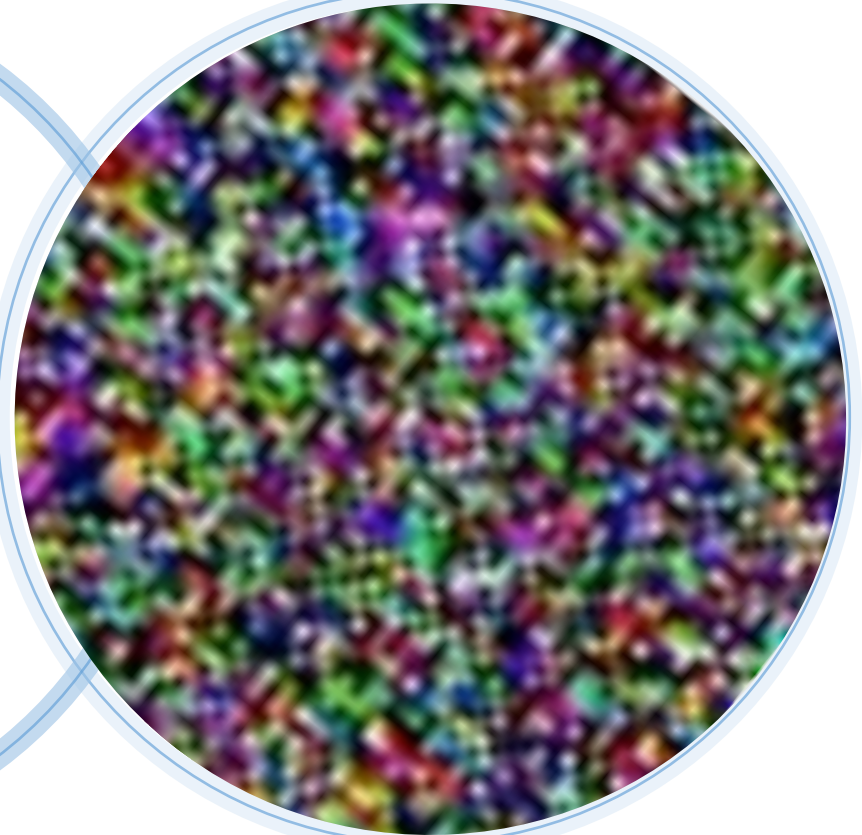
Intuition behind each component

Prompt: "a handsome cat  
dressed like Lincoln,  
trending art."

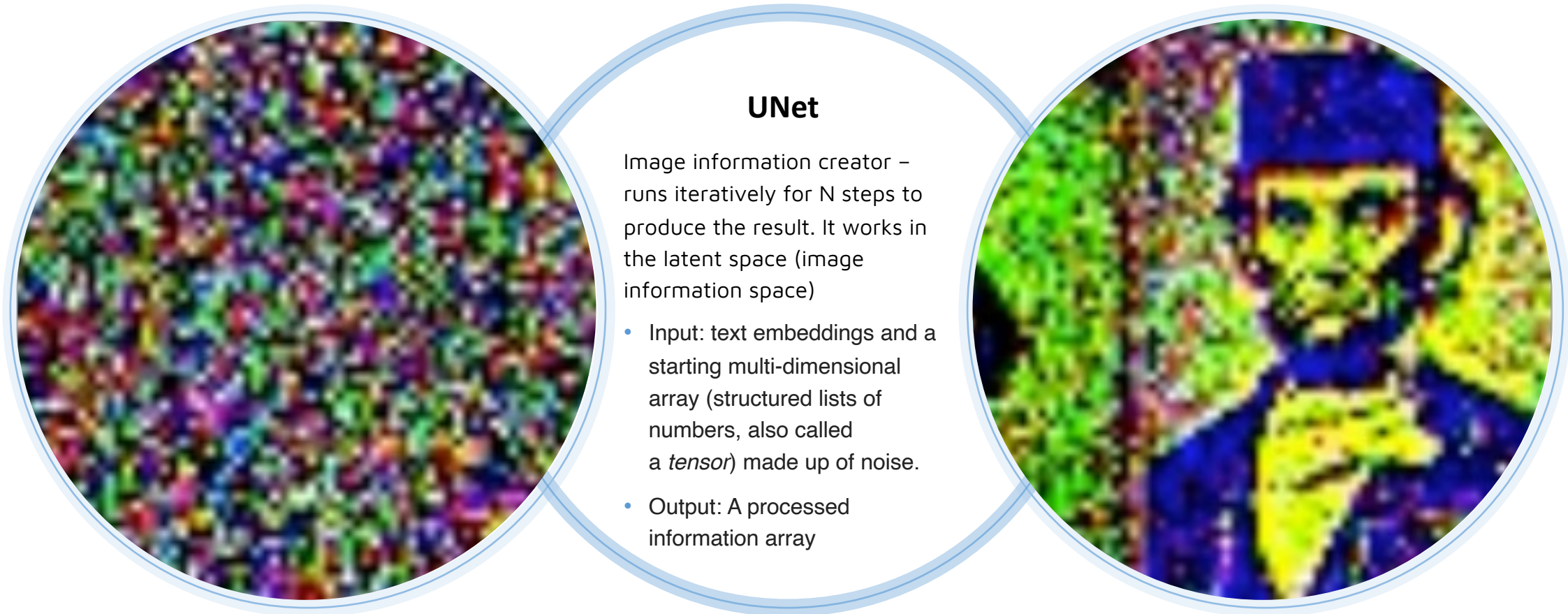
### ClipText

Text encoder used to translate  
prompt into information  
guiding the process of noise  
removal

- Input: text.
- Output: 77 token embeddings vectors, each in 768 dimensions.



