

GenAI in 2024

Overview of recent advances in generative AI for vision

Presented by Kanchana Ranasinghe | June 2024

Acknowledgements

- Content borrowed from:
<https://cvpr2023-tutorial-diffusion-models.github.io>
<https://cvpr2022-tutorial-diffusion-models.github.io>
- Lazebnik Slides:
https://slazebni.cs.illinois.edu/spring17/lec11_gan.pdf

Overview

Deep Generative Learning

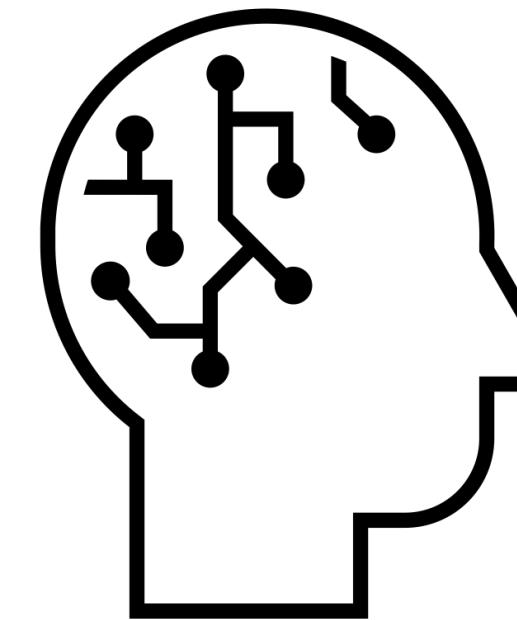
Learning to generate data



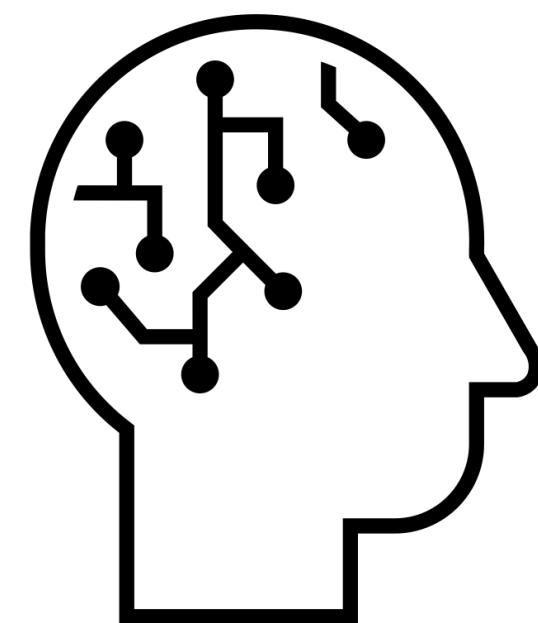
Samples from a Data Distribution

Train

A thick green arrow pointing from the cloud of cat images to the neural network icon.



Neural Network



Sample

A thick green arrow pointing from the neural network icon to a generated cat image.



Generative Learning Applications

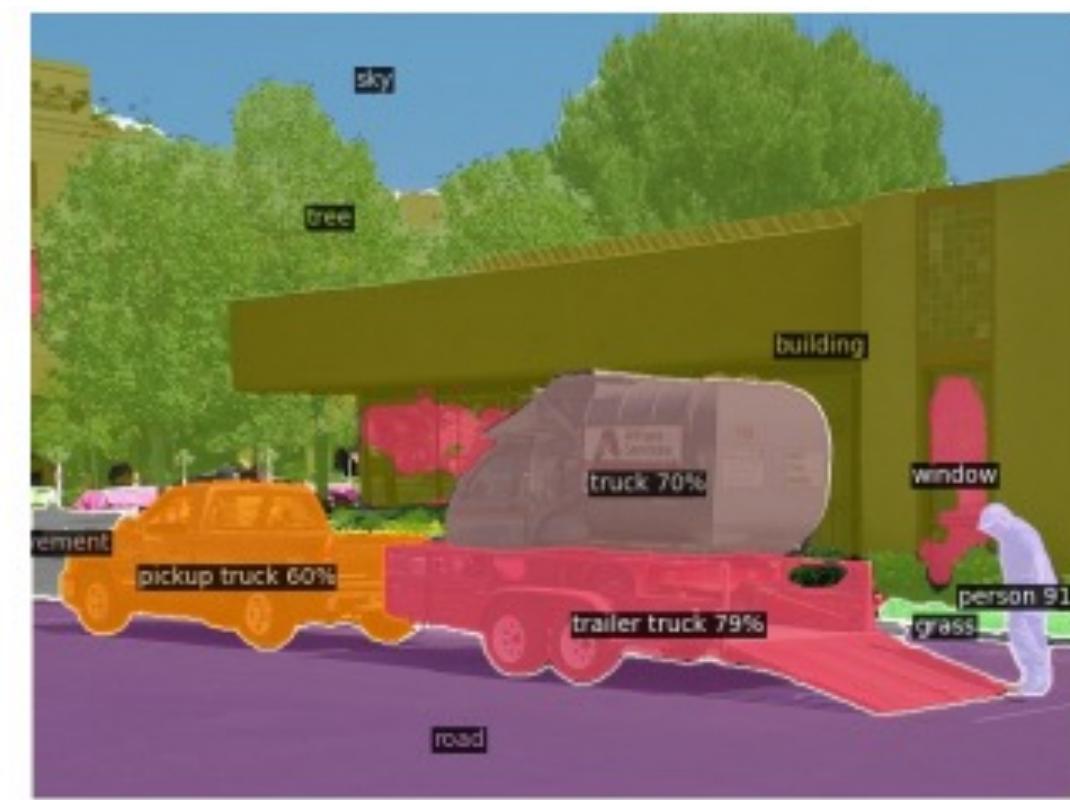
Art & Design



Content Generation



Representation Learning



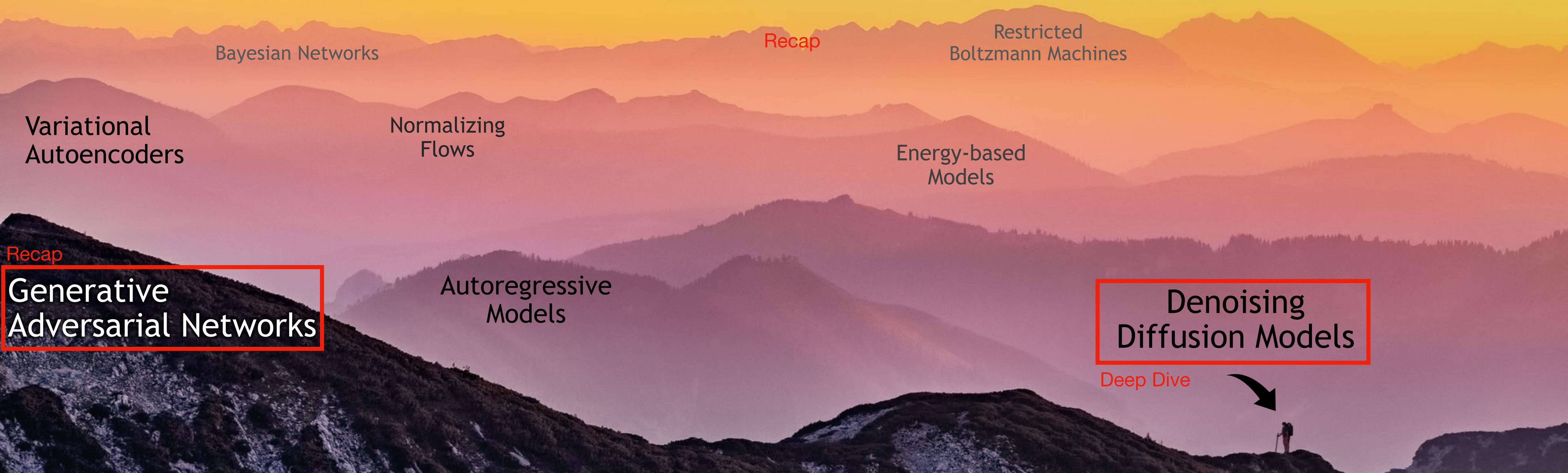
Entertainment



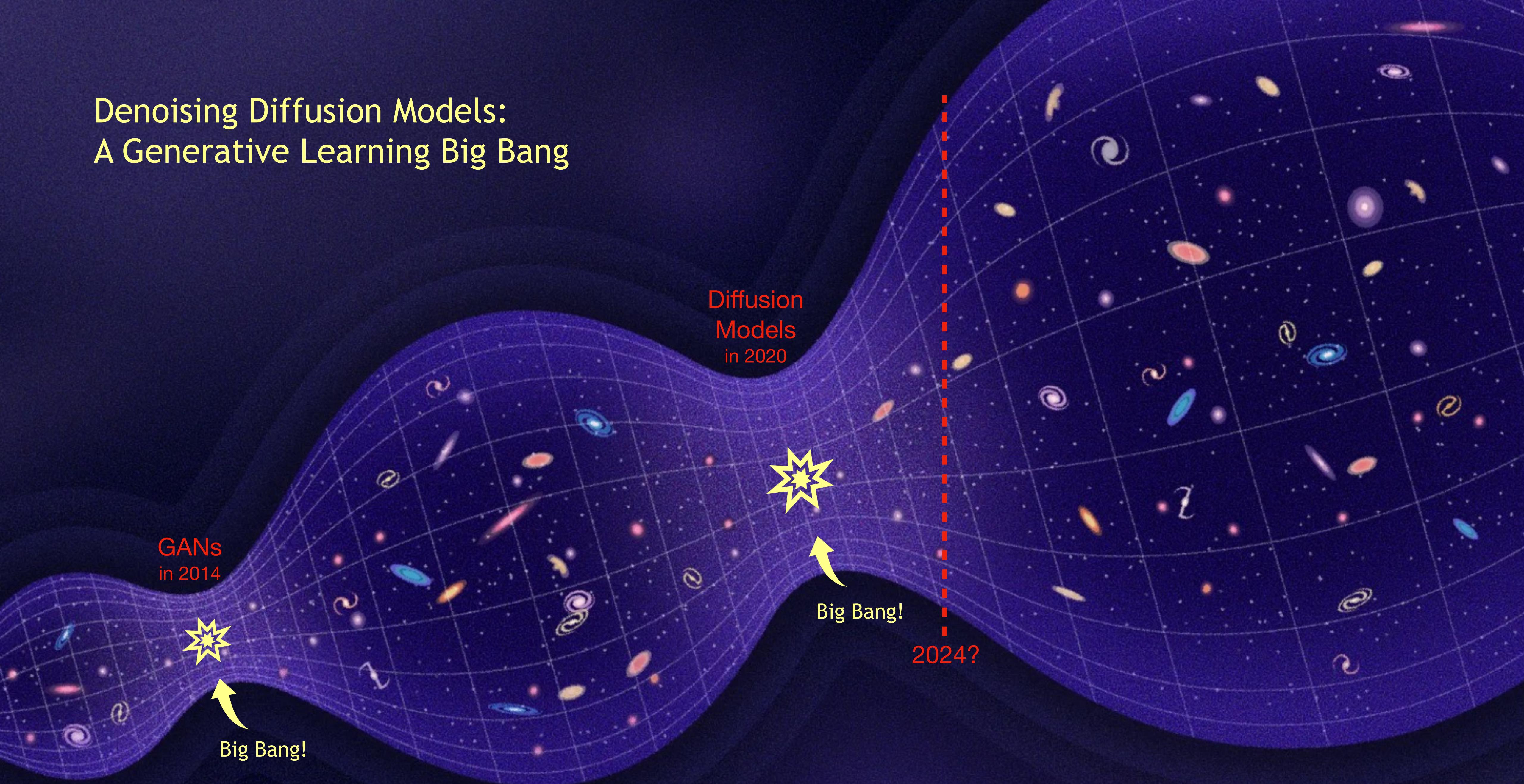
The Landscape of Deep Generative Learning



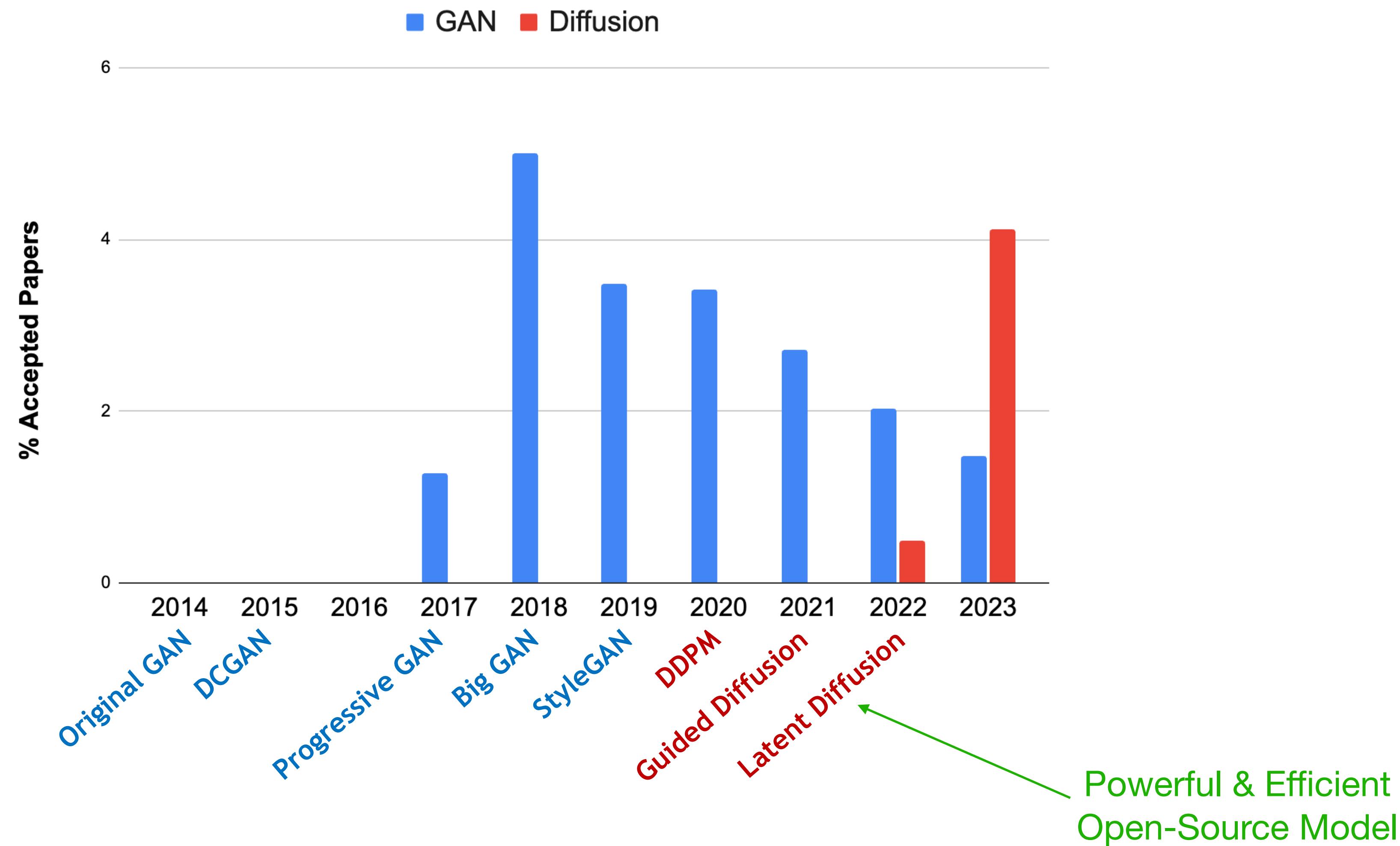
The Landscape of Deep Generative Learning



