

CMSC 491/691

Lecture 17: Image Synthesis

Input Text

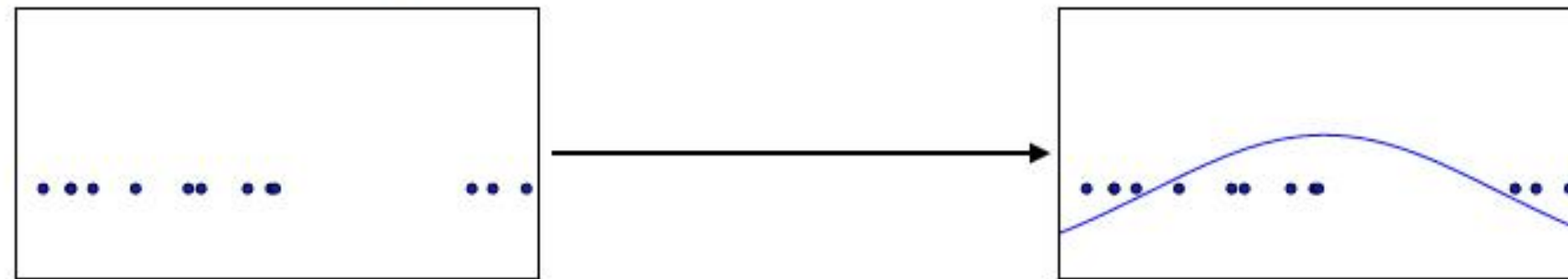
Struggling for meme ideas

Output Meme



Generative Modeling

- Density estimation



- Sample generation



Training examples

Model samples

Generative Adversarial Networks (GAN)

Why Generative Models?

- **We've only seen discriminative models so far**
 - Given an image \mathbf{X} , predict a label \mathbf{Y}
 - Estimates $\mathbf{P}(\mathbf{Y}|\mathbf{X})$
- **Discriminative models have several key limitations**
 - Can't model $\mathbf{P}(\mathbf{X})$, i.e. the probability of seeing a certain image
 - Thus, can't sample from $\mathbf{P}(\mathbf{X})$, i.e. **can't generate new images**
- **Generative models (in general) cope with all of above**
 - Can model $\mathbf{P}(\mathbf{X})$
 - Can generate new images

Generative Adversarial Networks

[Goodfellow et al., 2014]

Problem: Want to sample from complex, high-dimensional training distribution. There is no direct way to do this!

Solution: Sample from a simple distributions, e.g., random noise. Learn transformation to the training distribution

Question: What can we use to represent complex transformation function?

Output: Sample from training distribution

Input: Random noise

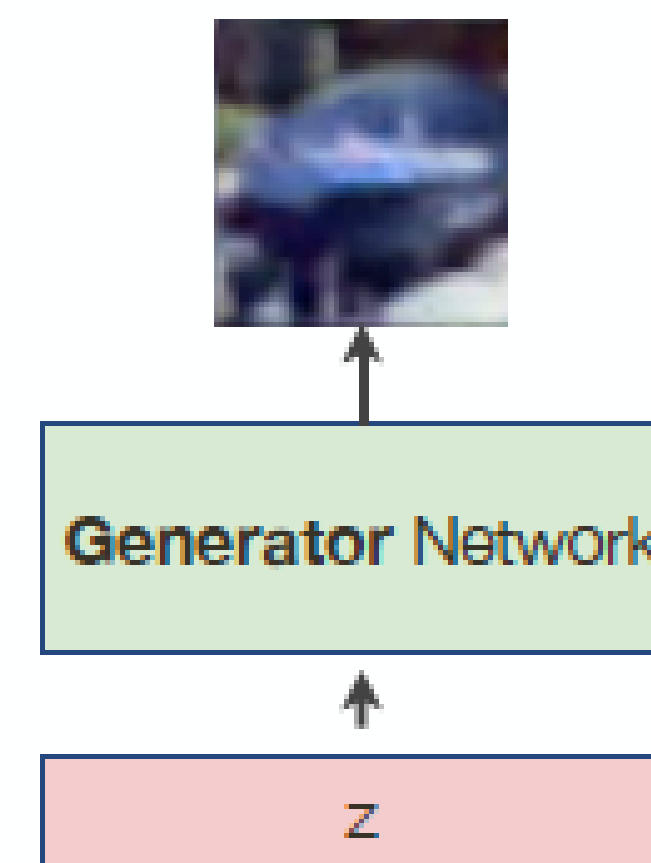
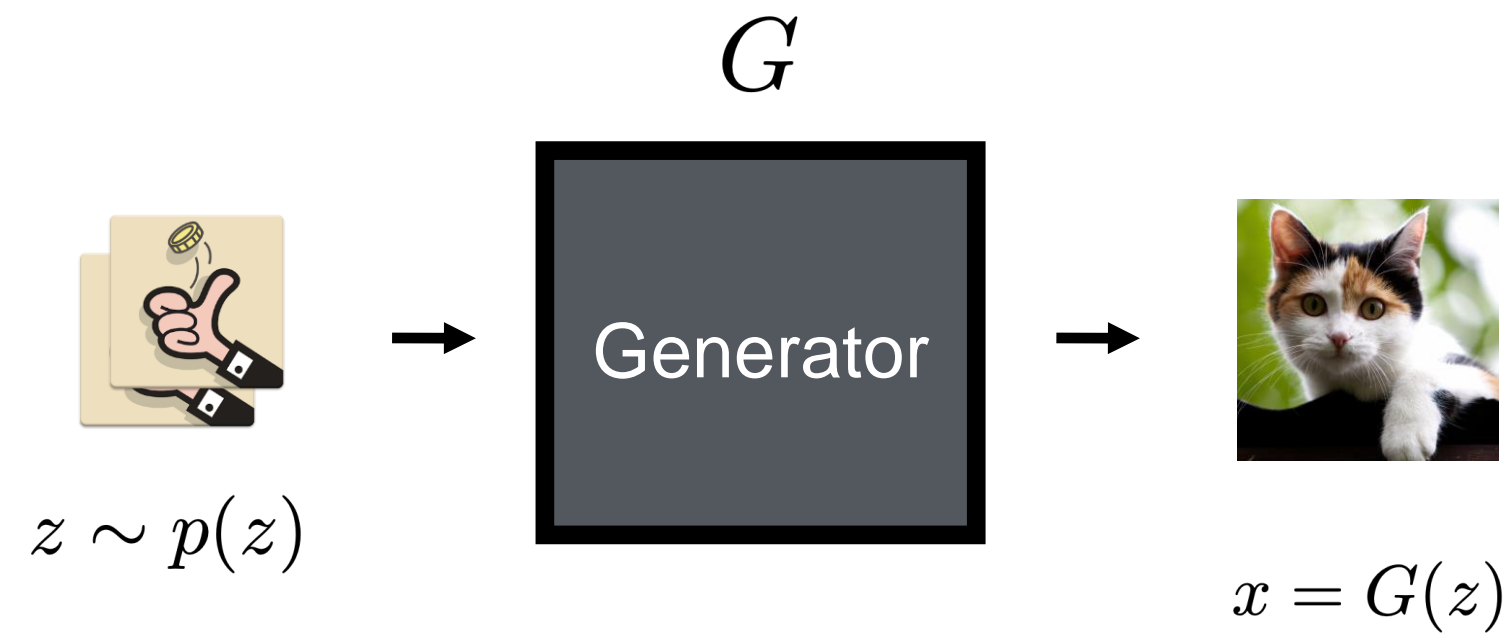


Image synthesis from “noise”



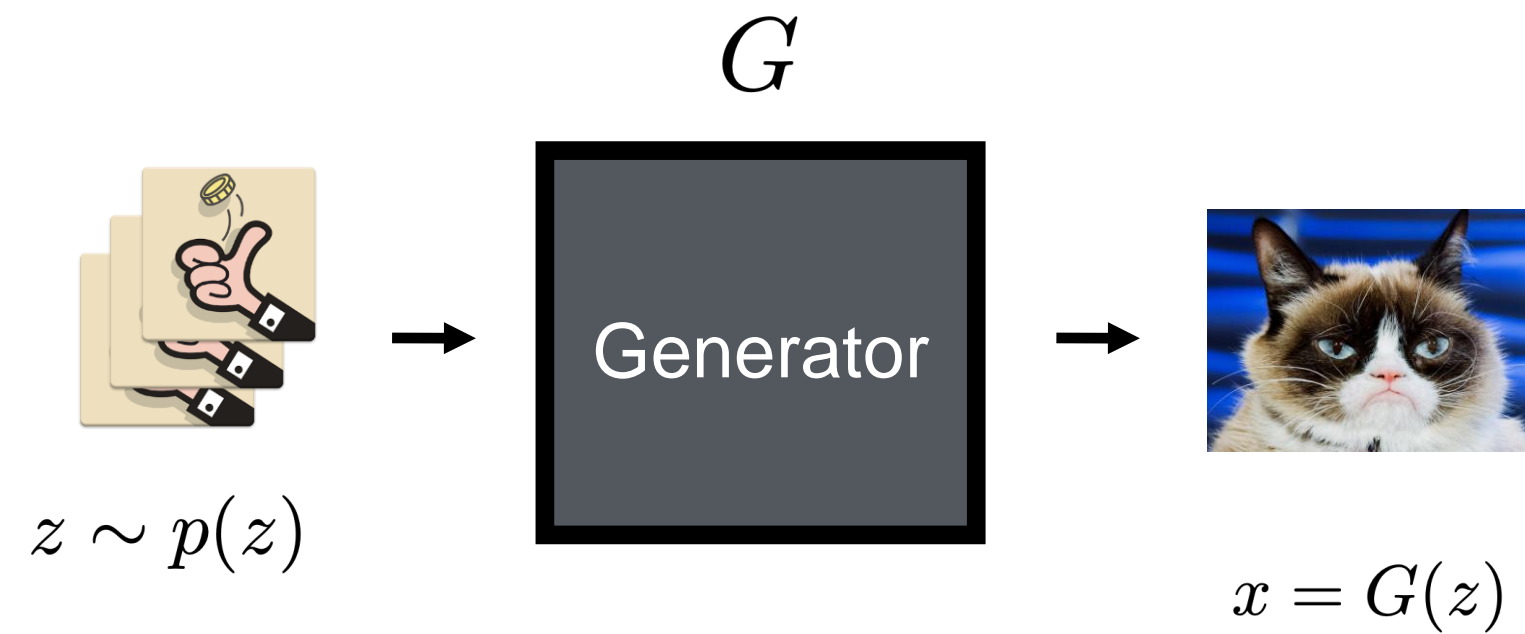
Sampler

$$G : \mathcal{Z} \rightarrow \mathcal{X}$$

$$z \sim p(z)$$

$$x = G(z)$$

Image synthesis from “noise”

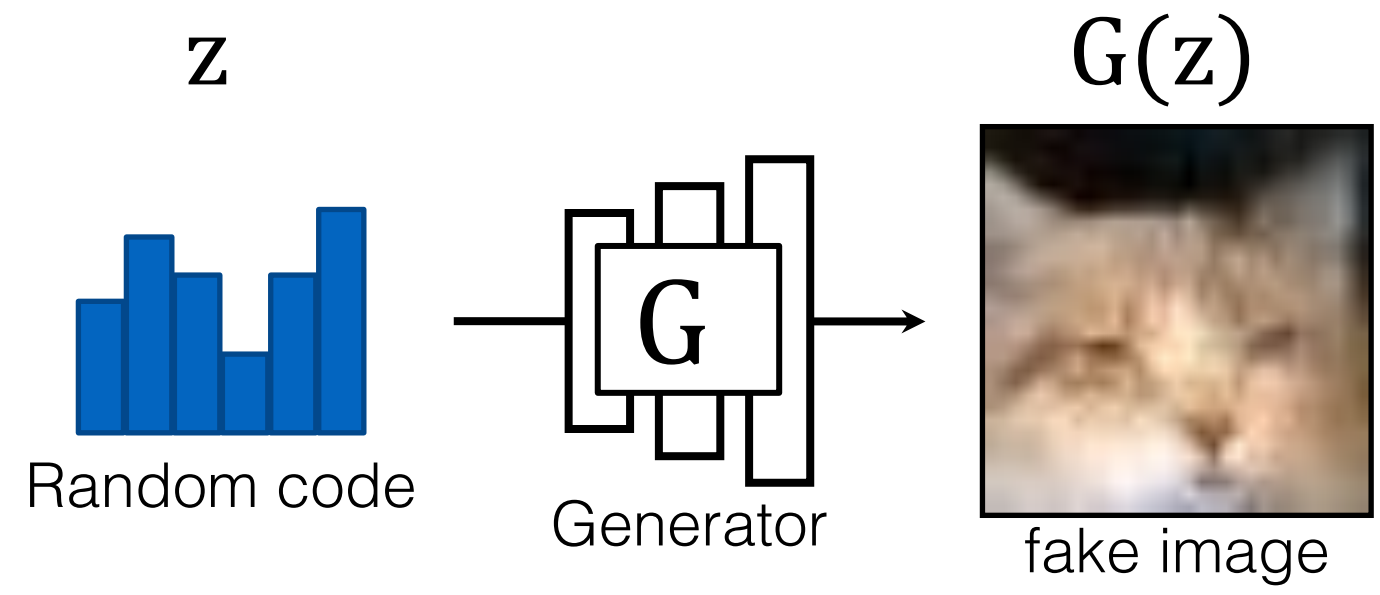


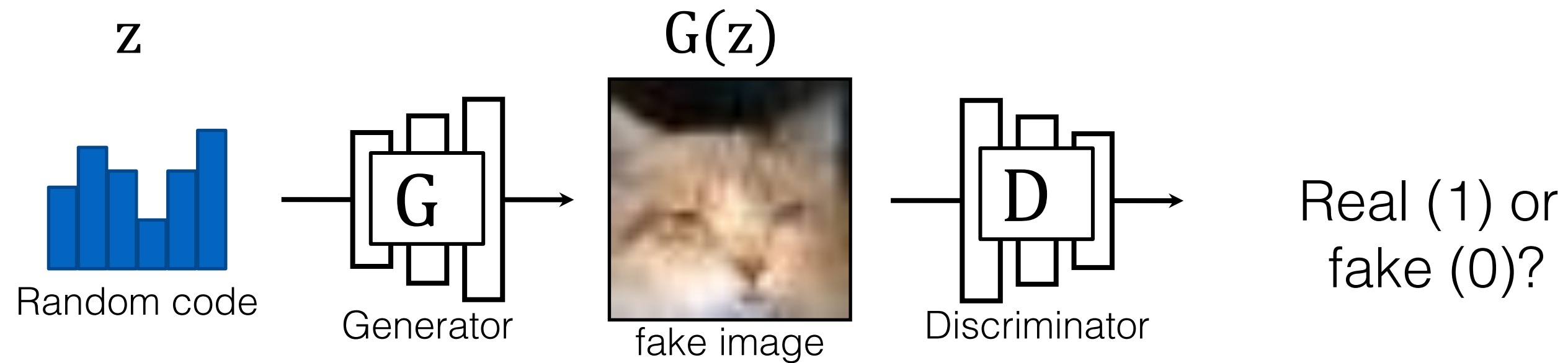
Sampler

$$G : \mathcal{Z} \rightarrow \mathcal{X}$$

$$z \sim p(z)$$

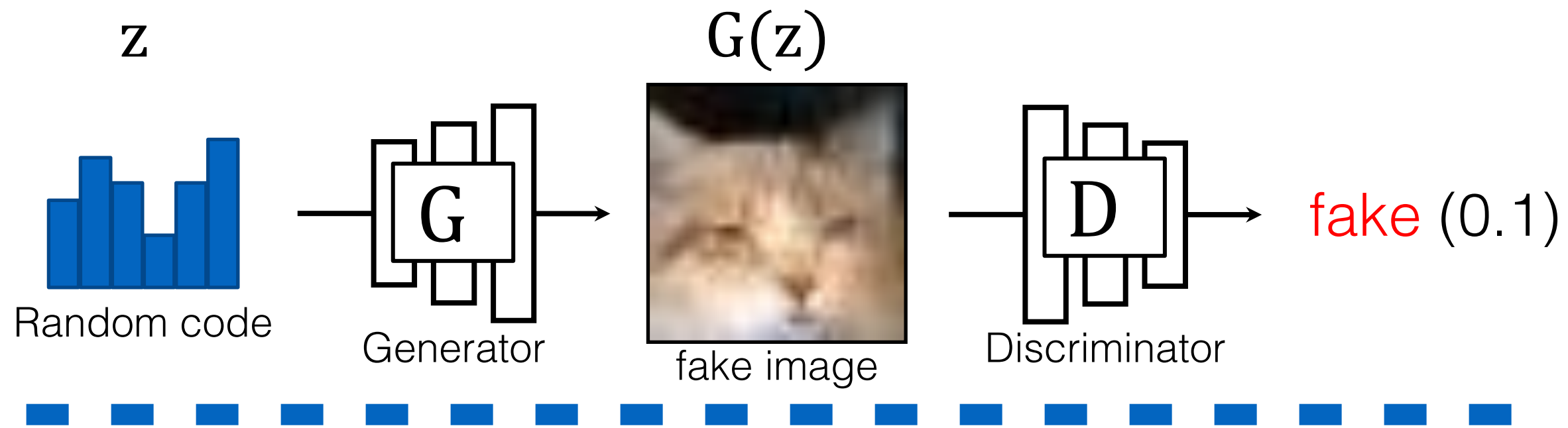
$$x = G(z)$$





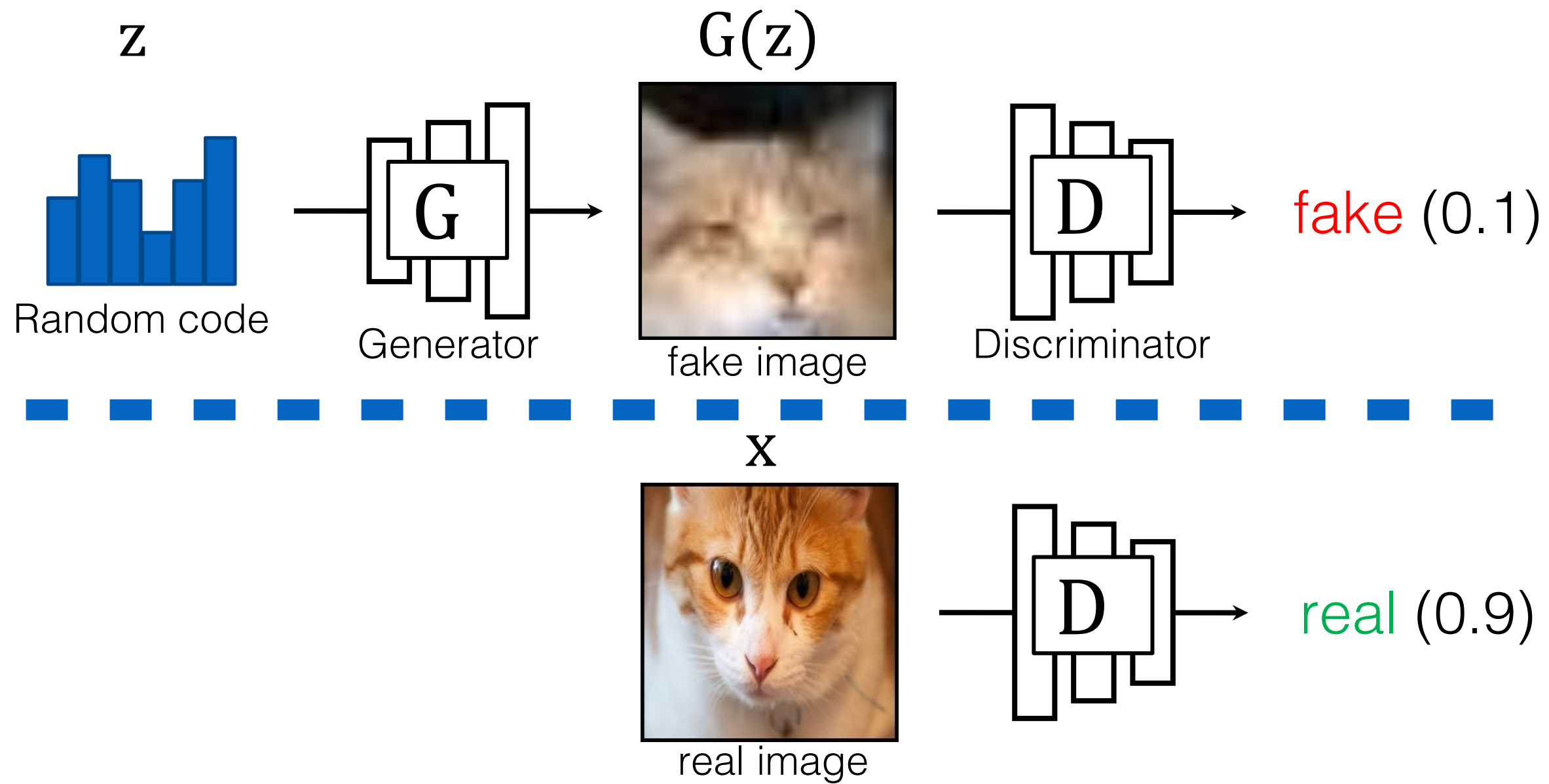
A two-player game:

- G tries to generate fake images that can fool D .
- D tries to detect fake images.



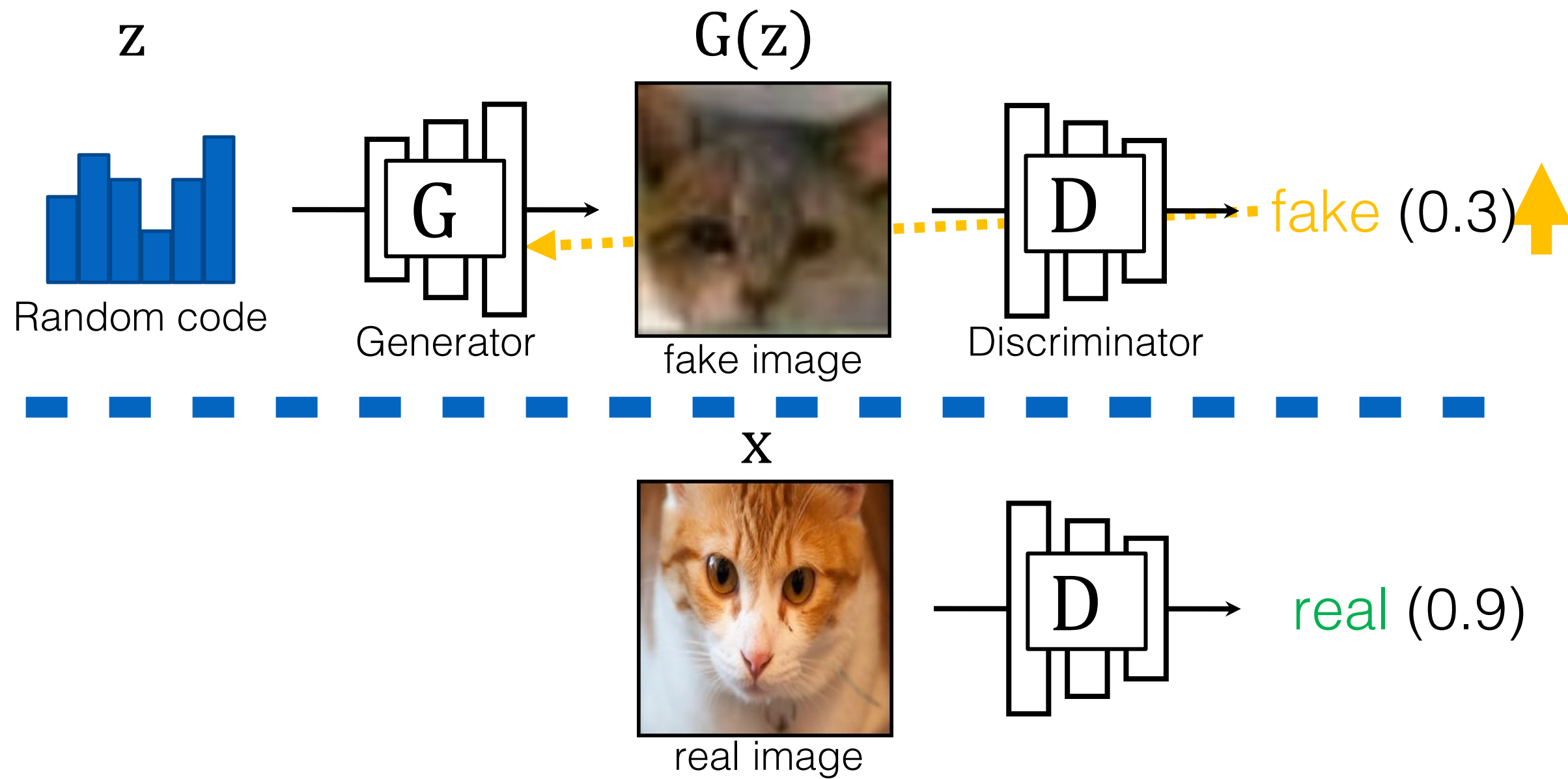
Learning objective (GANs)

$$\min_G \max_D \mathbb{E}_z [\log(1 - D(G(z)))]$$



Learning objective (GANs)

$$\min_G \max_D \mathbb{E}_z [\log(1 - D(G(z)))] + \mathbb{E}_x [\log D(x)]$$



Learning objective (GANs)

$$\min_G \max_D \mathbb{E}_z [\log(1 - D(G(z)))] + \mathbb{E}_x [\log D(x)]$$

GAN Training Breakdown

- From the discriminator D 's perspective:
 - binary classification: real vs. fake.
 - Nothing special: similar to 1 vs. 7 or cat vs. dog

$$\max_D \mathbb{E}[\log(1 - D(\text{cat}))] + \mathbb{E}[\log D(\text{cat})]$$

GAN Training Breakdown

- From the discriminator D 's perspective:
 - binary classification: real vs. fake.
 - Nothing special: similar to 1 vs. 7 or cat vs. dog

$$\max_D \mathbb{E}[\log(1 - D(\text{dog}))] + \mathbb{E}[\log D(\text{cat})]$$

- From the generator G 's perspective:
 - Optimizing a loss that depends on a classifier D
 - We have done it before (Perceptual Loss)

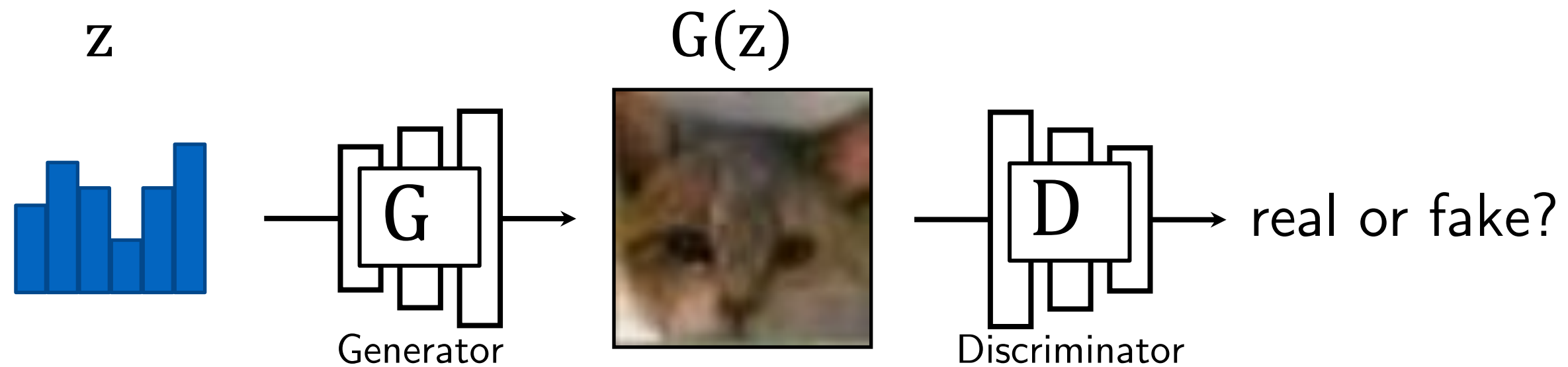
$$\min_G \mathbb{E}_z[\mathcal{L}_D(G(z))]$$

GAN loss for G

$$\min_G \mathbb{E}_{(x,y)} \|F(G(x)) - F(y)\|$$

Perceptual Loss for G

GAN Training Breakdown

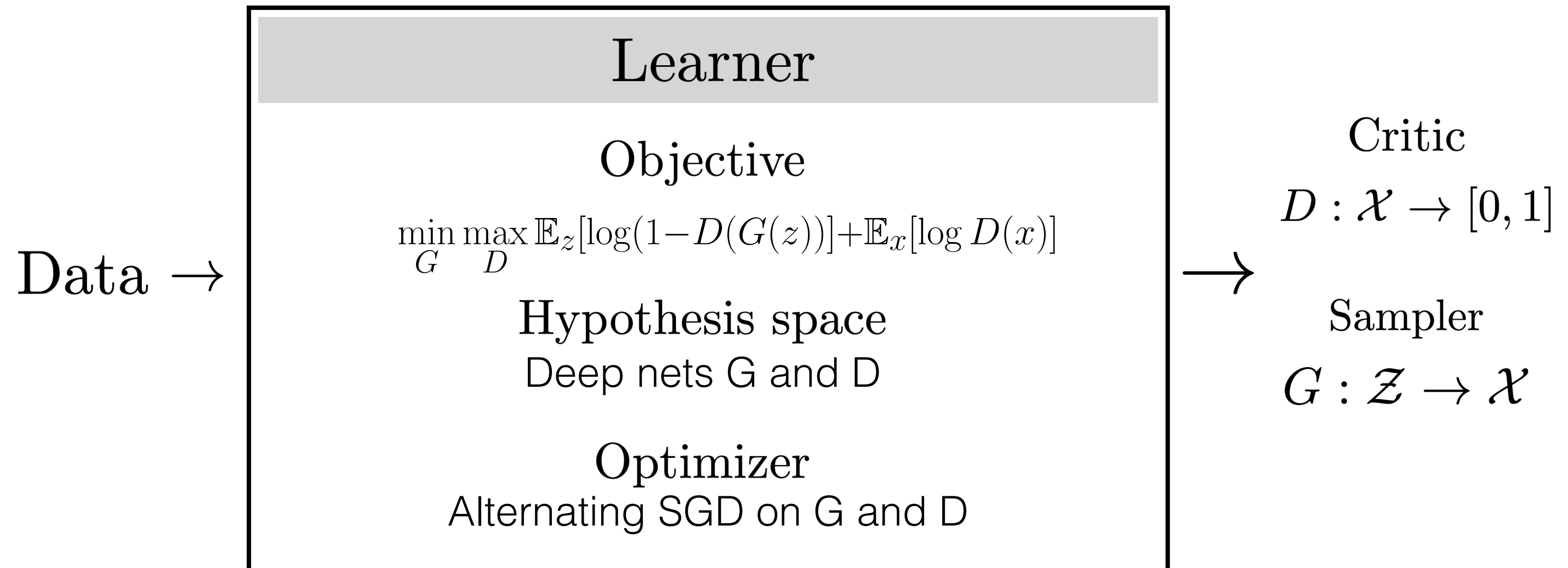


G tries to synthesize fake images that fool D

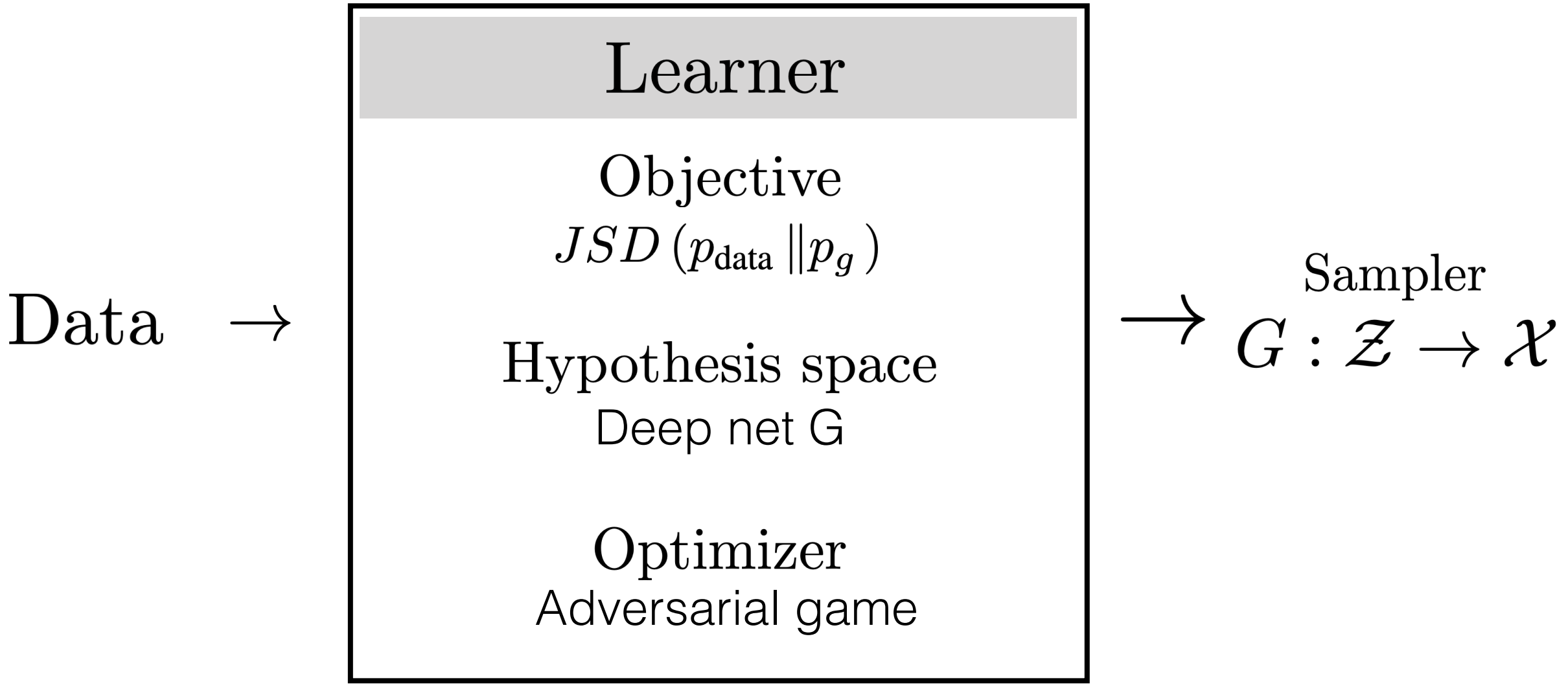
D tries to identify the fakes

- Training: iterate between training D and G with backprop.
- Global optimum when G reproduces data distribution.