Intuition behind each component



Intuition behind each component



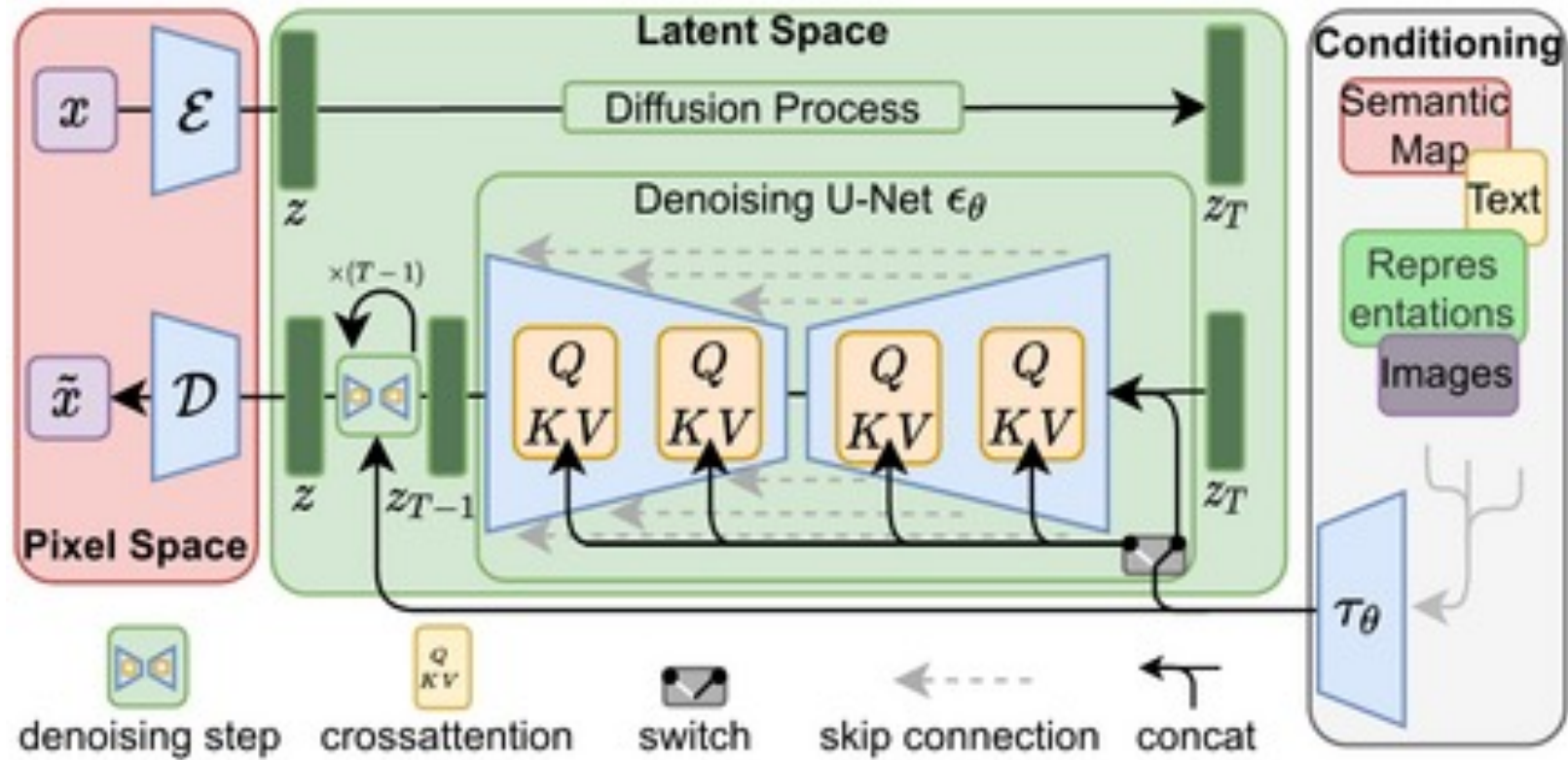
### VAE

Image Encoder/Decoder – runs once at the end of the process to create final image in pixel space

- Input / Output: The processed information array (dimensions: (4,64,64))
- Output / Input: The resulting image (dimensions: (3, 512, 512) which are (red/green/blue, width, height))

