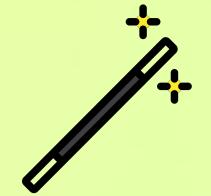




The Text2Image Magic



A high level overview of the technology's building blocks



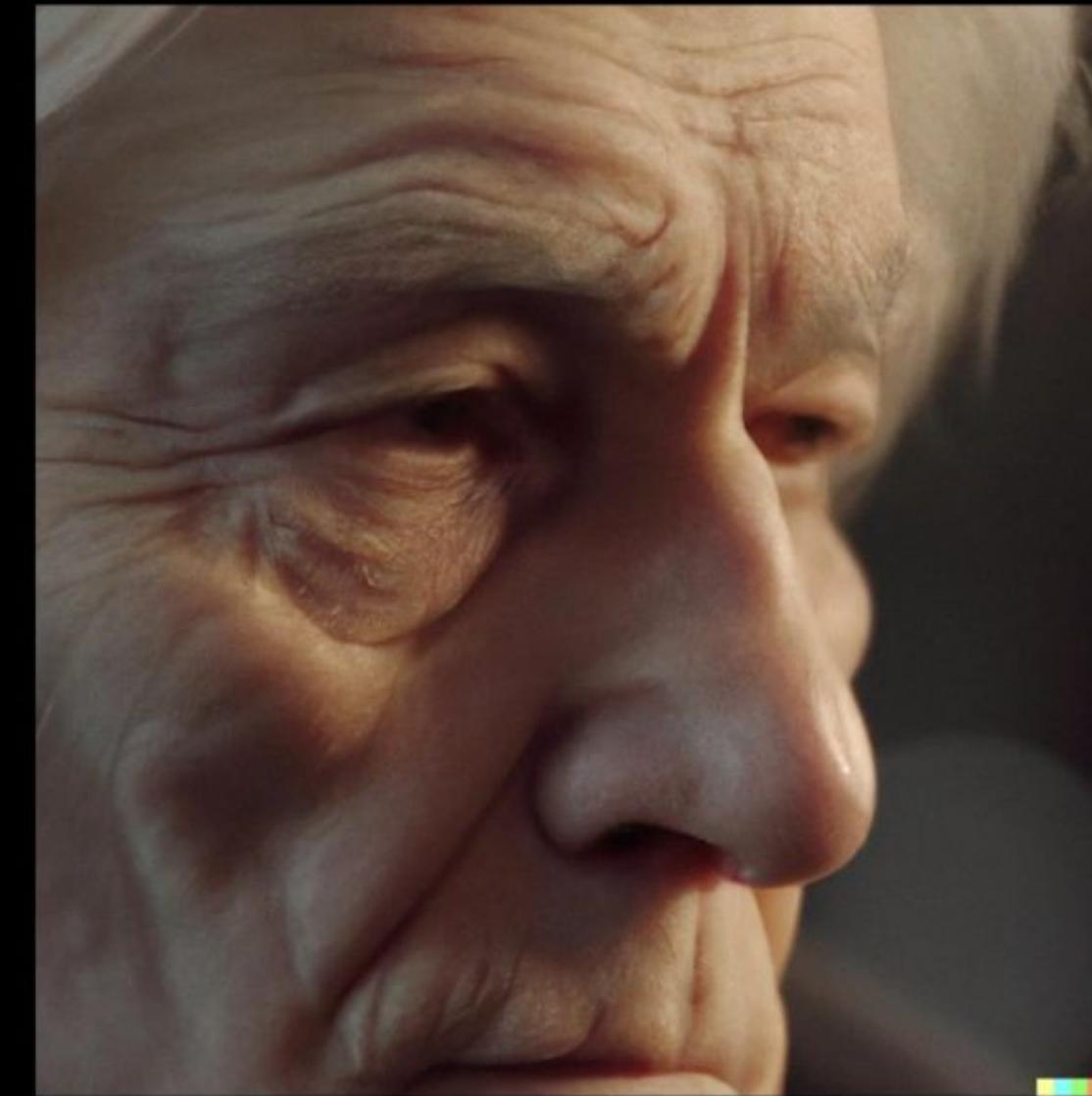
Naomi Ken Korem
Researcher



MIDJOURNEY



DALL-E 2



STABLEDIFFUSION



film still, portrait of an old man, wrinkles, dignified look, grey silver hair, peculiar nose, wise, eternal wisdom and beauty, incredible lighting and camera work, depth of field, bokeh, screenshot from a hollywood movie

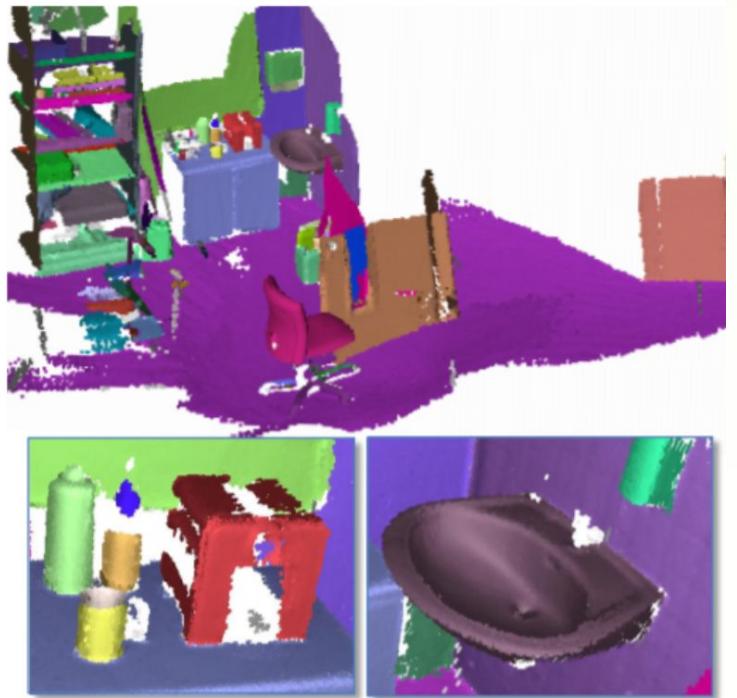


LAION
@laion_ai

...

Guiding Stable Diffusion with our CLIP H:
"Professional HDR photo of a polar bear holding a lollipop on a rooftop in Hong Kong looking up at a UFO in the night sky. A UFO flies above the polar bear. The polar bear holds a lollipop on a rooftop. The background shows Hong Kong."





Naomi Ken Korem

Lighticks Team



Harel Cain

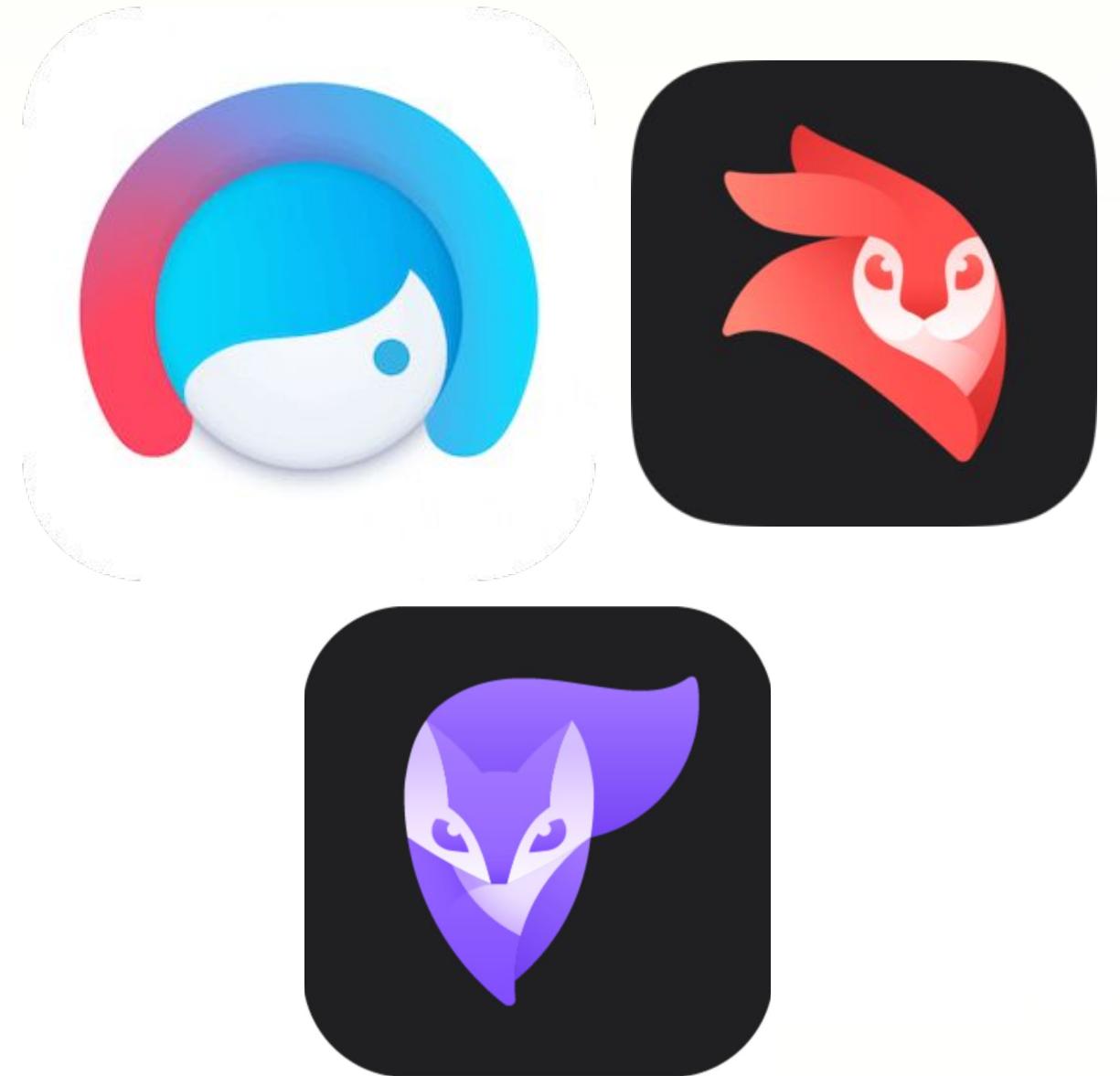


AJ Bruce



Shaked Dunskey

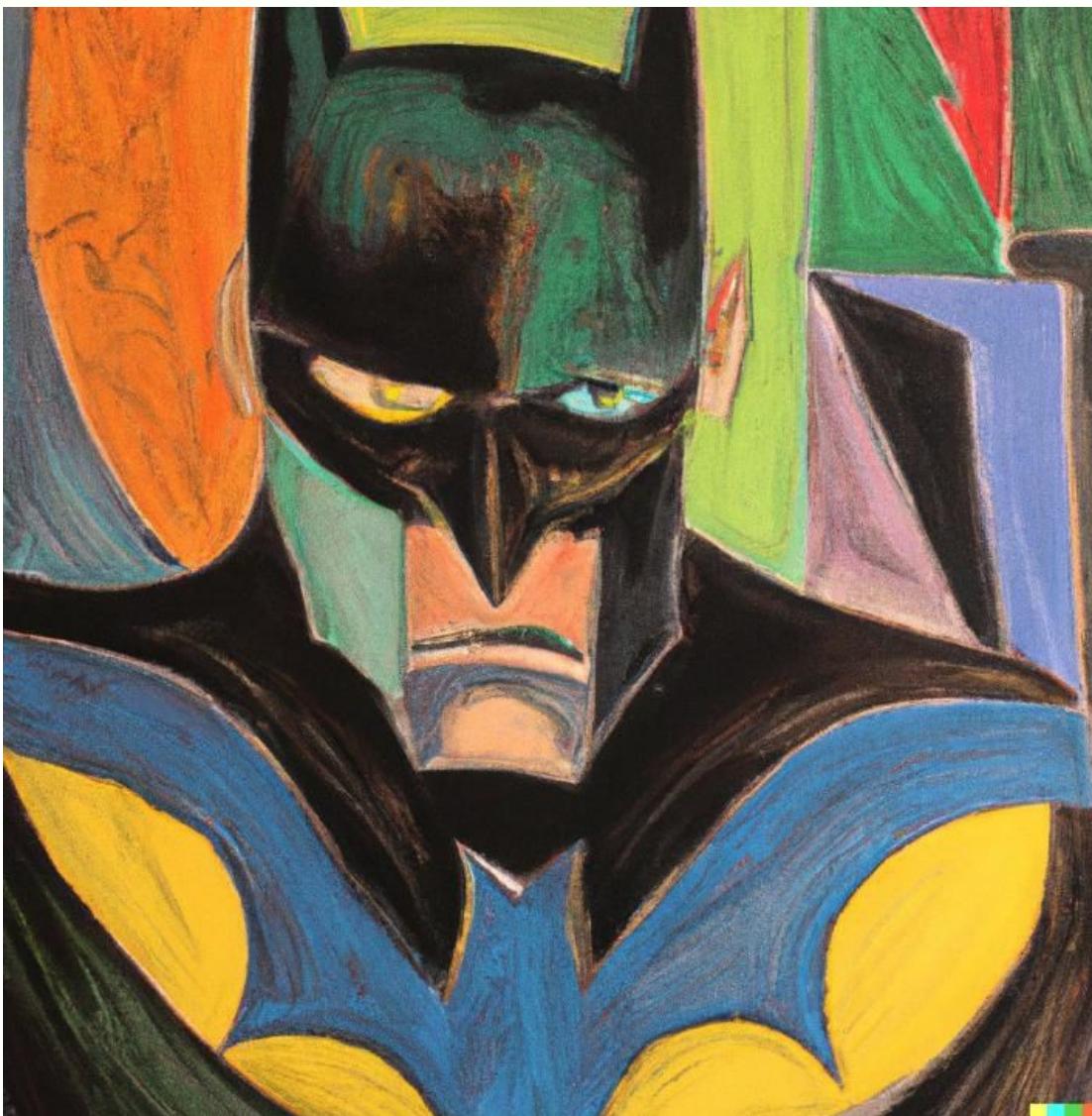
Lighticks



DALLE v2 (Apr 2022)



Academic researcher before a deadline, as painted by Edvard Munch.



Batman, as painted by Picasso.



Yoda angry at having lost money in the stock market.

Imagen (May 2022)

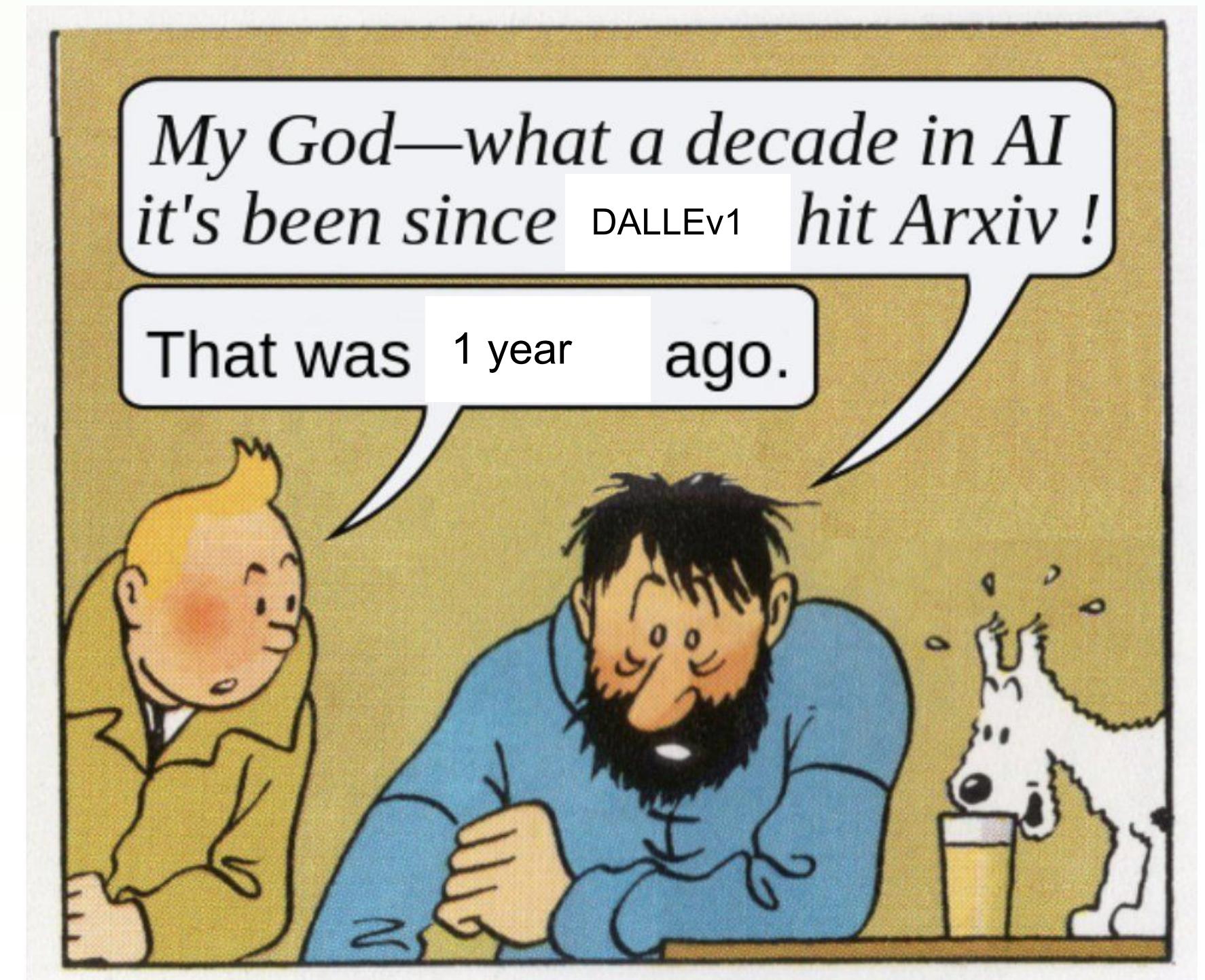


Android mascot made of bamboo.

A dog looking curiously in the mirror, seeing a cat.

A chrome-plated duck with a golden beak arguing with an angry turtle in a forest.

It all
happened so
fast



System 1



Fast, intuitive and
emotional

System 2



Slow, conscious
and effortful

Timeline

Neural style transfer and related domains and StyleGAN.

OpenAI publish DALL-E paper and open source CLIP model

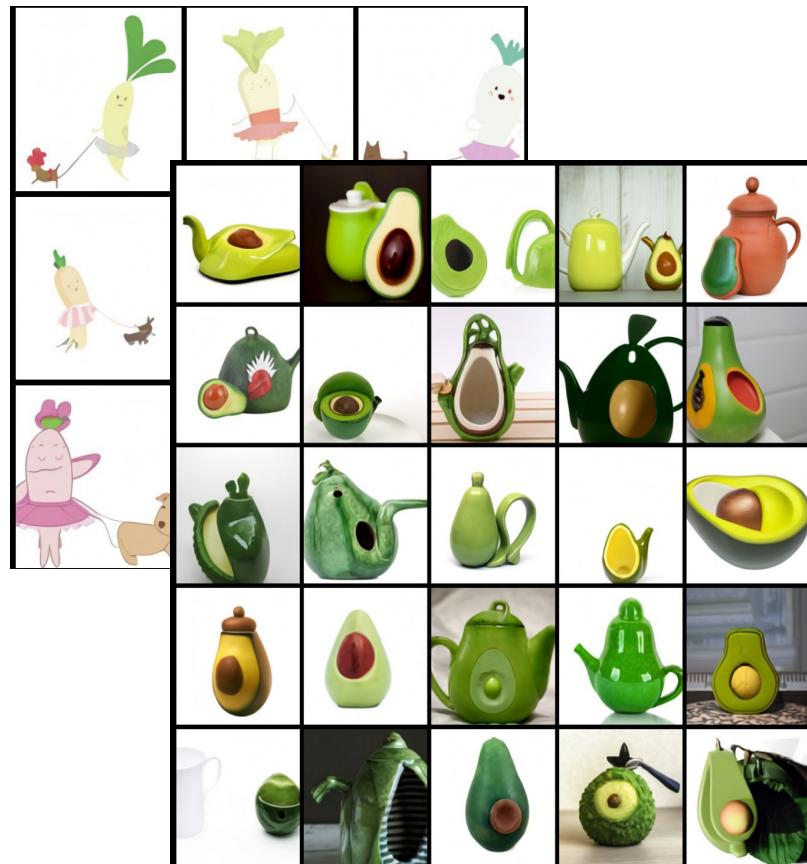
Diffusion models and Gans using CLIP

DALLE-2 is published by OpenAI and hype hits mainstream media

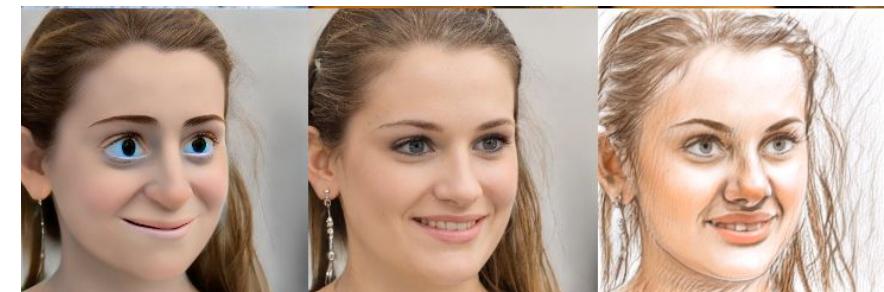
Google Brain releases Imagen

Stability.ai releases Stable Diffusion, open license

Pre 2021



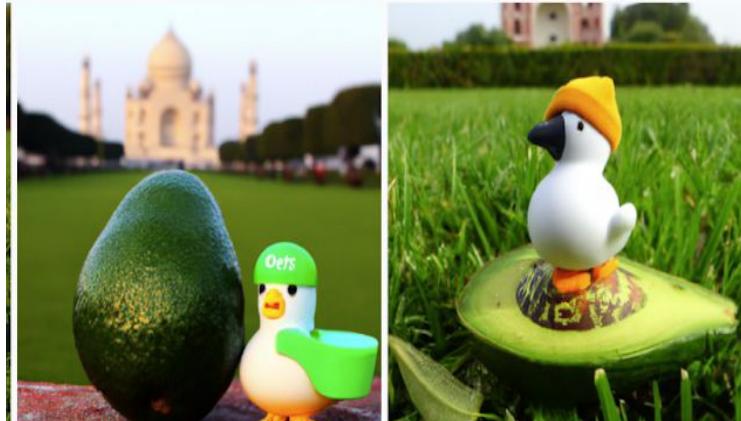
1-2021



2021-2022



04-2022



5-2022



8-2022



What cause this fast progress?

