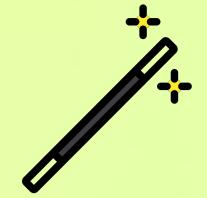




The Text2Image Magic



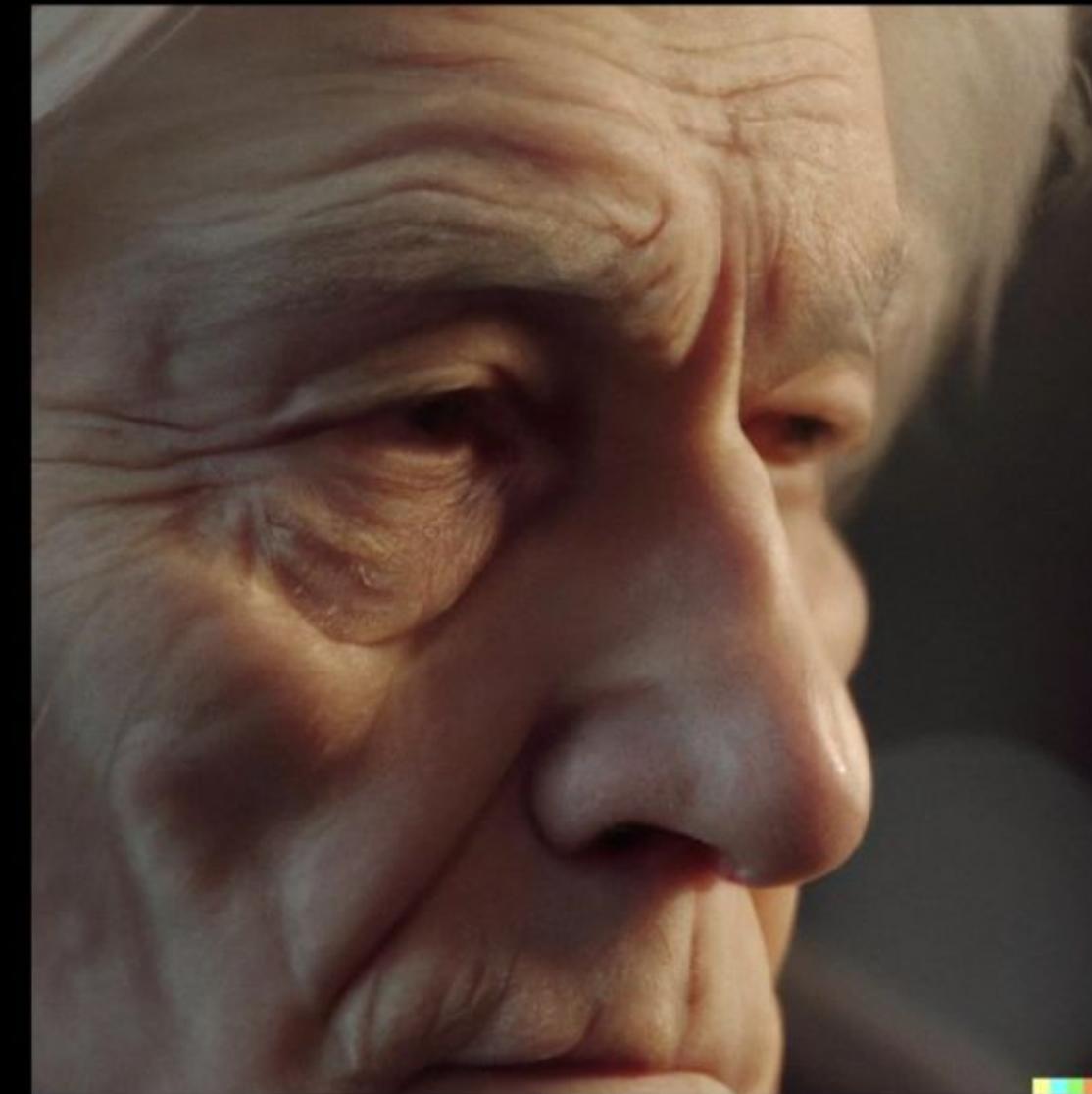
Naomi Ken Korem
Researcher



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

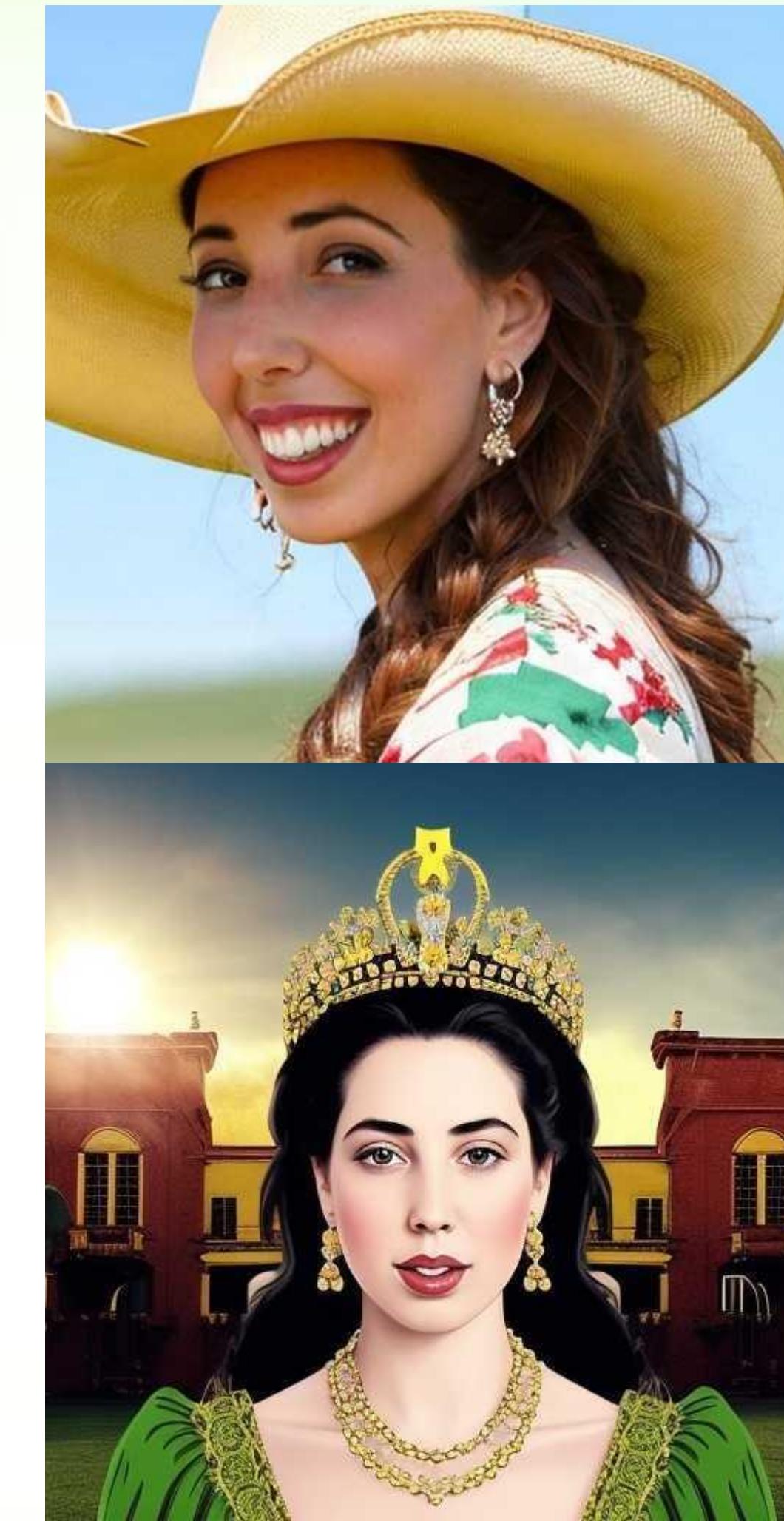
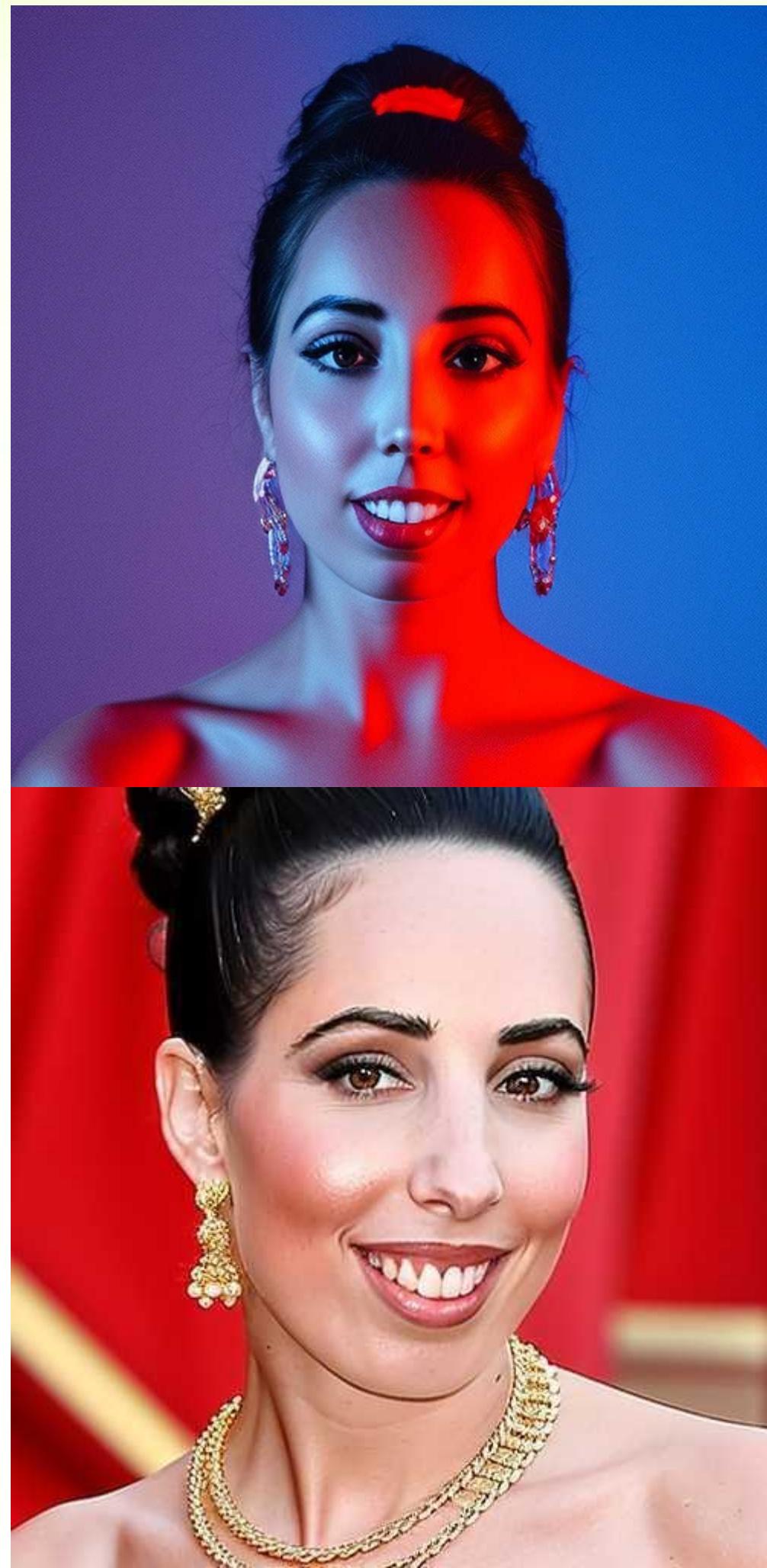


LAION
@laion_ai

...

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."





"Add fireworks to the sky"



"Replace the fruits with cake"

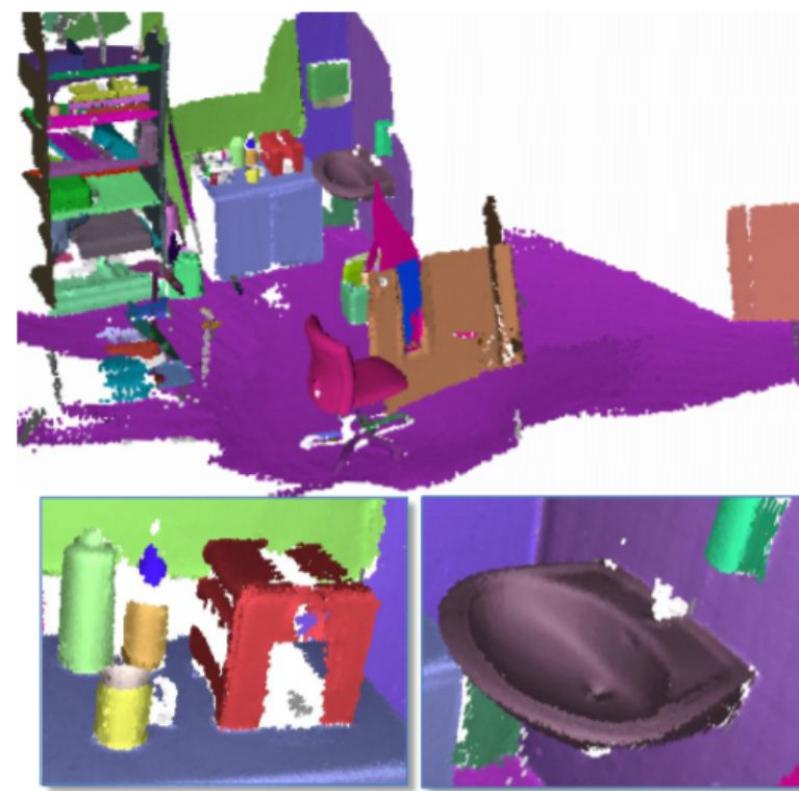


"Turn it into a still from a western"

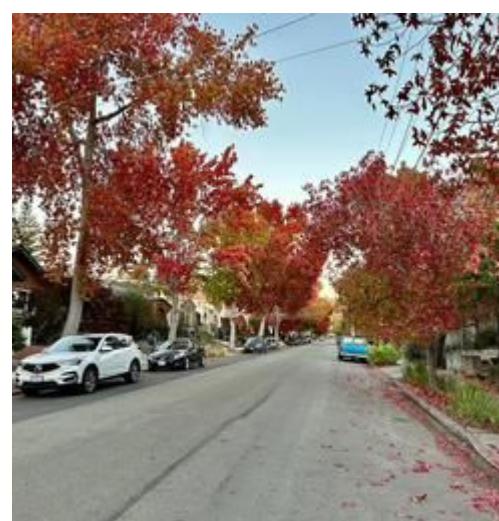


"Make his jacket out of leather"