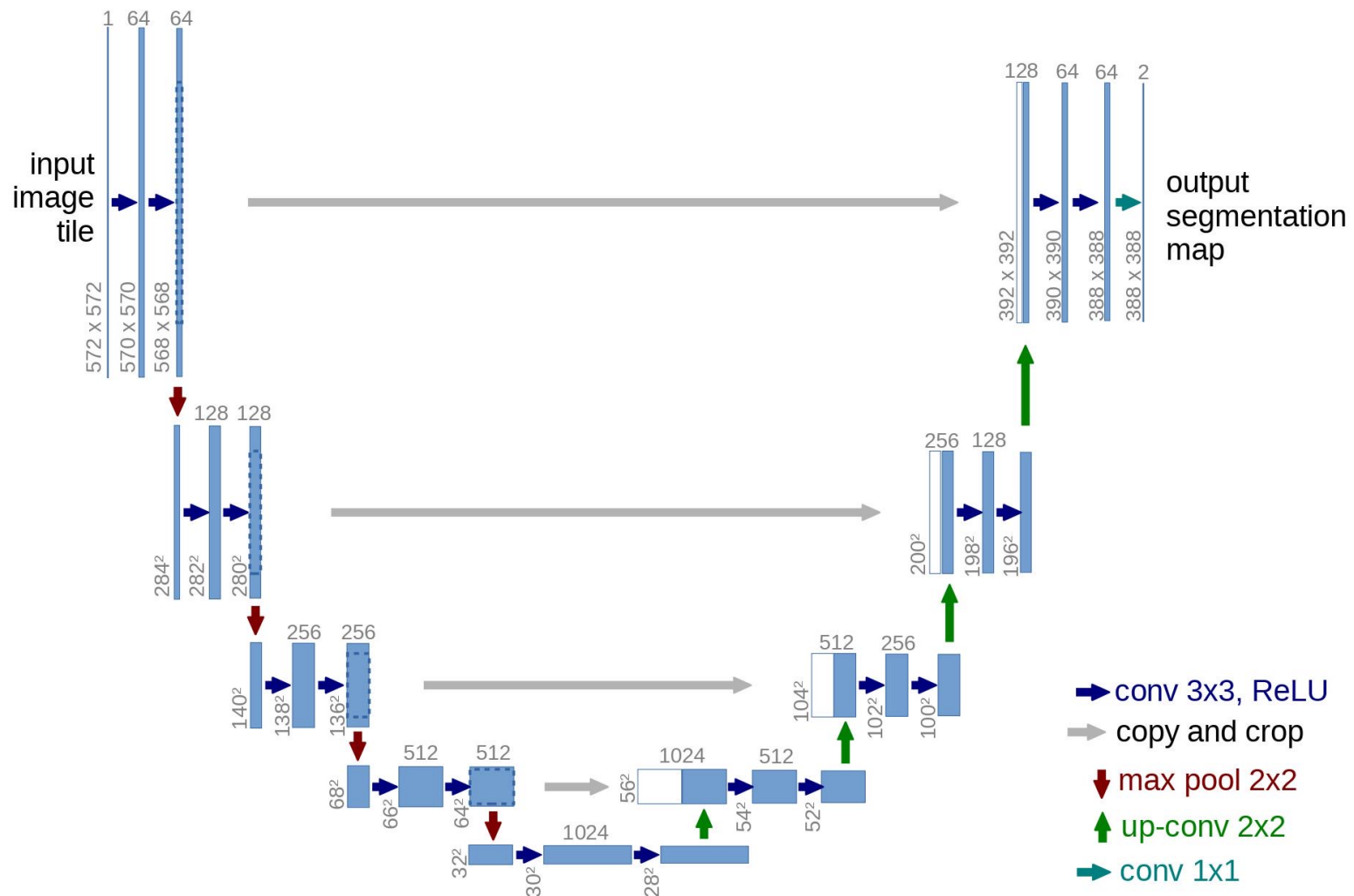


Stable Diffusion

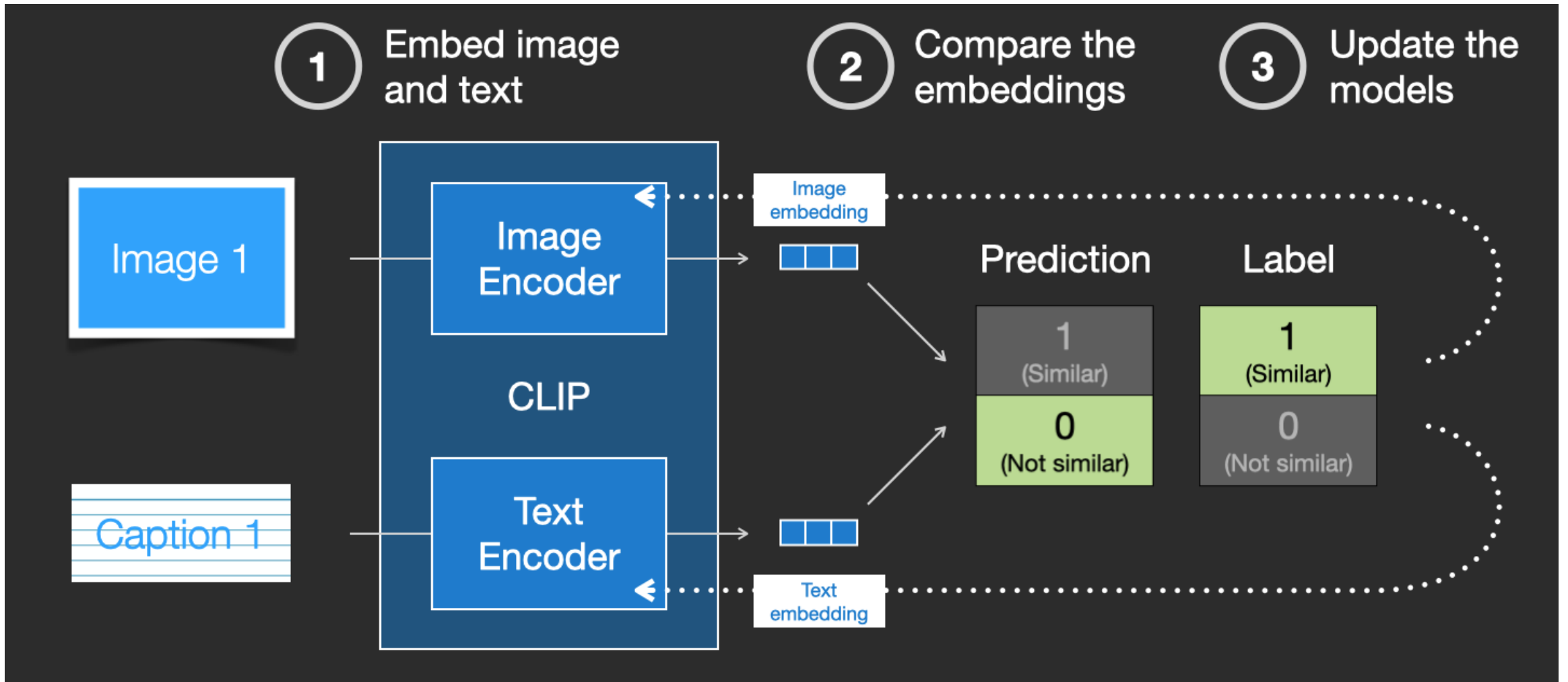


Source: Rombach & Blattmann, et al. "High-Resolution Image Synthesis with Latent Diffusion Models." CVPR 2022.

Denoising UNet

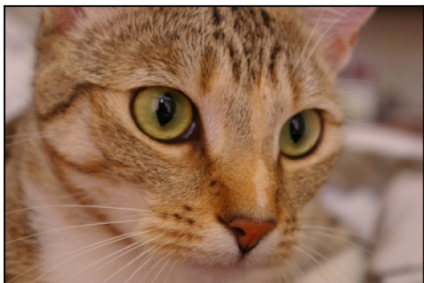


CLIP text model



## CLIP text model

chelsea.png  
a facial photo of a tabby cat



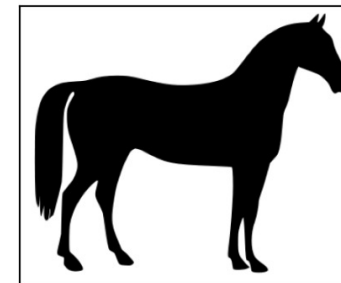
rocket.jpg  
a rocket standing on a launchpad



camera.png  
a person looking at a camera on a tripod



horse.png  
a black-and-white silhouette of a horse



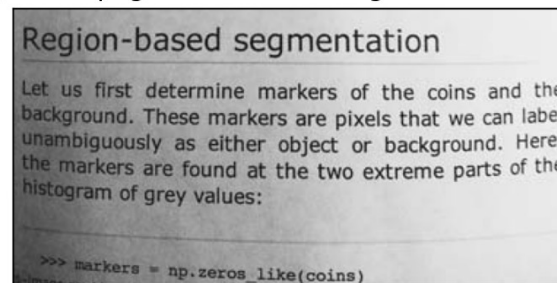
coffee.png  
a cup of coffee on a saucer



