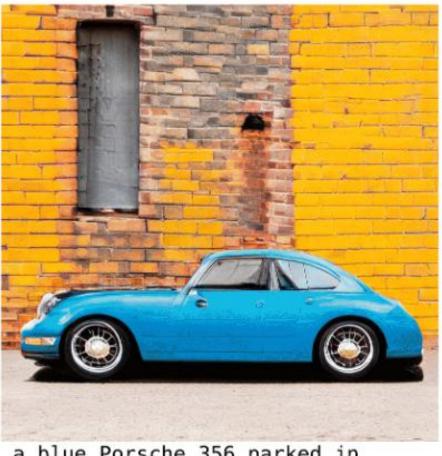
Text-to-Image (T2I)



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



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

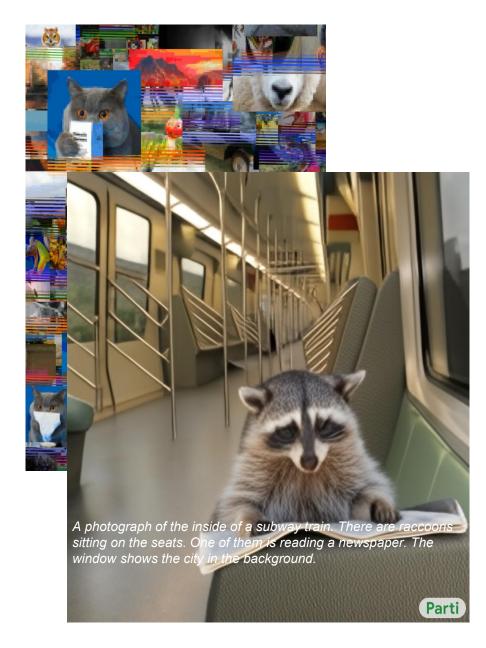


Eiffel Tower, landscape photography

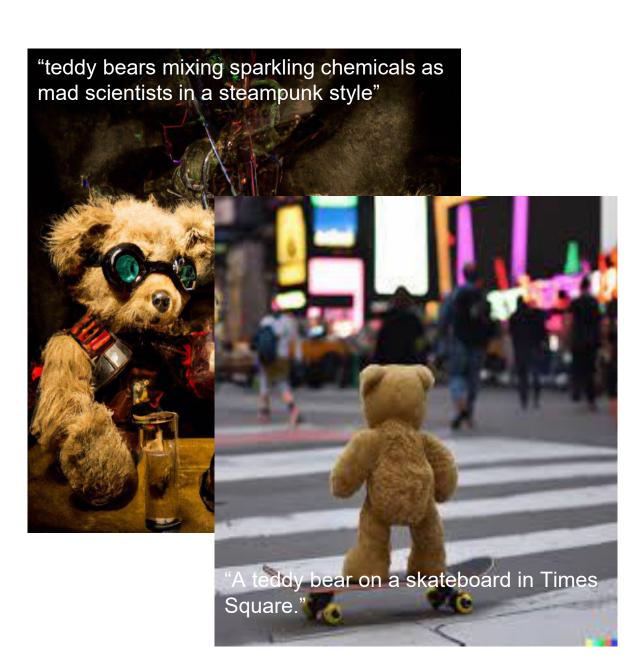


A painting of a majestic royal tall ship in Age of Discovery.

Text-to-Image Everywhere



Autoregressive models (Image GPT, Parti)

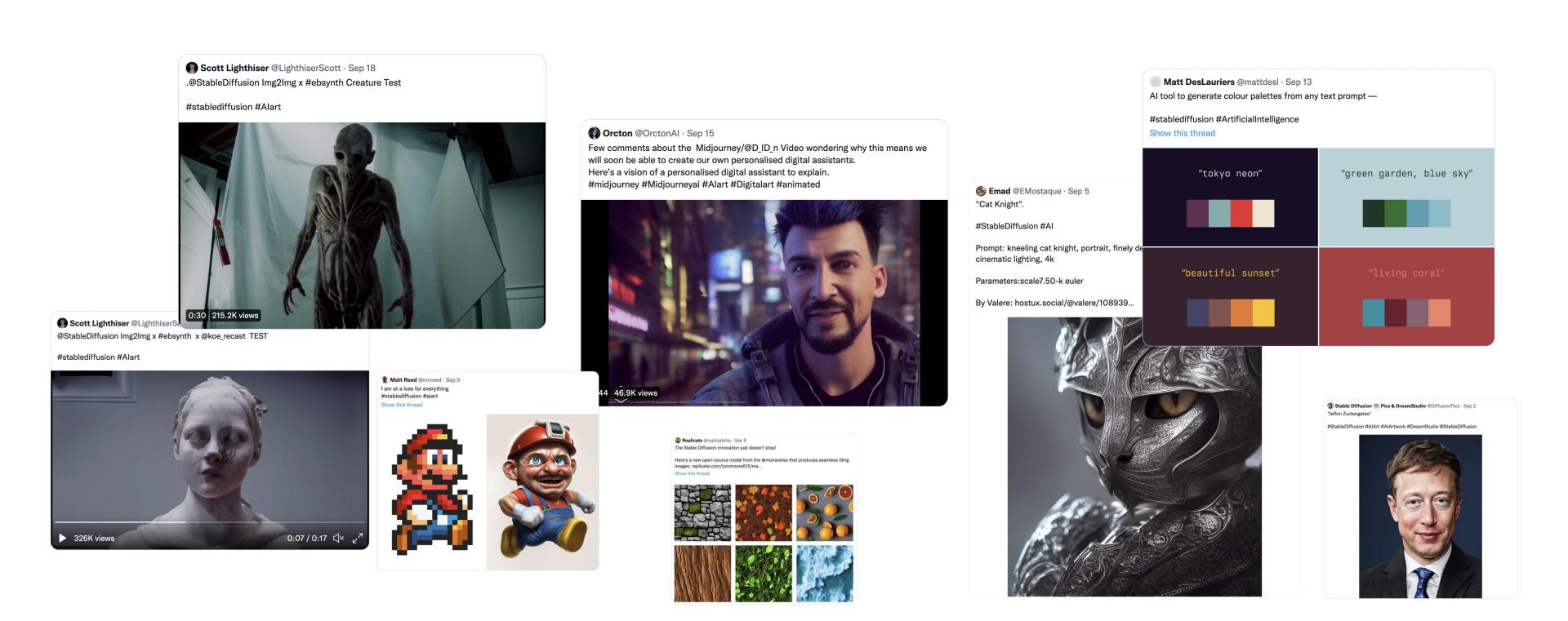


Diffusion models (DALL-E 2, Imagen)



GANs, Masked GIT (GigaGAN, MUSE)

Text-to-Image Everywhere



Where/when did it start?

First Text-to-Image System

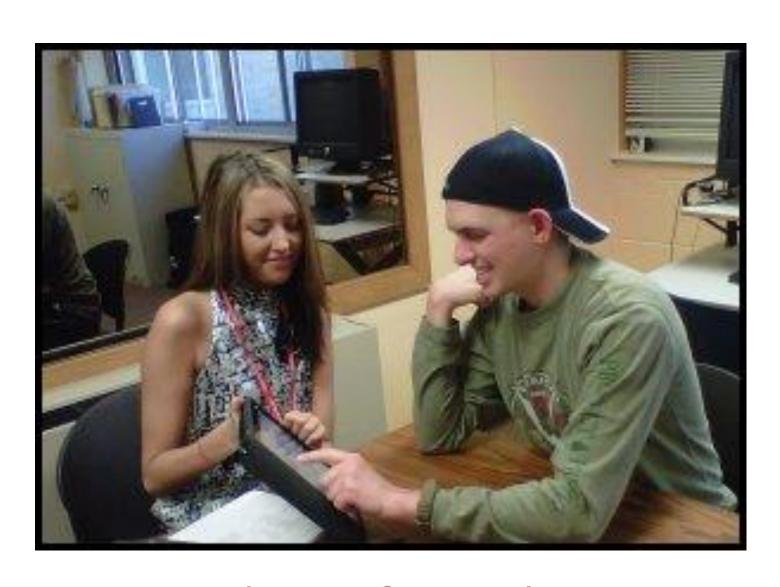


Step 1: Image Selection.

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



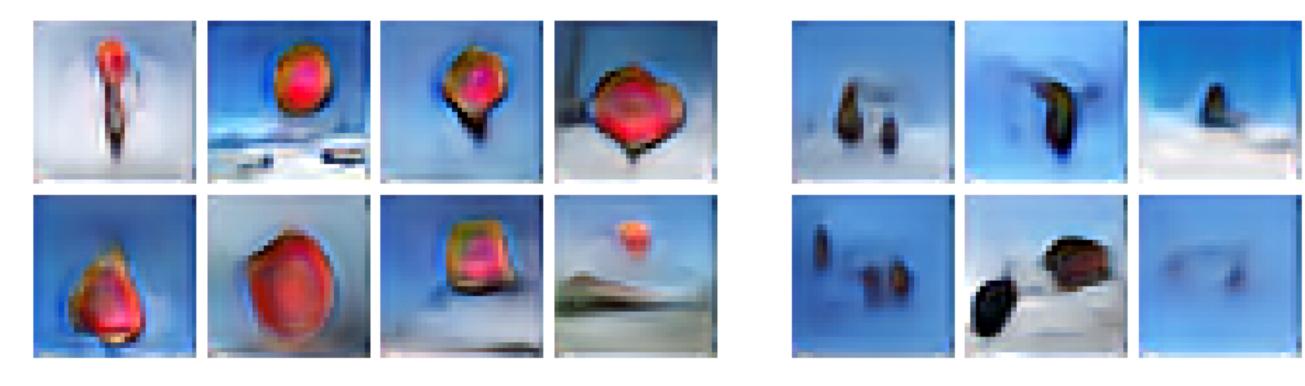
Therapy for people with communicative disorders



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 Strock. 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

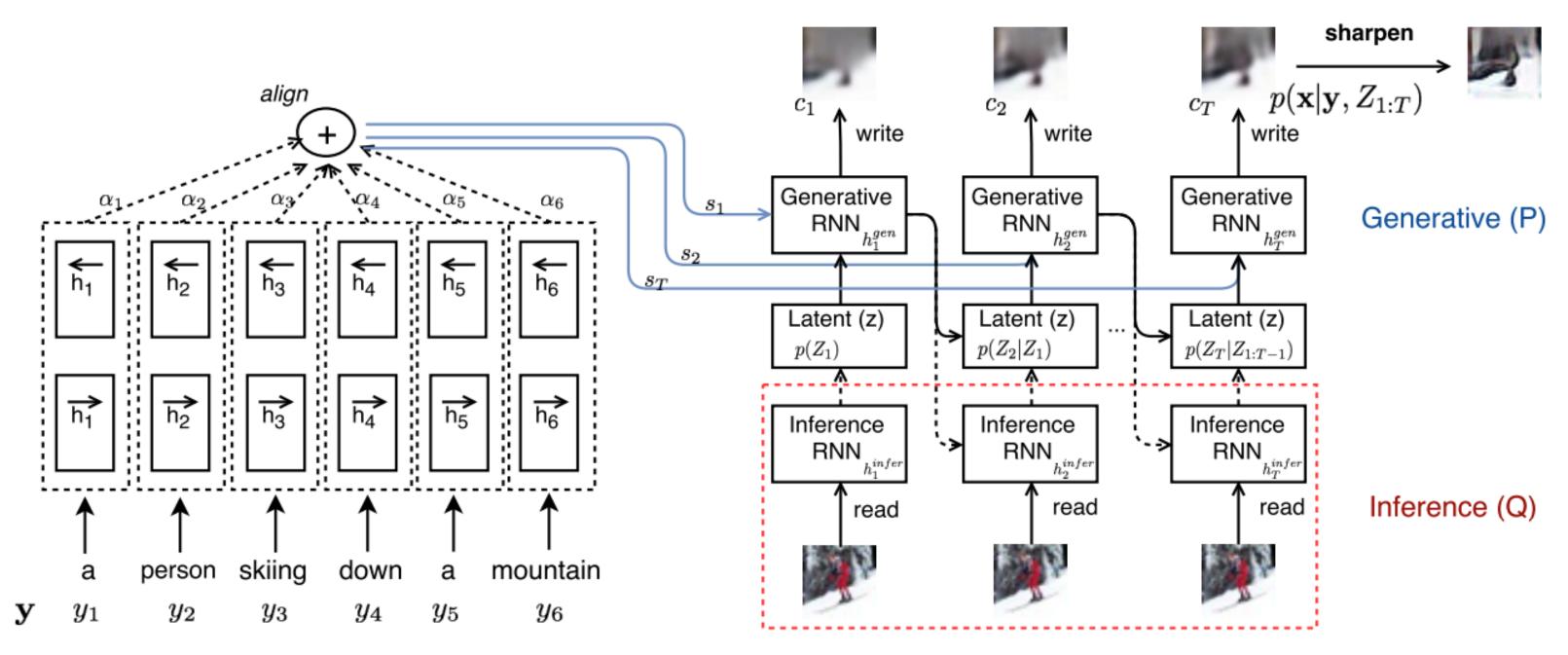


A toilet seat sits open in the grass field.