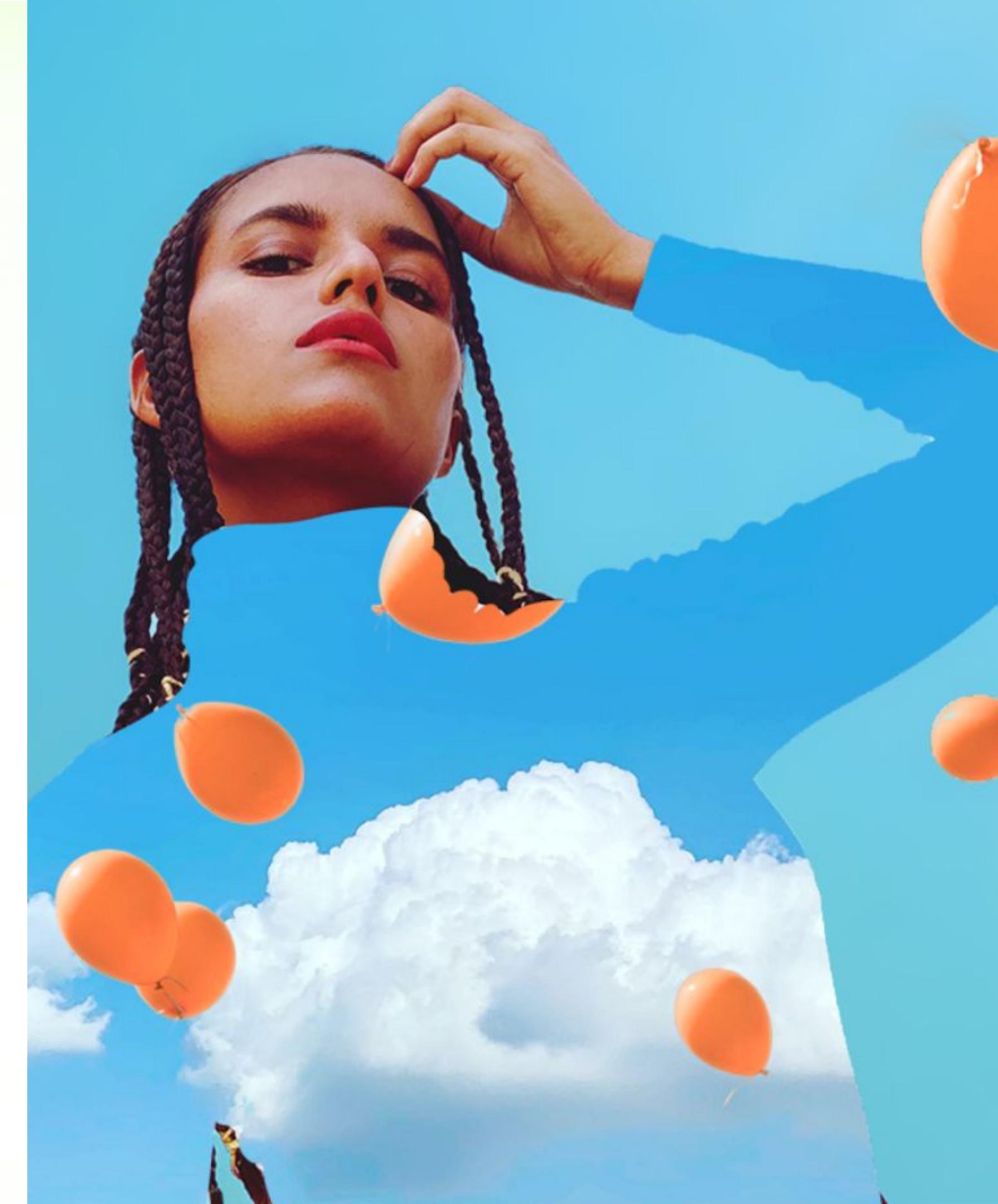
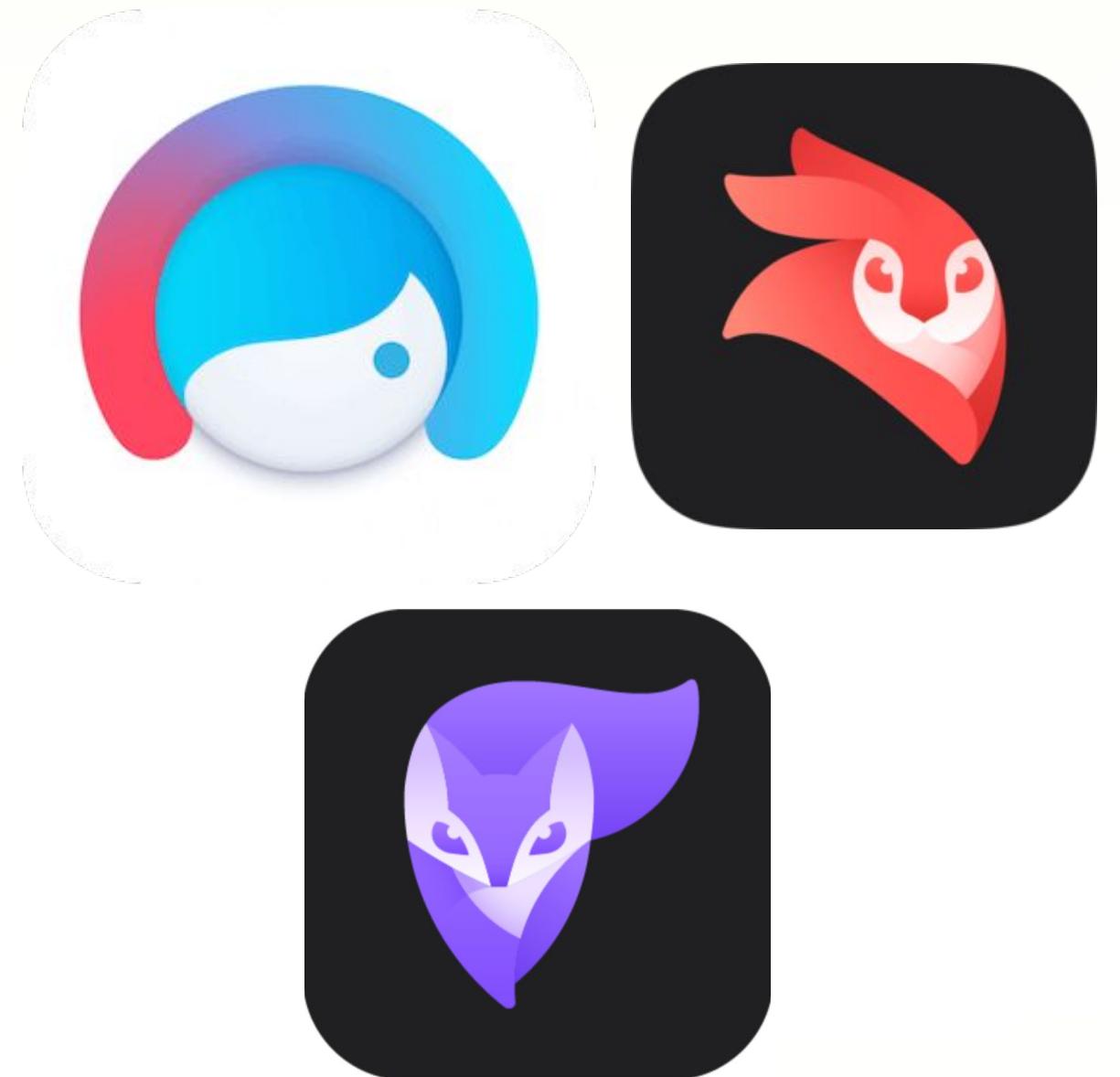




Naomi Ken Korem



Lighticks



Timeline

OpenAI publish DALL-E paper and open source CLIP model

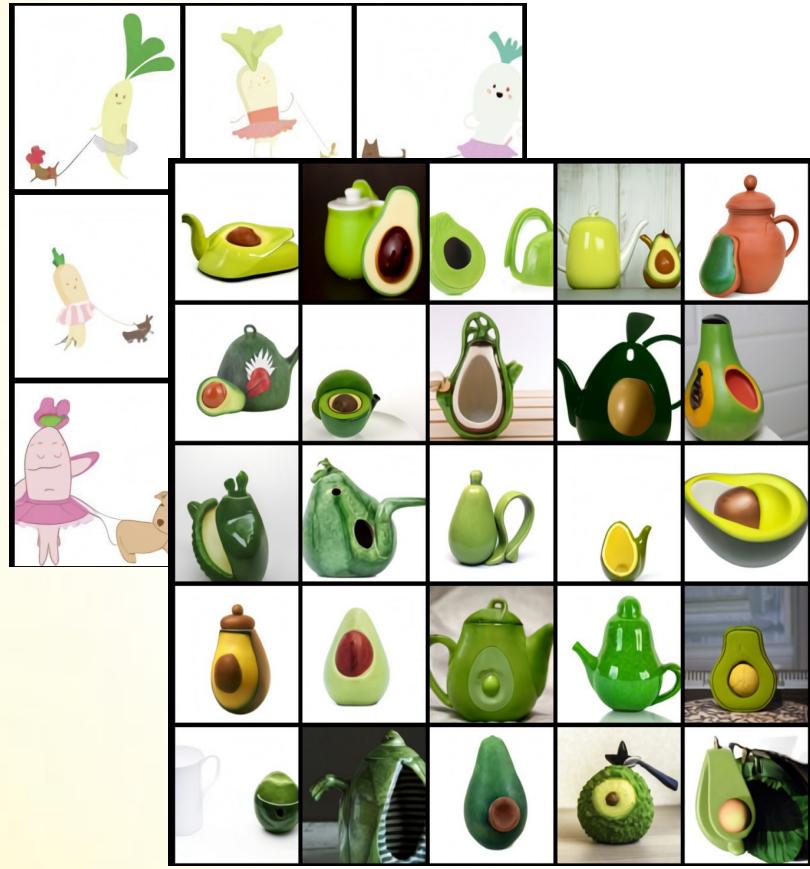
DALLE-2 is published by OpenAI and hype hits mainstream

Google Brain releases Imagen

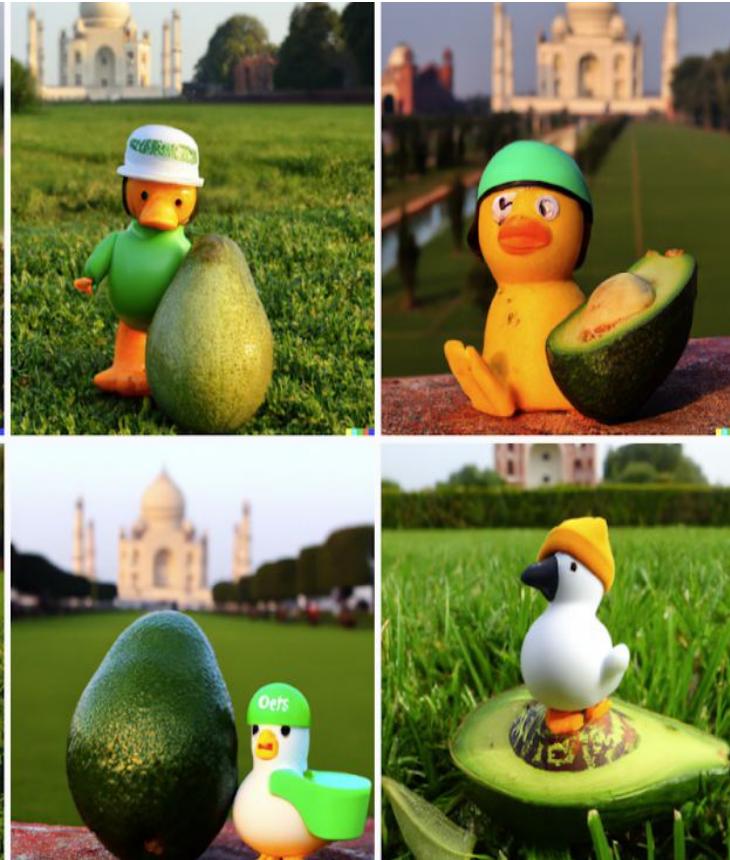
Stability.ai releases Stable Diffusion, open license

A LOT of works based on latest released models

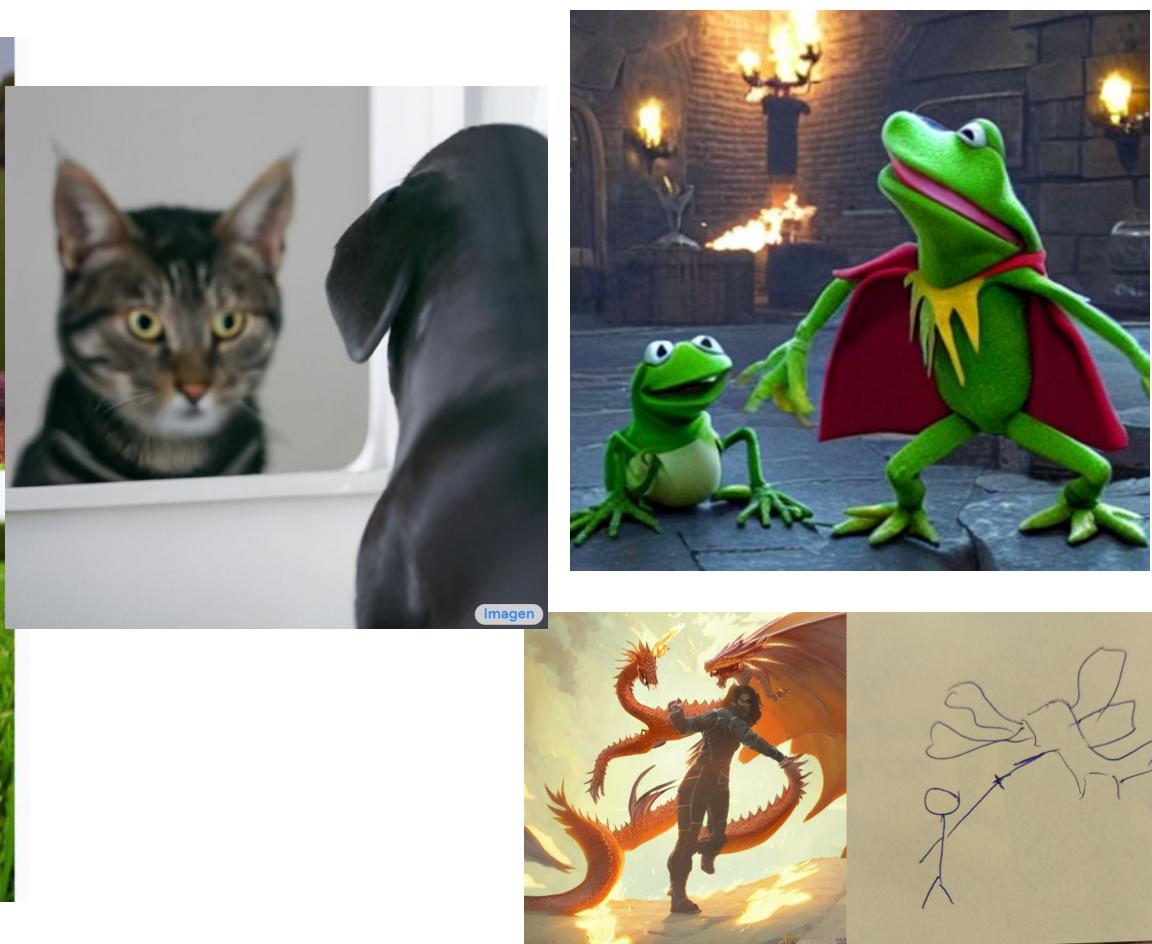
1-2021



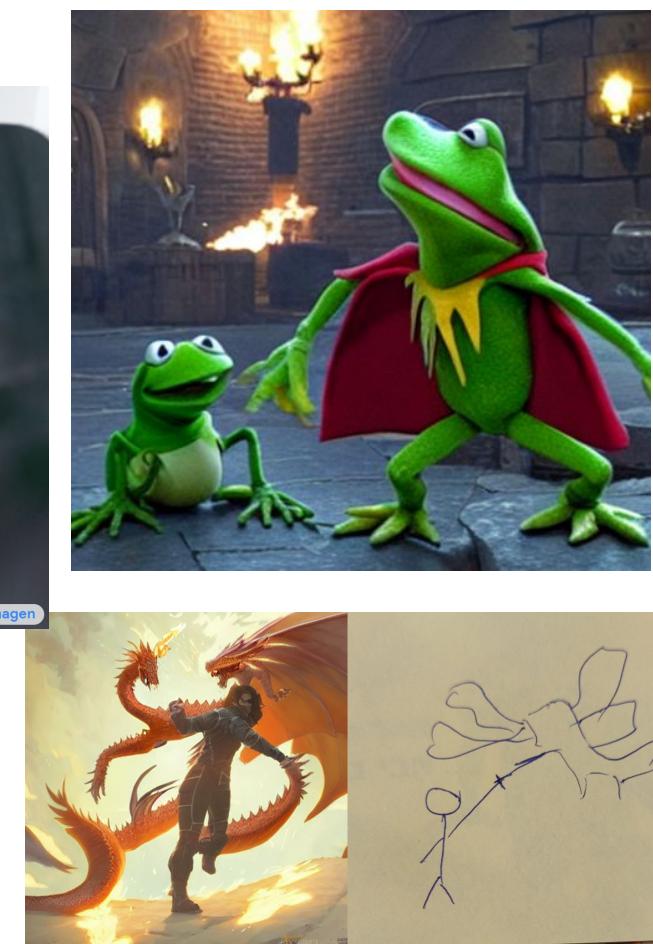
04-2022



5-2022



8-2022



Today



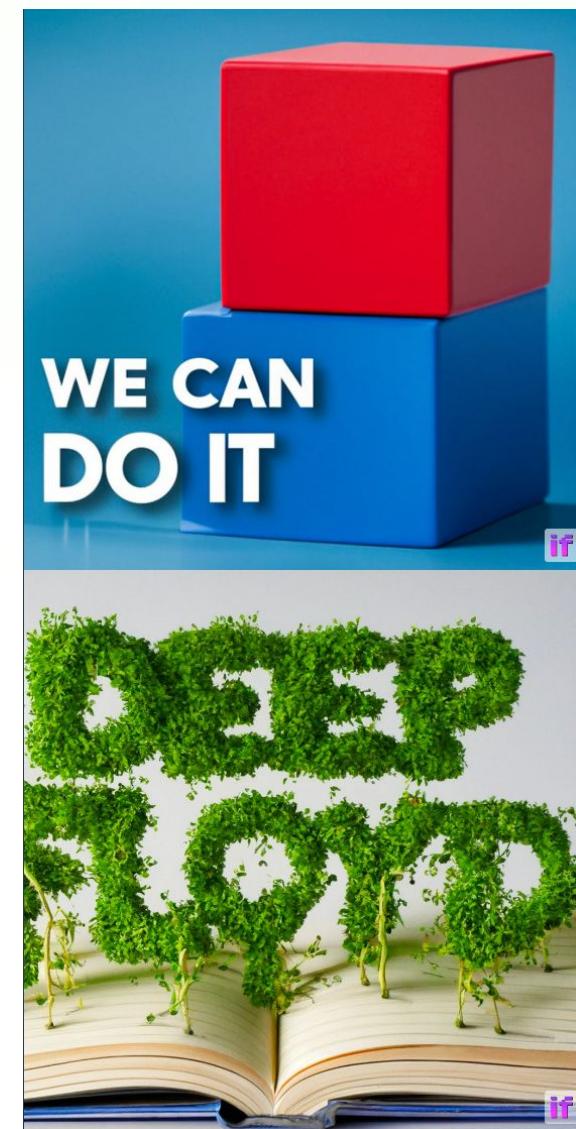
DeepFloyd (just came out)



capybara holding a neon sign with text that reads "capybara podcast", a professional photo of a capybara podcasting, capybara chimera animorph, transformer animal, anamorphic, 8k, 4k, 85 mm, f2. ... by rr.disc0



4 bottles of wine next to each other labeled



System 1



Fast, intuitive and
emotional

System 2



Slow, conscious
and effortful

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 -

1. Transformers
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.



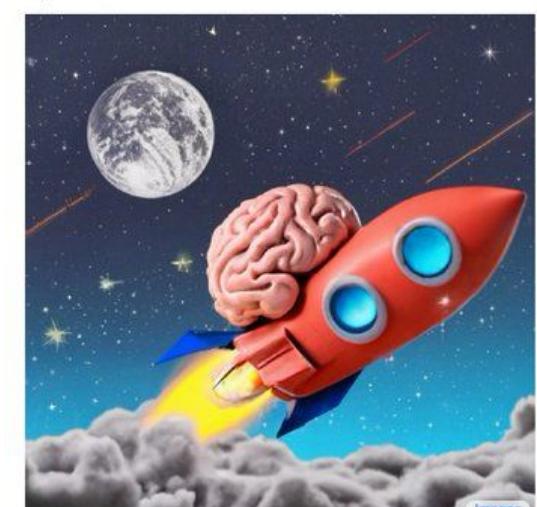
Teddy bears swimming at the Olympics 400m Butterfly event.



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



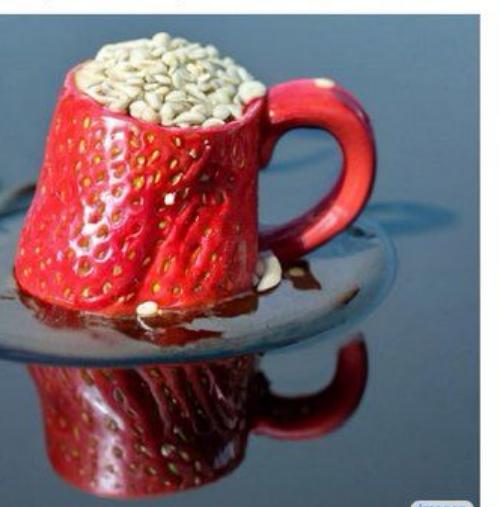
A cute sloth holding a small treasure chest. A bright golden glow is coming from the chest.