Generative Adversarial Network



Generative Adversarial Network



Generative Adversarial Nets

Generated Samples





Ian Goodfellow @goodfellow_ian · Jan 14



4.5 years of **GAN progress** on face generation. arxiv.org/abs/1406.2661

arxiv.org/abs/1511.06434 arxiv.org/abs/1606.07536 arxiv.org/abs/1710.10196

arxiv.org/abs/1812.04948



Samples from **StyleGAN2** [Karras et al., CVPR 2020]

Interpolation is impressive



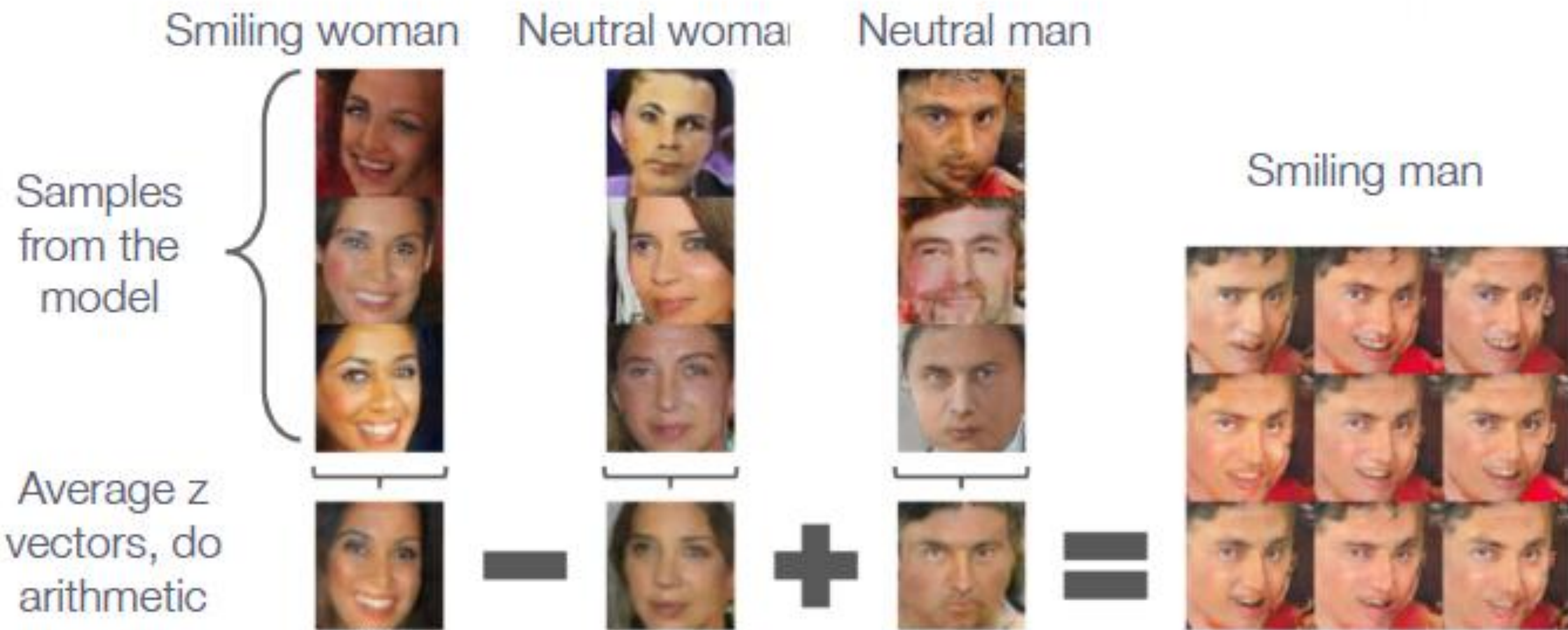
(c) Our results (128x128 with 128 filters)



(d) Mirror interpolations (our results 128x128 with 128 filters)

GANs: Interpretable Vector Math

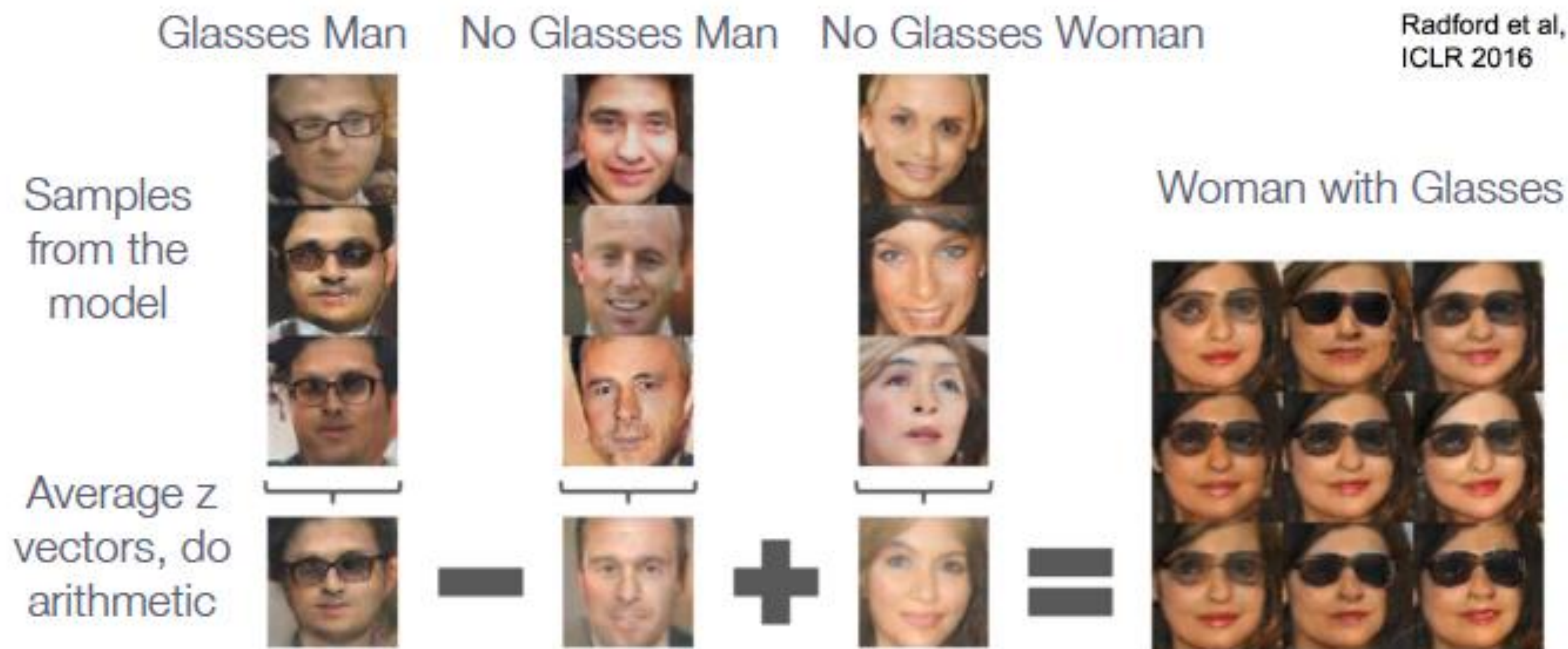
[Radford et al., 2016]



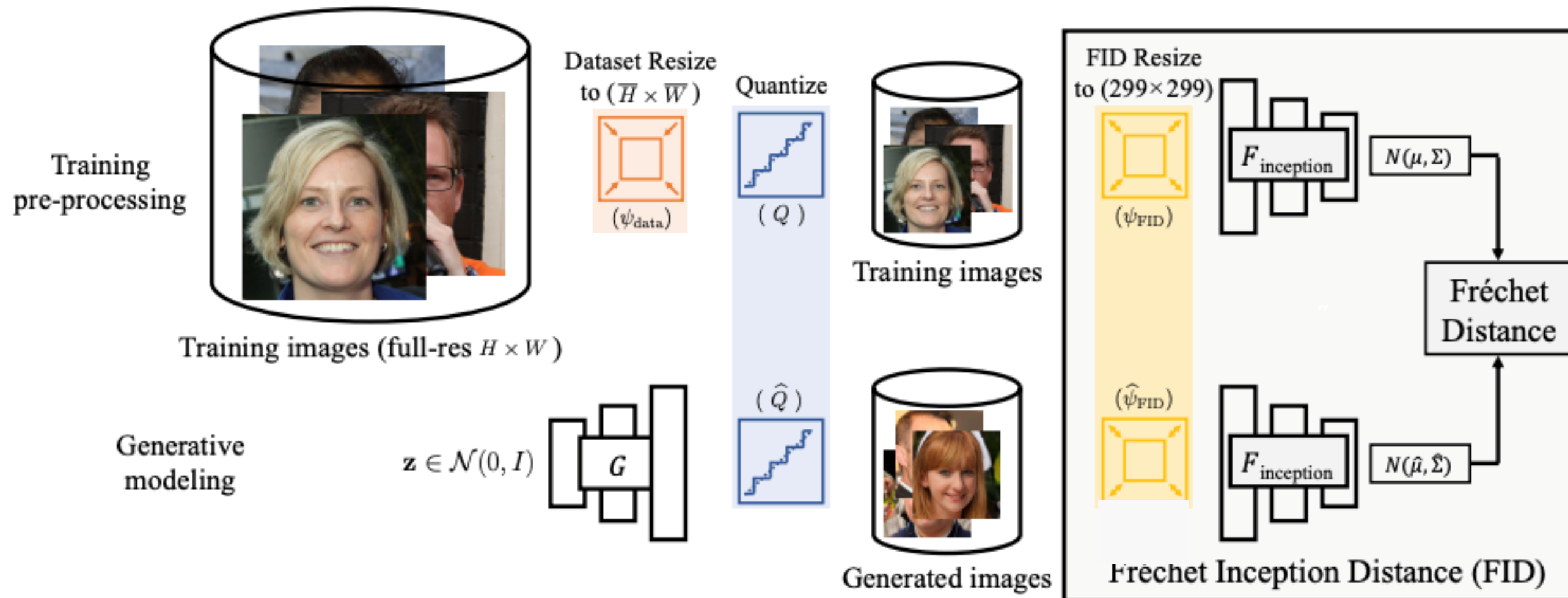
GANs: Interpretable Vector Math

[Radford et al., 2016]

Radford et al,
ICLR 2016



GANs evaluation (FID)



Fréchet Inception Distance (FID)

$$\mathbf{FID} = \|\mu - \hat{\mu}\|_2^2 + \text{Tr}(\Sigma + \hat{\Sigma} - 2(\Sigma\hat{\Sigma})^{1/2})$$

GANs evaluation (FID)

Clean-fid libraries for evaluating generative models

```
Python 3.7.10 (default, Feb 26 2021, 18:47:35)
[GCC 7.3.0] :: Anaconda, Inc. on linux
Type "help", "copyright", "credits" or "license" for more information.
>>> █
```

pip install clean-fid

Daily downloads (July, 2022): 100

Daily downloads (Feb, 2023) : 20, 000

Total downloads: 2, 600, 000

Better training and generation



(a) Church outdoor.

(b) Dining room.



(c) Kitchens.

(d) Conference room.

LSGAN. Mao et al. 2017.



BEGAN. Bertholet et al. 2017.

Source->Target domain transfer



horse → zebra

zebra → horse

apple → orange

→ summer Yosemite

→ winter Yosemite

CycleGAN. Zhu et al. 2017.

Text -> Image Synthesis

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.



Reed et al. 2017.

Many GAN applications



Pix2pix. Isola 2017. Many examples at <https://phillipi.github.io/pix2pix/>

Image-to-Image Translation

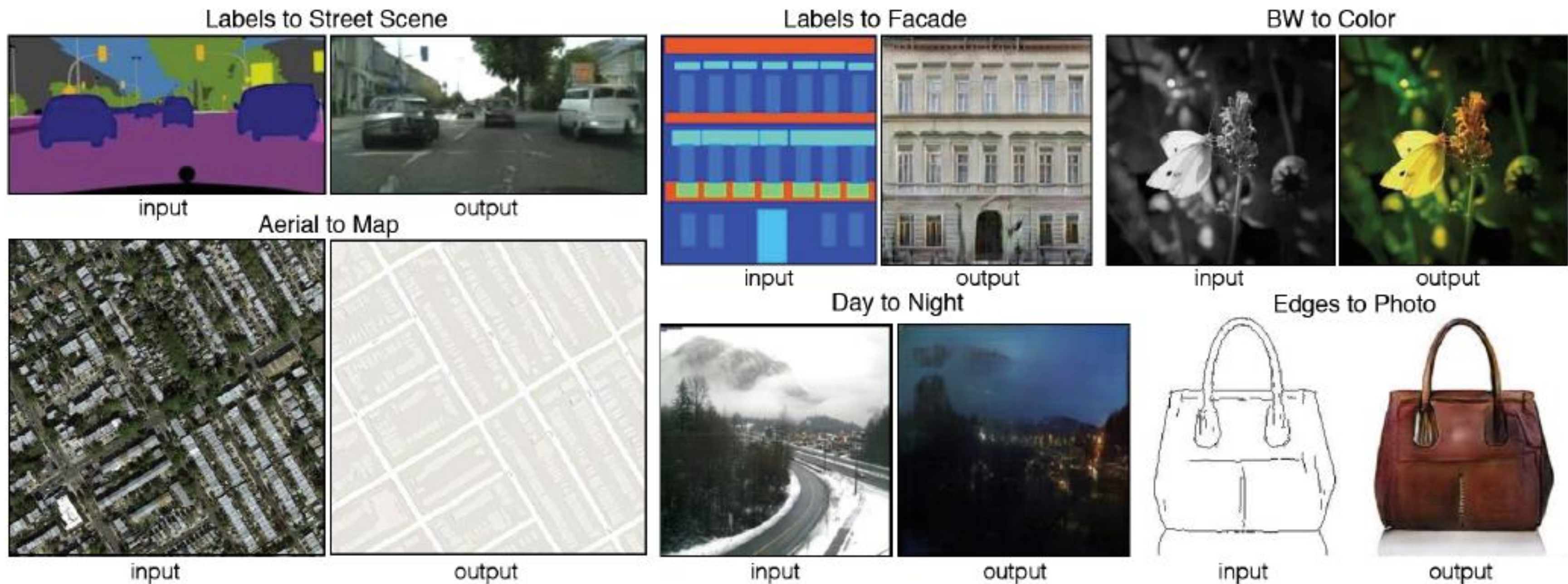


Figure 1 in the original paper.

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

Image-to-Image Translation

- Architecture: *DCGAN*-based architecture
- Training is conditioned on the images from the source domain.
- Conditional GANs provide an effective way to handle many complex domains without worrying about designing *structured loss* functions explicitly.

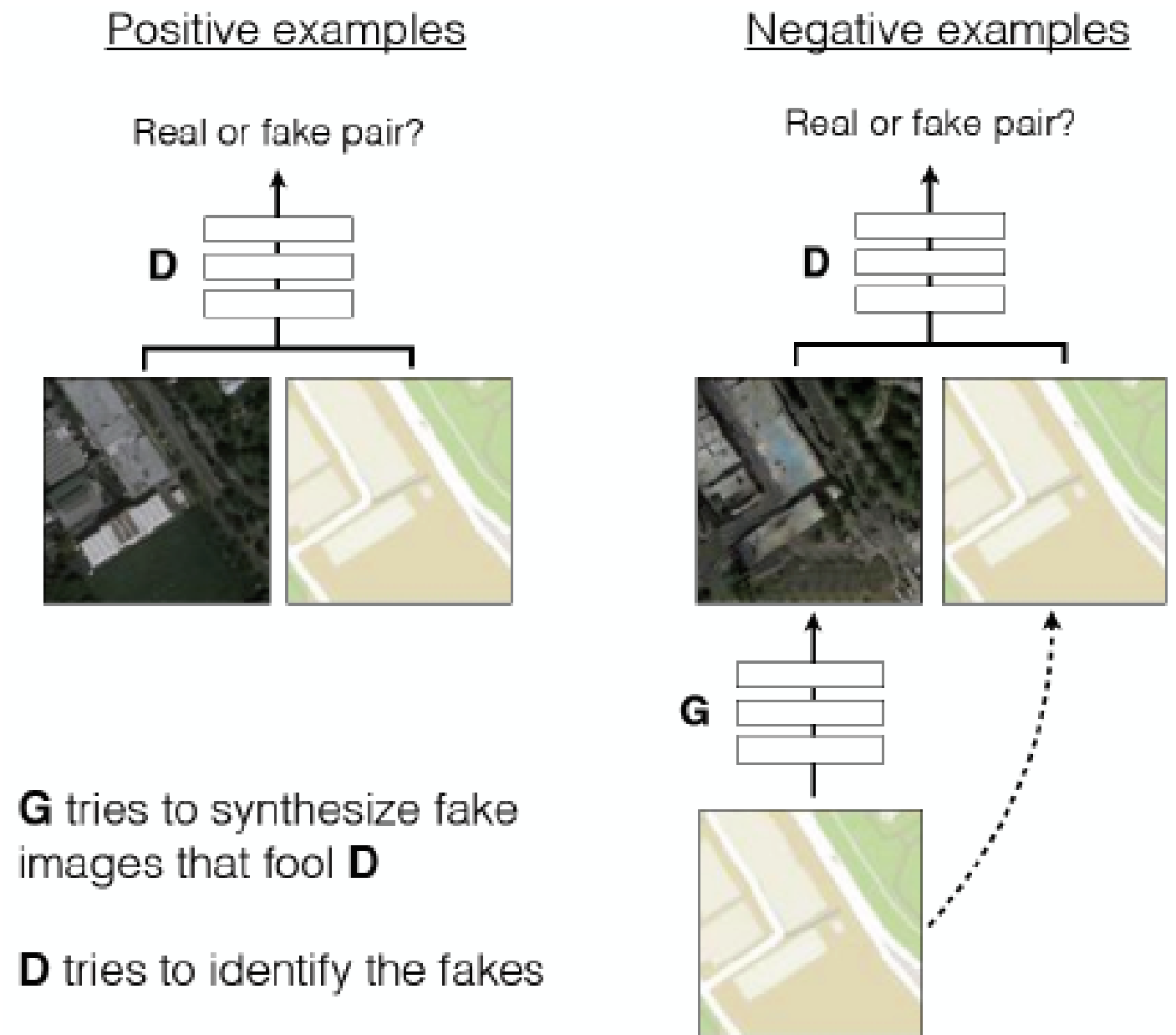


Figure 2 in the original paper.

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

- **Deep Learning models (in general) involve a single player**
 - The player tries to maximize its reward (minimize its loss).
 - Use SGD (with Backpropagation) to find the optimal parameters.
 - SGD has convergence guarantees (under certain conditions).
 - **Problem:** With non-convexity, we might converge to local optima.

$$\min_G L(G)$$

- **GANs instead involve two (or more) players**
 - Discriminator is trying to maximize its reward.
 - Generator is trying to minimize Discriminator's reward.

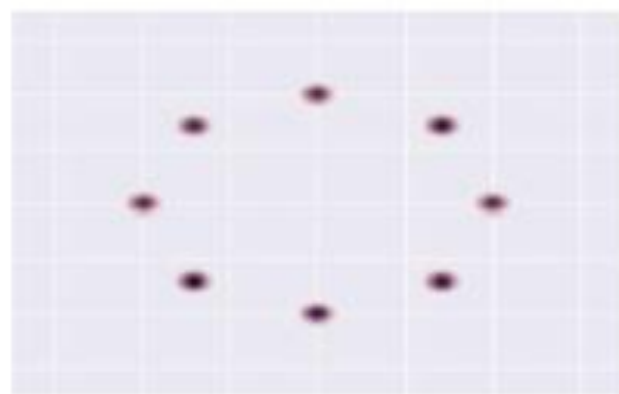
$$\min_G \max_D V(D, G)$$

- SGD was not designed to find the Nash equilibrium of a game.
- **Problem:** We might not converge to the Nash equilibrium at all.

Mode-Collapse

- Generator fails to output diverse samples

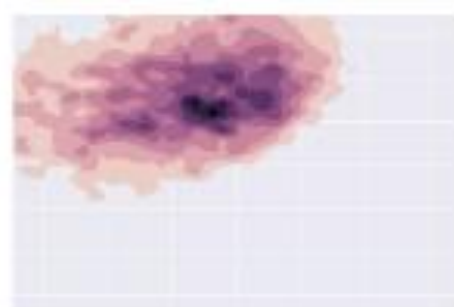
Target



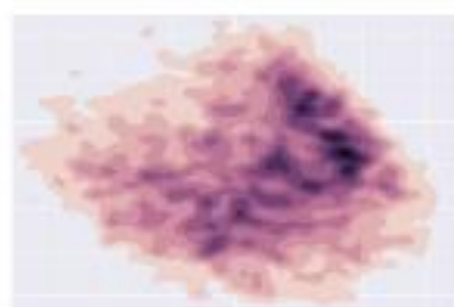
Expected



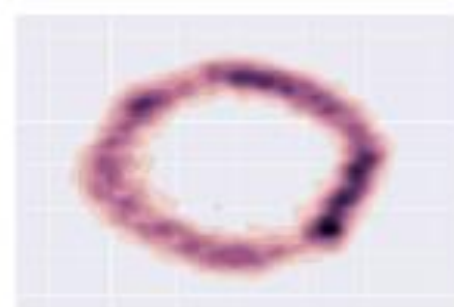
Step 0



Step 5k



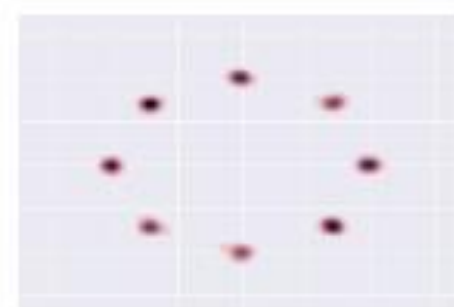
Step 10k



Step 15k

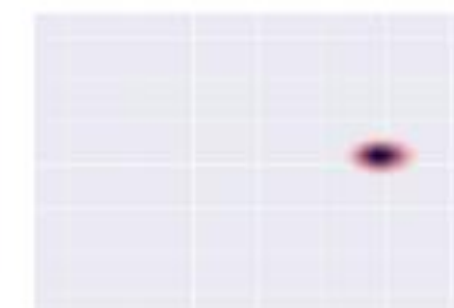
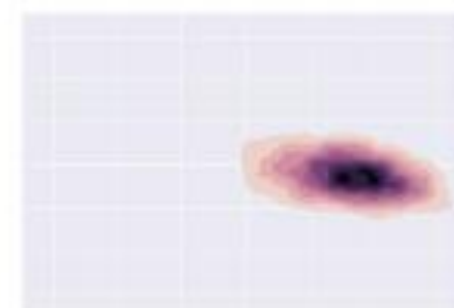


Step 20k



Step 25k

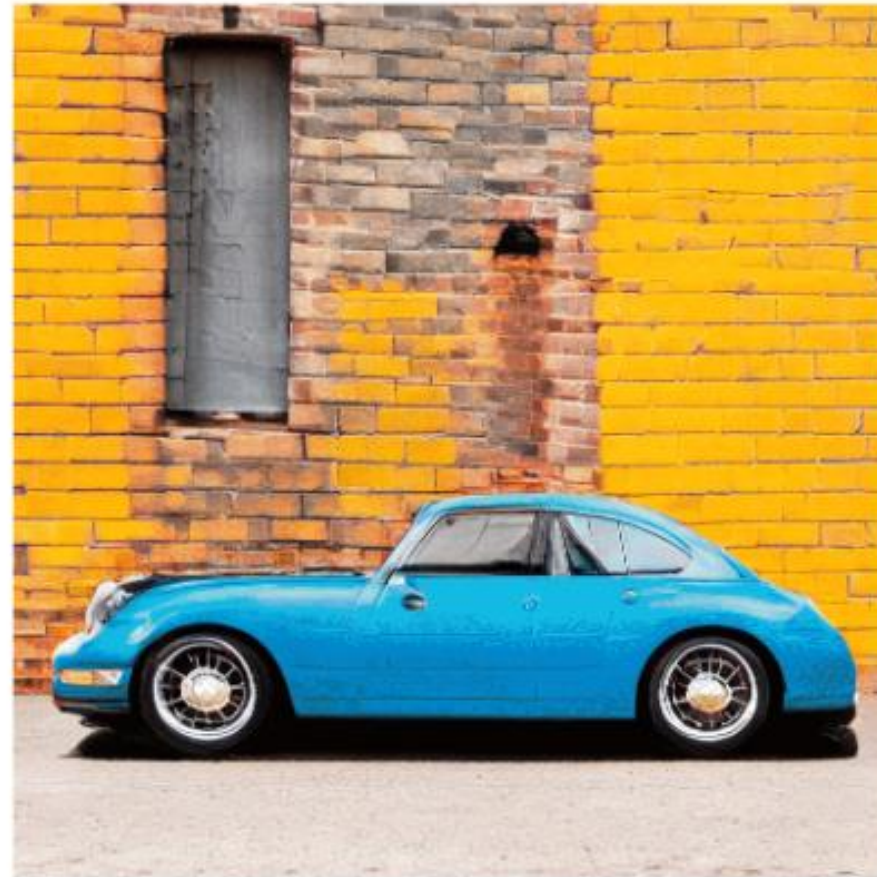
Output



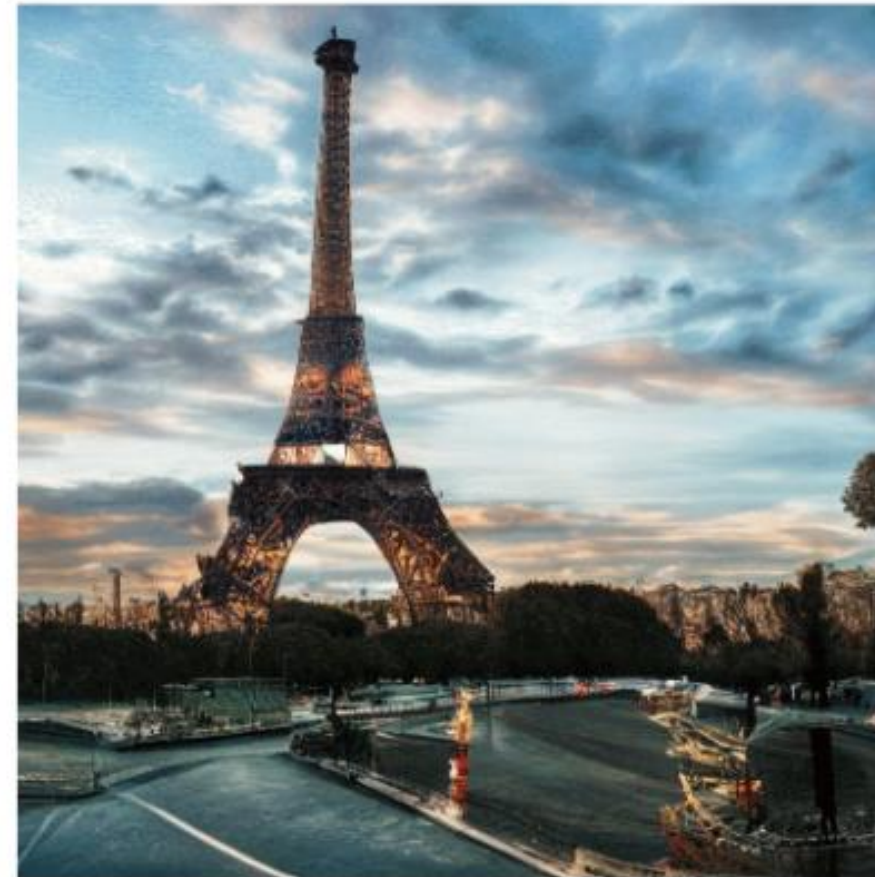
Text-to-Image (T2I)



A living room with a fireplace at a wood cabin. Interior design.



a blue Porsche 356 parked in front of a yellow brick wall.