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.

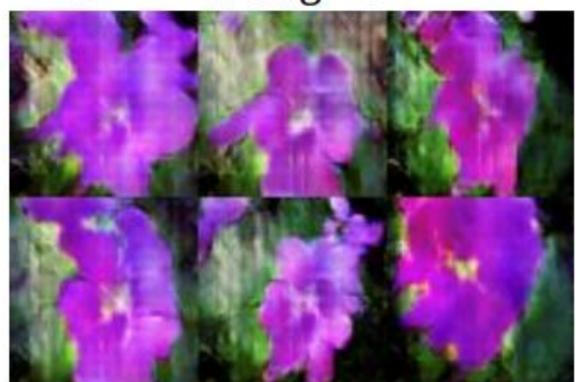




Generative Adversarial Text to Image Synthesis Scott Reed et al., ICML 2016

GAN-based Text-to-Image

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

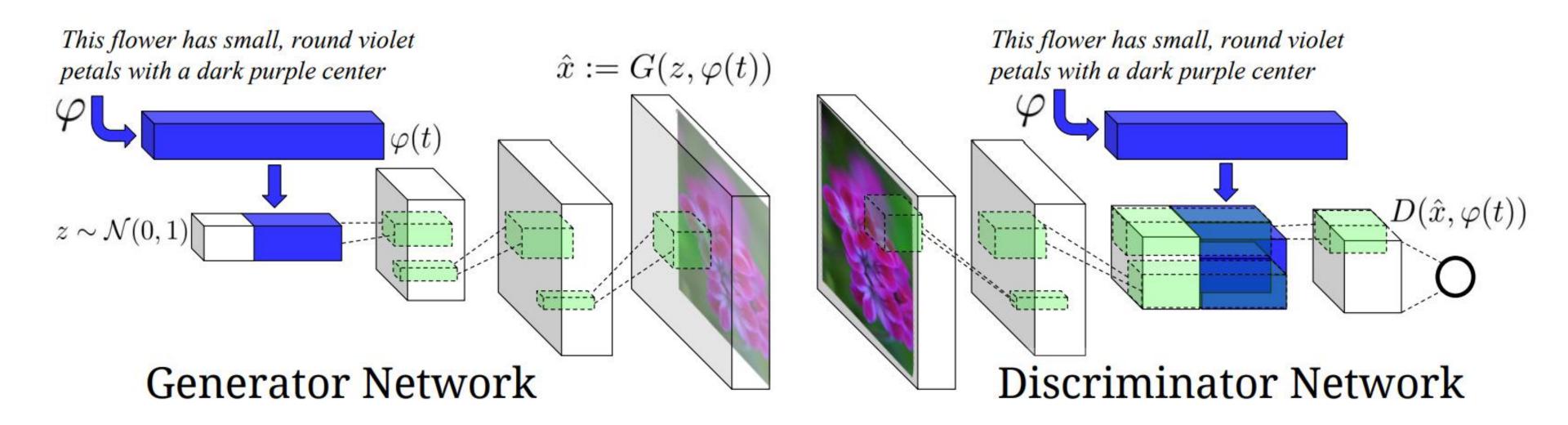


this white and yellow flower have thin white petals and a round yellow stamen



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

GAN-based Text-to-Image

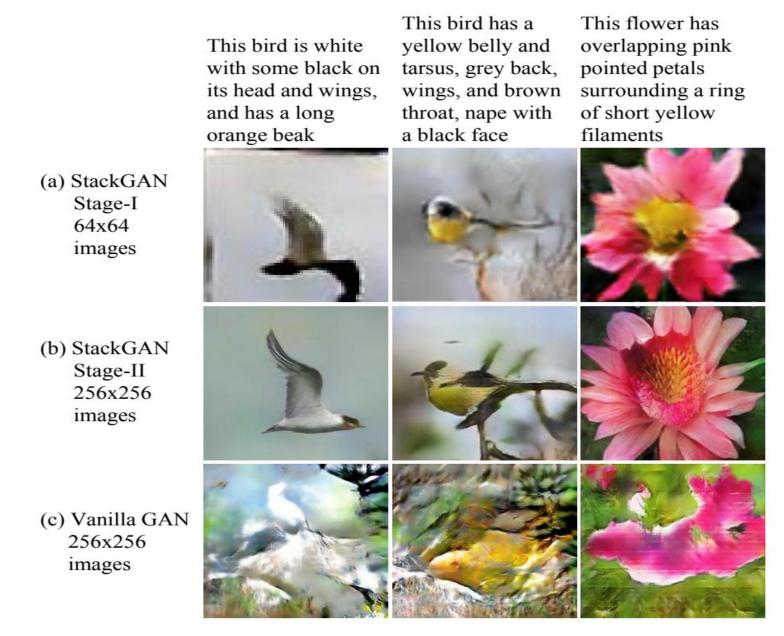


Conditional GAN + CNN + concatenation

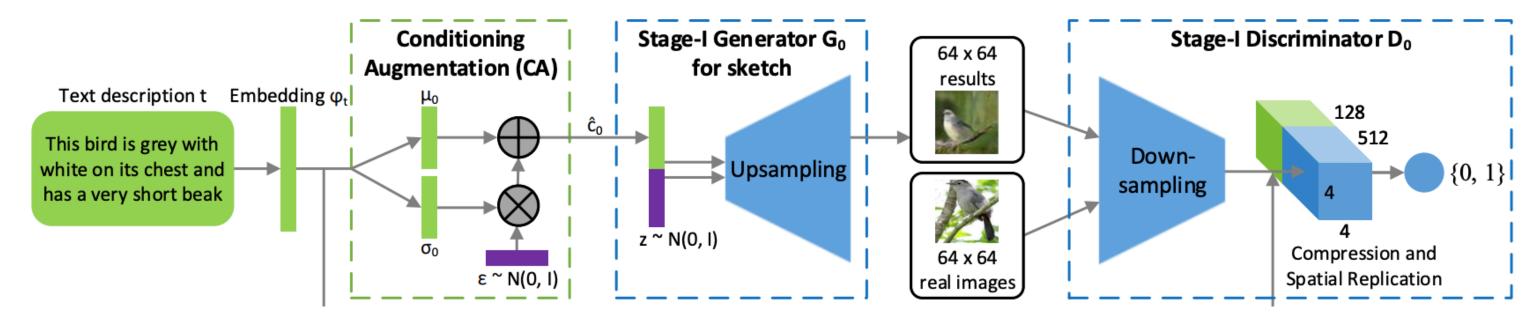
Generative Adversarial Text to Image Synthesis Scott Reed et al., ICML 2016



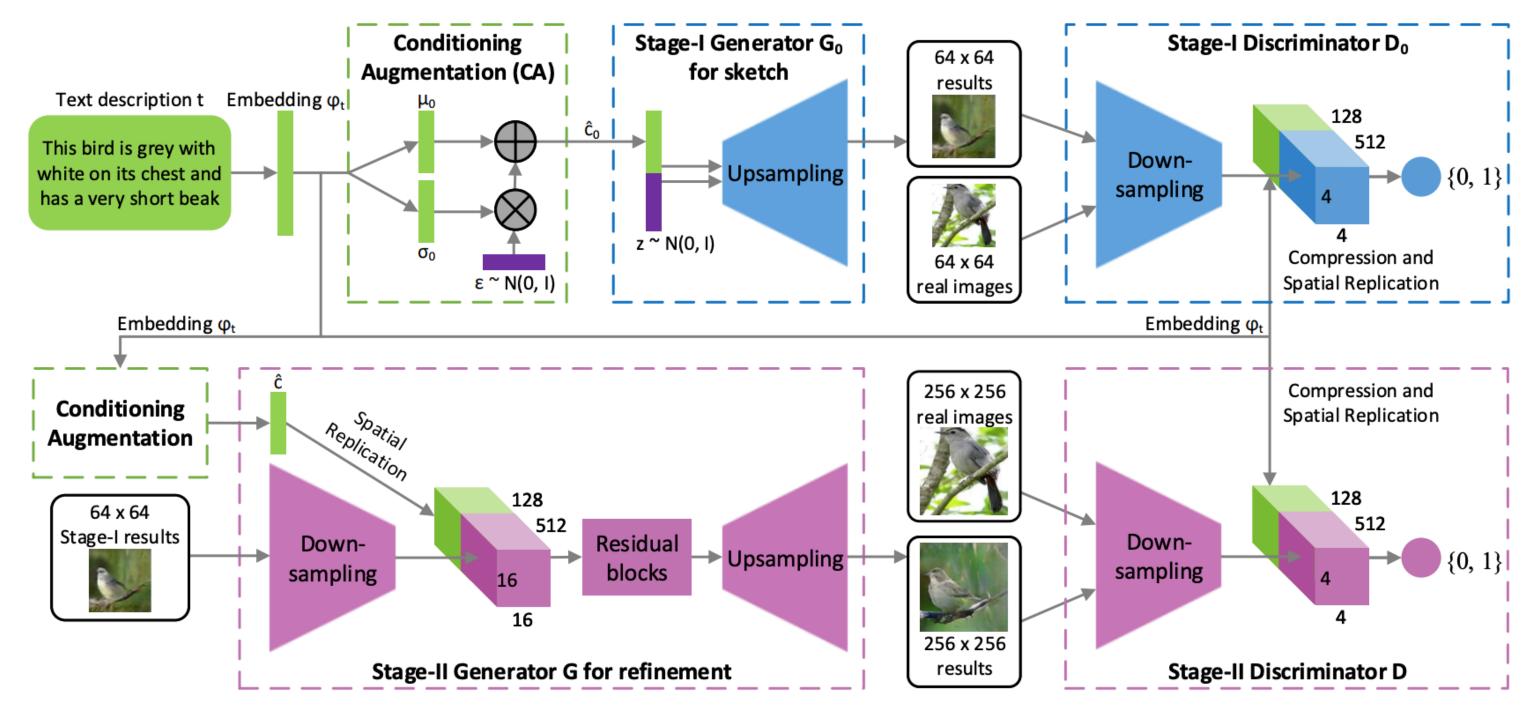
But these images are tiny ... How can we make them HD?



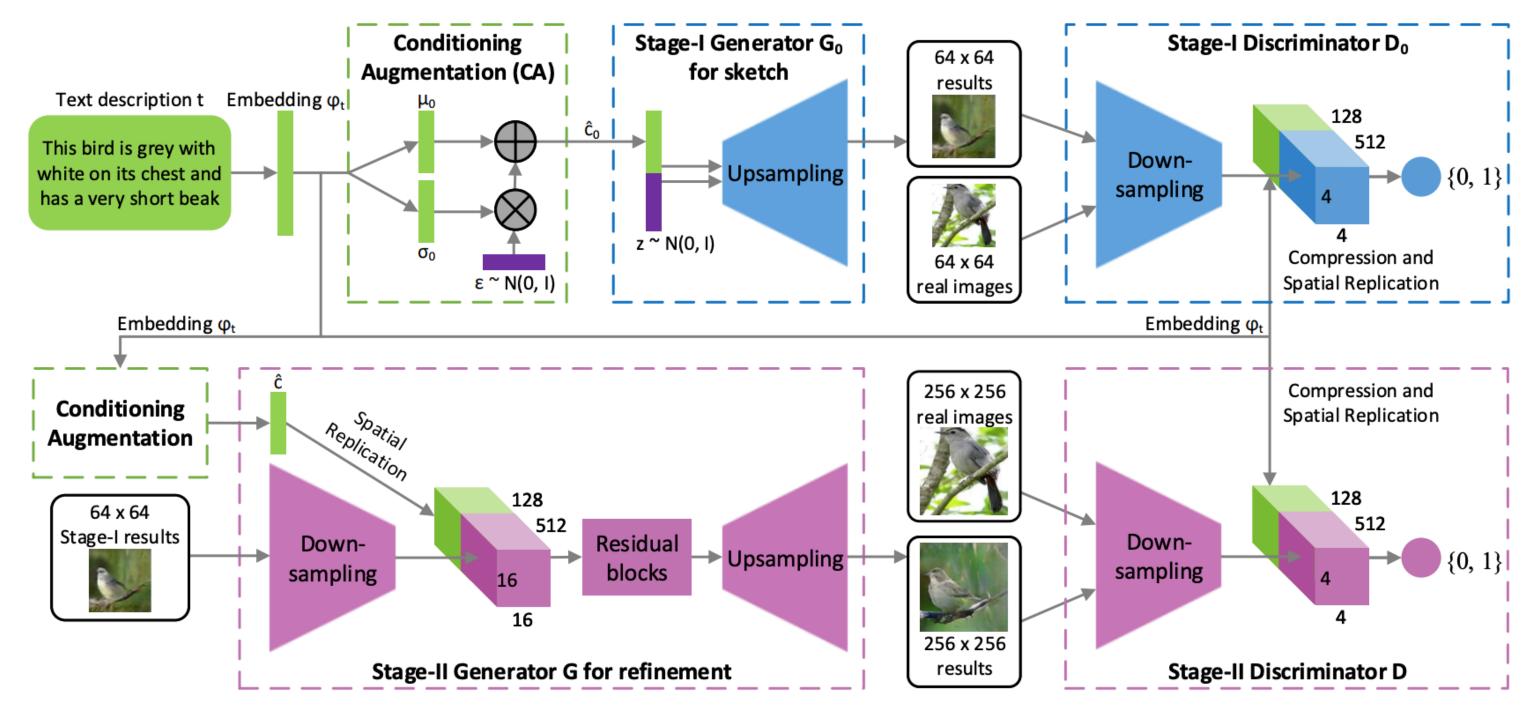
Two-stage Conditional GAN + CNN + concatenation