A brain riding a rocketship heading towards the moon.



A dragon fruit wearing karate belt in the snow.



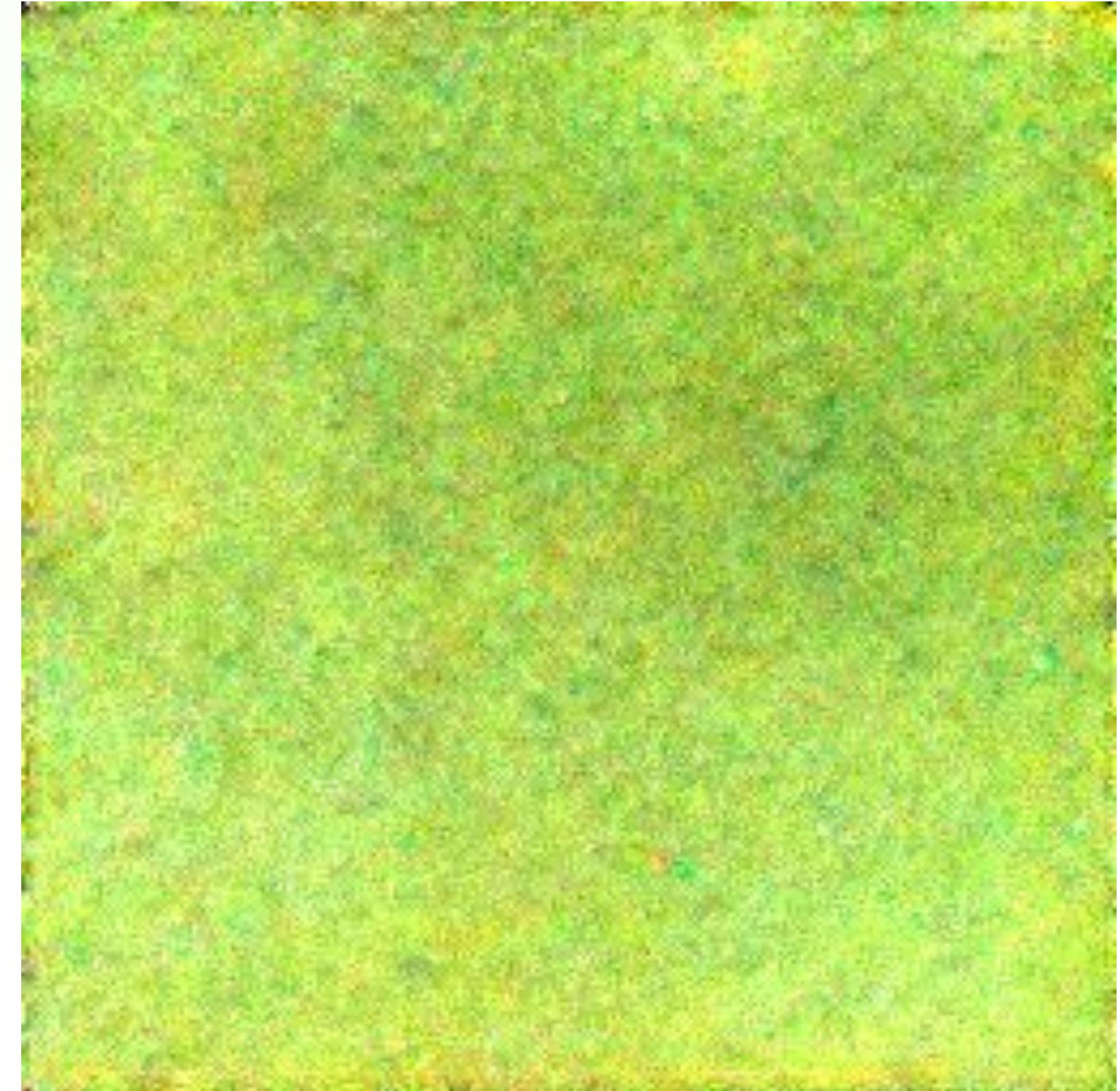
A strawberry mug filled with white sesame seeds. The mug is floating in a dark chocolate sea.

Agenda

1. Diffusion models
2. Hands on #1 training a diffusion model
3. Text 2 Image
4. Hands on #2 Stable diffusion pipelines

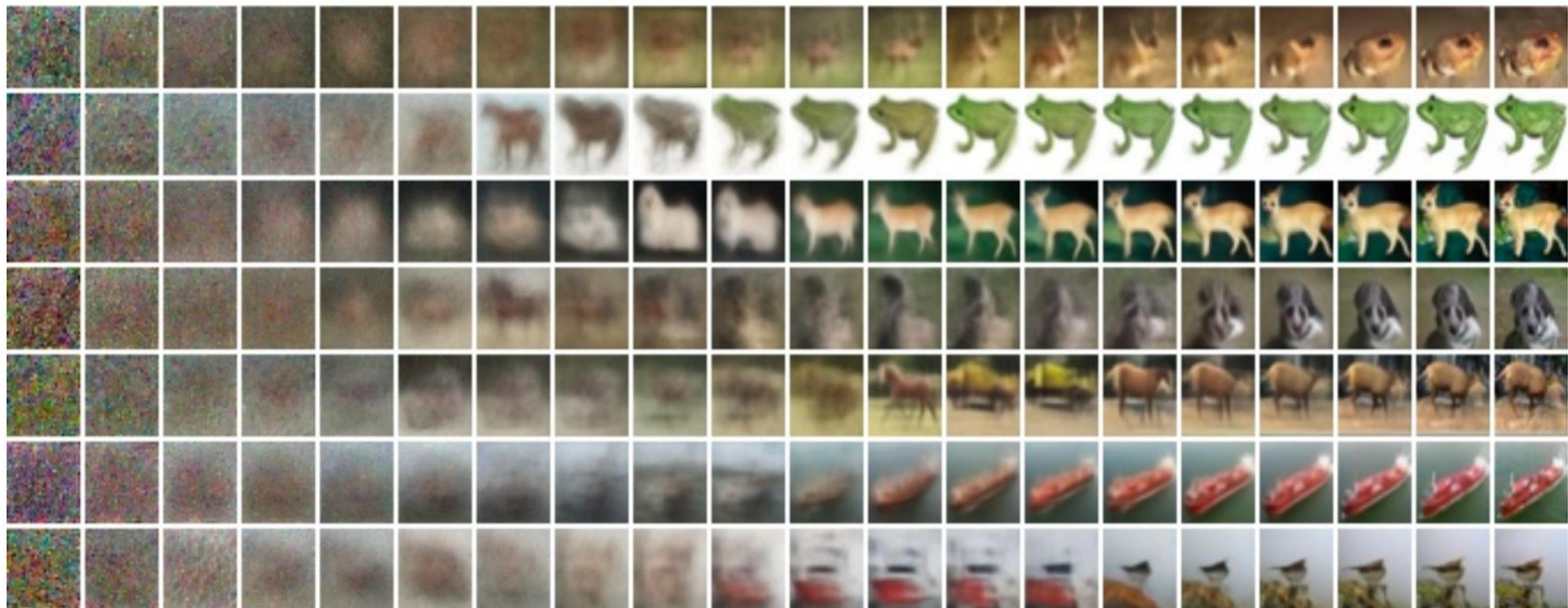
Diffusion models

- Model that learns to **Iteratively de-noise** until image is formed.
-
- 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.

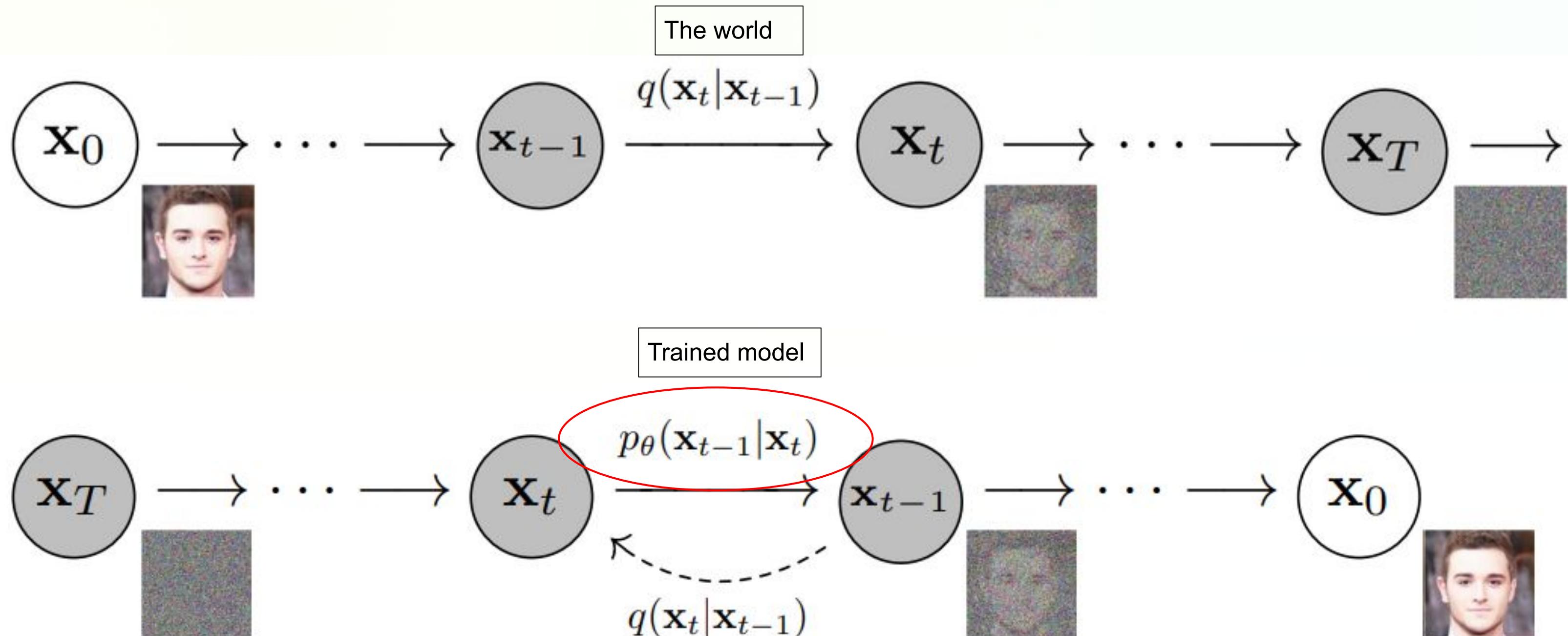


DiscoDiffusion v5, “Cheese and grapes next to a glass of wine”

Diffusion models- from noise to images in steps

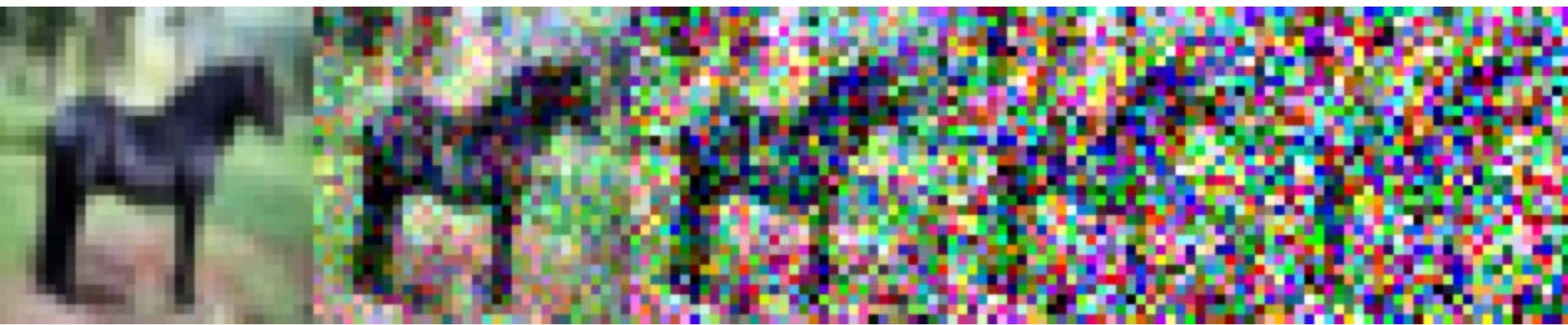


Training a diffusion model



Creating data for training

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



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

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

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

Diffusion model - The reverse step

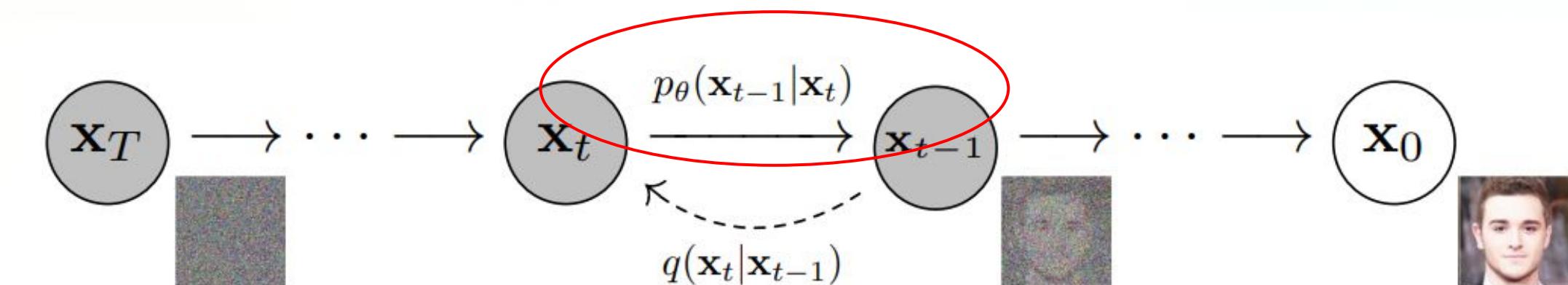
Diffusion model can be trained to predict \mathbf{x}_{t-1} given \mathbf{x}_t, t

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

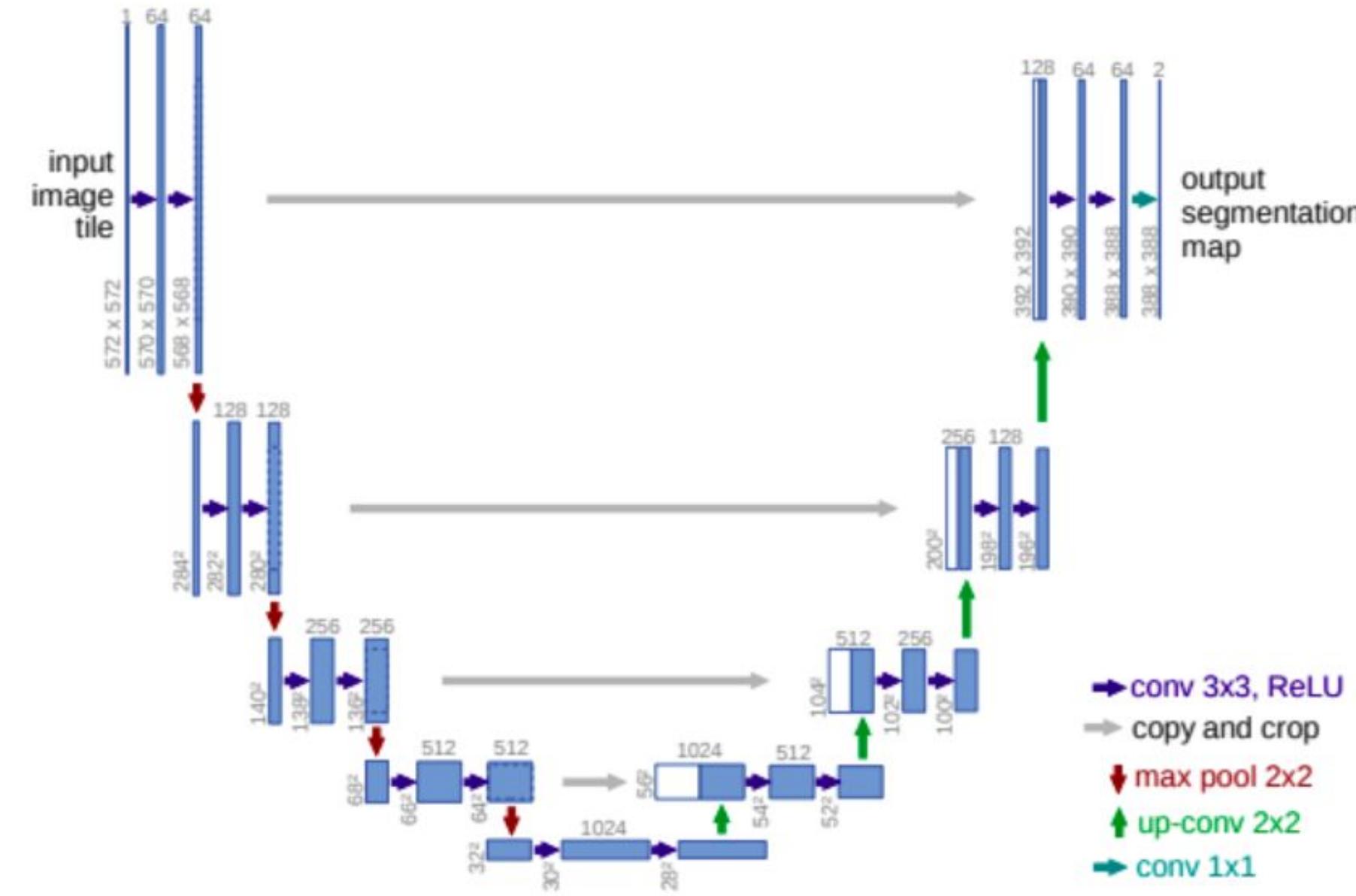


Dimensions:

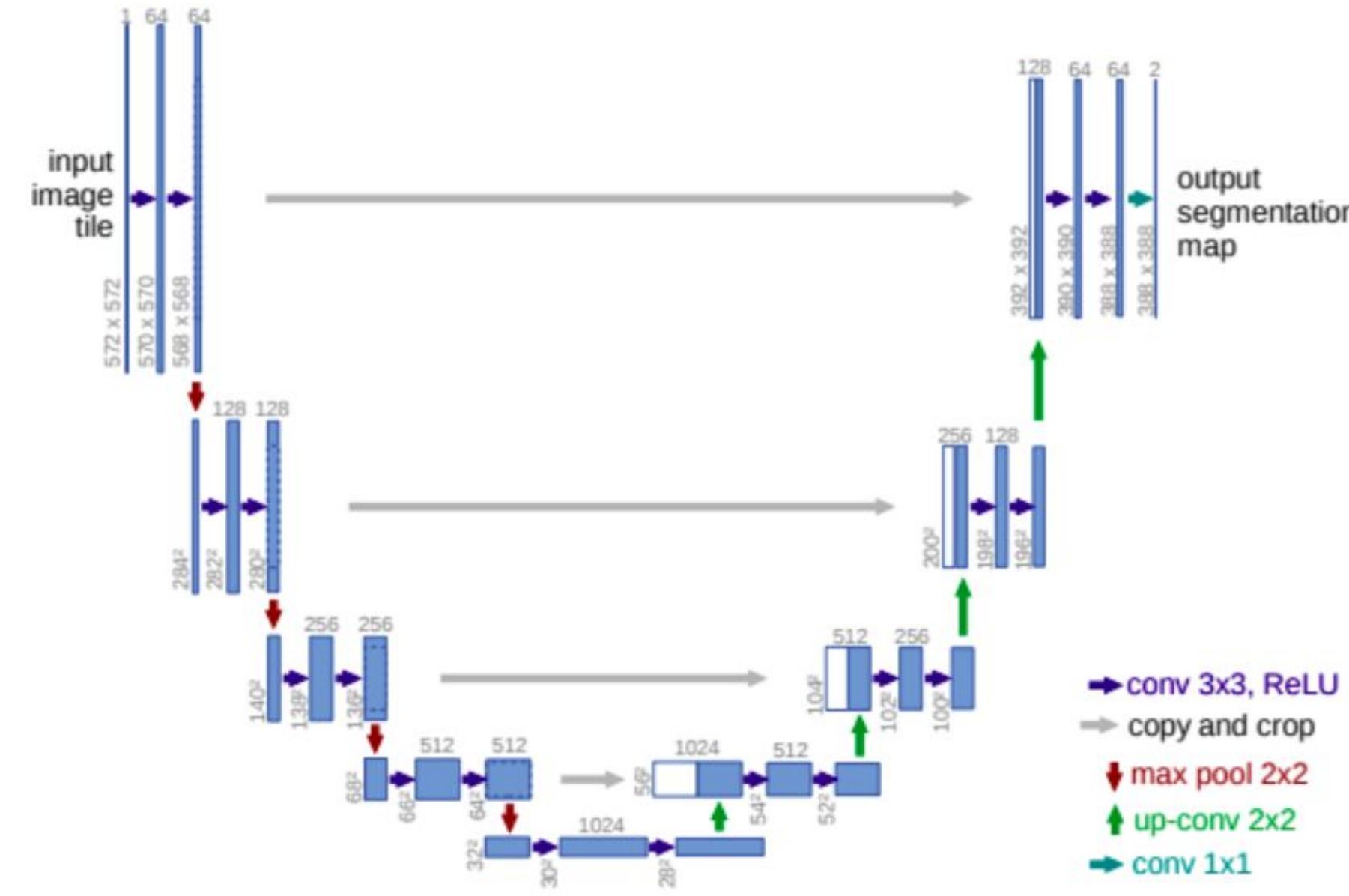
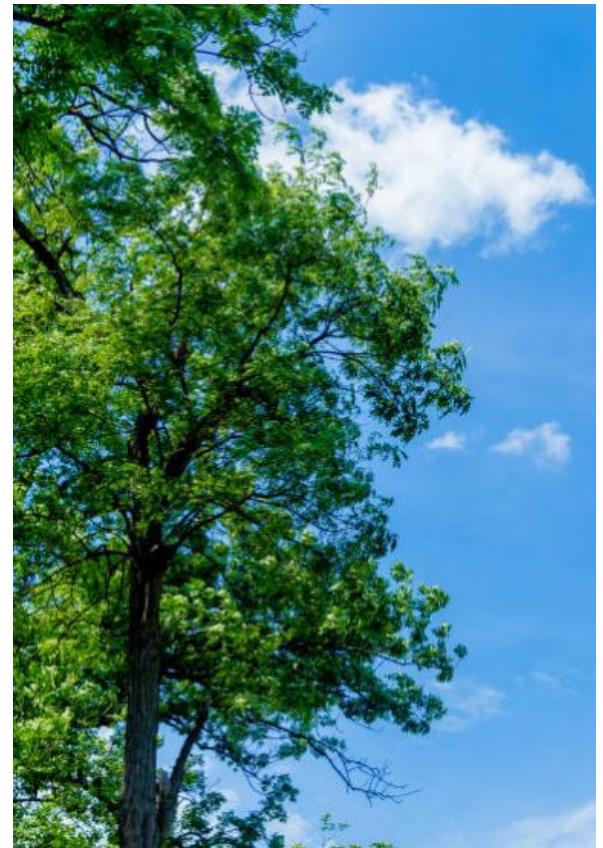
$\bar{\alpha}_t$ - scalar

\mathbf{x}_t Tensor 64x64x3

Unet Architecture

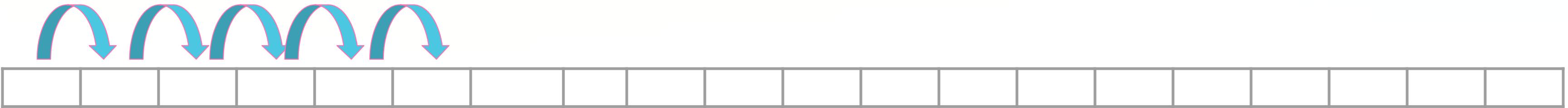
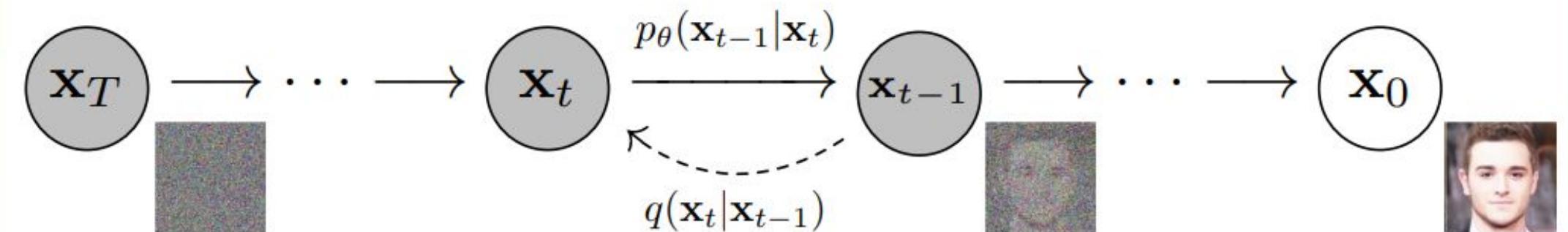


Unet Architecture



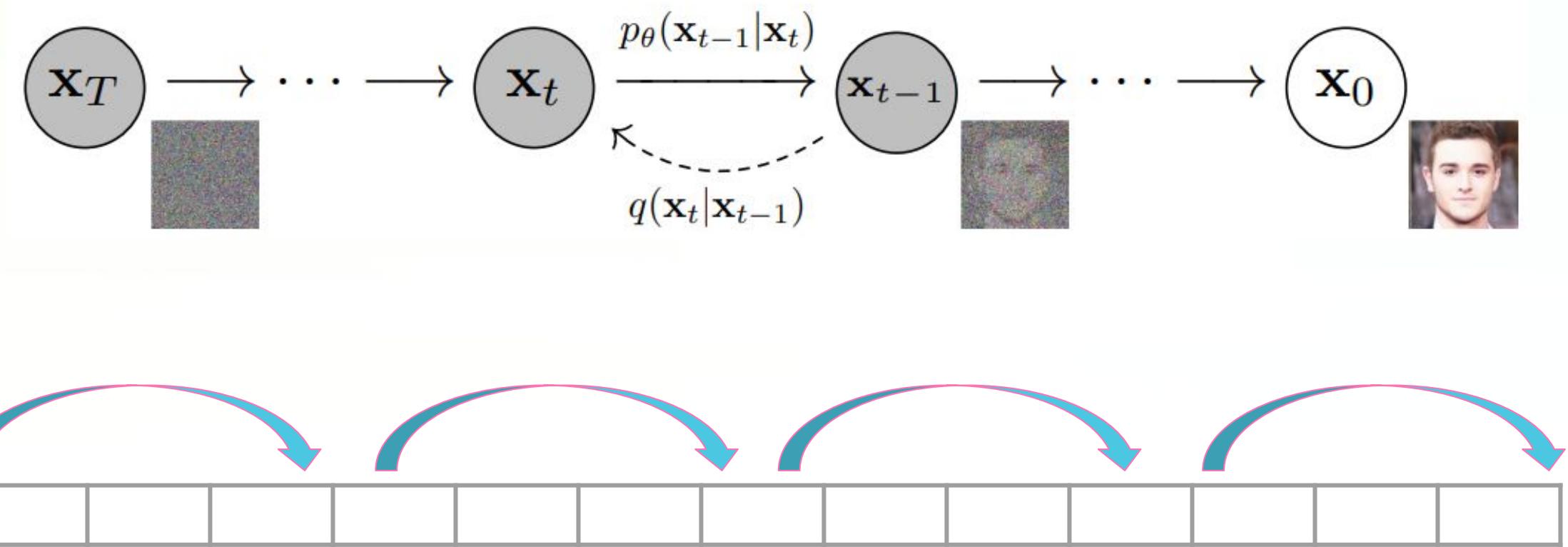
Inference

1000 steps??



Inference

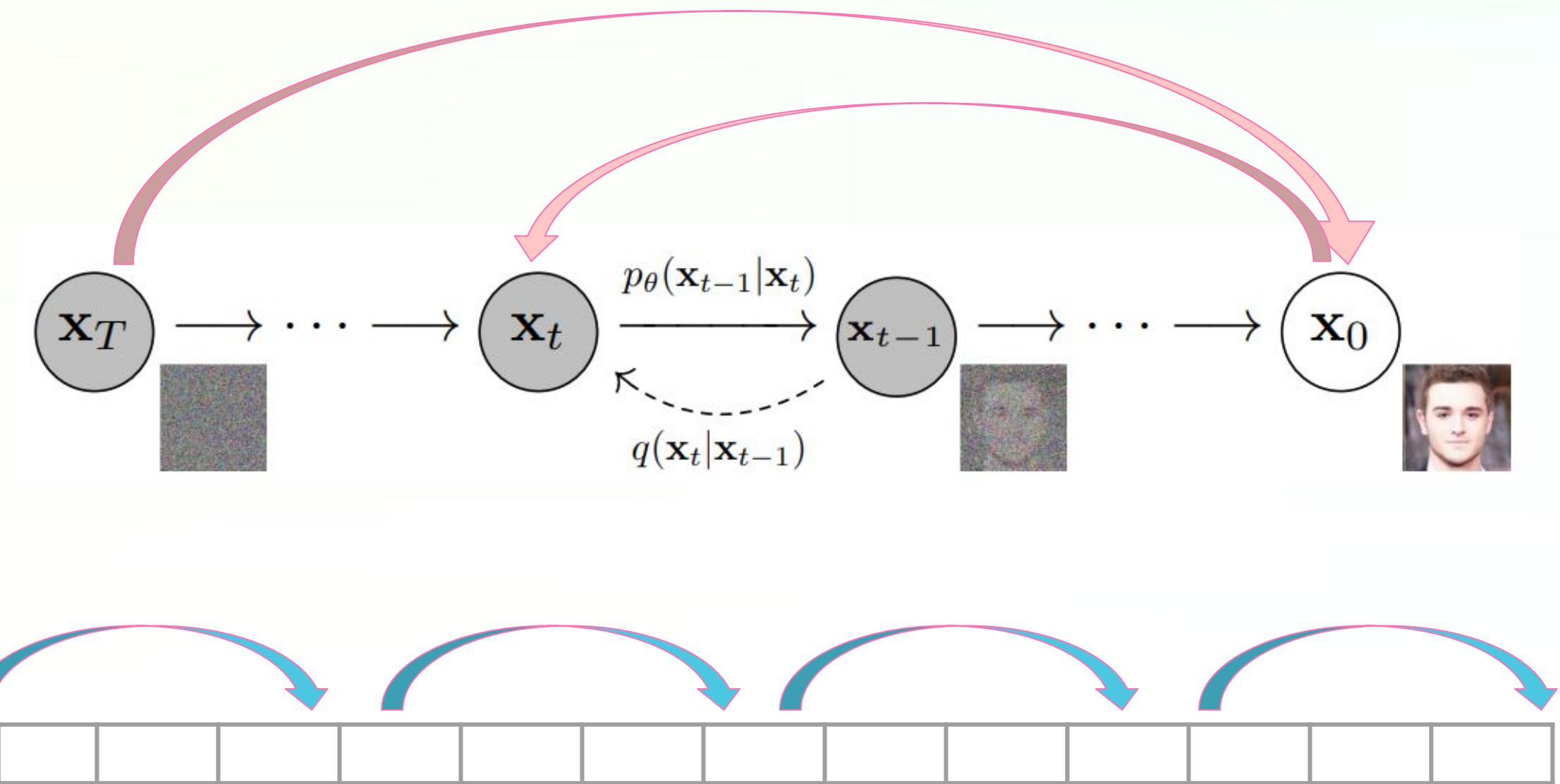
1000 steps??



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

Inference

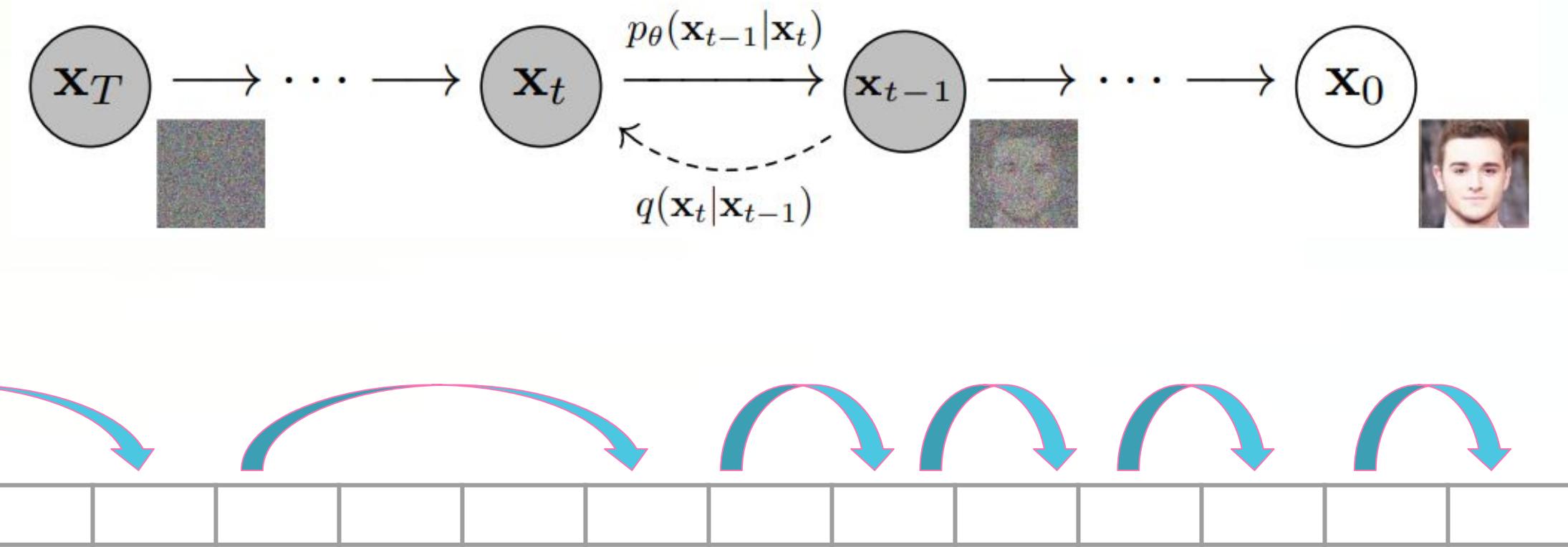
1000 steps??



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

Inference

1000 steps??



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

Questions?