astronaut.png  
a portrait of an astronaut with the American flag



page.png  
a page of text about segmentation



motorcycle\_right.png  
a red motorcycle standing in a garage



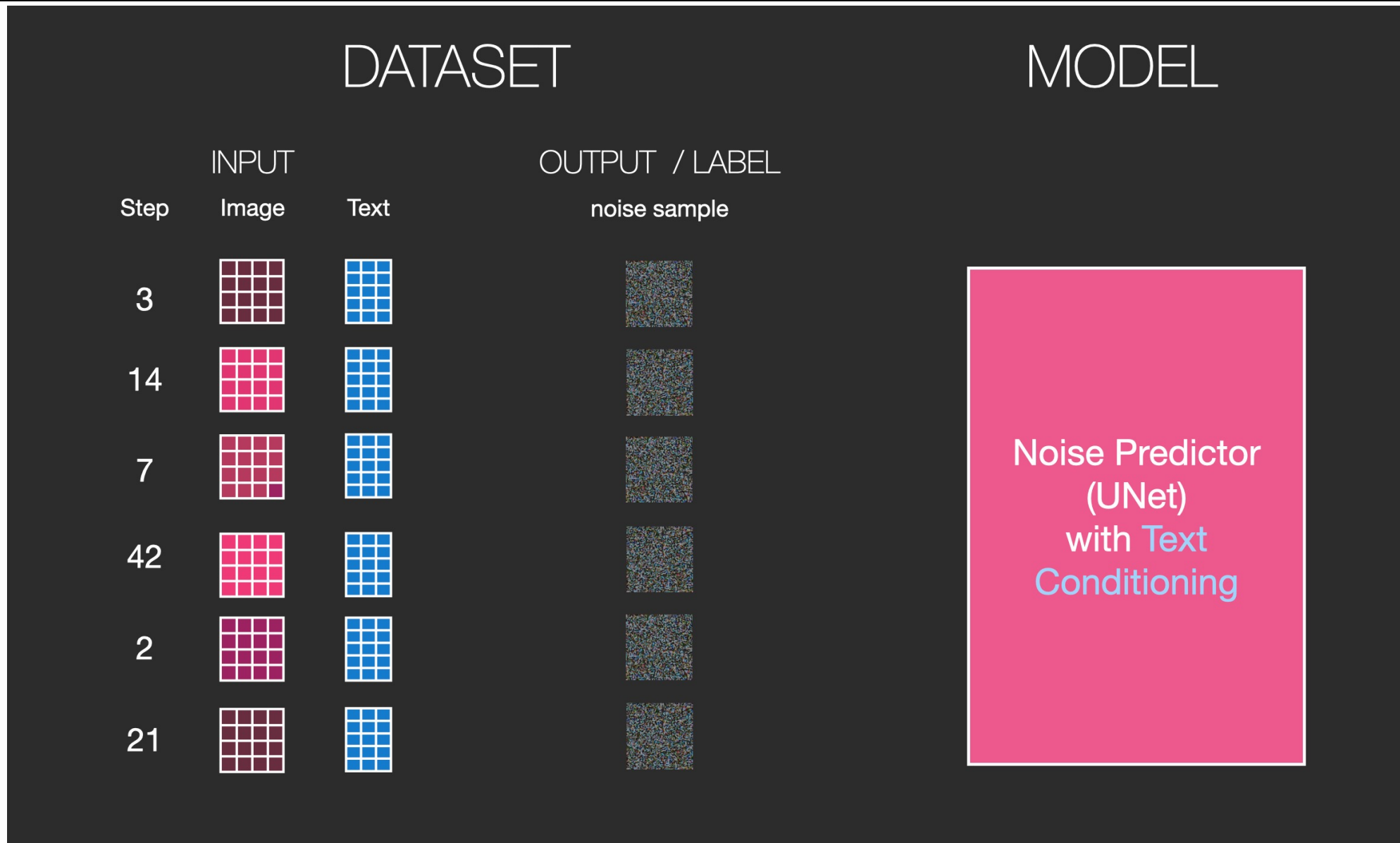


## CLIP text model



a facial photo of a tabby cat	0.31	0.12	0.16	0.15	0.17	0.12	0.12	0.12
a rocket standing on a launchpad	0.18	0.30	0.20	0.17	0.14	0.19	0.17	0.16
a person looking at a camera on a tripod	0.21	0.21	0.30	0.20	0.14	0.19	0.19	0.16
a black-and-white silhouette of a horse	0.15	0.15	0.21	0.35	0.15	0.11	0.17	0.17
a cup of coffee on a saucer	0.18	0.12	0.17	0.15	0.29	0.15	0.14	0.12
a portrait of an astronaut with the American flag	0.17	0.22	0.17	0.16	0.15	0.28	0.13	0.15
a page of text about segmentation	0.20	0.16	0.20	0.20	0.20	0.15	0.35	0.16
a red motorcycle standing in a garage	0.15	0.16	0.12	0.16	0.13	0.15	0.14	0.32

CLIP text model



Encoder / Decoder

Training:

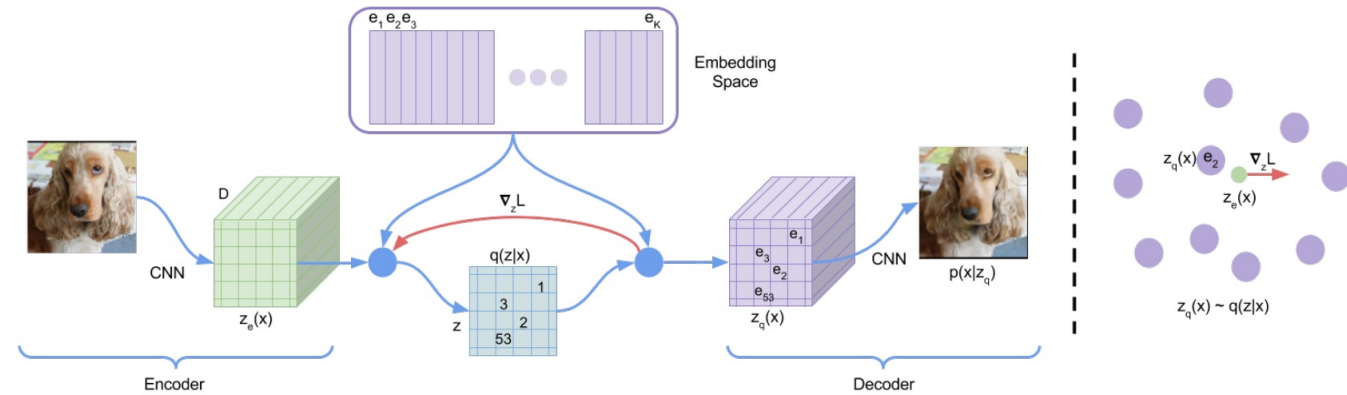
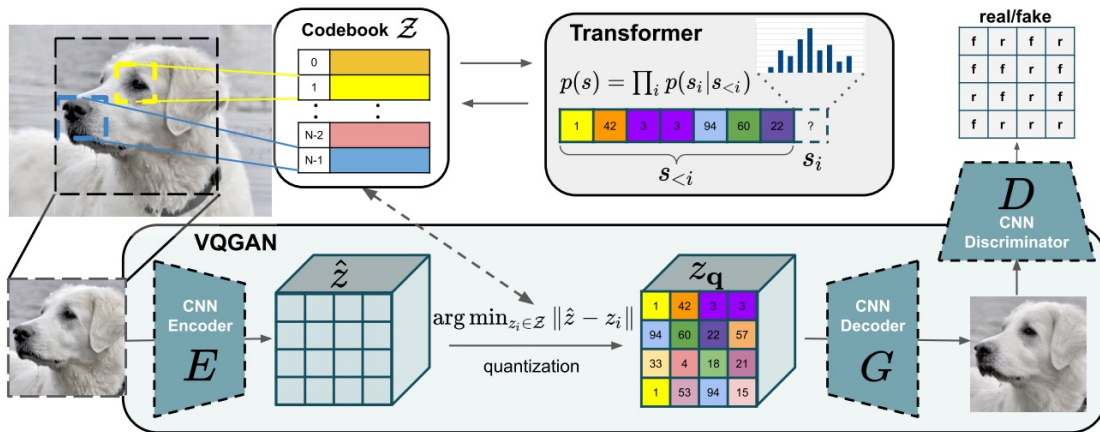


Figure 1: Left: A figure describing the VQ-VAE. Right: Visualisation of the embedding space. The output of the encoder  $z(x)$  is mapped to the nearest point  $e_2$ . The gradient  $\nabla_z L$  (in red) will push the encoder to change its output, which could alter the configuration in the next forward pass.

SD usage:

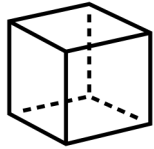
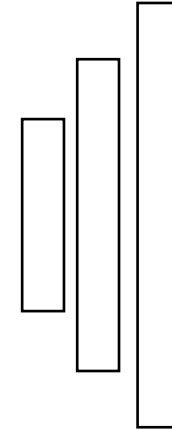
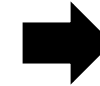
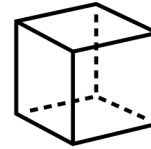
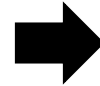
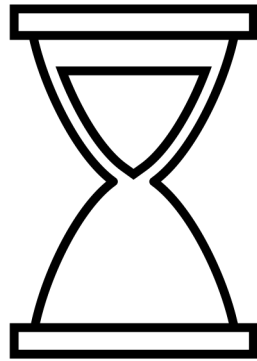
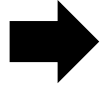


$x$  is (3, 512, 512) image tensor  
 $z$  is (4, 64, 64) latent tensor



Inference pipeline

Text prompt



We start in the latent space  $z$ . A random noise of dimension  $(4, 64, 64)$  is generated. This is controlled by seed.