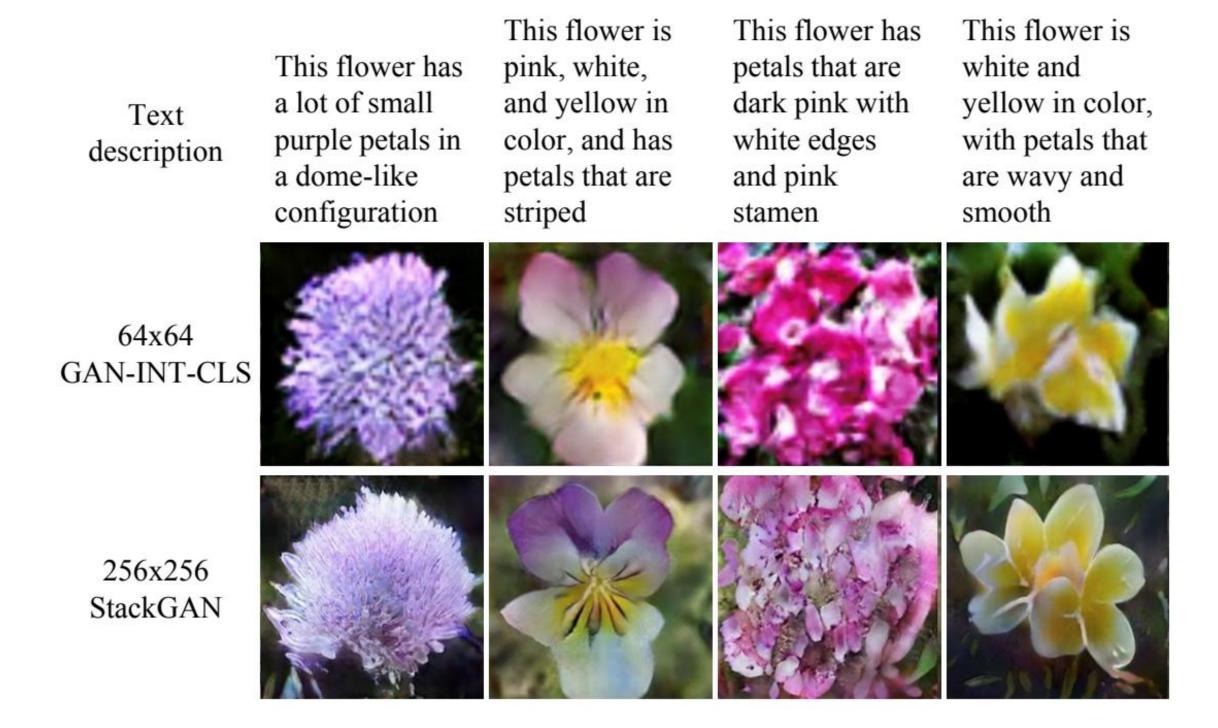
Two-stage Conditional GAN + CNN + concatenation



Two-stage Conditional GAN + CNN + concatenation



Two-stage Conditional GAN + CNN + concatenation

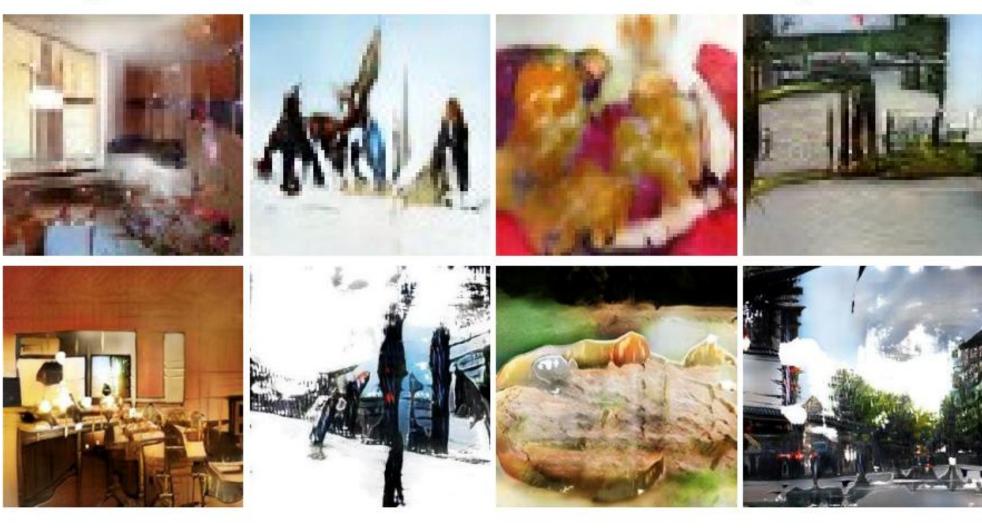


Text description

A picture of a very clean living room

A group of people on skis stand in the snow

256x256 StackGAN



Eggs fruit

candy nuts

and meat

served on

white dish

A street sign

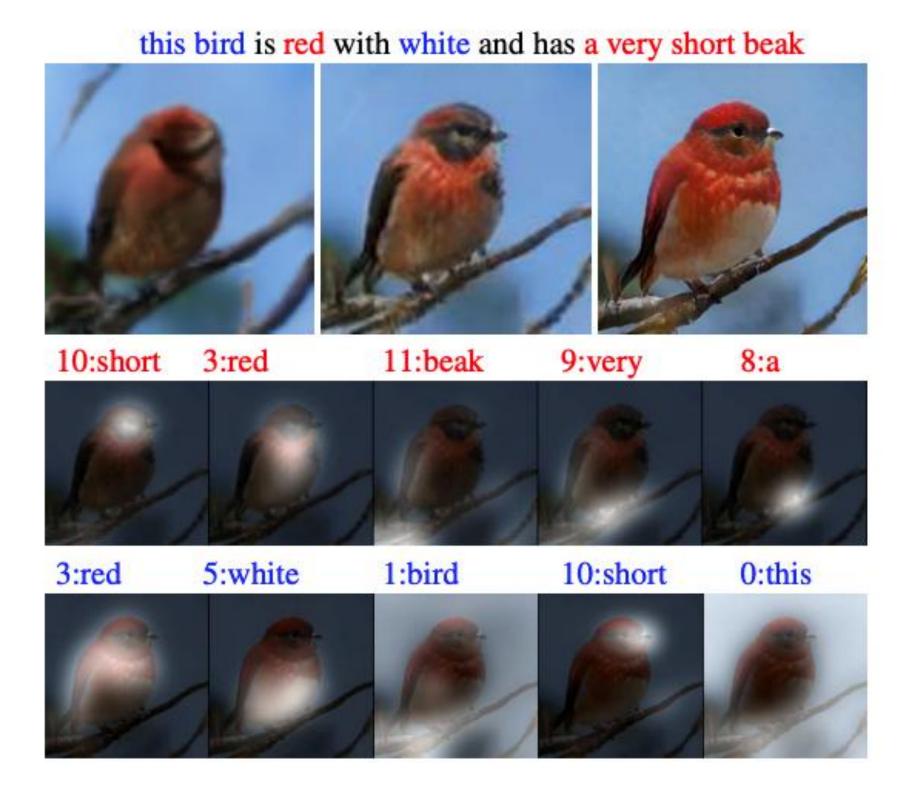
on a stoplight

pole in the

middle of a

day

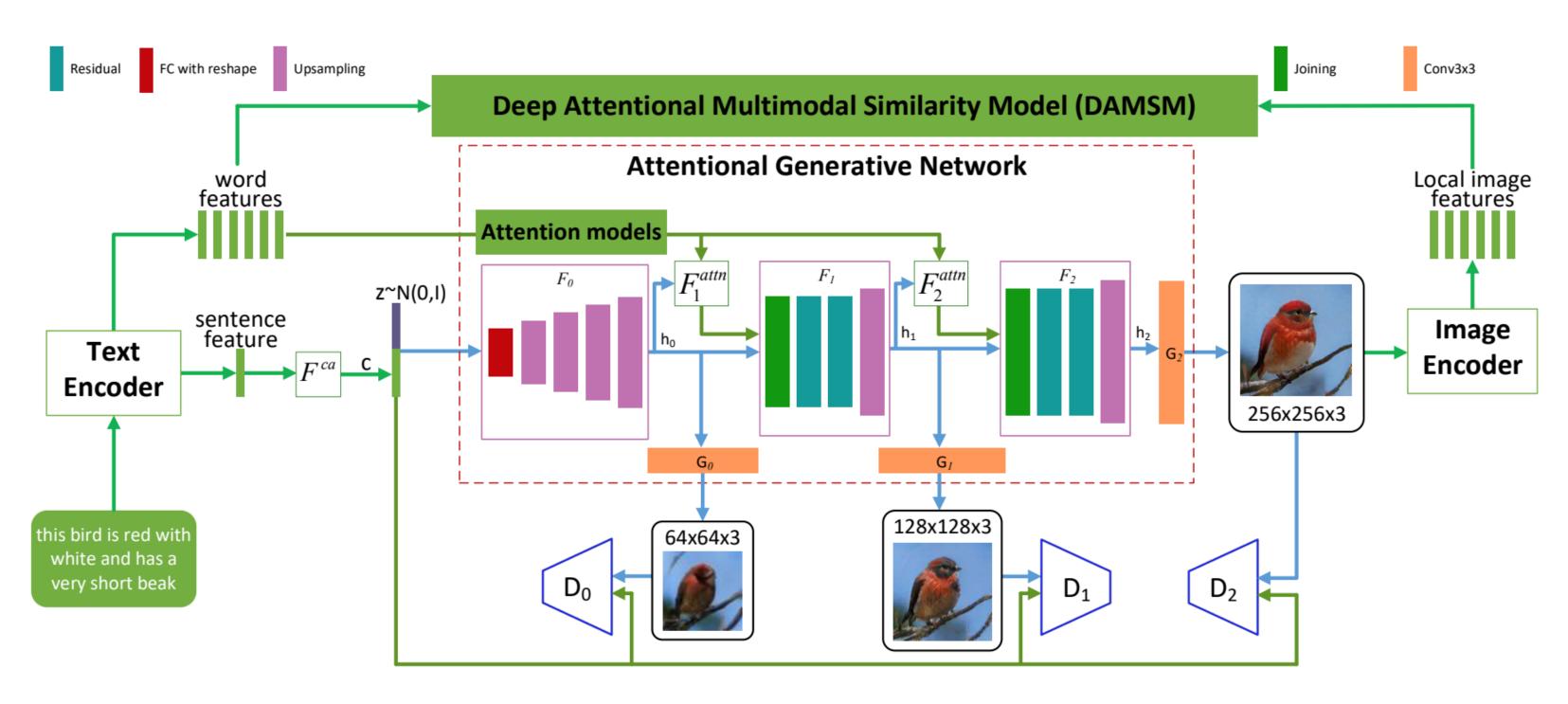
+ Cross-attention to connect Text and Image



AttnGAN: Fine-Grained Text to Image Generation with Attentional Generative Adversarial Networks