Denoising Diffusion Models: A Generative Learning Big Bang

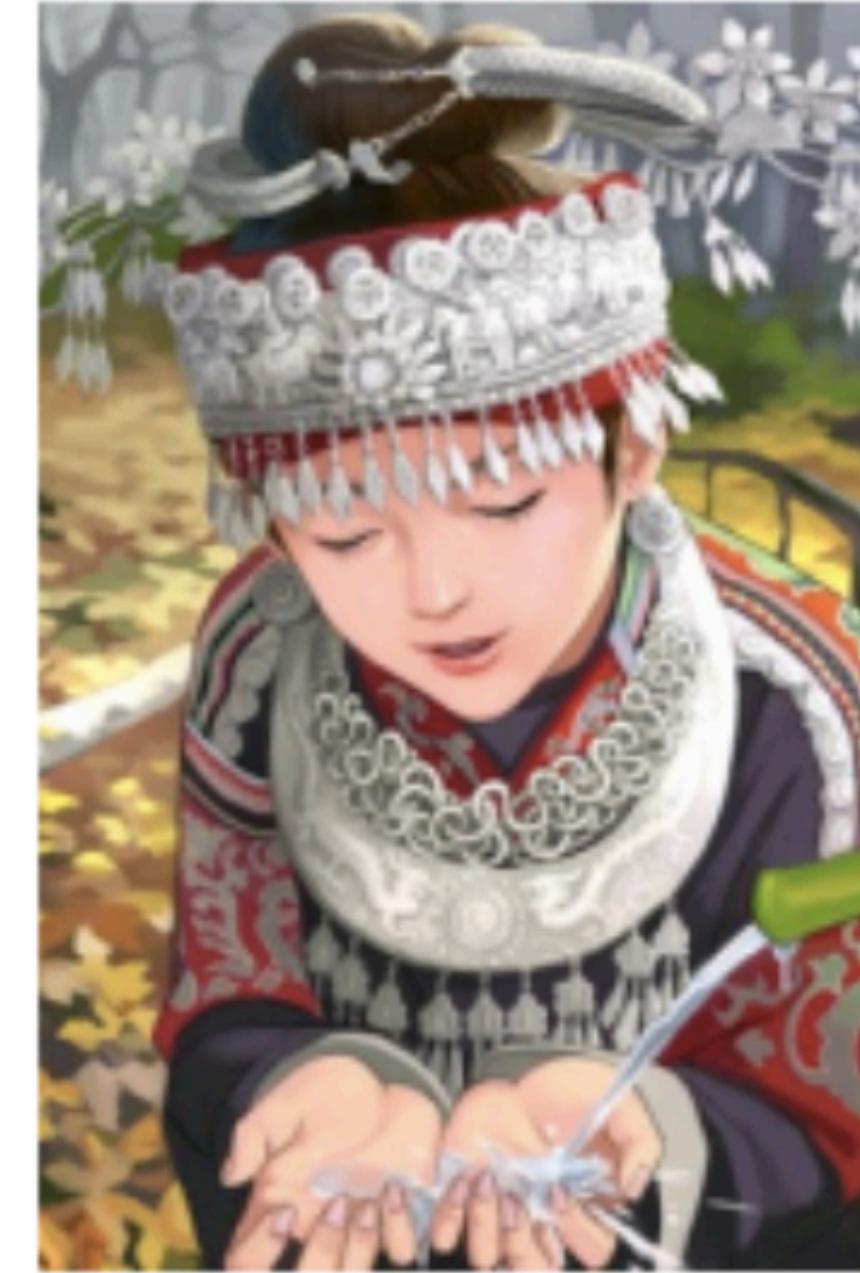


We May Not Know Cosmology, But We Know CVPR



Recap: GANs

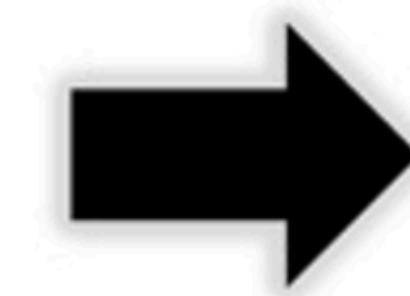
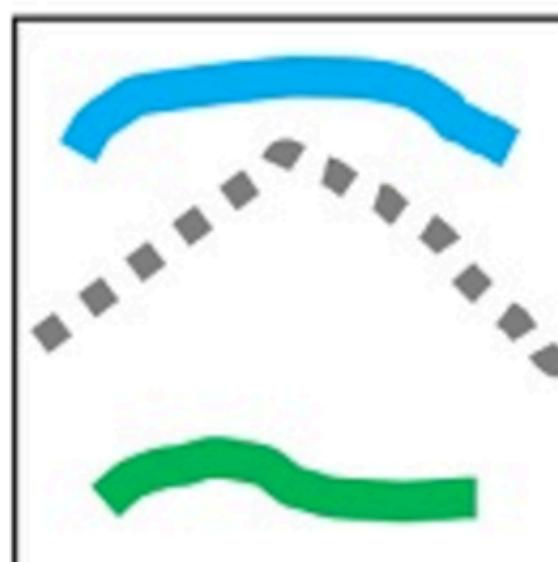
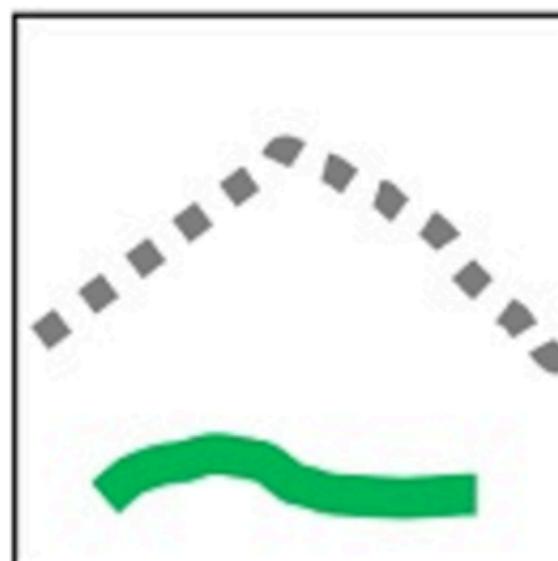
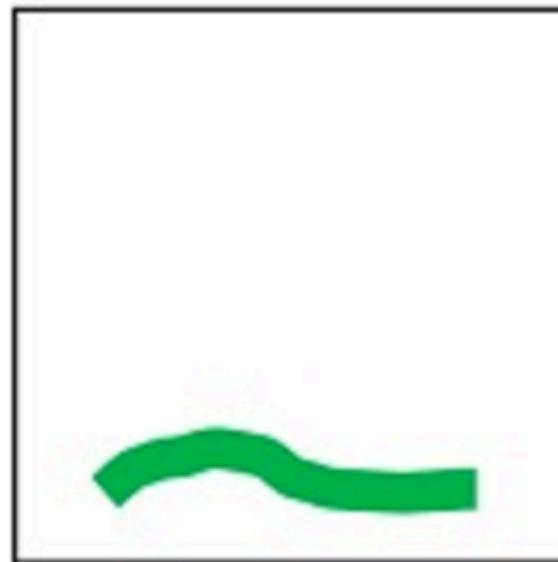
Which one is Computer generated?



Ledig, Christian, et al. "Photo-realistic single image super-resolution using a generative adversarial network." *arXiv preprint arXiv:1609.04802* (2016).

Magic of GANs...

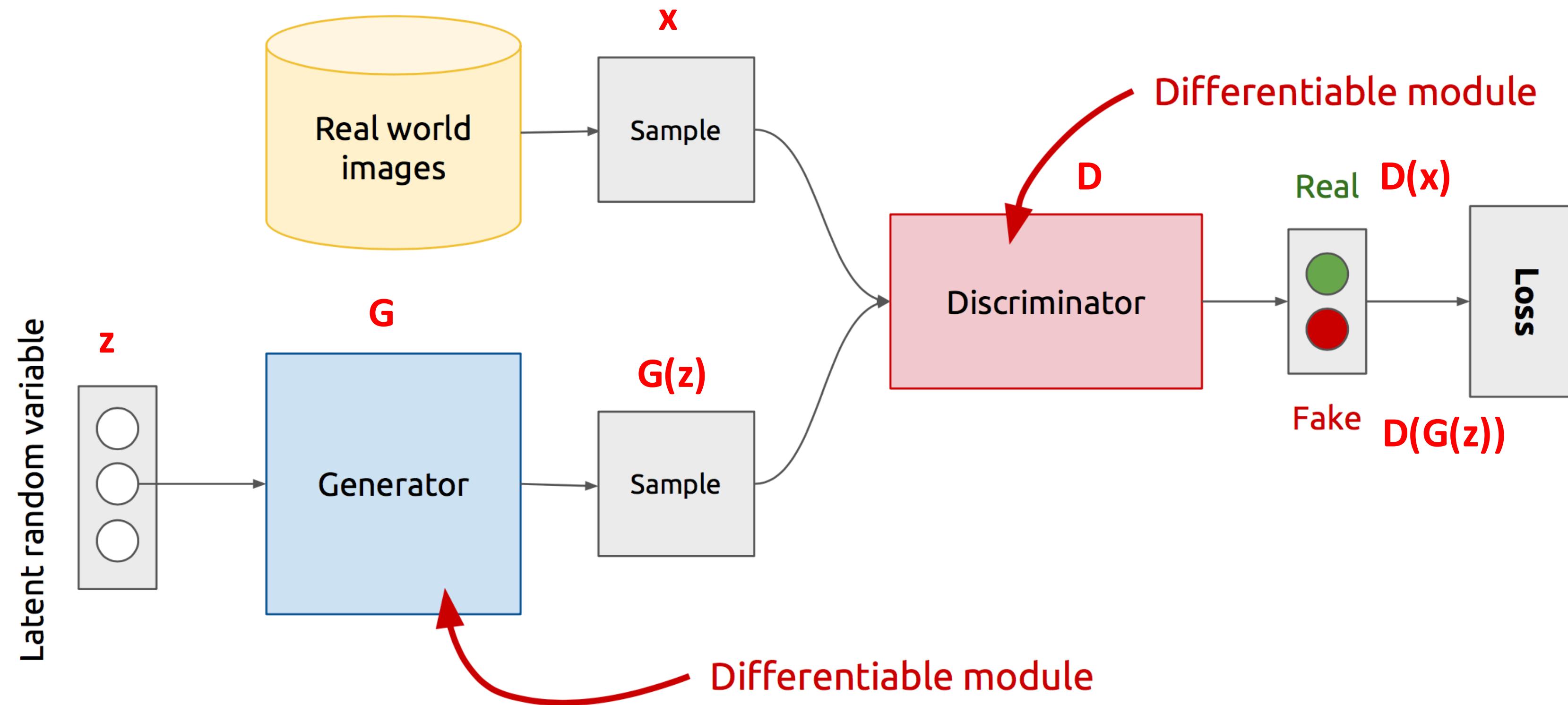
User edits



Generated images

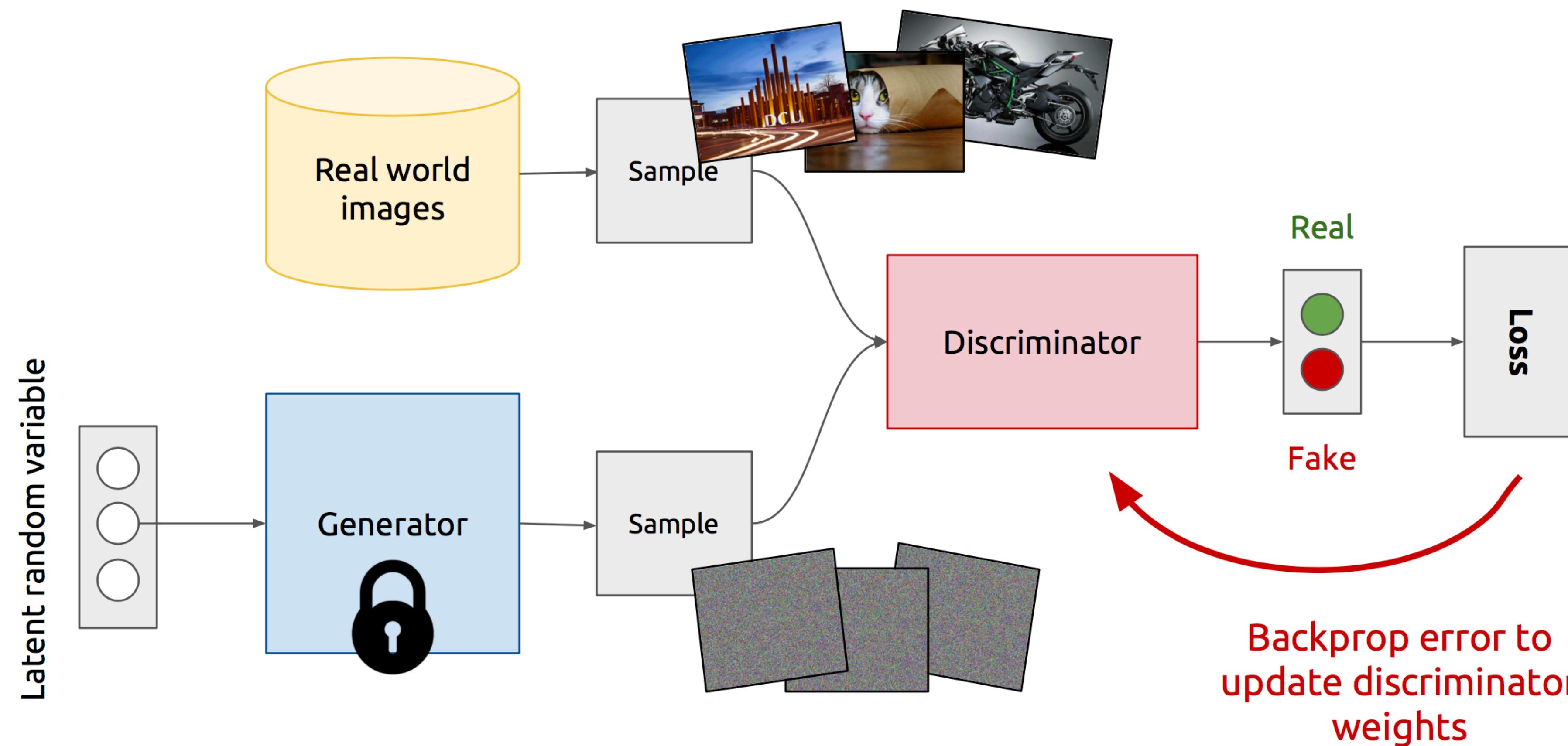


GAN's Architecture

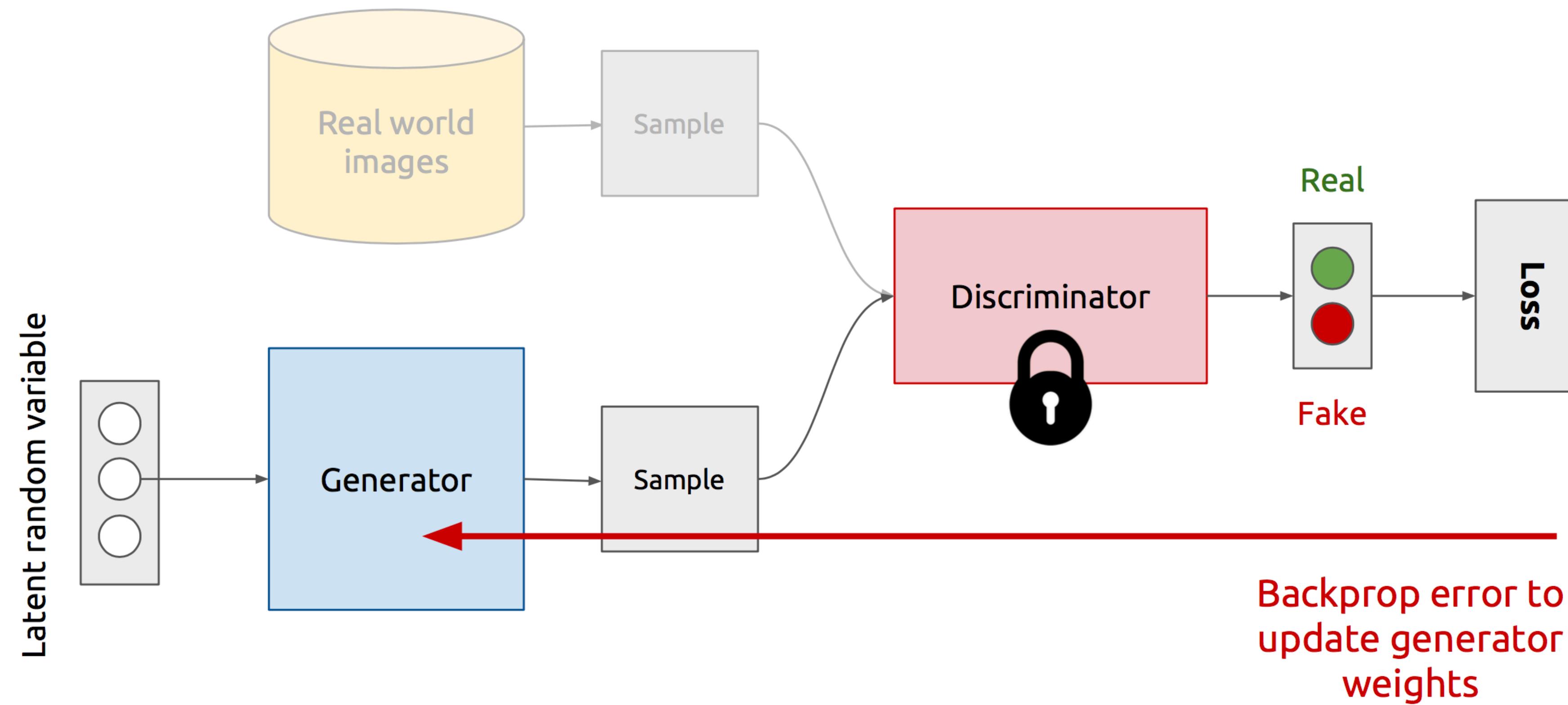


- Z is some random noise (Gaussian/Uniform).
- Z can be thought as the latent representation of the image.