1. **DATA**

5 Billion image, text pairs!

2. **Parallel training in scale -**

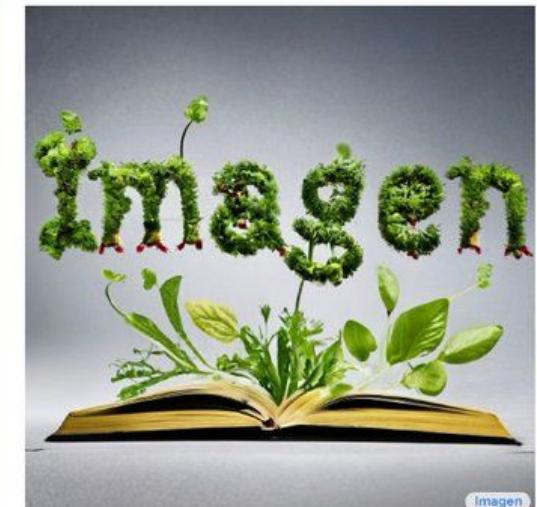
2.5M training steps of batch size 2048 in very short time. (256 TPU-v4 chips)

3. **Multi domain connections**

New architectures, layers (Bert, GPT3, StyleGAN, ViT)

The building blocks that got us there – our agenda for this workshop

1. Generators and StyleGAN
2. Diffusion models
3. Text conditioned diffusion models (Dalle2, glide)
4. Stable Diffusion



Sprouts in the shape of text 'Imagen' coming out of a fairytale book.



A photo of a Shiba Inu dog with a backpack riding a bike. It is wearing sunglasses and a beach hat.



A high contrast portrait of a very happy fuzzy panda dressed as a chef in a high end kitchen making dough. There is a painting of flowers on the wall behind him.



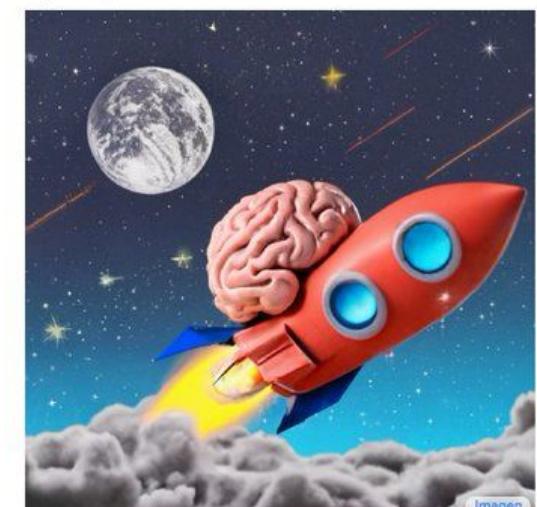
Teddy bears swimming at the Olympics 400m Butterfly event.



A cute corgi lives in a house made out of sushi.



A cute sloth holding a small treasure chest. A bright golden glow is coming from the chest.



A brain riding a rocketship heading towards the moon.



A dragon fruit wearing karate belt in the snow.



A strawberry mug filled with white sesame seeds. The mug is floating in a dark chocolate sea.

Generative Adversarial Network

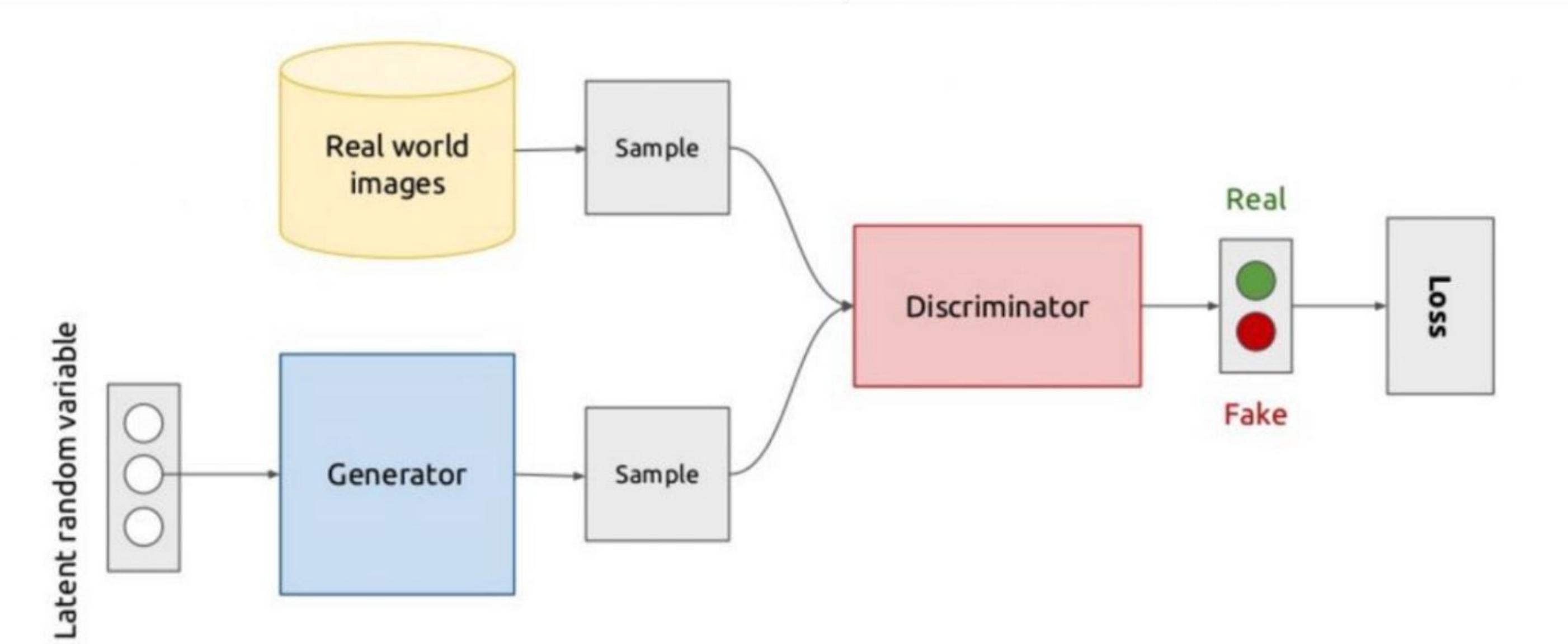
Input: random noise vector (dimensions - 1xd)

Output: Image (dimensions wxhx3)

GAN

(Generative adversarial network)

Back in 2014 argument in a bar led Ian Goodfellow to the idea of Generative Adversarial Network. Now many researchers believe that GAN is



or any other kind of neural networks. At the same time D needs to have about the same capacity as G or even greater.

[Generative Adversarial Networks](#)

StyleGAN

Nvidia, 2018



Analyzing and Improving the Image Quality of
StyleGAN

Neural StyleTransfer

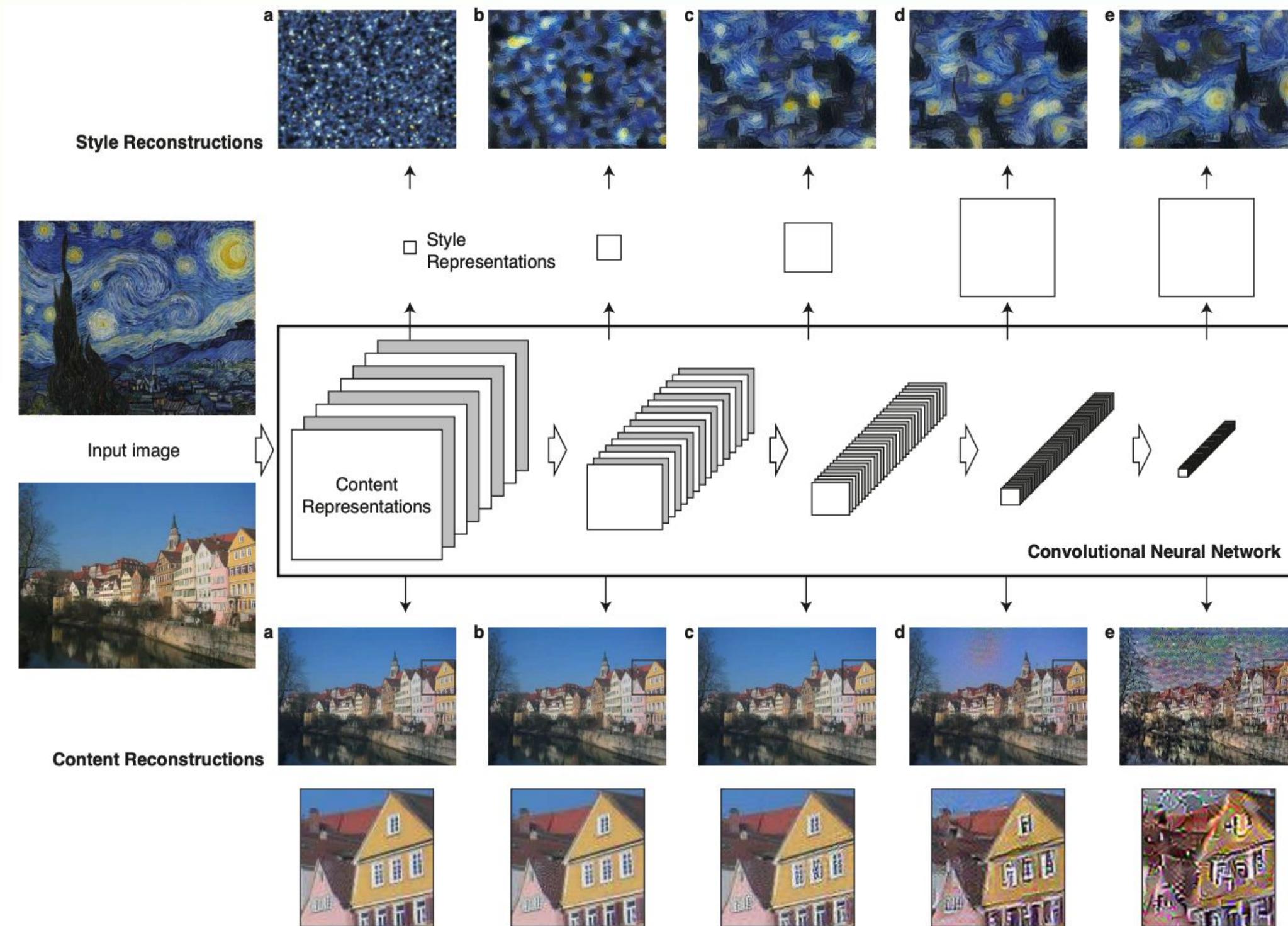
Match the correlation between the different filter response

Gram matrix:

$$G_{ij}^l = \sum_k F_{ik}^l F_{jk}^l.$$

Loss:

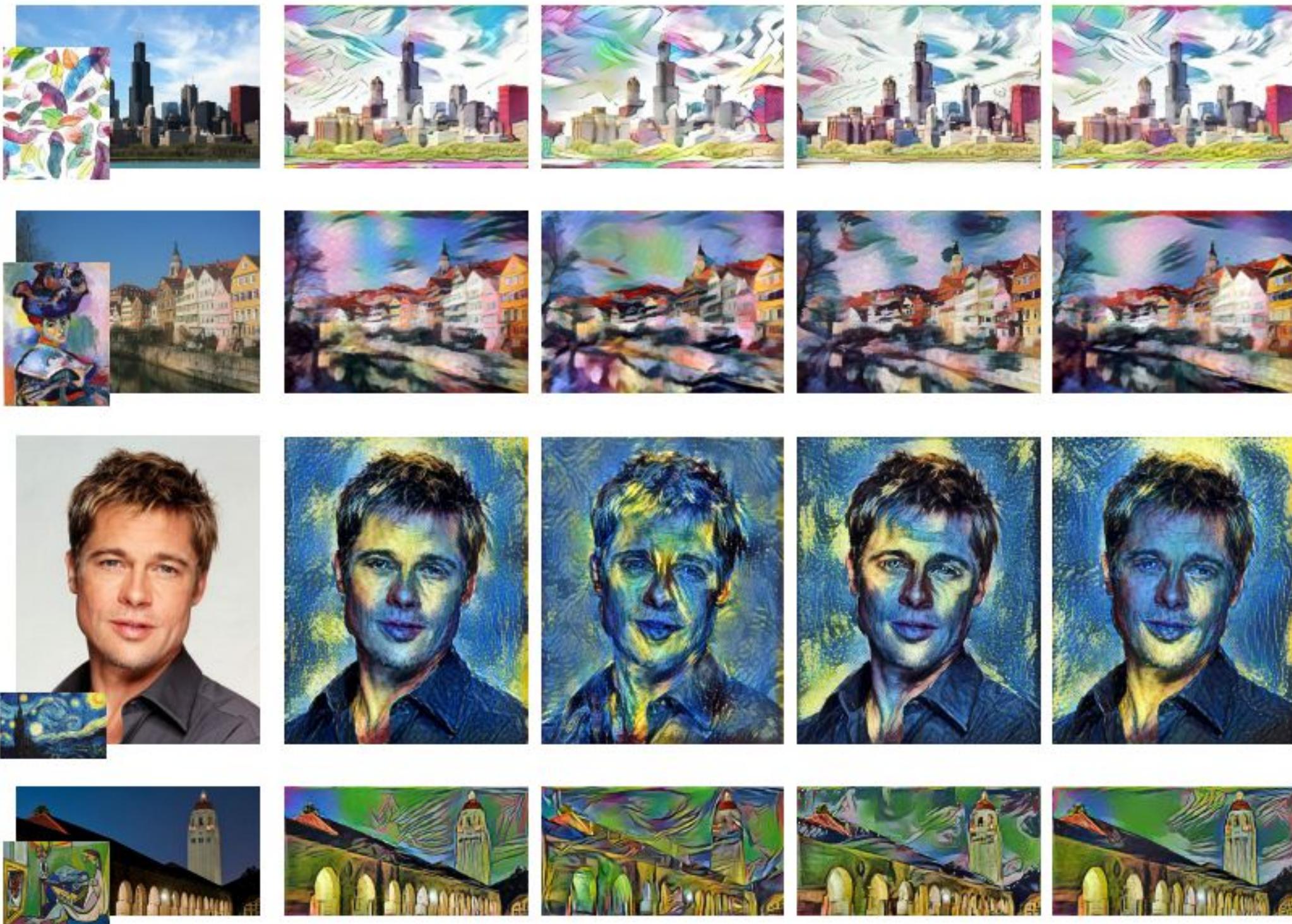
$$E_L = \sum (G^L - A^L)^2$$



[L. A. Gatys et al. Image Style Transfer Using Convolutional Neural Networks. CVPR 2016](#)

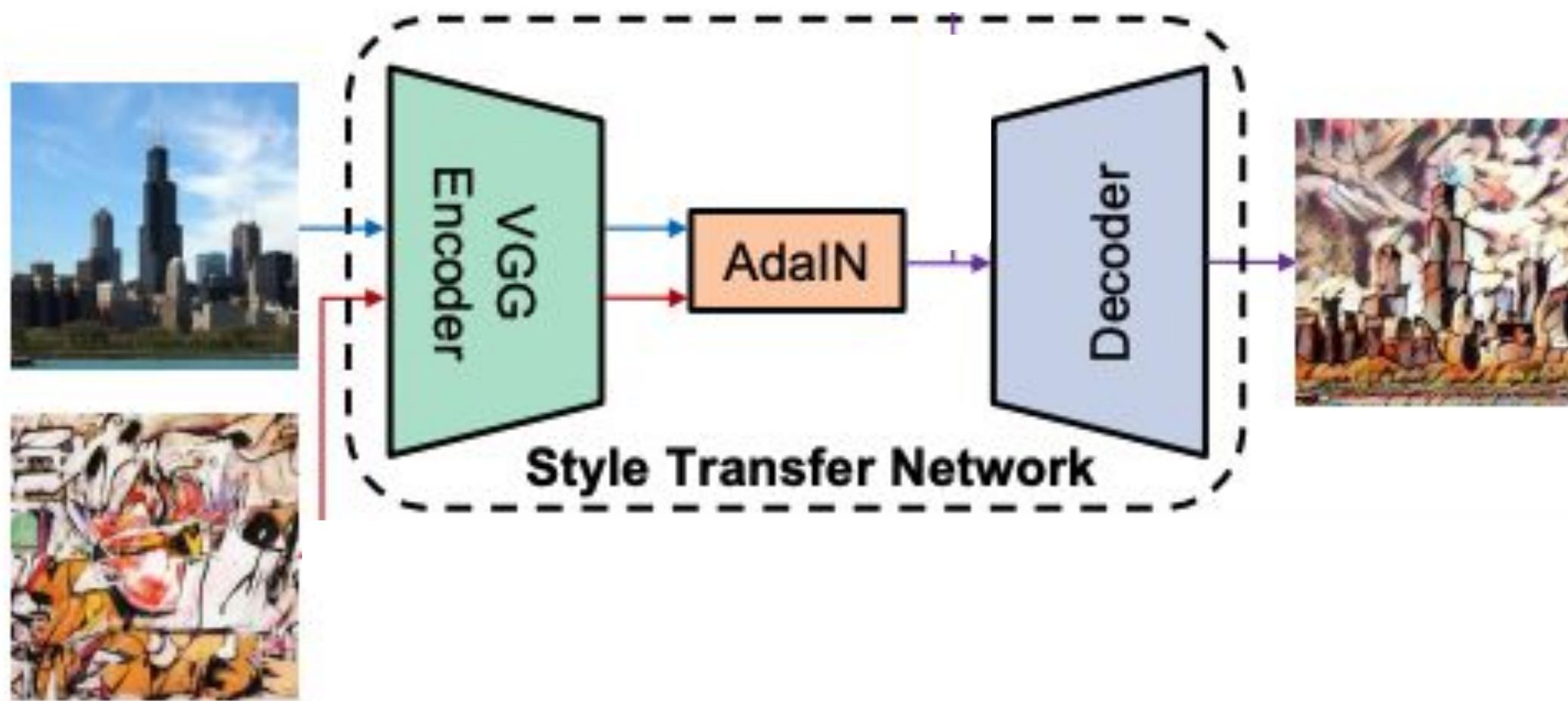
Matching other statistics

1. Style transfer as a distribution alignment
2. Matching other statistics, including channel-wise mean and variance



(a) Content / Style (b) linear (c) poly (d) Gaussian (e) BN

Introducing AdaIN Block

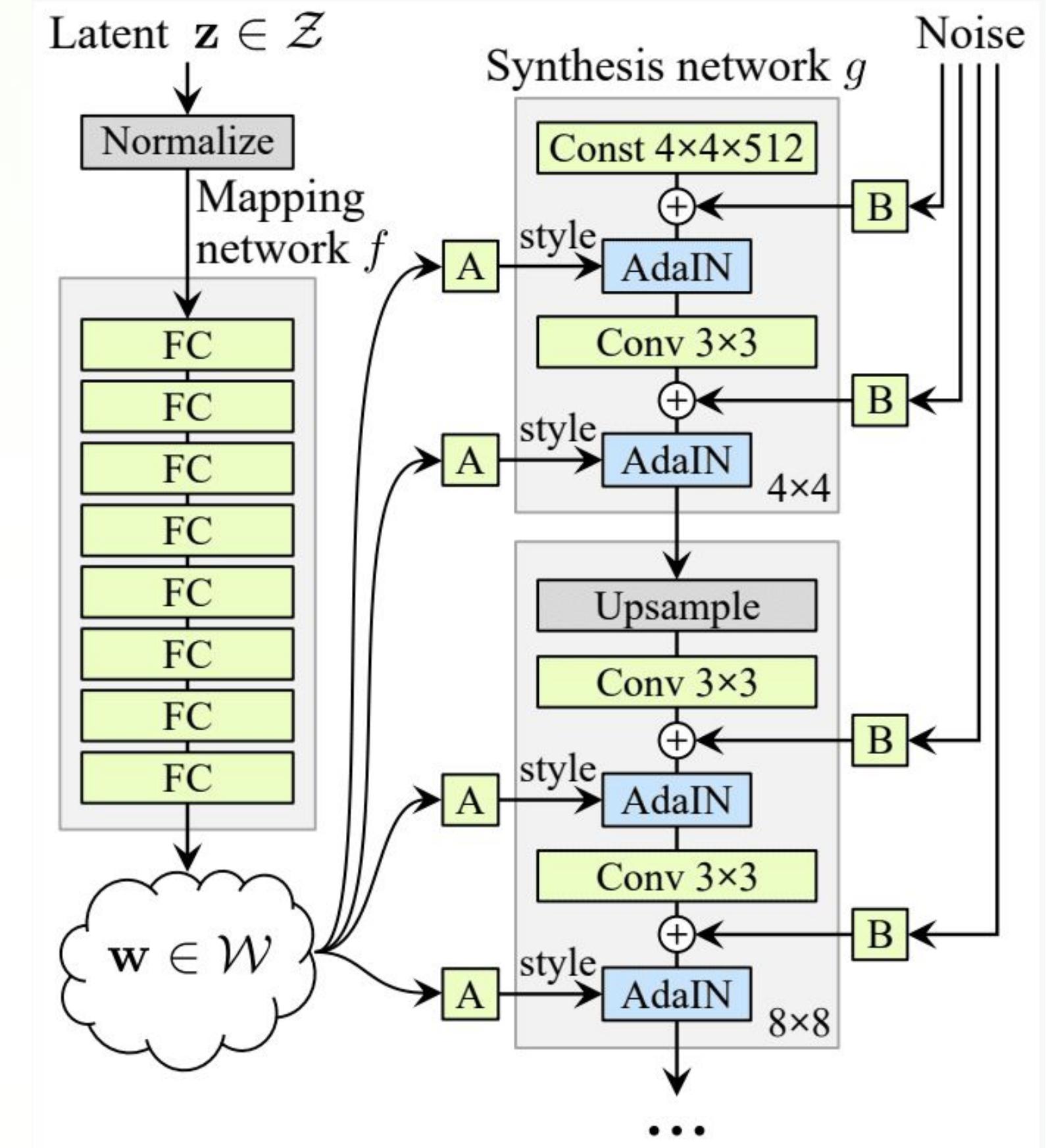


[X. Huang, S. Belongie, Arbitrary Style Transfer in Real-time with Adaptive Instance Normalization. ICCV 2017](#)