Tao Xu et al., CVPR 2018

+ Cross-attention to connect Text and Image

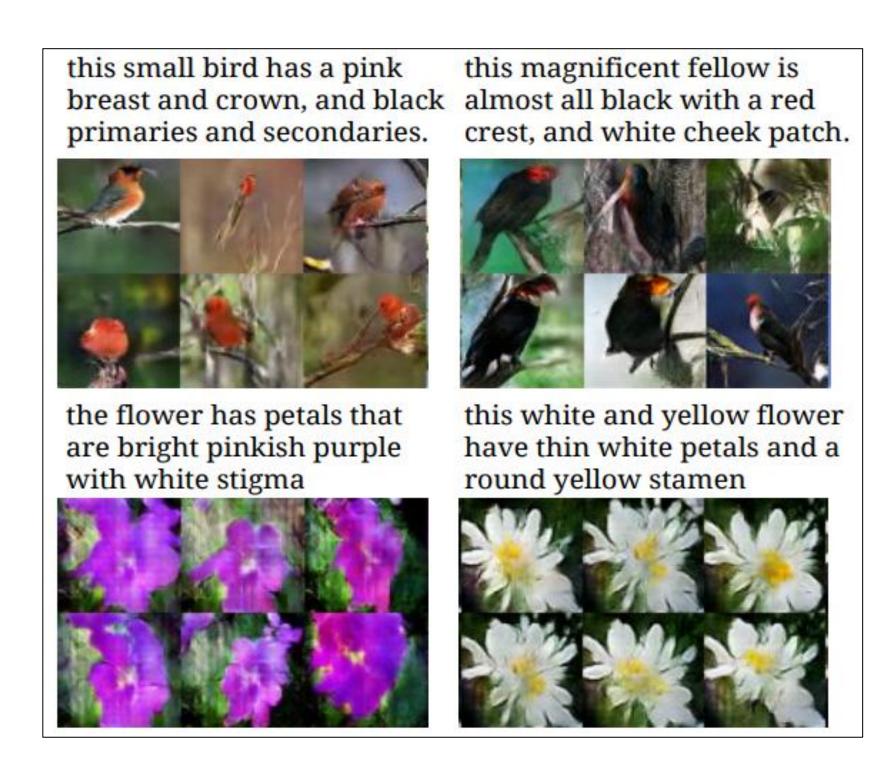


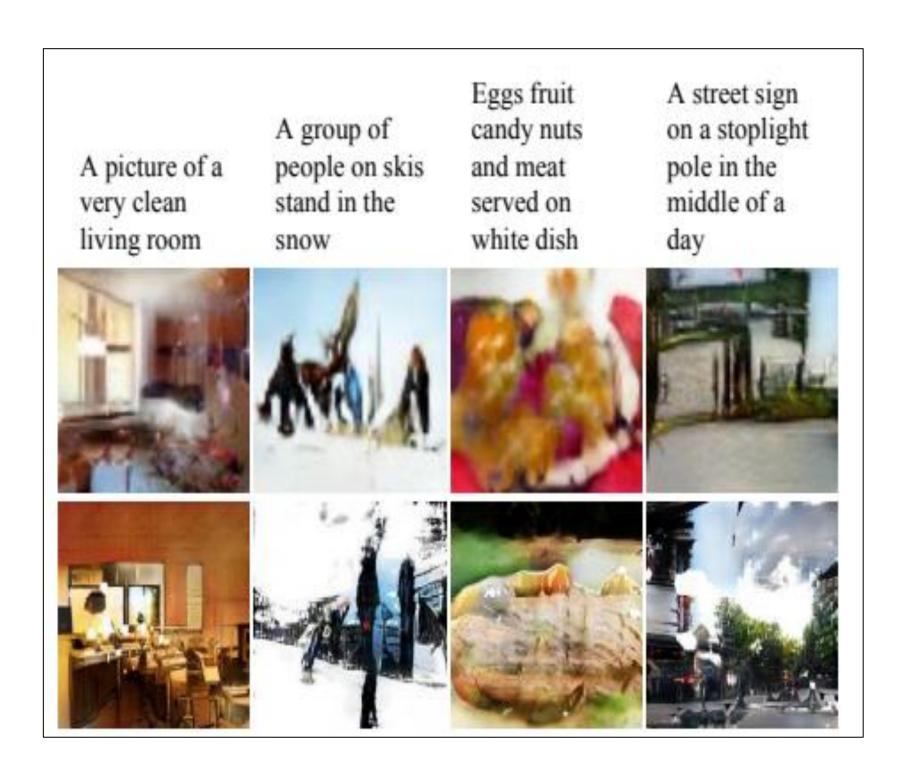
AttnGAN: Fine-Grained Text to Image Generation with Attentional Generative Adversarial Networks

Tao Xu et al., CVPR 2018

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

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





Can we synthesize images beyond single or a few categories

Taming Transformers for High-Resolution Image Synthesis

Björn Ommer Heidelberg Collaboratory for Image Processing, IWR, Heidelberg University, Germany *Both authors contributed equally to this work



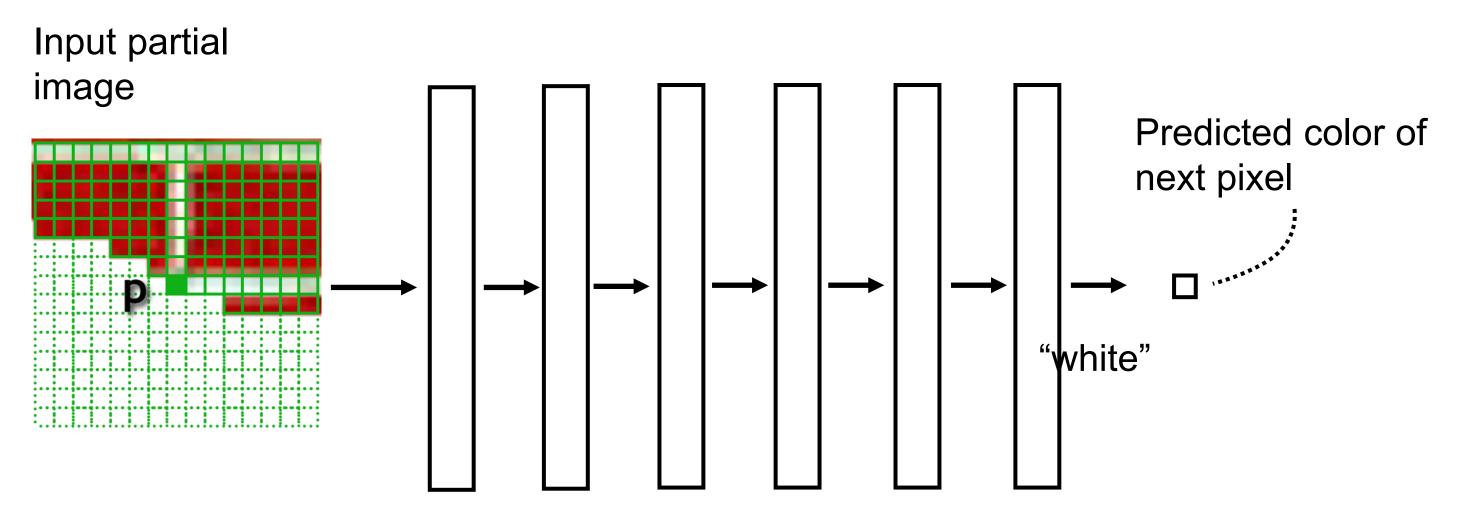
Figure 1. Our approach enables transformers to synthesize high-resolution images like this one, which contains 1280x460 pixels.

Abstract

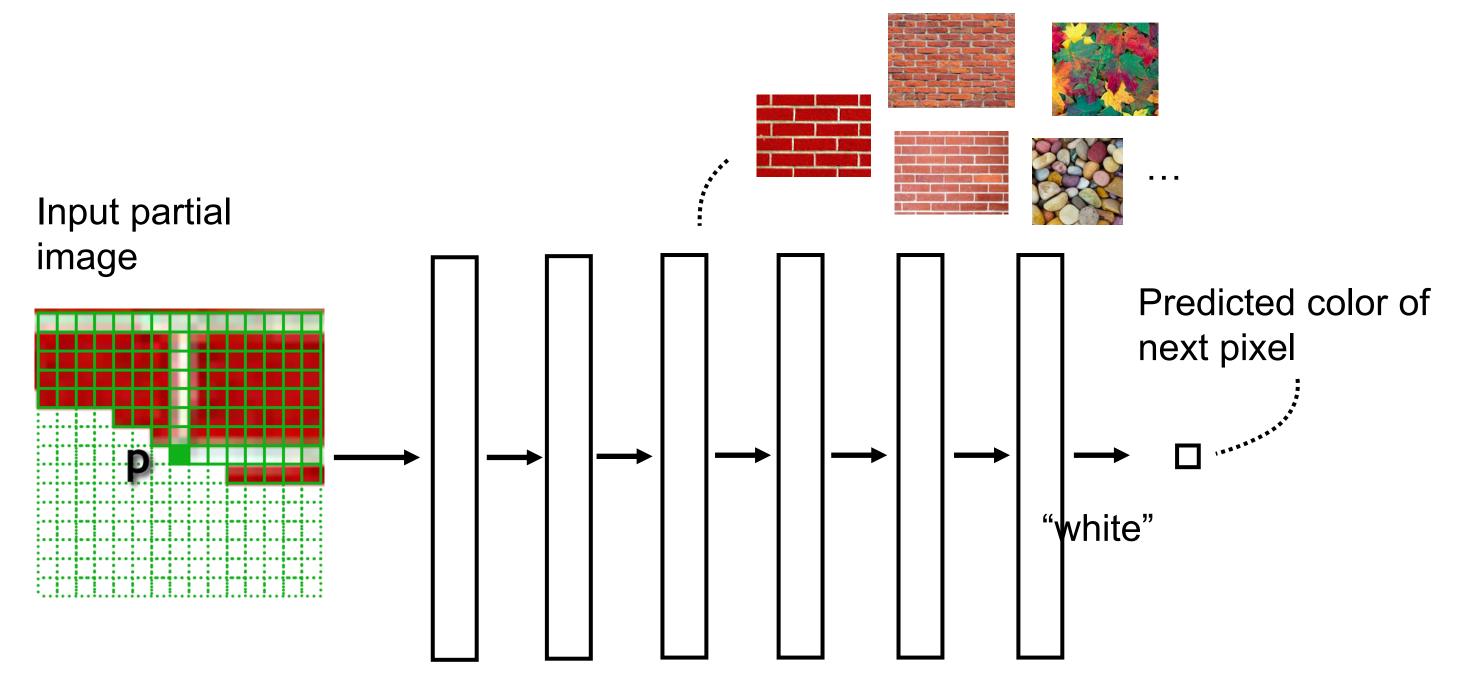
is and to learn long-range interactions on sequential

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 the locality of interactions and is therefore free

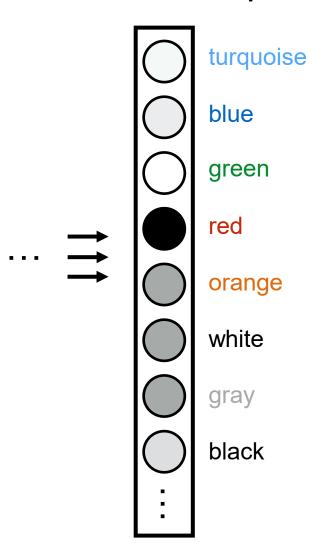
Autoregressive (AR) image synthesis



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

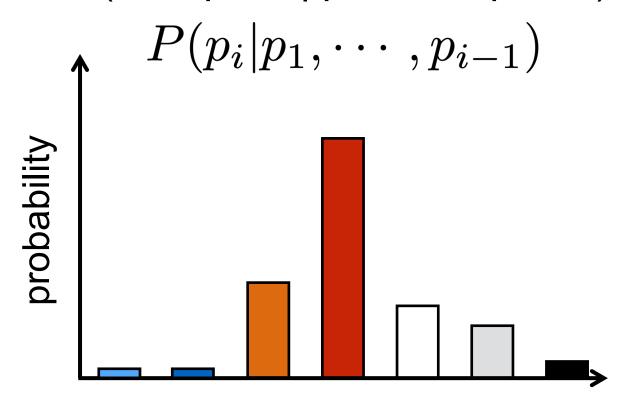


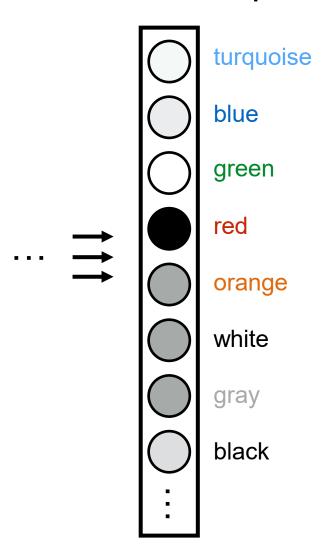
[PixelRNN, PixelCNN, van der Oord et al. 2016]