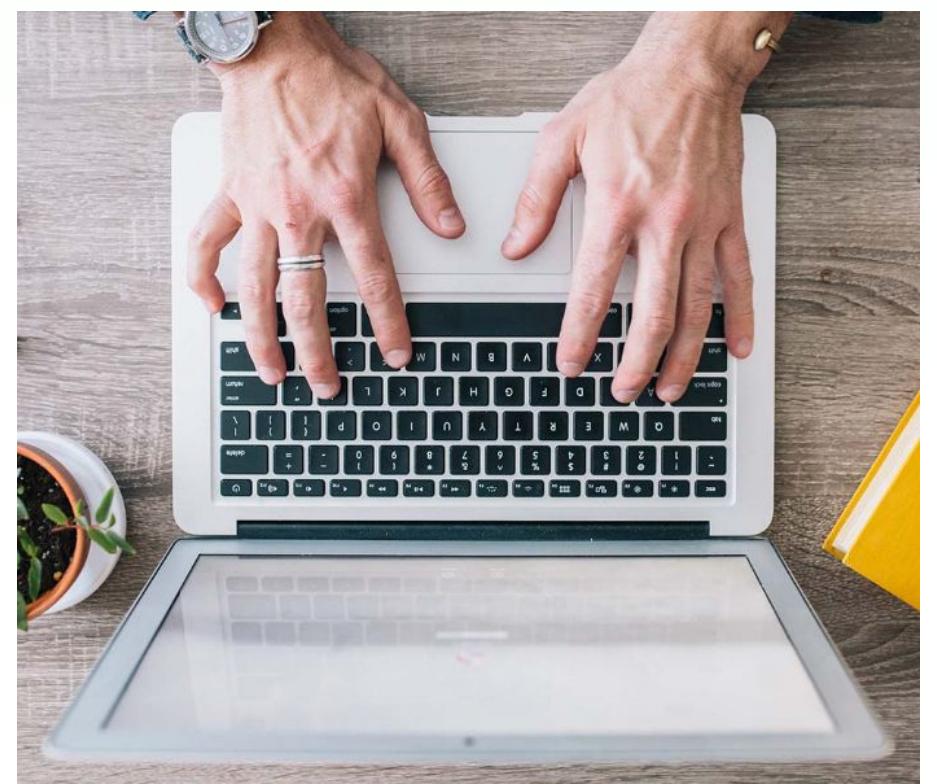
Hands on session

Hands on session (30 minutes)

Training a denoising model

1. Go to
<https://github.com/Naomi-Ken-Korem/text-to-image-presentation>

2. Open Diffusion models part 1.ipynb
(solution:
Diffusion models part 1_solution.ipynb)



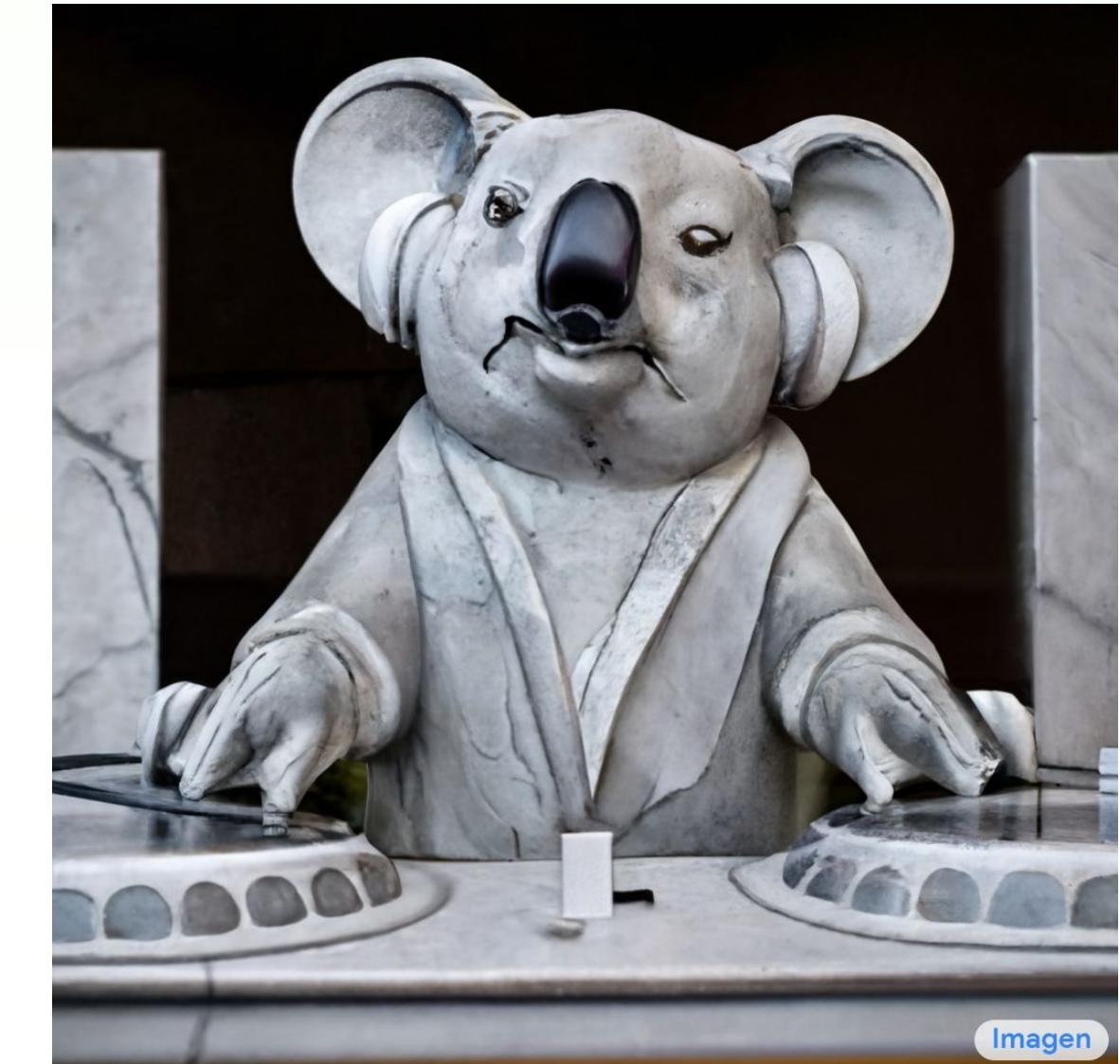
Part 2 - Text condition Diffusion models



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.

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.

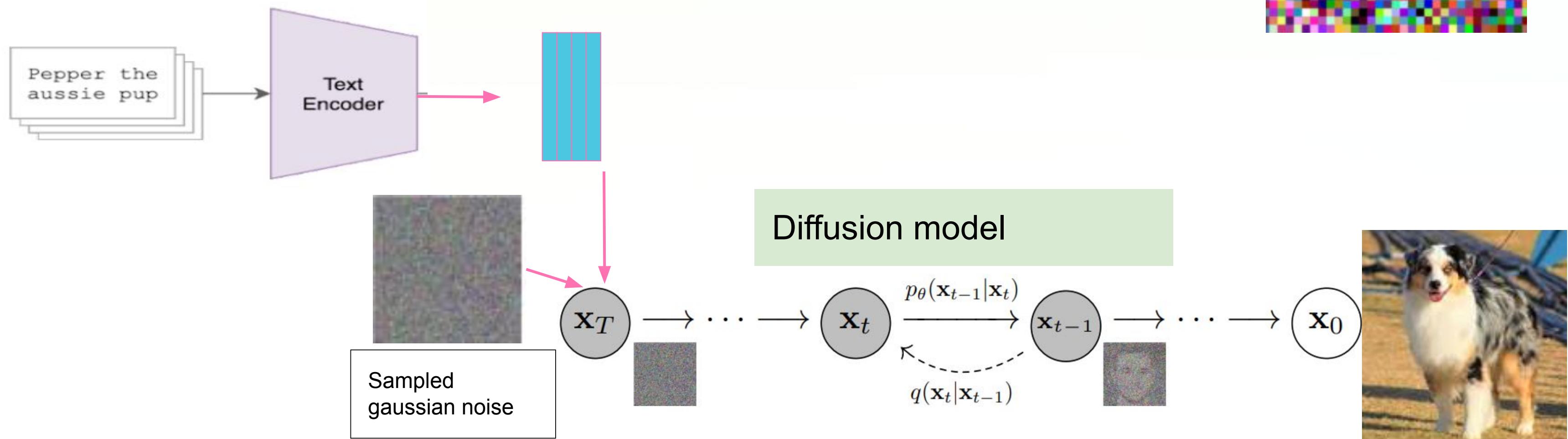


"a boat in the canals of venice"



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

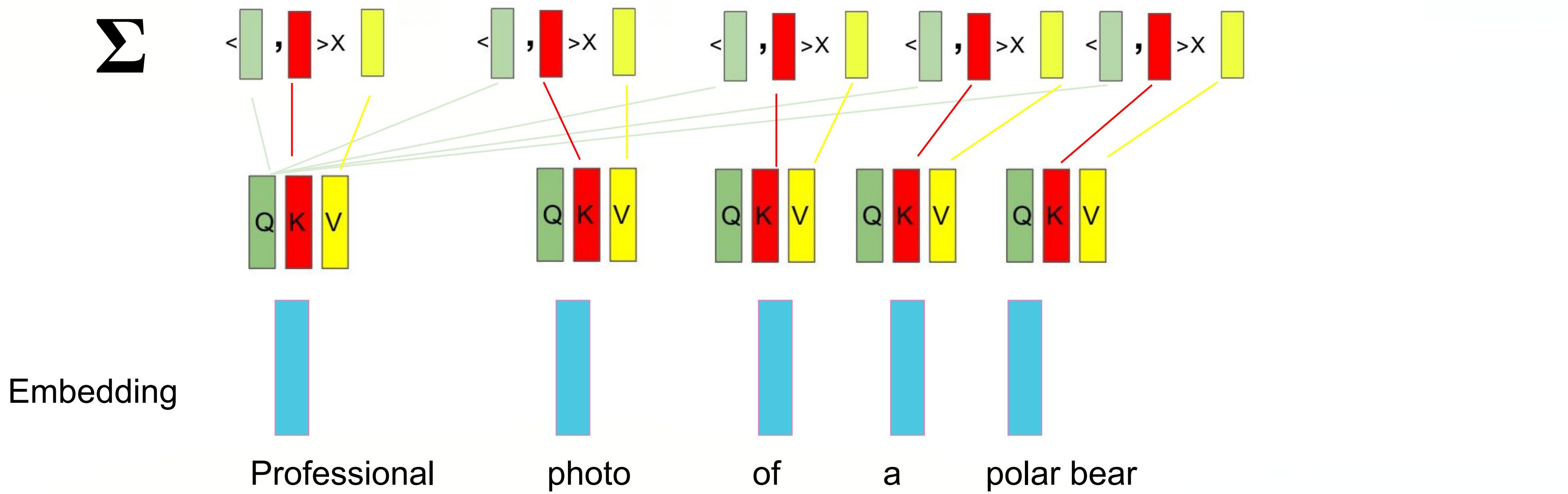
Glide text 2 image architecture



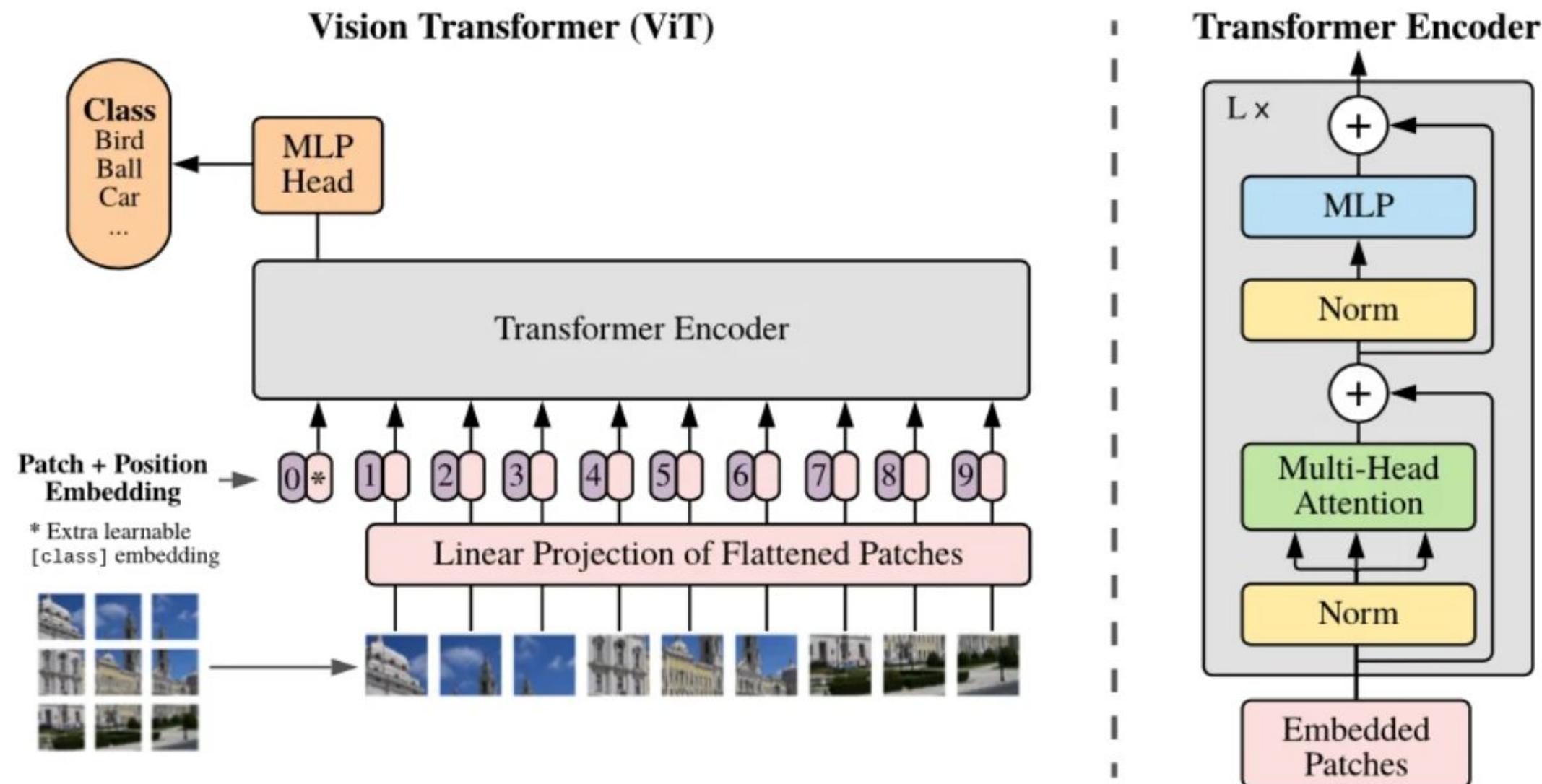
Transformers

Google, Dec 2017

Attention Layer

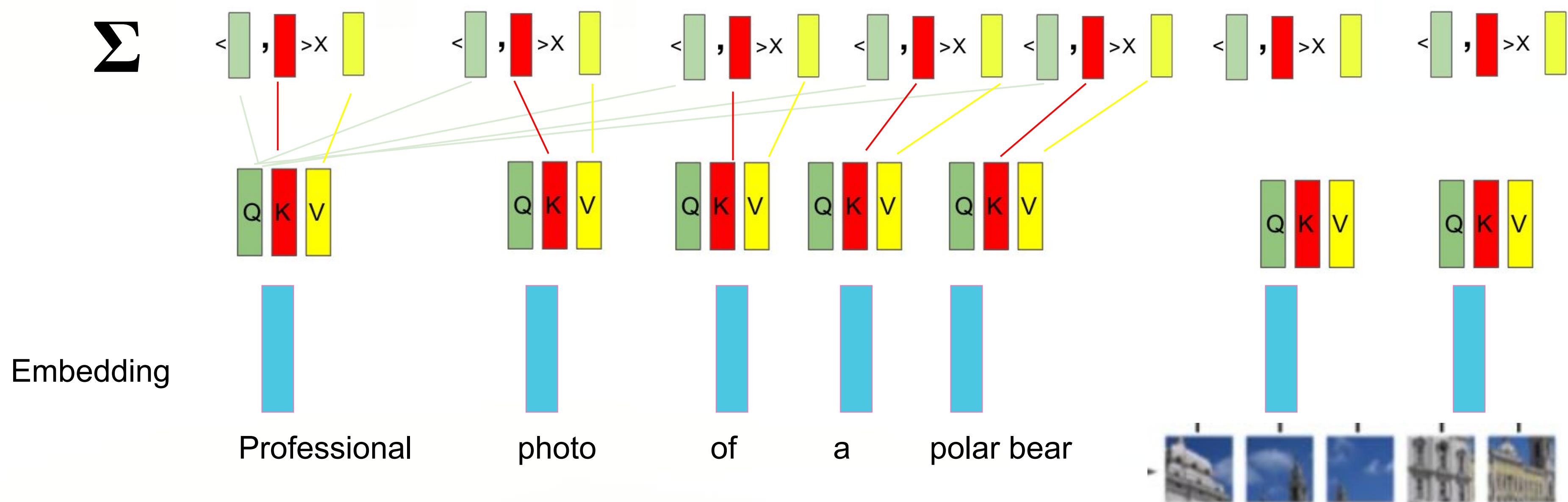


Visual Transformers

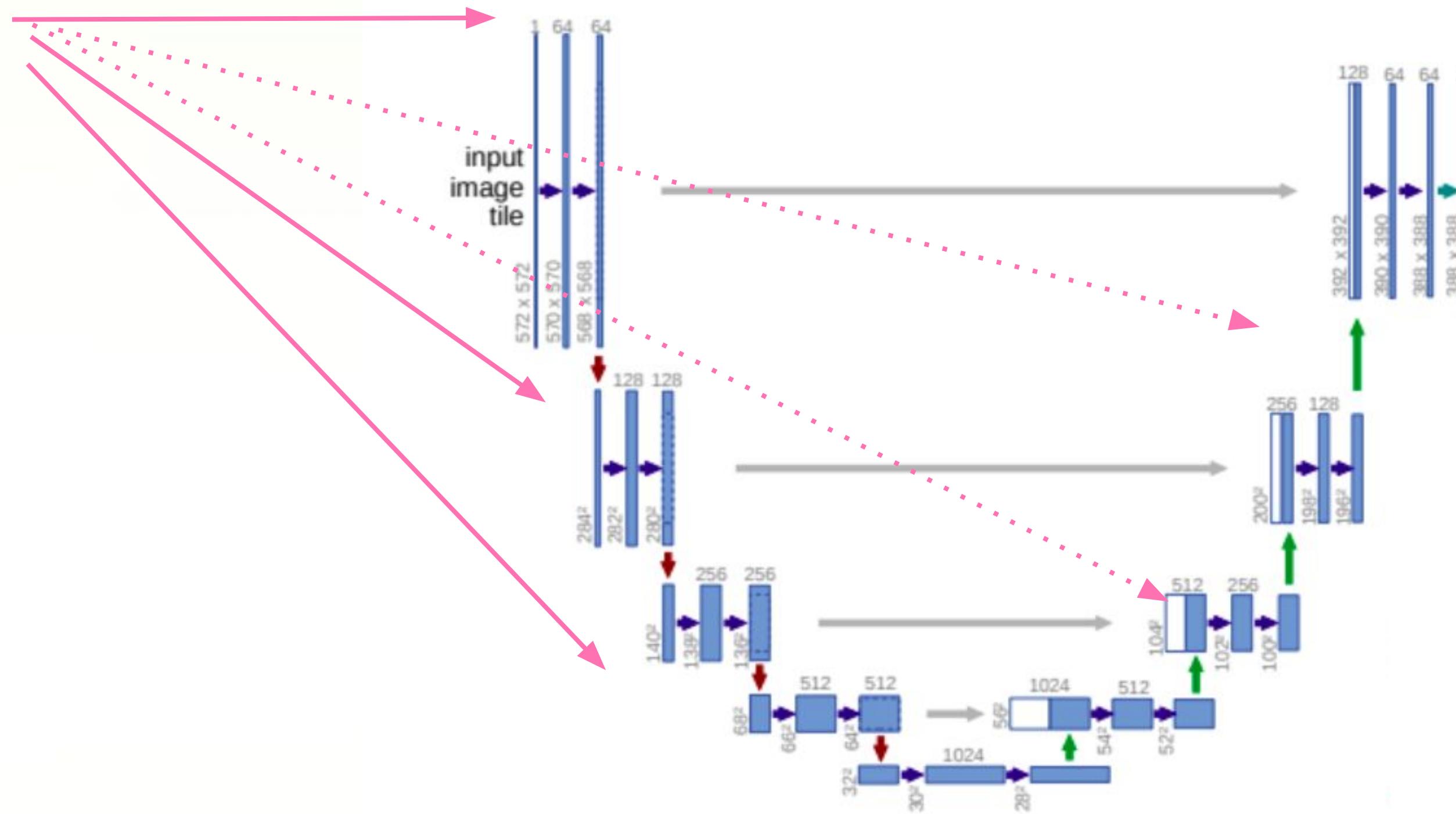


[An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale](#)