StyleGAN



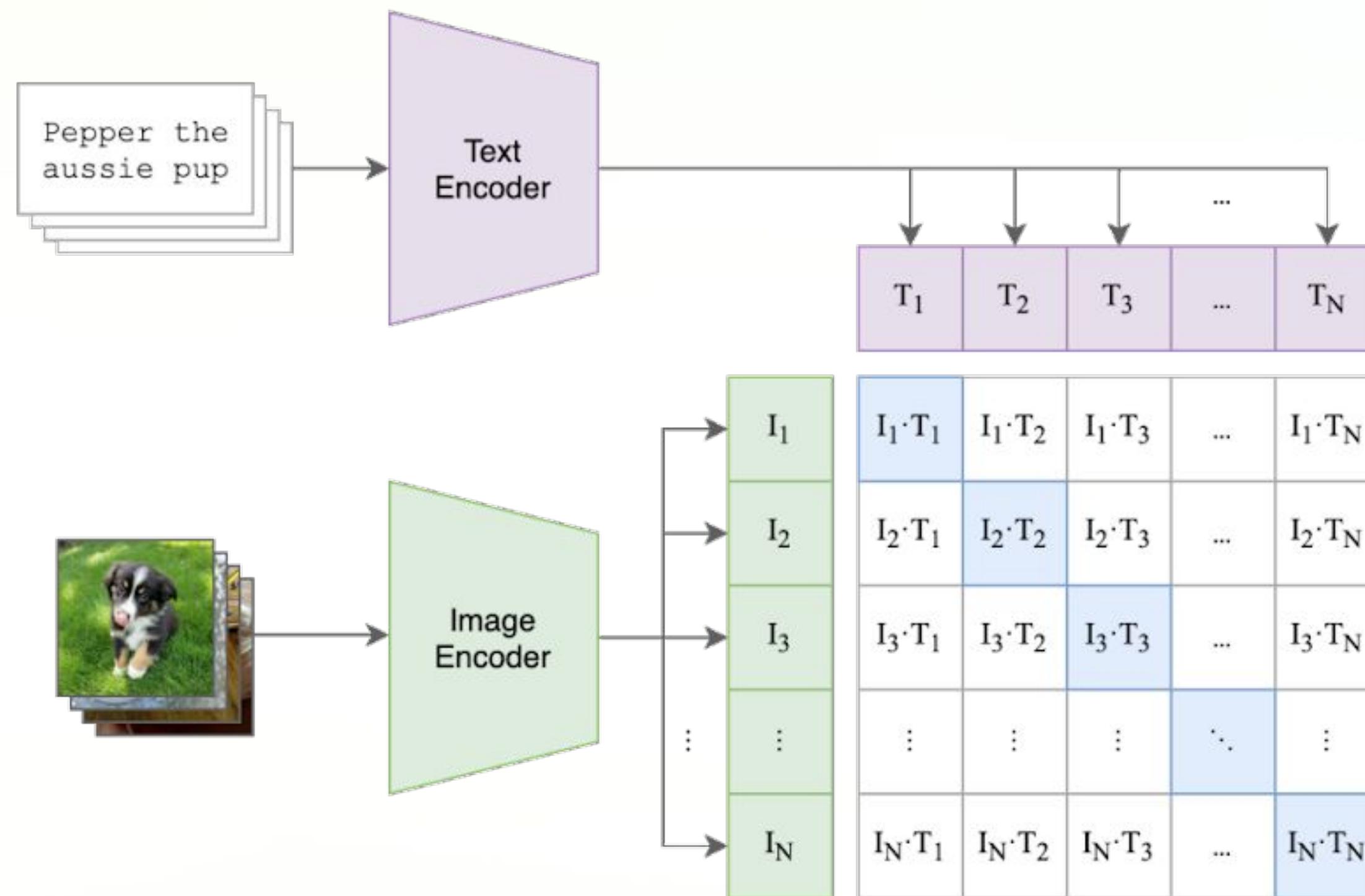
Analyzing and Improving the Image Quality of
StyleGAN





CLIP- Text and Image encoder

OpenAI, Jan 2021



Diffusion models

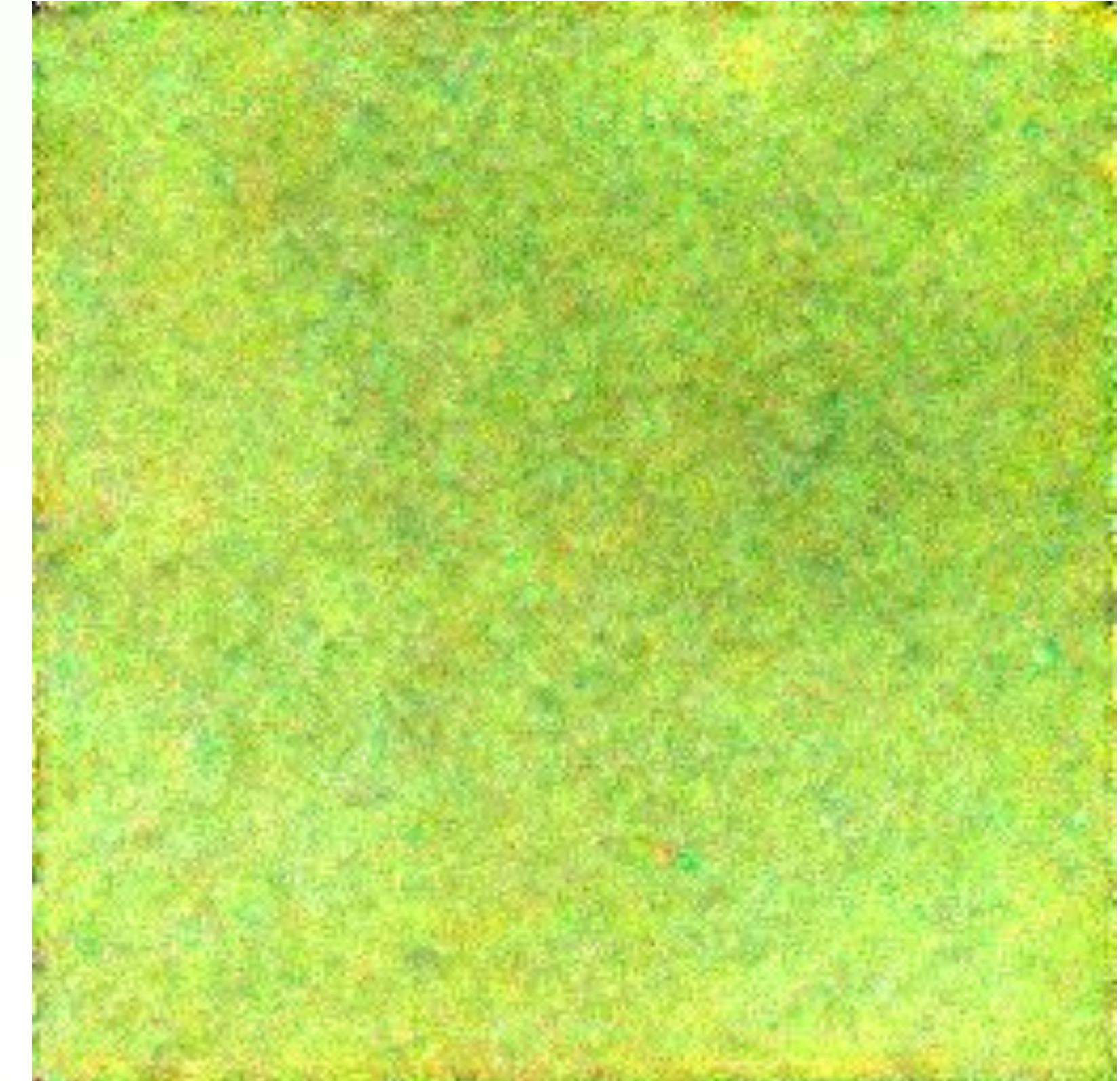
Guided image generation

UC Berkeley, June 2020

OpenAI, June 2021

Google Brain, July 2021

- Model that learns to iteratively de-noise until image is formed.
- A new form of image synthesis architectures.
- Can be conditioned on a class of images, text embedding, low-res image, etc.
- Slow to run and train, but inference becoming faster and faster these days.



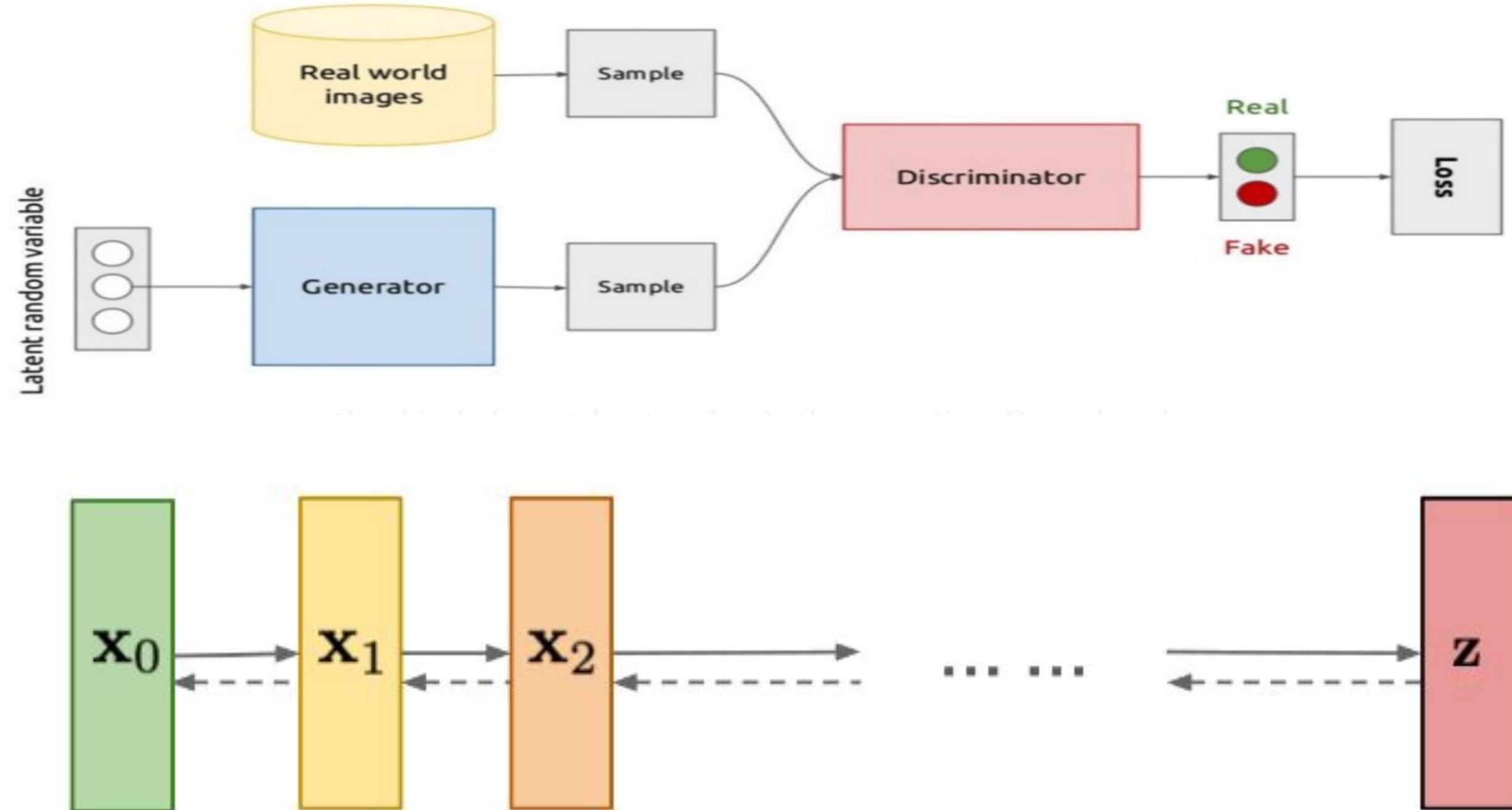
DiscoDiffusion v5, “Cheese and grapes next to a glass of wine”

Diffusion model

Input: random noise tensor (dimensions -
 $64 \times 64 \times 6$)

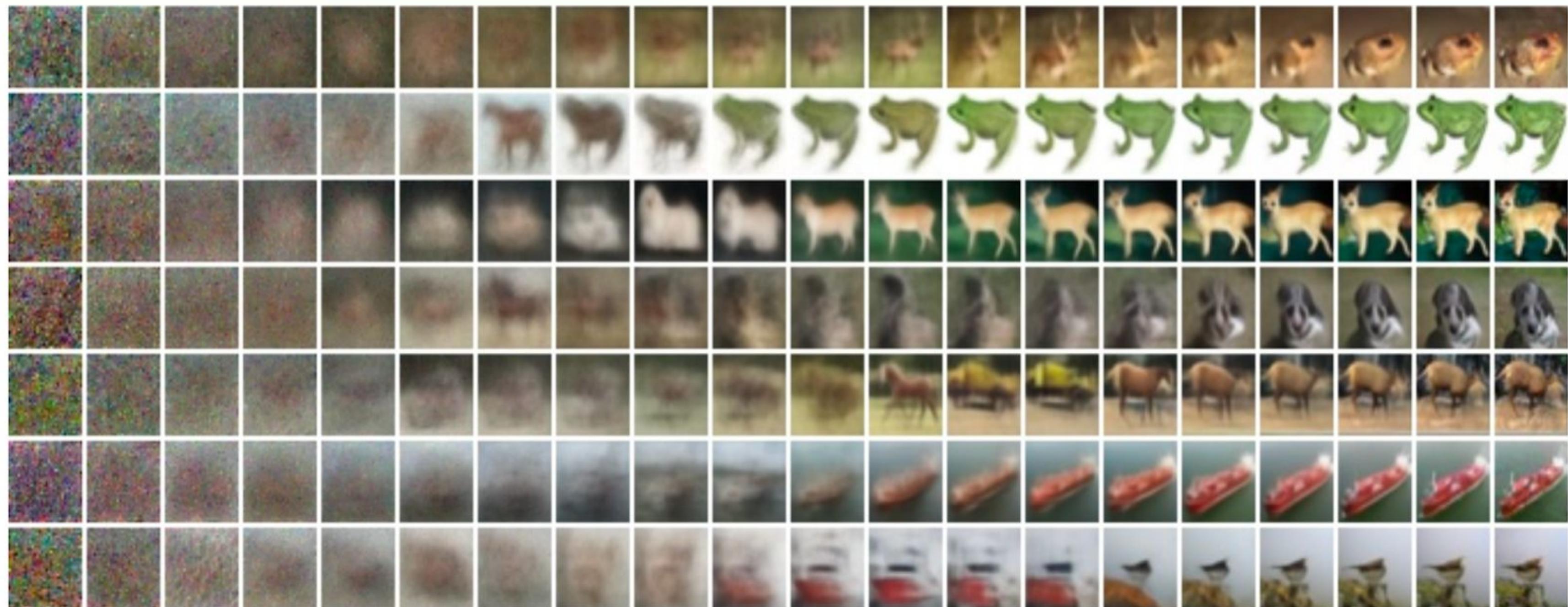
Output: Image (dimensions $64 \times 64 \times 3$)

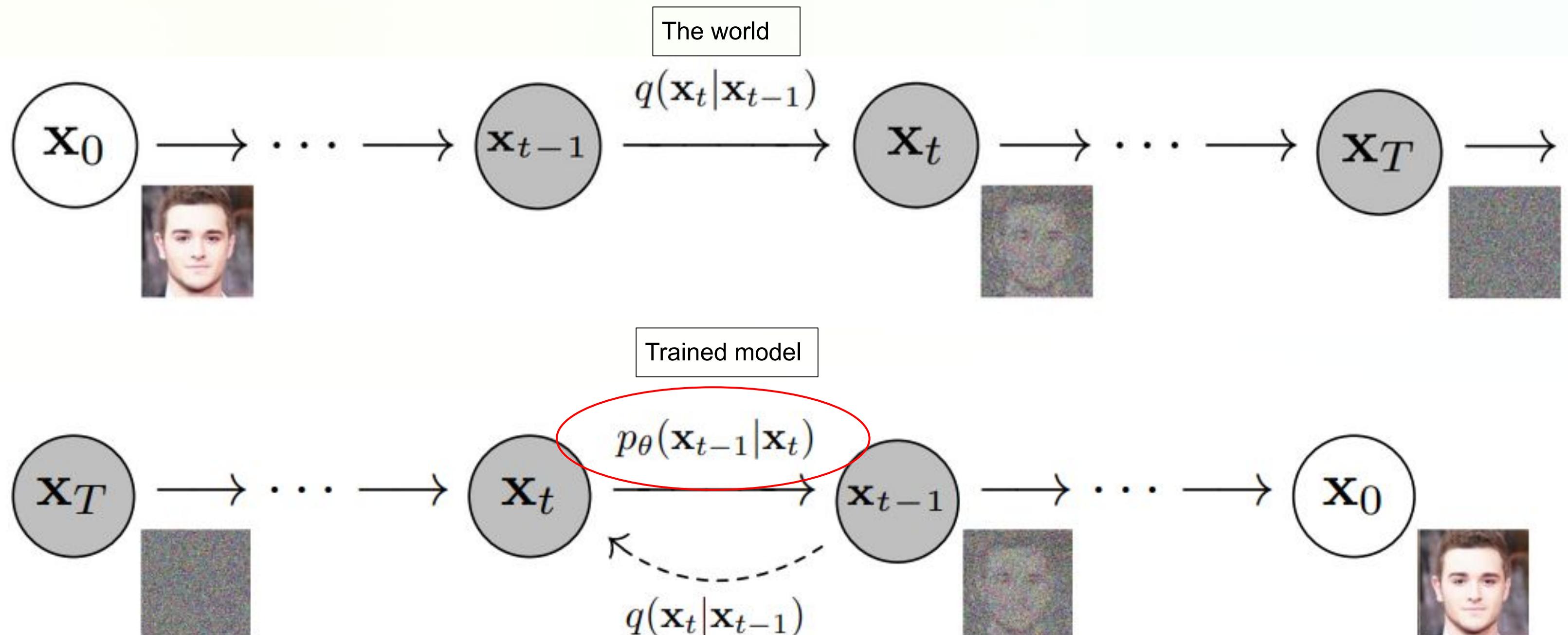
GAN: Adversarial training



Diffusion models:
Gradually add Gaussian noise and then reverse

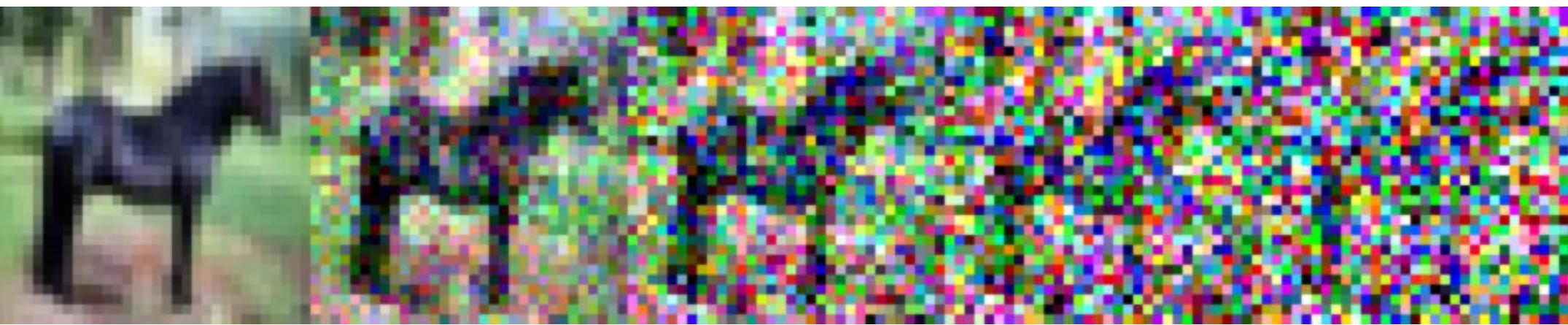
Diffusion models: from gaussian noise to images in small steps