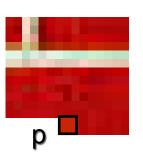


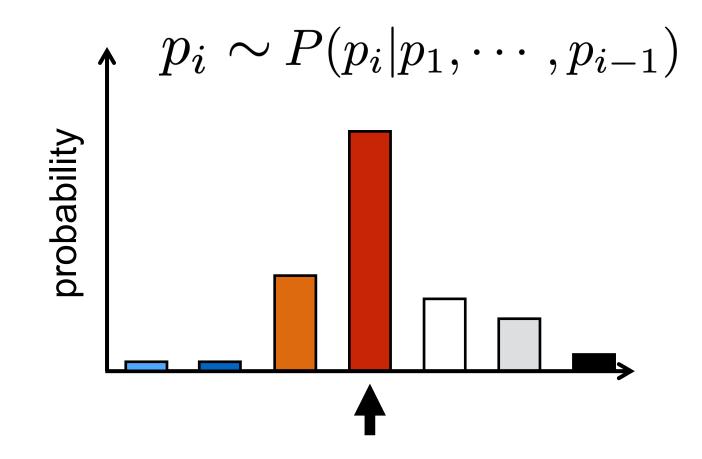


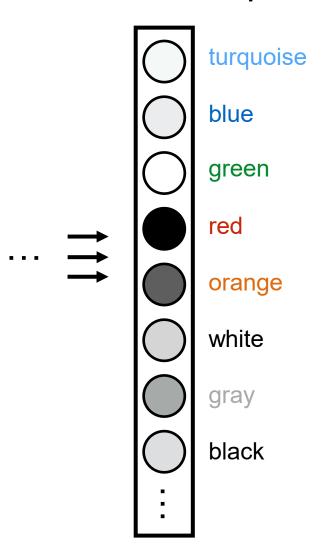
P(next pixel | previous pixels)



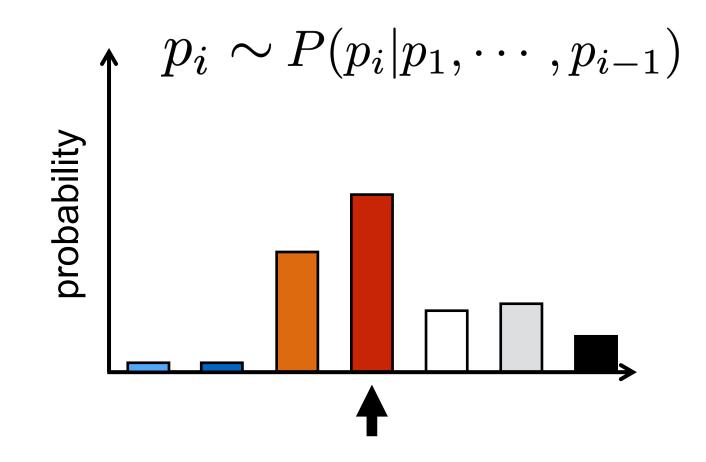


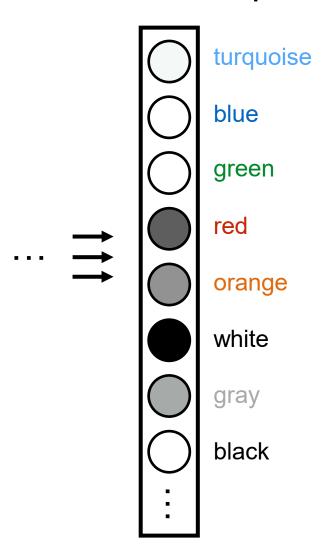




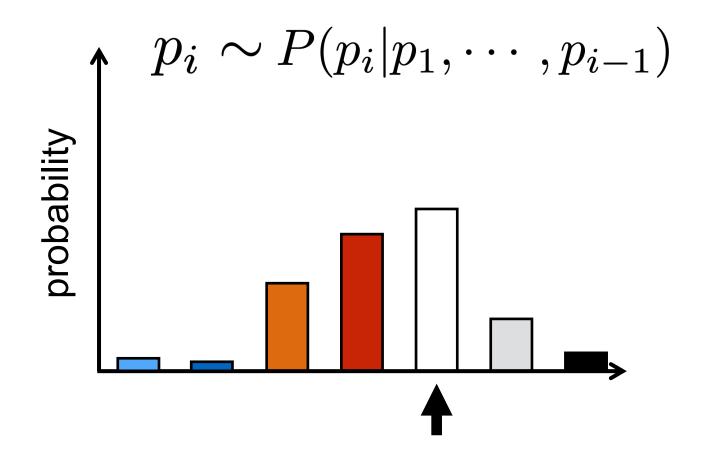


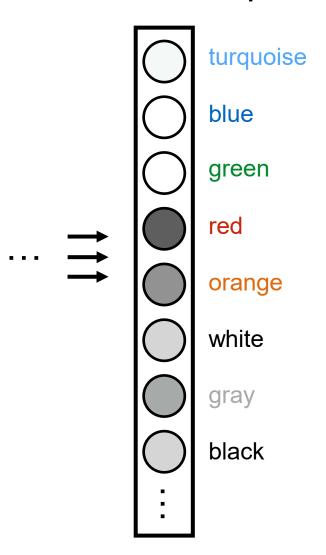




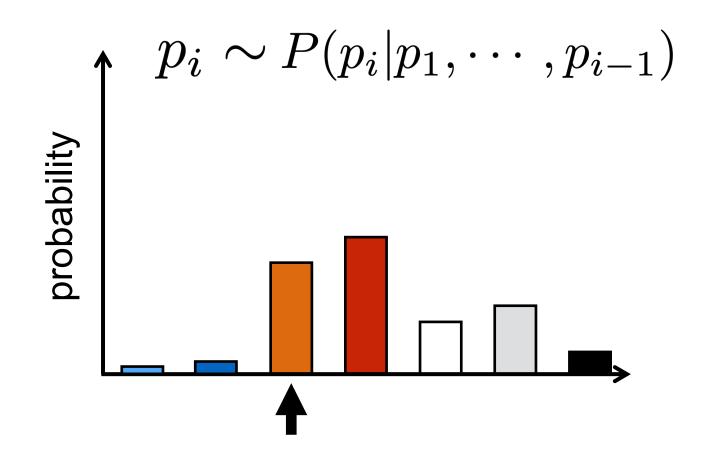




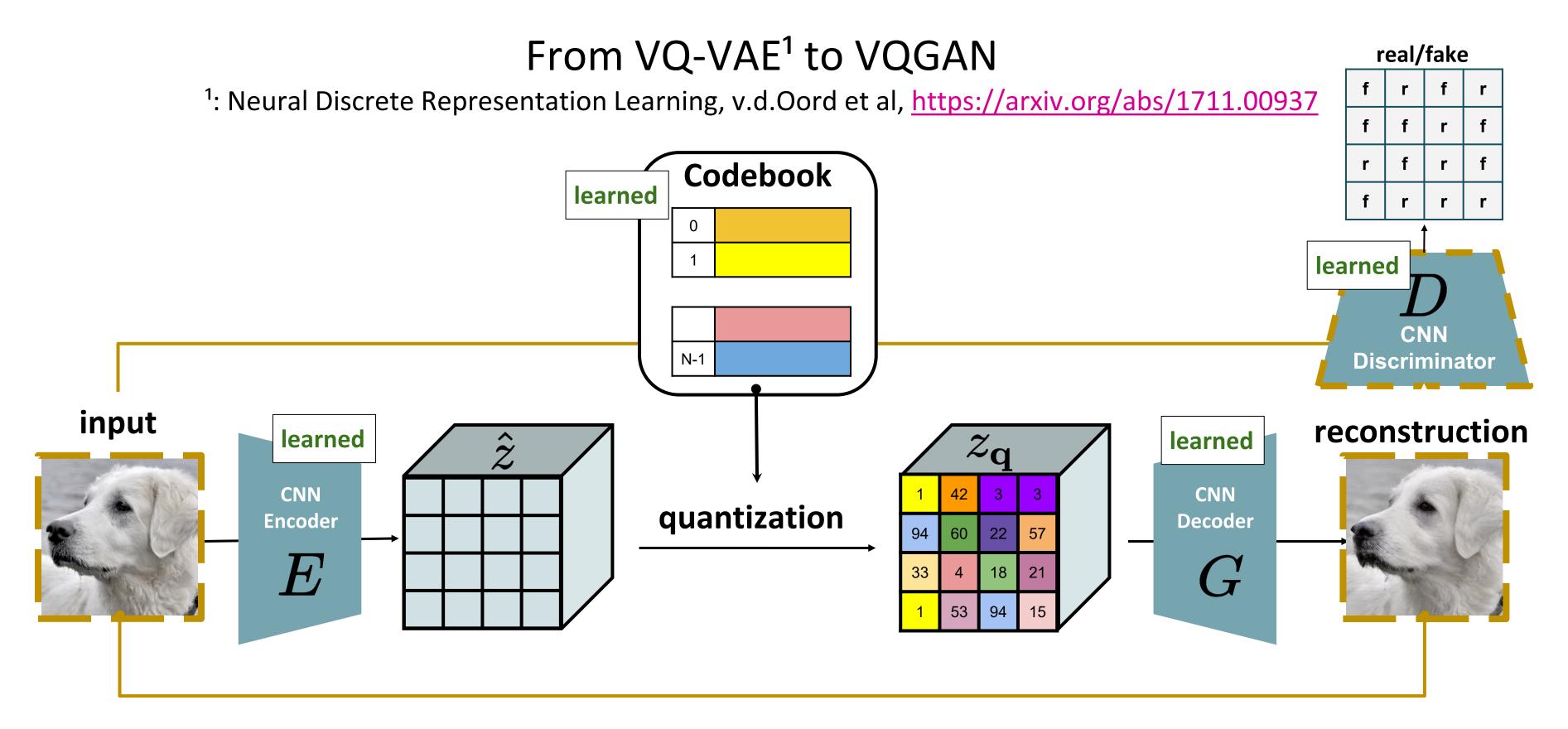








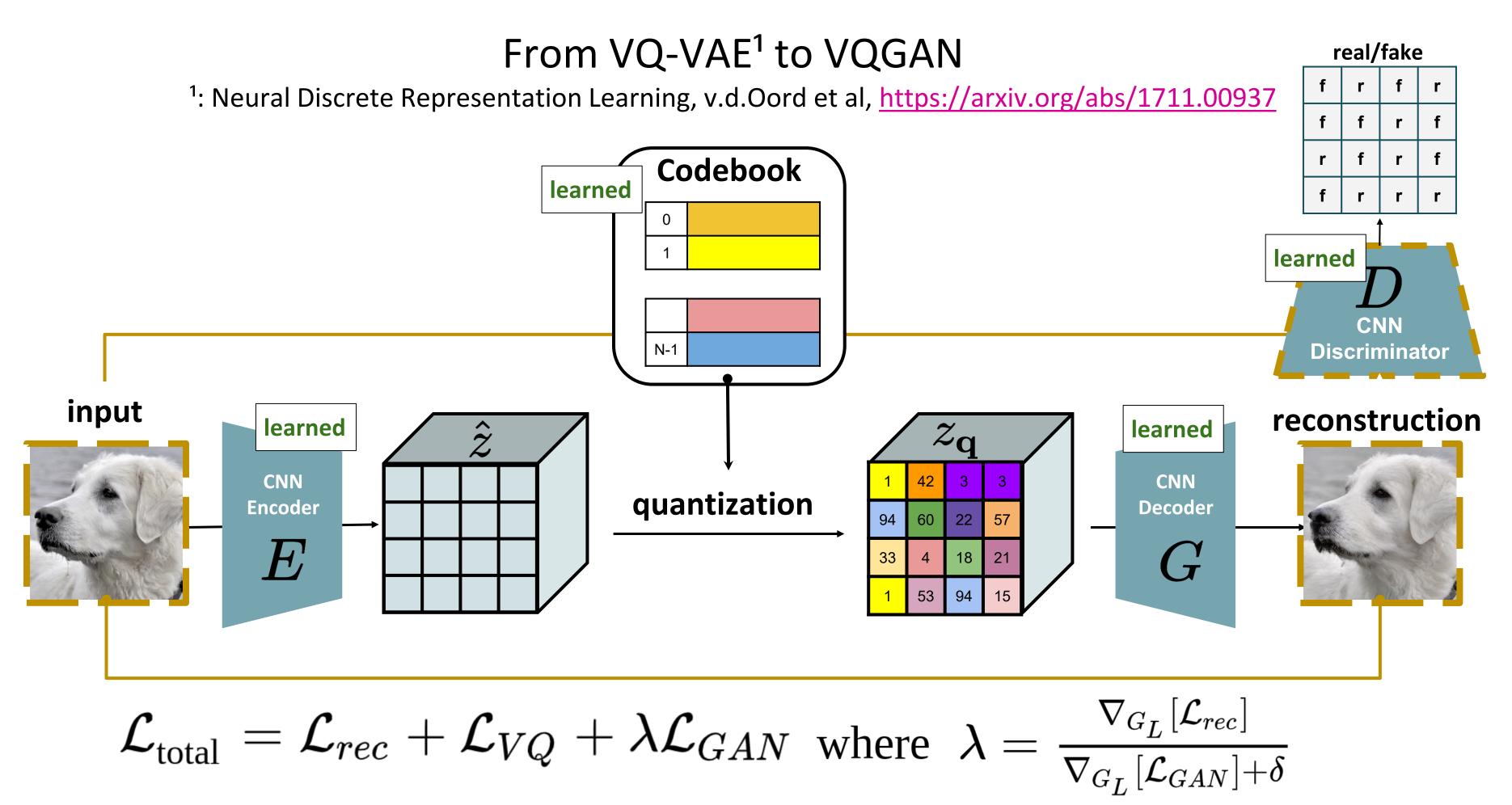
Generation is super slow? What should we do?