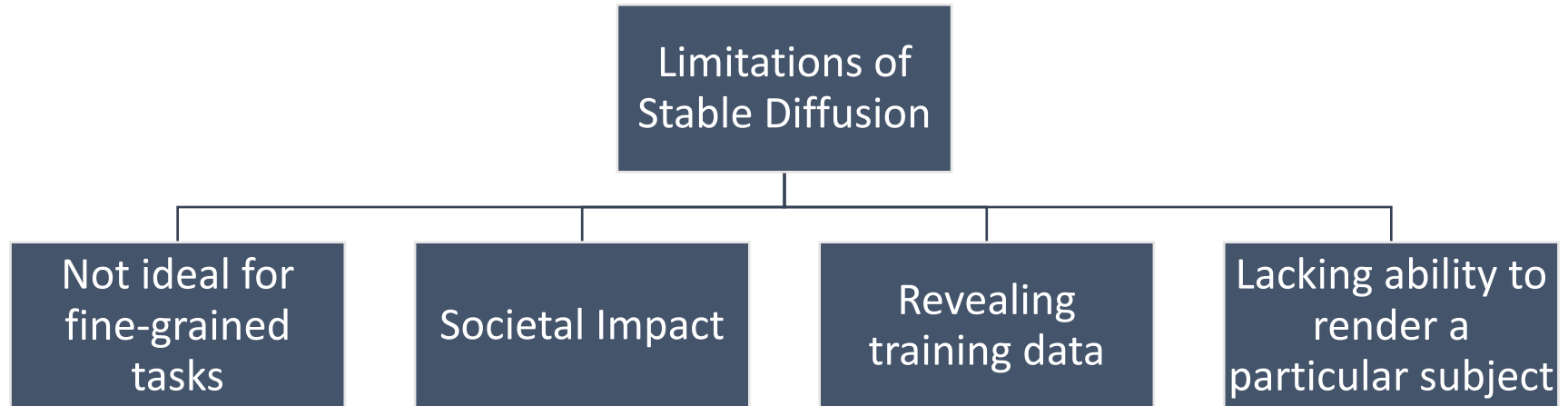
Unet takes latent noise and text prompt as input and predicts noise in latent space.

Latent space with noise iteratively subtracted

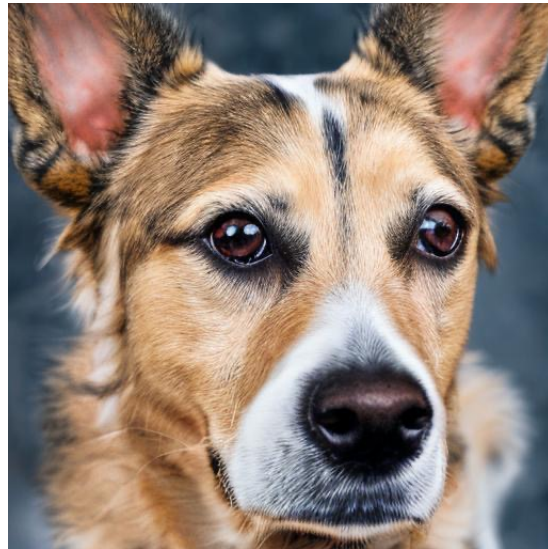
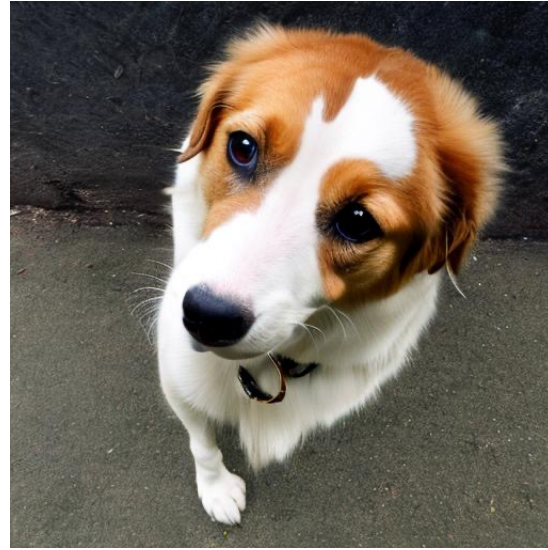
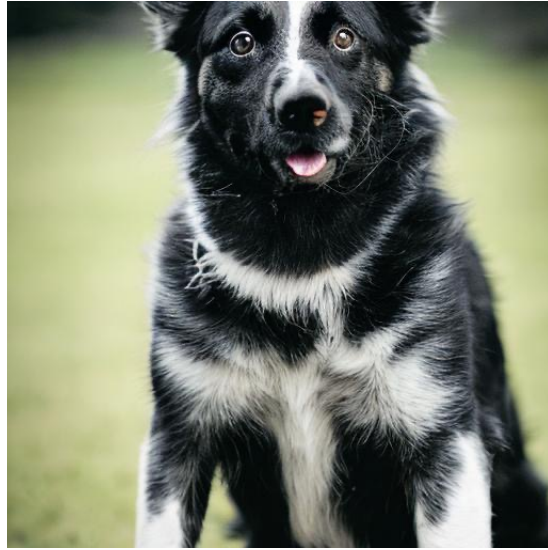
VEA decoder

# Limitations

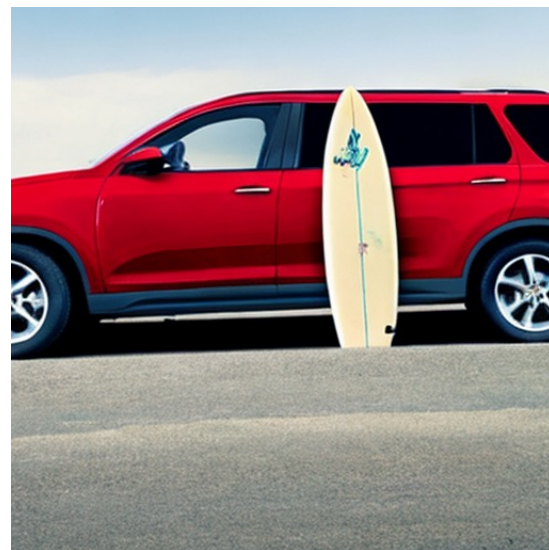
## Limitations



Exploration – a dog



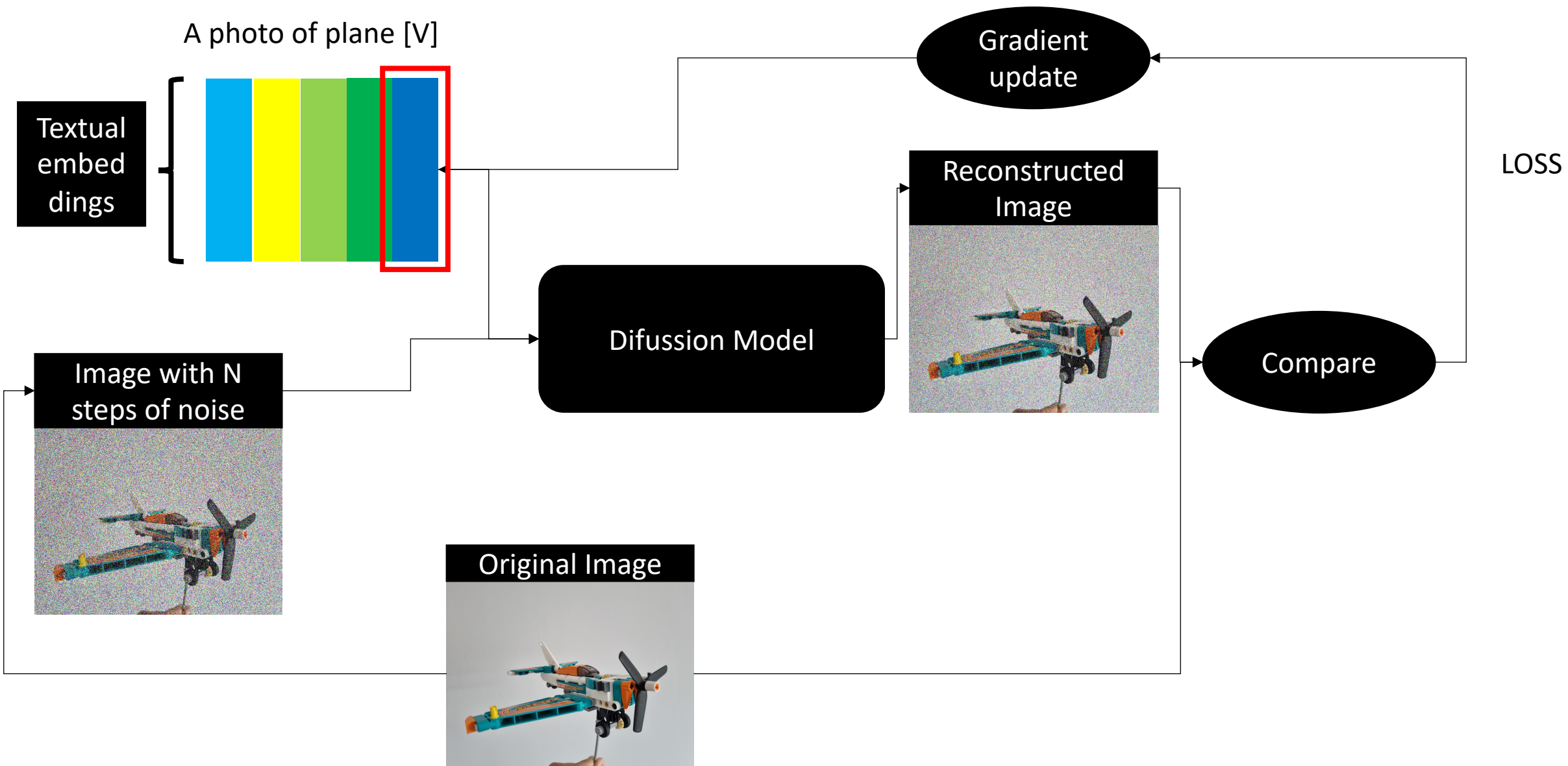
Adaptation – “A red SUV car with a surfing board painting on the side”