Training Discriminator

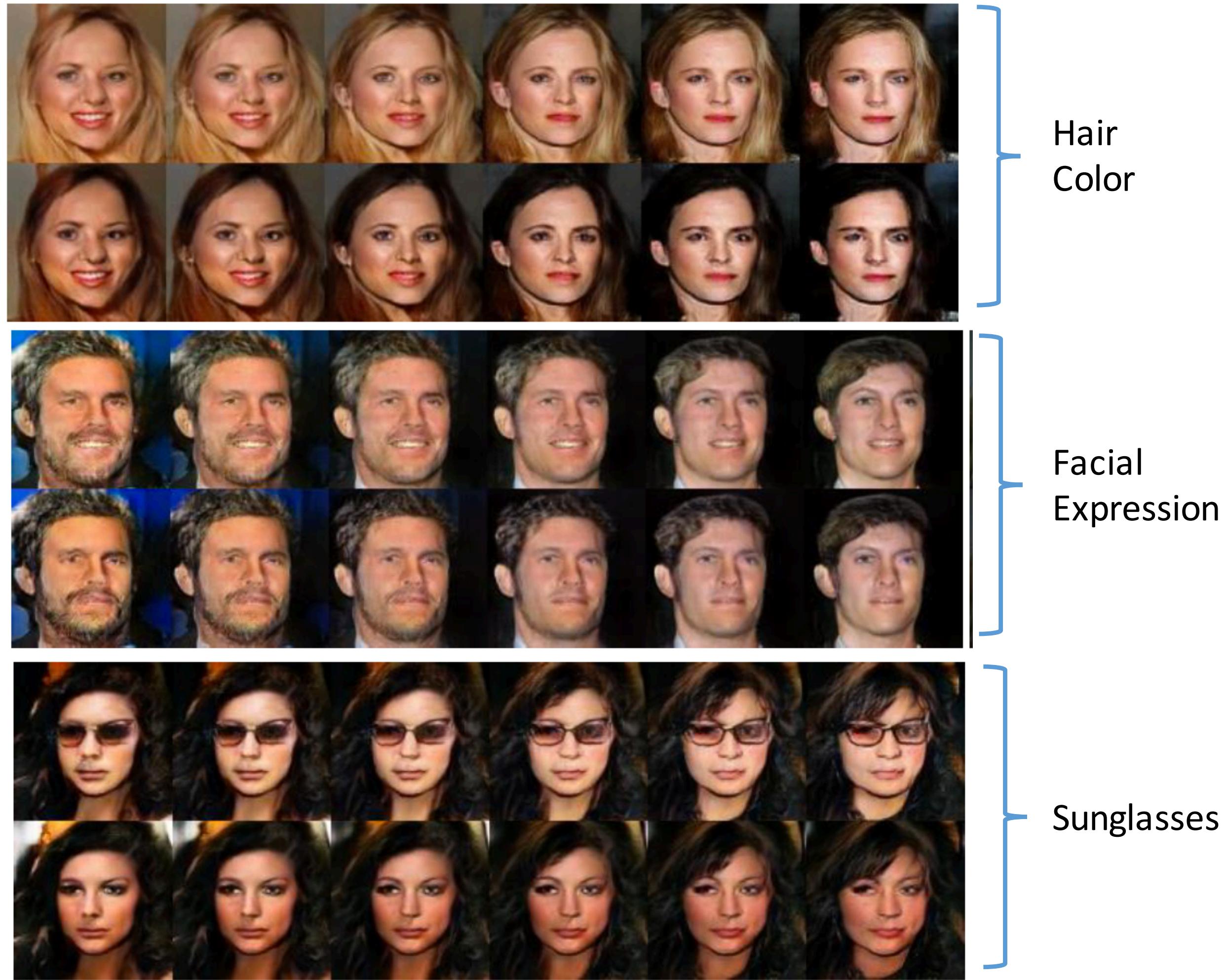


Training Generator



GANs in Action

Generating
Custom
Images



GANs in Action

Image-to-Image Translation

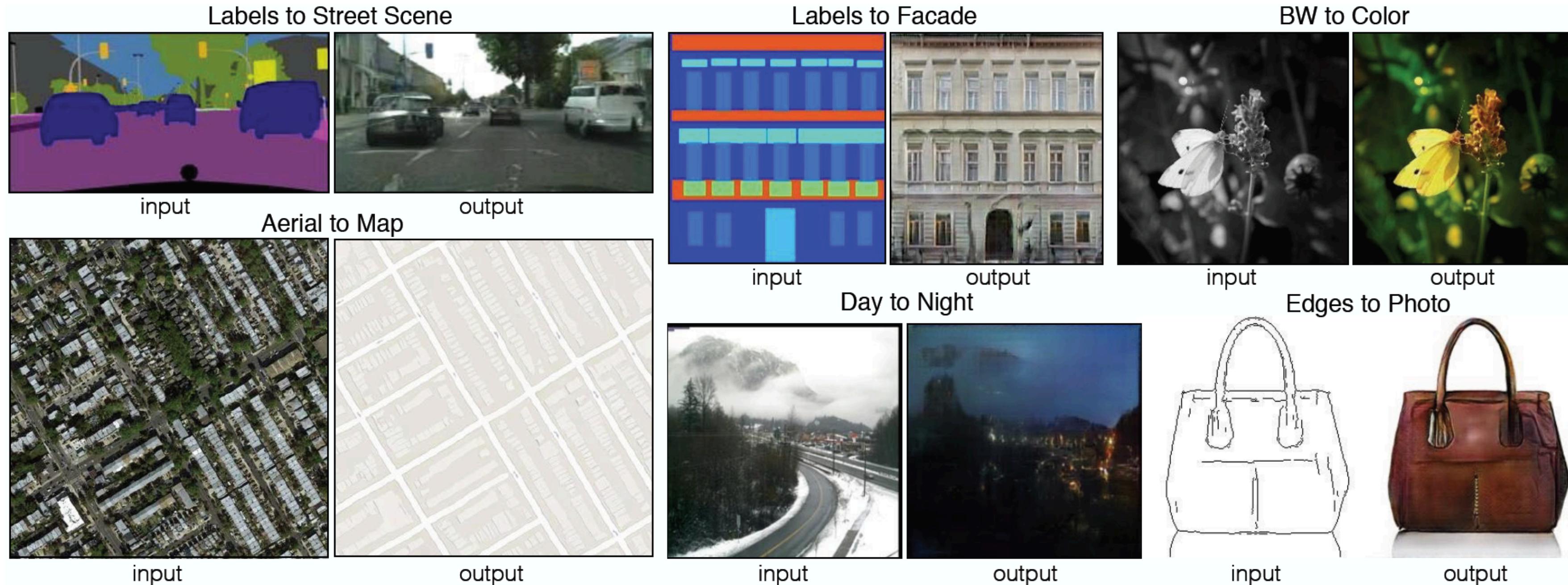


Figure 1 in the original paper.

[Link to an interactive demo of this paper](#)

GANs in Action

Text-to-Image Synthesis

Motivation

Given a text description, generate images closely associated.

Uses a conditional GAN with the generator and discriminator being condition on “dense” text embedding.

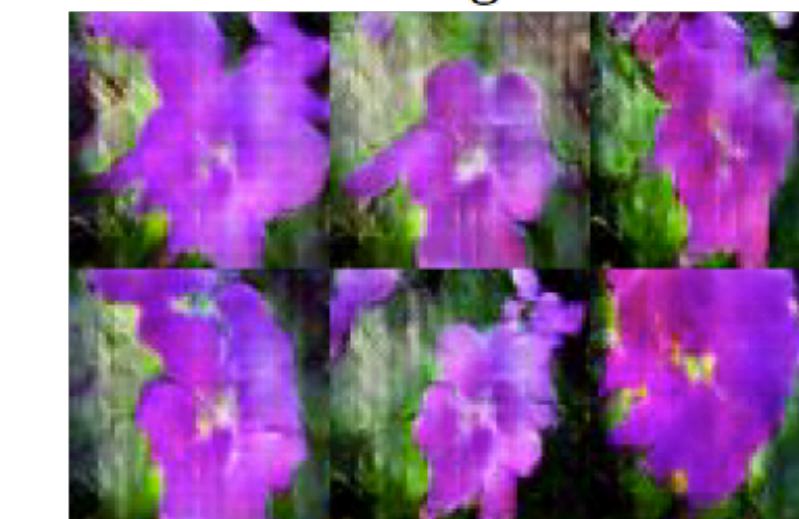
this small bird has a pink breast and crown, and black primaries and secondaries.



this magnificent fellow is almost all black with a red crest, and white cheek patch.



the flower has petals that are bright pinkish purple with white stigma



this white and yellow flower have thin white petals and a round yellow stamen

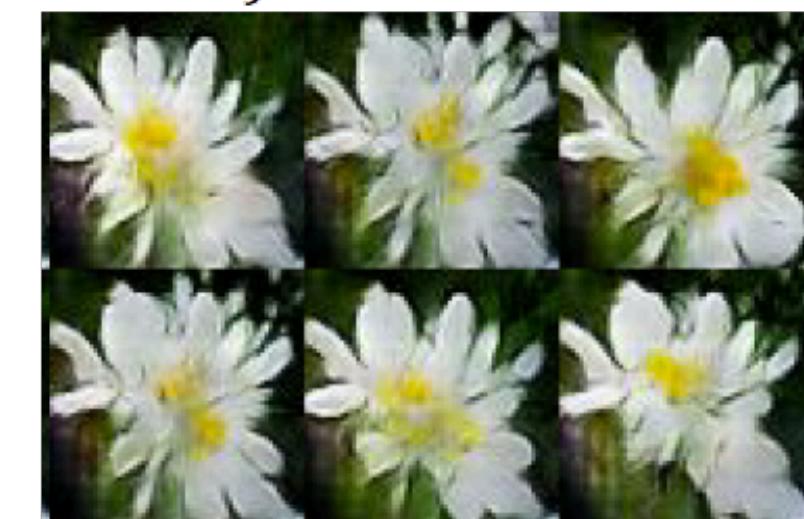


Figure 1 in the original paper.

GANs: Shortcomings

Problems with GANs

- **Probability Distribution is Implicit**

- Not straightforward to compute $P(X)$.
- Thus **Vanilla GANs** are only good for Sampling/Generation.

- **Training is Hard**

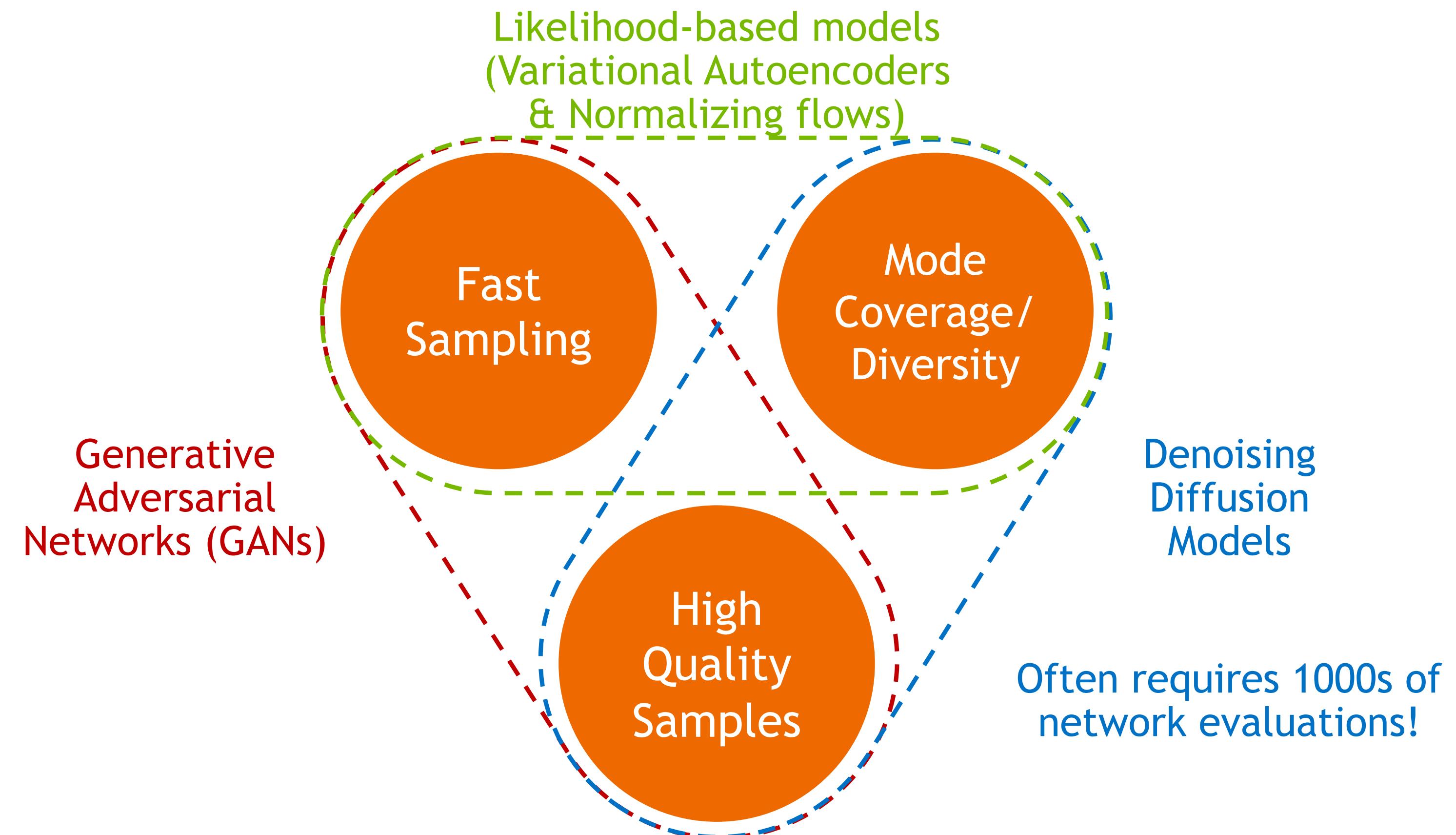
- Non-Convergence
- Mode-Collapse

Goodfellow, Ian. "NIPS 2016 Tutorial: Generative Adversarial Networks." arXiv preprint arXiv:1701.00160 (2016).

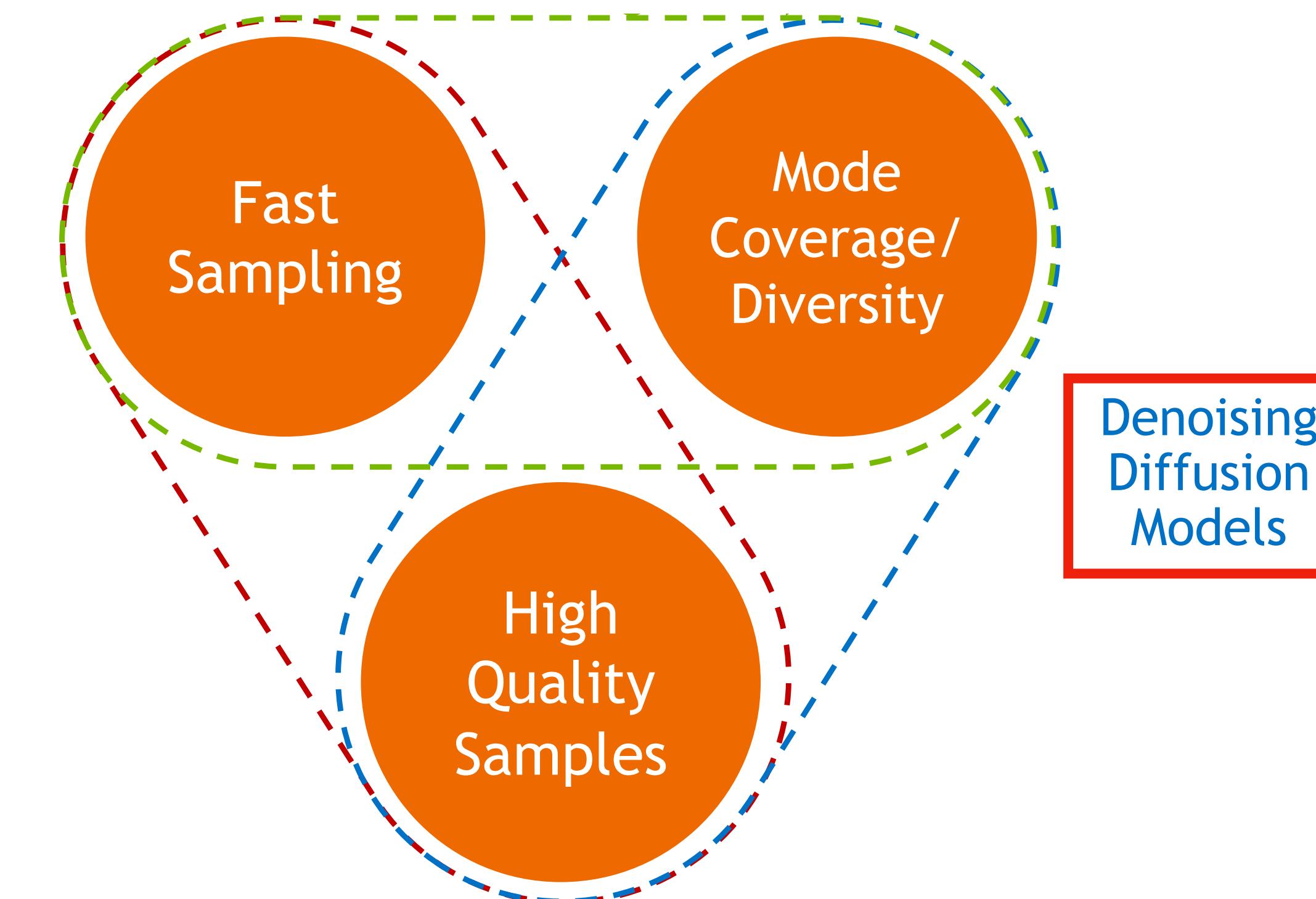
GANs: Shortcomings

- Training GANs is hard
 - Mode Collapse: all the images look similar
 - Non-convergence: does not learn anything good
- They are fast during inference though!

What makes a good generative model?
The generative learning trilemma



Deep Dive: Diffusion Models

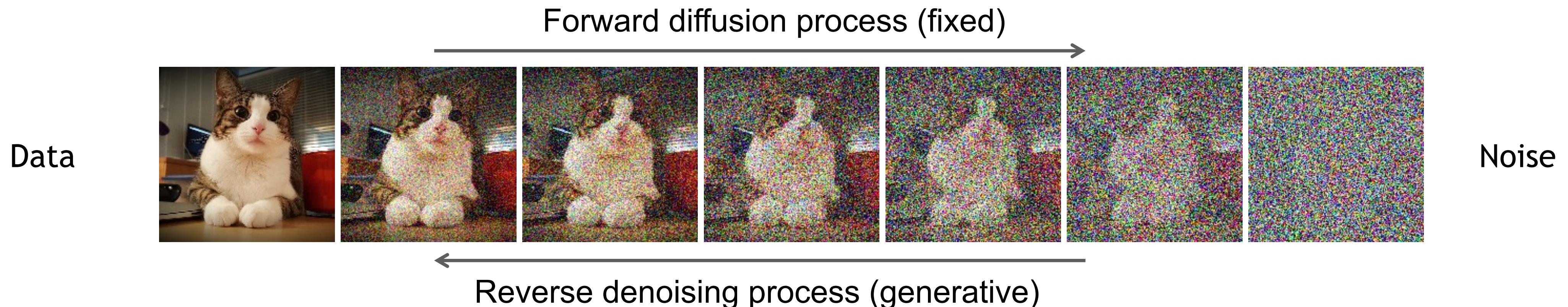


Denoising Diffusion Models

Learning to generate by denoising

Denoising diffusion models consist of two processes:

- Forward diffusion process that gradually adds noise to input
- Reverse denoising process that learns to generate data by denoising