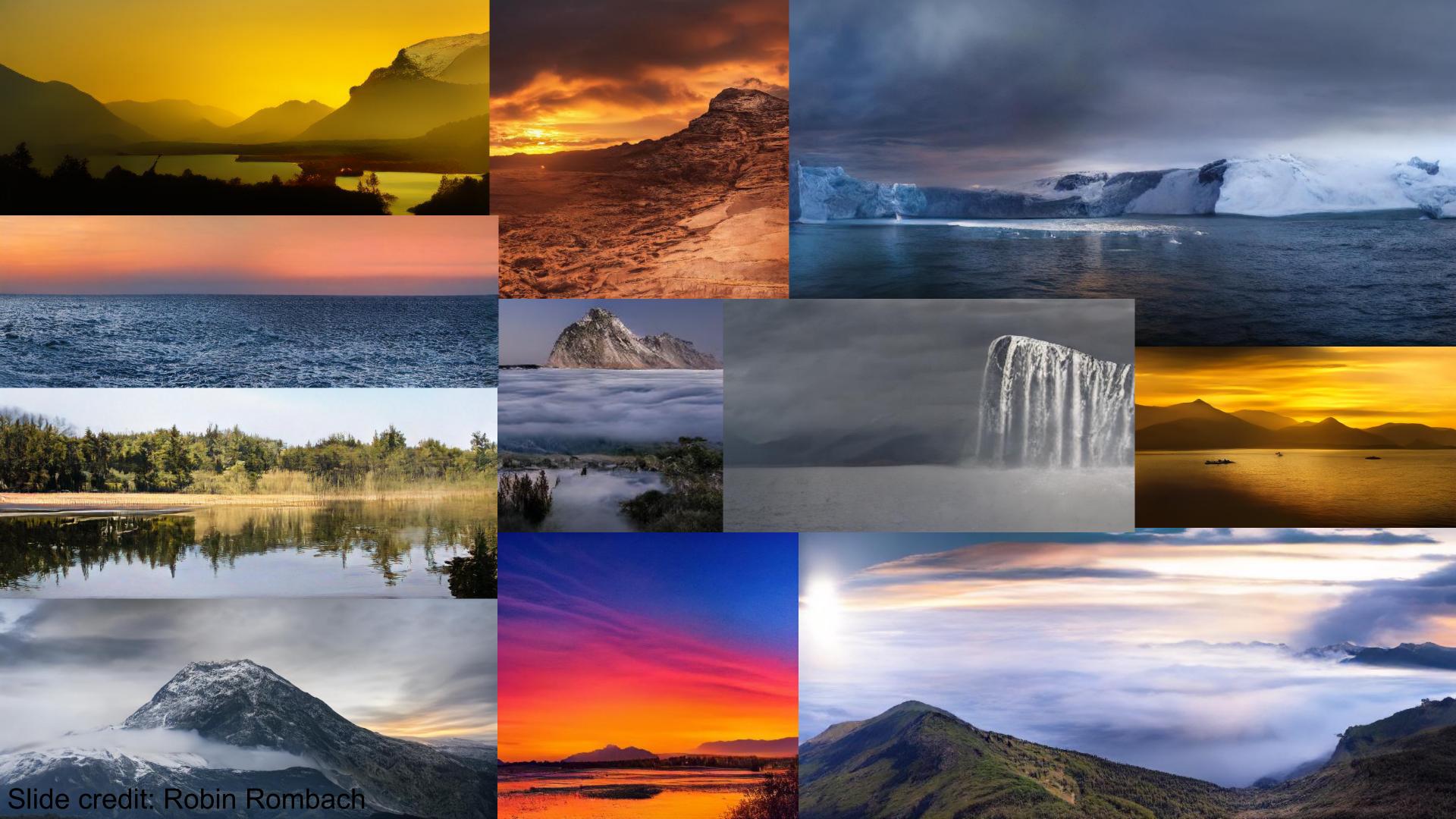


i) replace L2/L1 rec. loss with Perceptual loss (includes pixel-level)

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

Slide credit: Robin Rombach





Scaling VQGAN for Text-to-Image!

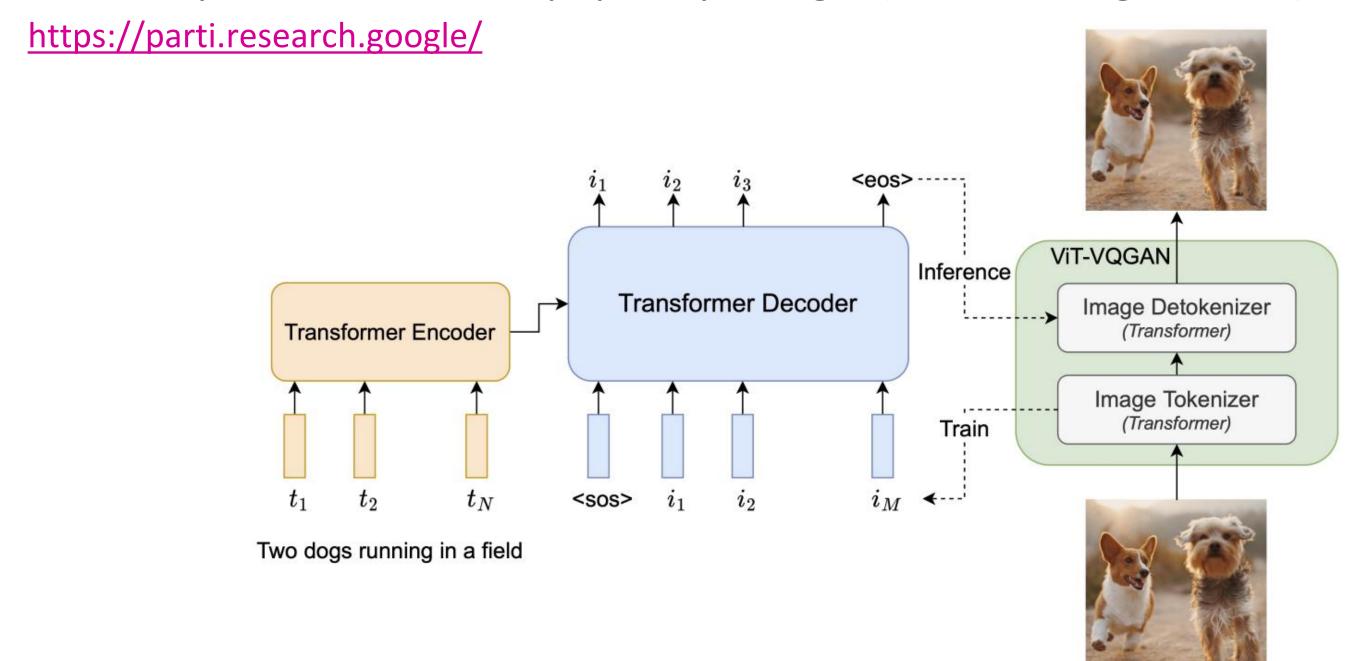
- see recently released "Parti" paper by Google (text-to-image model)
 - https://parti.research.google/



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)



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

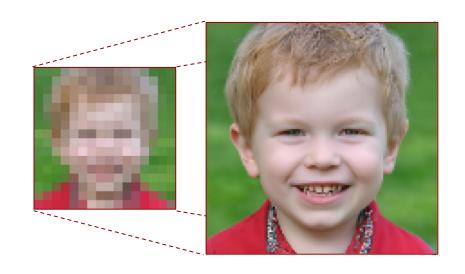
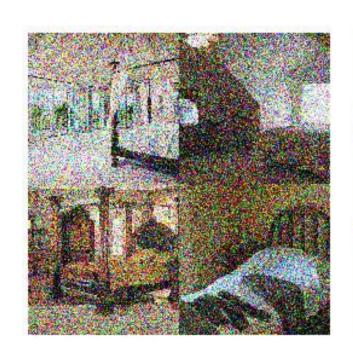


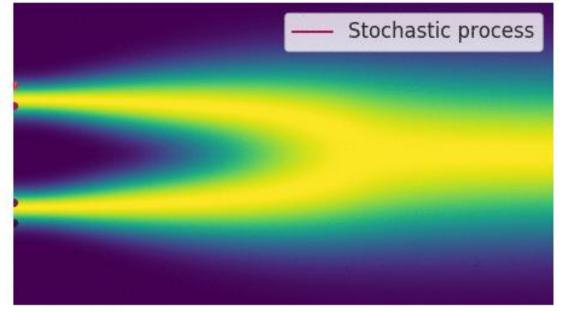
Image Super-Resolution via Iterative Refinement
Chitwan Saharia, et al

https://arxiv.org/abs/2104.07636

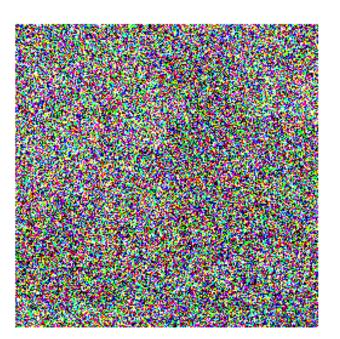
... but very expensive :(

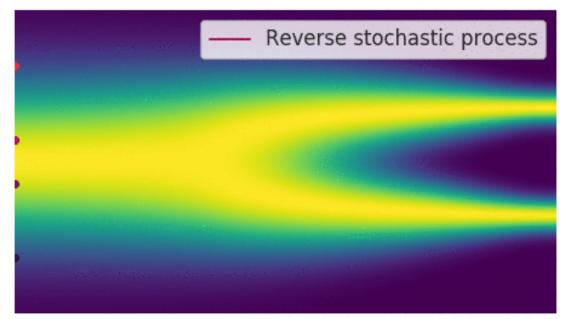
Brief Overview of Diffusion Models





"destroy" the data by gradually adding small amounts of gaussian noise





 "create" data by gradually denoising a noisy code from a stationary distribution

Animations from https://yang-song.github.io/blog/2021/score/

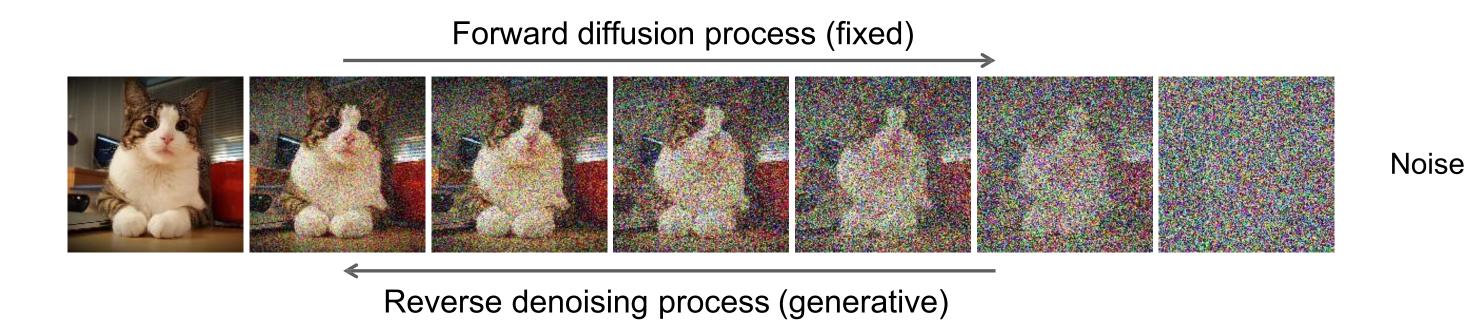
Denoising Diffusion Models

Learning to generate by denoising

Denoising diffusion models consist of two processes:

Data

- 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

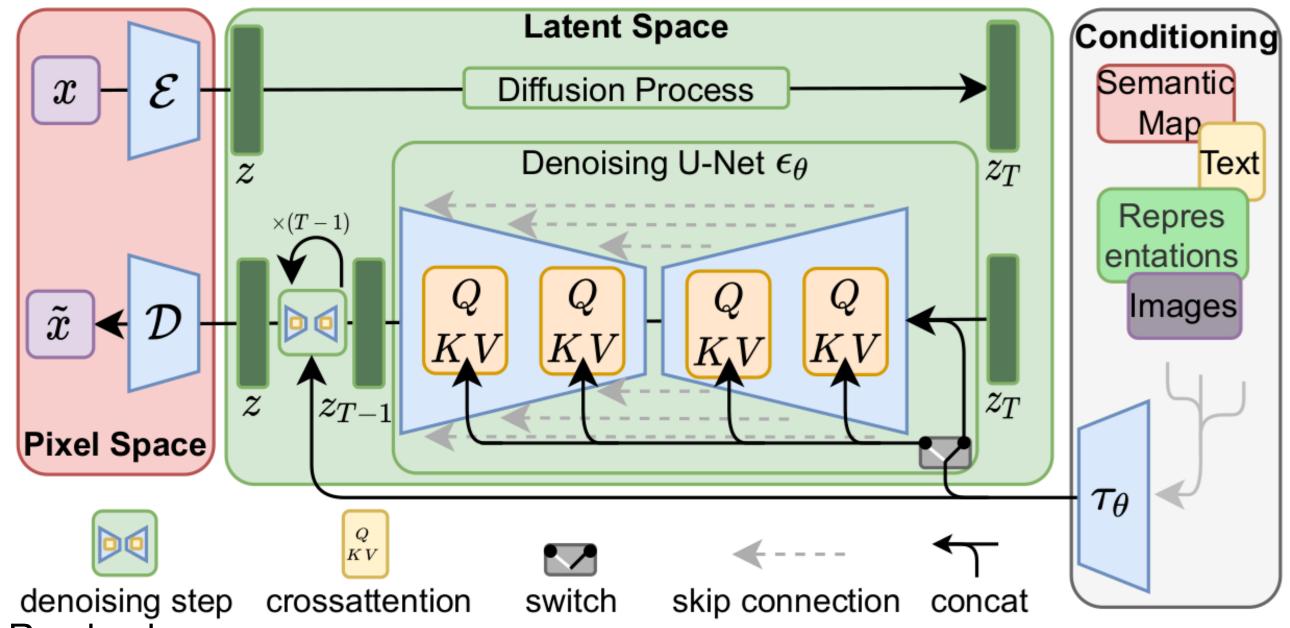
Latent Diffusion Modeling: Architecture

Autoencoder with KL or VQ regularization.

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

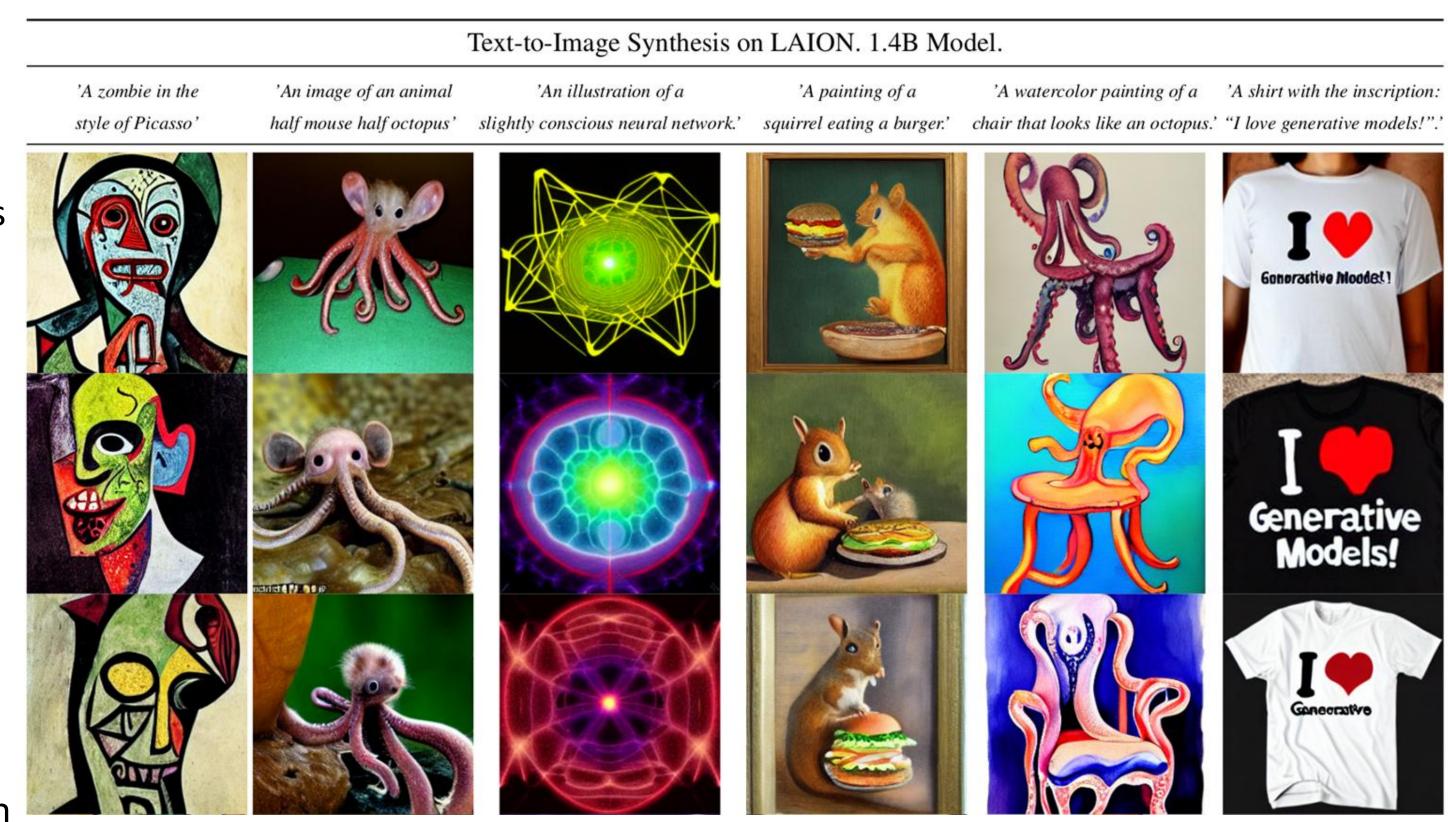
where $\lambda = rac{
abla_{G_L}[\mathcal{L}_{rec}]}{
abla_{G_L}[\mathcal{L}_{GAN}] + \delta}$

$$extsf{KL-reg.:} \quad \mathcal{L}_{ ext{total}} = \mathcal{L}_{rec} + eta \mathcal{L}_{KL} + \lambda \mathcal{L}_{GAN}$$



LDMs for Text-to-Image Synthesis

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



LDMs for Text-to-Image Synthesis

convolutional sampling (train on 256², generate on >256²)