Creating noised images

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



sample an arbitrary step of the noised latents directly conditioned on the input \mathbf{x}_0

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}) \text{ where } \bar{\alpha}_t = \prod_{i=1}^T \alpha_i \quad \alpha_t := 1 - \beta_t$$

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}$$

Training a diffusion model - The reverse step

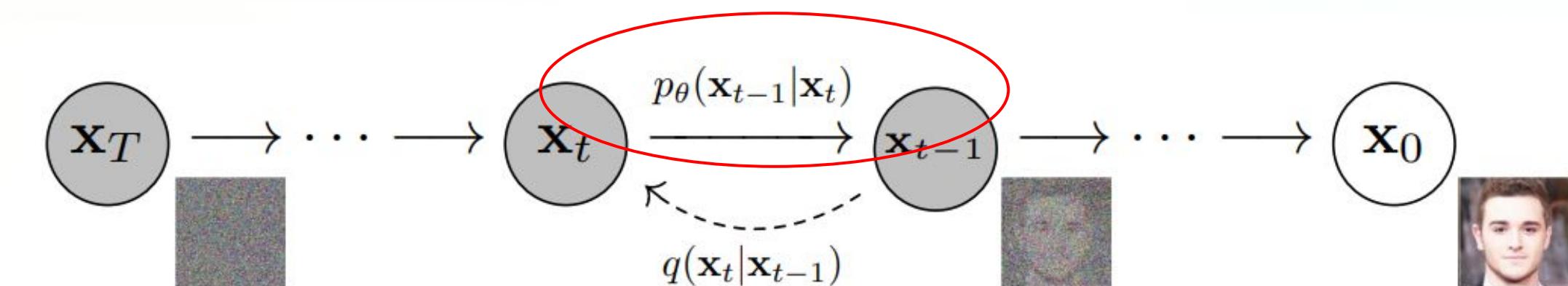
Diffusion model can be trained to predict \mathbf{x}_{t-1} given \mathbf{x}_t, t

But, In practice those models are trained to predict the added noise

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon$$

Loss: L2($noise_{\theta}(\mathbf{x}_t), \epsilon$)

Iteration step: $\mathbf{x}_{t-1} = \sqrt{\bar{\alpha}_t}(\mathbf{x}_t - noise_{\theta}(\mathbf{x}_t)) + \sqrt{(1 - \bar{\alpha}_t)}noise_{\theta}(\mathbf{x}_t)$

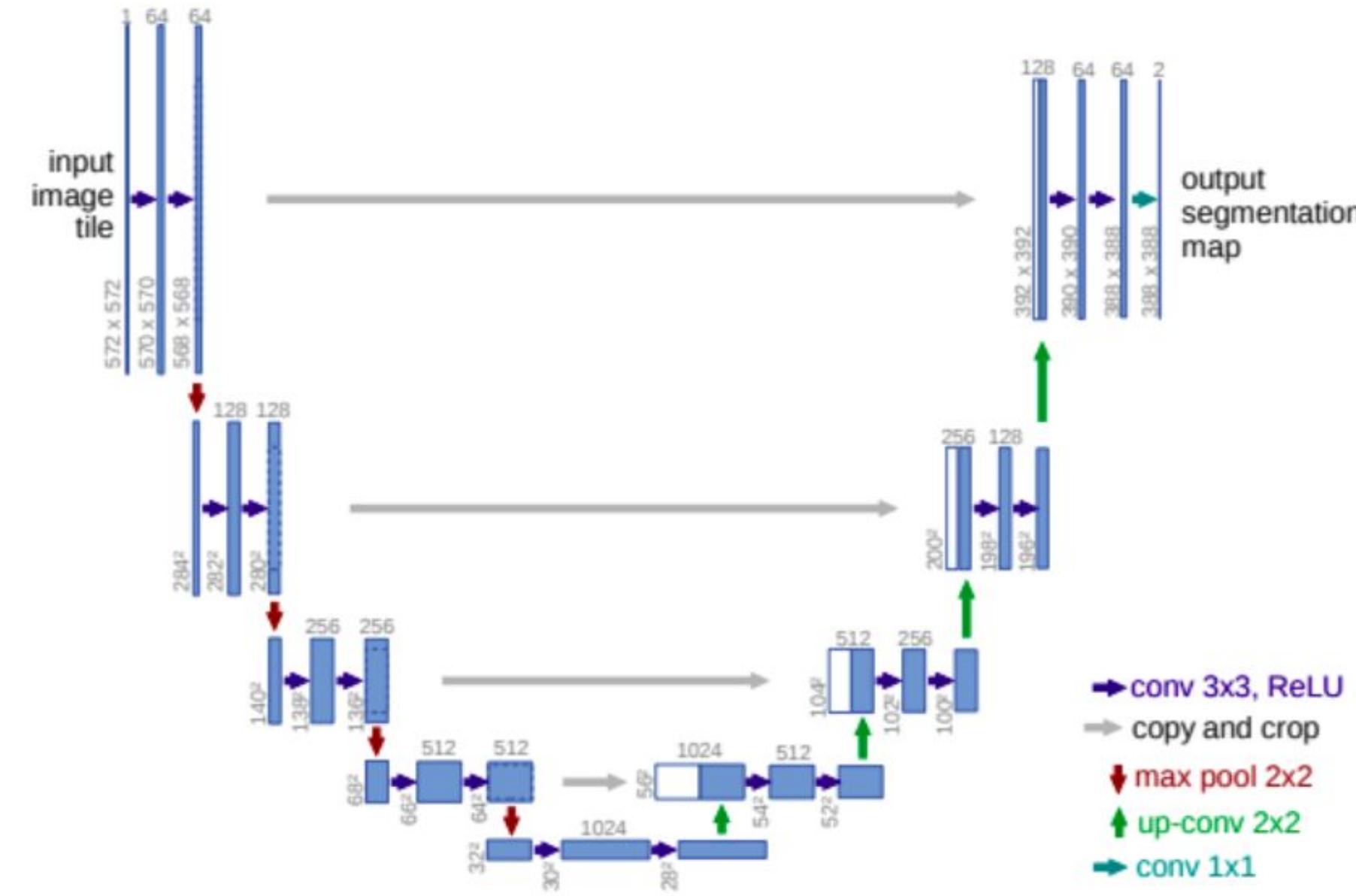


Dimensions:

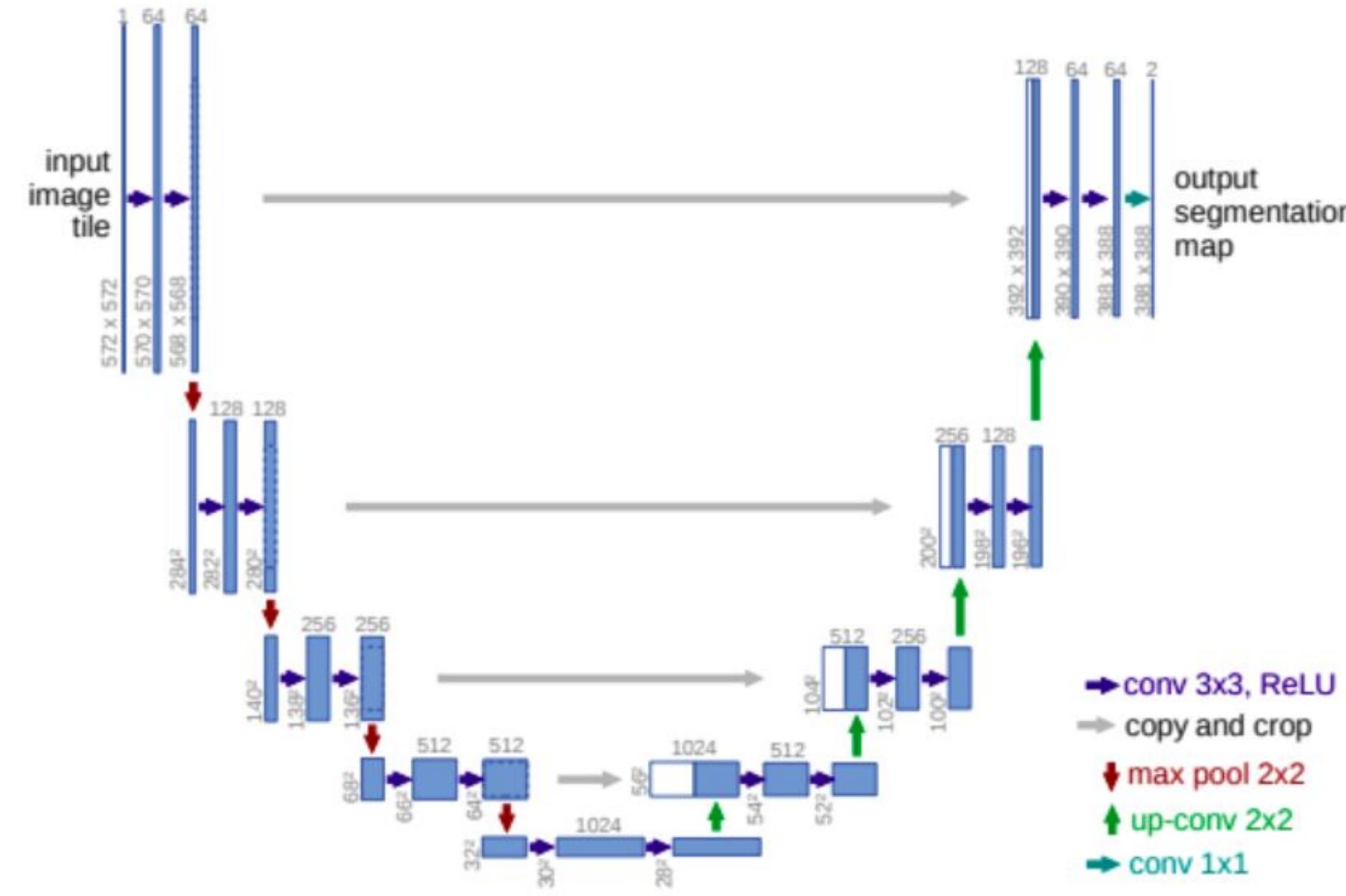
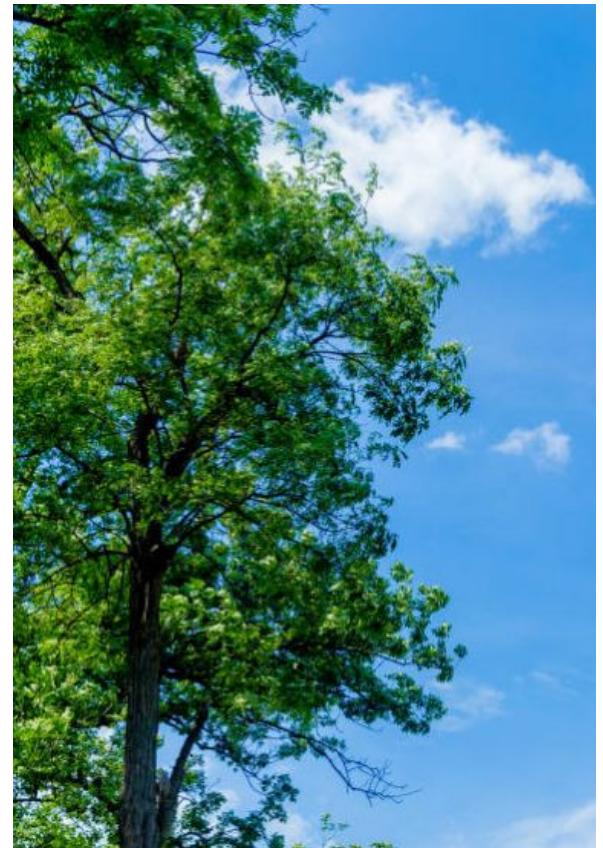
$\bar{\alpha}_t$ - scalar

\mathbf{x}_t Tensor 64x64x3

Unet Architecture



Unet Architecture



Questions?

Hands on session

Hands on session (30 minutes)

Training a denoising model

1. Go to

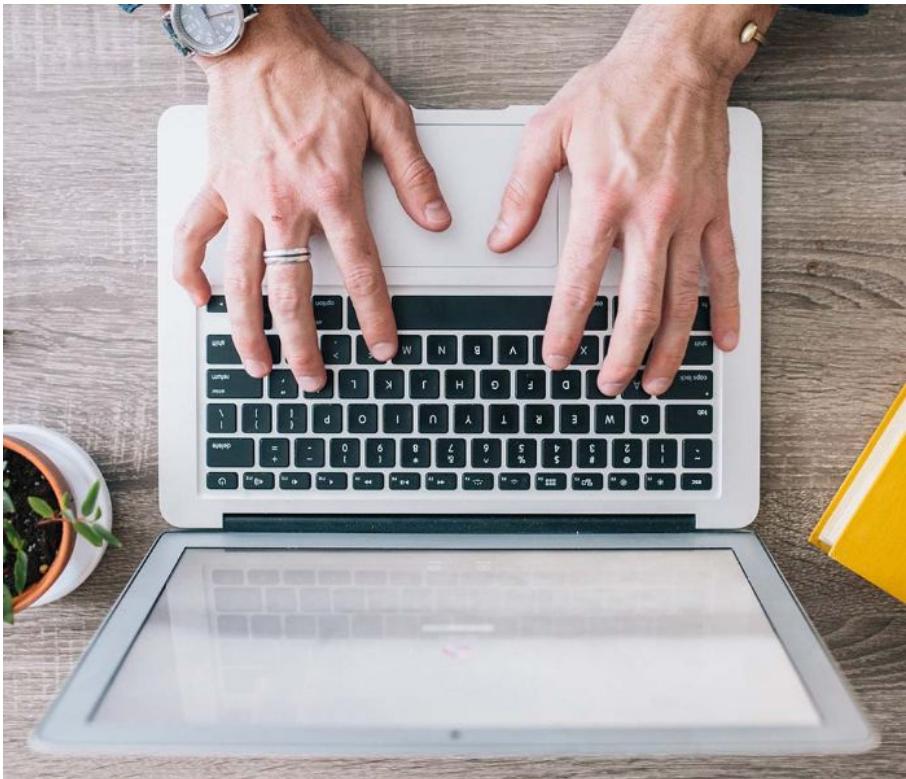
-<https://github.com/Naomi-Ken-Korem/text-to-image-datatalks>

2. Open Diffusion models part 1.ipynb

(solution:

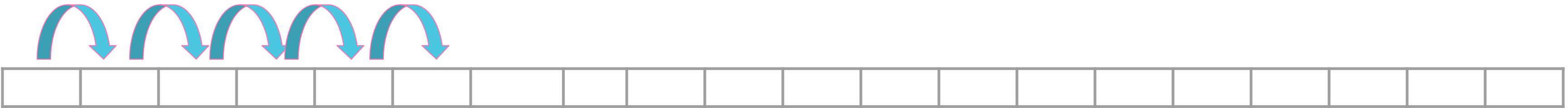
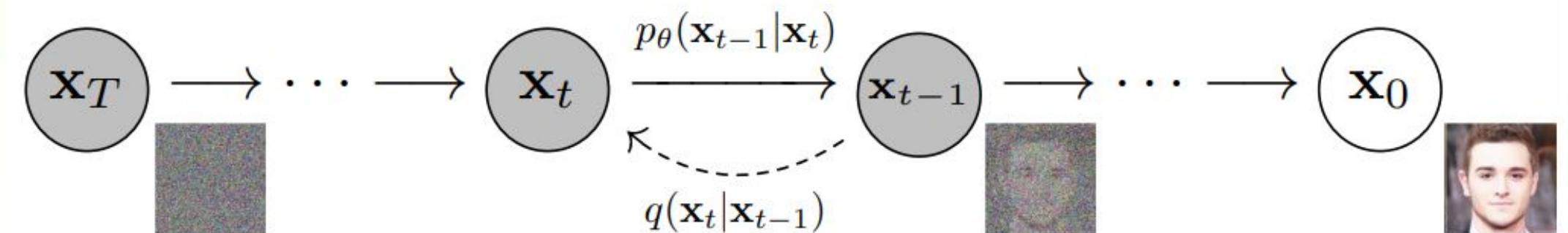
Diffusion models part 1_solution.ipynb)

3. Solve notebook till the training cell, and we will continue from there.



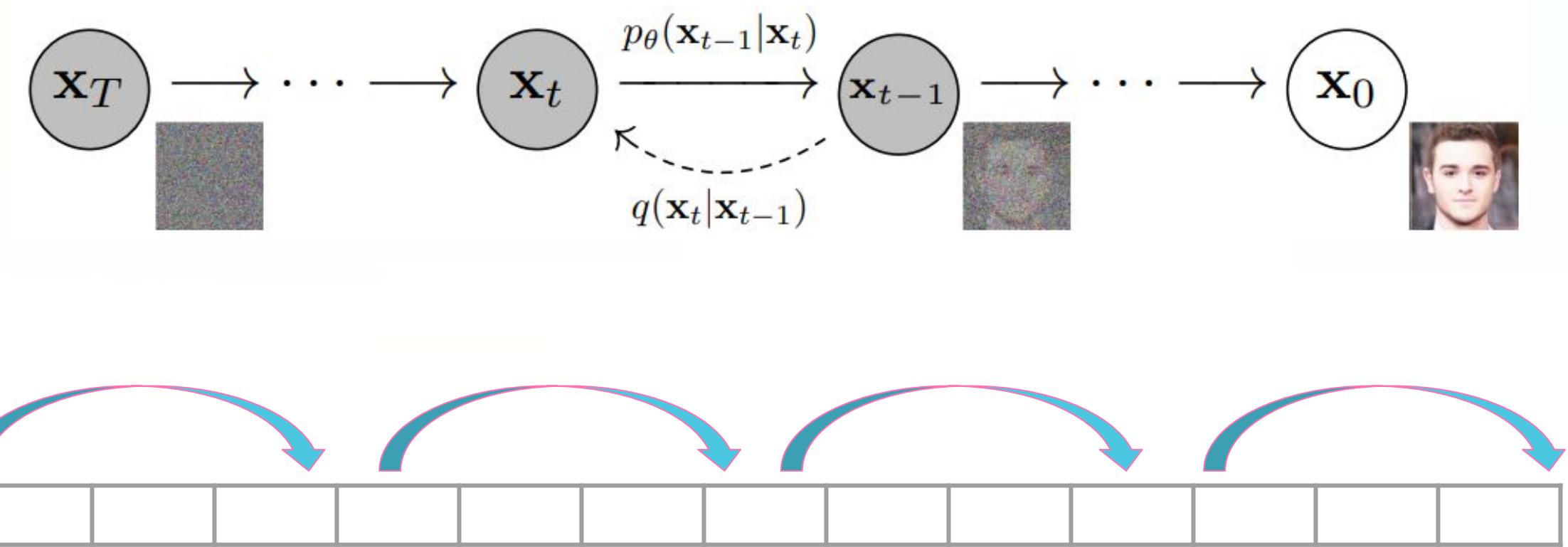
Inference

1000 steps??



Inference

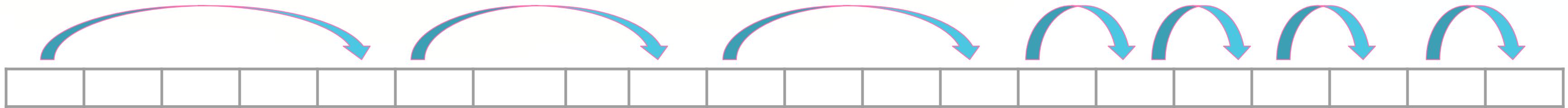
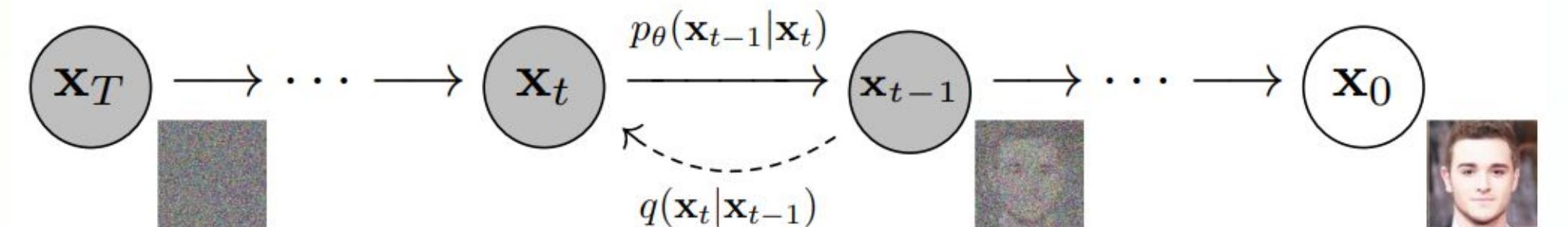
1000 steps??