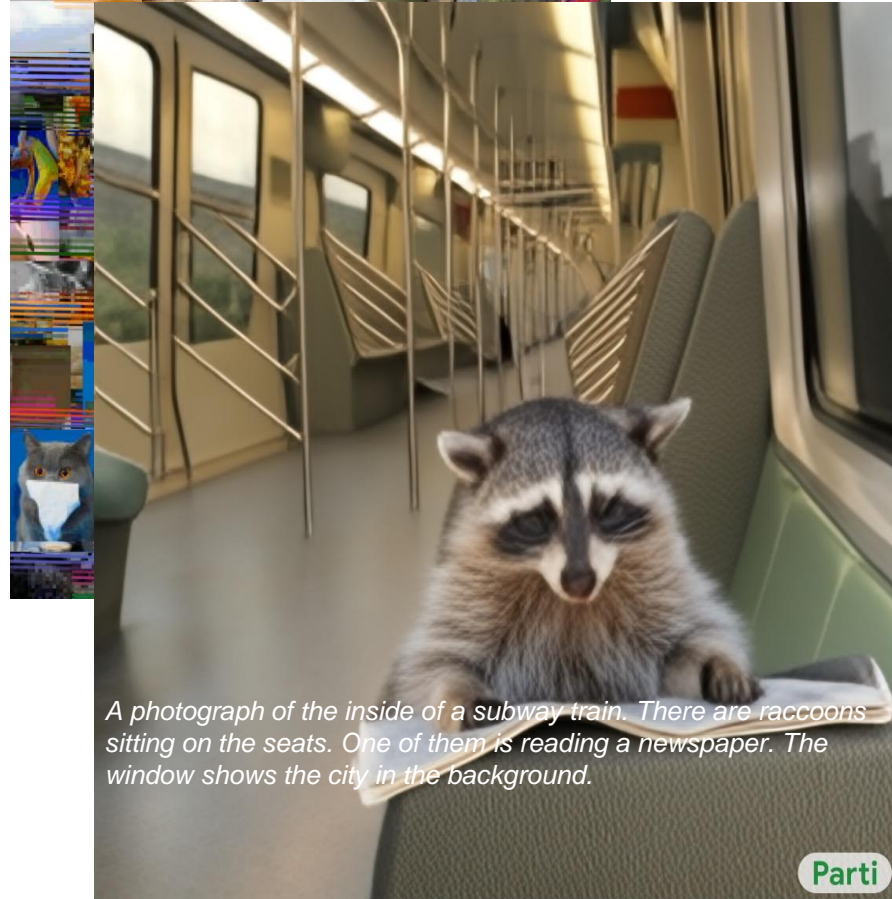
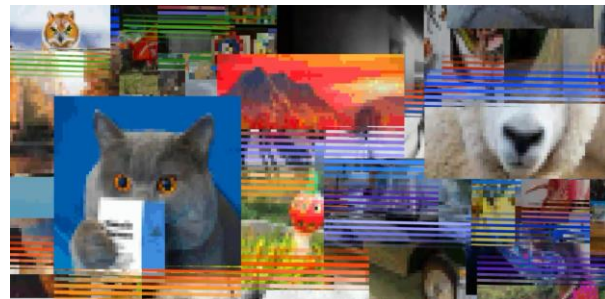


Eiffel Tower, landscape photography

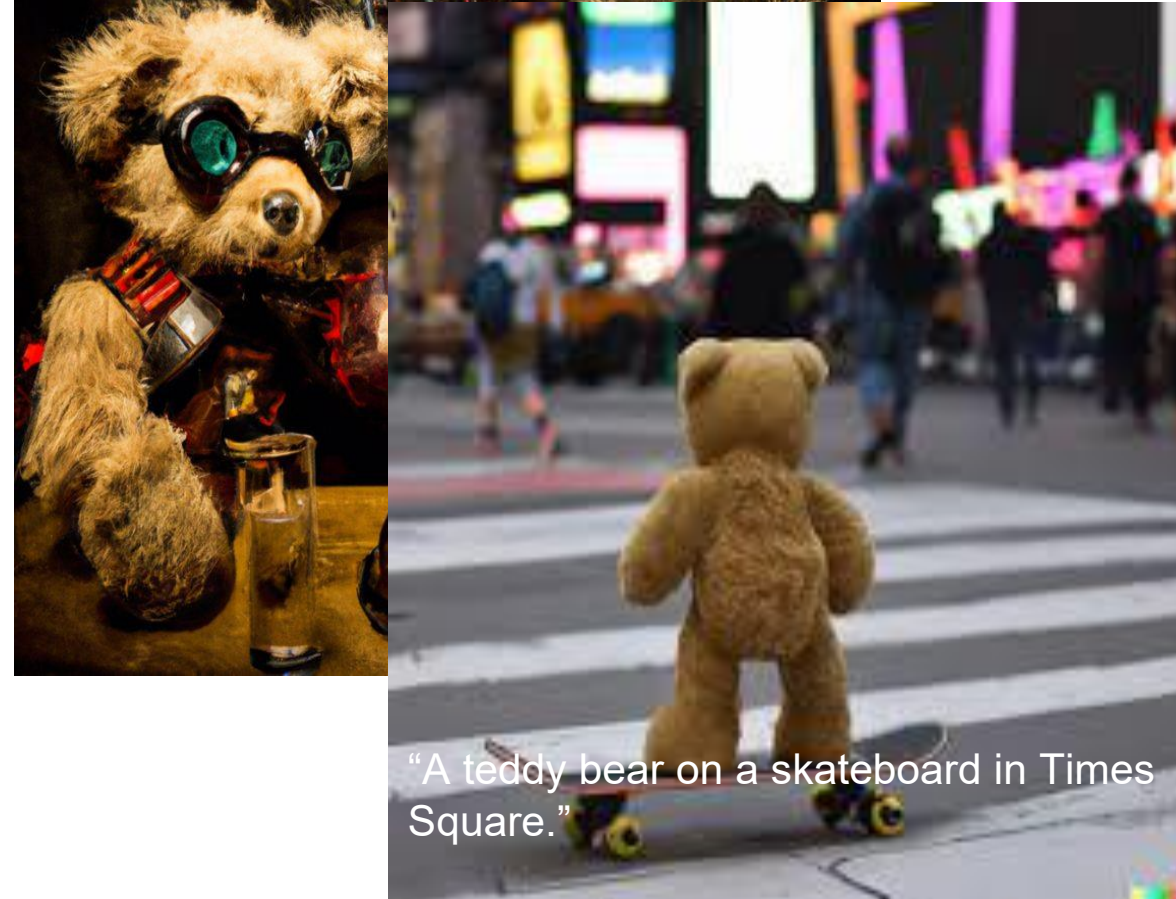
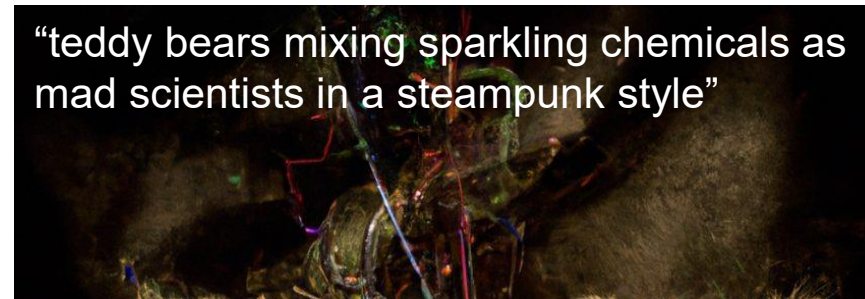


A painting of a majestic royal tall ship in Age of Discovery.

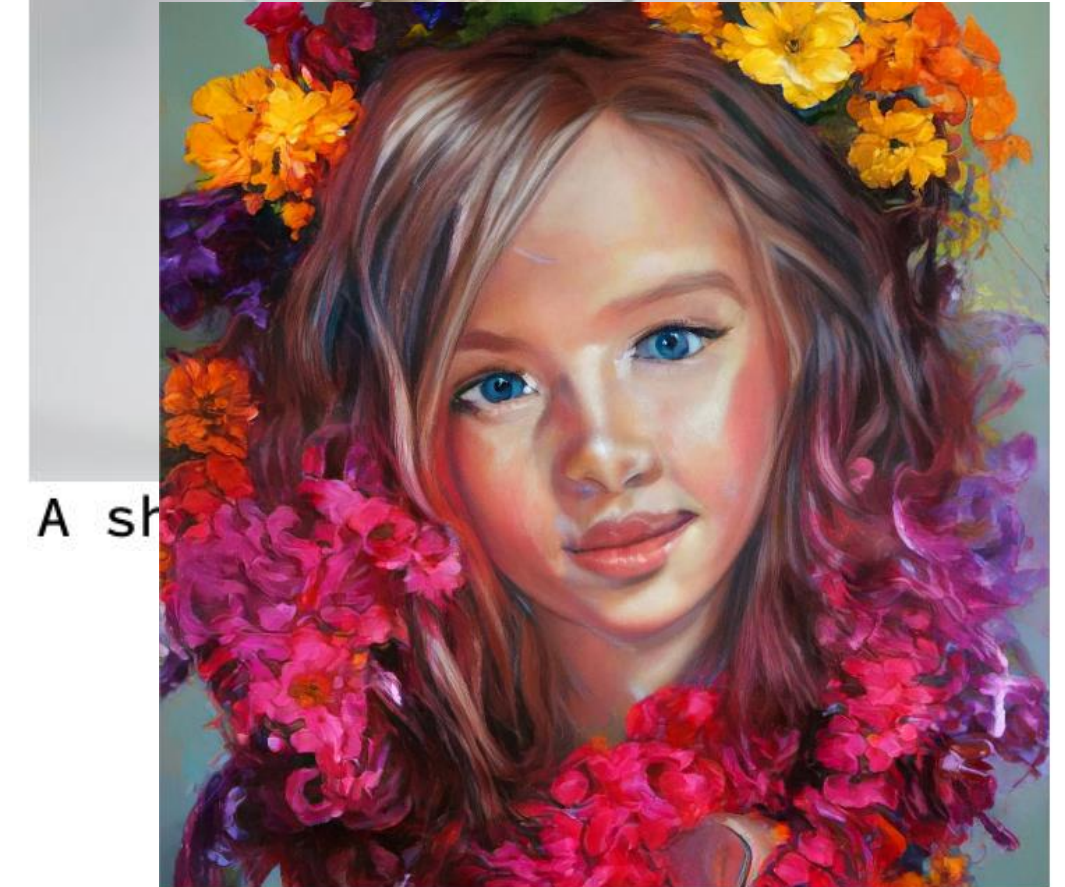
Text-to-Image Everywhere



Autoregressive models
(Image GPT, Parti)

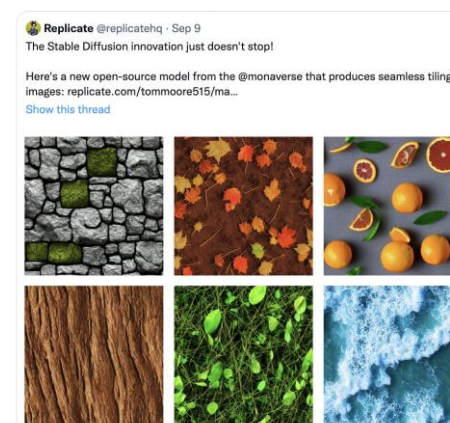
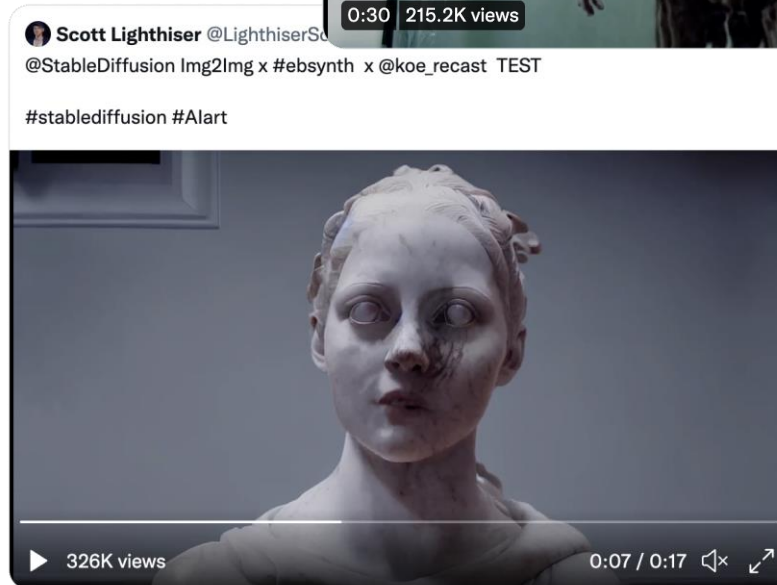


Diffusion models
(DALL-E 2, Imagen)



GANs, Masked GIT
(GigaGAN, MUSE)

Text-to-Image Everywhere



Where/when did it start?

First Text-to-Image System

First the
farmer gives
hay to the
goat. Then
the farmer
gets milk
from the
cow.



Step 1: Image Selection.

Step 2: Layout Optimization (Minimum overlap, Centrality, Closeness)

A Text-to-Picture Synthesis System for Augmenting Communication

Xiaojin Zhu, Andrew Goldberg, Mohamed Eldawy, Charles Dyer, and Bradley Strock. AAAI 2007

First Text-to-Image System



Therapy for people
with communicative disorders

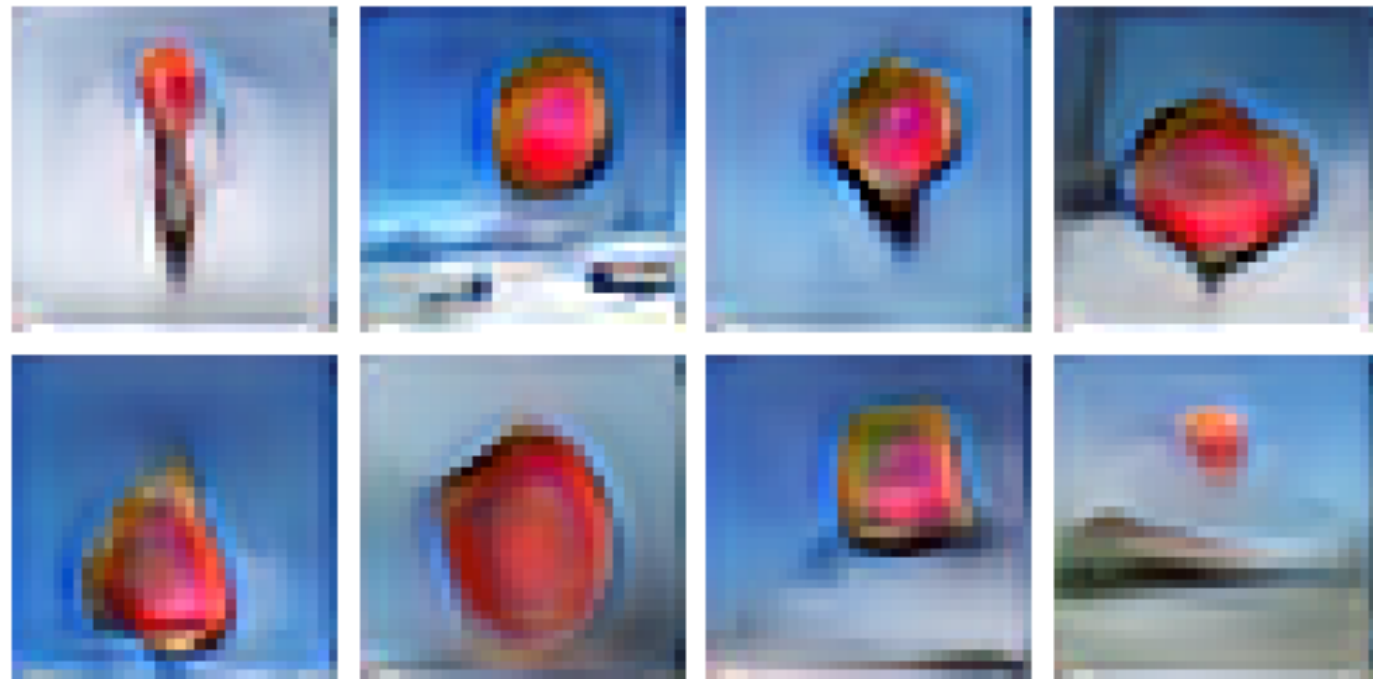


Math learning and reading comprehension
for young children

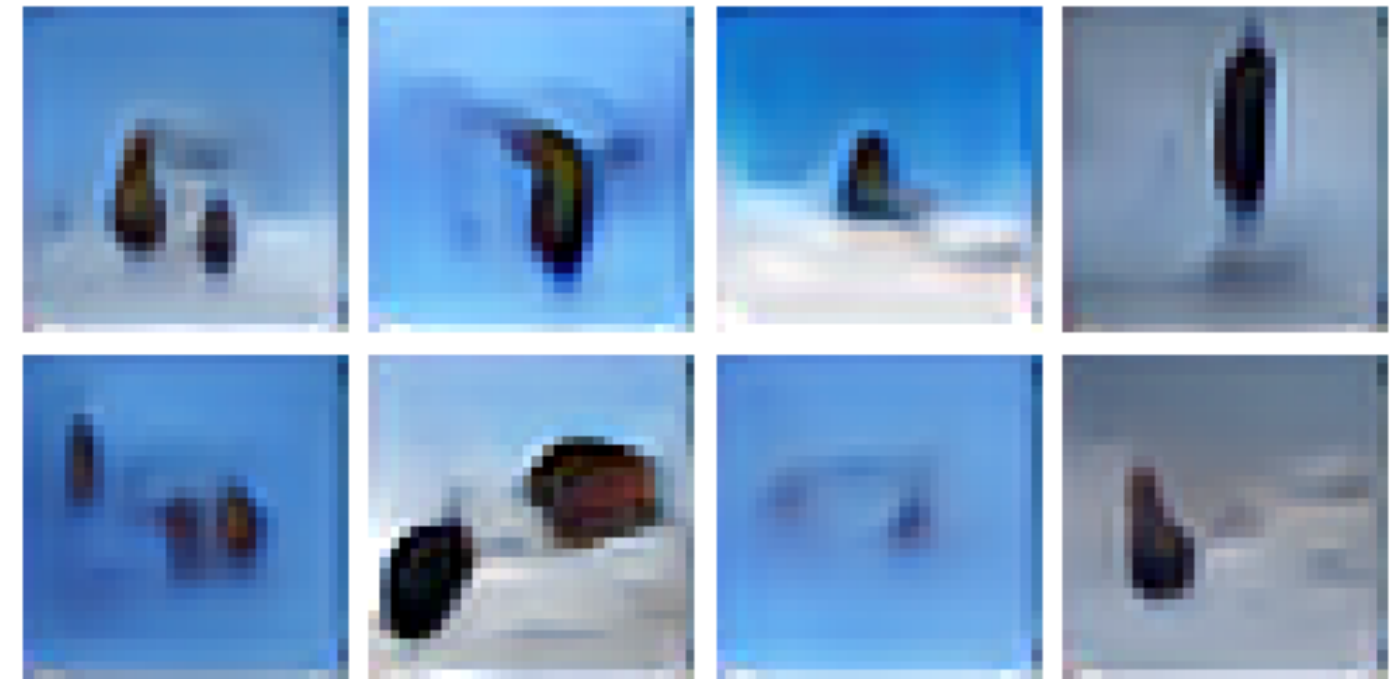
A Text-to-Picture Synthesis System for Augmenting Communication

Xiaojin Zhu, Andrew Goldberg, Mohamed Eldawy, Charles Dyer, and Bradley Stroock. AAAI 2007

First Deep Learning Work



A stop sign is flying in blue skies.



A herd of elephants flying in the blue skies.

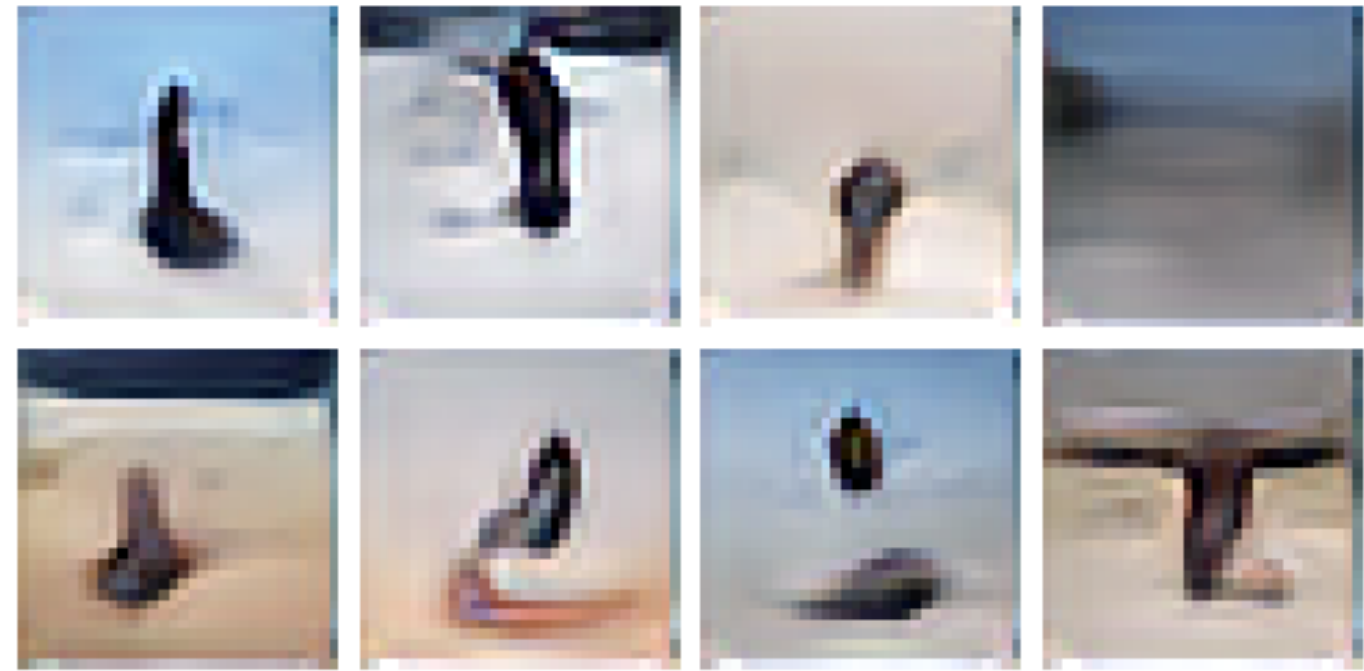
Generating Images from Captions with Attention.

Elman Mansimov, Emilio Parisotto, Jimmy Lei Ba, Ruslan Salakhutdinov. ICLR 2016.

First Deep Learning Work



A toilet seat sits open in the grass field.

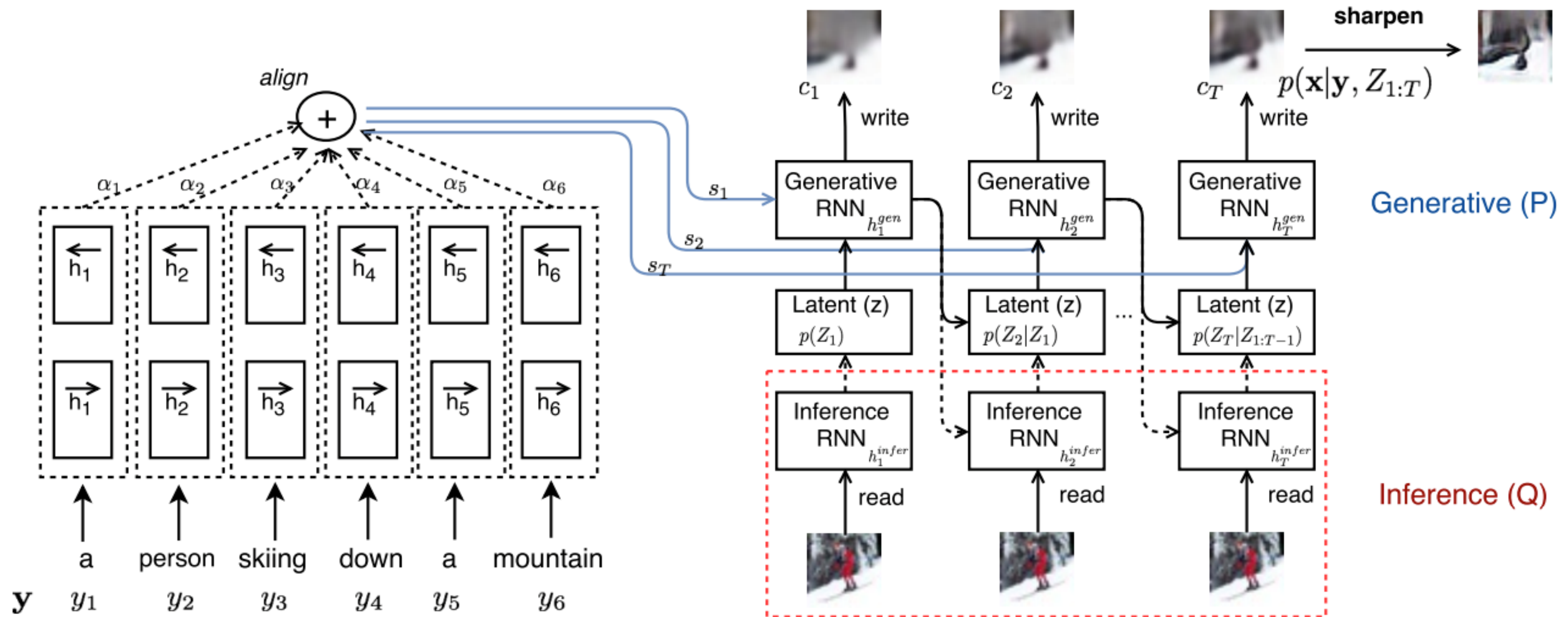


A person skiing on sand clad vast desert.

Generating Images from Captions with Attention.

Elman Mansimov, Emilio Parisotto, Jimmy Lei Ba, Ruslan Salakhutdinov. ICLR 2016.

First Deep Learning Work



VAES + RNN+ cross-attention

Generating Images from Captions with Attention.

Elman Mansimov, Emilio Parisotto, Jimmy Lei Ba, Ruslan Salakhutdinov. ICLR 2016.

Can we improve it?

How can we improve it?

- Better generative modeling techniques.
- Better text encoders.
- Better generator architectures.
- Better ways to connect text and image.
- Bigger data + more GPU/TPU computing.
- Bigger model sizes.

GAN-based Text-to-Image

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.



GAN-based Text-to-Image

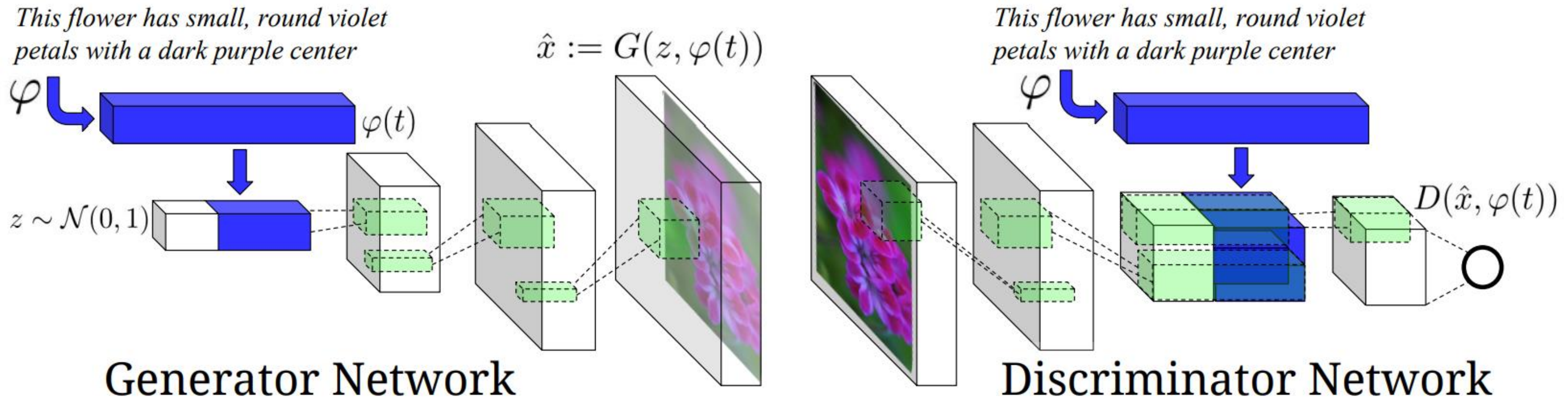
the flower has petals that are bright pinkish purple with white stigma



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



GAN-based Text-to-Image

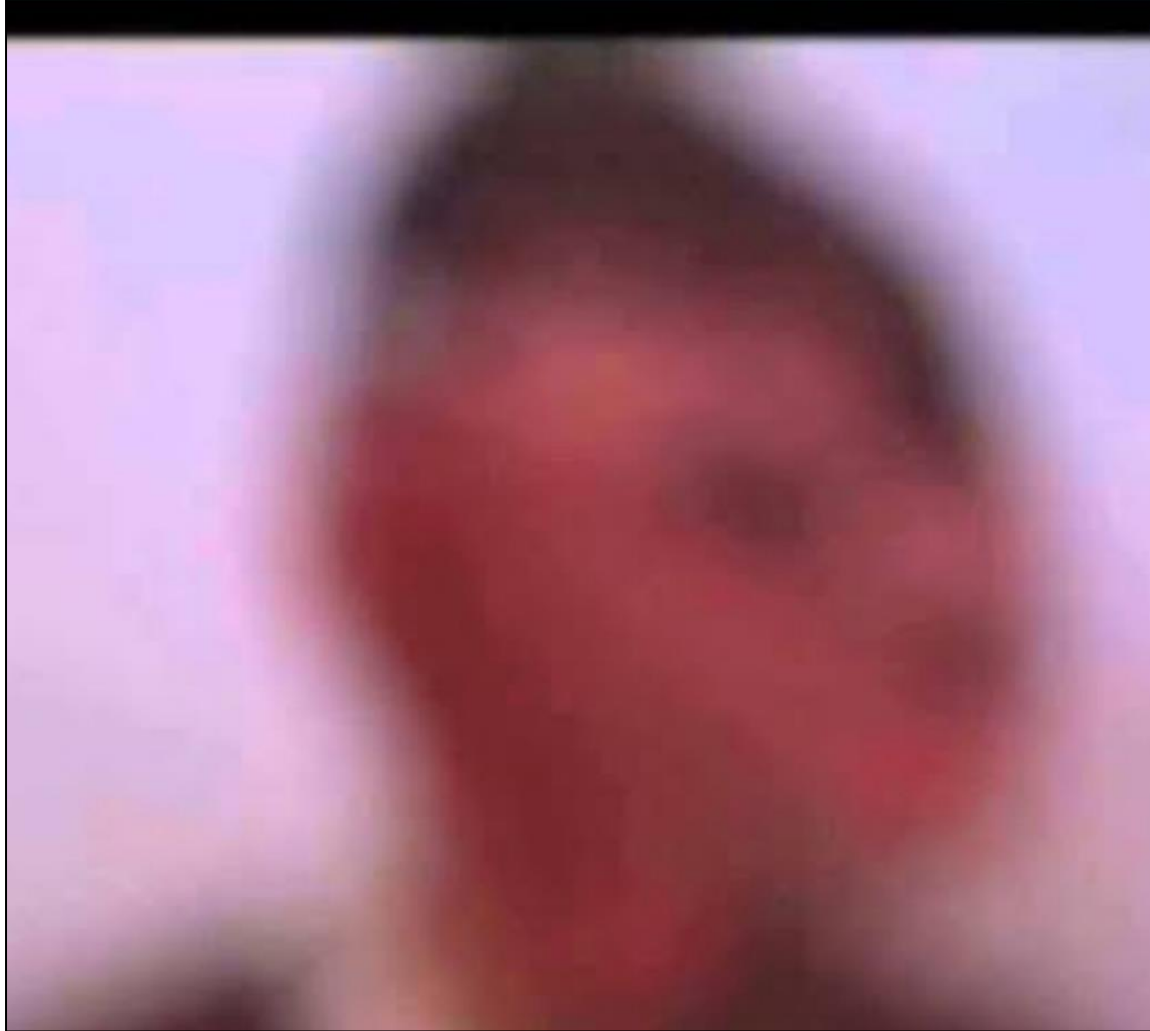


Conditional GAN + CNN + concatenation

Generative Adversarial Text to Image Synthesis
Scott Reed et al., ICML 2016

Video title :- "Real footage of aliens caught on tape 100 percent authentic"

The video quality :-



4K



144p



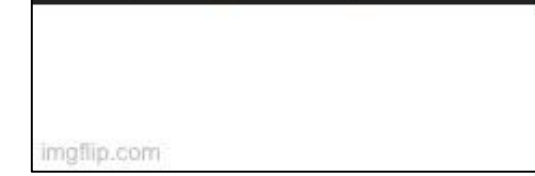
Gamepedia · 1h ago
its too high in quality lower it



Quality for current video · 144p

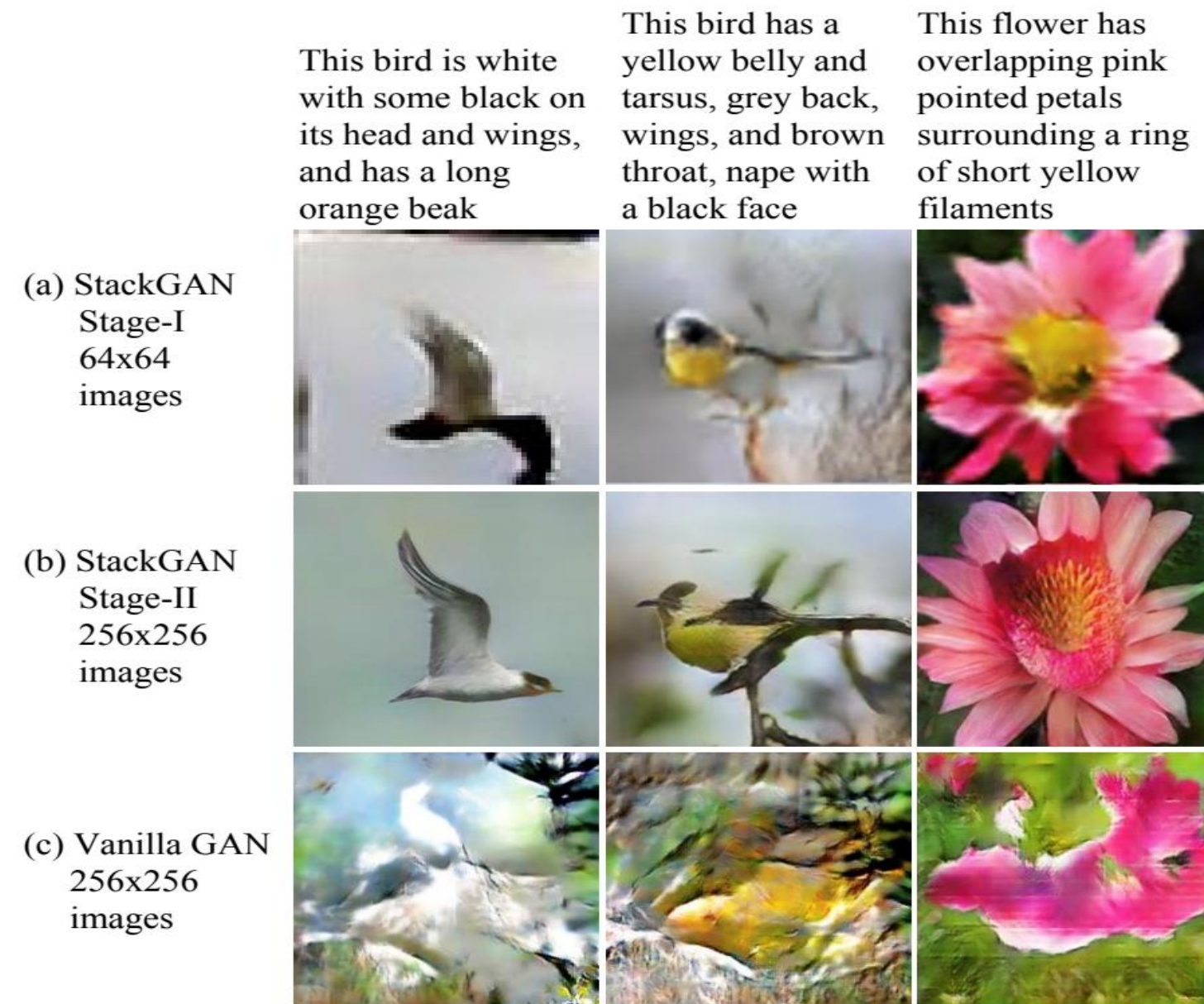
- 360p
- 240p
- ✓ 144p

This selection only applies to the current video. For all videos, go to Settings > Video quality preferences.