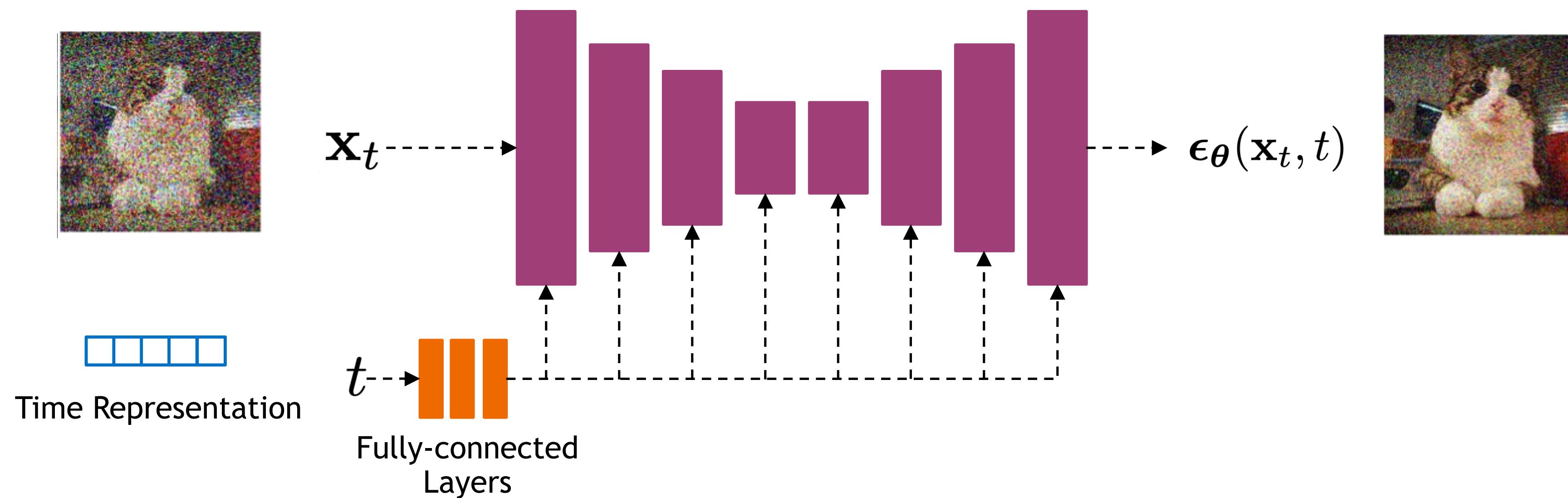


[Sohl-Dickstein et al., Deep Unsupervised Learning using Nonequilibrium Thermodynamics, ICML 2015](#)

[Ho et al., Denoising Diffusion Probabilistic Models, NeurIPS 2020](#)

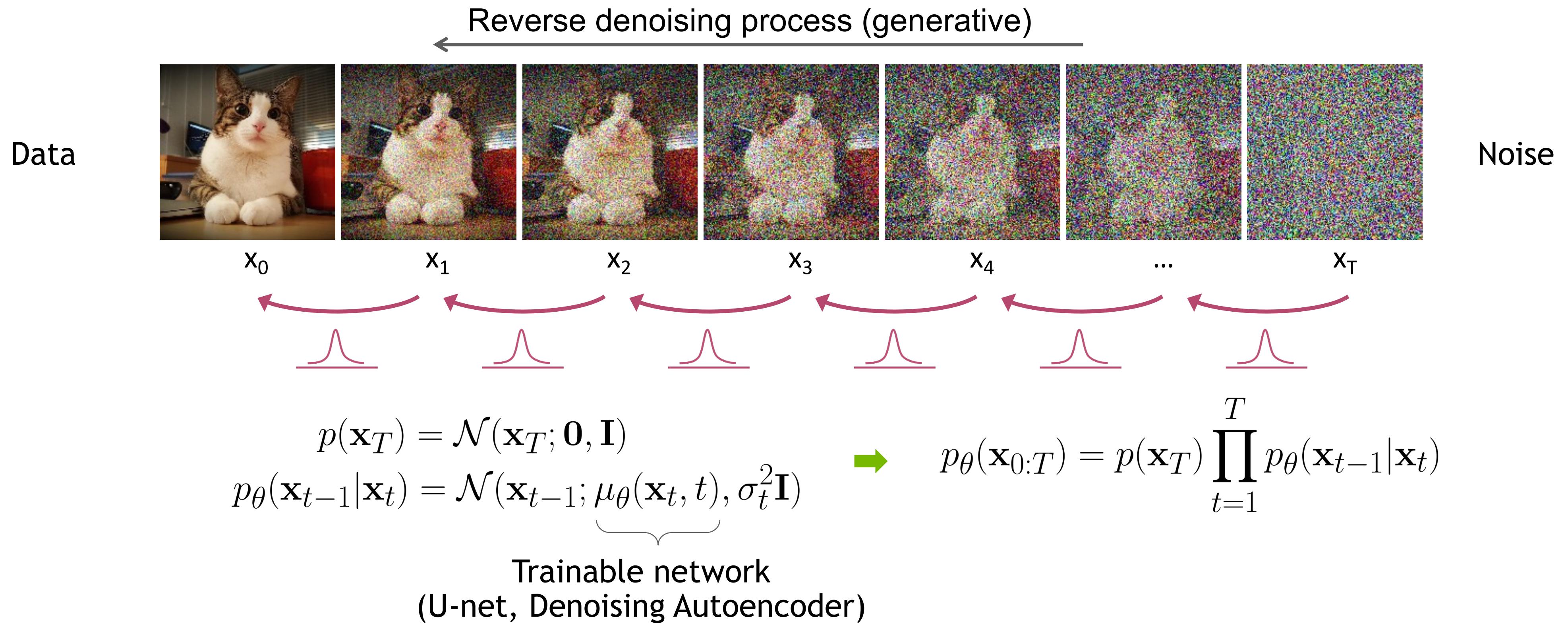
[Song et al., Score-Based Generative Modeling through Stochastic Differential Equations, ICLR 2021](#)

Reverse Denoising Process



Reverse Denoising Process

Formal definition of forward and reverse processes in T steps:



Training and Sample Generation

Algorithm 1 Training

```
1: repeat
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$ 
4:    $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on
       
$$\nabla_{\theta} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t)\|^2$$

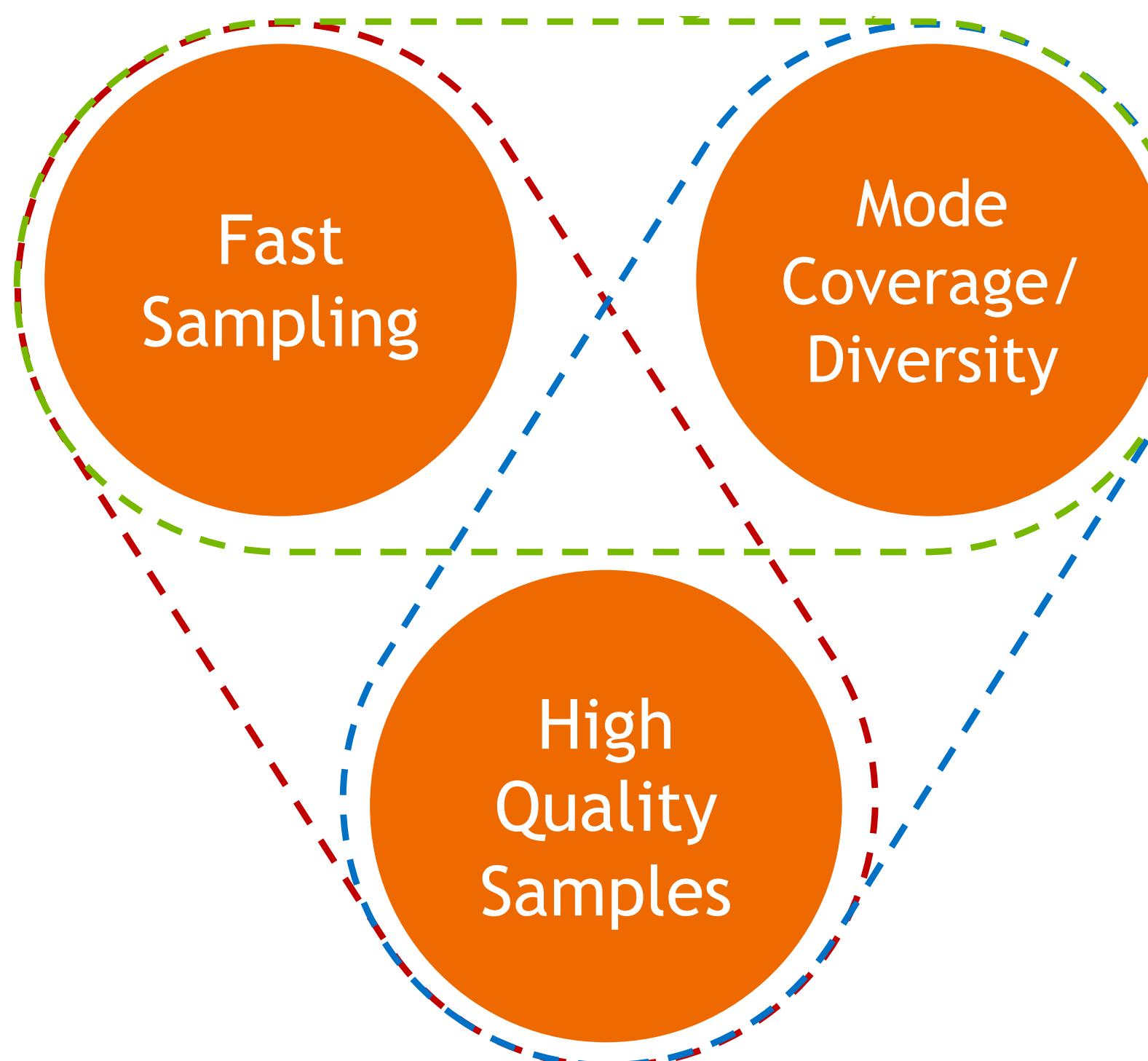
6: until converged
```

Algorithm 2 Sampling

```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
4:   
$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$$

5: end for
6: return  $\mathbf{x}_0$ 
```

Diffusion Model Inference is relative slow
(T network iterations for each generation)



Algorithm 2 Sampling

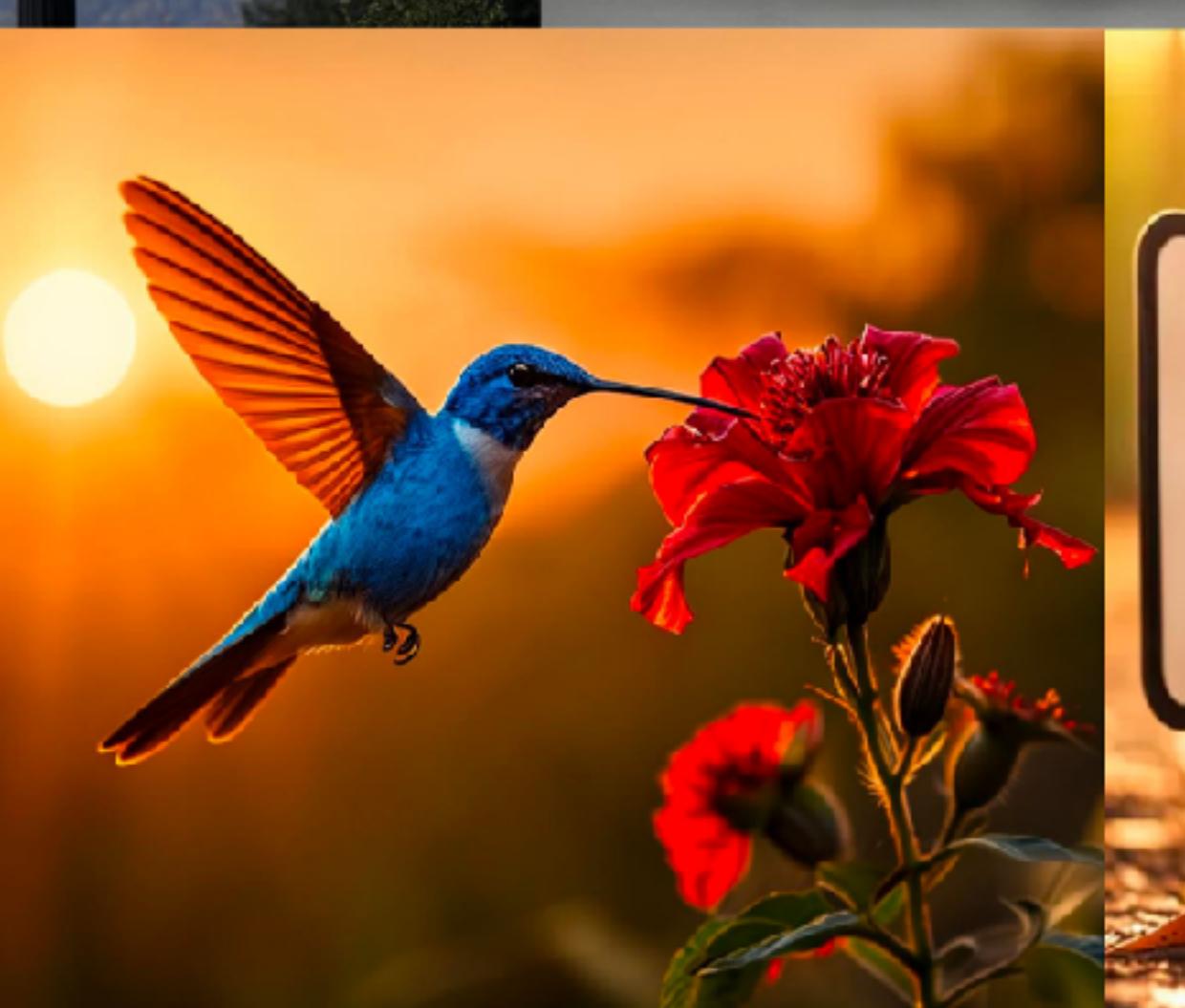
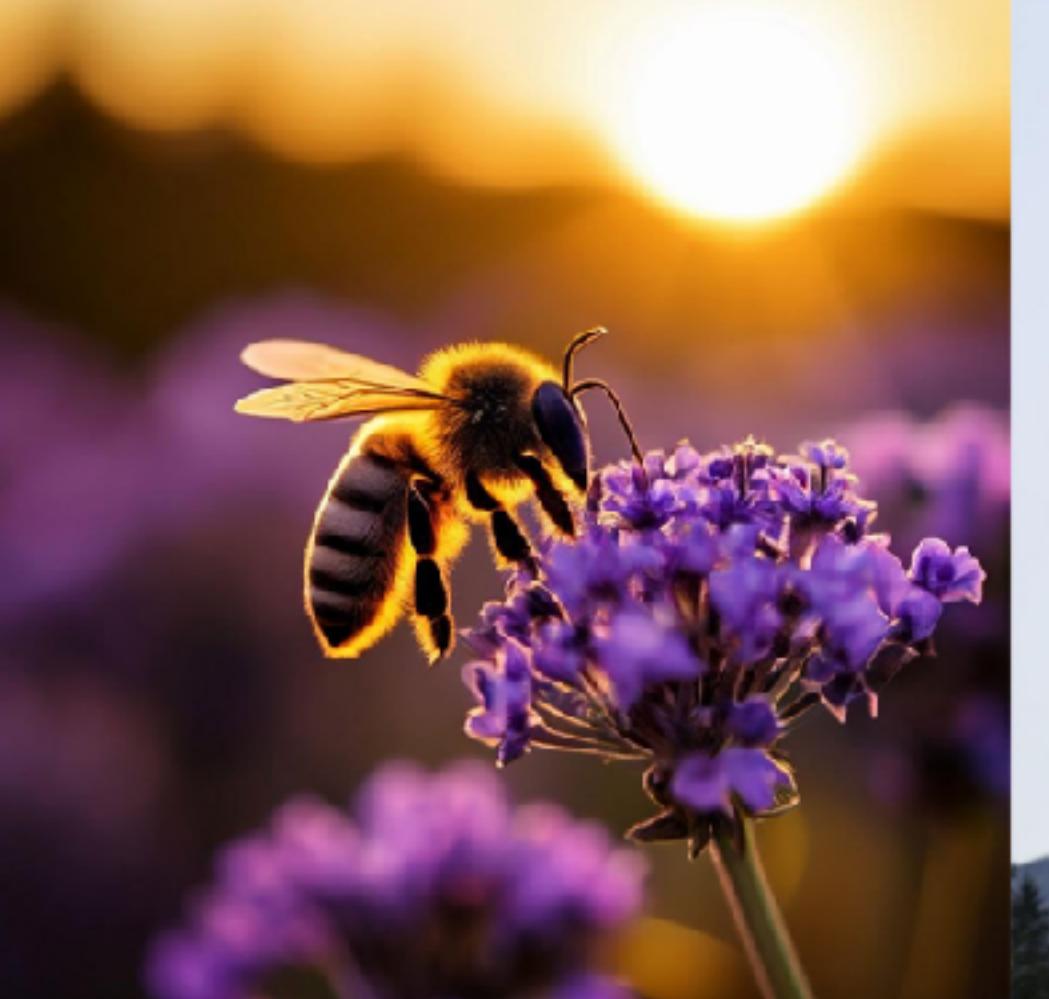
```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
```

Diffusion Models: The Fancy

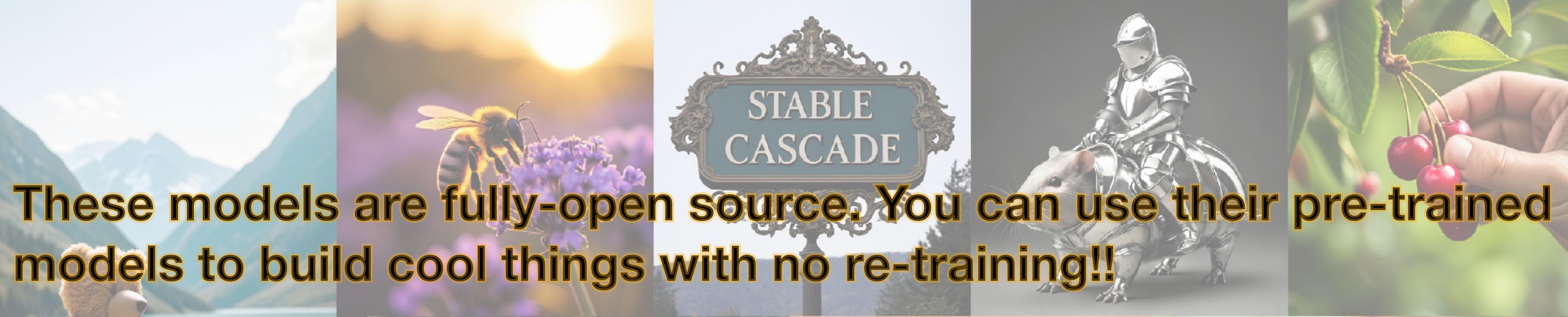


What can you do with Diffusion Models?