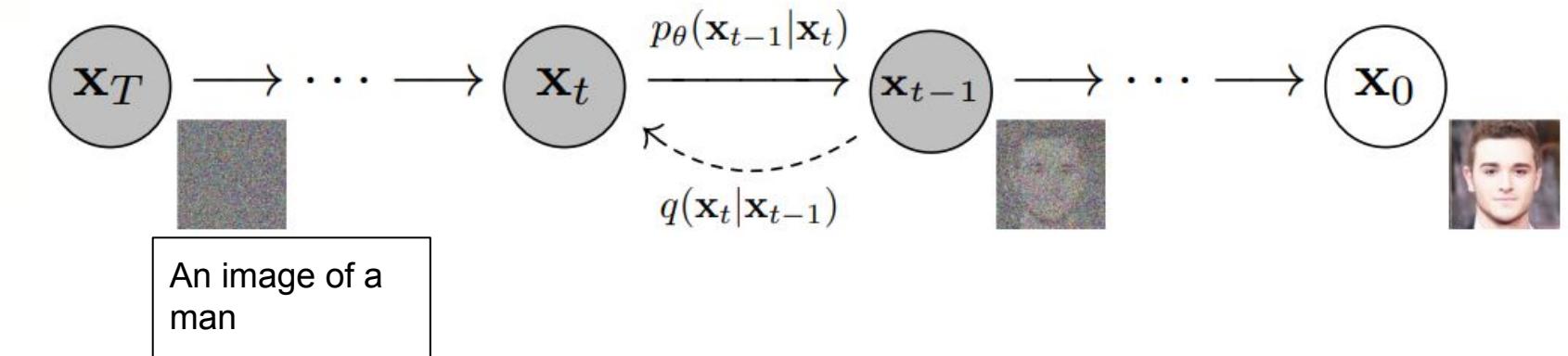
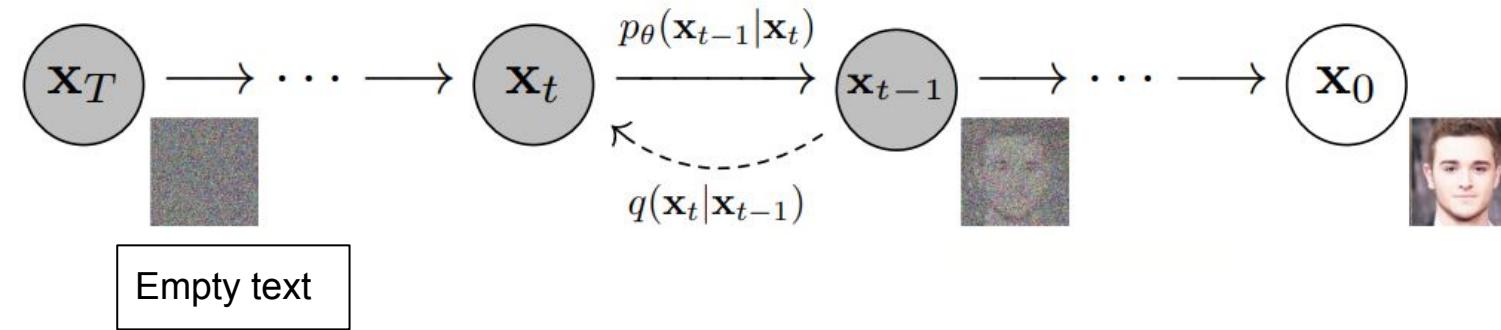
Cross Attention Layer



Adding textual input into the diffusion model



Guidance scale trick- extrapolation in the direction of the textual concept



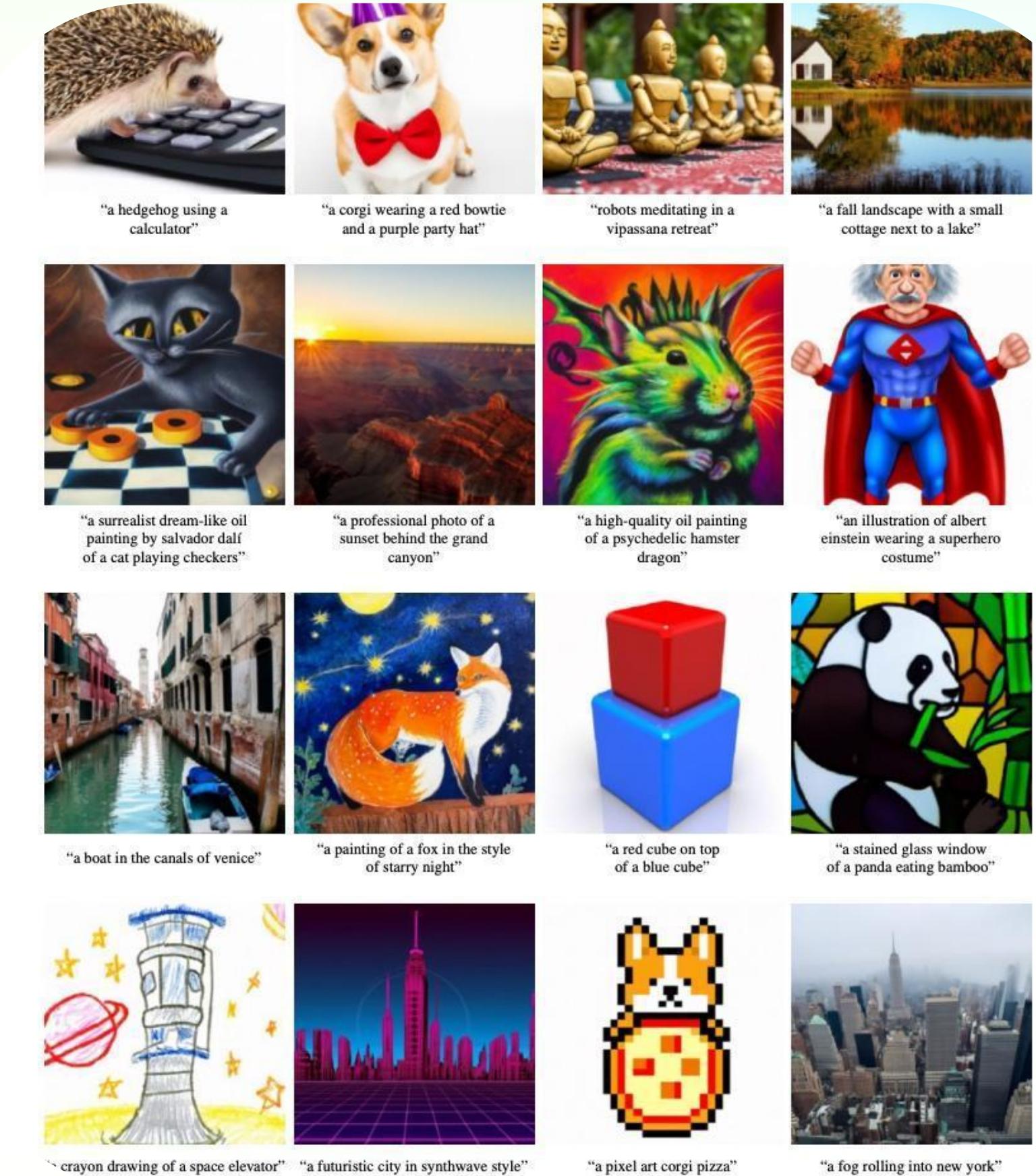
$$\mathbf{x} = \mathbf{x}_{empty} + guidancescale * (\mathbf{x}_{prompt} - \mathbf{x}_{empty})$$

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

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

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 (and deepFloyd) architecture

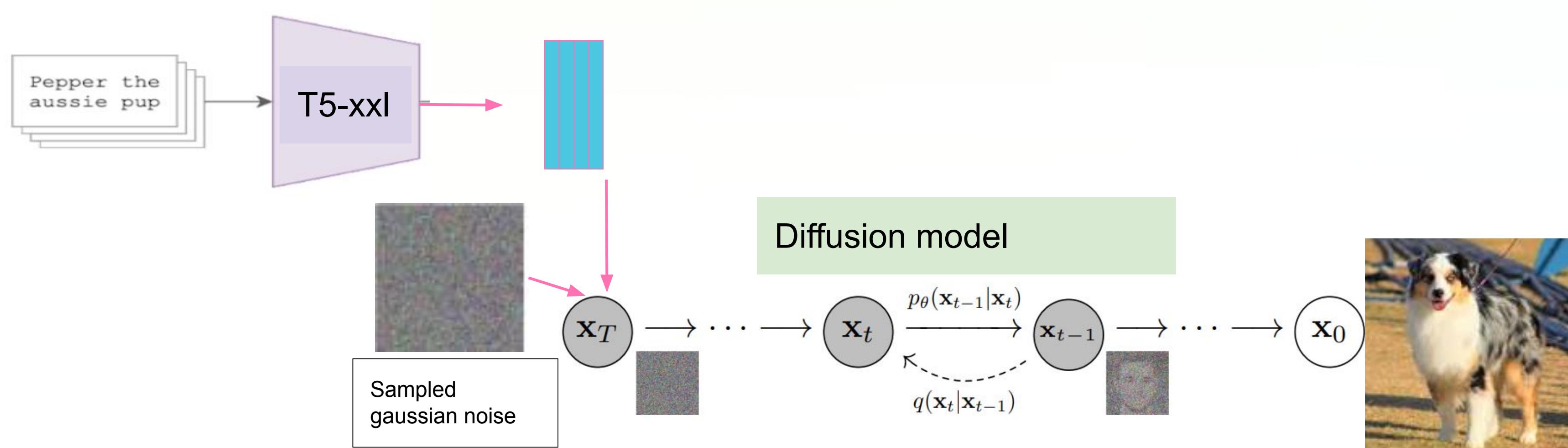
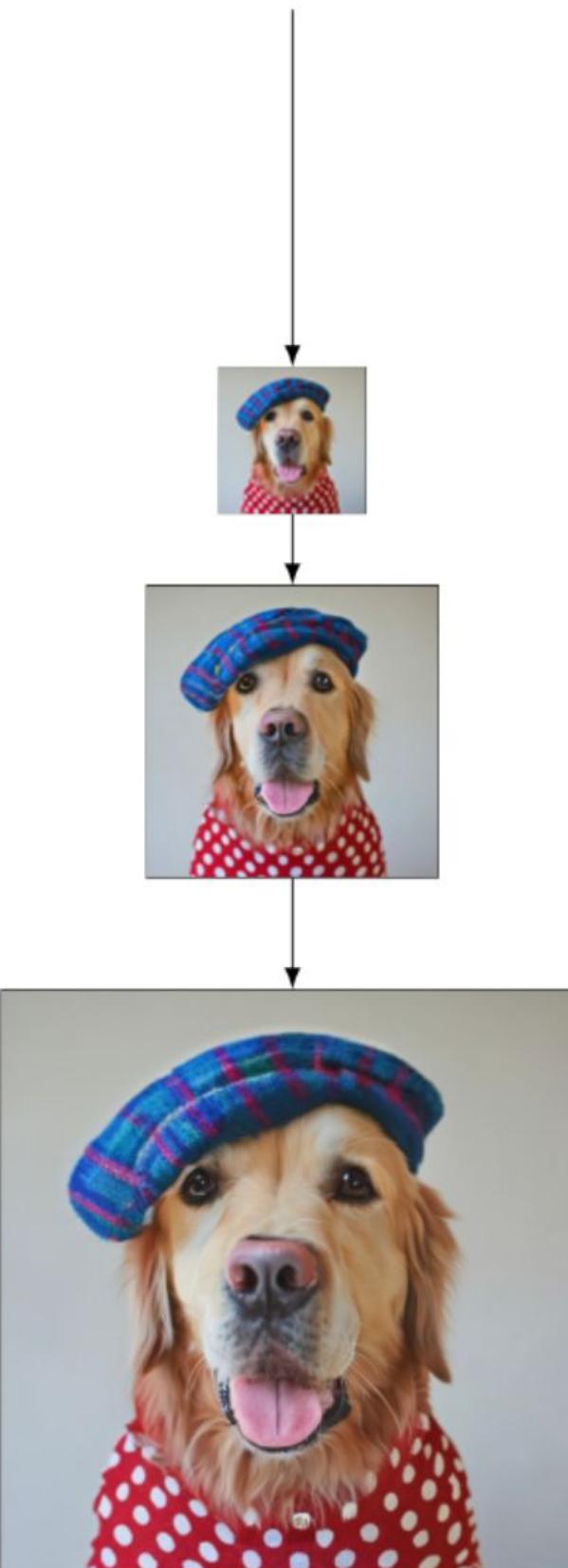


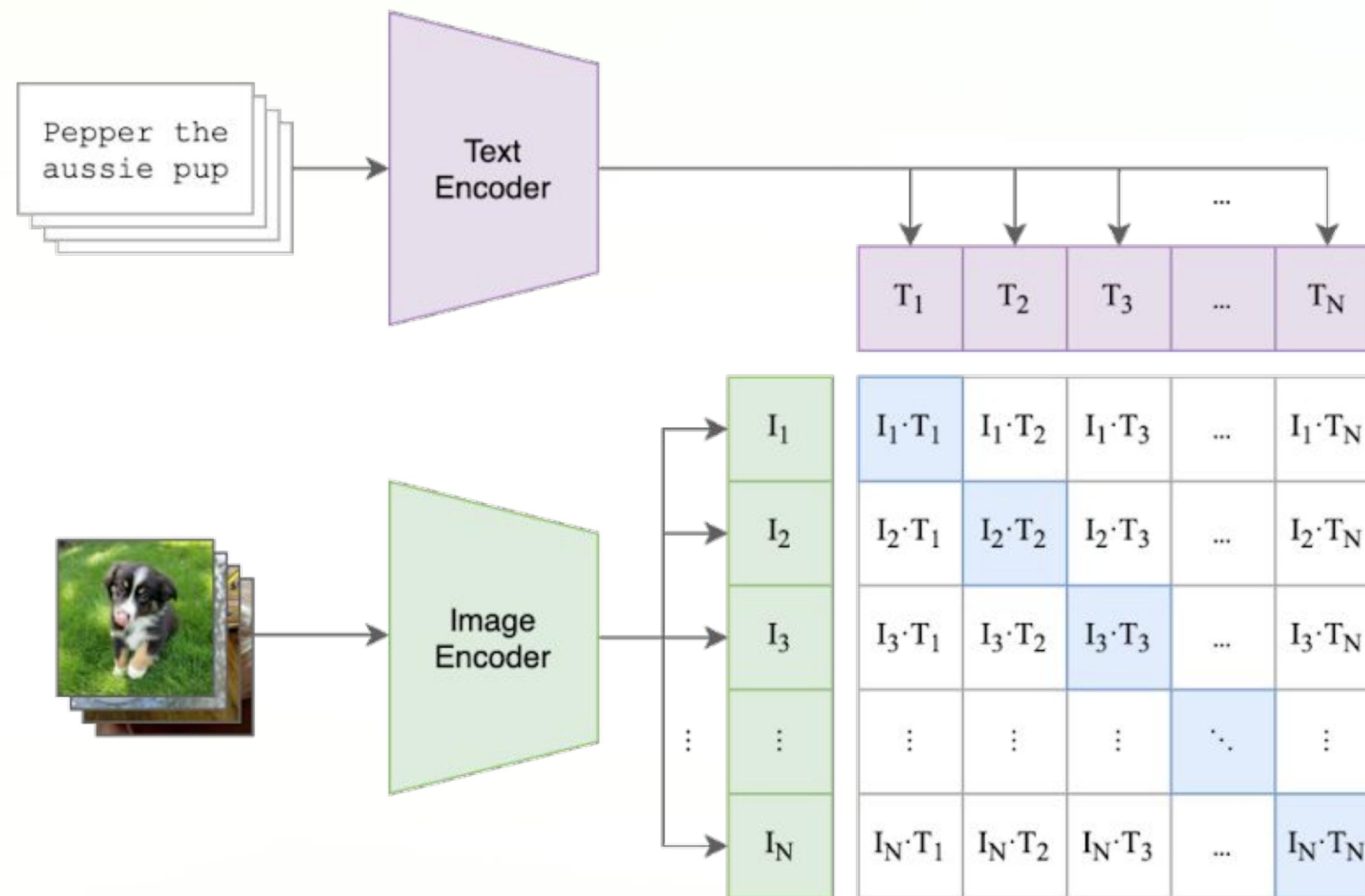
Imagen: Super Resolution models

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



CLIP- Text and Image encoder

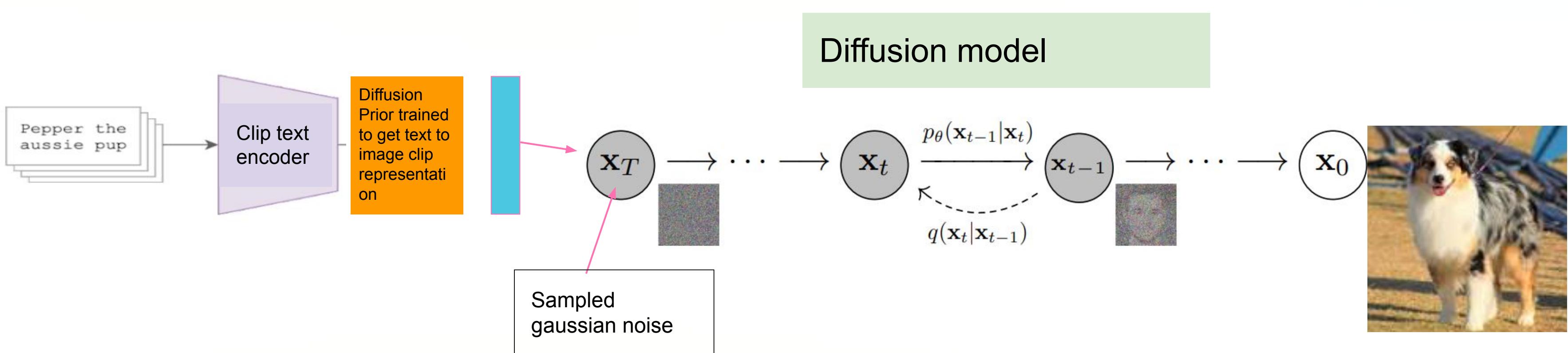
OpenAI, Jan 2021



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.