$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}$$

Inference

1000 steps??



$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon$$

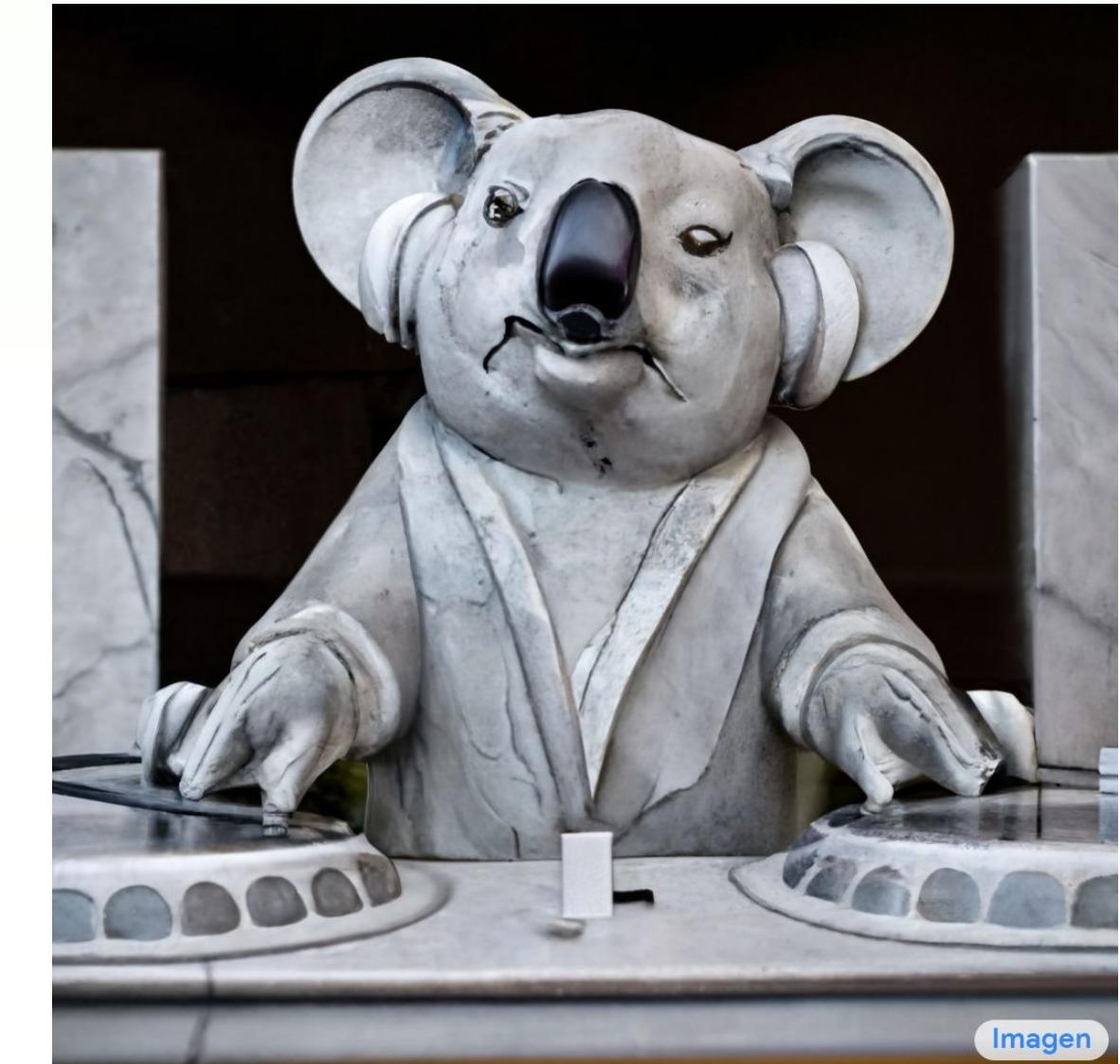
Part 2 - Text condition Diffusion models



A dragon fruit wearing karate belt in the snow.



An art gallery displaying Monet paintings. The art gallery is flooded. Robots are going around the art gallery using paddle boards.

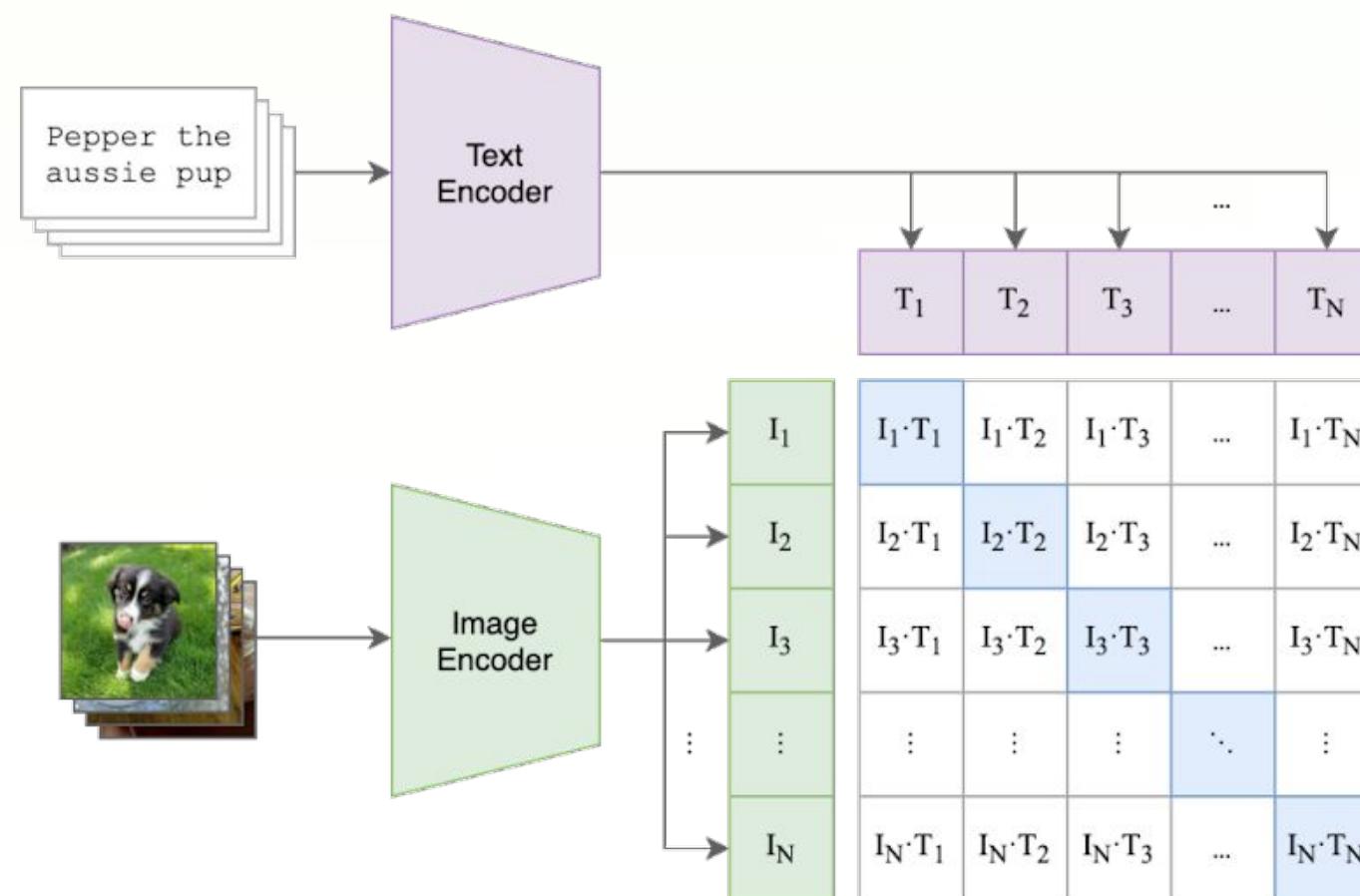


A marble statue of a Koala DJ in front of a marble statue of a turntable. The Koala is wearing large marble headphones.

Recap



StyleGAN



CLIP



Diffusion models

GLIDE

OpenAI, Dec 2021

Given a text prompt + random noise image, generates a related image using a text transformer model based on the model of "Attention is all you need" (Bert) and a conditioned diffusion model.

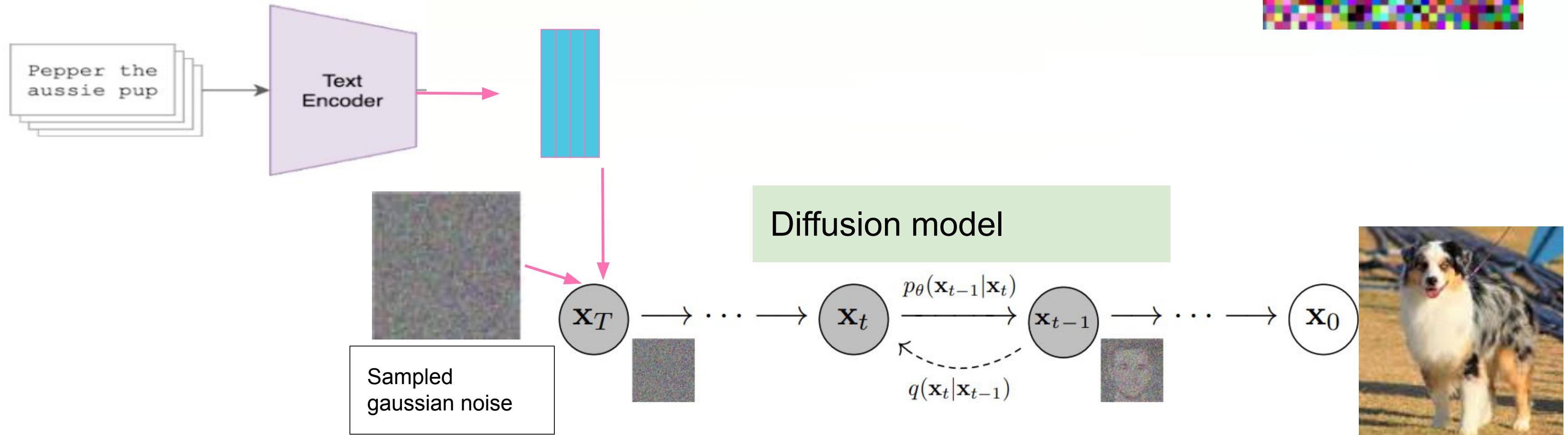


"a boat in the canals of venice"



"a painting of a fox in the style of starry night"

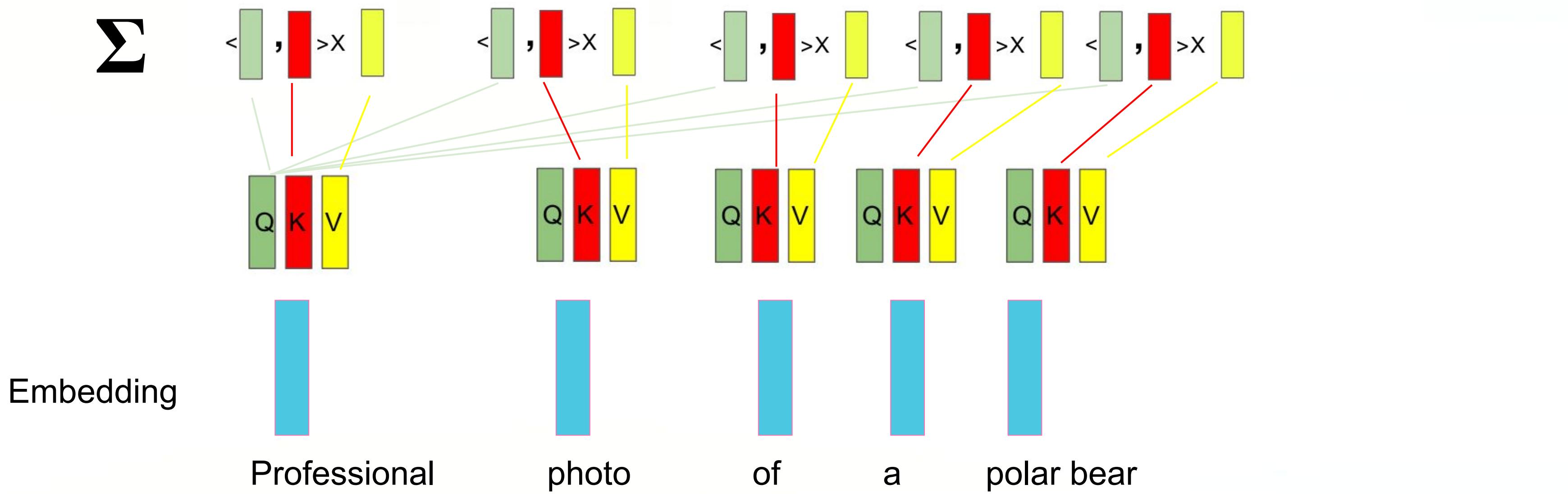
Glide text 2 image architecture



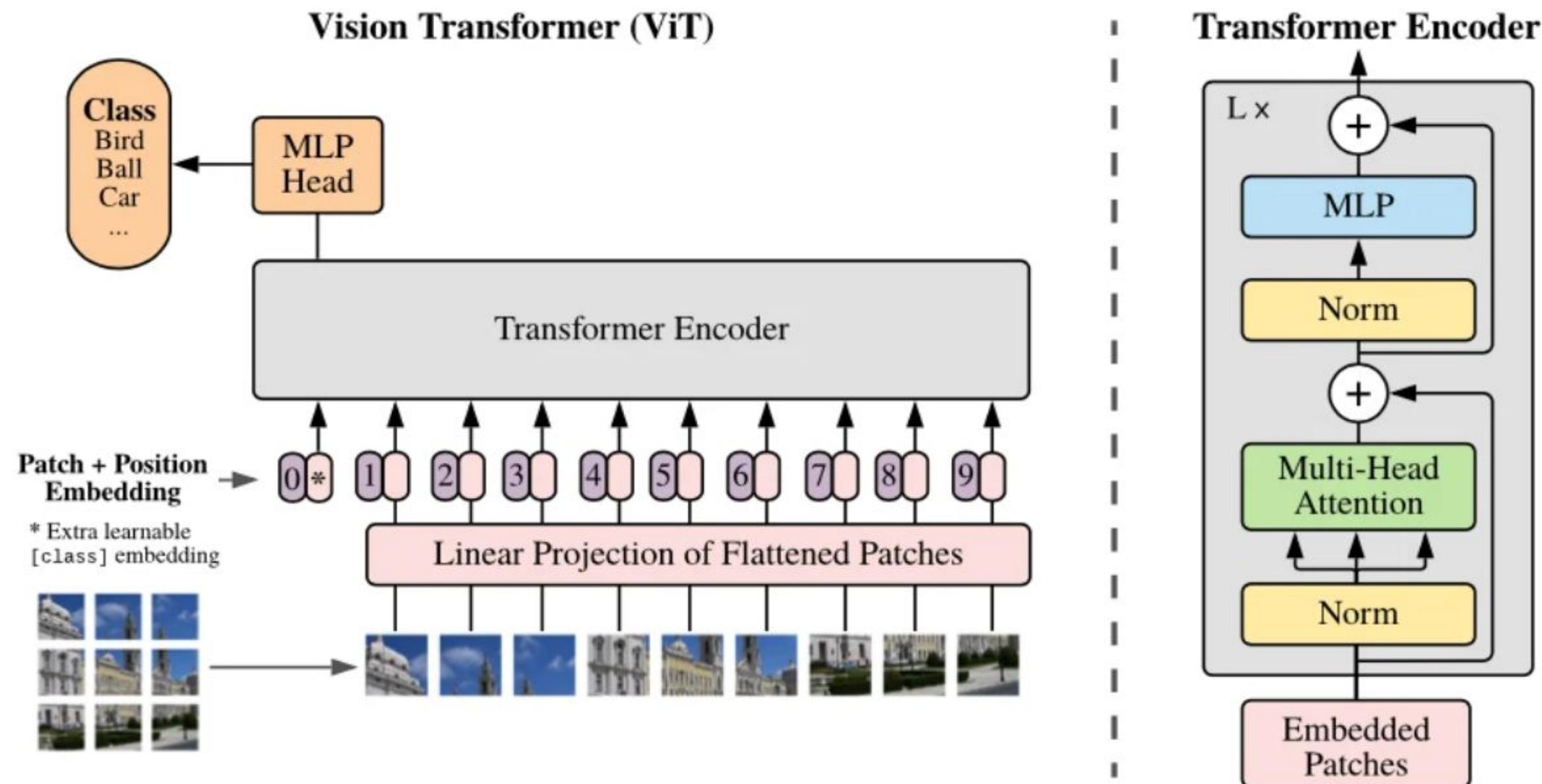
Transformers

Google, Dec 2017

Attention Layer

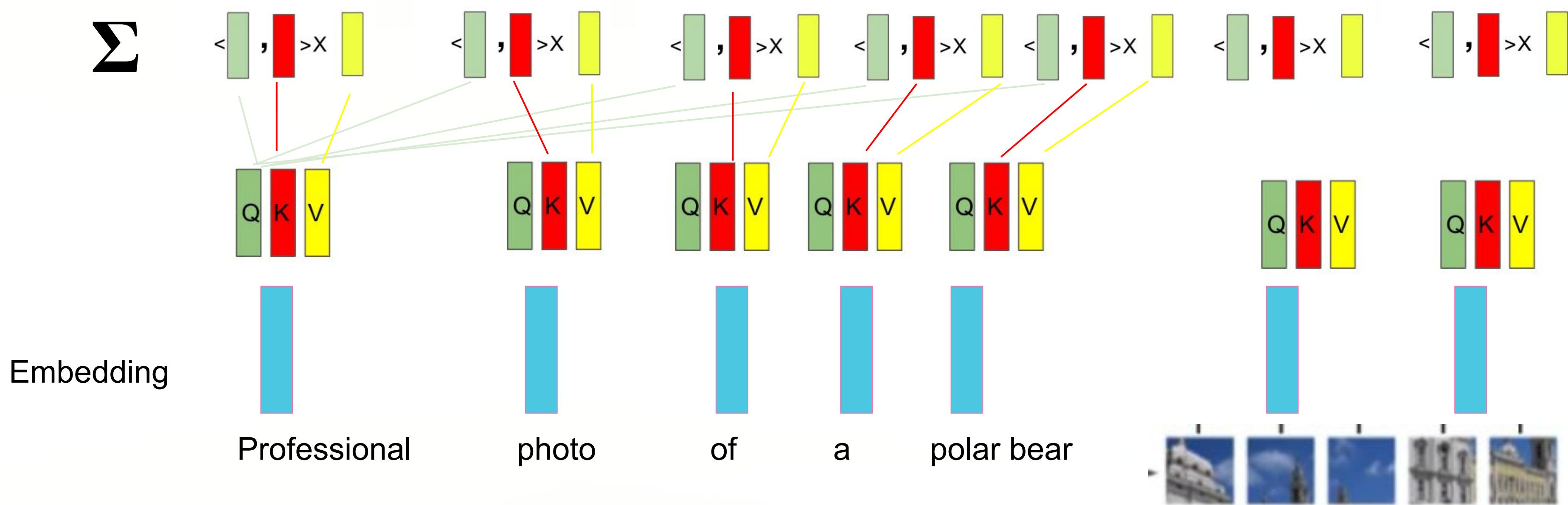


Visual Transformers

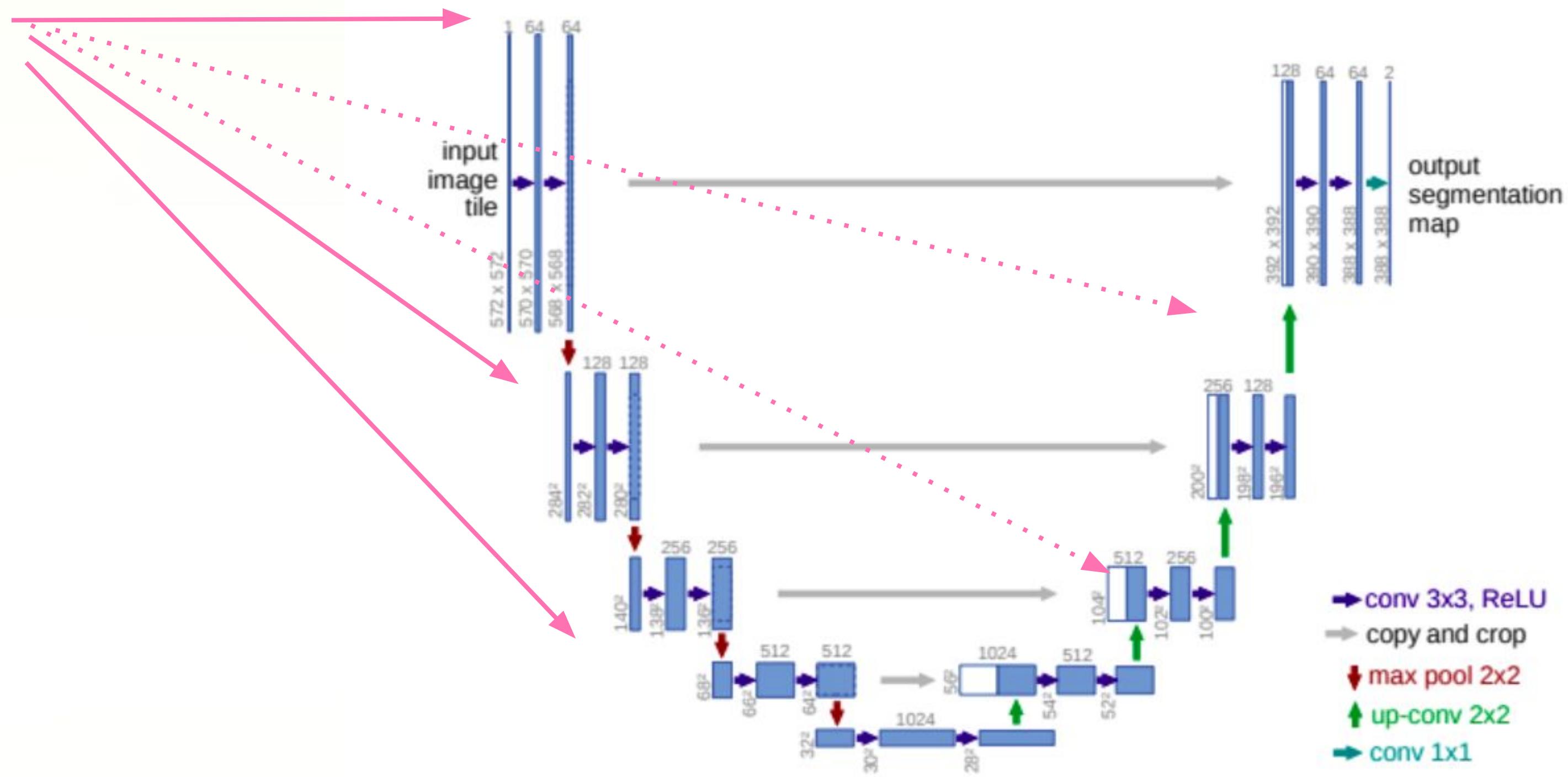


[An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale](#)