But these images are tiny ...
How can we make them HD?

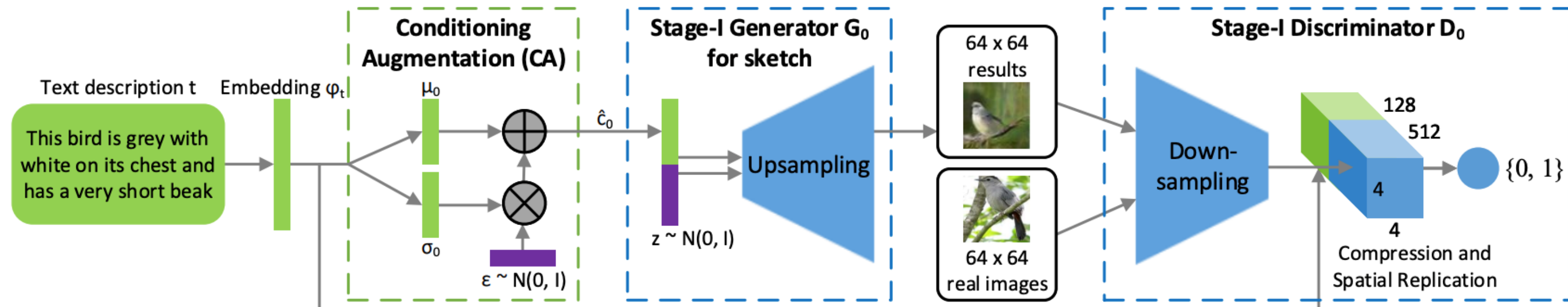
+Two-stage Models



Two-stage Conditional GAN + CNN + concatenation

StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks
Han Zhang et al., ICCV 2017

+Two-stage Models

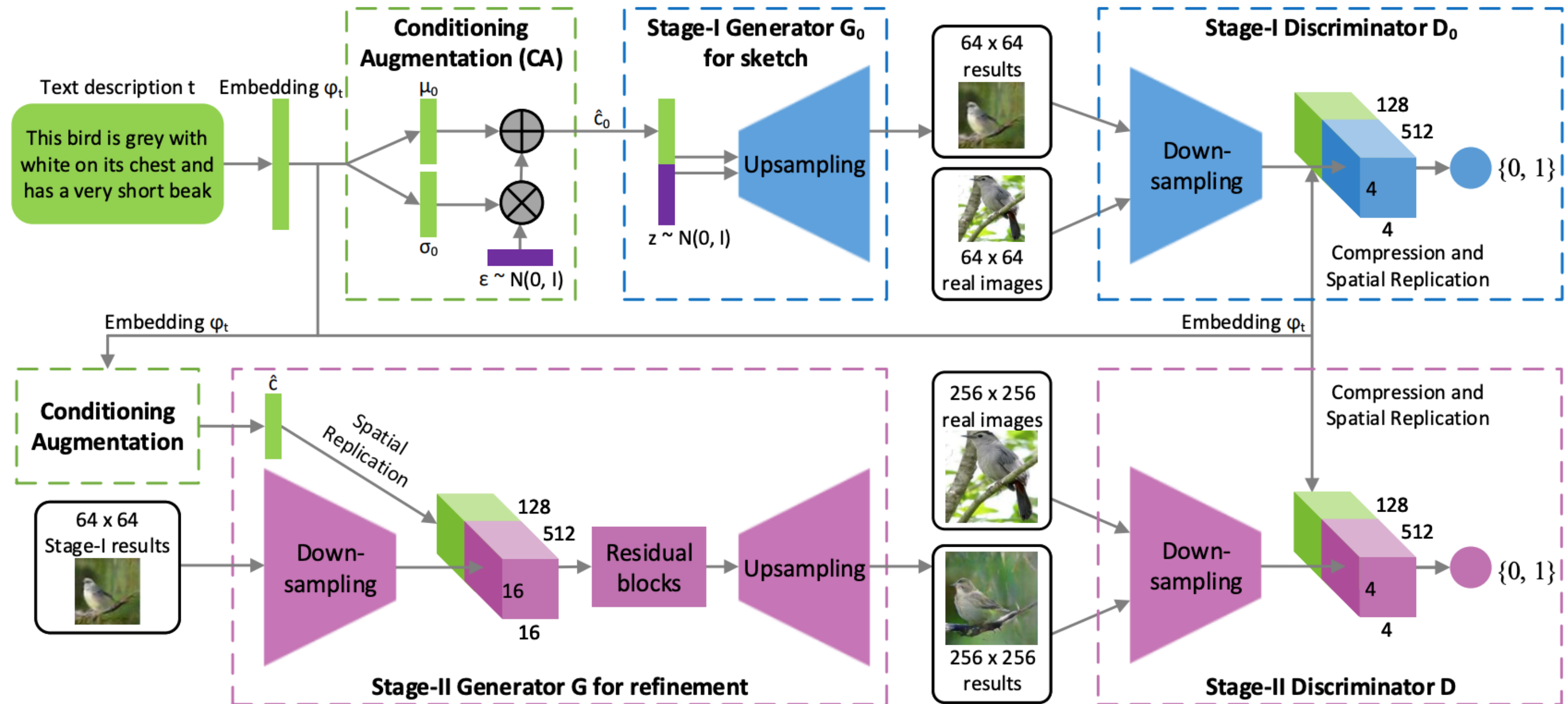


Two-stage Conditional GAN + CNN + concatenation

StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks

Han Zhang et al., ICCV 2017

+Two-stage Models

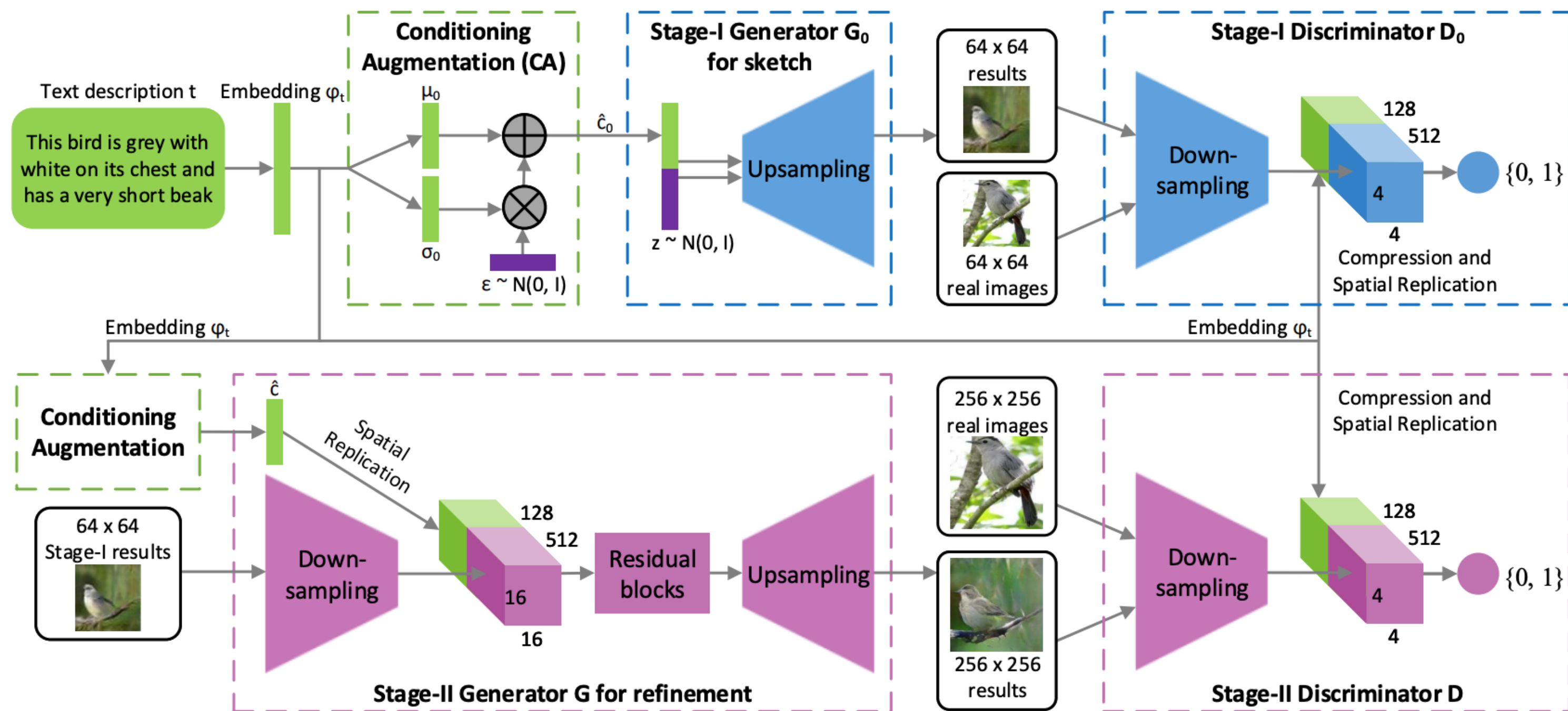


Two-stage Conditional GAN + CNN + concatenation

StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks

Han Zhang et al., ICCV 2017





+Two-stage Models



Two-stage Conditional GAN + CNN + concatenation

StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks
Han Zhang et al., ICCV 2017

+Two-stage Models

Text description	This flower has a lot of small purple petals in a dome-like configuration	This flower is pink, white, and yellow in color, and has petals that are striped	This flower has petals that are dark pink with white edges and pink stamen	This flower is white and yellow in color, with petals that are wavy and smooth
64x64 GAN-INT-CLS				
256x256 StackGAN				

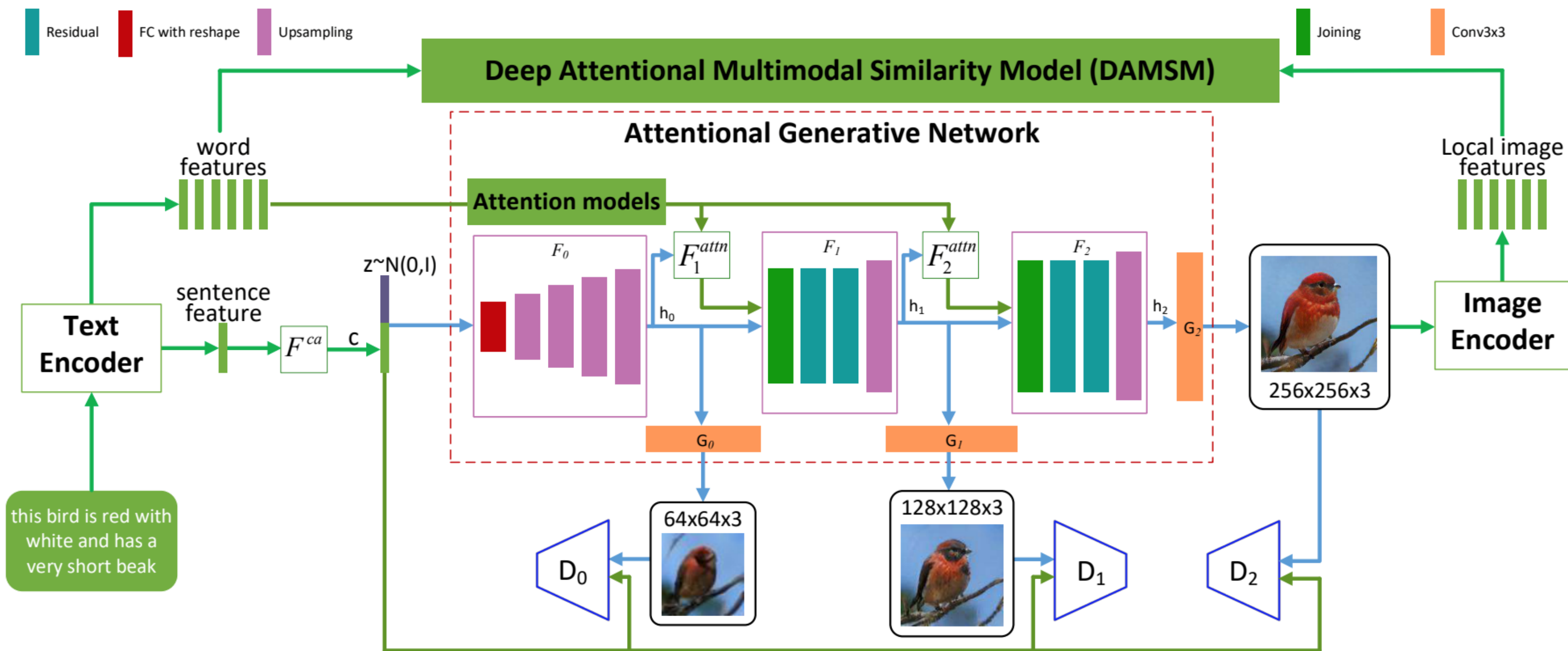
+Two-stage Models



+ Cross-attention to connect Text and Image



+ Cross-attention to connect Text and Image



AttnGAN: Fine-Grained Text to Image Generation with Attentional Generative Adversarial Networks
Tao Xu et al., CVPR 2018

**Got Stuck in 2018-2020
(Birds, MS COCO)**

Got Stuck in 2018-2020 (Birds, MS COCO)

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



A picture of a very clean living room



A group of people on skis stand in the snow



Eggs fruit candy nuts and meat served on white dish



A street sign on a stoplight pole in the middle of a day



Can we synthesize images
beyond single or a few categories

Taming Transformers for High-Resolution Image Synthesis

Patrick Esser* Robin Rombach* Björn Ommer
Heidelberg Collaboratory for Image Processing, IWR, Heidelberg University, Germany

*Both authors contributed equally to this work



Figure 1. Our approach enables transformers to synthesize high-resolution images like this one, which contains 1280x460 pixels.

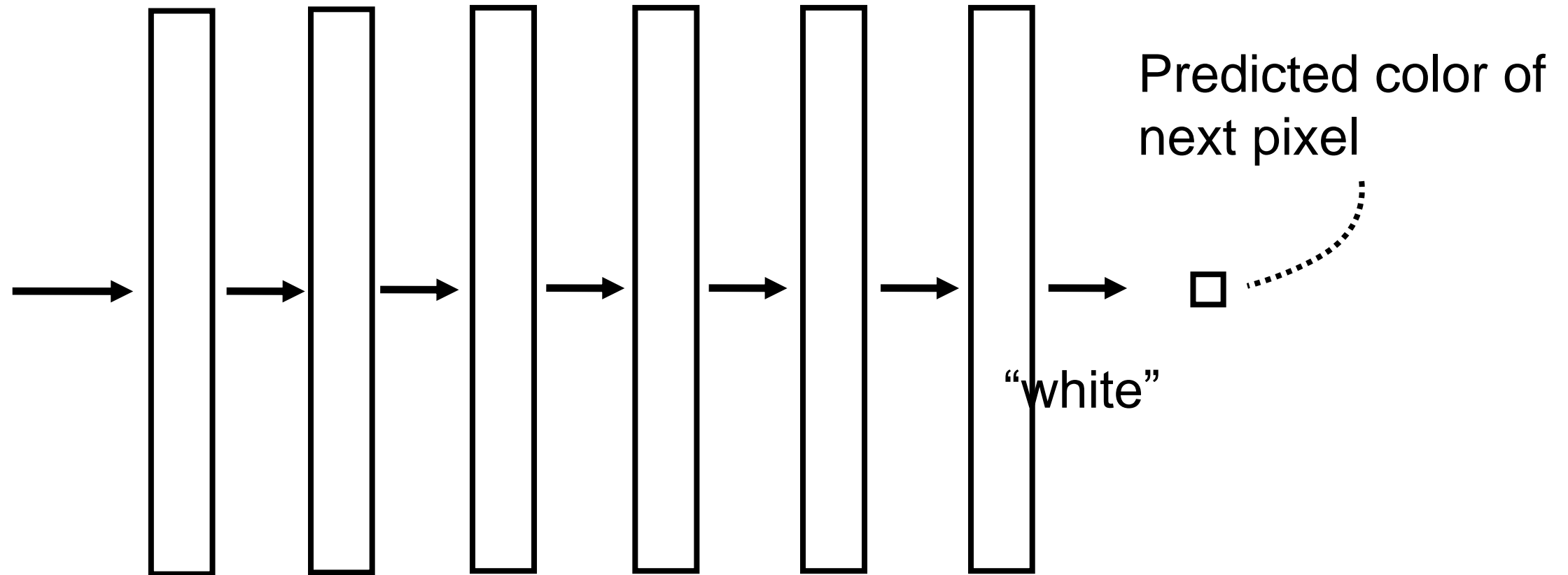
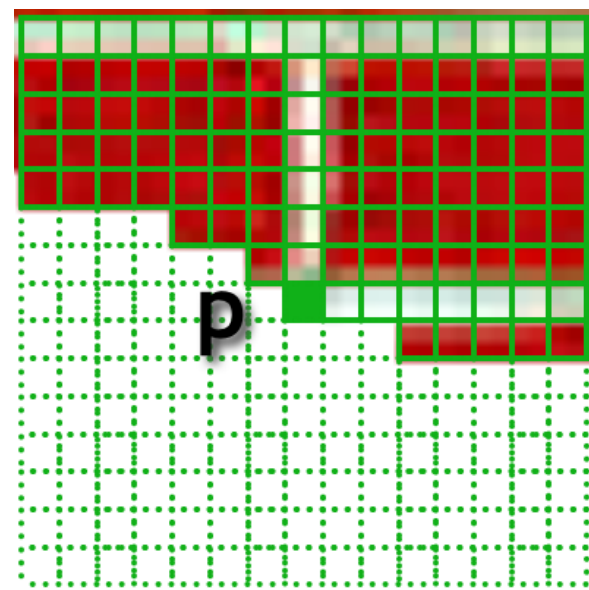
Abstract

Designed to learn long-range interactions on sequential data, transformers have recently achieved state-of-the-art results

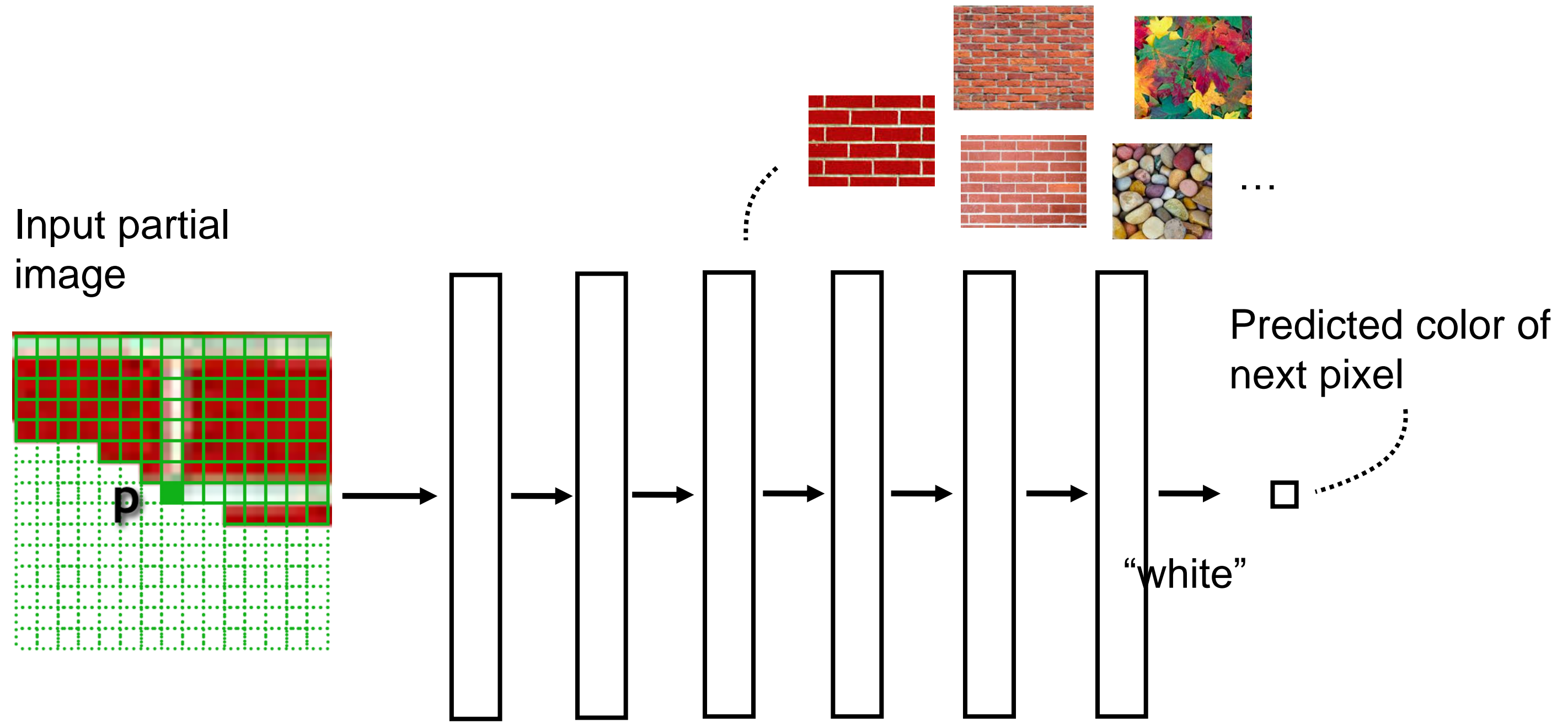
and are increasingly adapted in other areas such as audio [12] and vision [8, 16]. In contrast to the predominant vision architecture, convolutional neural networks (CNNs), the transformer architecture contains no built-in inductive bias, the locality of interactions and is therefore free to learn long-range dependencies from its inputs. However,

Autoregressive (AR) image synthesis

Input partial
image

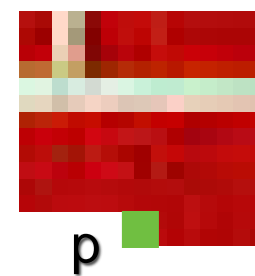
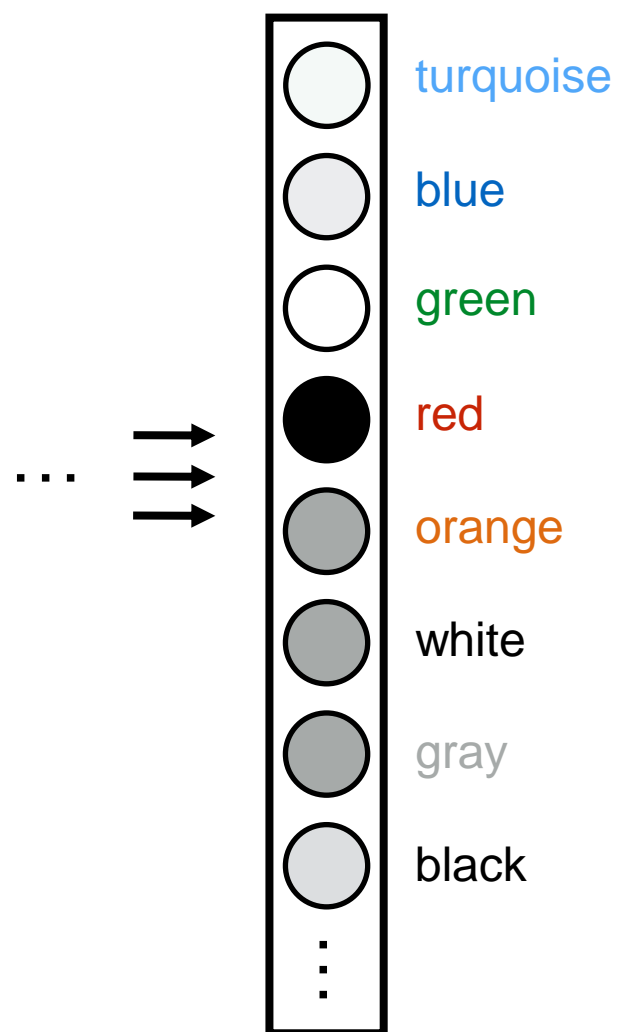


[PixelRNN, PixelCNN, van der Oord et al. 2016]



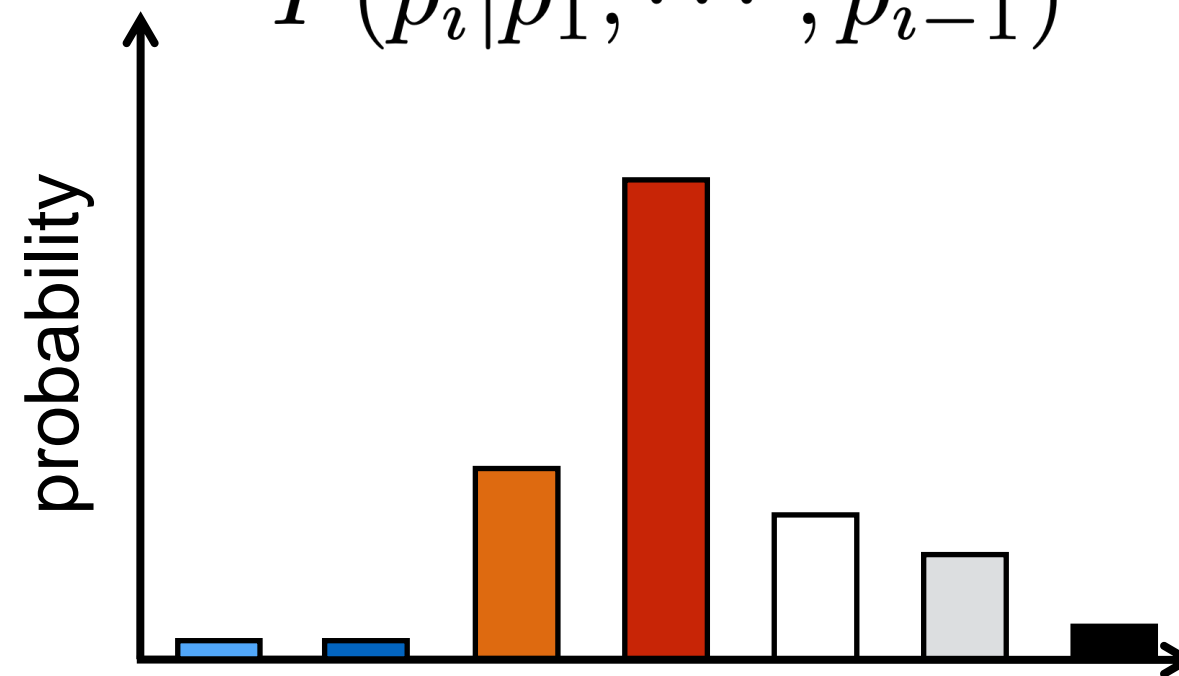
[PixelRNN, PixelCNN, van der Oord et al. 2016]