- A. Stable Diffusion
- B. SDEdit
- C. Imagic
- D. InstructPix2Pix
- E. DreamBooth
- F. Textual Inversion



- Stable Diffusion: <https://github.com/Stability-AI/stablediffusion>
- Stable Cascade: <https://github.com/Stability-AI/StableCascade>

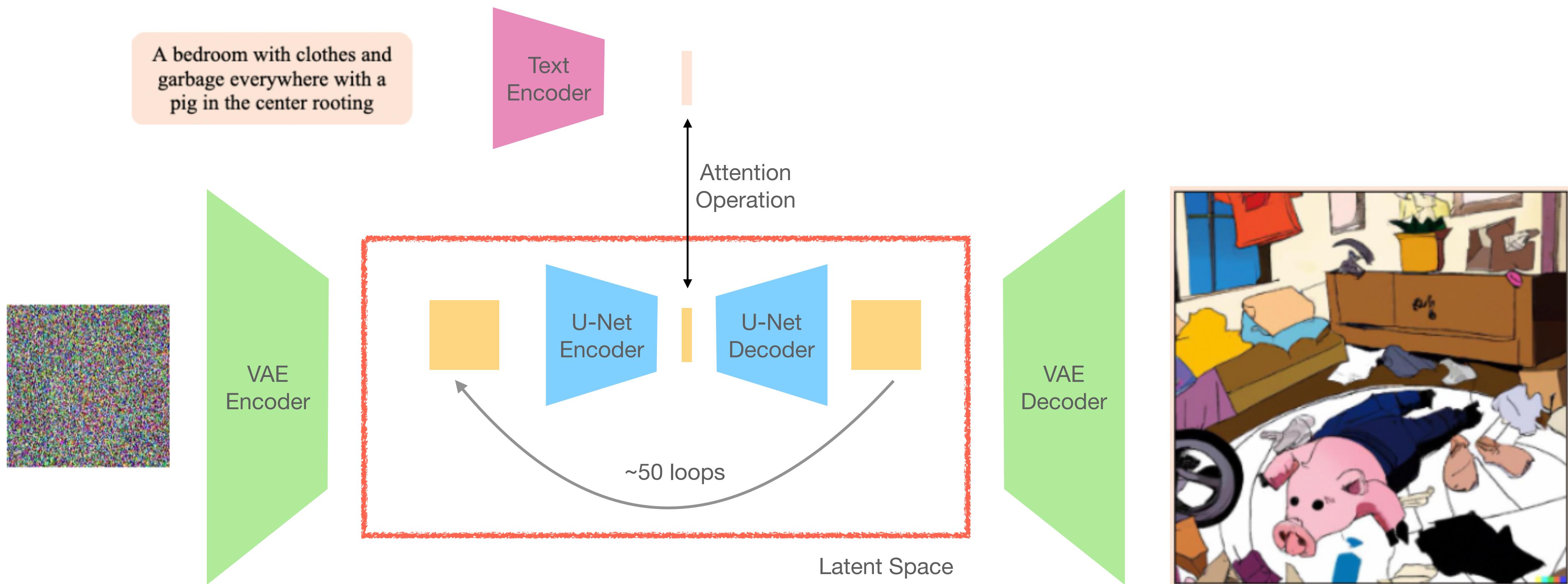


These models are fully-open source. You can use their pre-trained models to build cool things with no re-training!!



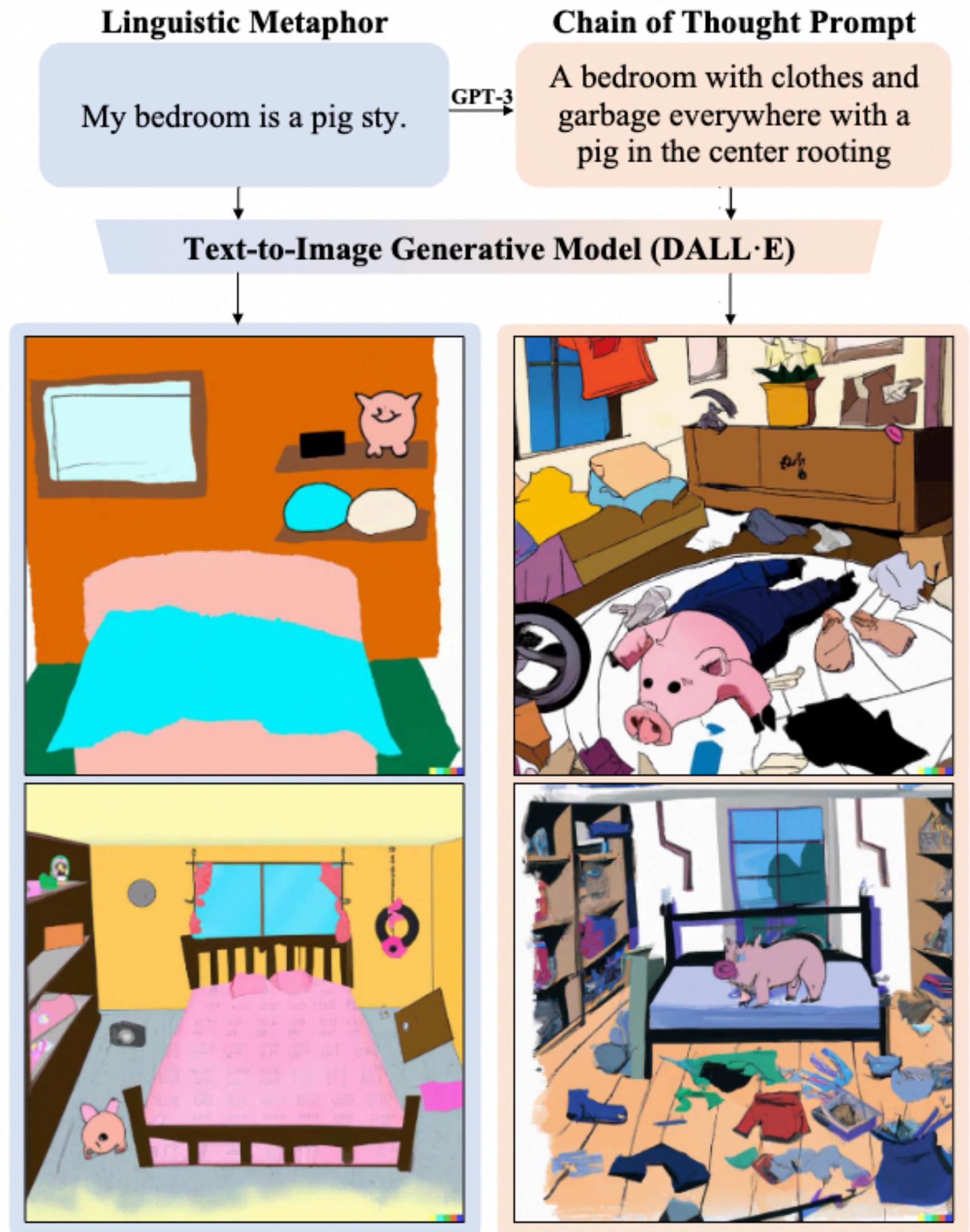
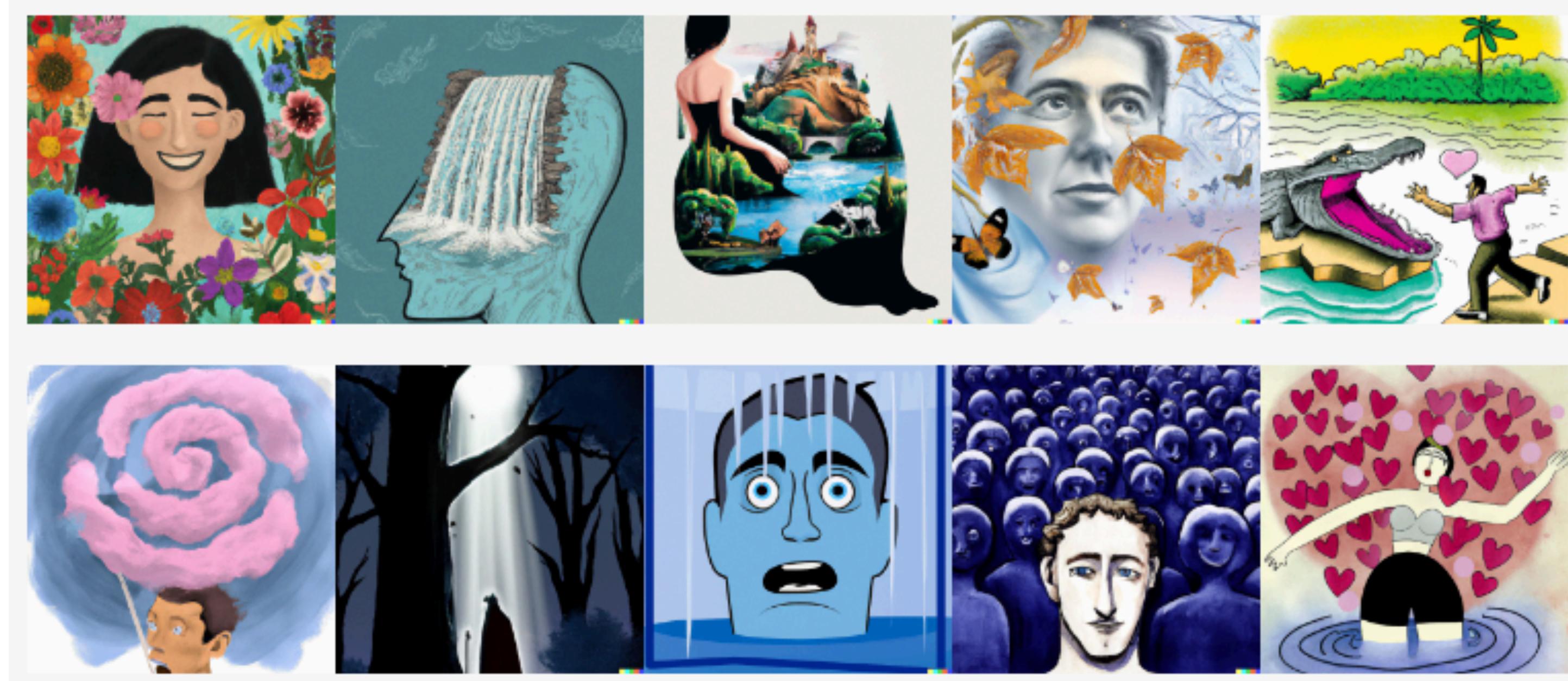
Stable Diffusion: insides?

A Latent Diffusion Model



Using Stable Diffusion Artistic Projects?

- Combine with Large Language Models
- Generate emotions & metaphors in visual form



Using Stable Diffusion

Personal Project: See <https://diffusionillusions.com>

- Minimal modifications to Stable Diffusion allows us to generate Visual Illusions
- Outstanding Demo at CVPR (top vision conference)
- All experiments on a single GPU machine
- Inference works on colab (~6 minutes per generation)

Flip Illusion



180°
rotate



Using Stable Diffusion

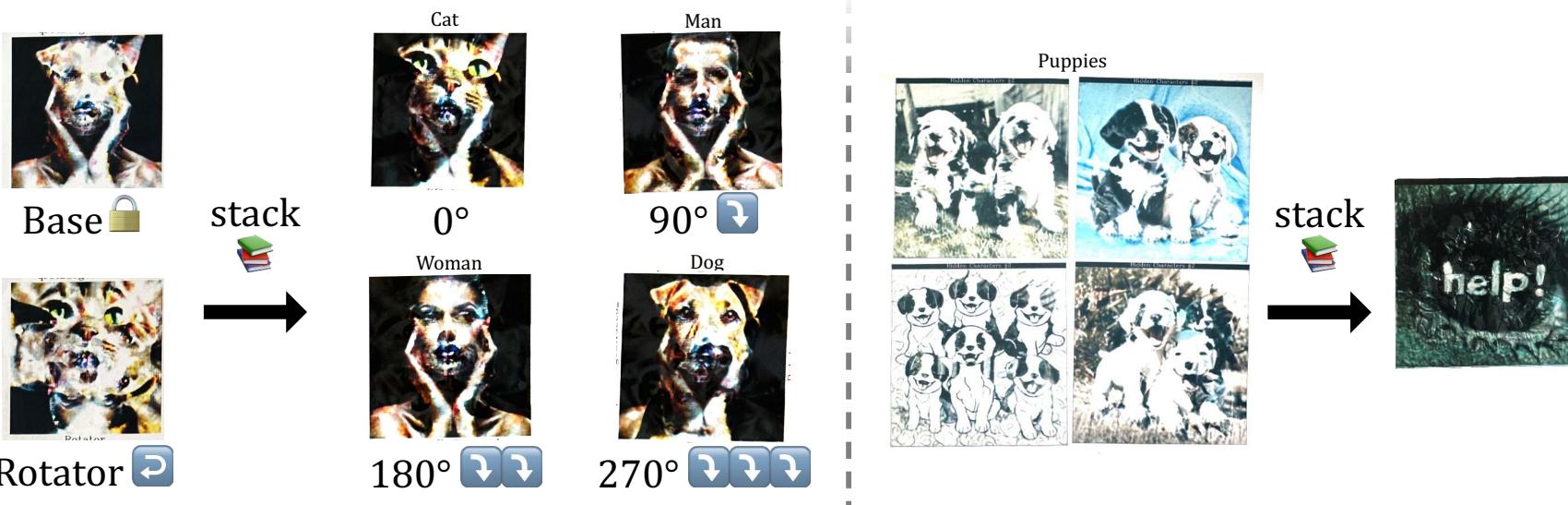
Personal Projects? See <https://diffusionillusions.com>

Flip Illusion

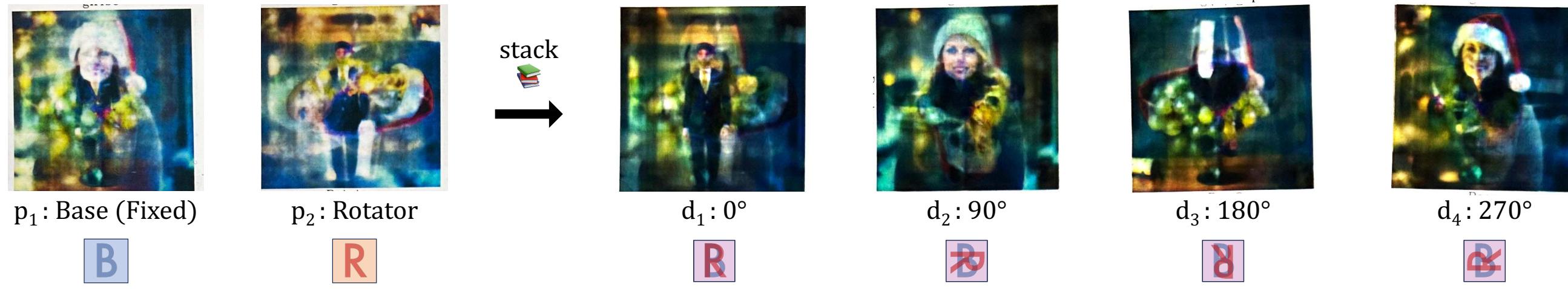


180°
rotate
↻

Rotation Overlay Illusion



Another Rotation Overlay Illusion Example



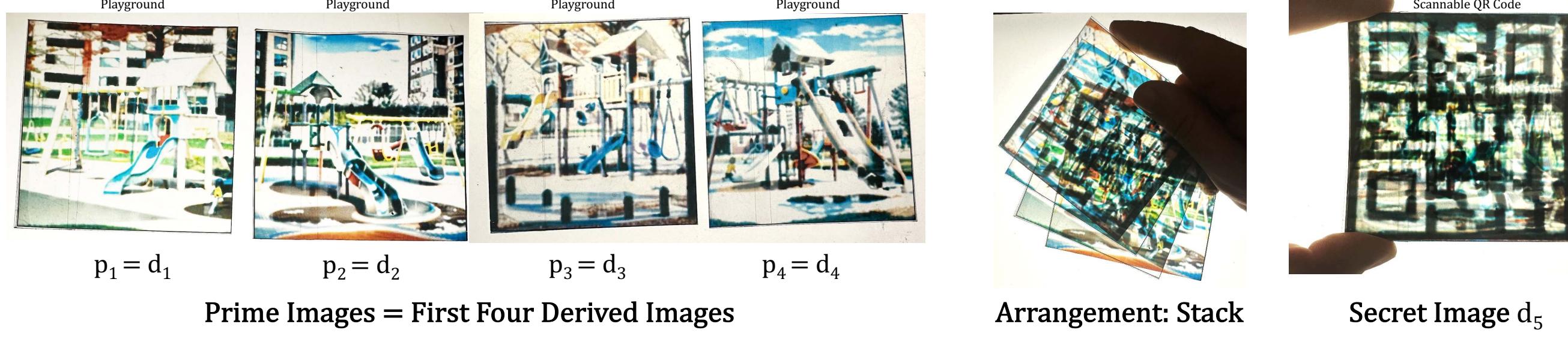
p₁: Base (Fixed)