Cross Attention Layer

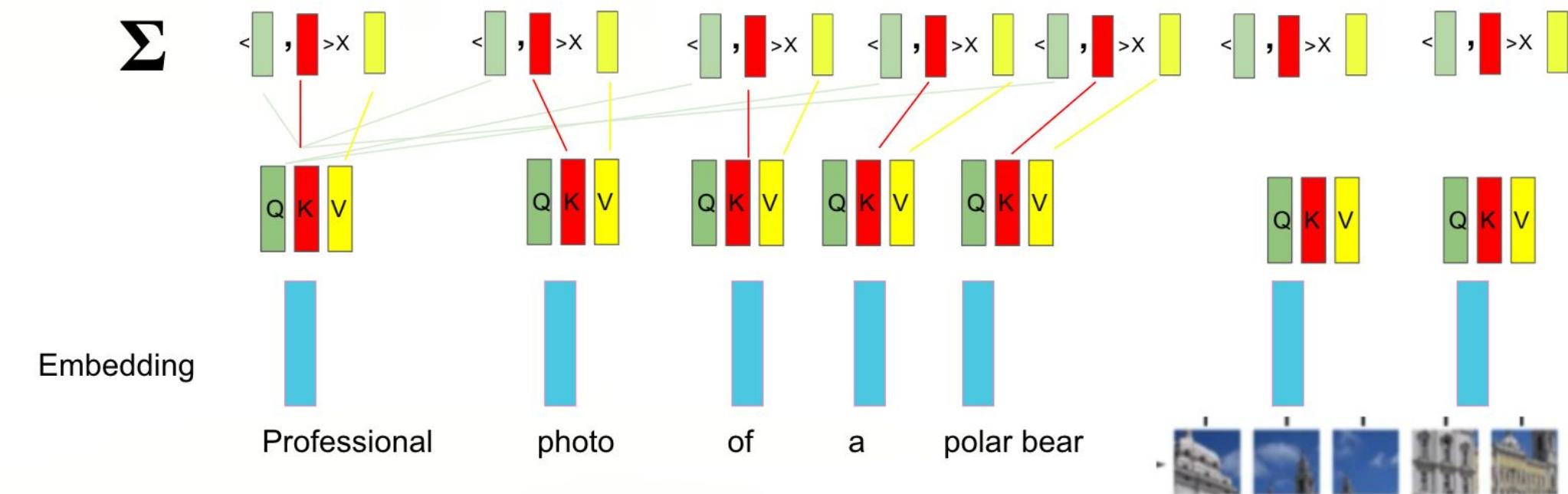
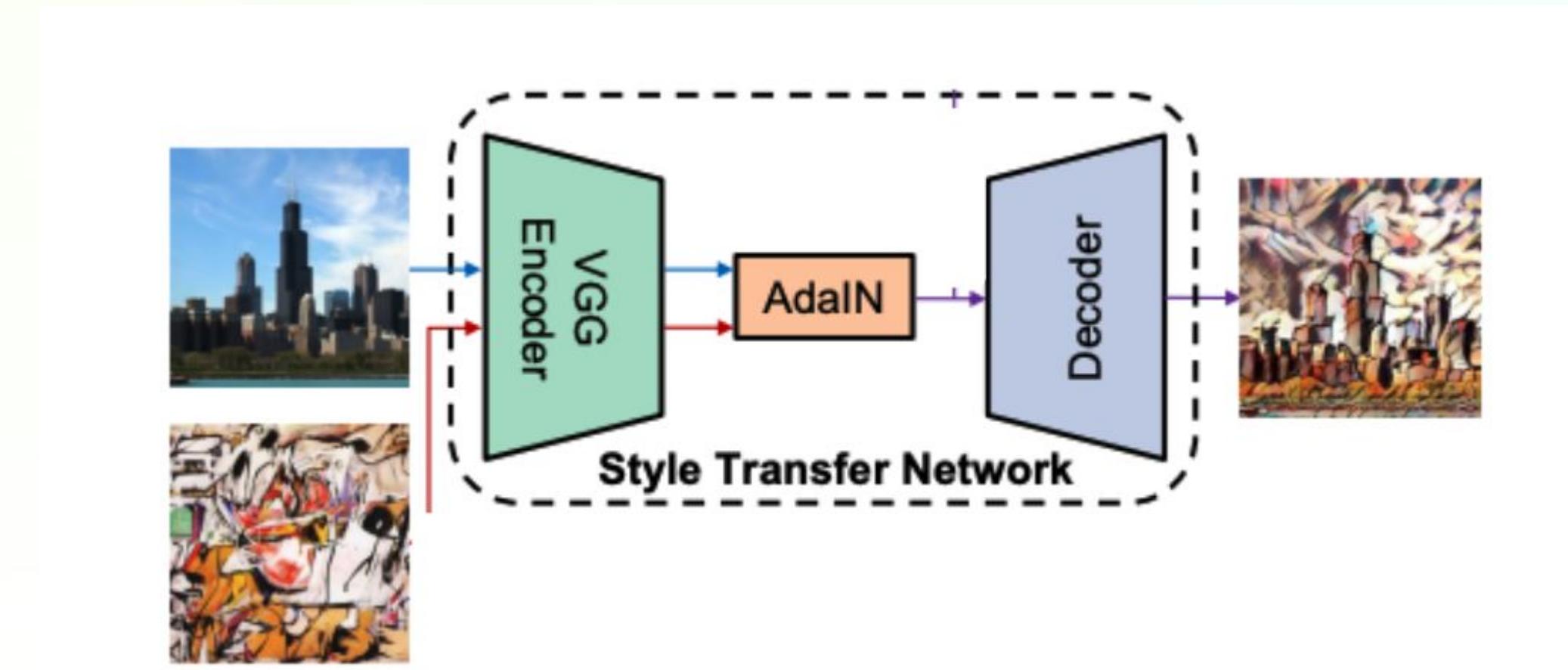


Adding textual input into the diffusion model

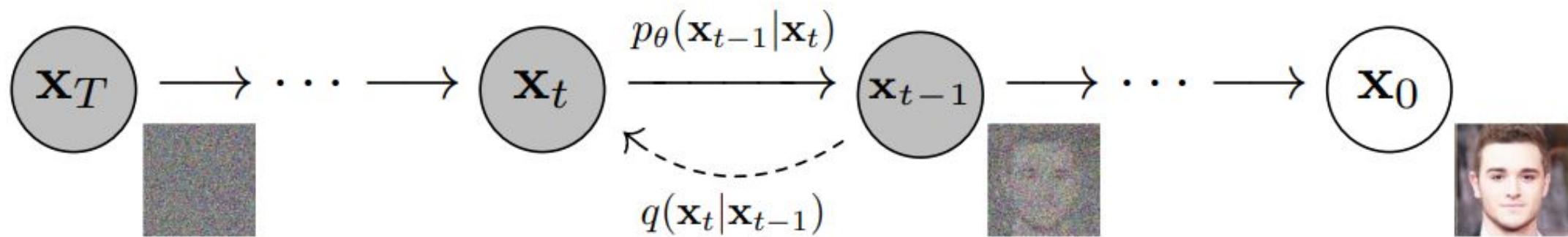


Multimodal layers

1. Style/Text injection via channel statistics
2. Style/Text injection via Cross attention



Class free guidance - extrapolation toward the textual concept



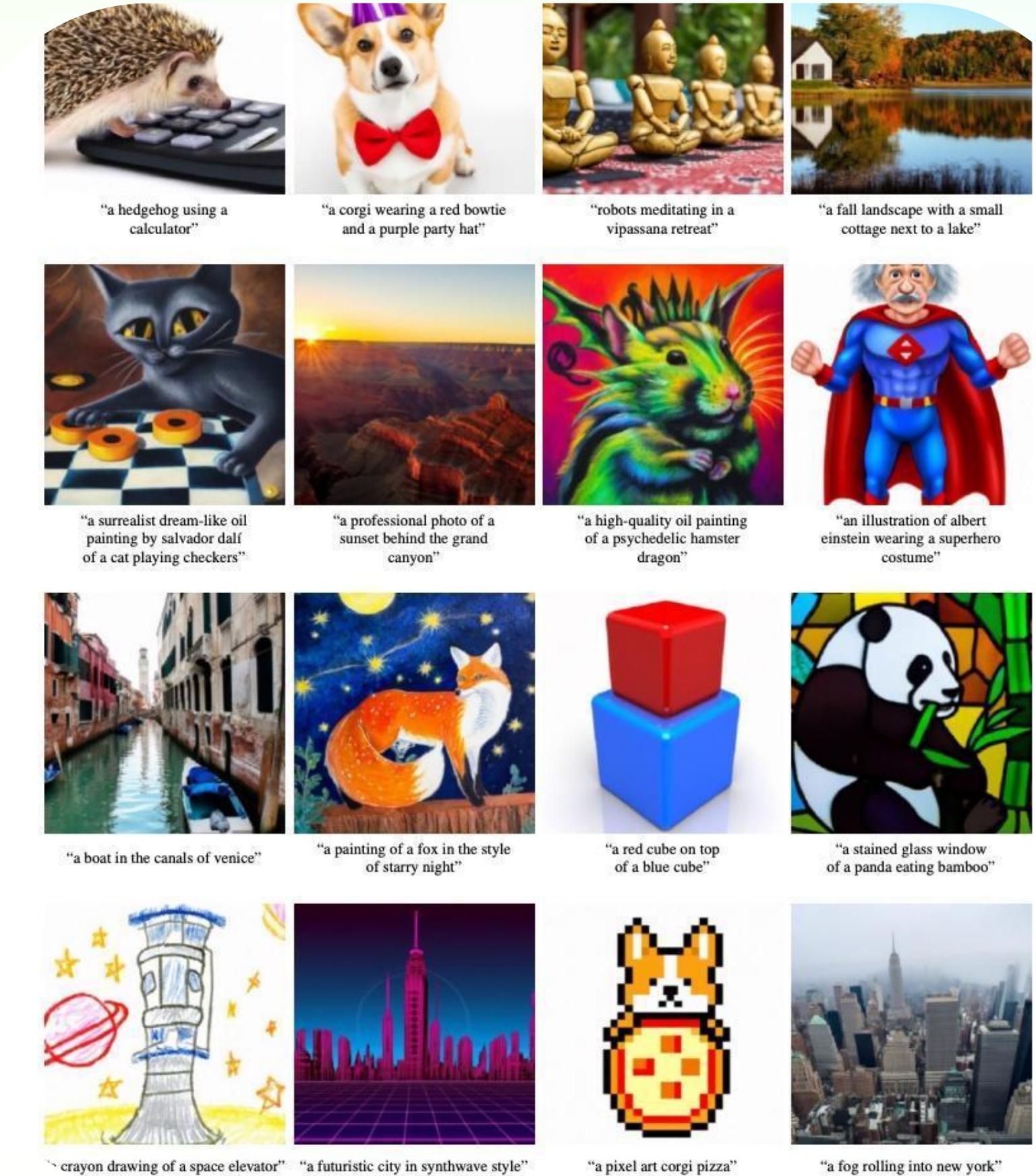
$$\mathbf{x} = \mathbf{x}_{empty} + guidancescale * (\mathbf{x}_{prompt} - \mathbf{x}_{empty})$$

$$\mathbf{x}_{empty:t-1} = P(x_t, "")$$

$$\mathbf{x}_{prompt:t-1} = P(x_t, " prompt ")$$

Seminal publications

Now we can
understand Imagen
and Dalle-2



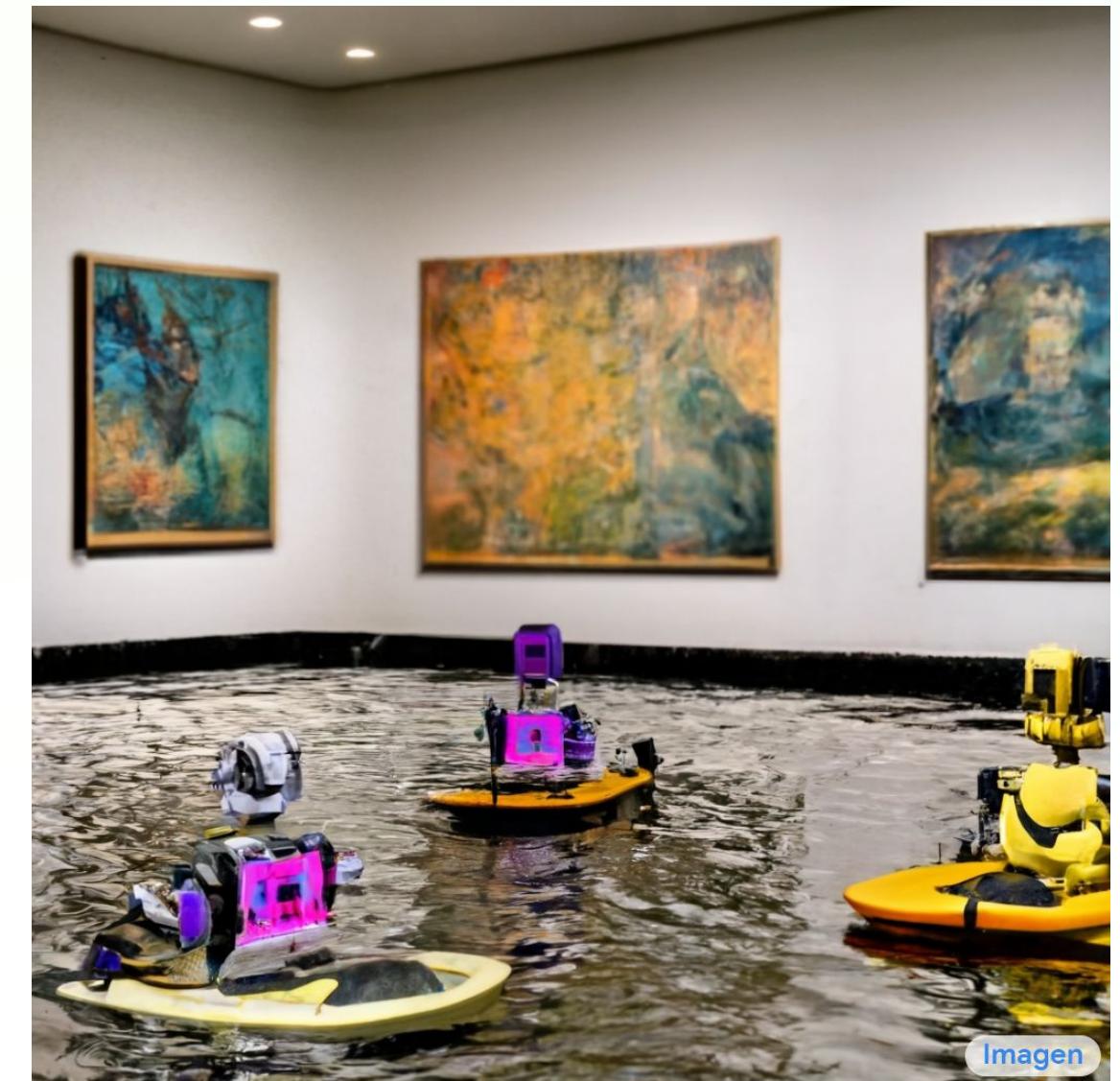
Imagen

Google Brain, May 2022

Huge text model (T5-XXL)

Given text embedding generates an image using a guided diffusion model.

Learn two 4x super-resolution steps (16x in total).



An art gallery displaying Monet paintings.
The art gallery is flooded. Robots are going around the art gallery using paddle boards.

Imagen architecture

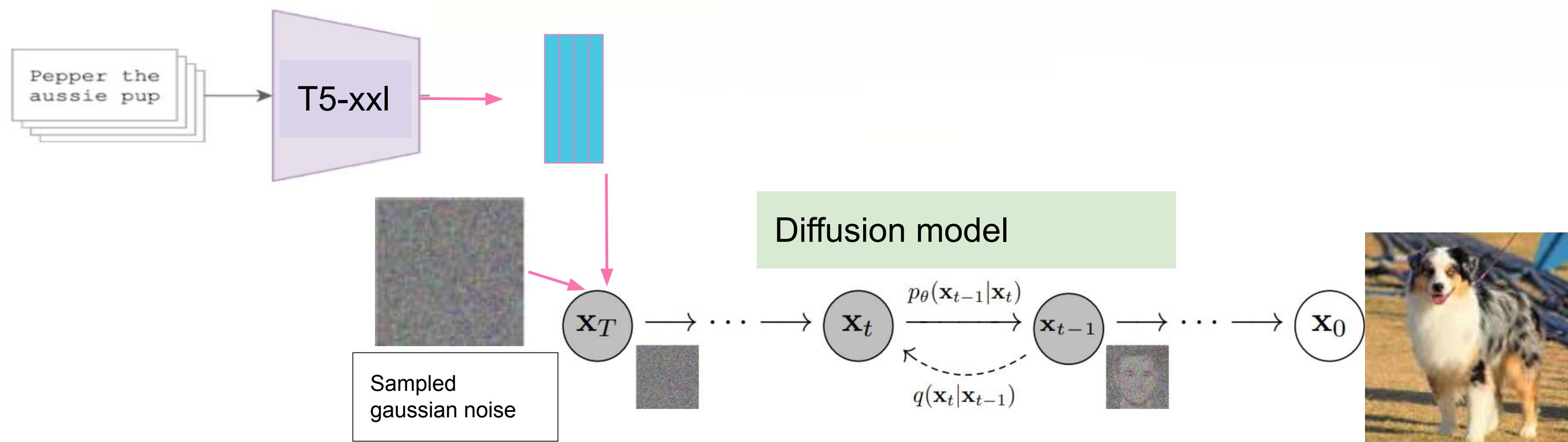
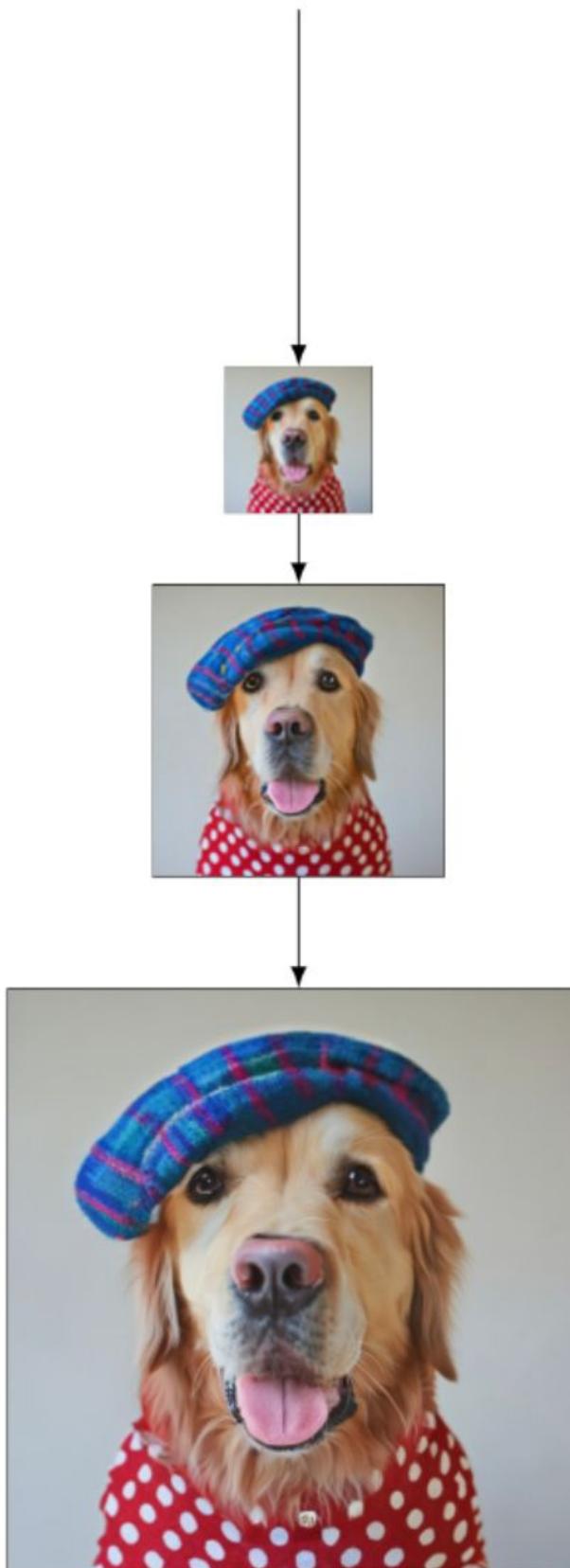


Imagen: Super Resolution models

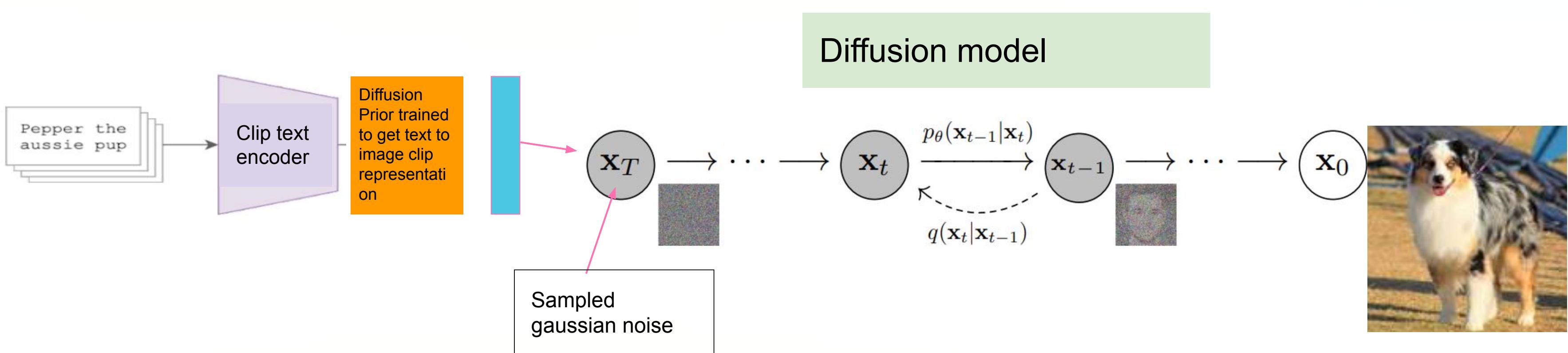
“A Golden Retriever dog wearing a blue checkered beret and red dotted turtleneck.”