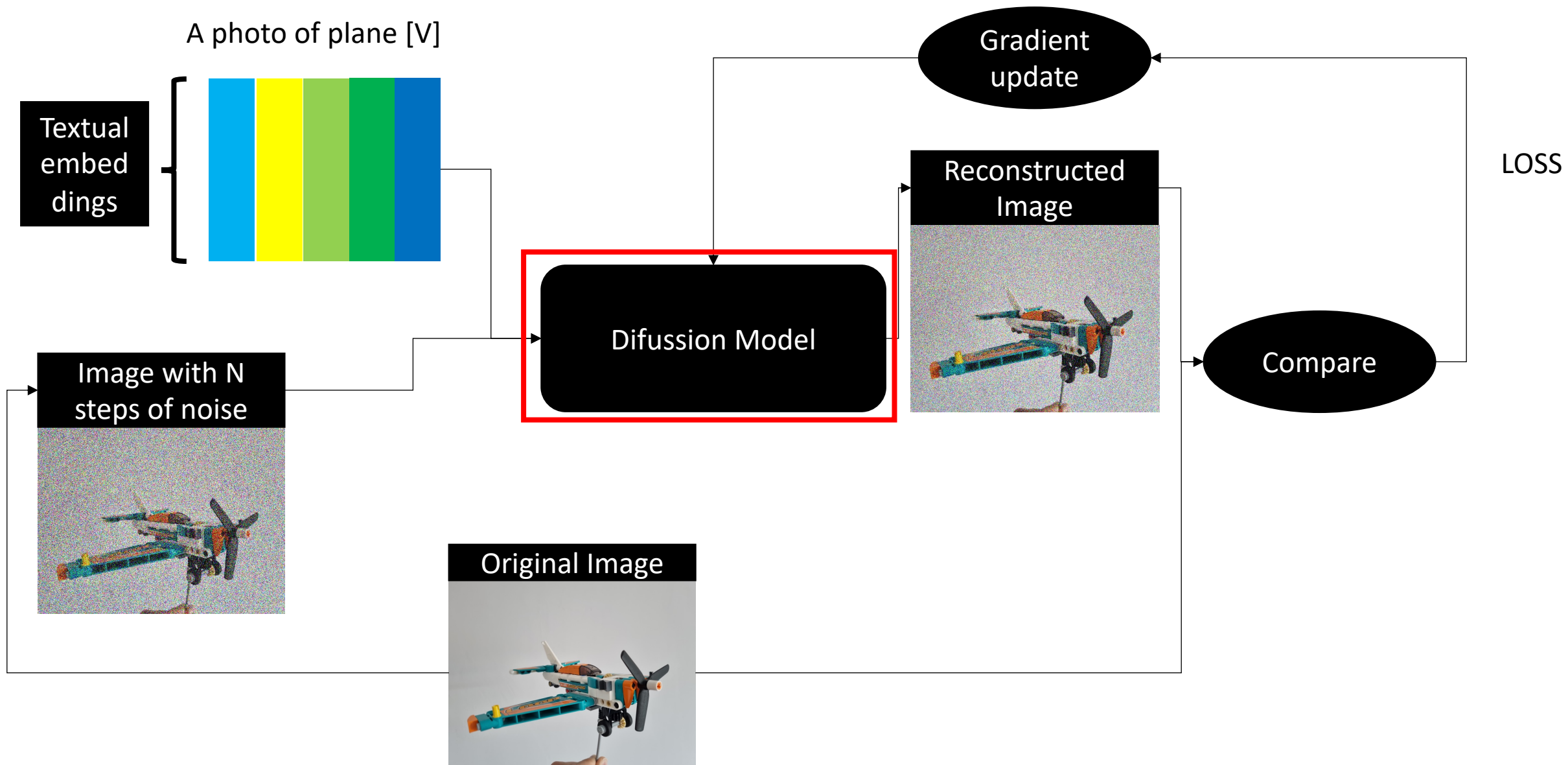


# Fine-tuning approach

Textual Inversion



DreamBooth

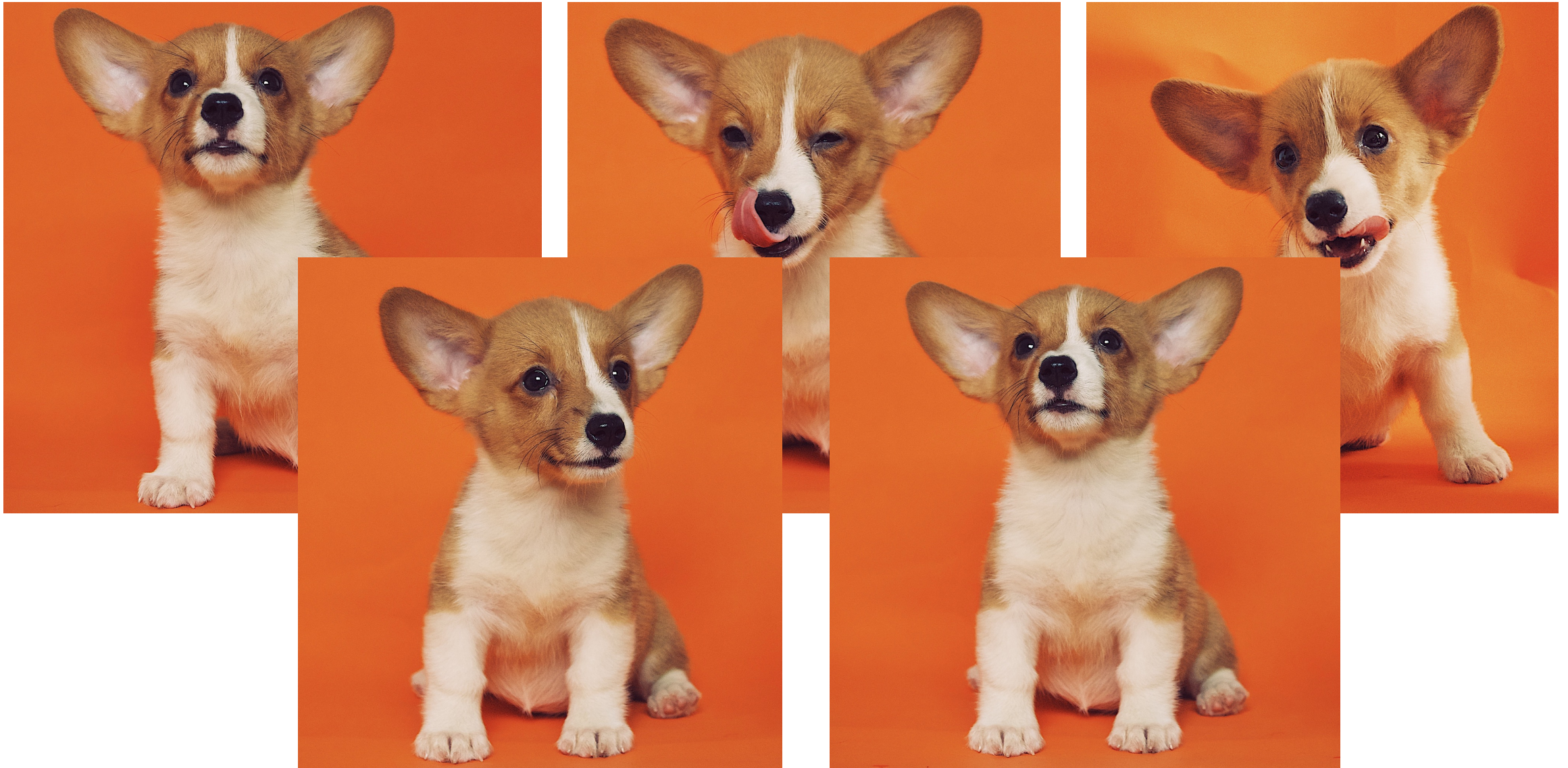


## Textual Inversion vs DreamBooth

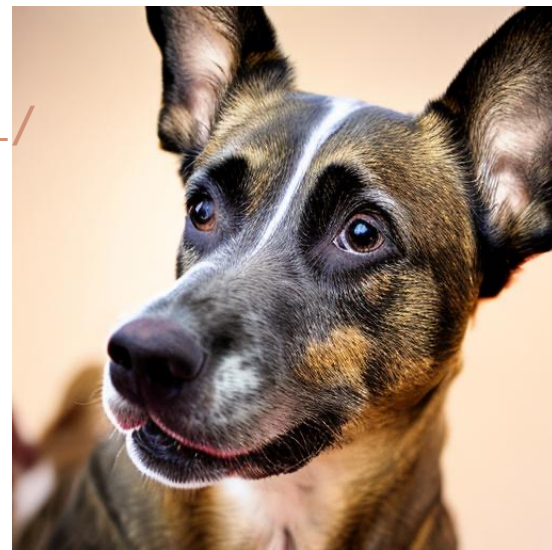
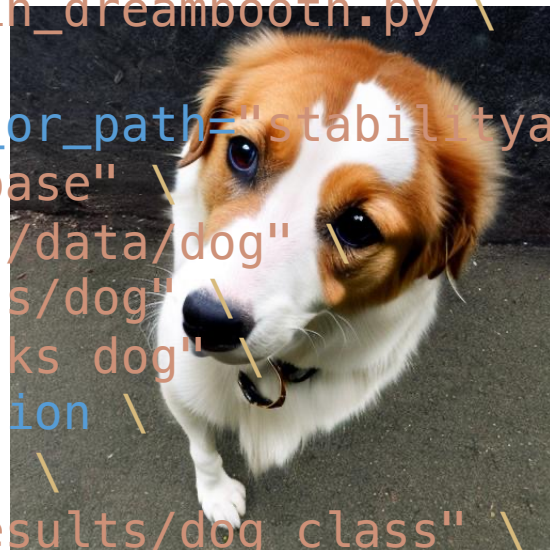
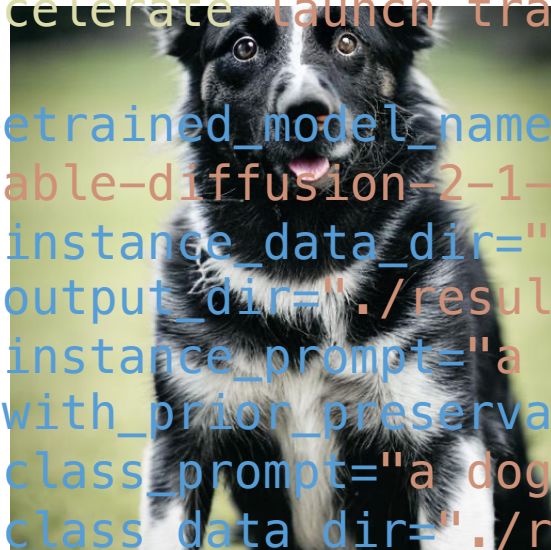
	<b>Textual Inversion</b>	<b>DreamBooth</b>
Effectiveness	Quite good	Subjectively better than Textual Inversion