p₂: Rotator

B

R

Another Hidden Overlay Illusion Example



p₁ = d₁

p₂ = d₂

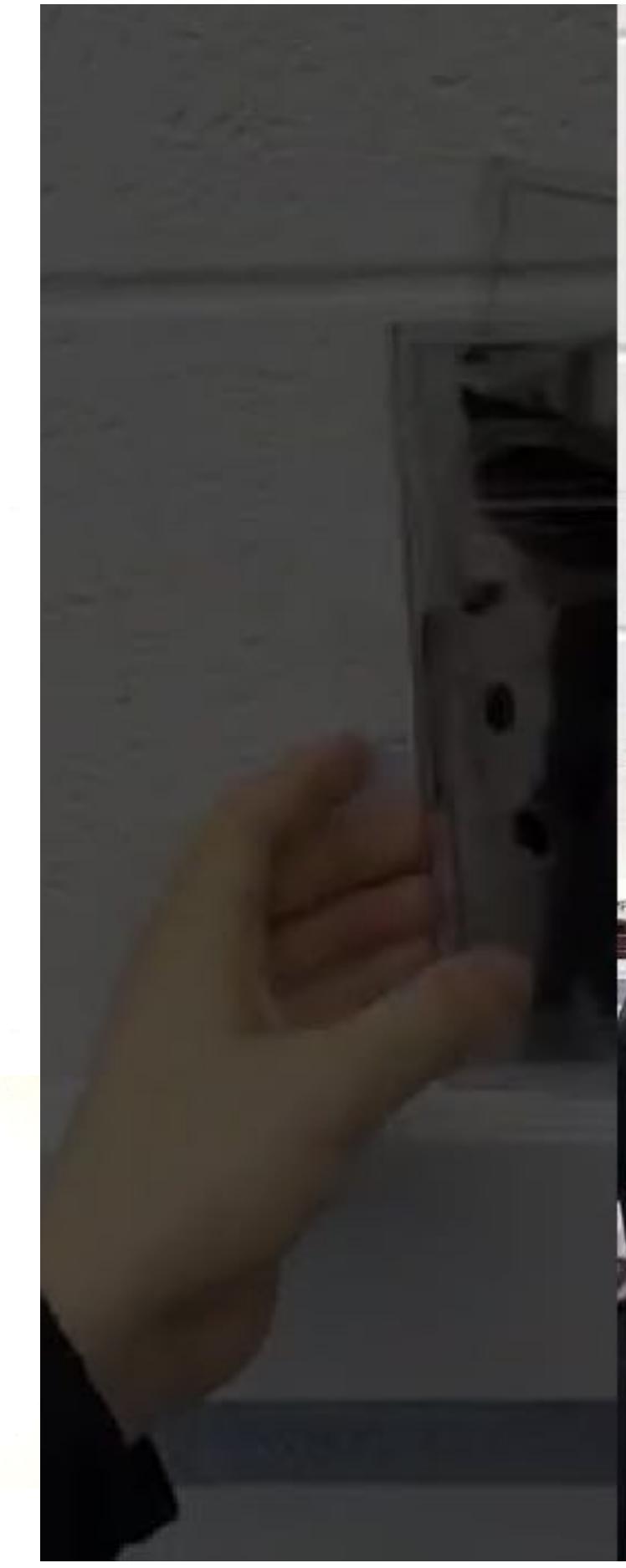
p₃ = d₃

p₄ = d₄

Prime Images = First Four Derived Images

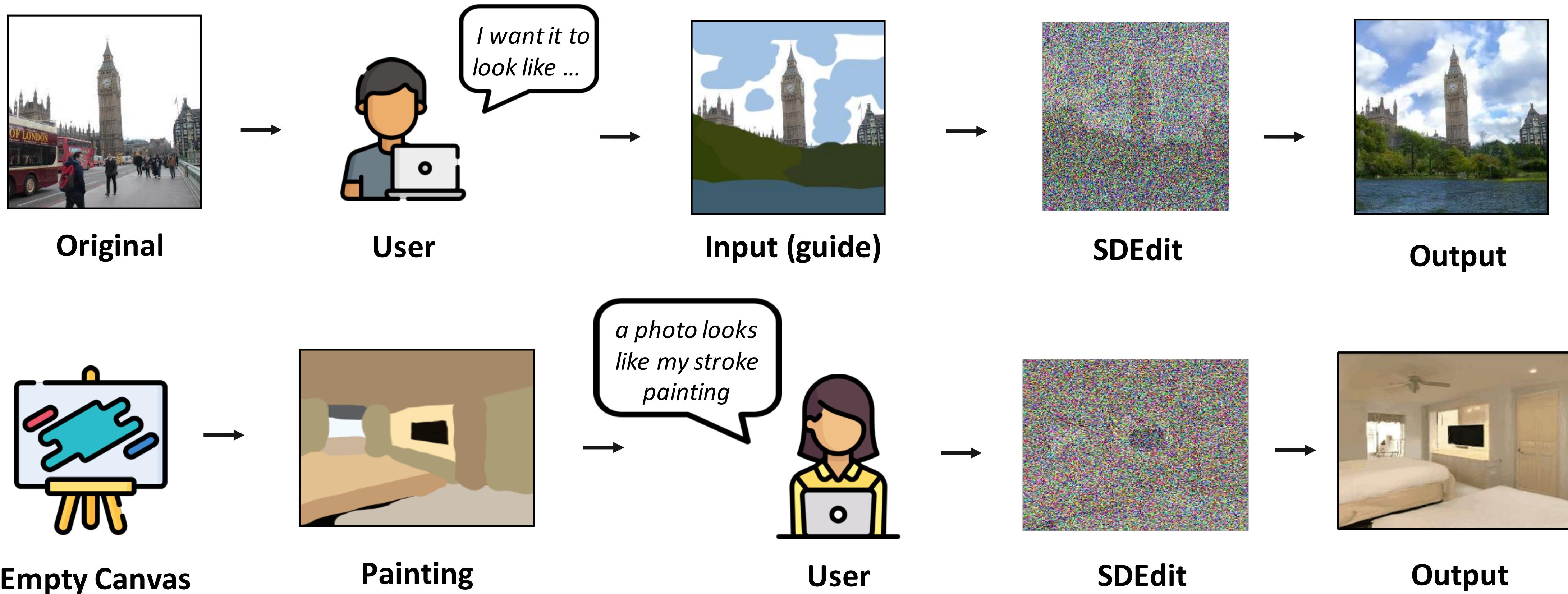
Arrangement: Stack

Secret Image d₅

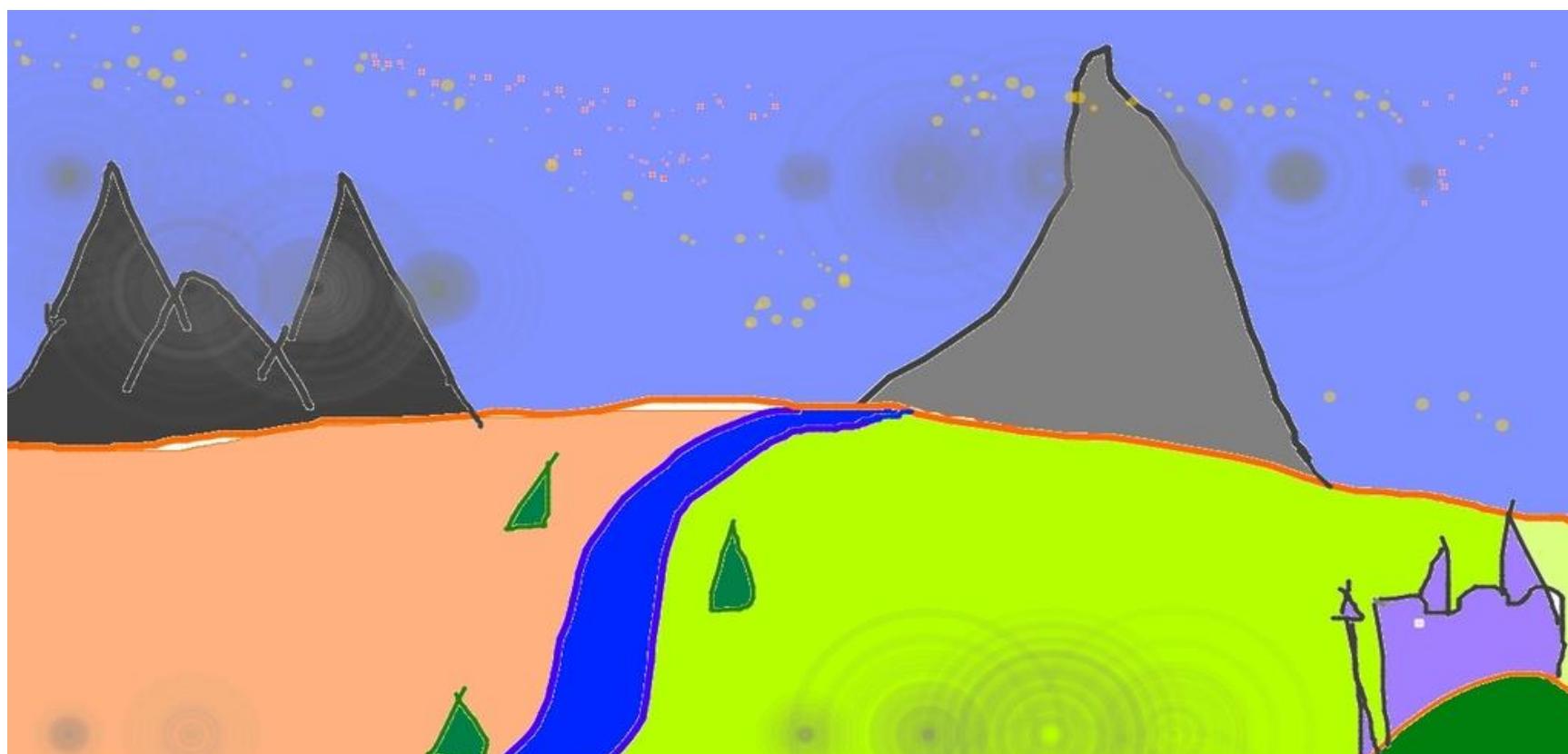


diffusionillusions.com

B. SDEdit

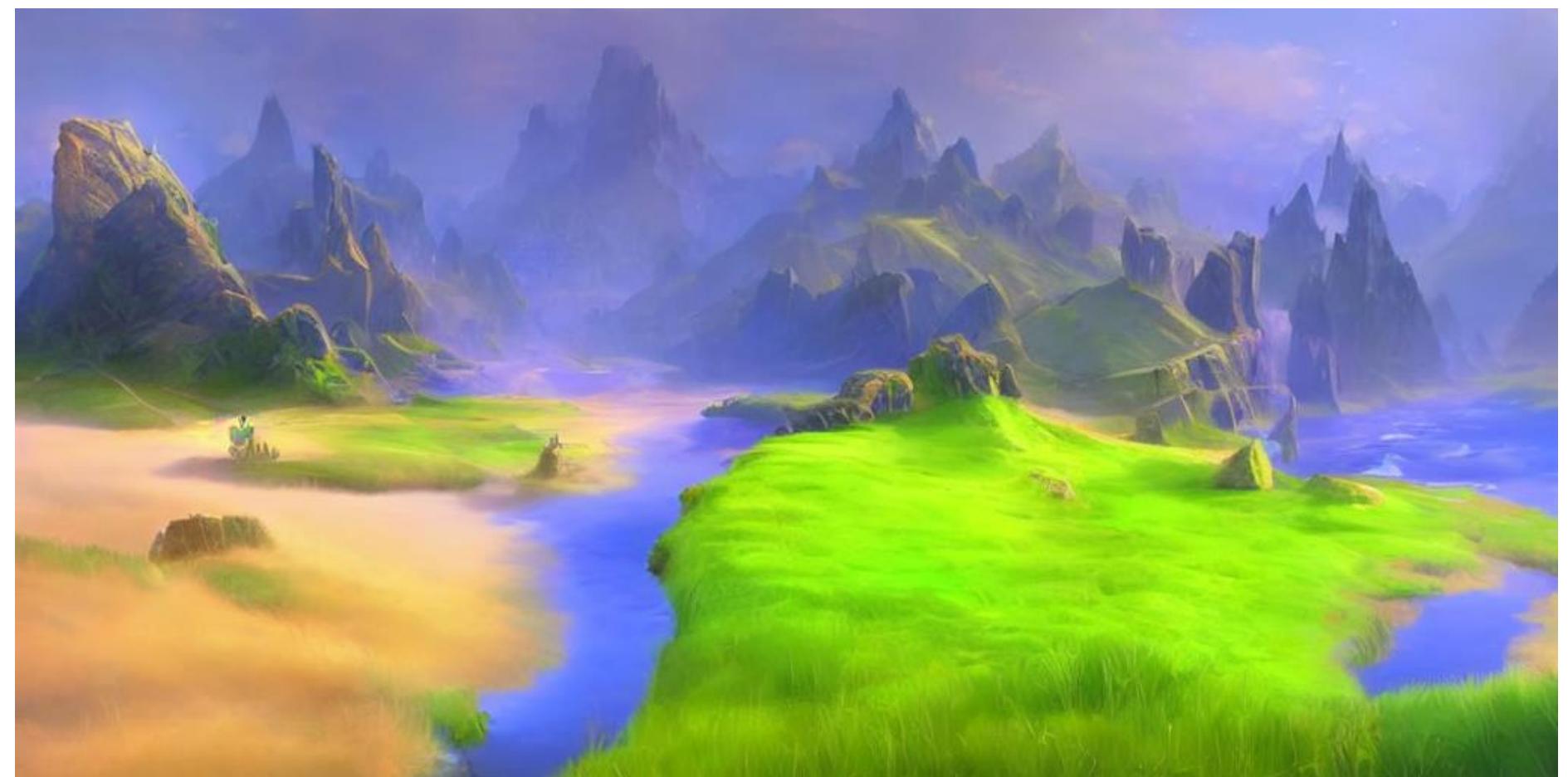


Application to Stable Diffusion (img2img)



Input

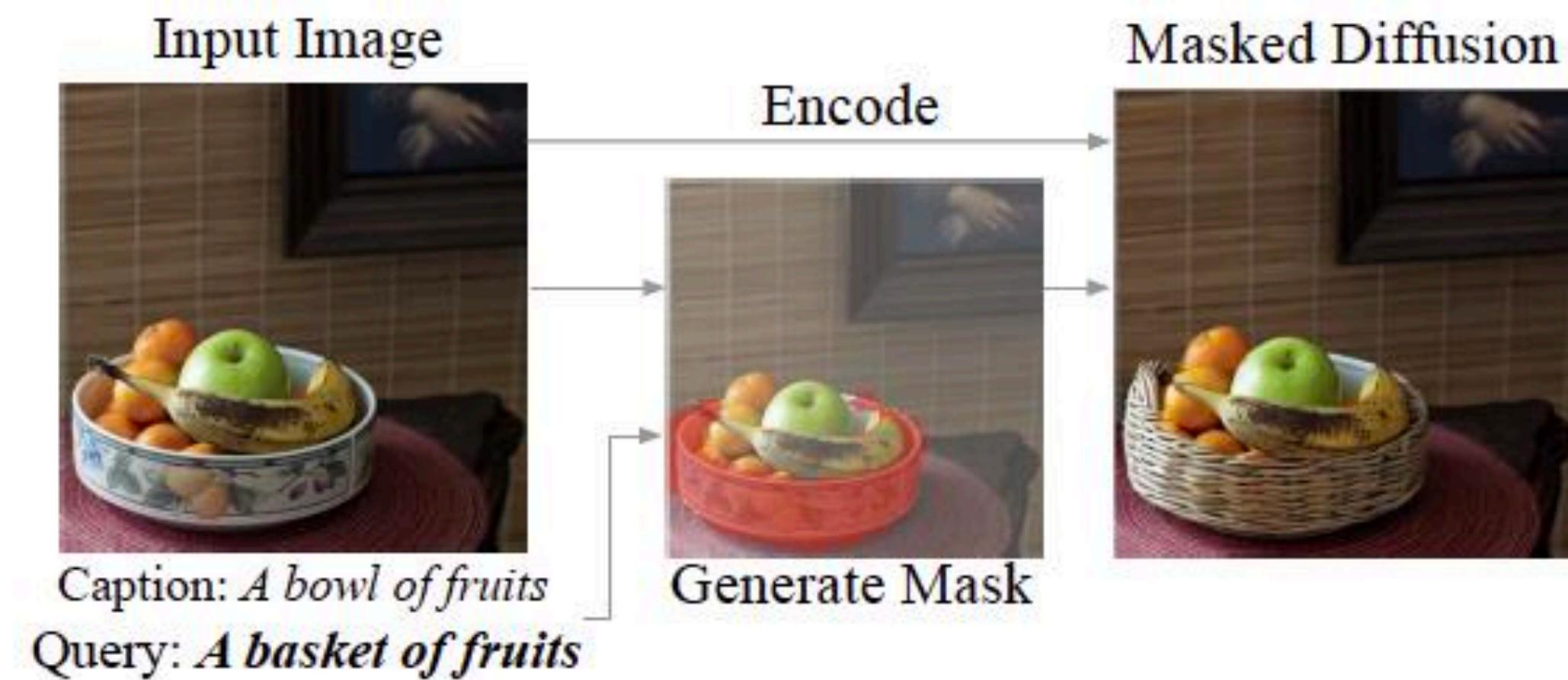
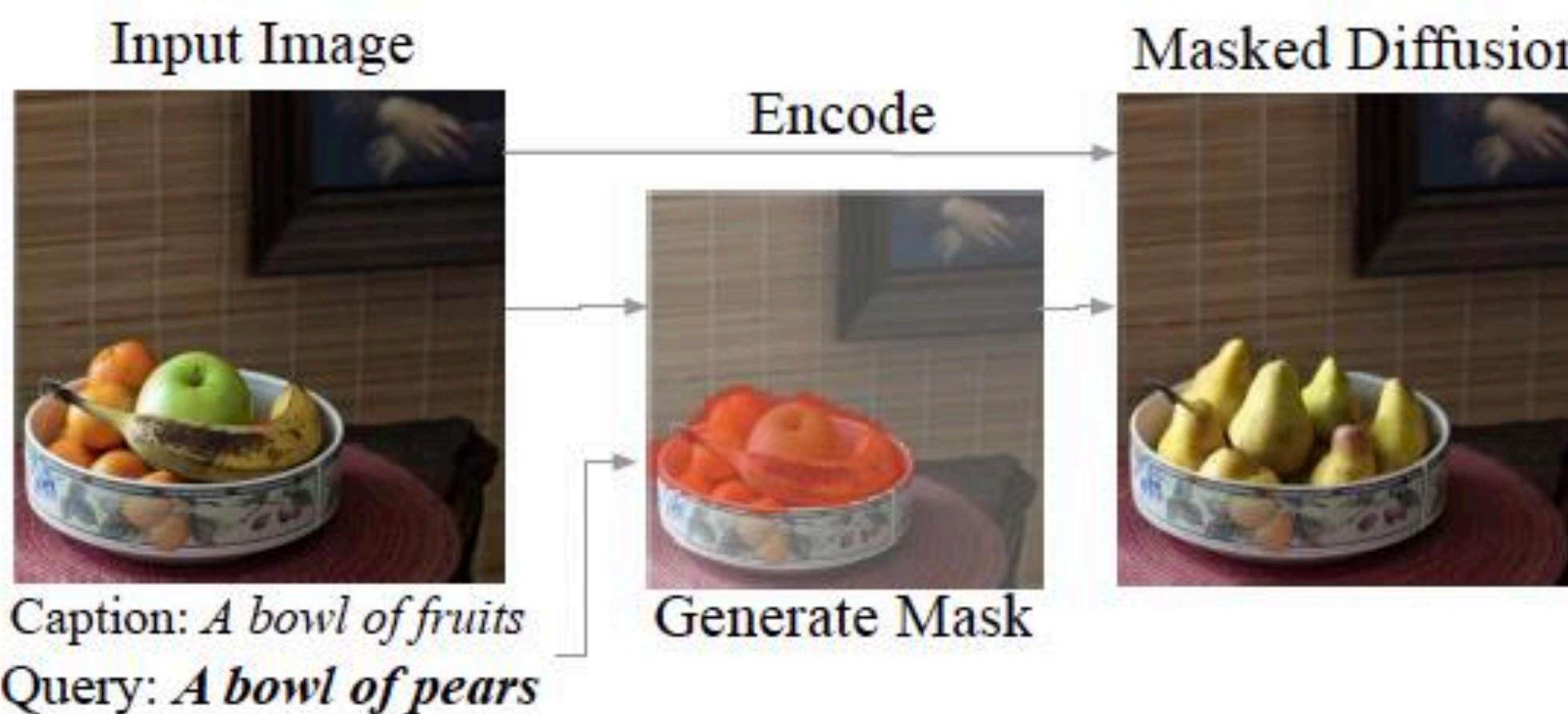
SDEdit/img2img
→