Network output

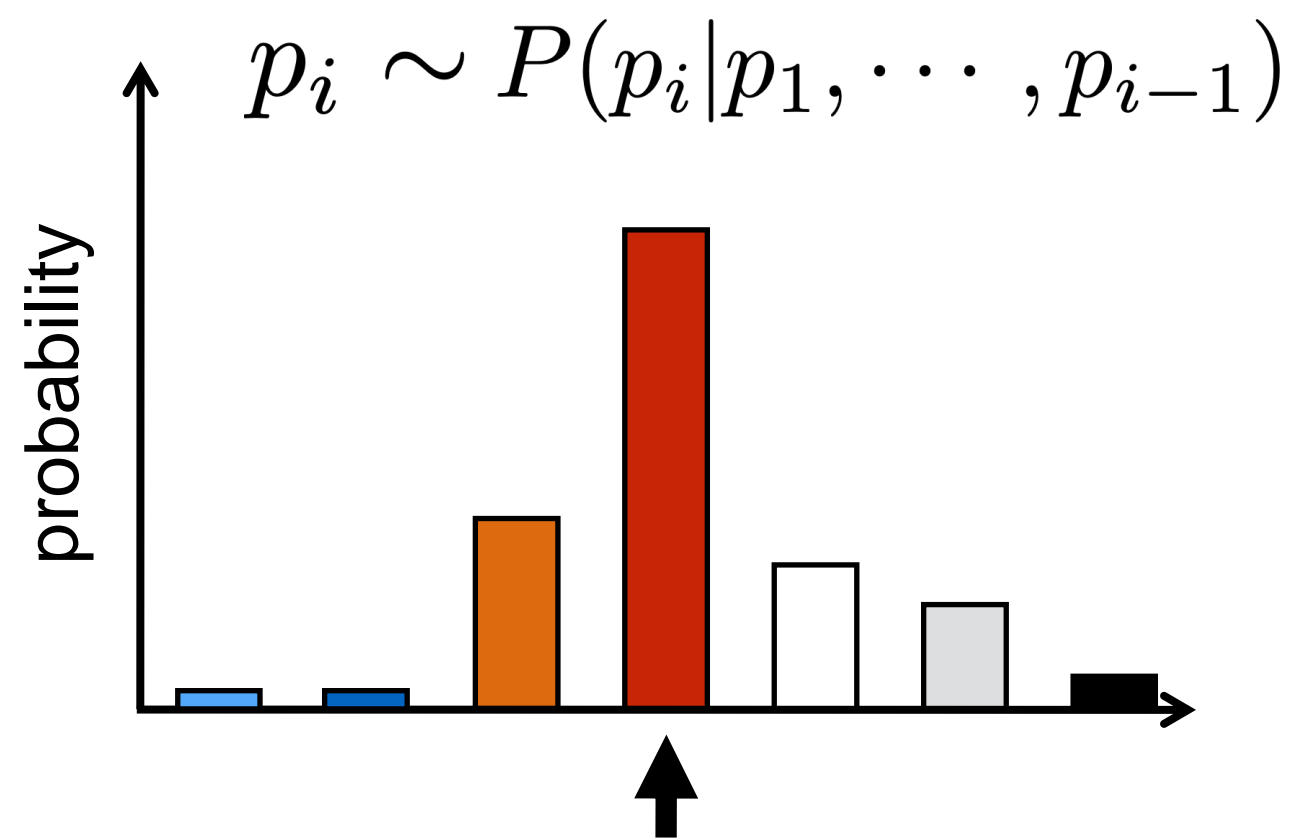
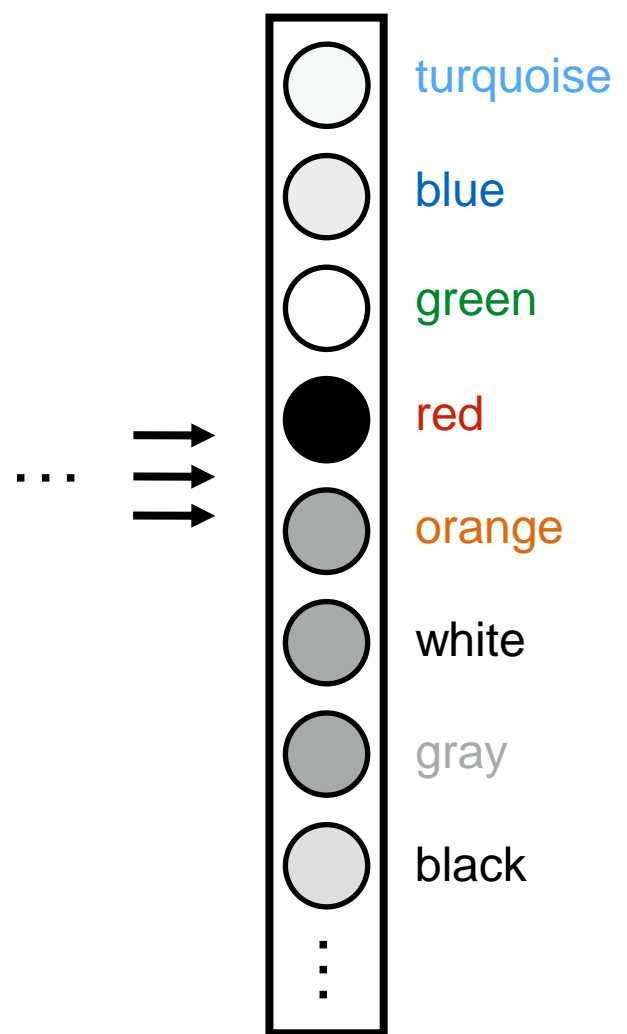


P(next pixel | previous pixels)

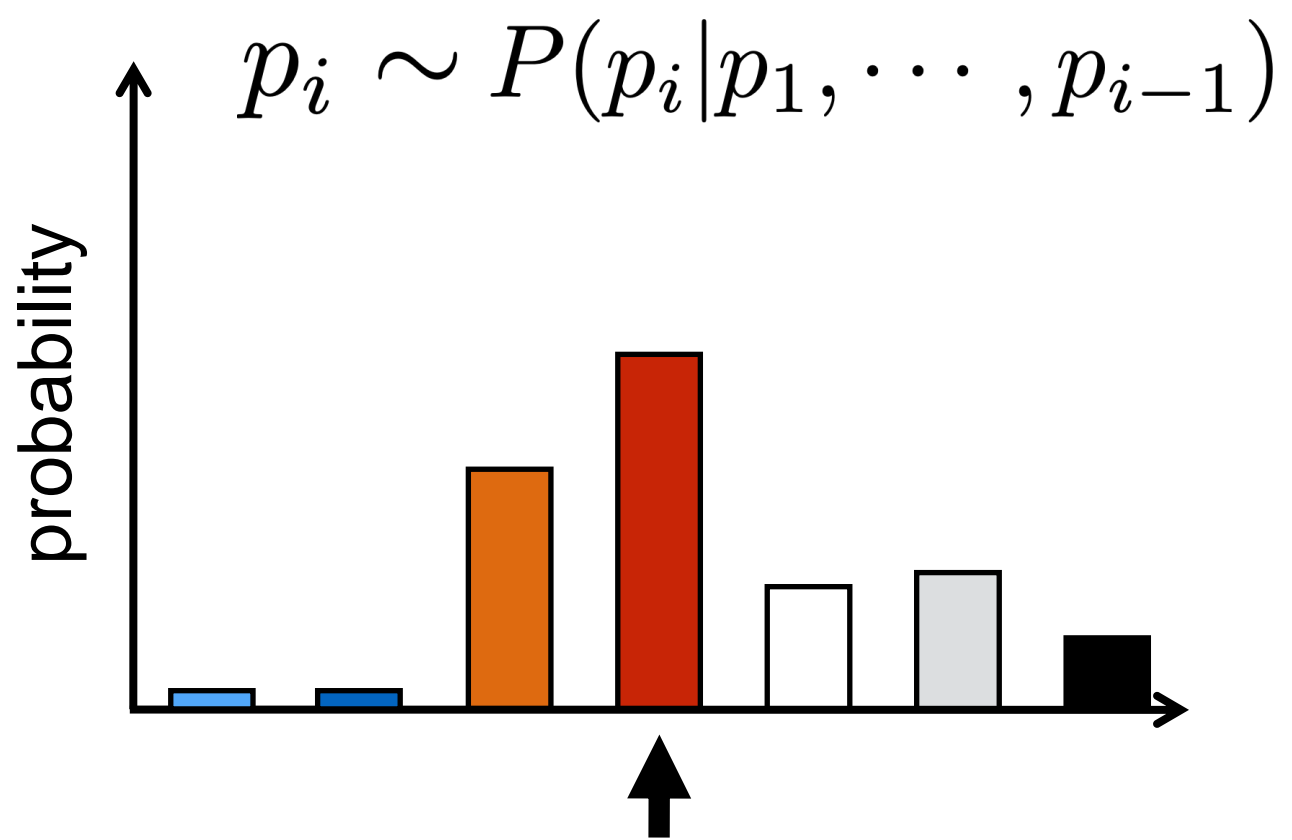
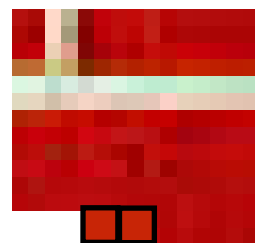
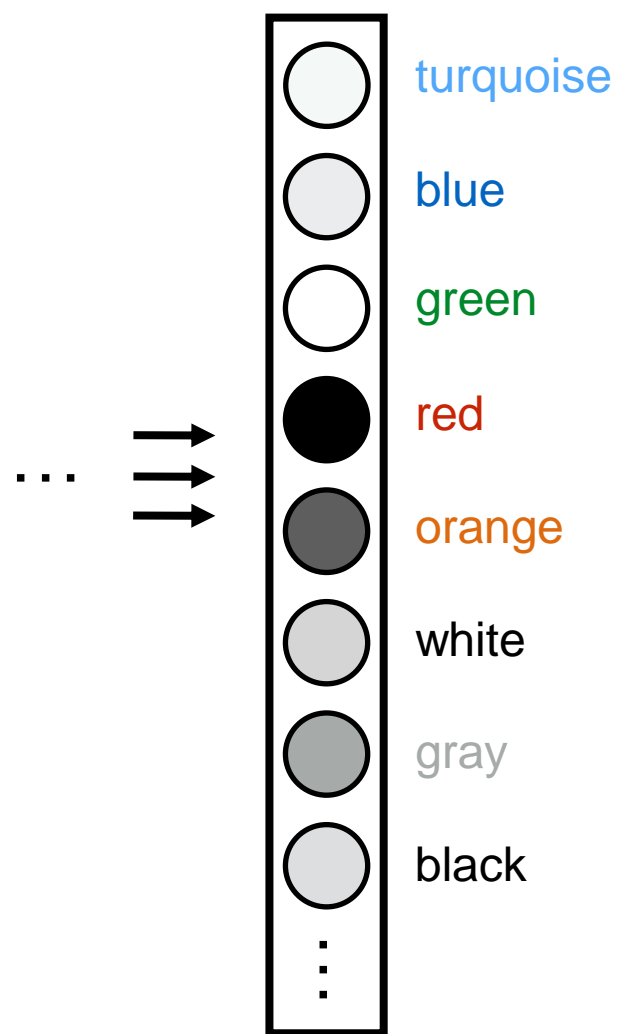
$$P(p_i | p_1, \dots, p_{i-1})$$



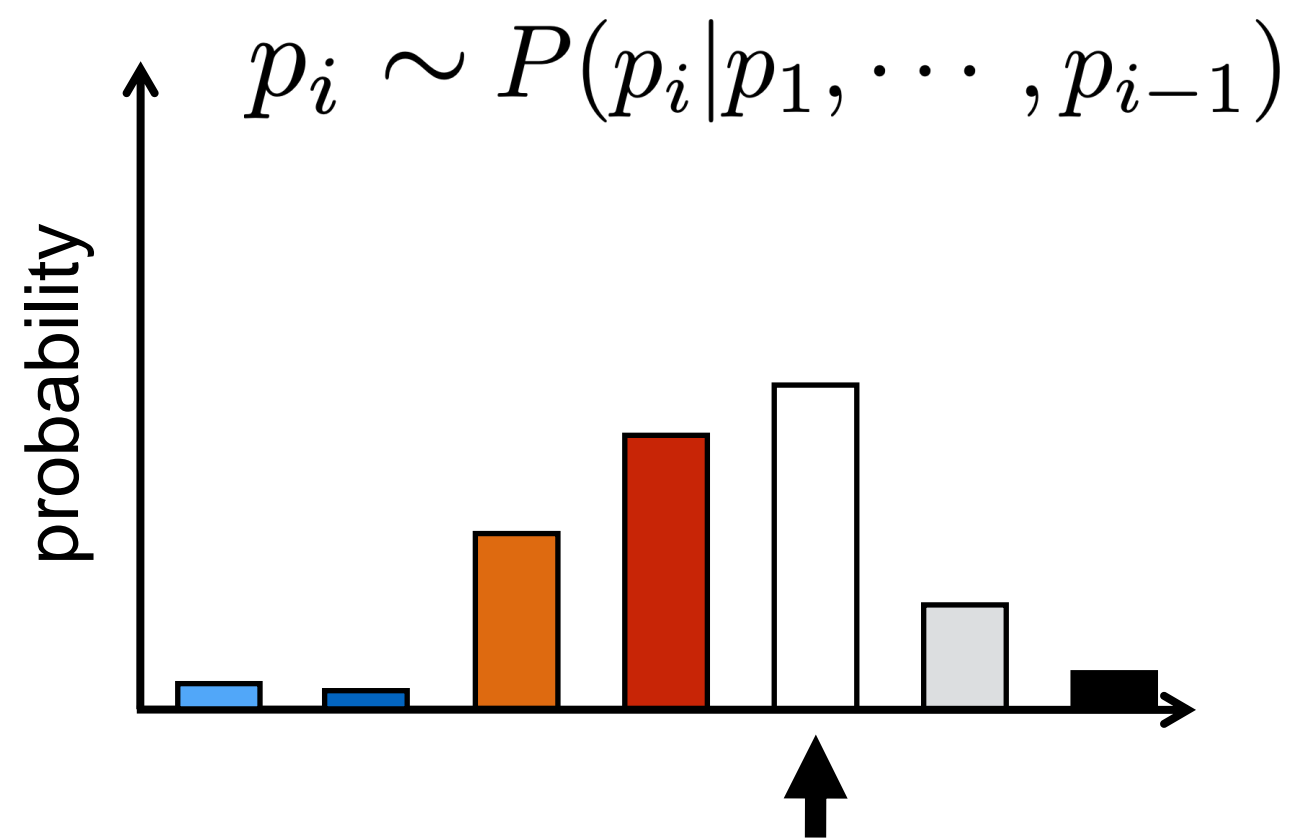
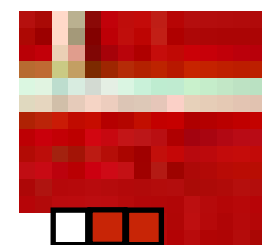
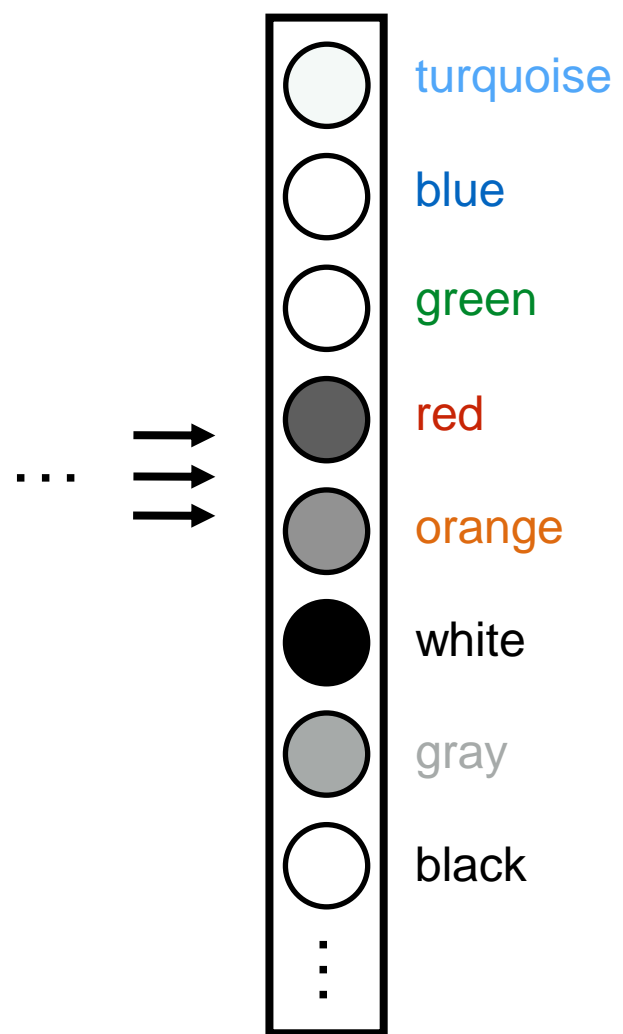
Network output



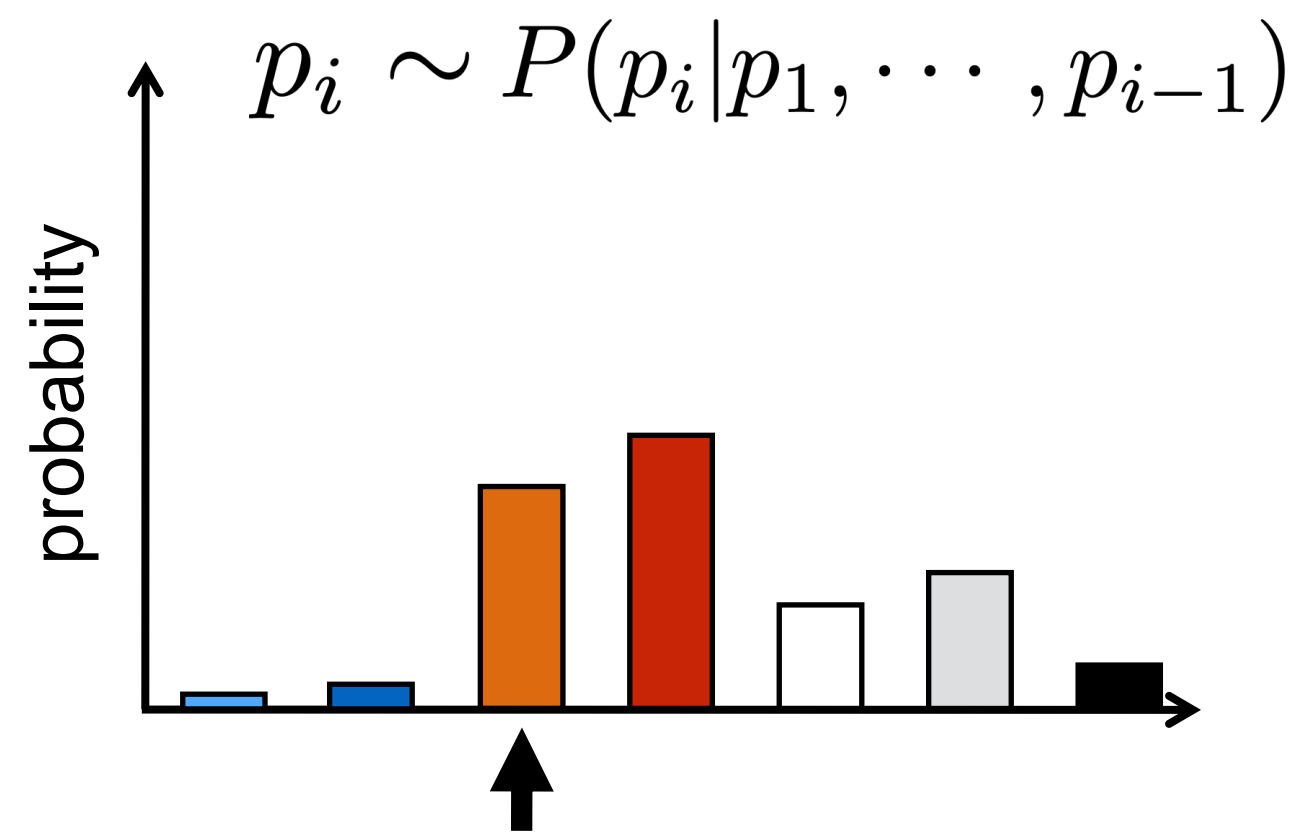
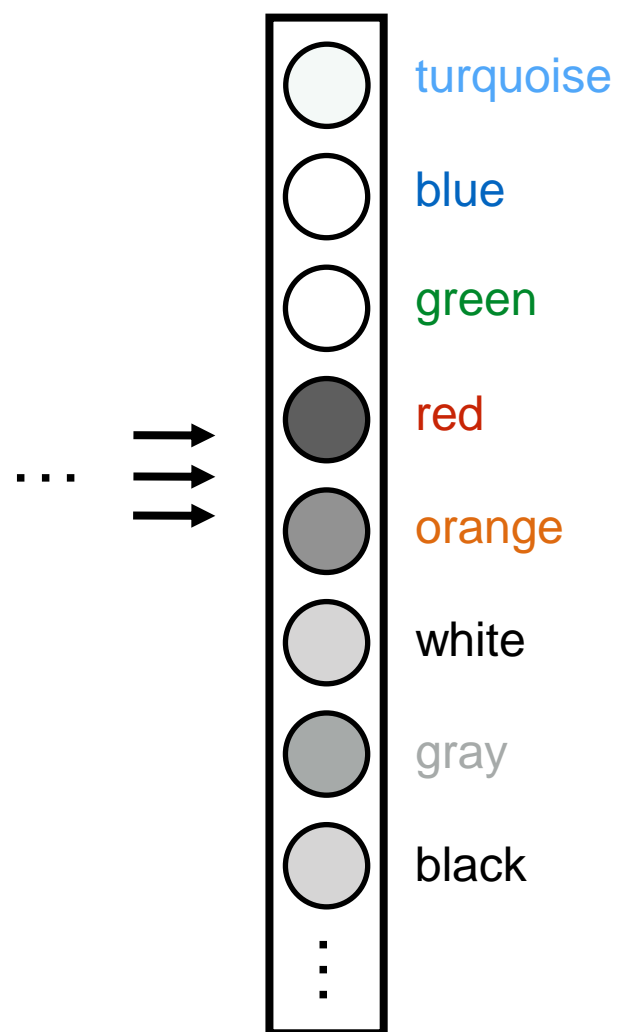
Network output



Network output



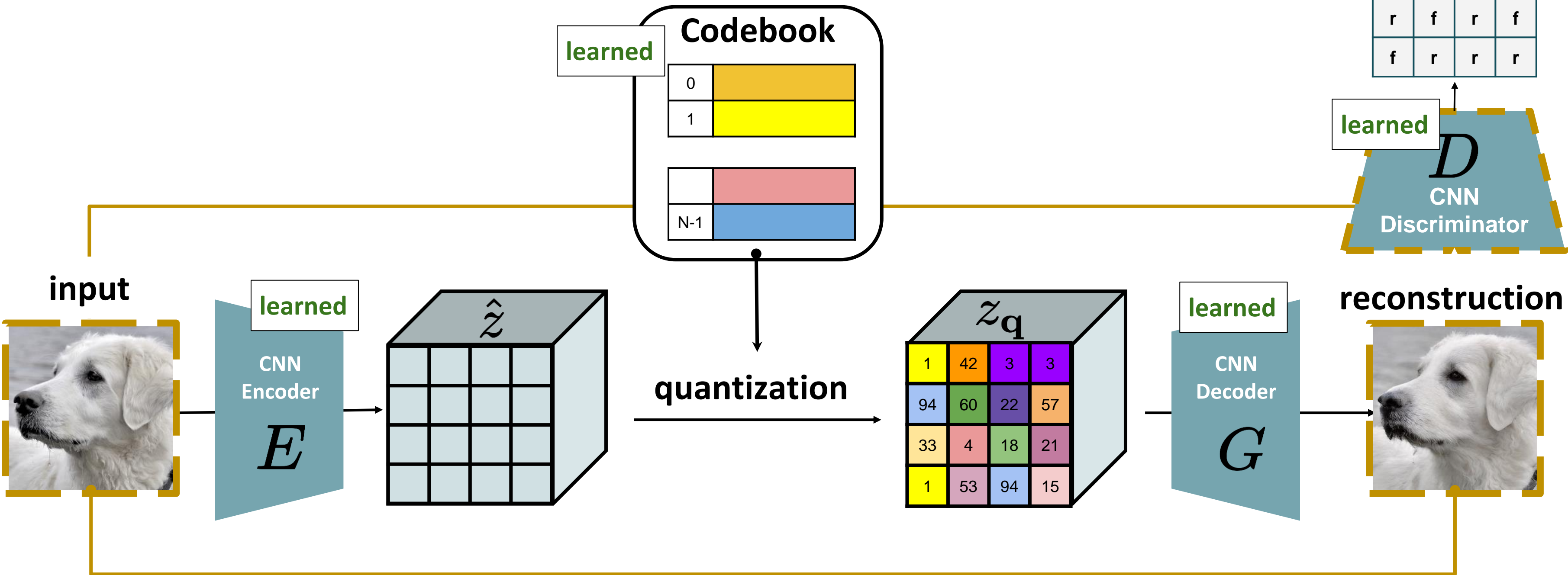
Network output



Generation is super slow?
What should we do?

From VQ-VAE¹ to VQGAN

¹: Neural Discrete Representation Learning, v.d.Oord et al, <https://arxiv.org/abs/1711.00937>

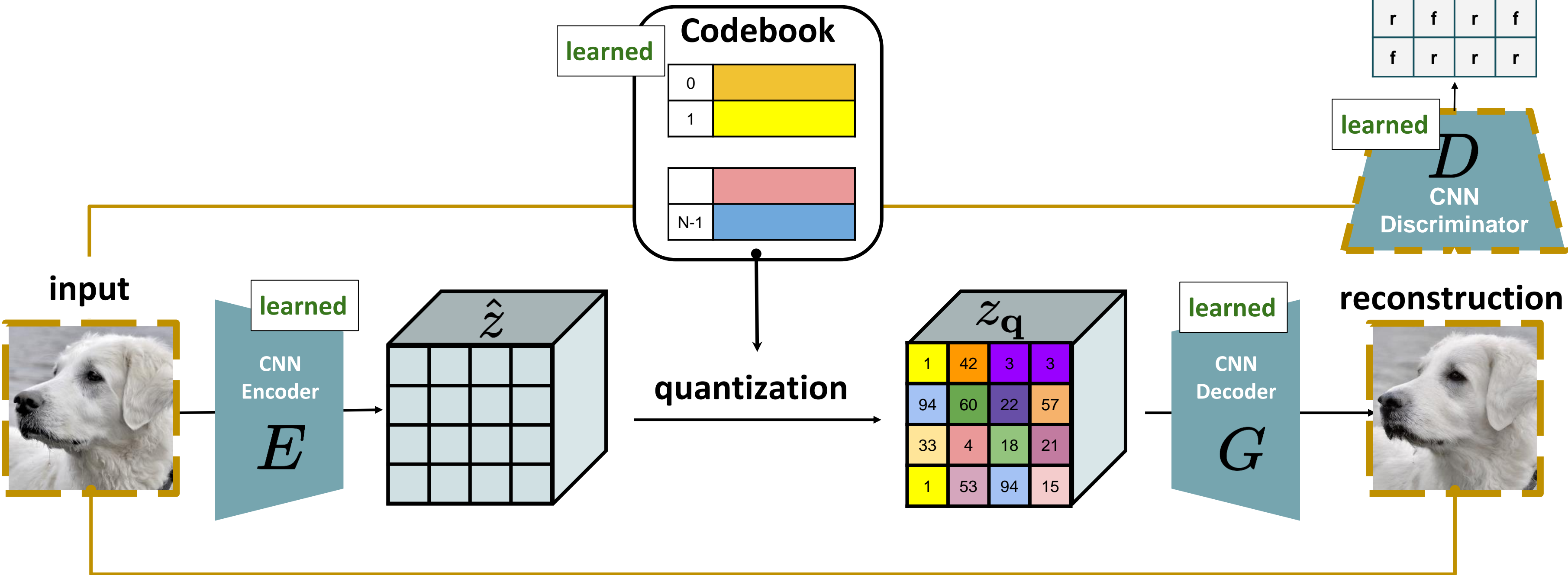


i) replace L2/L1 rec. loss with Perceptual loss (includes pixel-level)

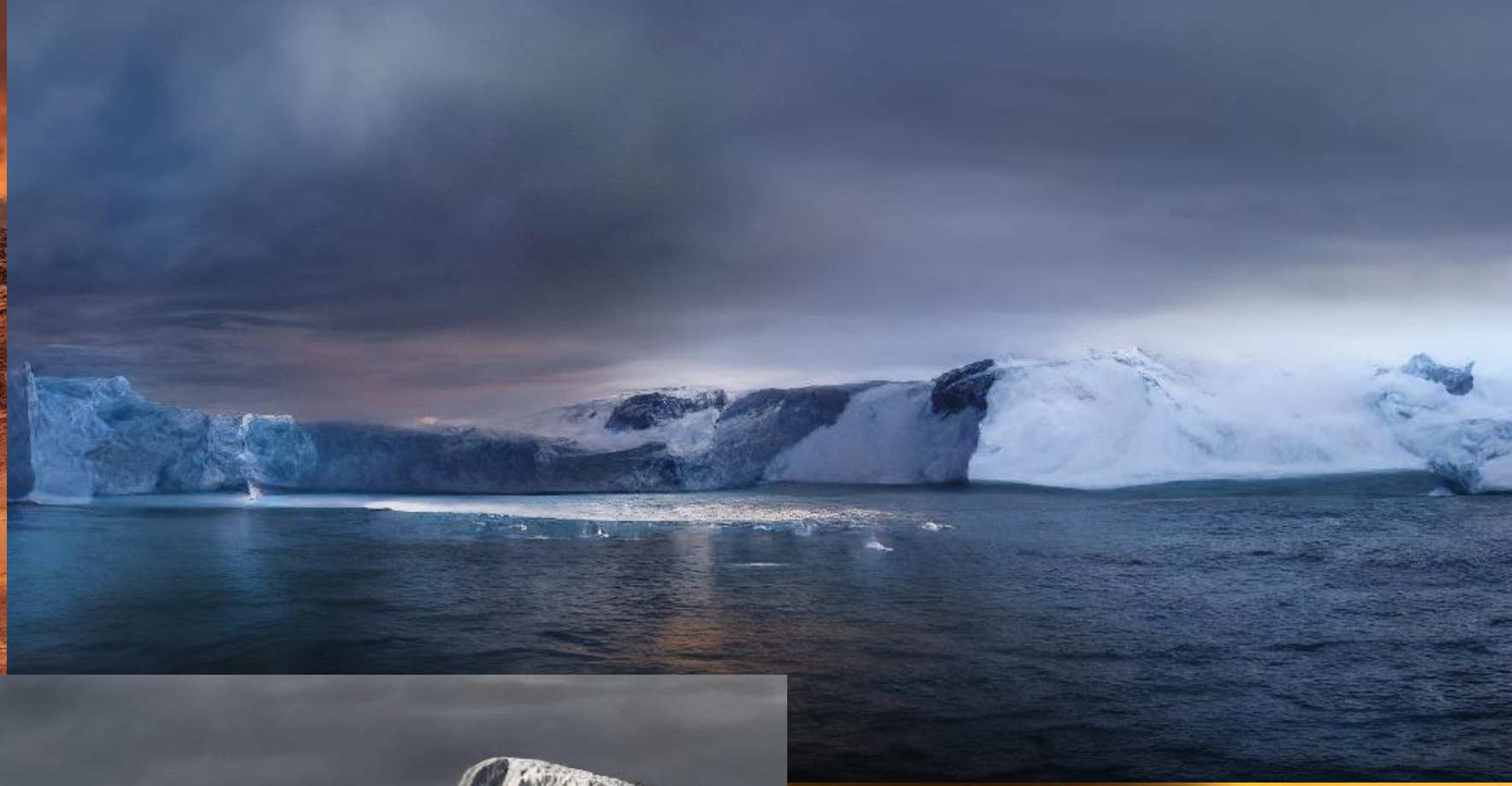
ii) add (patch-wise) Discriminator to favor realism over perfect reconstruction

From VQ-VAE¹ to VQGAN

¹: Neural Discrete Representation Learning, v.d.Oord et al, <https://arxiv.org/abs/1711.00937>



$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{rec}} + \mathcal{L}_{\text{VQ}} + \lambda \mathcal{L}_{\text{GAN}} \text{ where } \lambda = \frac{\nabla_{G_L} [\mathcal{L}_{\text{rec}}]}{\nabla_{G_L} [\mathcal{L}_{\text{GAN}}] + \delta}$$



Slide credit: Robin Rombach

Scaling VQGAN for Text-to-Image!