Storage	Efficient when it comes to storage management since base model is not changed and it is only the small vector representing the learned concept that needs to be stored.	For every new concept new big model files is created, so in case many concepts need to be learned we end up with bunch of models.

# Example – DreamBooth

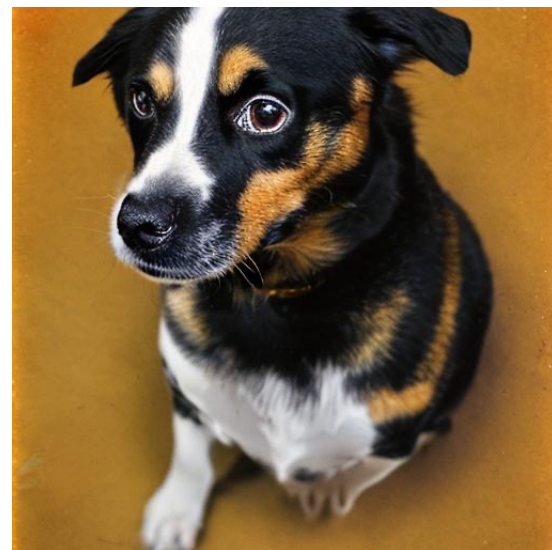
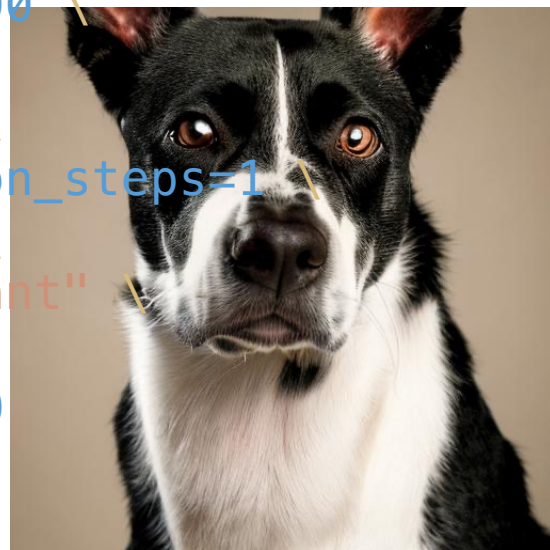
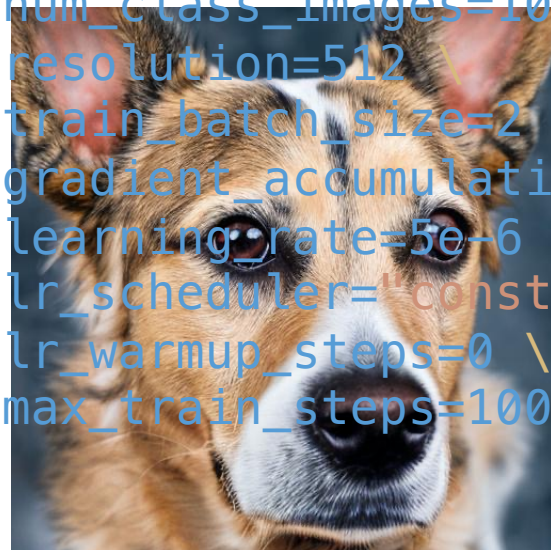
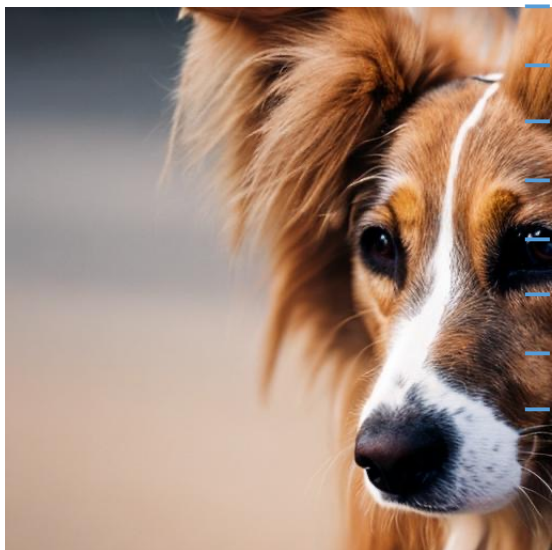
DreamBooth – experiment 2