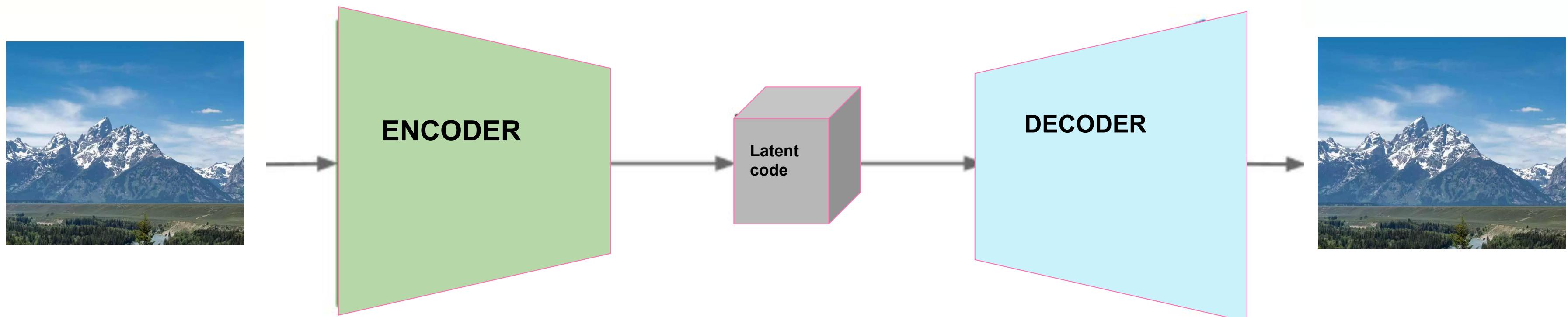
Stable Diffusion (Image to Image pipe)



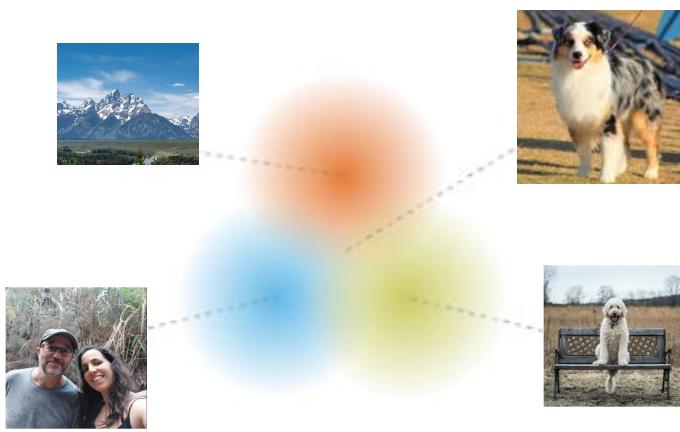
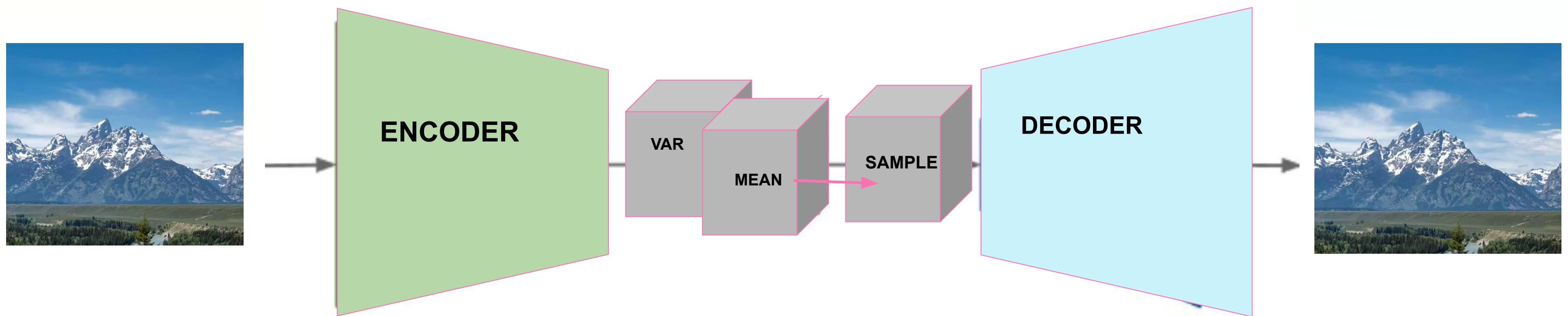
High-Resolution Image Synthesis with Latent Diffusion Models