- see recently released “Parti” paper by Google (text-to-image model)
 - <https://parti.research.google/>

350M

750M

3B

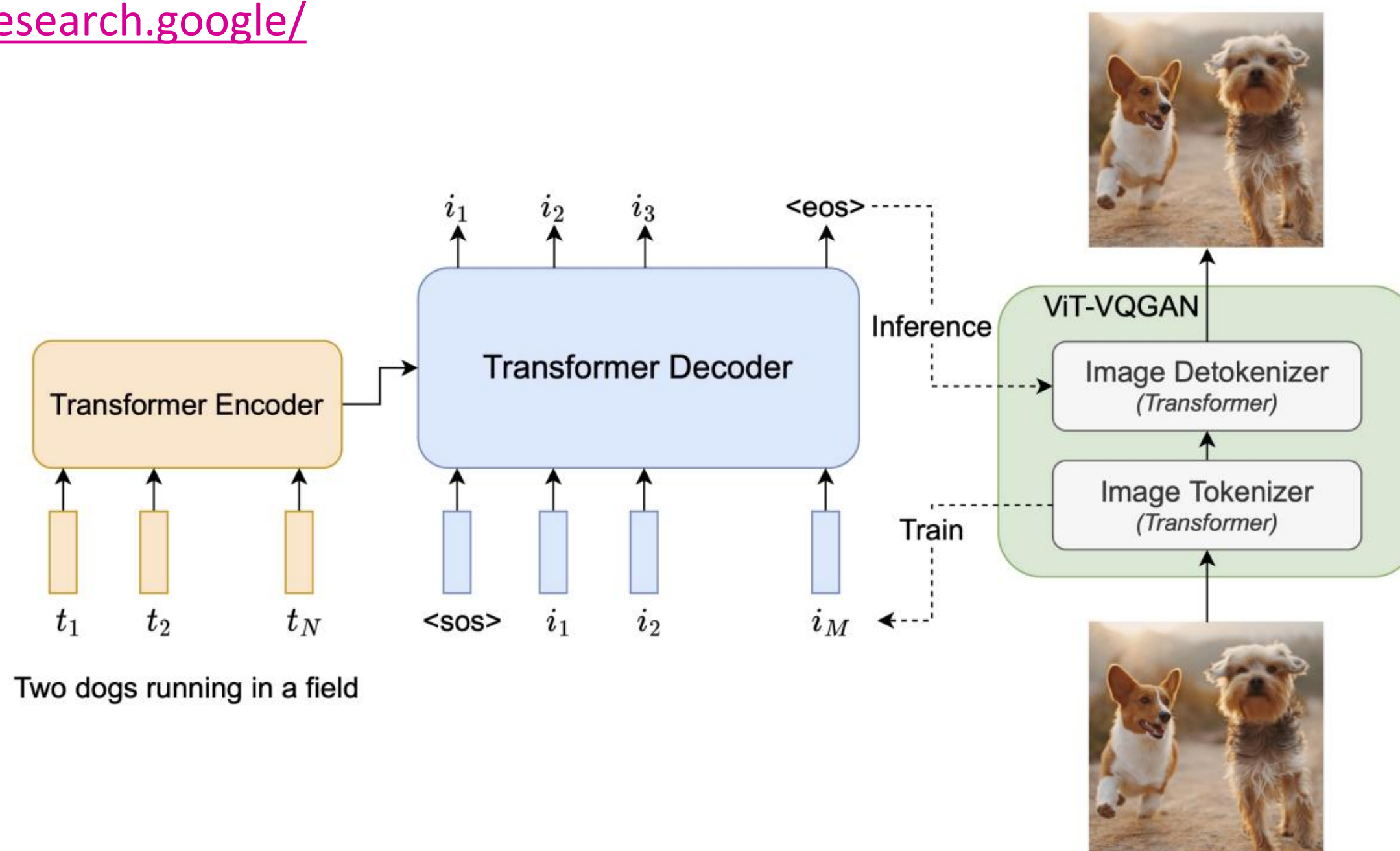
20B



A portrait photo of a kangaroo wearing an orange hoodie and blue sunglasses standing on the grass in front of the Sydney Opera House holding a sign on the chest that says Welcome Friends!

Scaling VQGAN for Text-to-Image!

- see recently released “Parti” paper by Google (text-to-image model)
 - <https://parti.research.google/>



Transformer-based Encoder/Decoder + Transformer-based Autoregressive models

Another Approach: Diffusion Models!

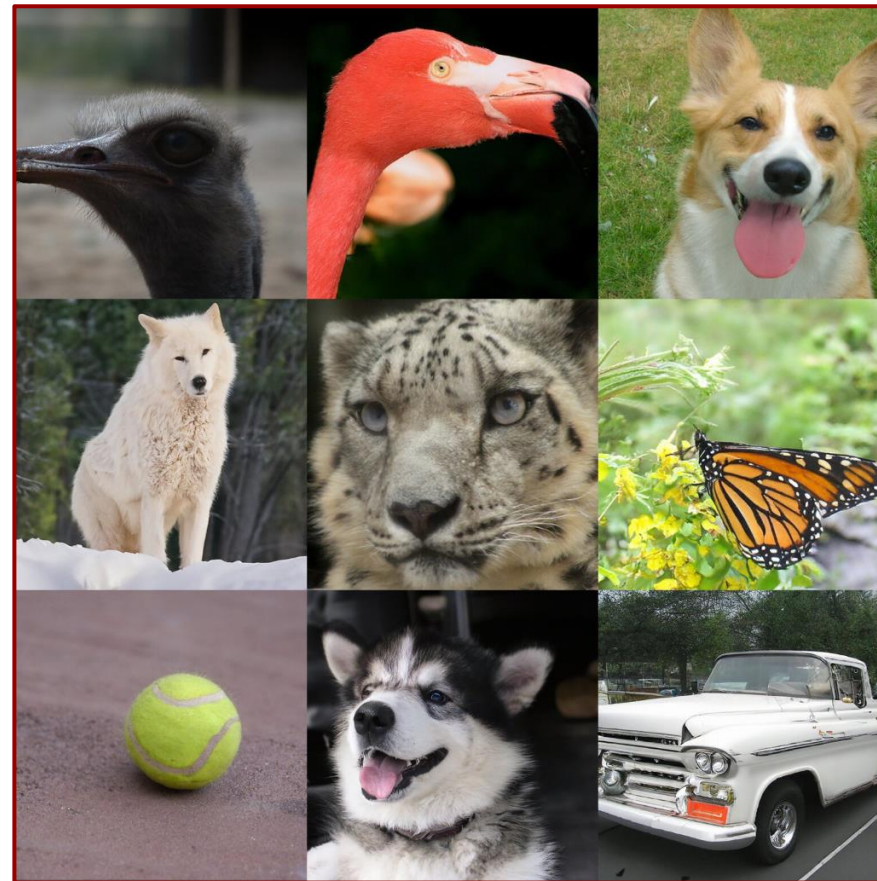
great results for image synthesis



Denoising Diffusion Probabilistic Models

Jonathan Ho, Ajay Jain, et al

<https://arxiv.org/abs/2006.11239>



Diffusion Models beat GANs on Image Synthesis

Prafulla Dhariwal, Alex Nichol

<https://arxiv.org/abs/2105.05233>

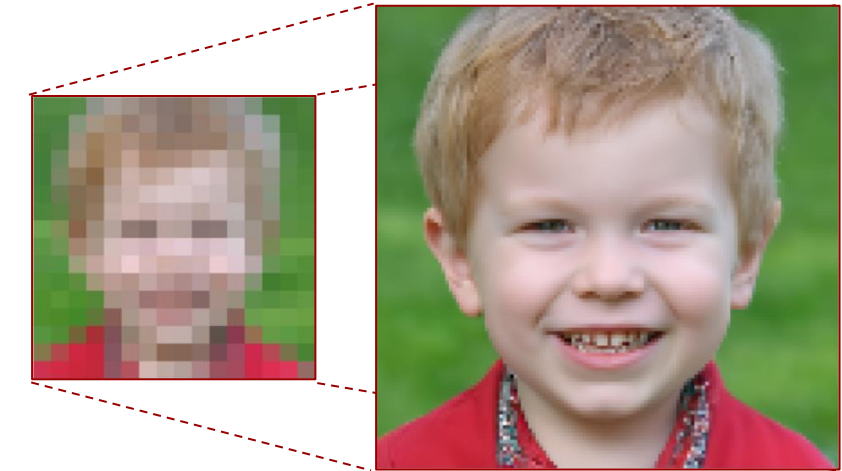


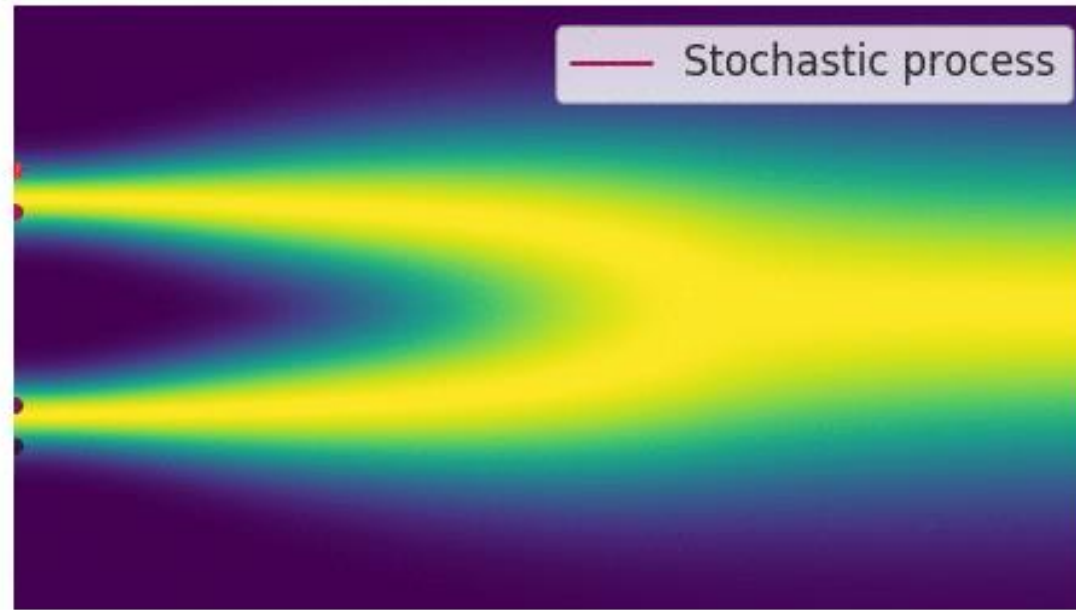
Image Super-Resolution via Iterative Refinement

Chitwan Saharia, et al

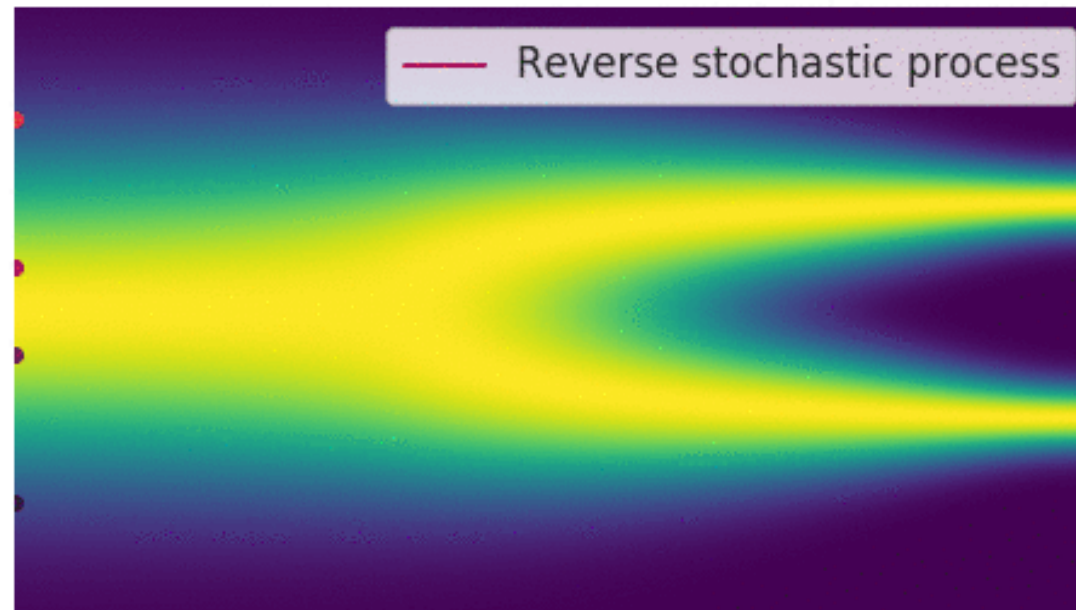
<https://arxiv.org/abs/2104.07636>

... but very expensive :(

Brief Overview of Diffusion Models



- “destroy” the data by gradually adding small amounts of gaussian noise



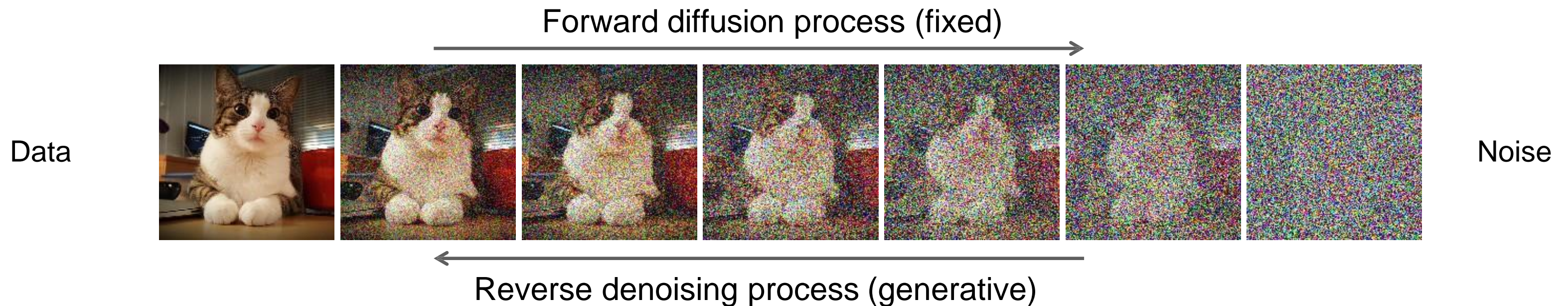
- “create” data by gradually denoising a noisy code from a stationary distribution

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](#)

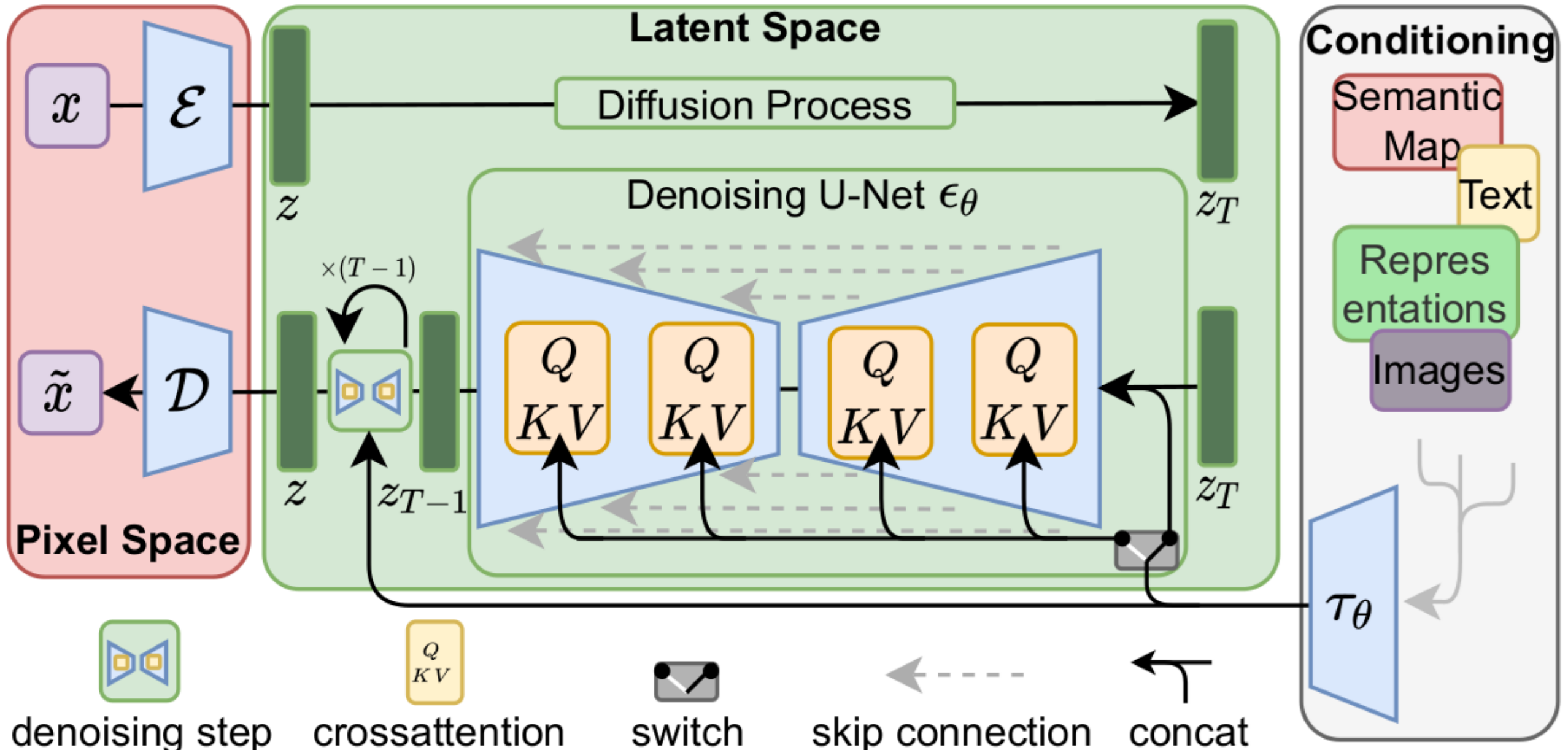
Latent Diffusion Modeling: Architecture

Autoencoder with KL or VQ regularization.

VQ-reg.: $\mathcal{L}_{total} = \mathcal{L}_{rec} + \mathcal{L}_{VQ} + \lambda \mathcal{L}_{GAN}$

where $\lambda = \frac{\nabla_{G_L}[\mathcal{L}_{rec}]}{\nabla_{G_L}[\mathcal{L}_{GAN}] + \delta}$



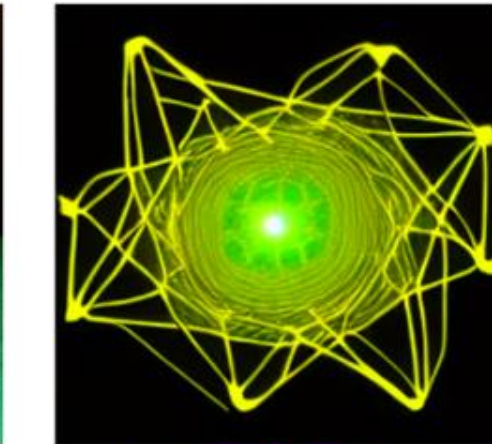









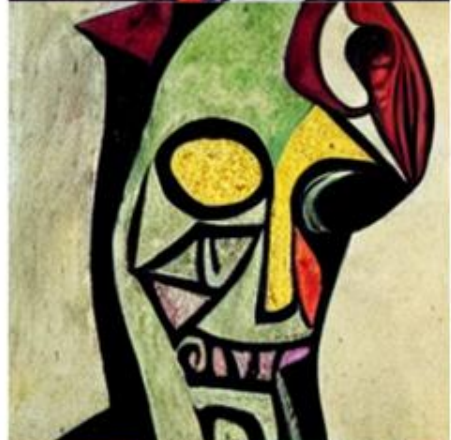





KL-reg.: $\mathcal{L}_{total} = \mathcal{L}_{rec} + \beta \mathcal{L}_{KL} + \lambda \mathcal{L}_{GAN}$



LDMs for Text-to-Image Synthesis

- 32x32 cont. space
- 600M Transformer
- 800M UNet
- 400M Image/Text Pairs

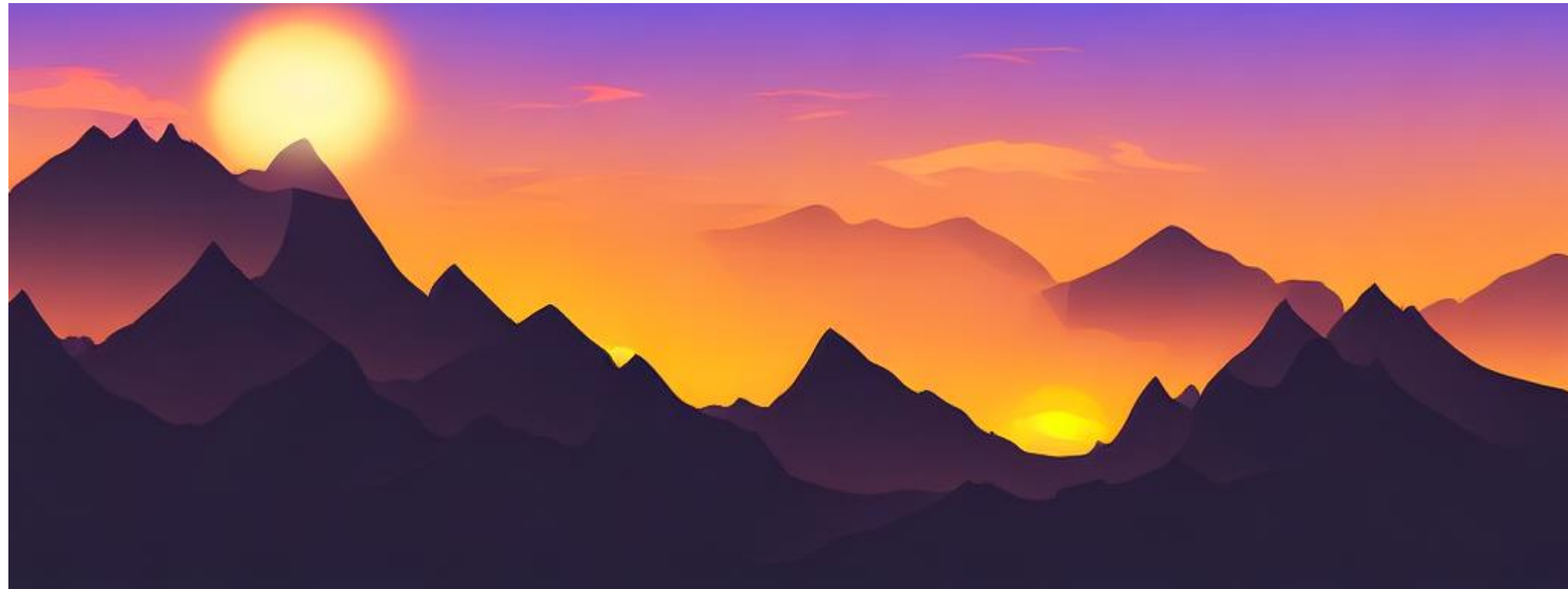
Text-to-Image Synthesis on LAION. 1.4B Model.

<i>'A zombie in the style of Picasso'</i>	<i>'An image of an animal half mouse half octopus'</i>	<i>'An illustration of a slightly conscious neural network.'</i>	<i>'A painting of a squirrel eating a burger.'</i>	<i>'A watercolor painting of a chair that looks like an octopus.'</i>	<i>'A shirt with the inscription: "I love generative models!"'</i>
					
					
					

LDMs for Text-to-Image Synthesis

convolutional sampling (train on 256^2 , generate on $>256^2$)

“A sunset over a mountain range, vector image”



“A sunset over a mountain range, oil on canvas”





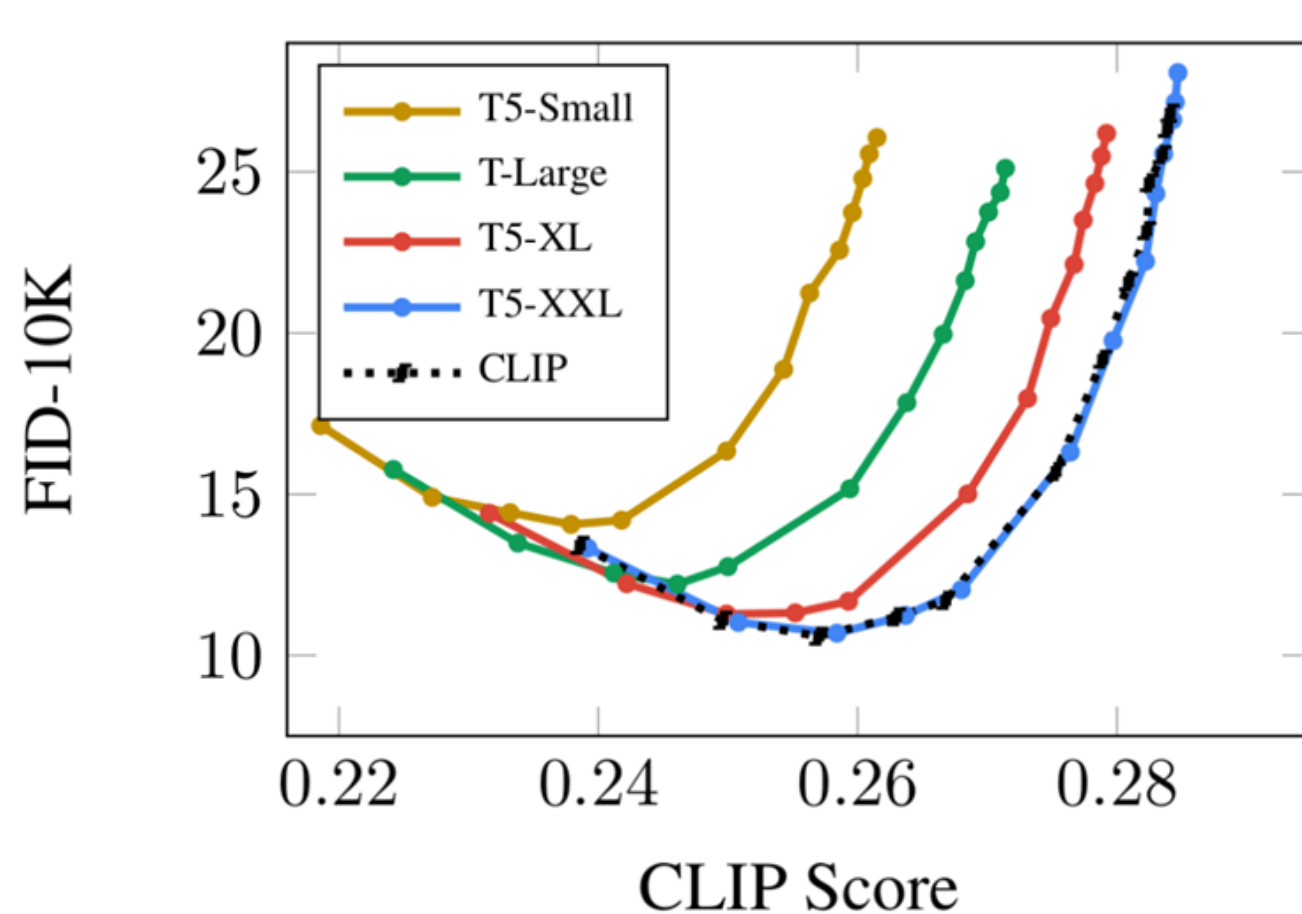
Stable Diffusion

Latent Diffusion ++

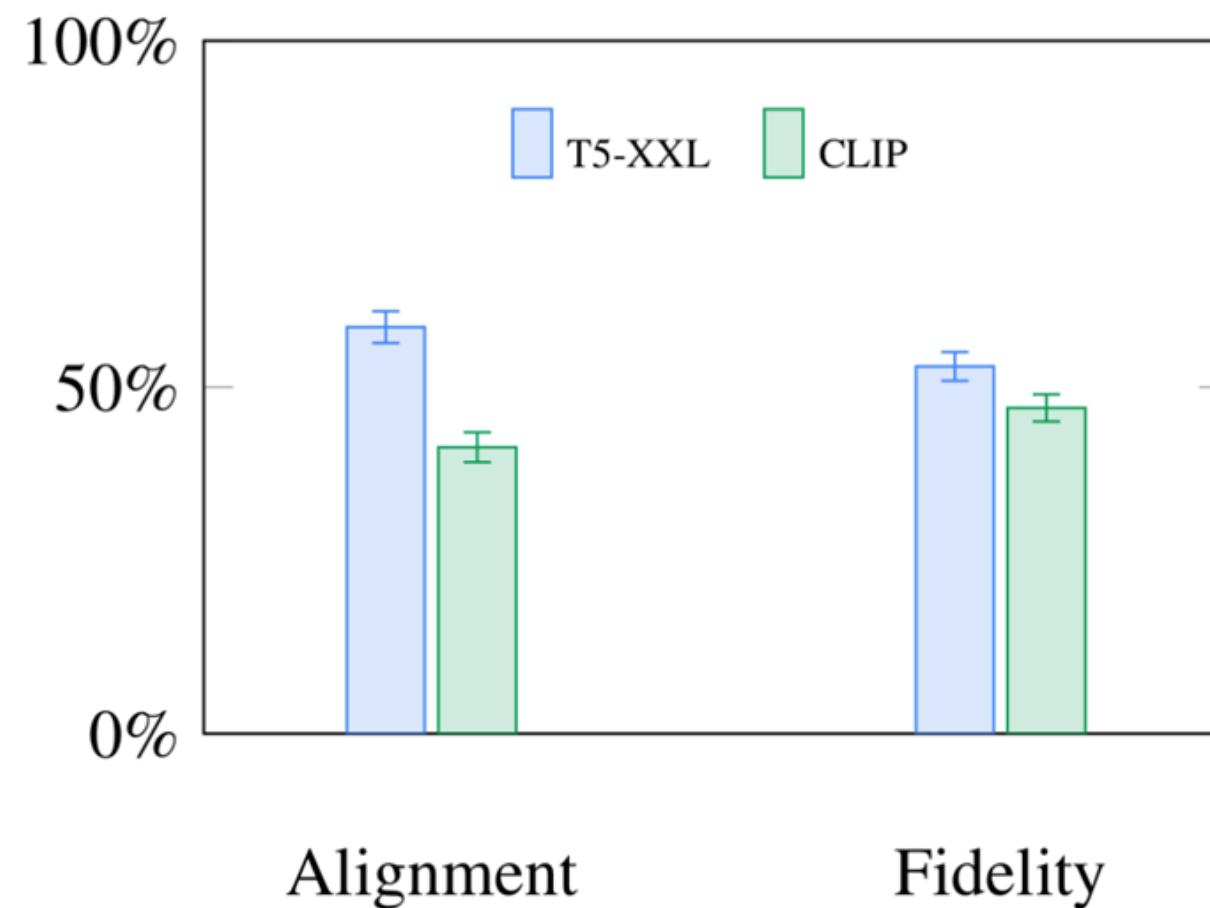


From Latent to Stable Diffusion

- goal: achieve a small model that people can actually run locally on “small” GPUs (~10GB VRAM)
- progressive training: pretrain on 256x256, then continue on 512x512
- fix text encoder (as in Imagen)
- → choose CLIP (ViT-L/14) since performance/size tradeoff seems significant



(a) Pareto curves comparing various text encoders.

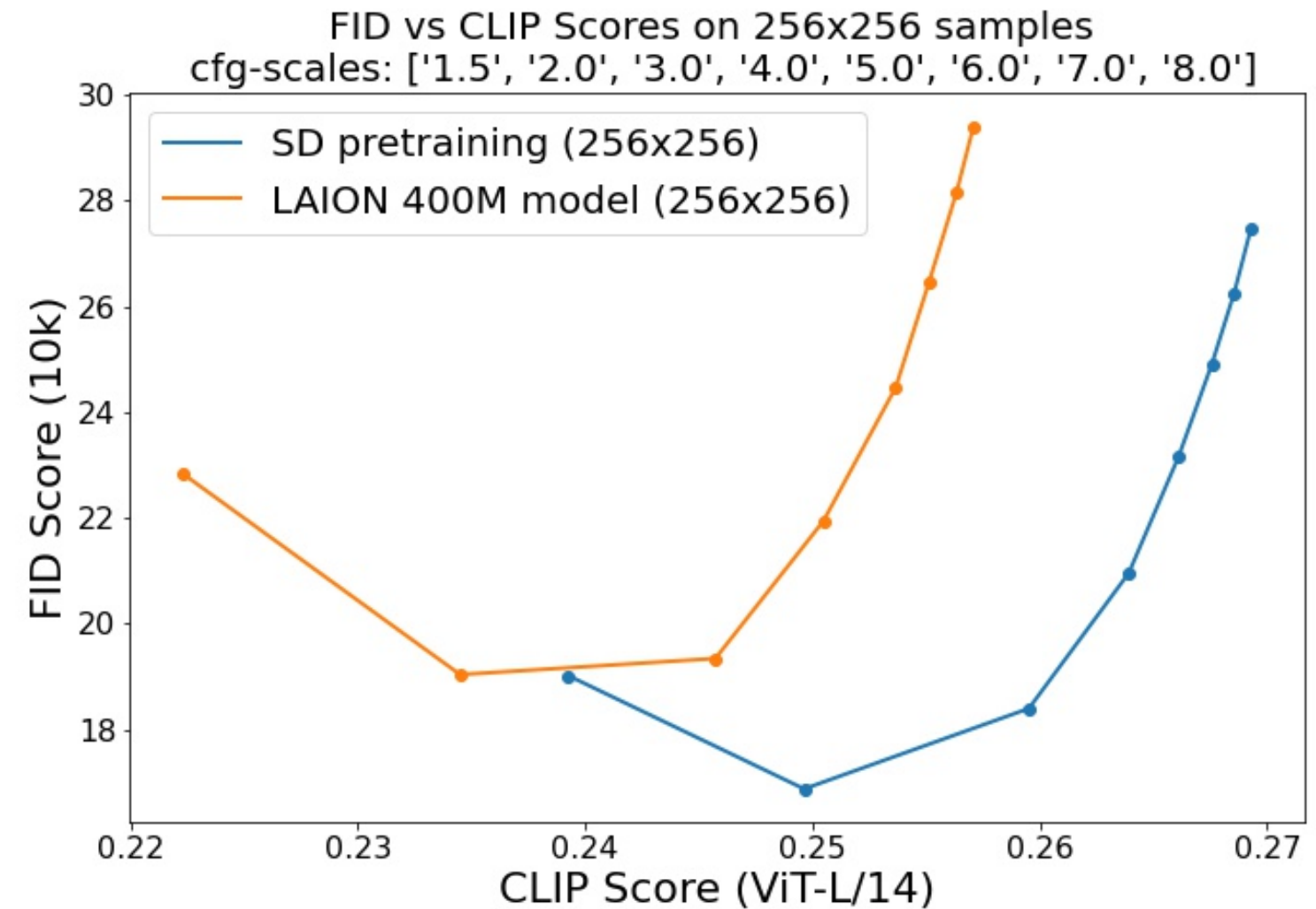


(b) Comparing T5-XXL and CLIP on DrawBench.