DreamBooth – experiment 2



```
accelerate launch train_dreambooth.py \  
--pretrained_model_name_or_path="stabilityai/\  
stable-diffusion-2-1-base" \  
--instance_data_dir="./data/dog" \  
--output_dir="./results/dog" \  
--instance_prompt="a sks dog" \  
--with_prior_preservation \  
--class_prompt="a dog" \  
--class_data_dir="./results/dog_class" \  
--num_class_images=1000 \  
--resolution=512 \  
--train_batch_size=2 \  
--gradient_accumulation_steps=1 \  
--learning_rate=5e-6 \  
--lr_scheduler="constant" \  
--lr_warmup_steps=0 \  
--max_train_steps=1000
```



## DreamBooth – experiment 2

```
accelerate launch train_dreambooth.py \  
--  
pretrained_model_name_or_path="stabilityai/  
stable-diffusion-2-1-base" \  
--instance_data_dir="./data/dog" \  
--output_dir="./results/dog" \  
--instance_prompt="a sks dog" \  
--with_prior_preservation \  
--class_prompt="a dog" \  
--class_data_dir="./results/dog_class" \  
--num_class_images=1000 \  
--resolution=512 \  
--train_batch_size=2 \  
--gradient_accumulation_steps=1 \  
--learning_rate=5e-6 \  
--lr_scheduler="constant" \  
--lr_warmup_steps=0 \  
--max_train_steps=1000
```

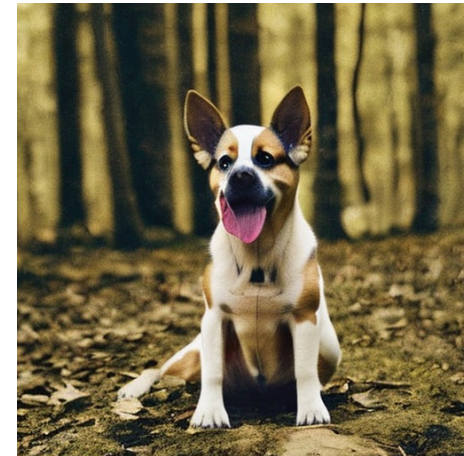
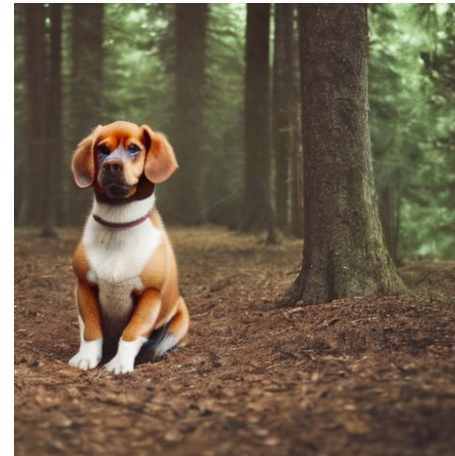


DreamBooth – experiment 2

**Prompt:** SKS dog in a forest



**Prompt:** dog in a forest



DreamBooth – experiment 2

**Prompt:** SKS dog over pyramids



**Prompt:** dog over pyramids

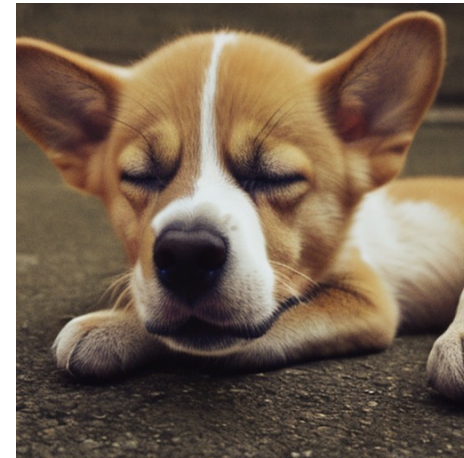
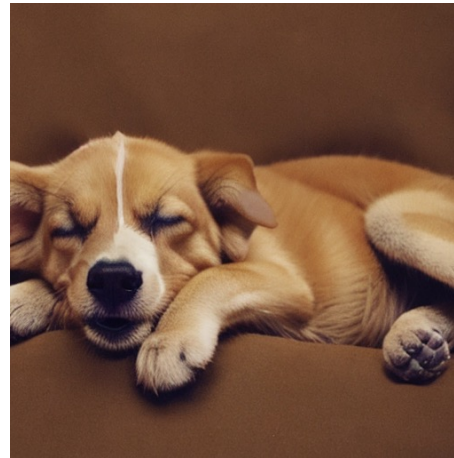


DreamBooth – experiment 2

**Prompt:** SKS dog underwater



**Prompt:** SKS dog sleeping



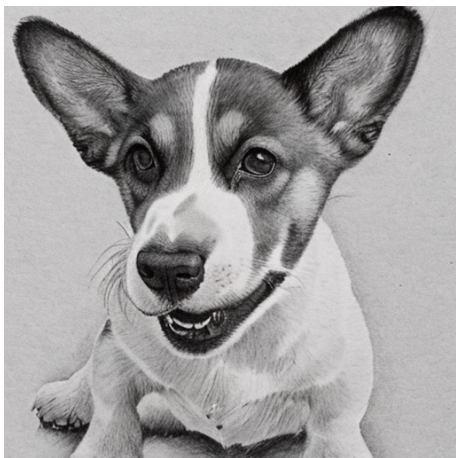
DreamBooth – experiment 2

**Prompt:** SKS dog in red dress



DreamBooth – experiment 2

**Prompt:** sketch of SKS dog



**Prompt:** SKS dog by Picasso



**Prompt:** dog by Picasso