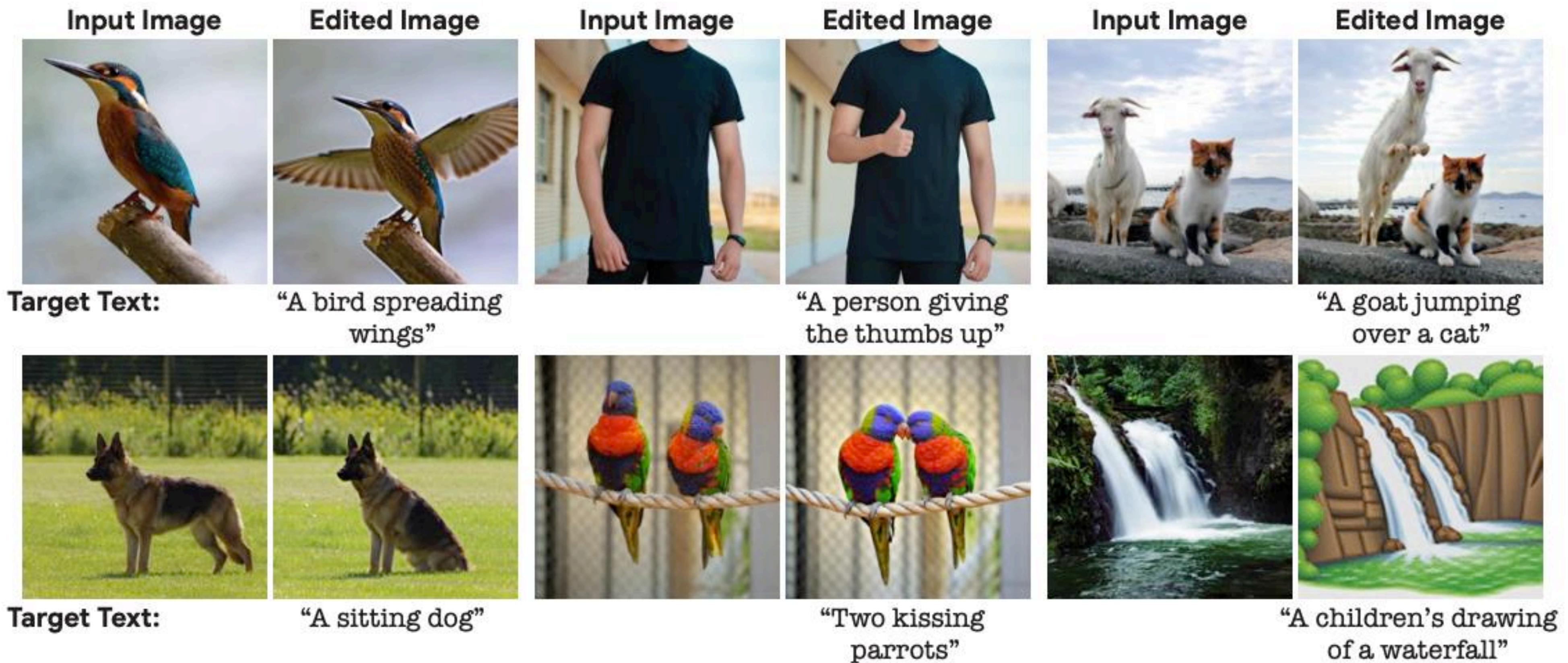
Outputs

DiffEdit: Diffusion-based semantic image editing with mask guidance

Instead of asking users to provide the mask, the model will generate the mask itself based on the caption and query.



C. Imagic



Imagic: Text-Based Real Image Editing with Diffusion Models



D. InstructPix2Pix



D. InstructPix2Pix

Training Data Generation

(a) Generate text edits:

Input Caption: "photograph of a girl riding a horse" → 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" → Stable Diffusion + Prompt2Prompt → Edited Caption: "photograph of a girl riding a dragon" → 

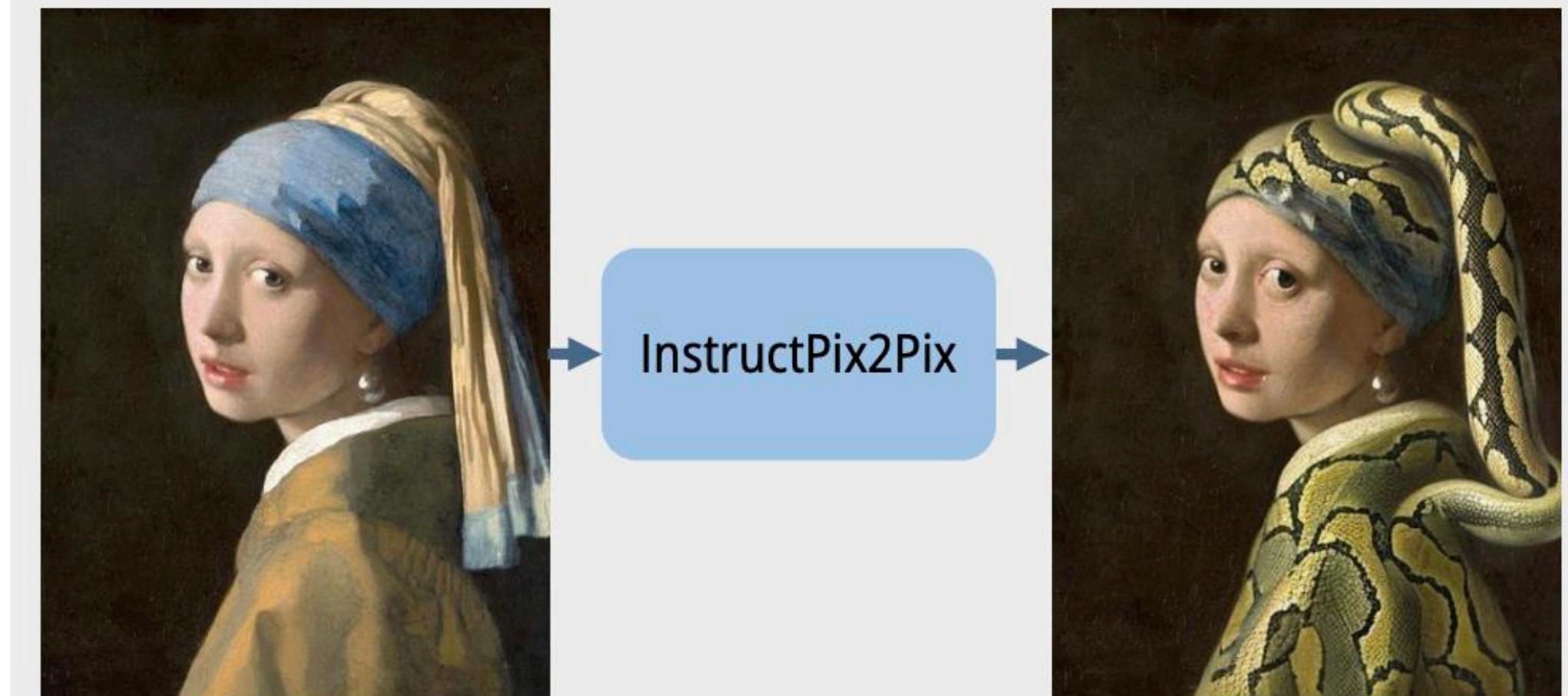
(c) Generated training examples:



Instruction-following Diffusion Model

(d) Inference on real images:

"turn her into a snake lady"



F. DreamBooth: Personalization!



Input images

Make it generate only
images of your dog!!



in the Acropolis



swimming



in a doghouse



sleeping



in a bucket

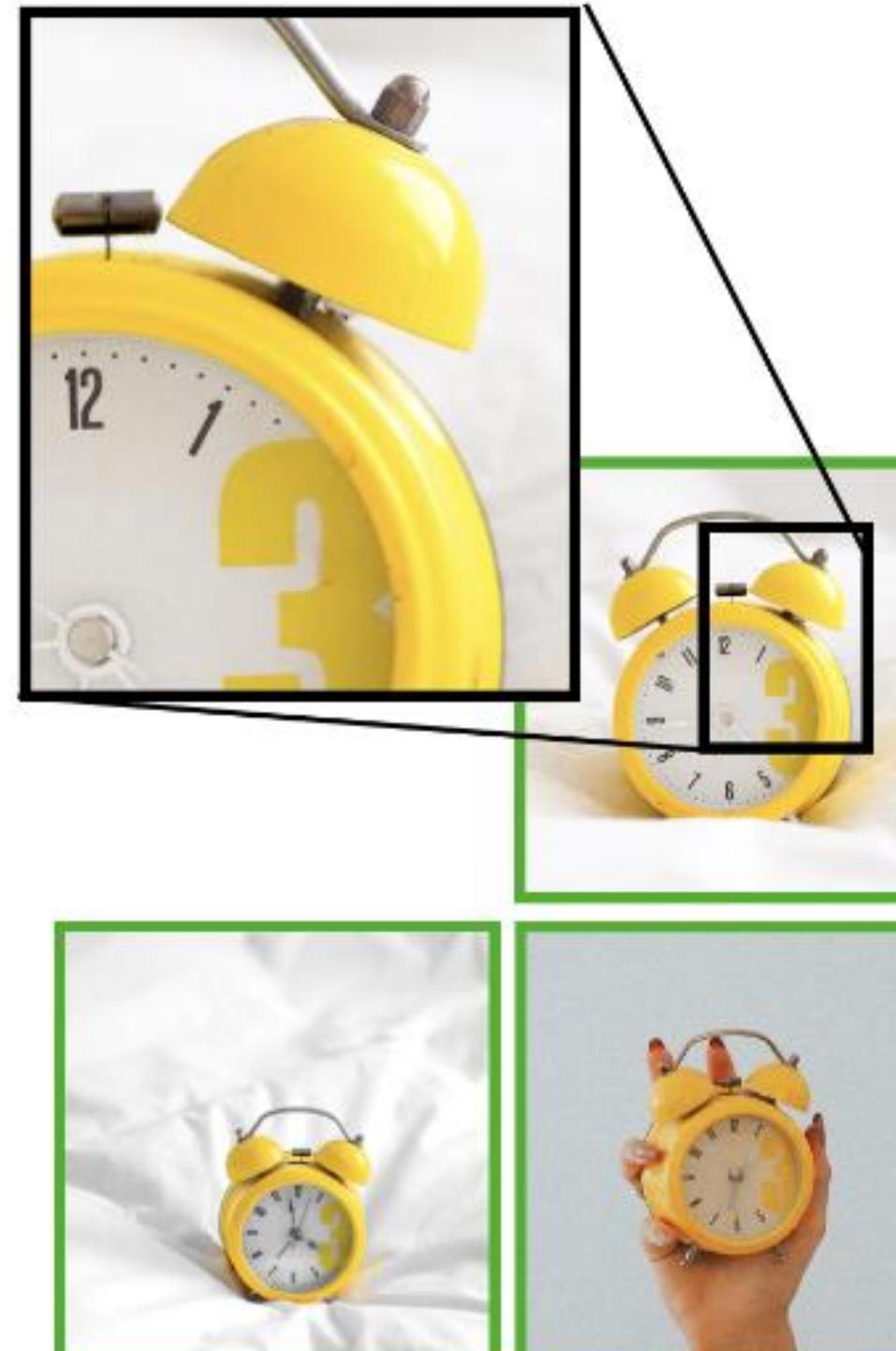


getting a haircut

Generated images

F. DreamBooth: Personalization!

Efficient Fine Tuning to Personalize a Diffusion Model



Input Images



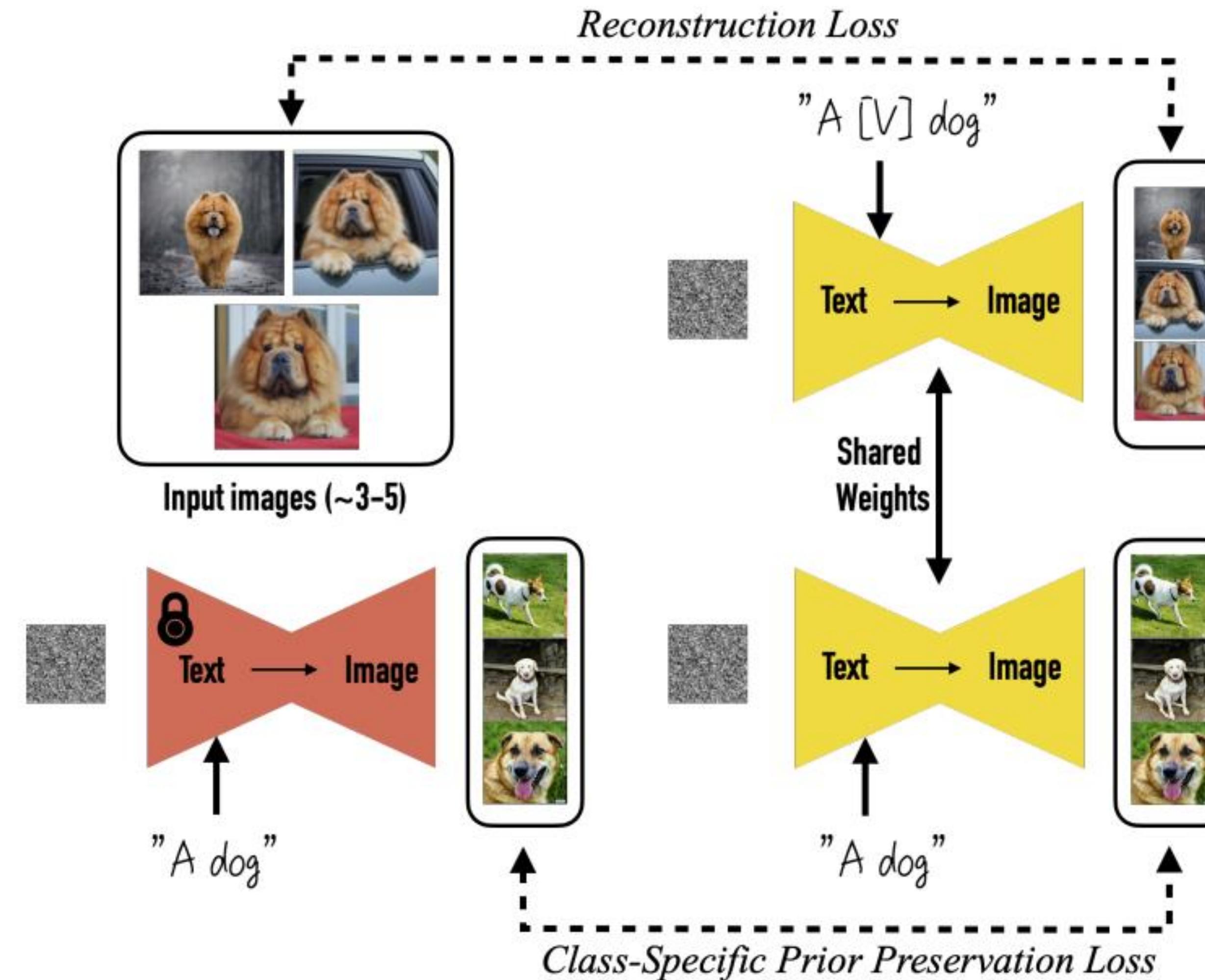
Image-guided, DALL-E2



Text-guided, Imagen

F. DreamBooth: Personalization!

Efficient Fine Tuning to Personalize a Diffusion Model



F. DreamBooth: Personalization!



Input images



A [V] backpack in the
Grand Canyon



A wet [V] backpack
in water



A [V] backpack in Boston



A [V] backpack with the
night sky



Input images



A [V] teapot floating
in milk



A transparent [V] teapot
with milk inside



A [V] teapot
pouring tea



A [V] teapot floating
in the sea

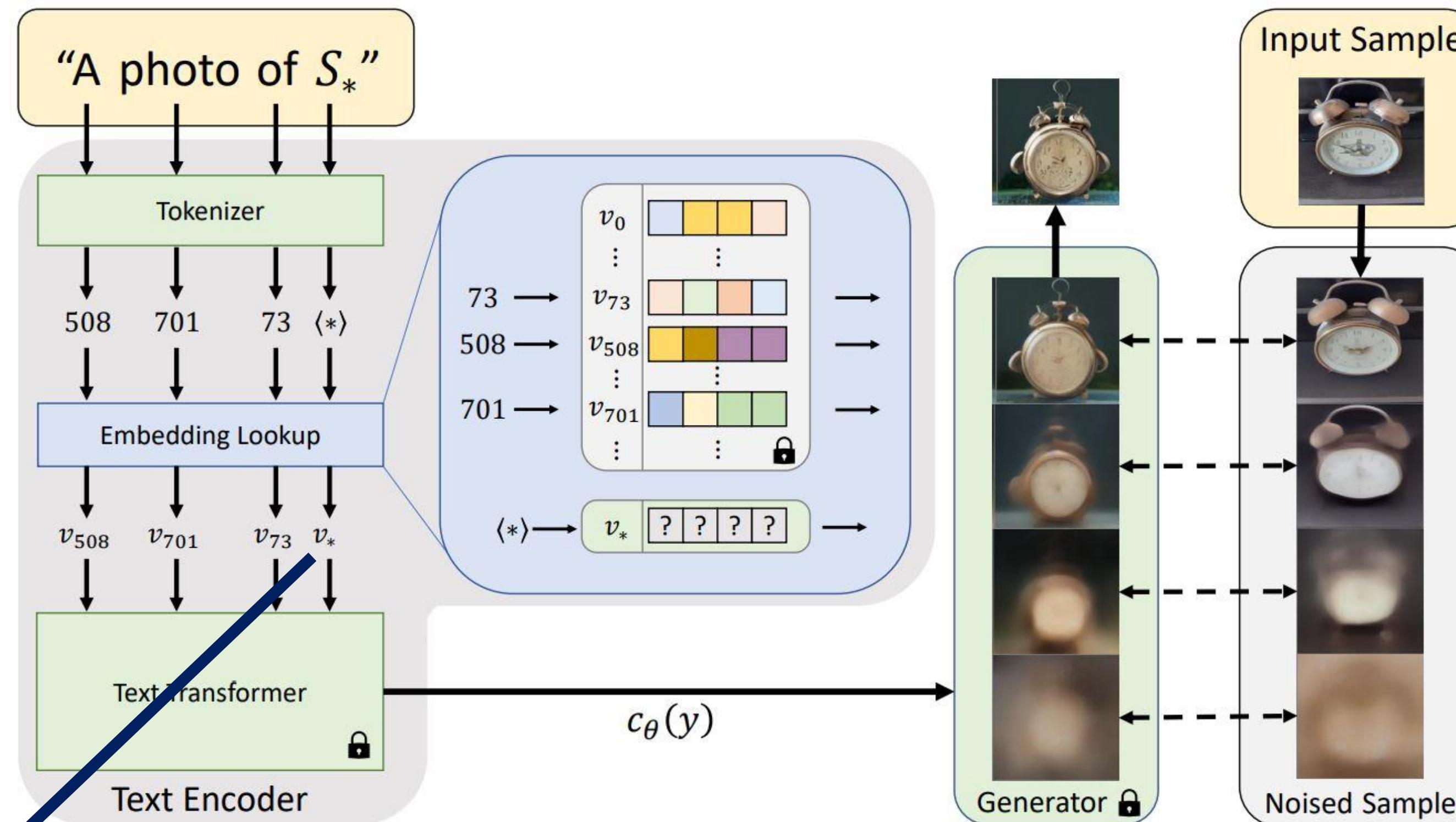
G. Textual Inversion: Personalization!