Dalle-2

openAI april 2022

Unclip – from clip image vector back into image





Delighted to announce the public open source release
of [#StableDiffusion!](#)

Please see our release post and retweet!
[stability.ai/blog/stable-di...](https://stability.ai/blog/stable-diffusion-public-release/)

Proud of everyone involved in releasing this tech that is
the first of a series of models to activate the creative
potential of humanity

[תרגום את הציג](#)

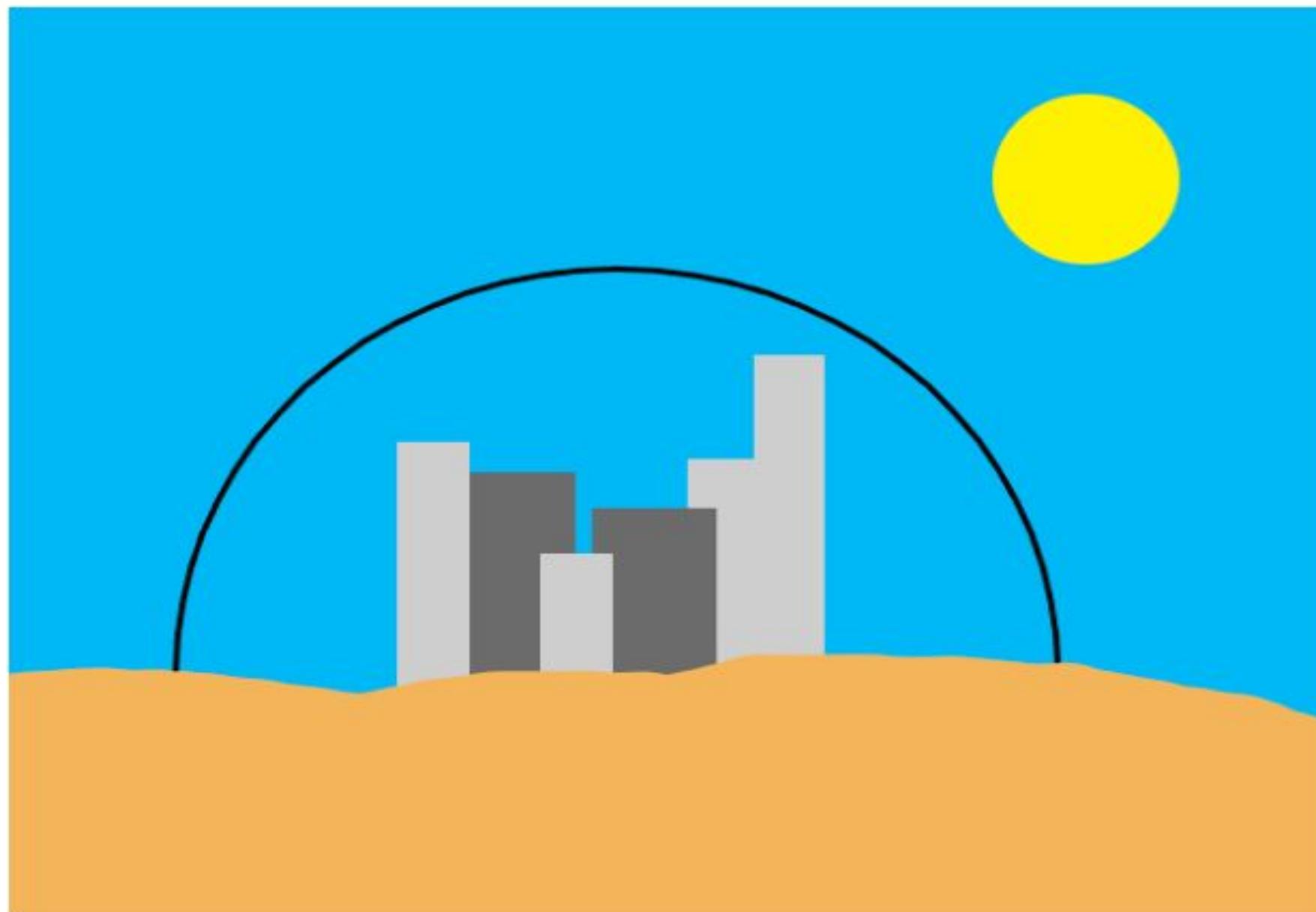
stability.ai

Stable Diffusion Public Release — Stability.Ai

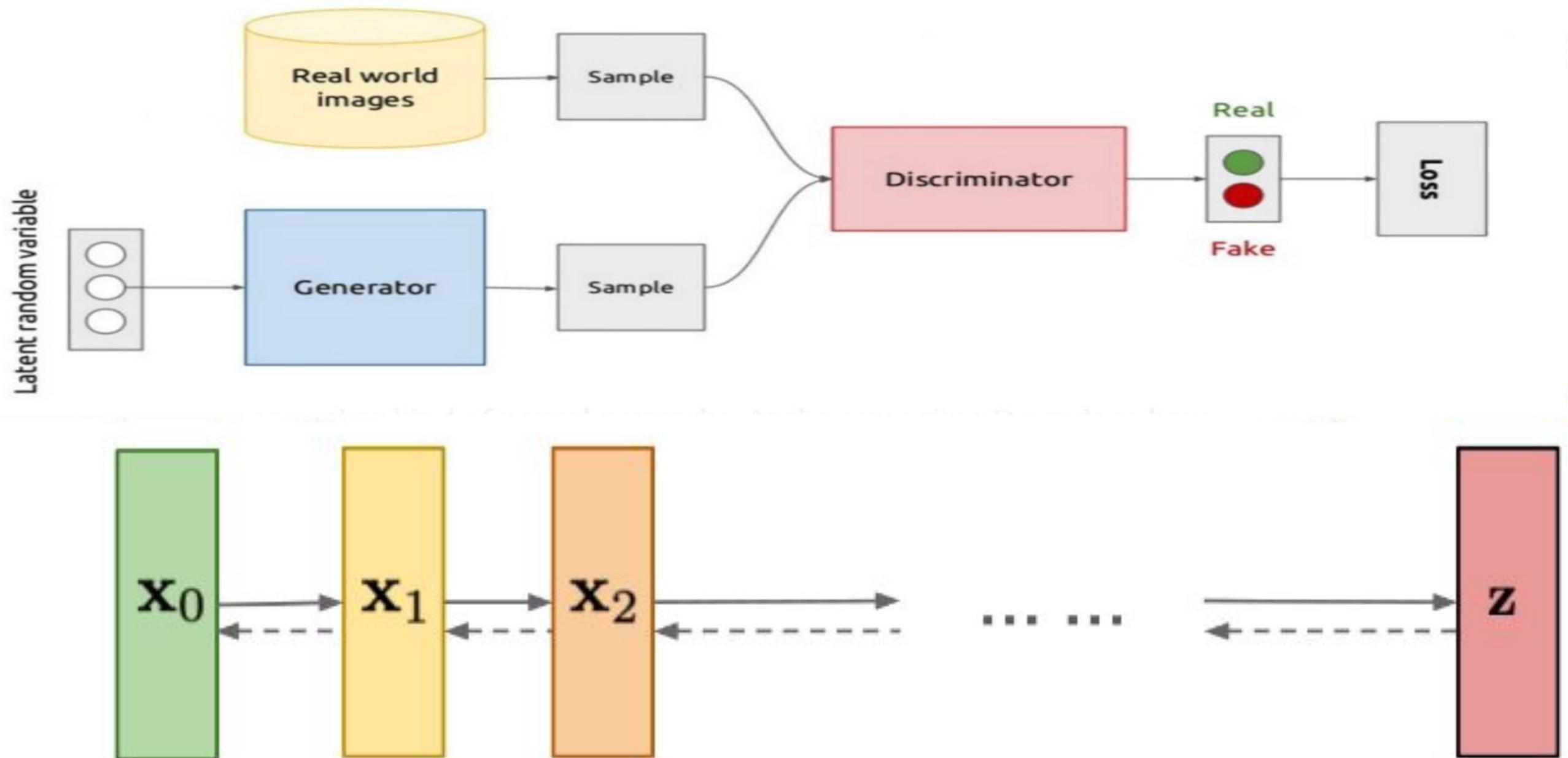
We are delighted to announce the public release of Stable
Diffusion and the launch of DreamStudio Lite.



Stable Diffusion (Image to Image pipe)

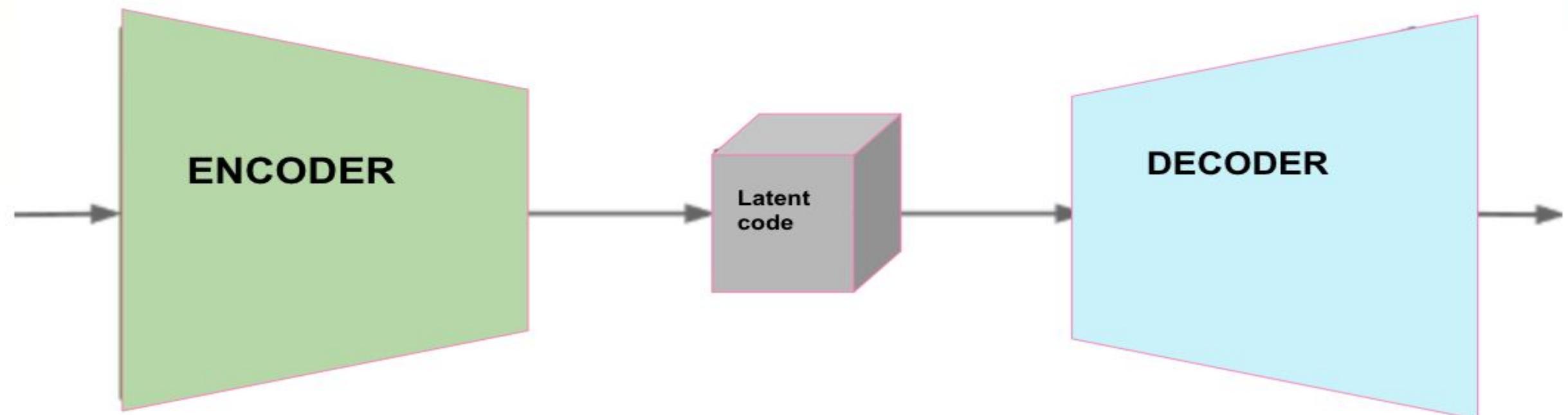


GAN: Adversarial training

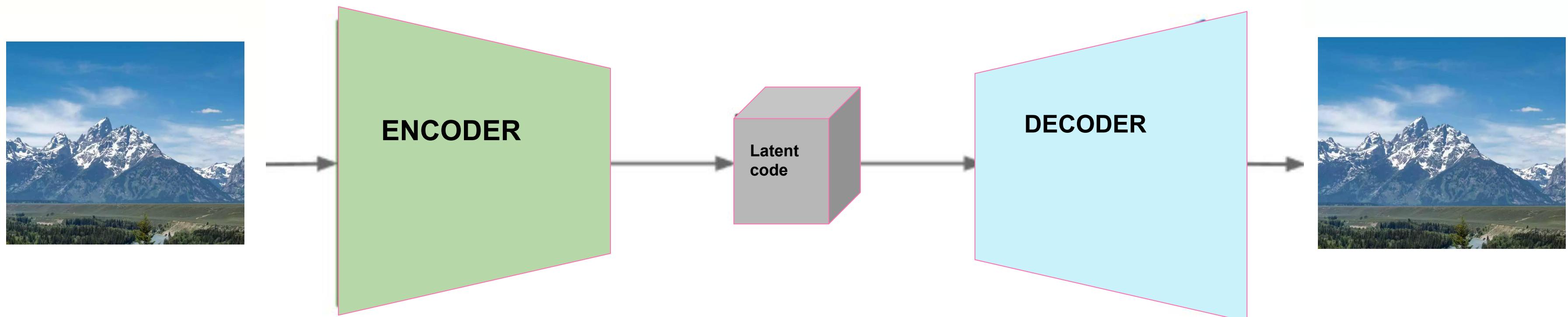


Diffusion models:
Gradually add Gaussian noise and then reverse

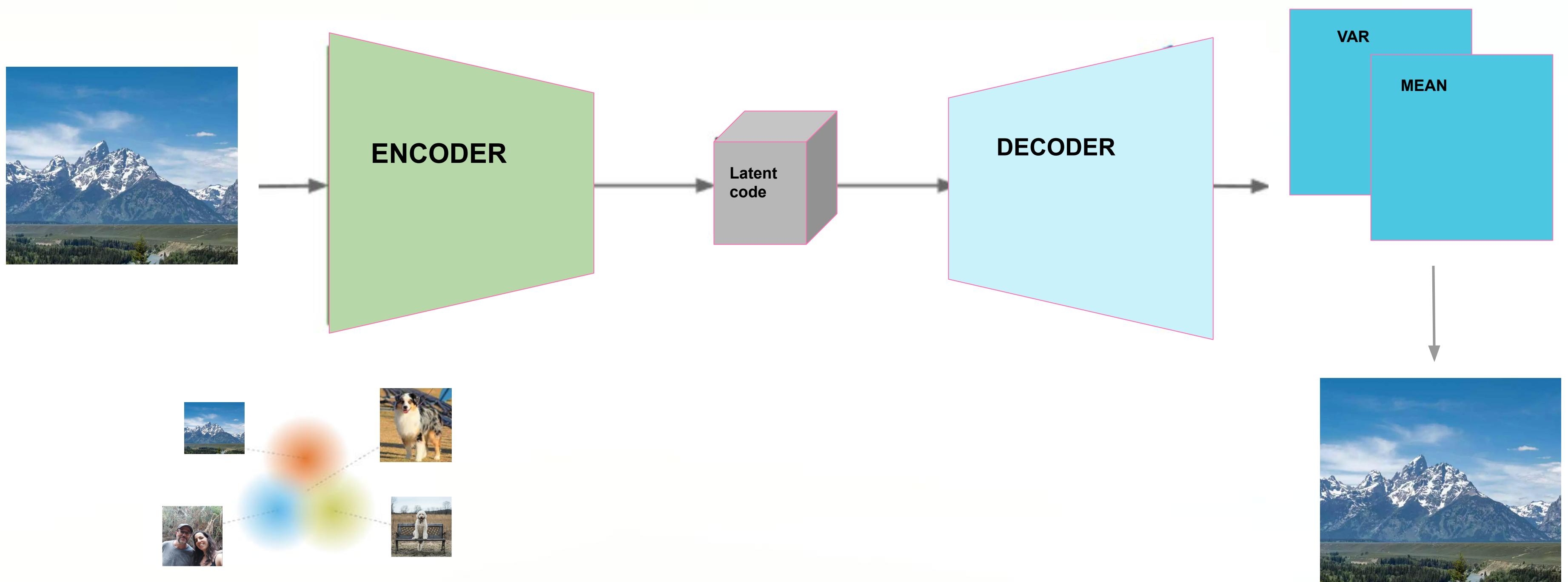
Variational Autoencoder:



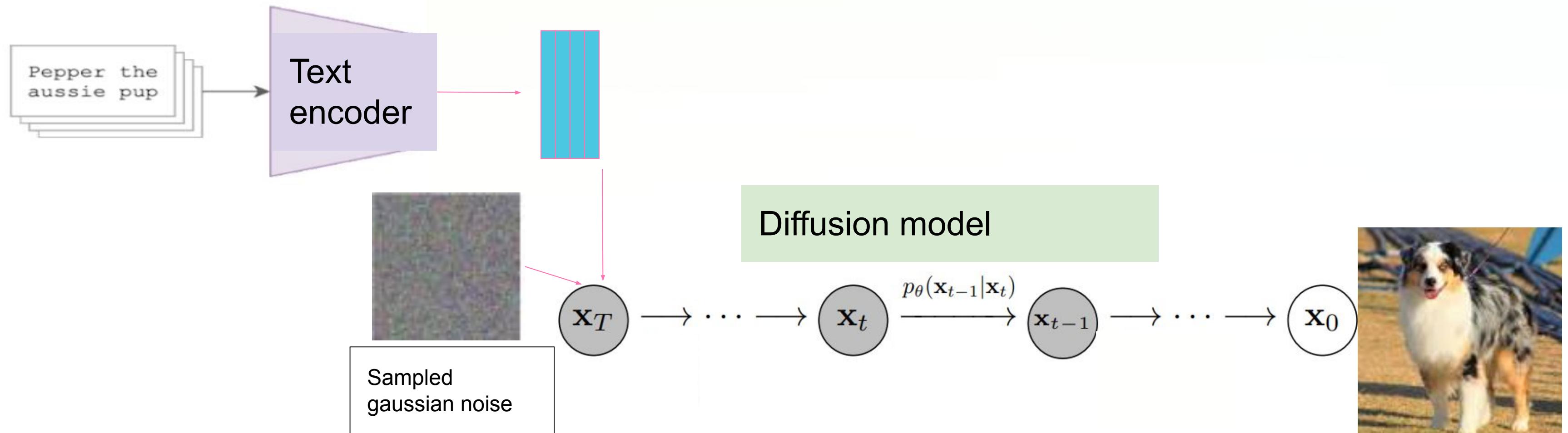
AutoEncoder



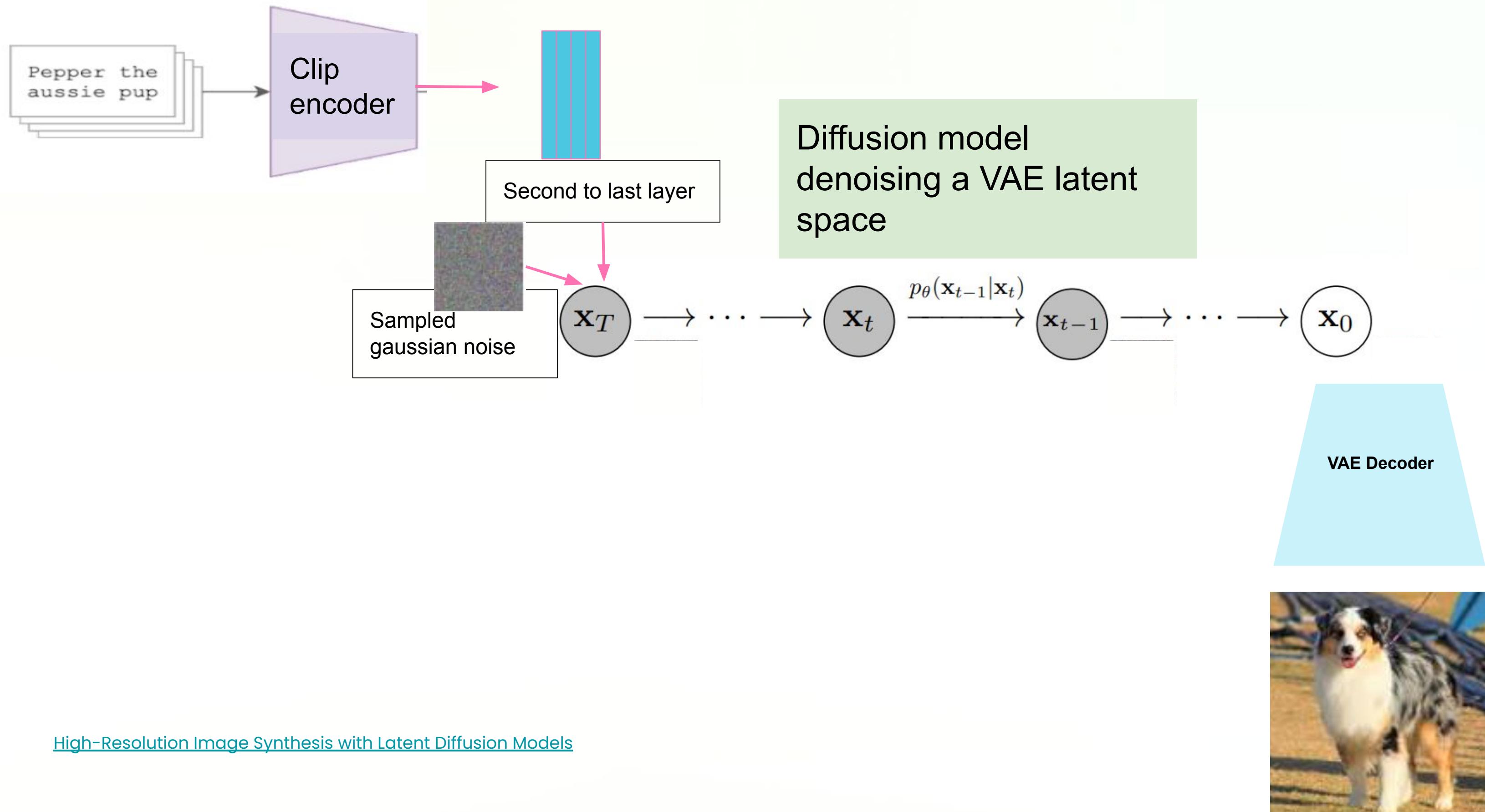
Variational AutoEncoder



Original Text to Image pipe

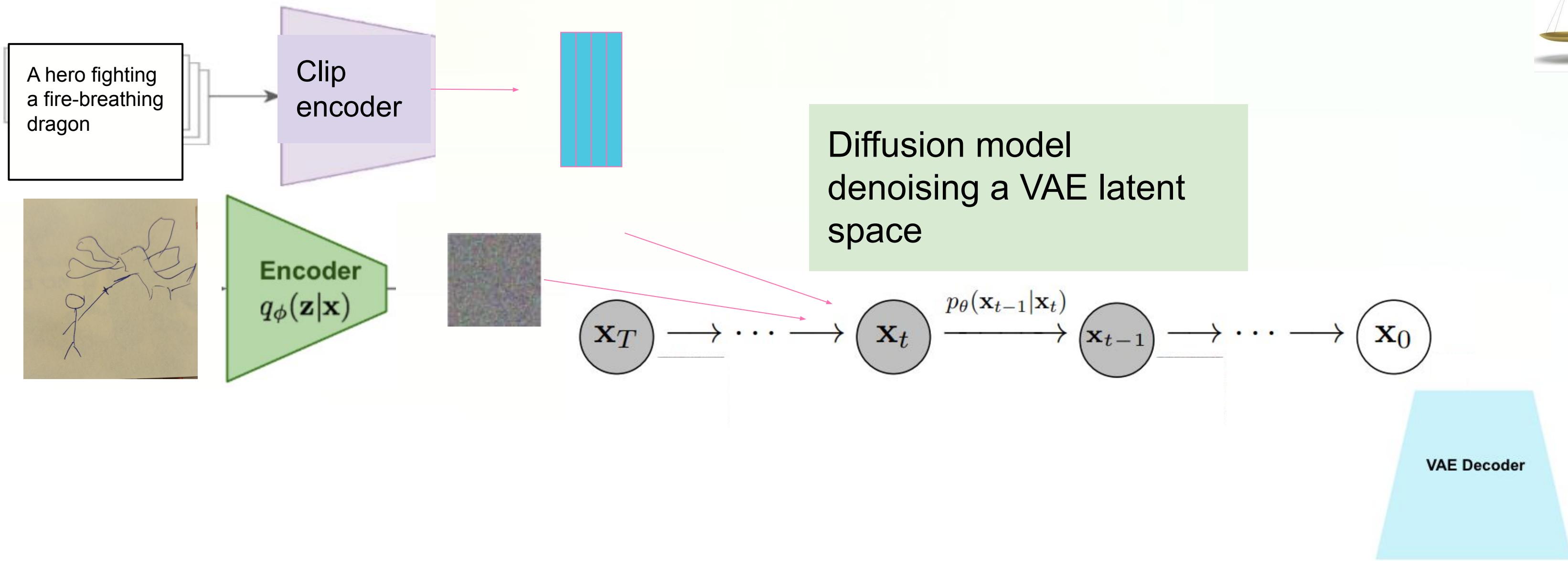


Stable Diffusion Text to Image pipe

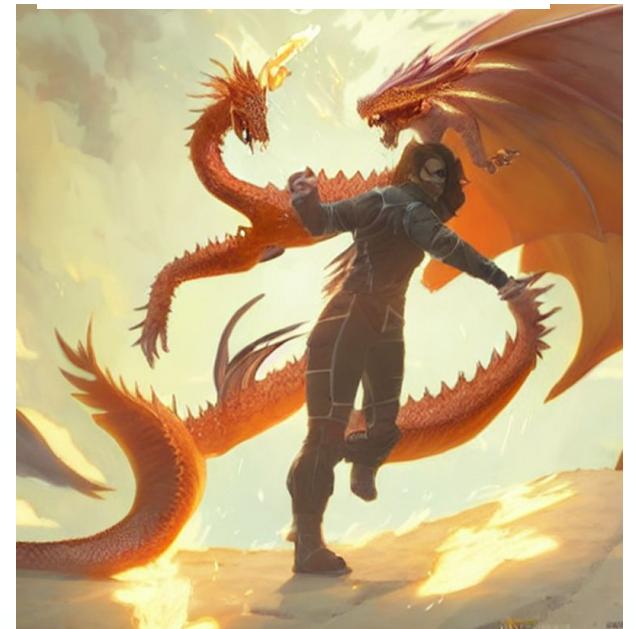


Editing pipelines

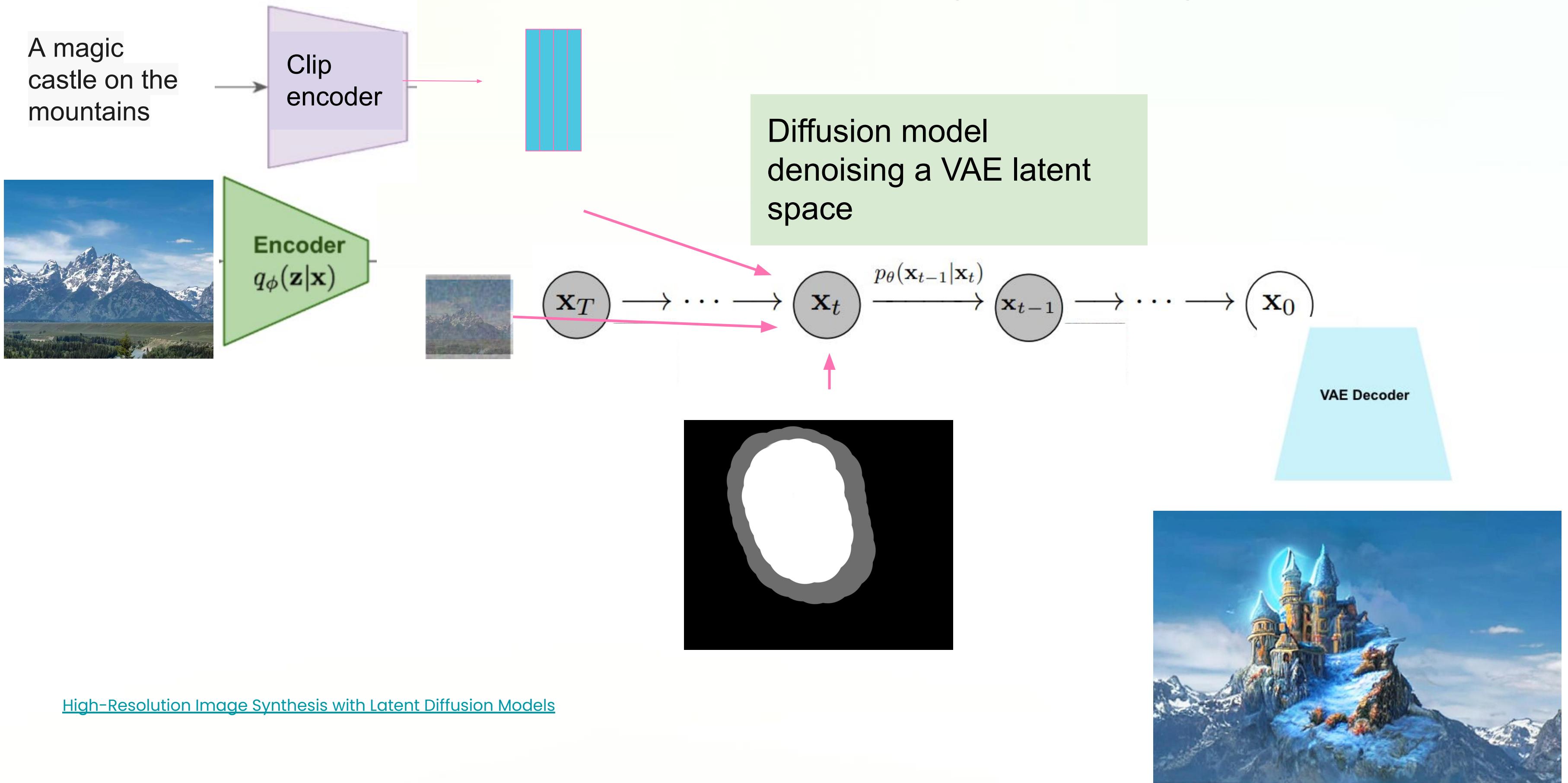
Stable Diffusion Image to Image pipe



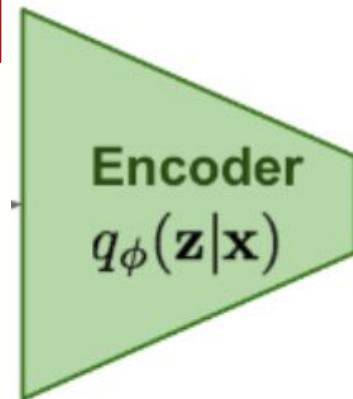
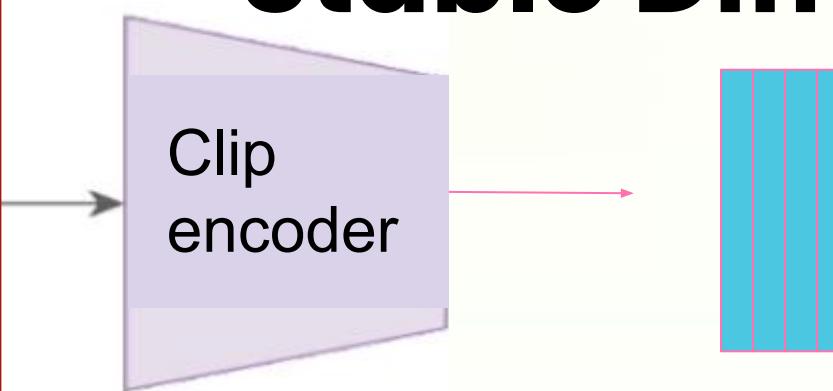
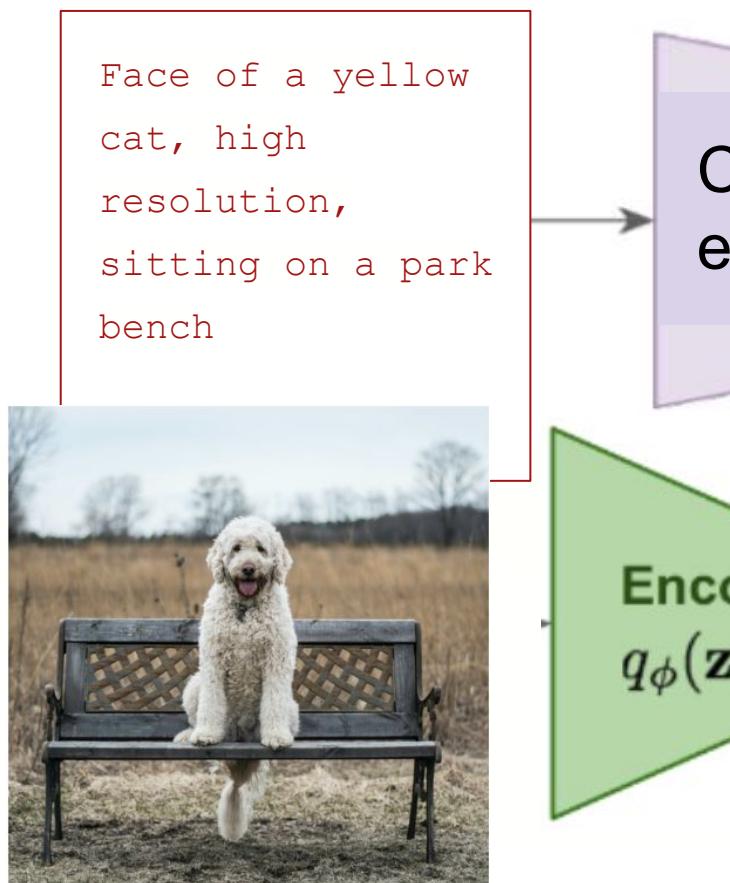
[High-Resolution Image Synthesis with Latent Diffusion Models](#)



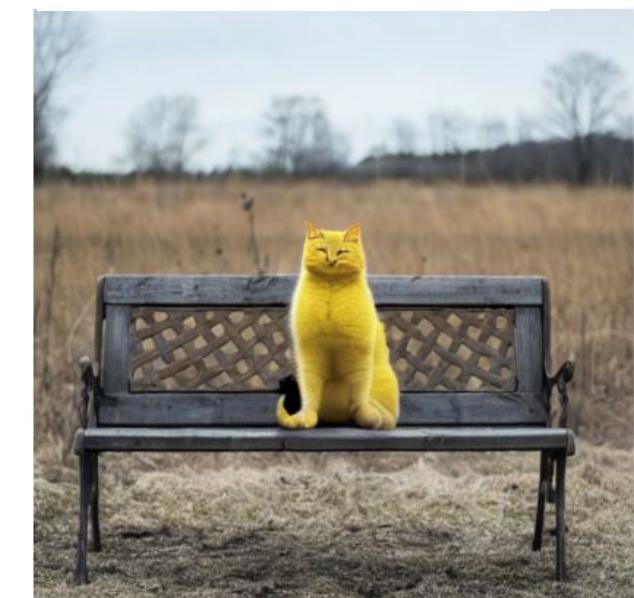
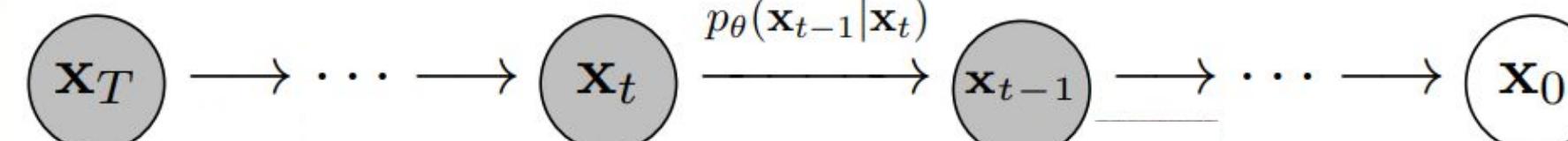
Stable Diffusion Inpainting pipe (legacy)



Stable Diffusion Inpainting pipe

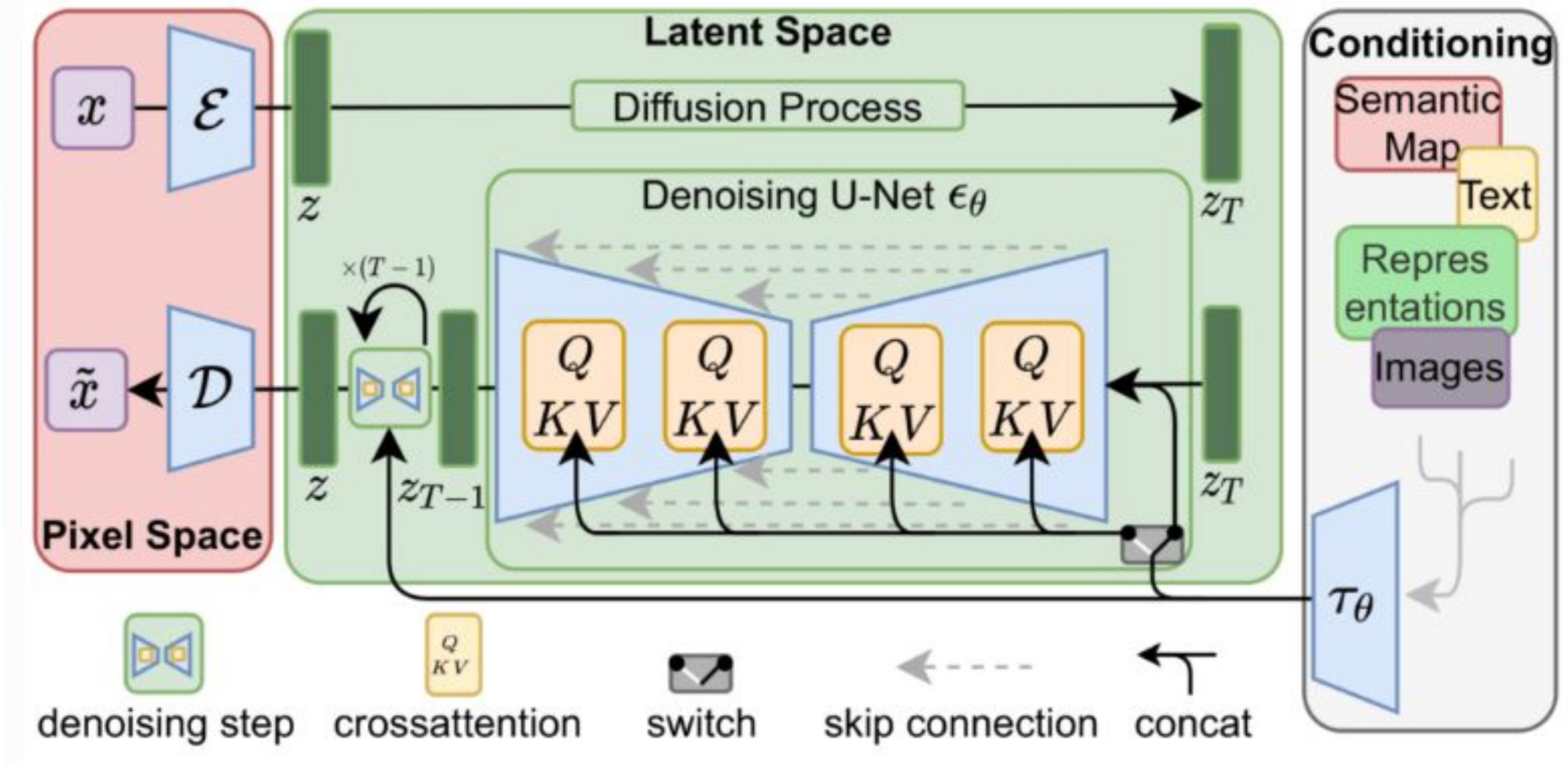


Inpainting finetuned diffusion model denoising a VAE latent space



Inpainting





Useful gits to follow

<https://github.com/huggingface/diffusers>

<https://github.com/AUTOMATIC1111/stable-diffusion-webui>

Good image 2 image example:

https://colab.research.google.com/github/patal-suraj/Notebooks/blob/master/image_2_image_using_diffusers.ipynb

Questions?

Hands on workshop

1. Go to :
<https://github.com/Naomi-Ken-Korem/text-to-image-datatalks>
2. Optional: Class free guidance trick-
Glide_text2image.ipynb
3. Stable diffusion editing pipelines
Stable_diffusion_editing_pipes.ipynb
OR stable_diffusion_pipes_full_version.ipynb



Takeaway

1. Multimodel domain inherits ideas from computer vision, and nlp.
2. We covered a lot of architectures, and technical ideas. Combining those existing ideas leads to great progress
3. When you read papers and tweets on the topic try to look for the building blocks that we mentioned today



Thank you!