Robin Rombach*, Andreas Blattmann*, Dominik Lorenz, Patrick Esser, Björn Ommer

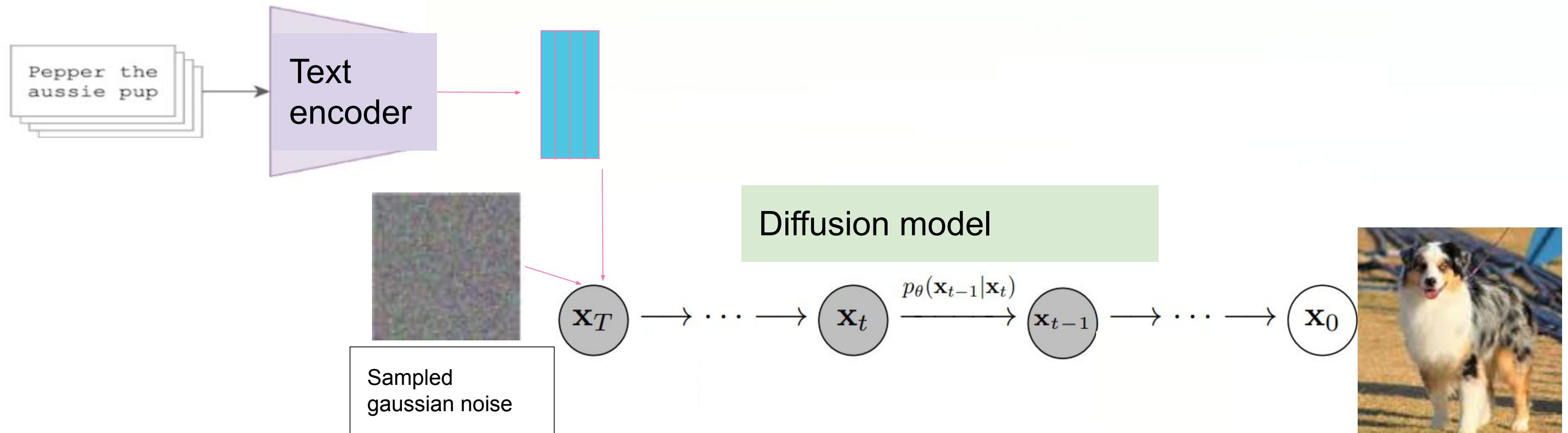
AutoEncoder



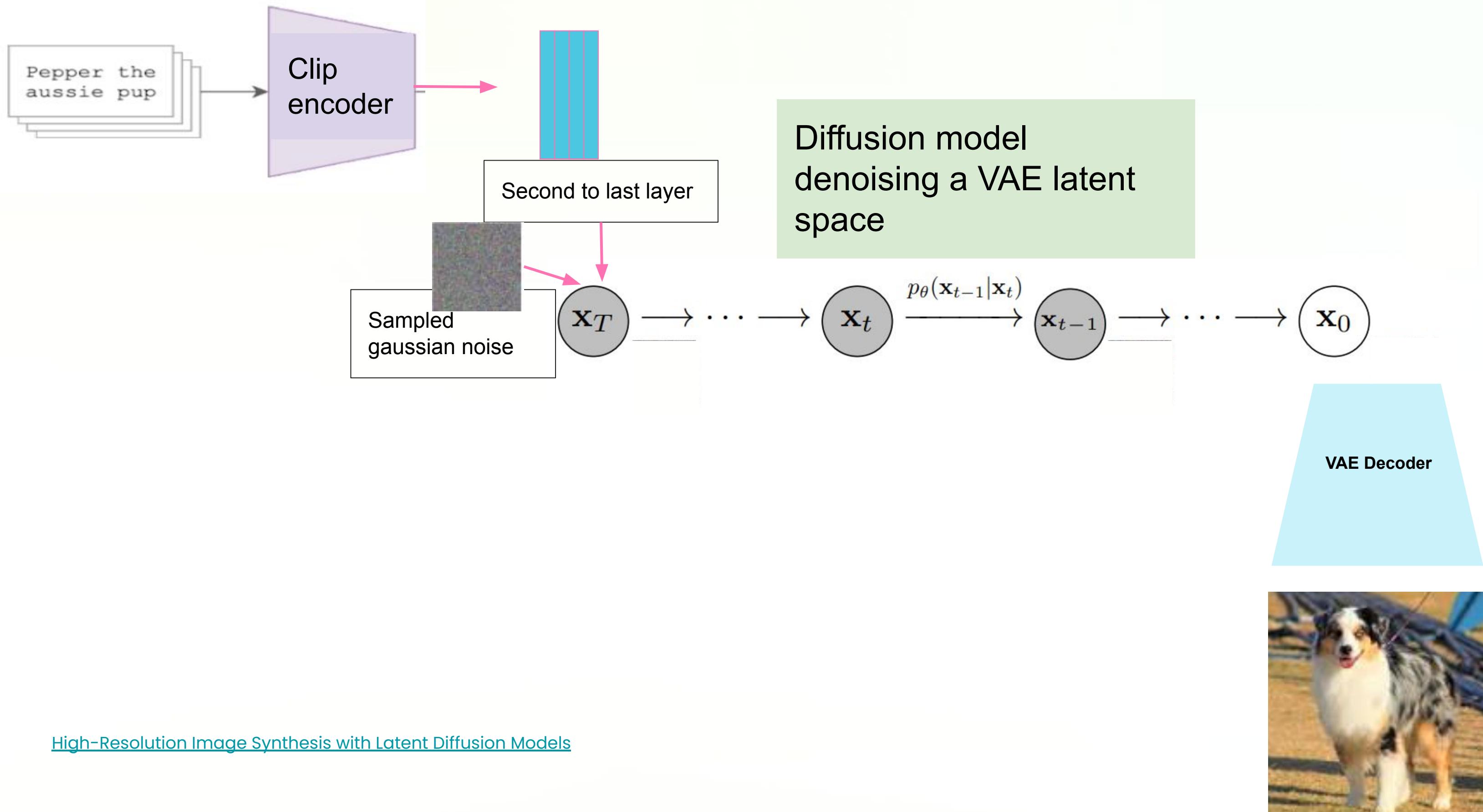
Variational AutoEncoder



Original Text to Image pipe

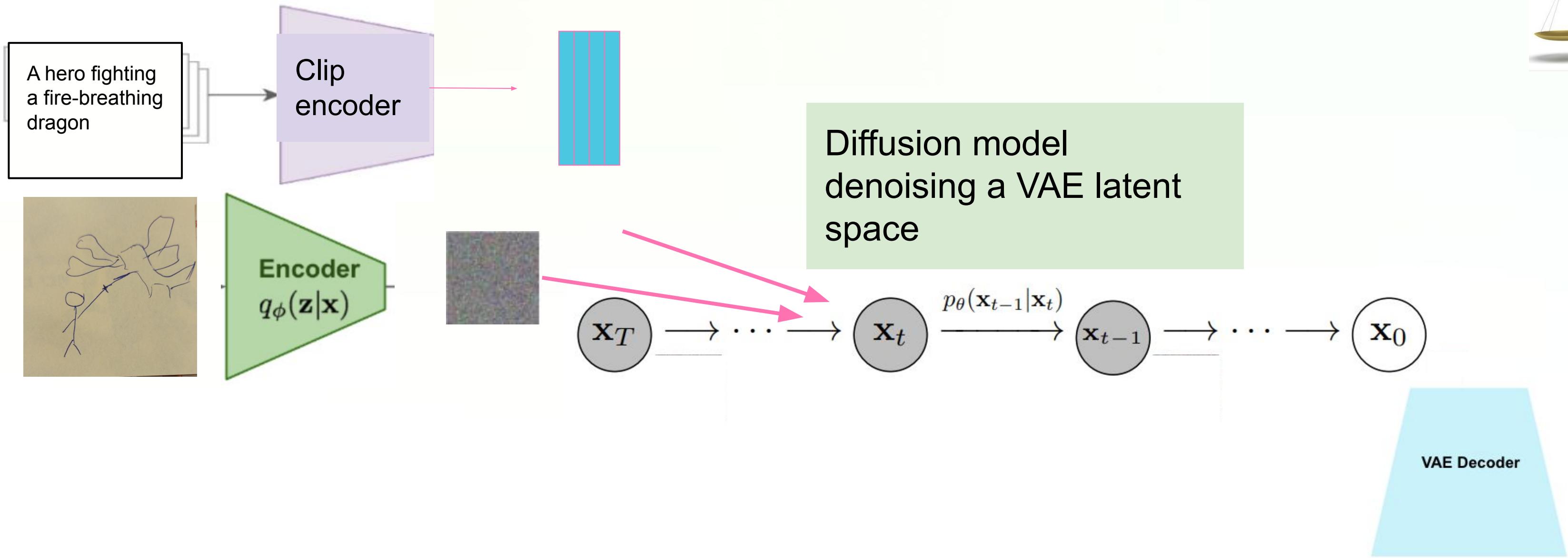


Stable Diffusion Text to Image pipe

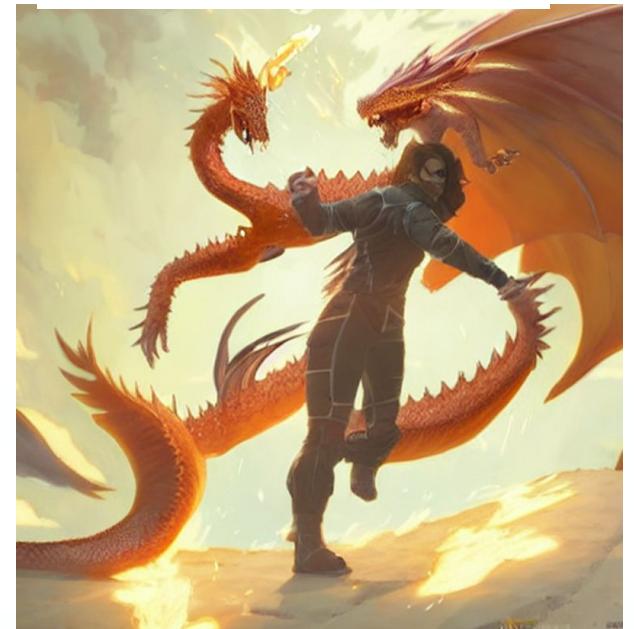


Editing pipelines

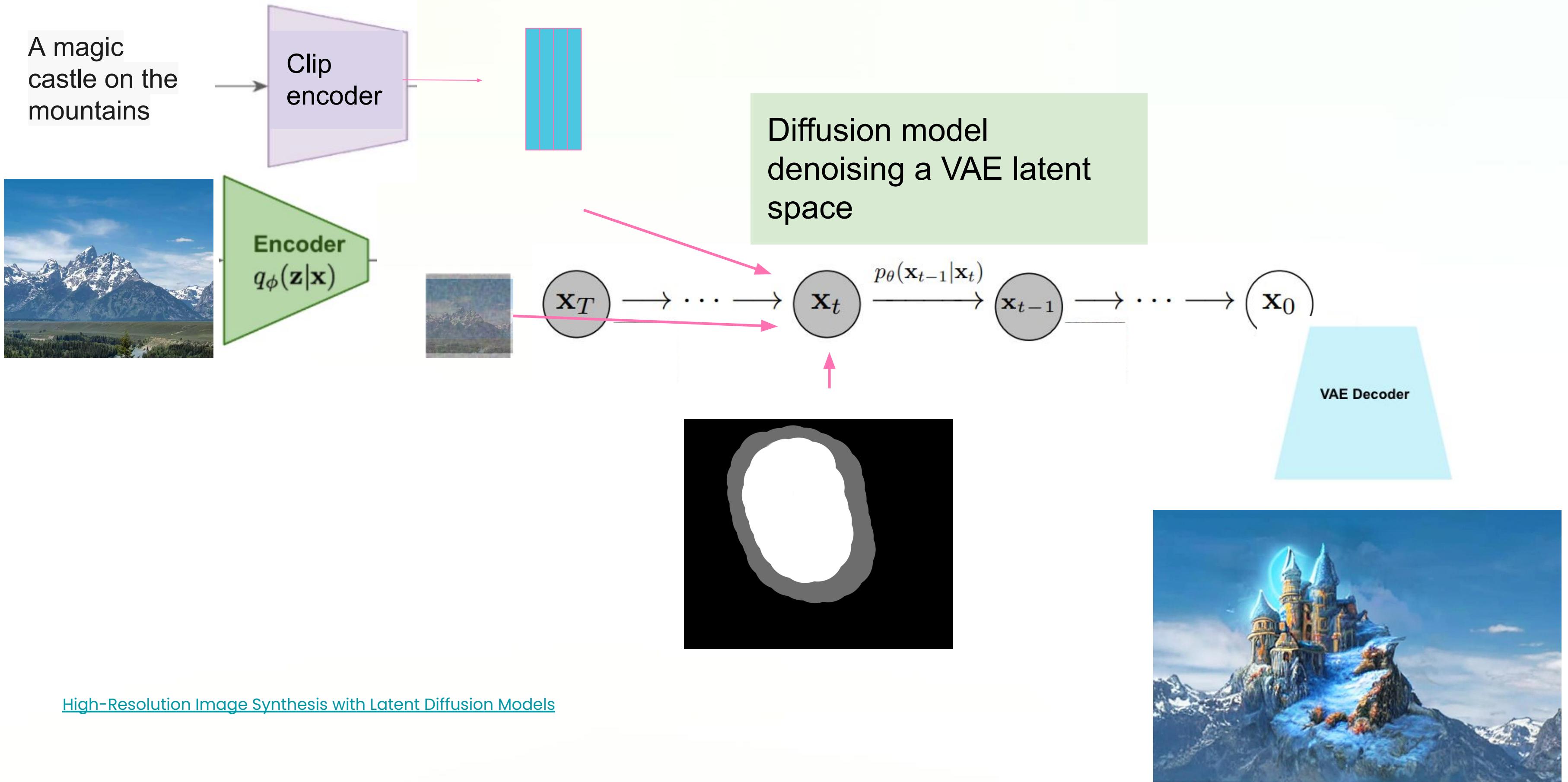
Stable Diffusion Image to Image pipe



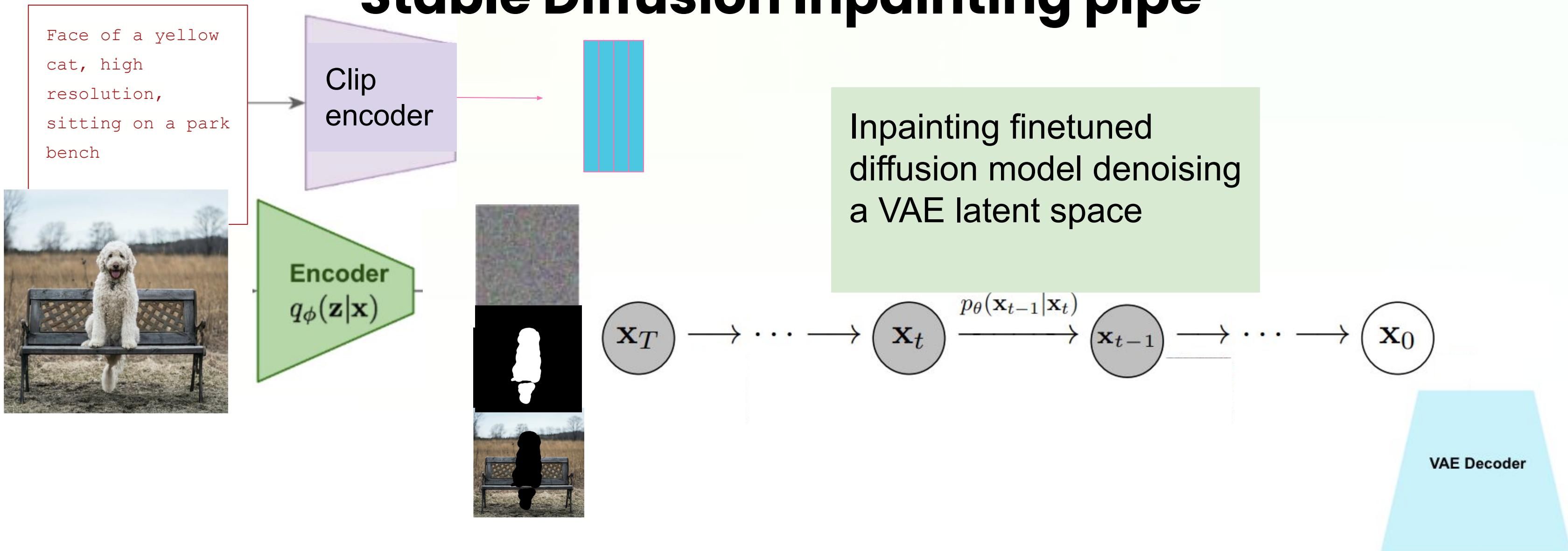
[High-Resolution Image Synthesis with Latent Diffusion Models](#)



Stable Diffusion Inpainting pipe (legacy)

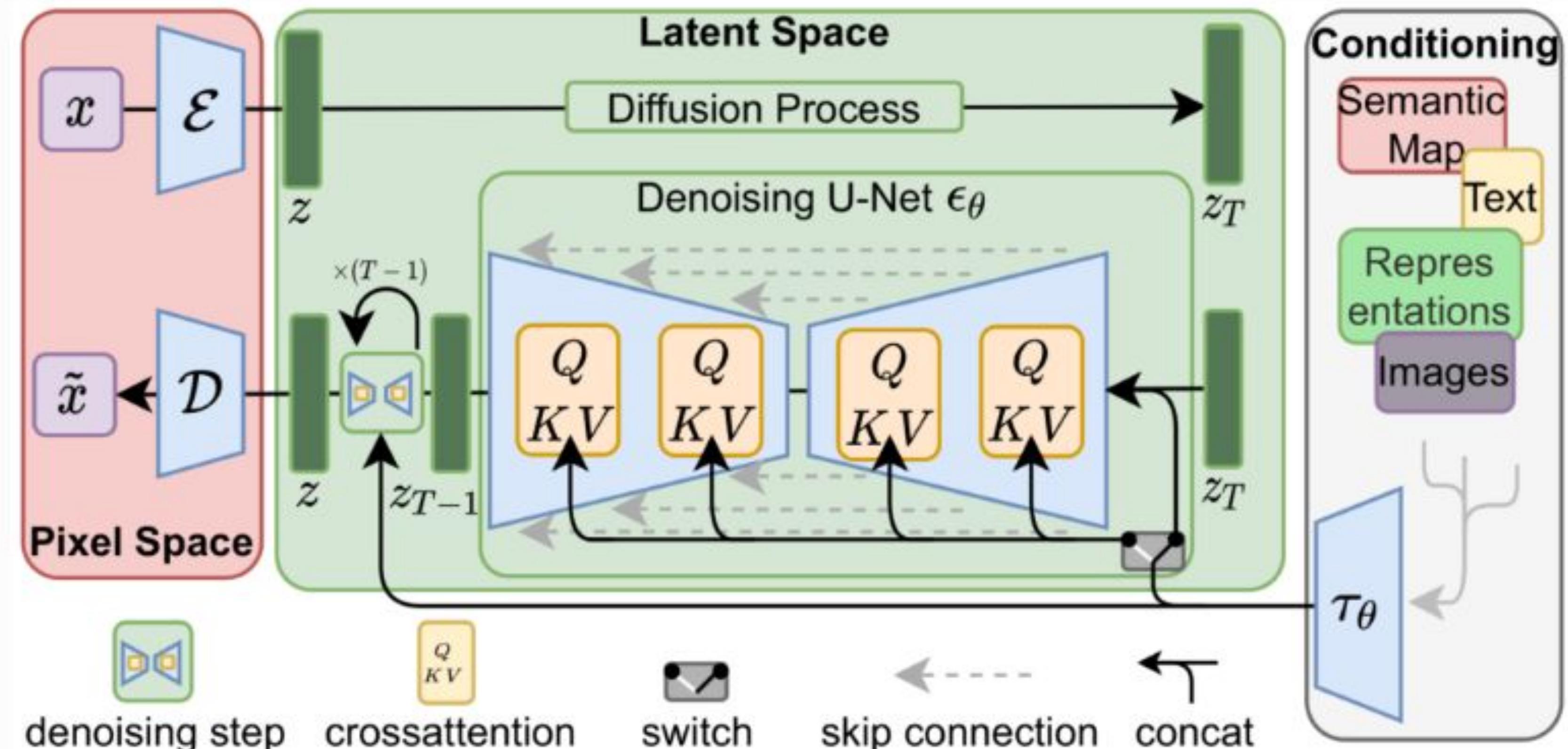


Stable Diffusion Inpainting pipe



Inpainting





Useful gits to follow

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

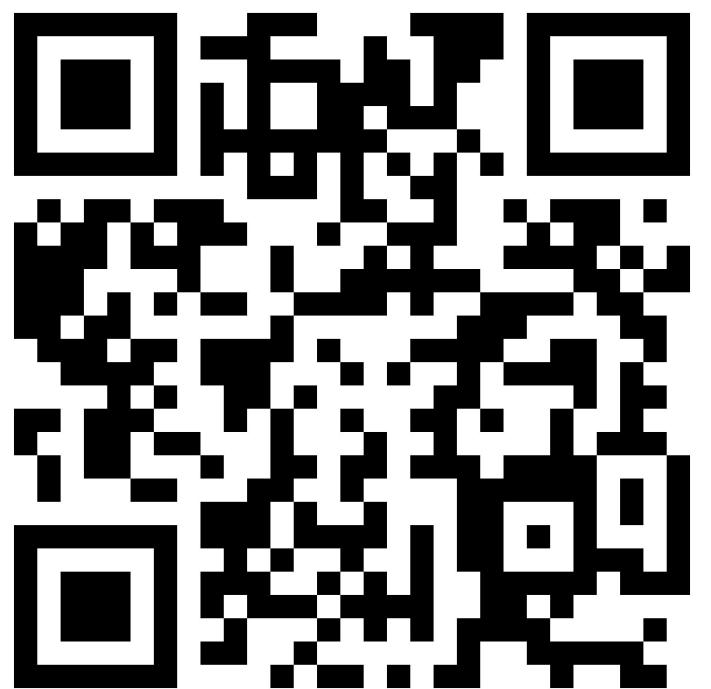
<https://github.com/AUTOMATIC1111/stable-diffusion-webui>

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 session (40 minutes)



Stable_diffusion_editing_pipes.ipynb



Takeaway

1. You now know what is it diffusion models
2. We covered a lot of architectures, and technical ideas. Combining those existing ideas leads to great progress
3. It should be pretty easy for you to start playing with diffusers git



Thank you!

Instruct pix2pix

Tim Brooks* Aleksander Holynski* Alexei A.
Efros University of California, Berkeley
Jan, 2023

"Swap sunflowers with roses"

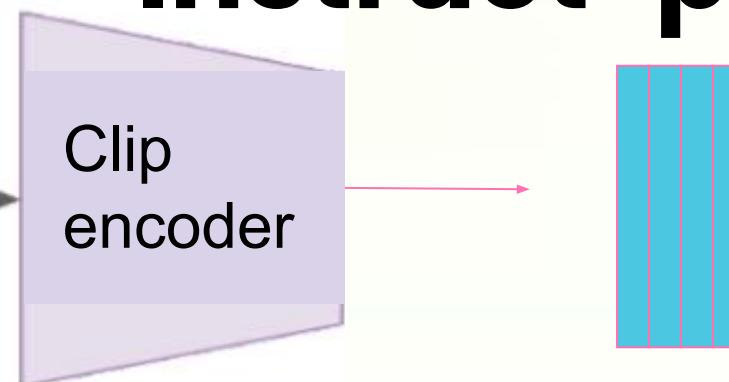


"Make his jacket out of leather"

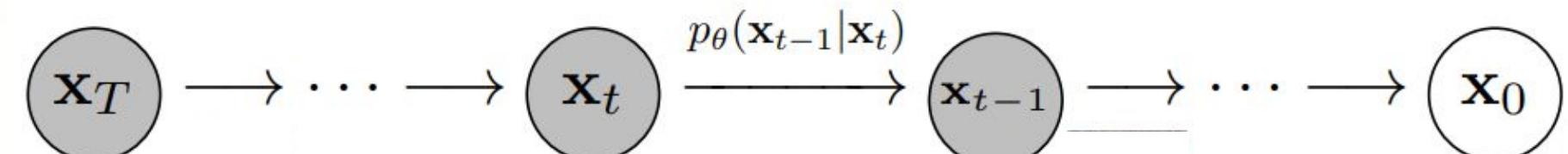
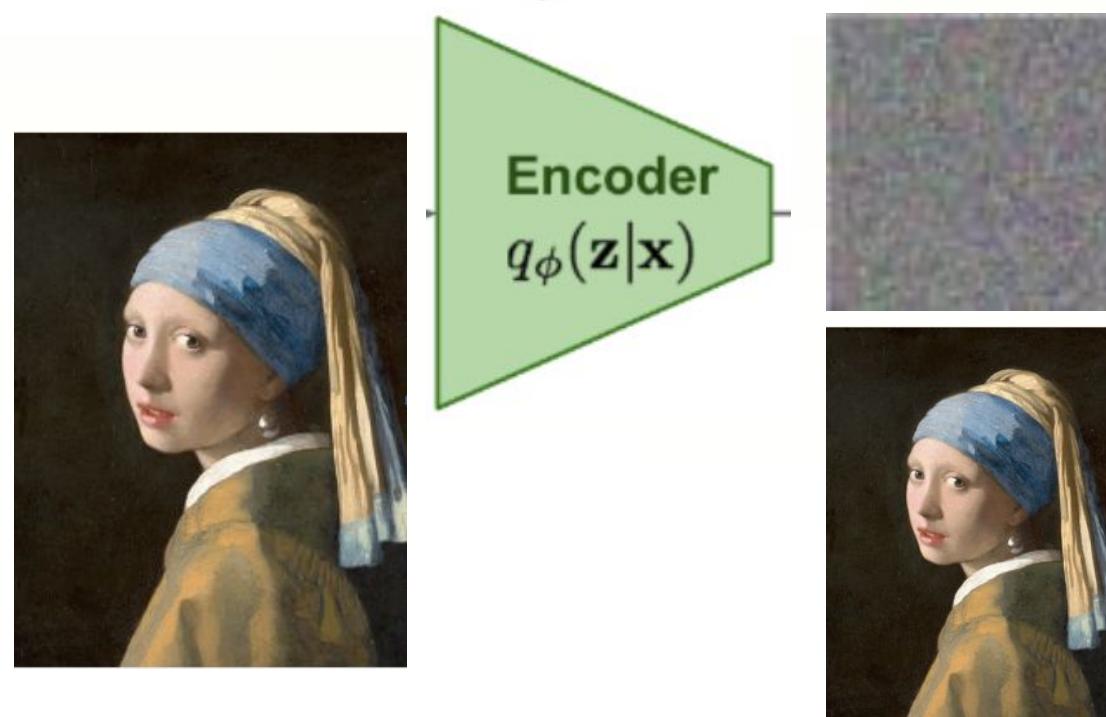


Instruct-pix2pix pipe

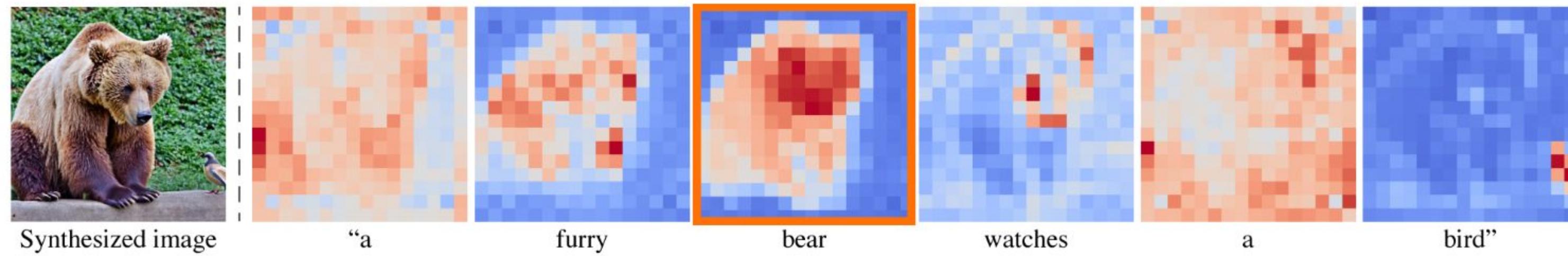
Turn her into a
snake lady



Instruct pix2pix model



Cross attention visualization



Training Data Generation

(a) Generate text edits:

Input Caption: "photograph of a girl riding a horse" →
700 samples

GPT-3

Instruction: "have her ride a dragon"

Edited Caption: "photograph of a girl riding a dragon"

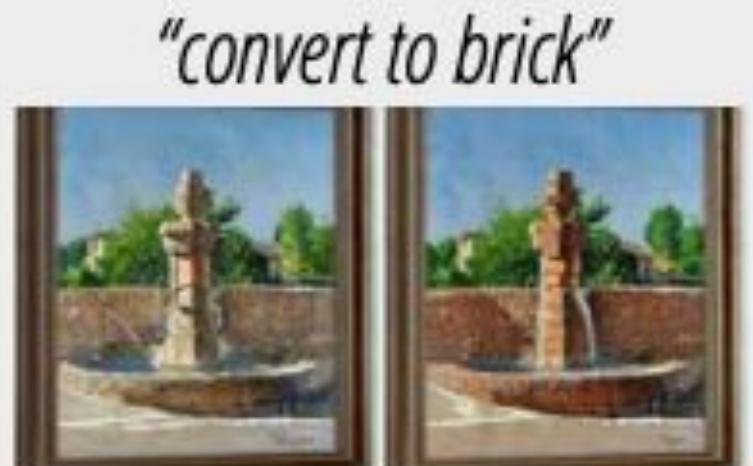
(b) Generate paired images:

Input Caption: "photograph of a girl riding a horse"
Edited Caption: "photograph of a girl riding a dragon"

Stable Diffusion
+ Prompt2Prompt



(c) Generated training examples:



...