From Latent Diffusion to Stable Diffusion

Stage 1: Pretraining @256x256

- 237k steps at resolution 256x256 on LAION 2B(en)
- batch-size = 2048
- ~ 64 A100 GPUs



10k random COCO val captions / 50 decoding steps

From Latent Diffusion to Stable Diffusion

Stage 2: Training @512x512. batch-size=2048, #gpus=256

part 1 (v1.1):

- 194k steps at resolution 512x512 on laion-high-resolution (170M examples from LAION-5B with resolution $\geq 1024 \times 1024$).

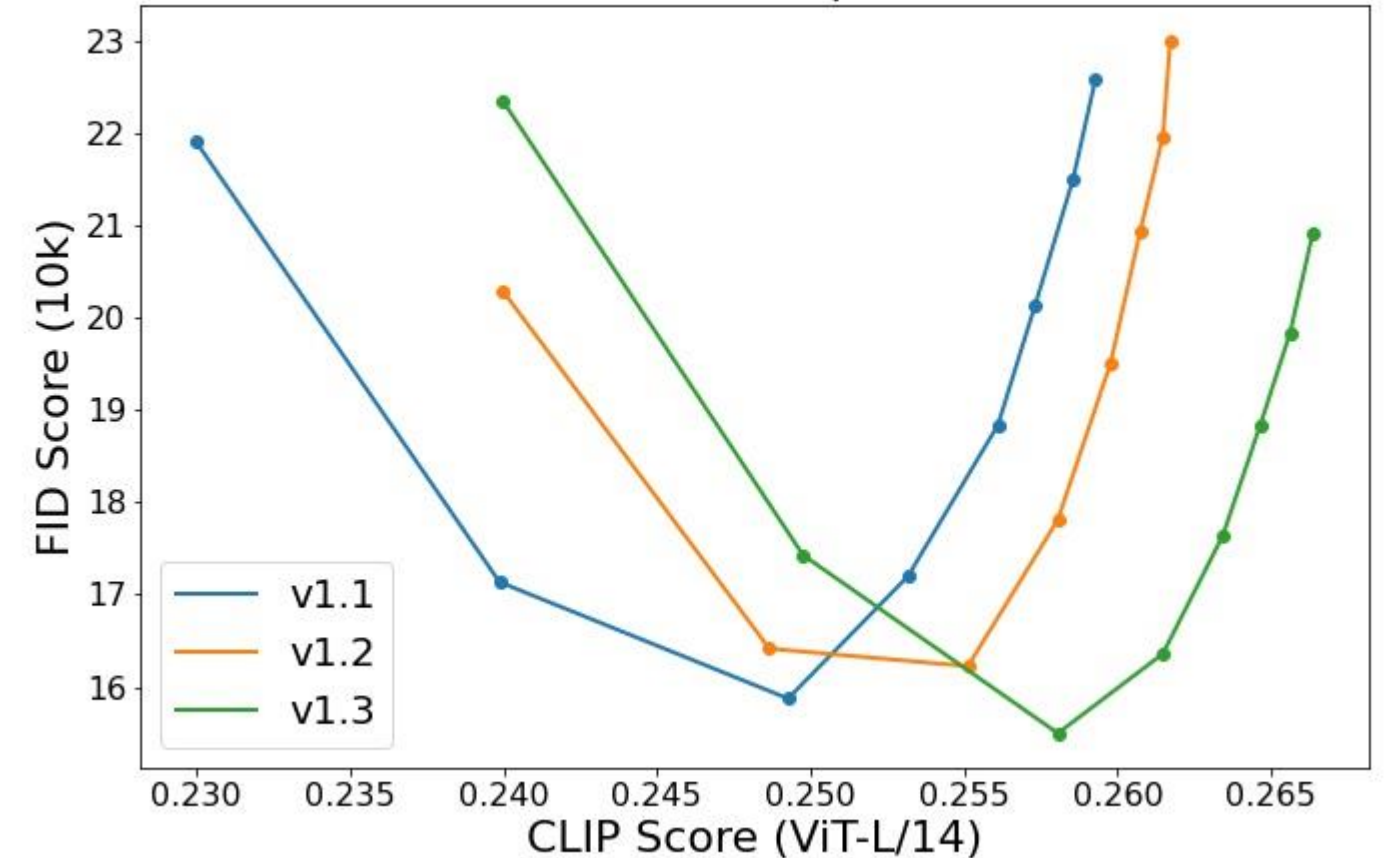
part 2 (v1.2):

- 515k steps at resolution 512x512 on "laion-improved-aesthetics" (a subset of laion2B-en, filtered to images with an original size $\geq 512 \times 512$, estimated aesthetics score > 5.0 , and an estimated watermark probability < 0.5)

part 3/4 (v1.3/v1.4):

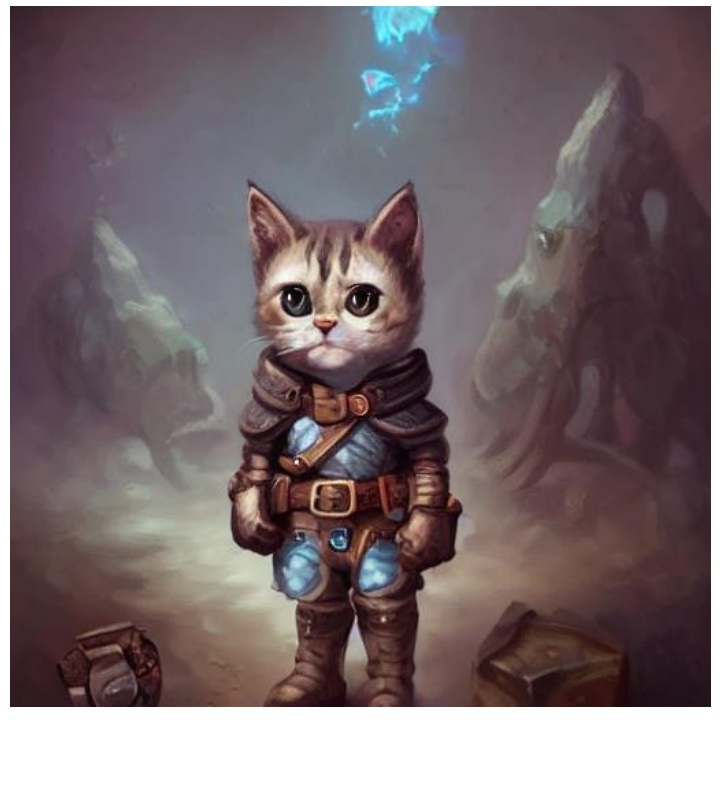
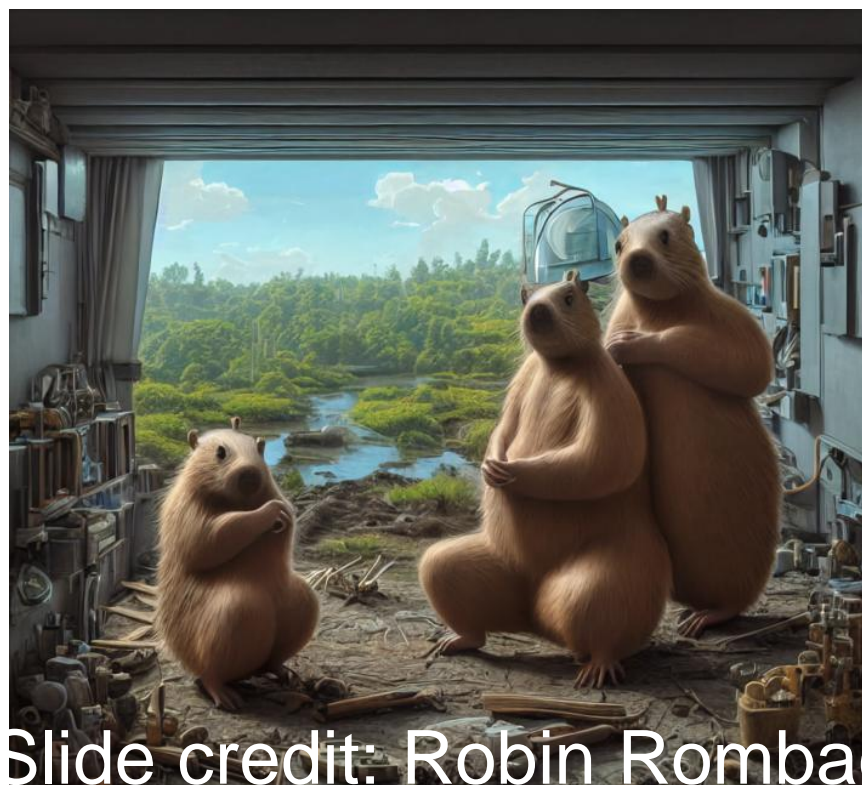
- 195k/225k steps at resolution 512x512 on "laion-improved-aesthetics" and 10% dropping of the text-conditioning

FID vs CLIP Scores on 512x512 samples for different v1-versions



10k random COCO val captions / 50 decoding steps

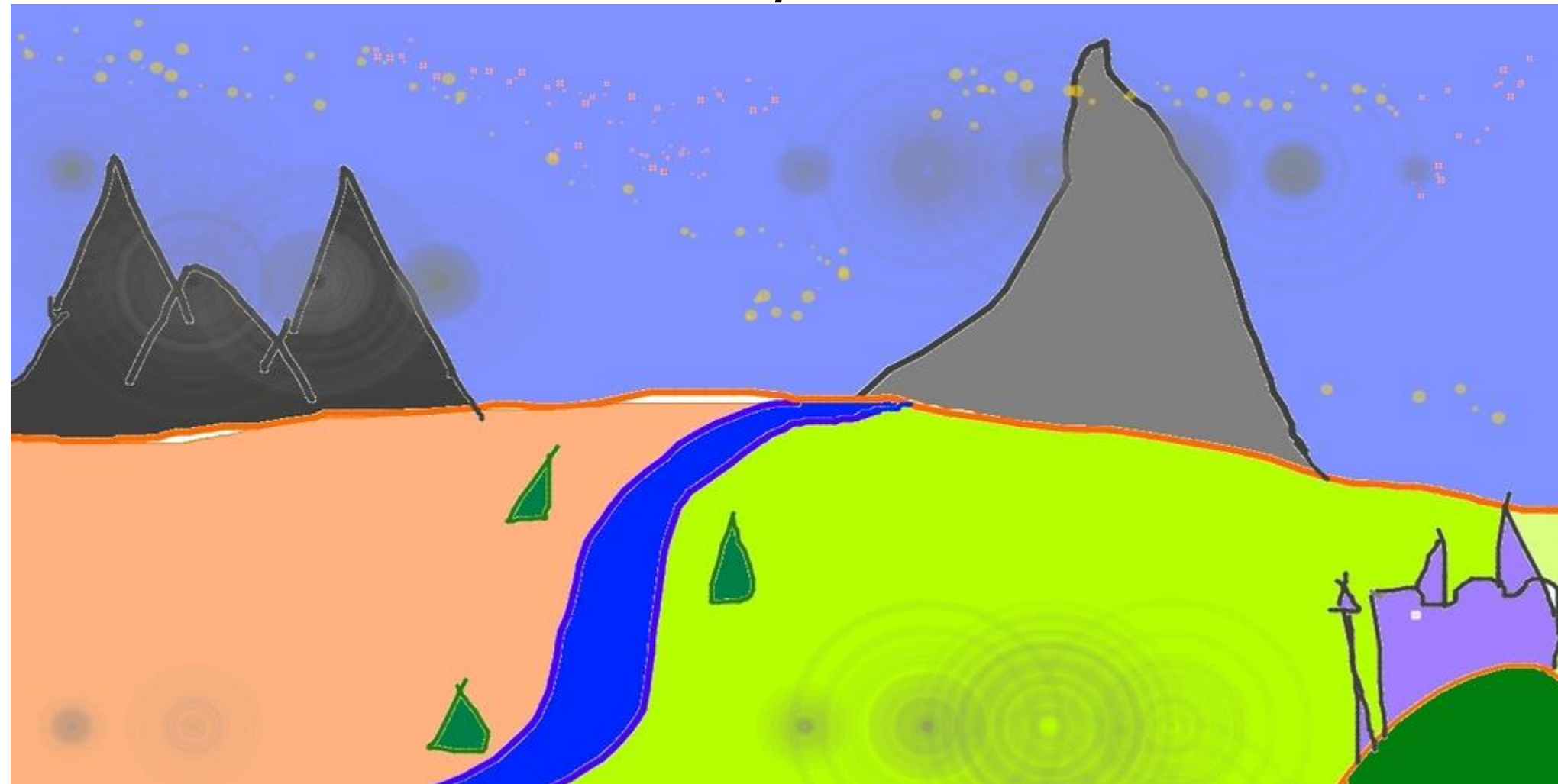
→ 4.2 GB checkpoint (EMA only, fp32)



Slide credit: Robin Rombach

Text-Guided Image-to-Image

input



“a fantasy landscape, watercolor painting”



“a fantasy landscape, trending on artstation”

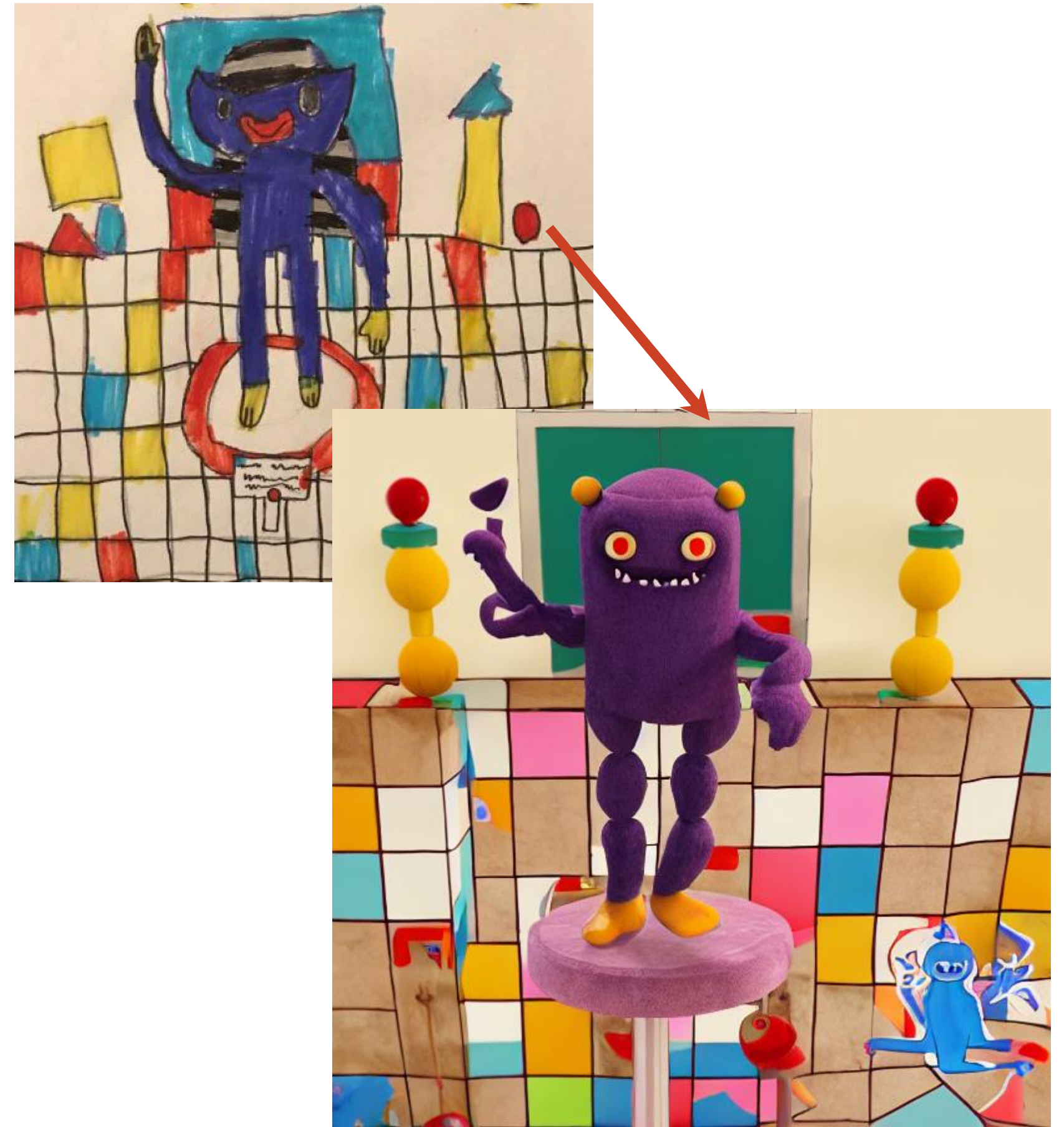


“a fantasy landscape, by Simon Stalenhag”

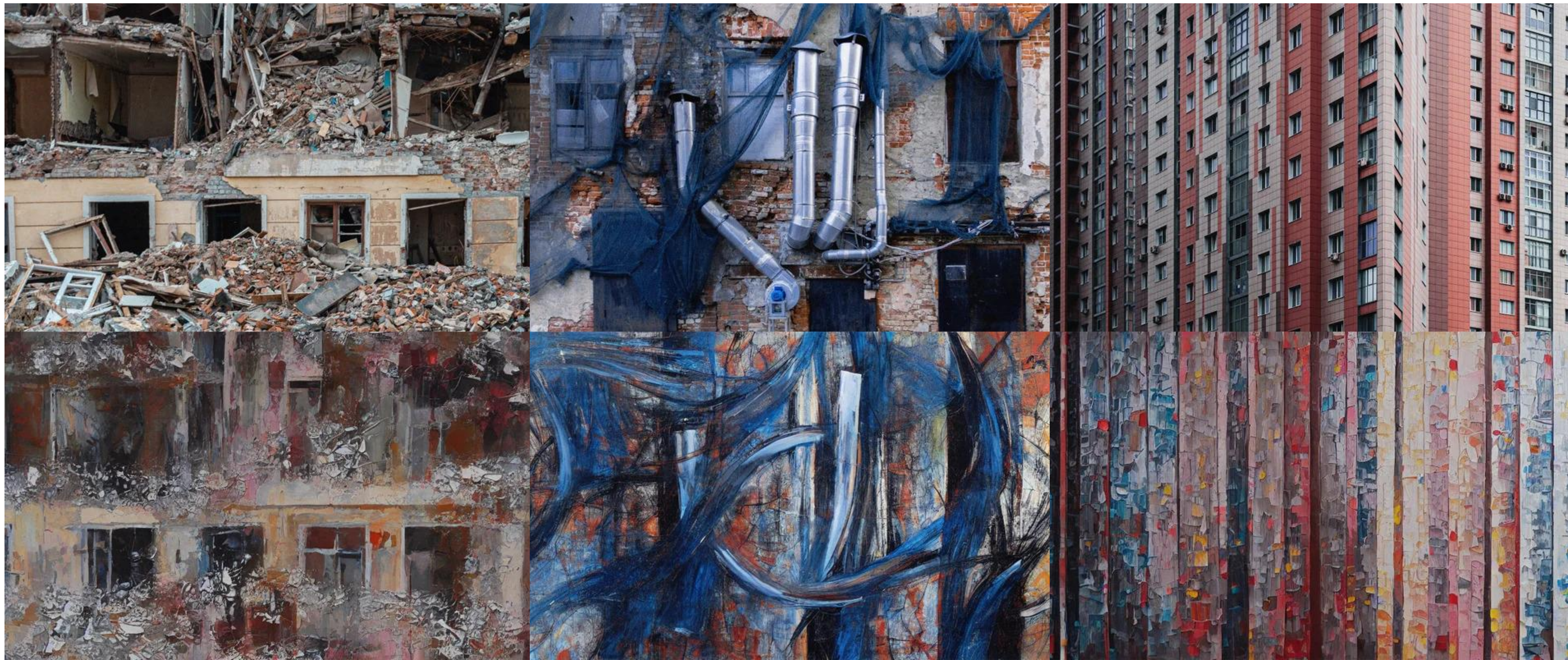


“Upgrade” your child’s artwork

original post: https://www.reddit.com/r/StableDiffusion/comments/wyq04v/using_img2img_to_upgrade_my_sons_artwork/



abstract art from photos



original post by [u/Pereulkov](https://www.reddit.com/u/Pereulkov)

https://www.reddit.com/r/StableDiffusion/comments/xhhyad/i_made_abstract_art_from_my_photos/

Video Synthesis



Stable Diffusion (img2img) + EBSynth by Scott Lightsier:

<https://twitter.com/LighthiserScott/status/1567355079228887041?t=kXXCAVtuO5IJCgro3Ma3A&s=19>

EBSynth: single-frame video stylization app: <https://ebsynth.com/>

Prompt Marketplace (promptbase.com)

DALL-E, GPT-3, Midjourney, Stable Diffusion, ChatGPT Prompt Marketplace

Find top prompts, produce better results, save on API costs, sell your own prompts.

Find a prompt

Sell a prompt



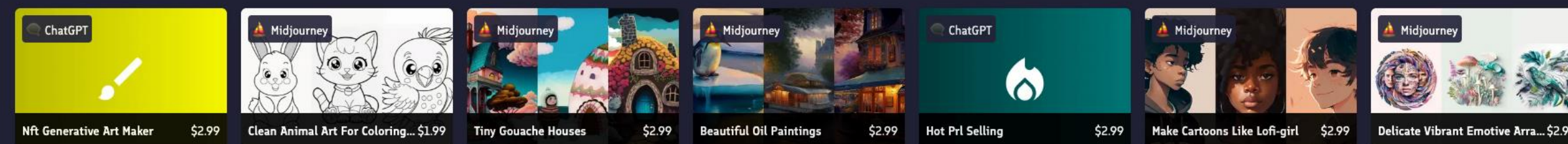
Featured in

TechCrunch THE VERGE WIRED FASTCOMPANY FINANCIAL TIMES Atlantic yahoo!finance WSJ

Featured Prompts



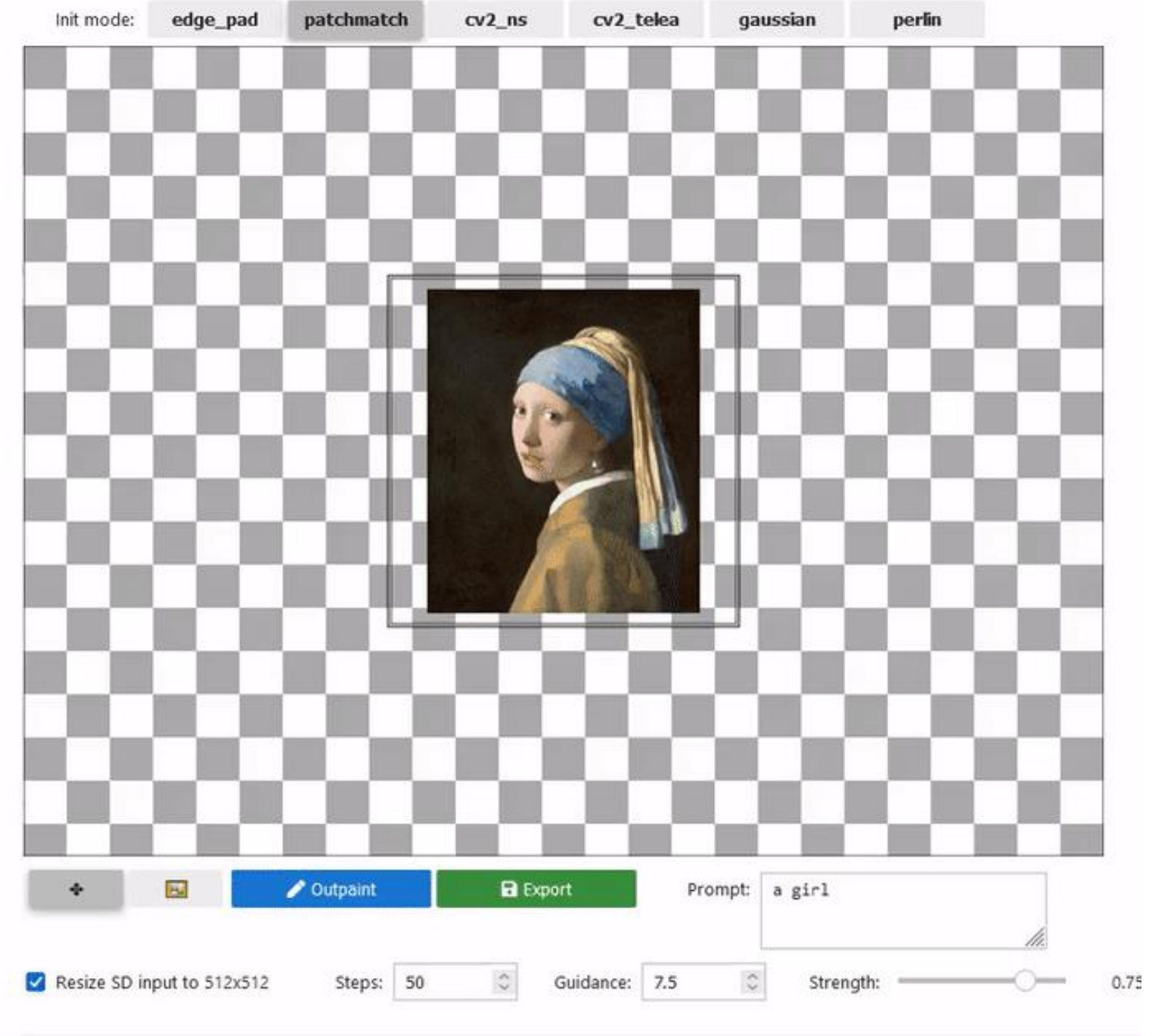
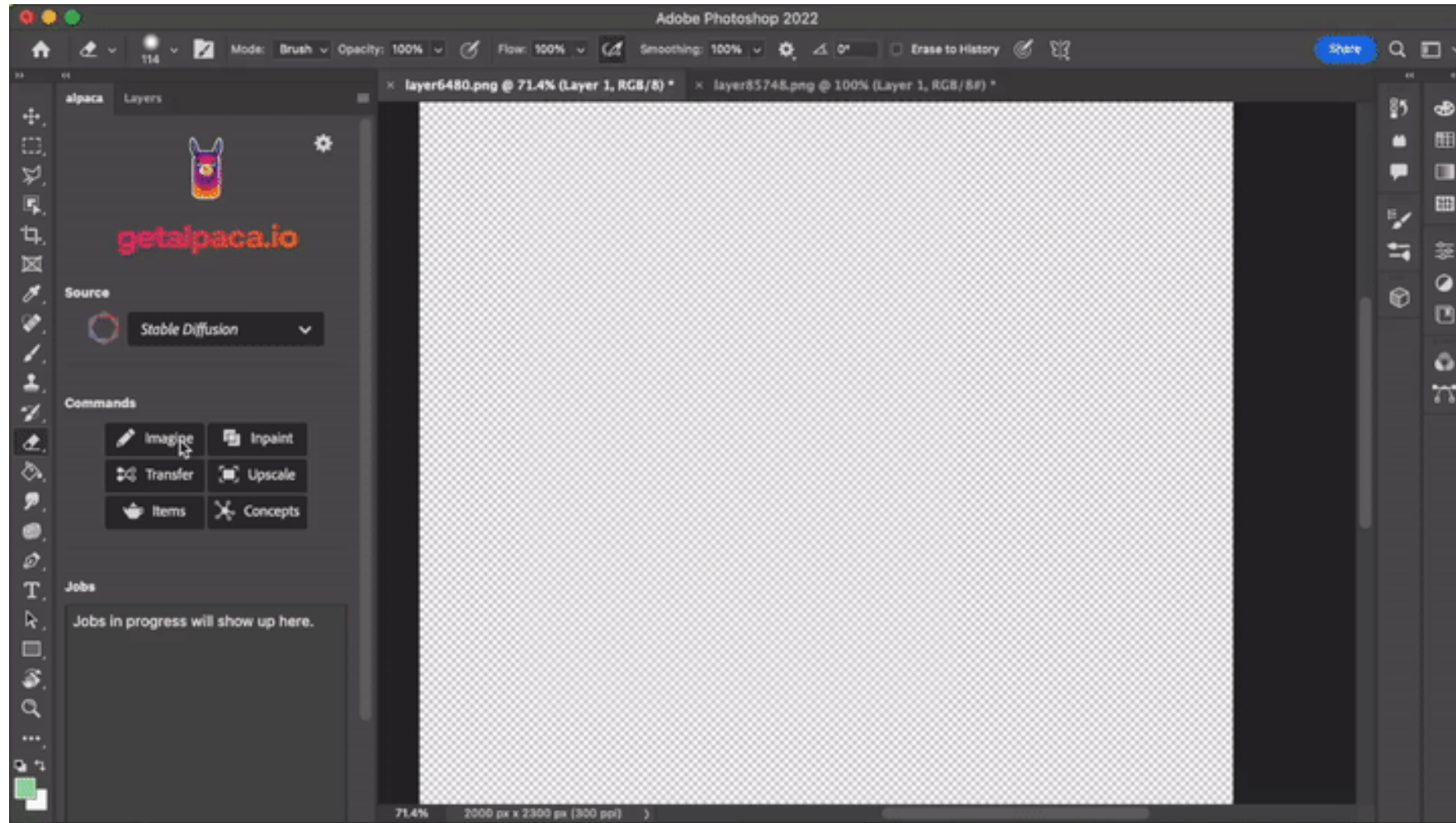
Hottest Prompts



Newest Prompts



UIs / Plug-Ins for Photoshop, GIMP etc

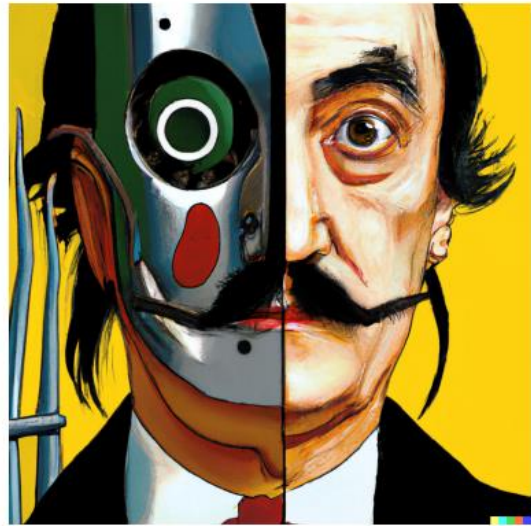


<https://twitter.com/wbuchw/status/1563162131024920576>

<https://github.com/lkwq007/stablediffusion-infinity>

What if you have 1,000+ GPUs/TPUs

DALL-E 2, Imagen



vibrant portrait painting of Salvador Dalí with a robotic half face



a shiba inu wearing a beret and black turtleneck



a close up of a handpalm with leaves growing from it



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.



an espresso machine that makes coffee from human souls, artstation



panda mad scientist mixing sparkling chemicals, artstation



a corgi's head depicted as an explosion of a nebula



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.