Learn A Text Embedding to Personalize a Diffusion Model



G. Textual Inversion: Personalization!

Learning (or optimizing) the Personalized Text Embedding



$$v_* = \arg \min_v \mathbb{E}_{z \sim \mathcal{E}(x), y, \epsilon \sim \mathcal{N}(0, 1), t} \left[\|\epsilon - \epsilon_\theta(z_t, t, c_\theta(y))\|_2^2 \right]$$

G. Textual Inversion: Some Results



Input samples

→



“ S_* sports car”



“ S_* made of lego”



“ S_* onesie”



“da Vinci sketch of S_* ”



Input samples

→



“Manga drawing of a steaming S_* ”



“A S_* watering can”



“ S_* Death Star”



“A poster for the movie ‘The Teapot’ starring S_* ”

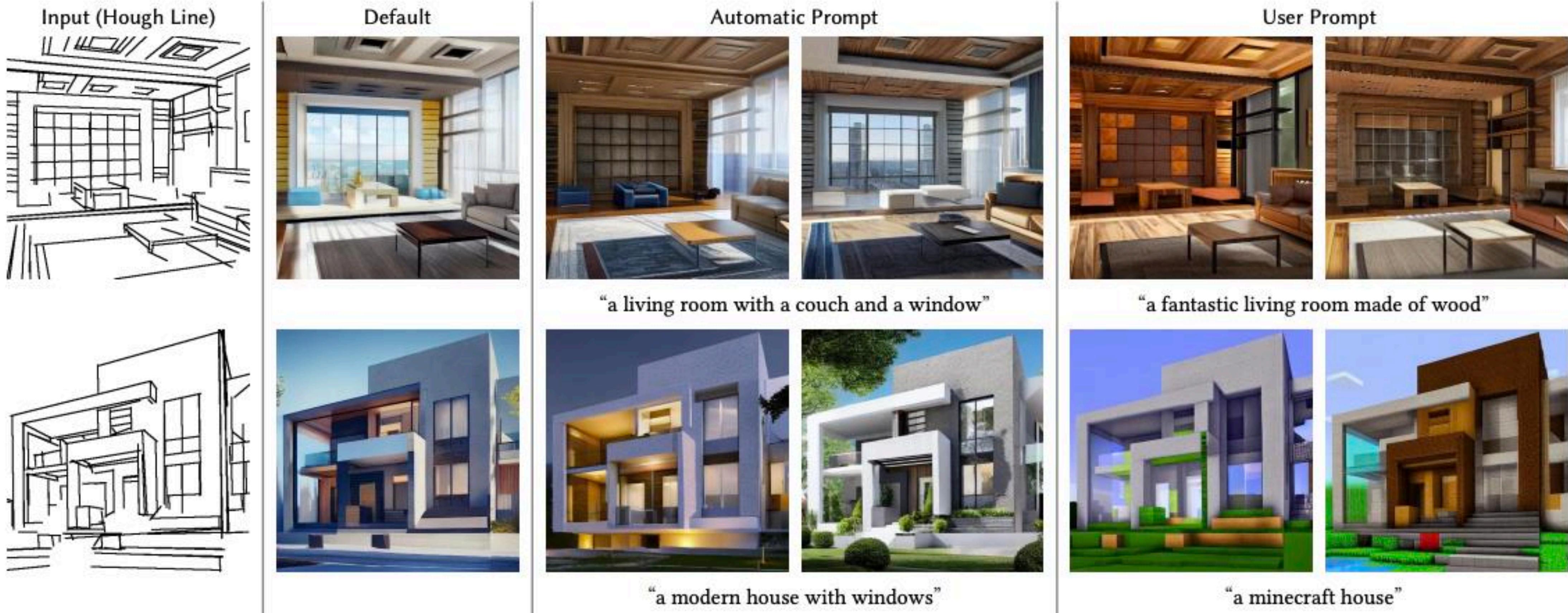
G. Textual Inversion: Personalization!

Can even learn an artistic style



ControlNet

Controlling Diffusion Model Generation
ICCV '23 Best Paper (Marr Prize)

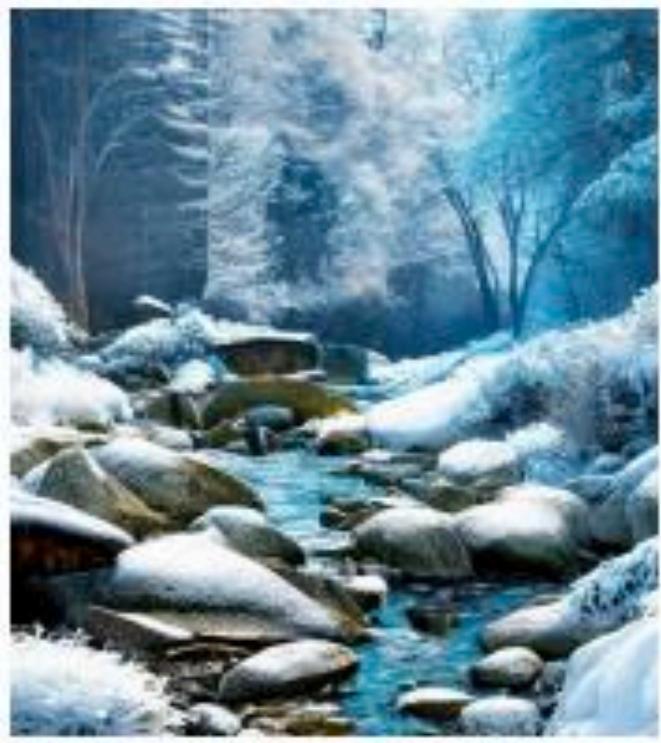
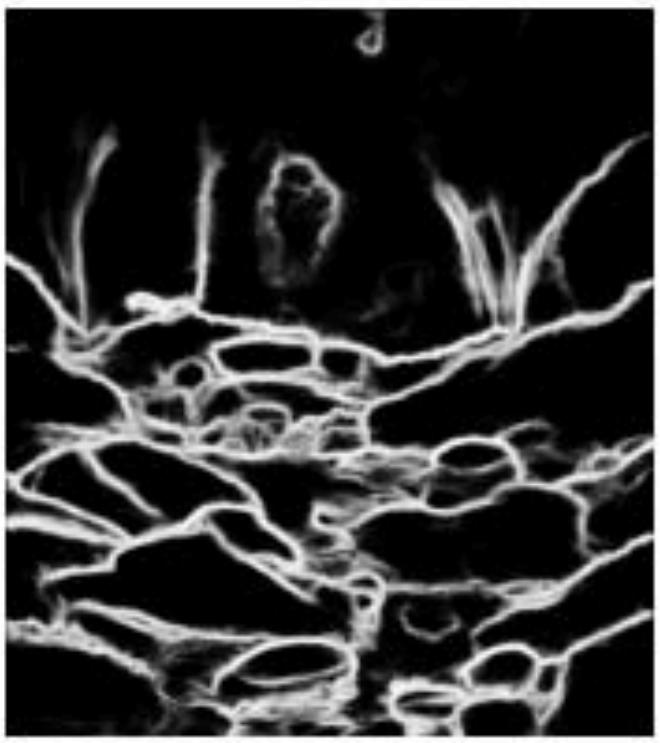


ControlNet

Controlling Diffusion Model Generation
ICCV '23 Best Paper (Marr Prize)



“a bird on a branch of a tree”



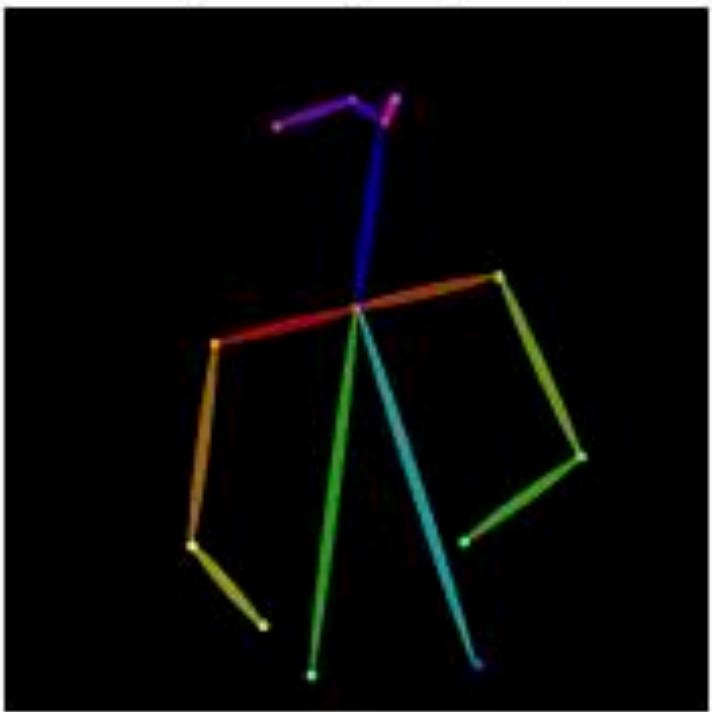
“ a stream running through a forest”

“river in forest, winter, snow”

ControlNet

Artistic Applications in Fashion and Design?

Input (openpose)



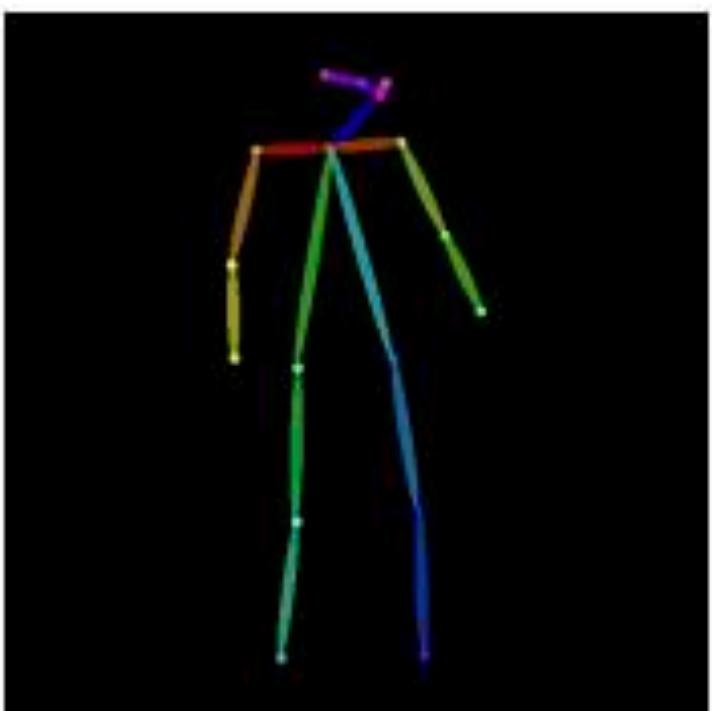
Default



User Prompt



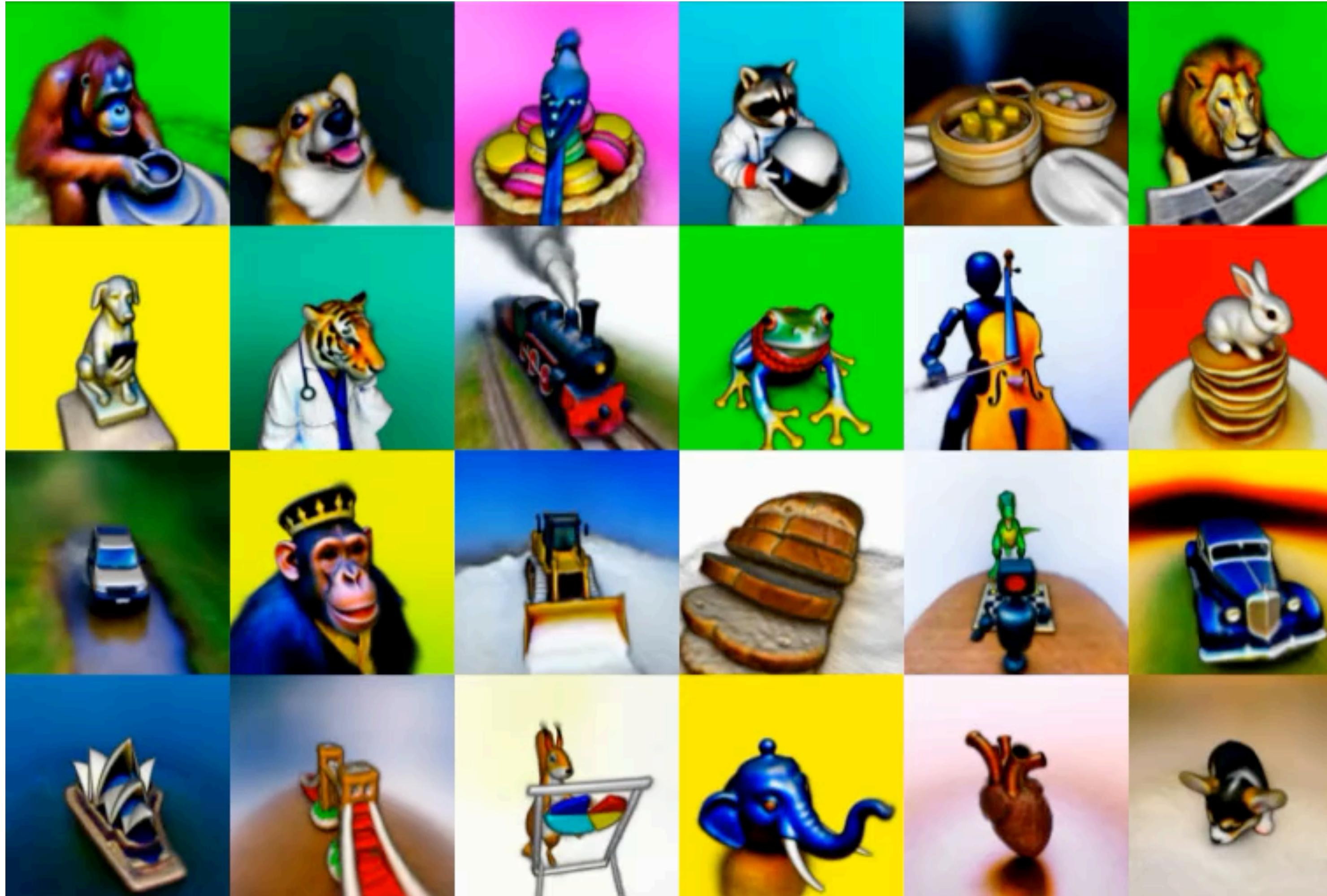
“chef in the kitchen”



“astronaut”

Vision Transcends Beyond Images

Generating 3D objects, videos, and more!



DreamFusion

OpenSource Version:
[https://github.com/
ashawkey/stable-
dreamfusion](https://github.com/ashawkey/stable-dreamfusion)

Vision Transcends Beyond Images

OpenSource Video Generation



Prompt: “*Darth Vader is surfing on waves.*”

Cerspense Models

OpenSource:
[https://huggingface.co/
cerspense/
zeroscope v2 576w](https://huggingface.co/cerspense/zeroscope_v2_576w)

Getting Started?

- Fast open-source diffusion models:
 1. <https://github.com/Stability-AI/StableCascade>
 2. <https://huggingface.co/ByteDance/SDXL-Lightning>
- Open-source video generation models:
 1. [https://huggingface.co/cerspense/zeroscope v2 576w](https://huggingface.co/cerspense/zeroscope_v2_576w)
 2. <https://makepixelsdance.github.io>
- Animation Models:
 1. <https://showlab.github.io/magicanimate/>

Getting Started? HF Diffusers Library!

- Video Generation with Cerscape model: example code

```
import torch  
  
from diffusers import DiffusionPipeline, DPMsolverMultistepScheduler  
from diffusers.utils import export_to_video  
  
pipe = DiffusionPipeline.from_pretrained("cerspense/zeroscope_v2_576w", torch_dtype=torch.float16)  
pipe.scheduler = DPMsolverMultistepScheduler.from_config(pipe.scheduler.config)  
pipe.enable_model_cpu_offload()  
  
prompt = "Darth Vader is surfing on waves"  
video_frames = pipe(prompt, num_inference_steps=40, height=320, width=576, num_frames=24).frames  
video_path = export_to_video(video_frames)
```