- Pixel-based Diffusion (No encoder-decoder)
- pre-trained text encoder (CLIP, t5)
- Diffusion model + classifier-free guidance
- Cascaded models: 64->128->512

<https://cdn.openai.com/papers/dall-e-2.pdf>

<https://arxiv.org/abs/2205.11487>

But what about ...

Evaluation?

Robustness ?

Reasoning Abilities?

Efficiency?

Do T2I Models Generate Accurate Spatial Relationships?

A chair above a knife



✘ above(chair, knife)

A teddy bear below a bed



✘ below(teddy bear, bed)

A fork to the left of a carrot



✘ left(fork, carrot)

A person to the right of a truck



✘ right(person, truck)



Benchmarking Spatial Relationships in Text-to-Image Generation

Tejas Gokhale ^{1*}
Eric Horvitz ²

Hamid Palangi ²
Ece Kamar ²

Besmira Nushi ²
Chitta Baral ¹

Vibhav Vineet ²
Yezhou Yang ¹

¹Arizona State University

²Microsoft Research

A chair above a knife



✘ above(chair, knife)

A teddy bear below a bed



✘ below(teddy bear, bed)

A fork to the left of a carrot



✘ left(fork, carrot)

A person to the right of a truck



✘ right(person, truck)



Figure 1: We benchmark T2I models on their competency with generating appropriate spatial relationships in their visual renderings. Although text inputs may explicitly mention these spatial relationships, T2I models lack such spatial understanding.

[visort2i.github.io](https://github.com/microsoft/VISOR)

<https://github.com/microsoft/VISOR>

VISOR reveals the ineffectiveness of T2I models in generating multiple objects with correct spatial relationships.

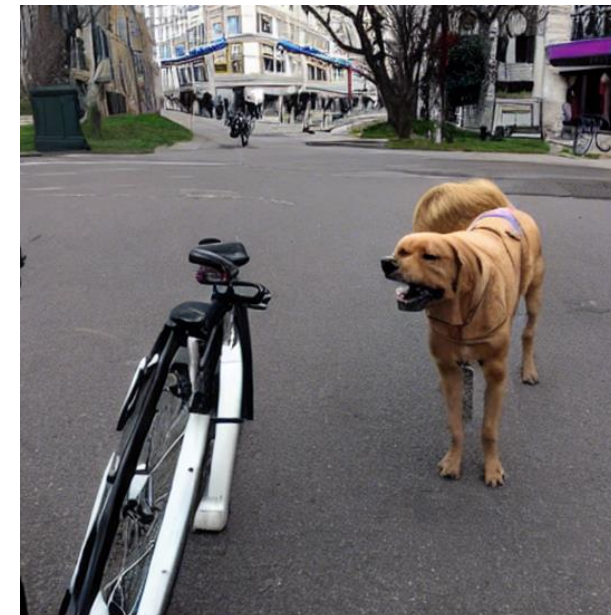
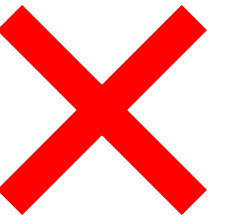
Attribute-Level Compositionality

Compositionality

an armchair in the shape of an avocado. an armchair imitating an avocado.



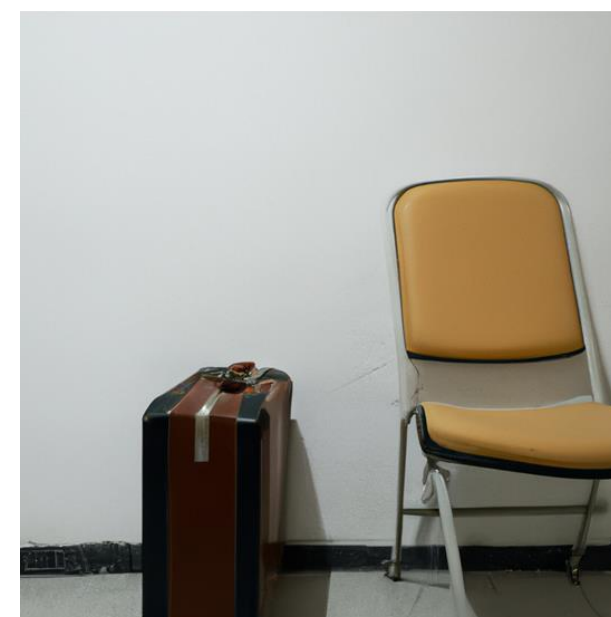
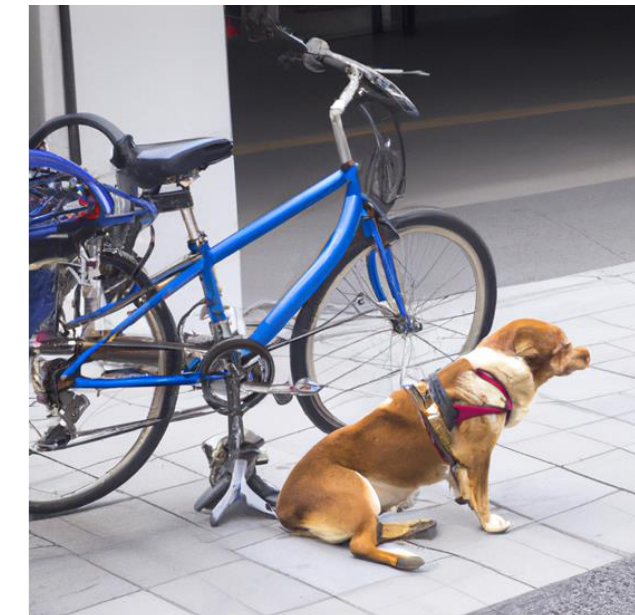
Object-Level / Spatial



"A dog to the left of a bicycle"



dog:
bicycle:
left(dog, bicycle)



"A suitcase above a chair"



suitcase:
chair:
above(suitcase, chair)



Follow-up (Method to Improve Spatial Reasoning in T2I)



Getting it *Right*: Improving Spatial Consistency in Text-to-Image Models

Agneet Chatterjee^{1,*}, Gabriela Ben Melech Stan^{2,*}, Estelle Aflalo², Sayak Paul³, Dhruba Ghosh⁴,
Tejas Gokhale⁵, Ludwig Schmidt⁴, Hannaneh Hajishirzi⁴, Vasudev Lal², Chitta Baral¹, Yezhou Yang¹

¹Arizona State University, ²Intel Labs, ³Hugging Face, ⁴University of Washington

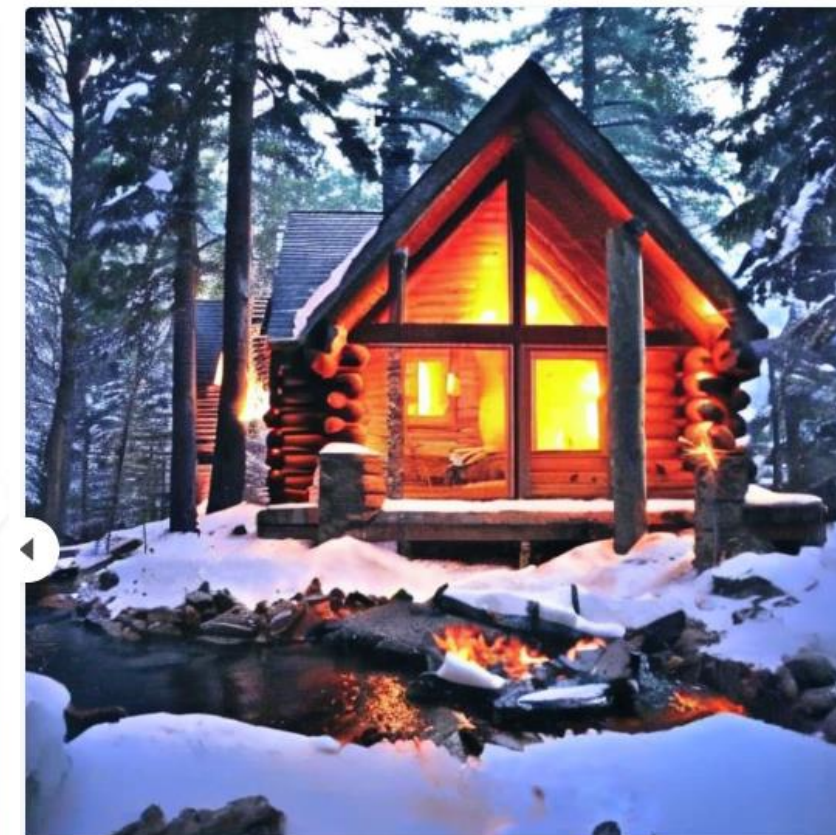
⁵University of Maryland, Baltimore County



A giraffe to the right of a truck.



A hair drier to the right of a wine glass.



A cozy cabin nestled in the woods, with a stream flowing in front and a fire burning in the fireplace inside.



A cat sitting on a chair with a lamp to the right and a window above, casting shadows on the floor below.

ConceptBed (AAAI 2024)



CONCEPTBED: Evaluating Concept Learning Abilities of Text-to-Image Diffusion Models

Maitreya Patel^{1*}, Tejas Gokhale², Chitta Baral¹, Yezhou Yang¹

¹ Arizona State University

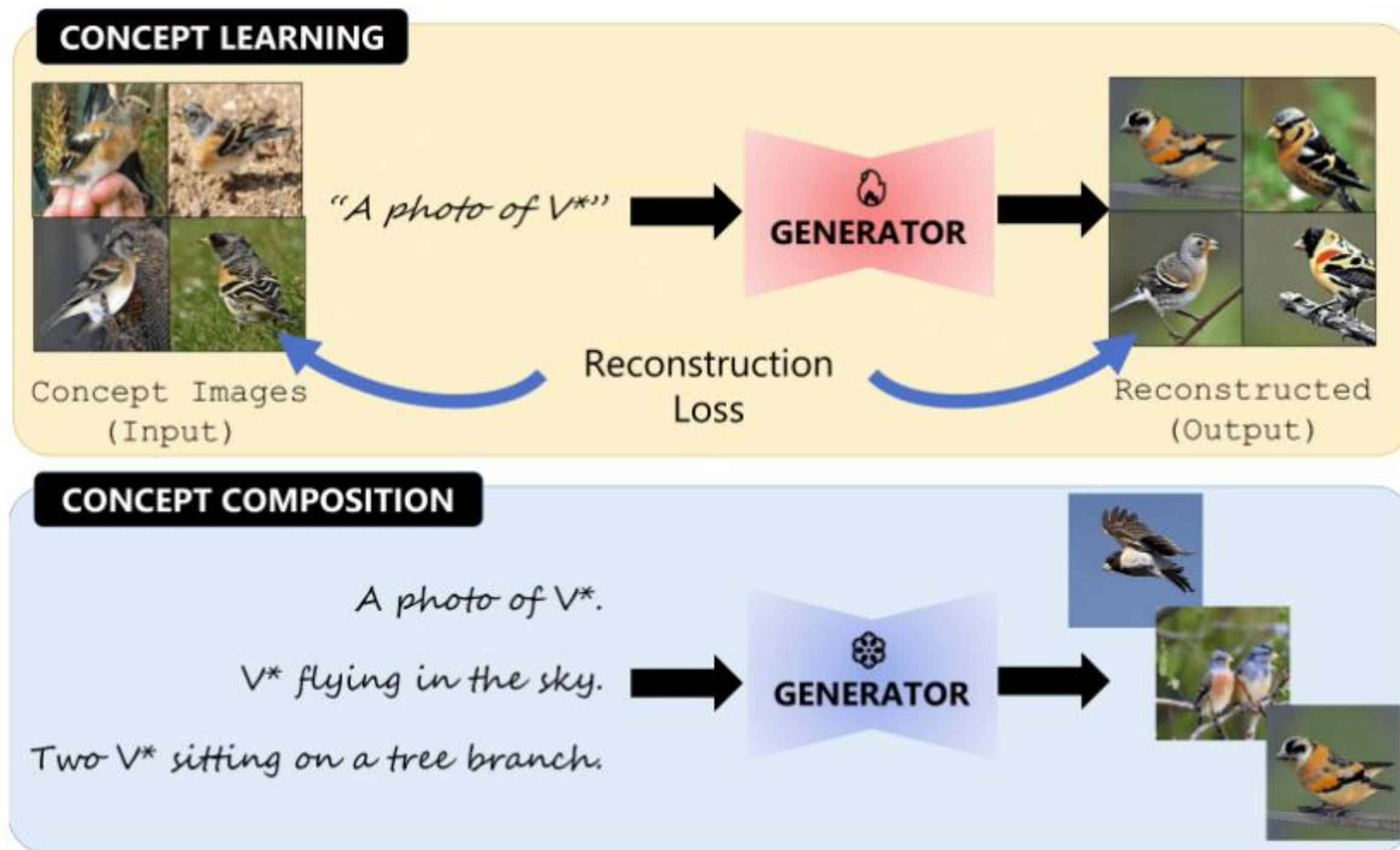
² University of Maryland Baltimore County

ConceptBed

Evaluating Concept Learning Abilities of Text-to-Image Diffusion Models

Workflow:

- Textual inversion models learn visual concepts from a few examples.
- These concepts “ V^* ” are stored as text embeddings.
- T2I models use the new concepts in novel compositions




ConceptBed

Evaluating Concept Learning Abilities of Text-to-Image Diffusion Models

Findings:

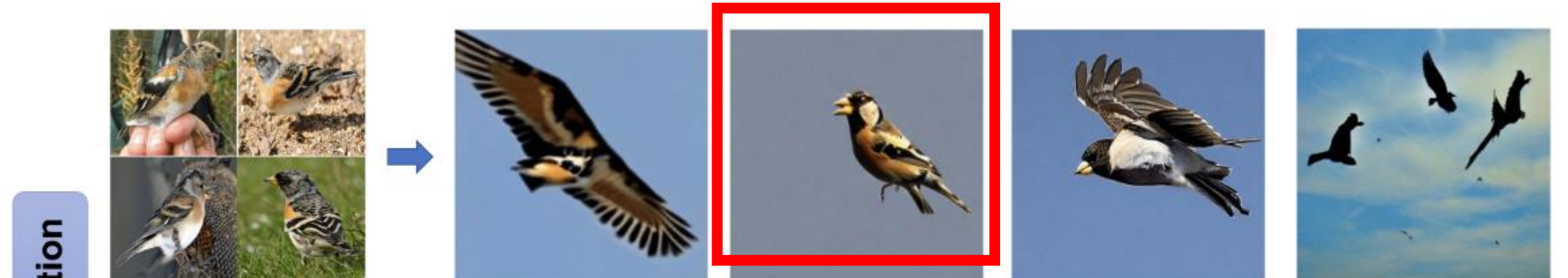
- Compositionality is hard!
- “flying”
 - where are the wings?
 - Would a bird float with that pose?
- Counting ...

Object




A photo of V*

Composition



V* flying in the sky.



Two V* sitting on a tree branch.

- Dataset
- Evaluation Metric: “Concept Confidence Deviation”

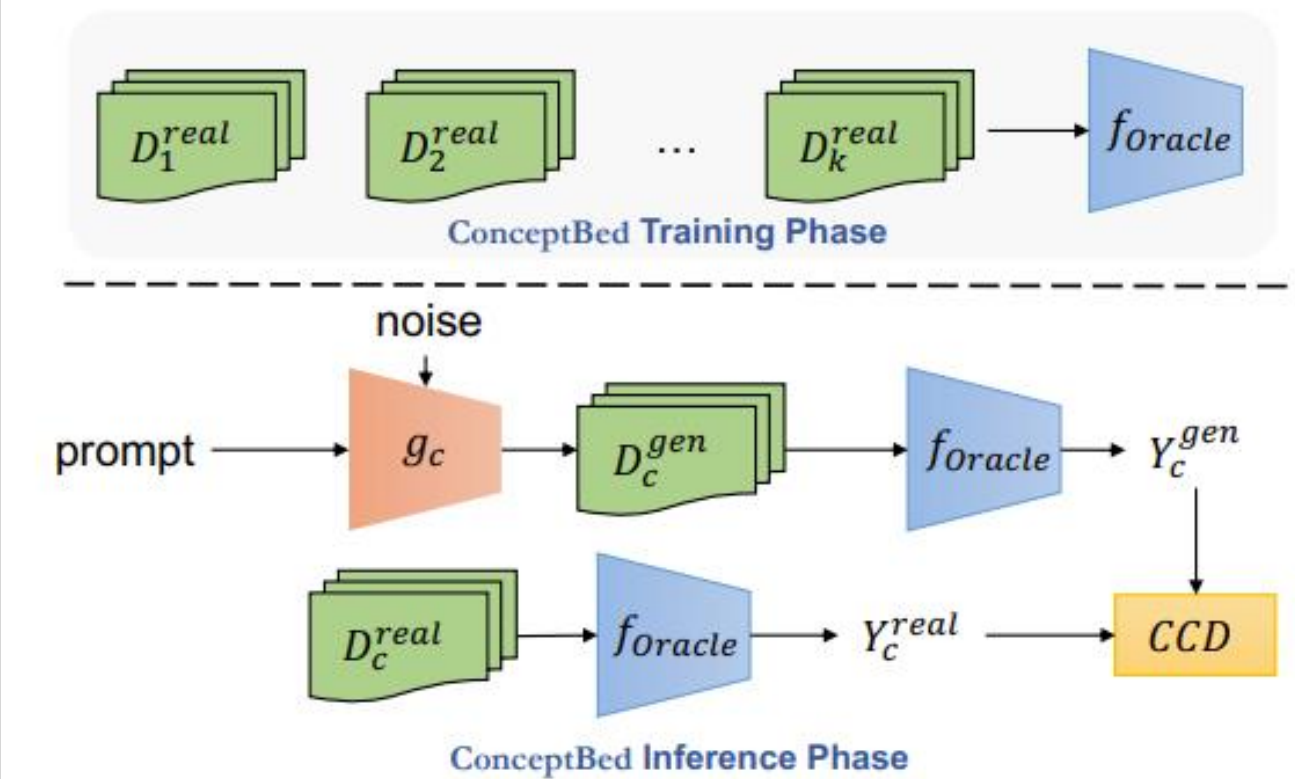
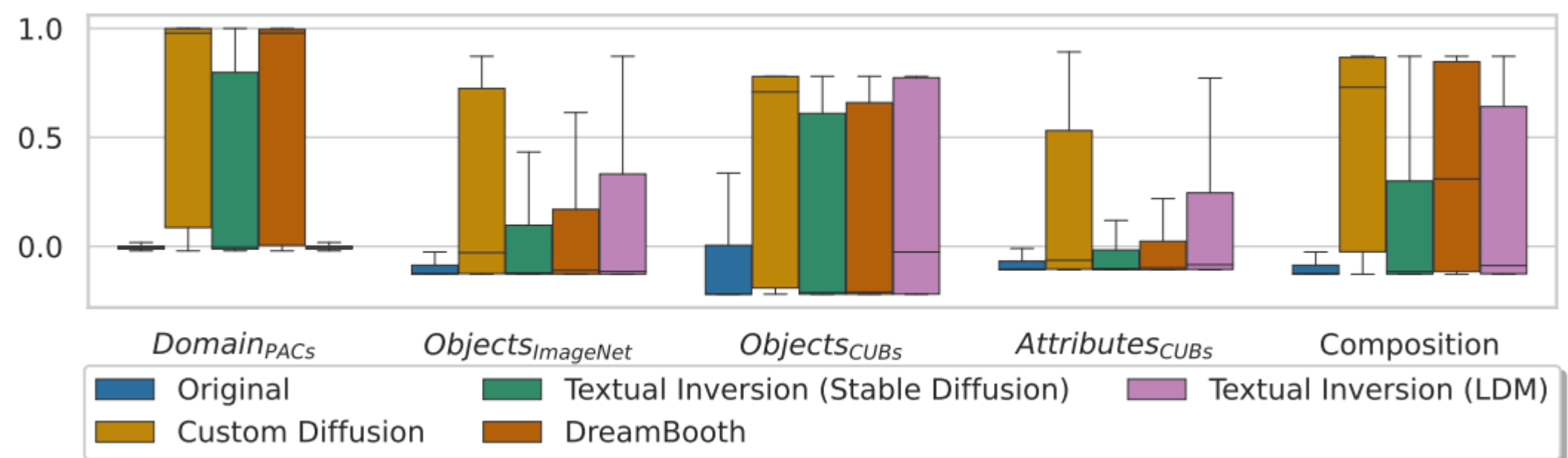
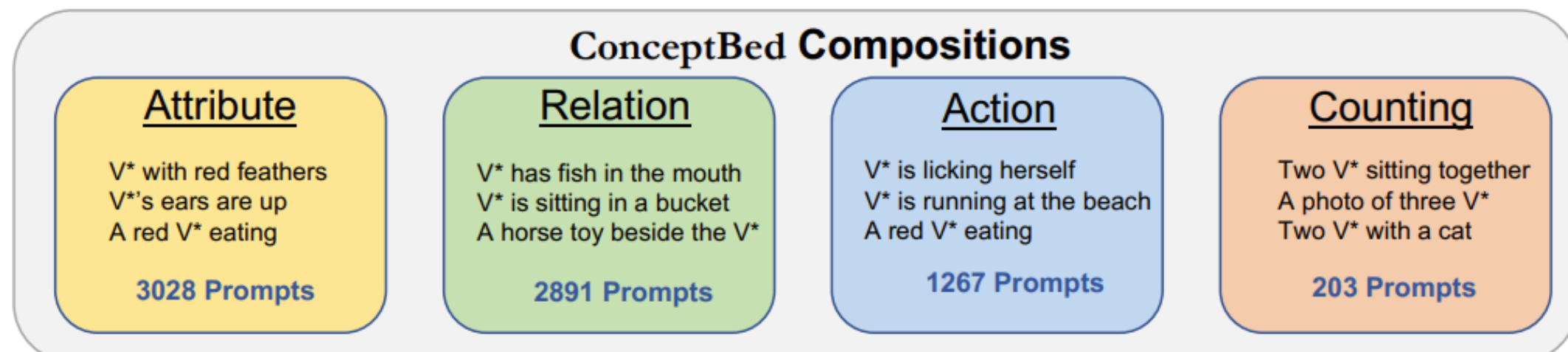
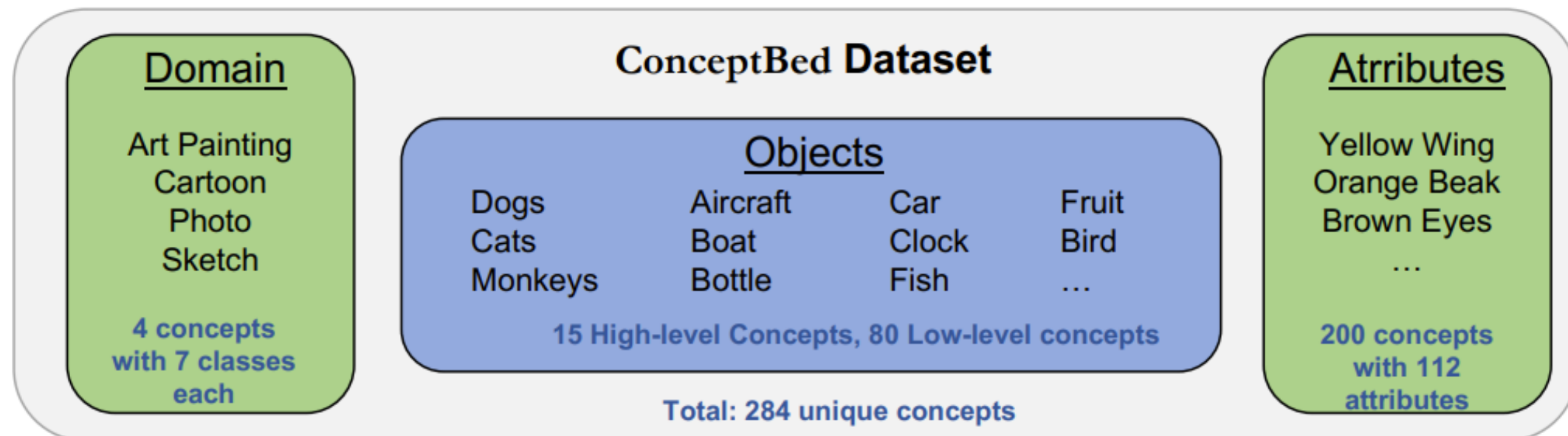
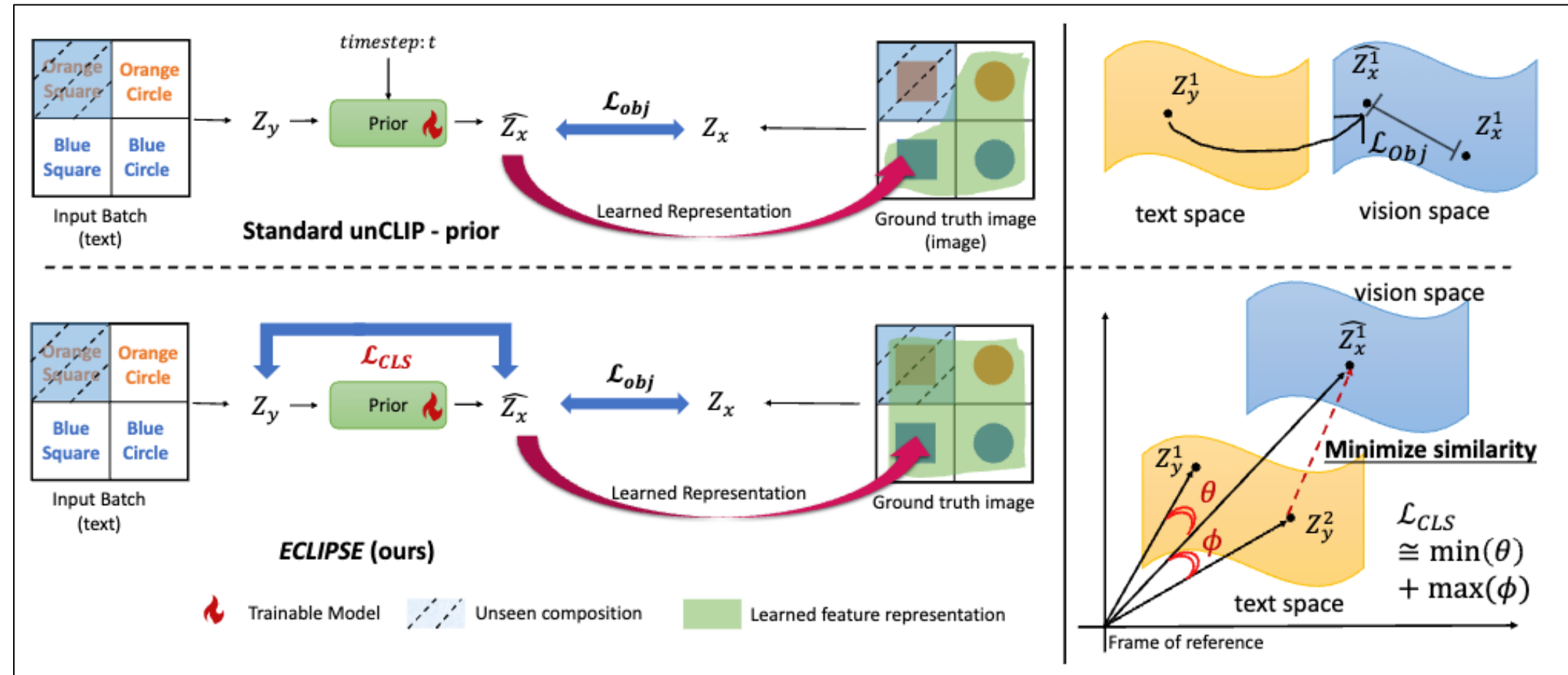
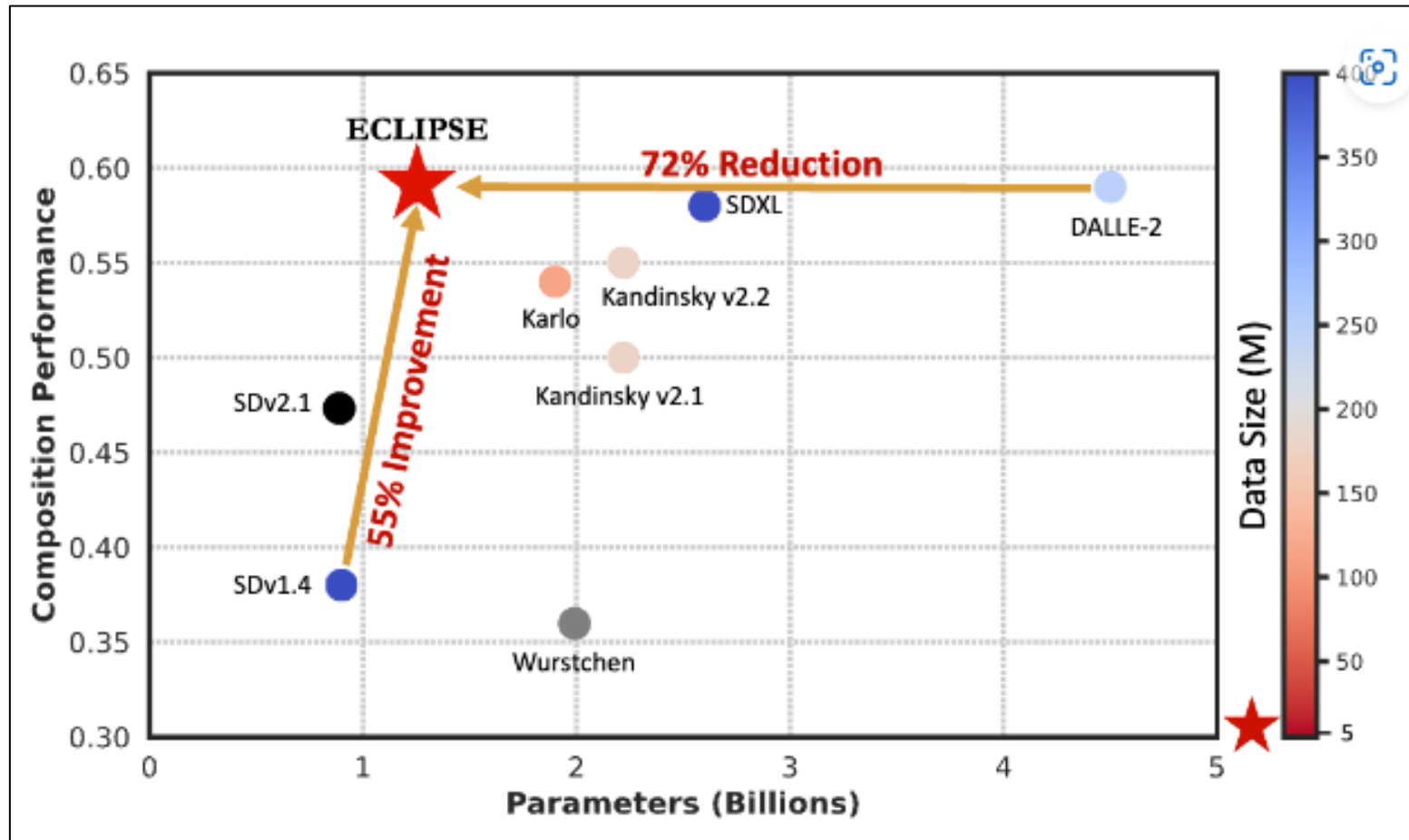


Figure 3: The outline of **CONCEPTBED** evaluation framework.

ECLIPSE: Resource-Efficient T2I Prior (CVPR'24)



- **ECLIPSE** leverages pre-trained vision-language models (e.g., CLIP) to distill the knowledge into the prior model.
- CLIP Contrastive Learning is enough to achieve state-of-the-art text-to-image prior **without the diffusion process**.
- This allows training model with only 33M parameters and 0.6M image-text pairs.