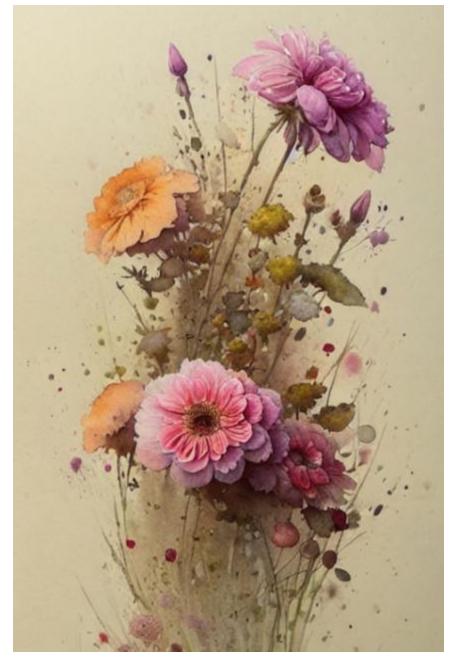




Stable Diffusion

Latent Diffusion ++













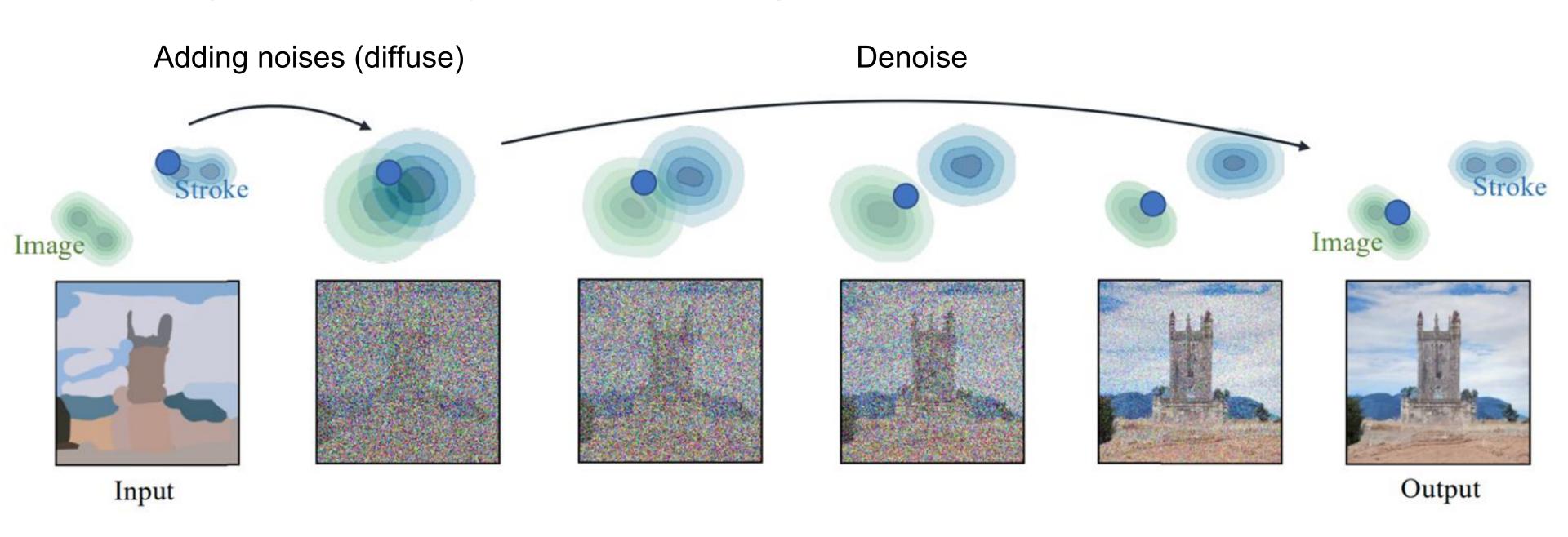




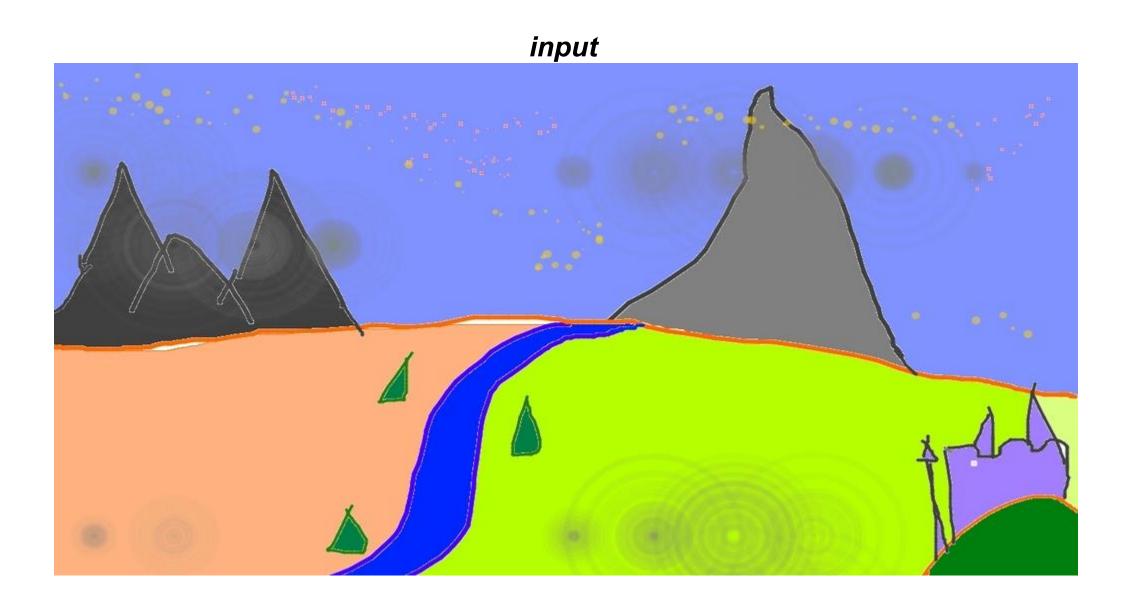


Text-Guided Image-to-Image

SDEdit (https://arxiv.org/abs/2108.01073) recipe: diffuse → denoise



Text-Guided Image-to-Image

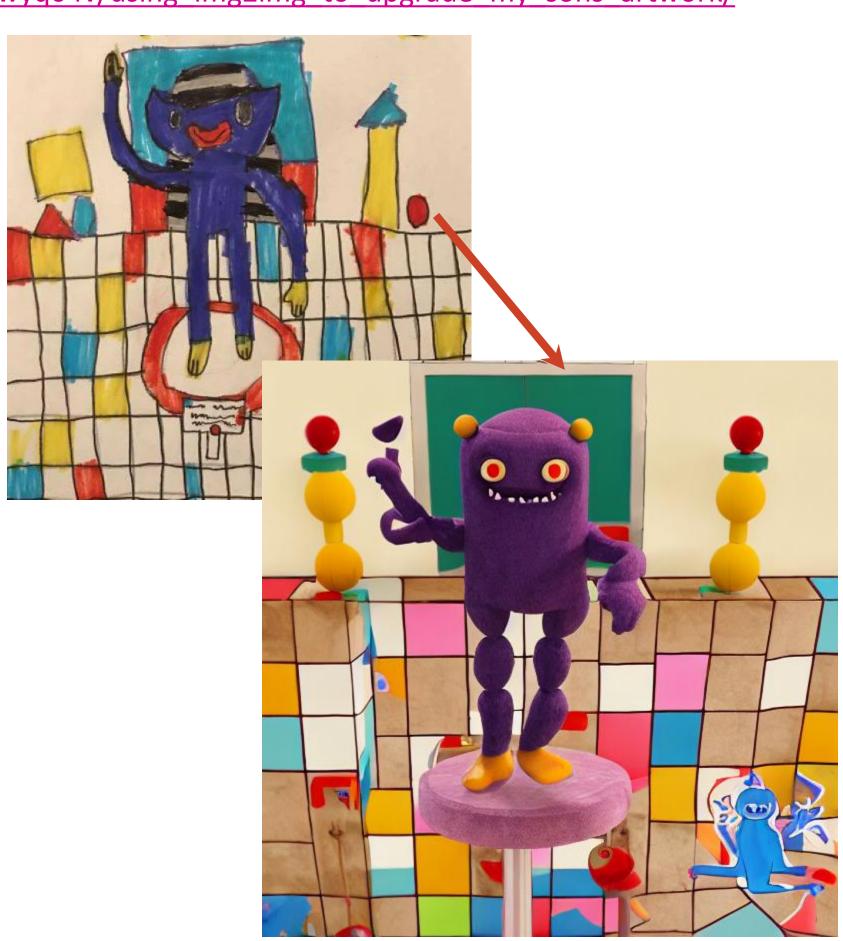


"a fantasy landscape, watercolor painting" "a fantasy landscape, trending on artstation" "a fantasy landscape, by Simon Stalenhag" Slide credit: Robin Rombach

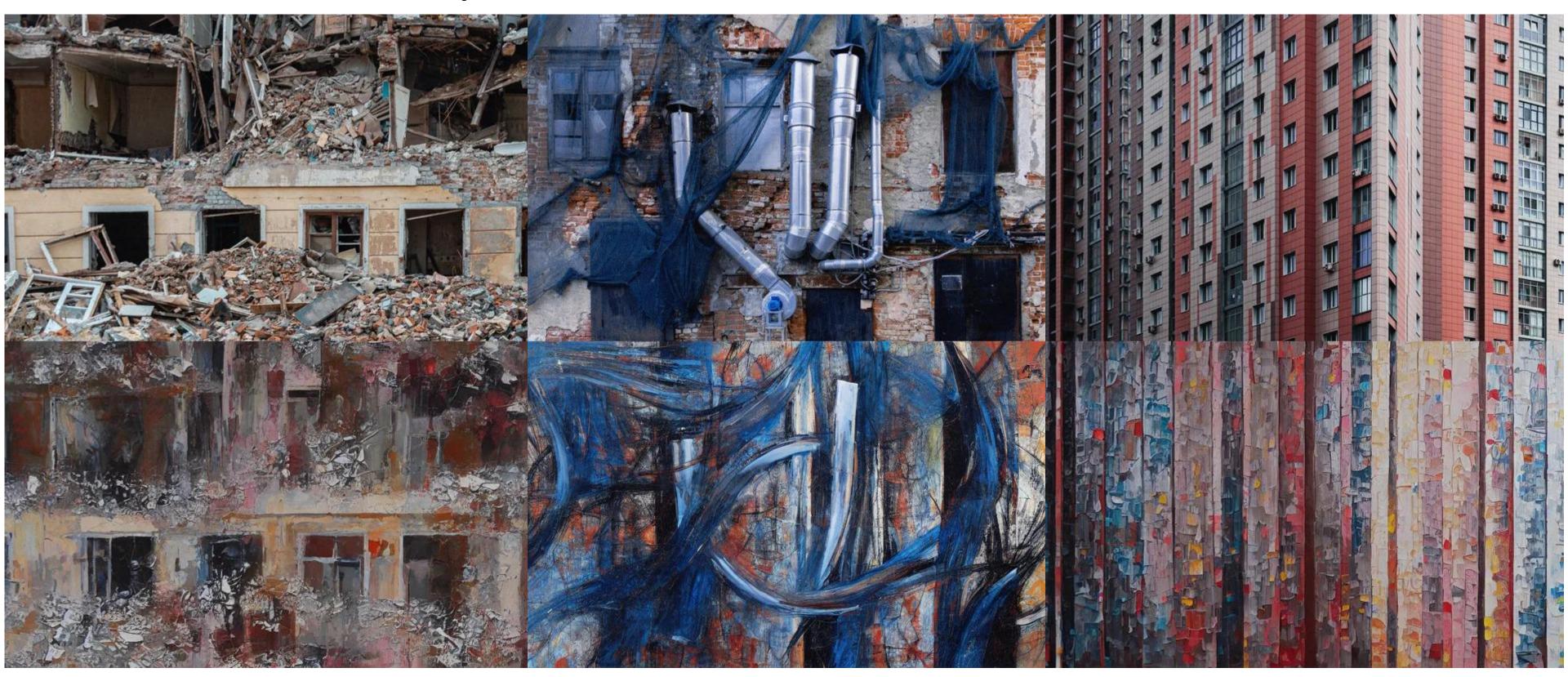
"Upgrade" your child's artwork

original post: https://www.reddit.com/r/StableDiffusion/comments/wyq04v/using-img2img-to-upgrade-my-sons-artwork/





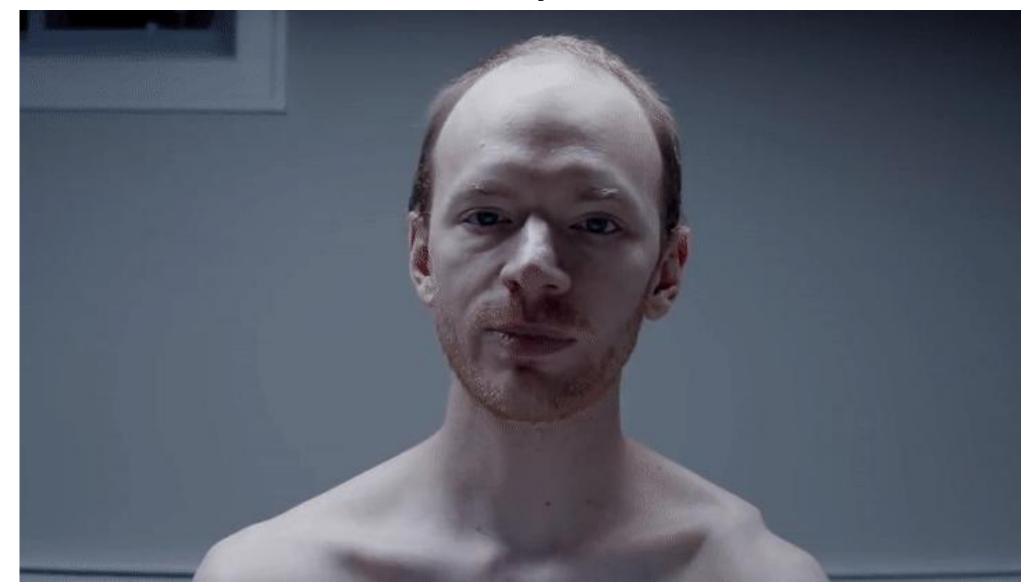
abstract art from photos



original post by <u>u/Pereulkov</u>

https://www.reddit.com/r/StableDiffusion/comments/xhhyad/i made abstract art from my photos/

Video Synthesis

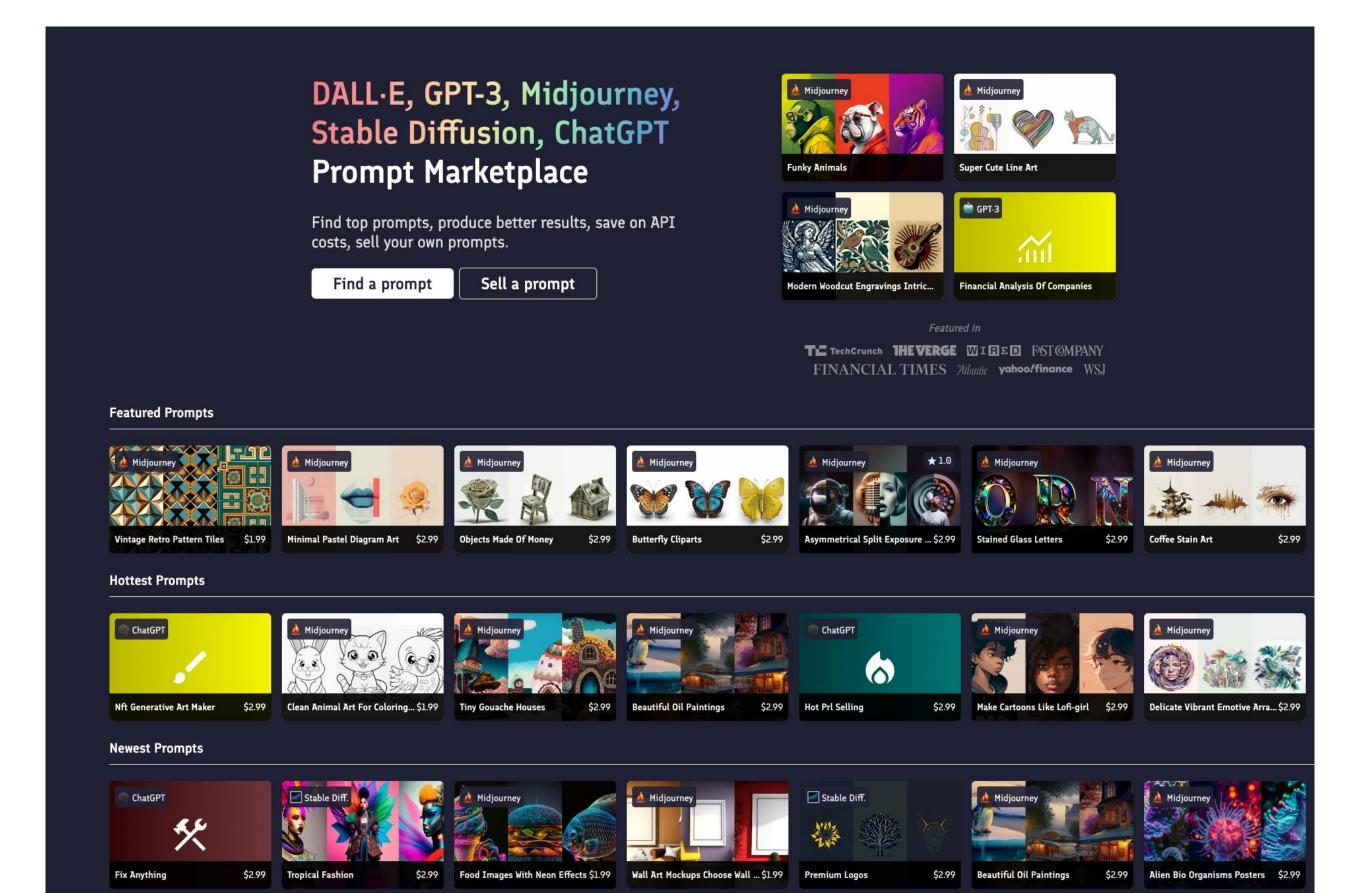


Stable Diffusion (img2img) + EBSynth by Scott Lightsier:

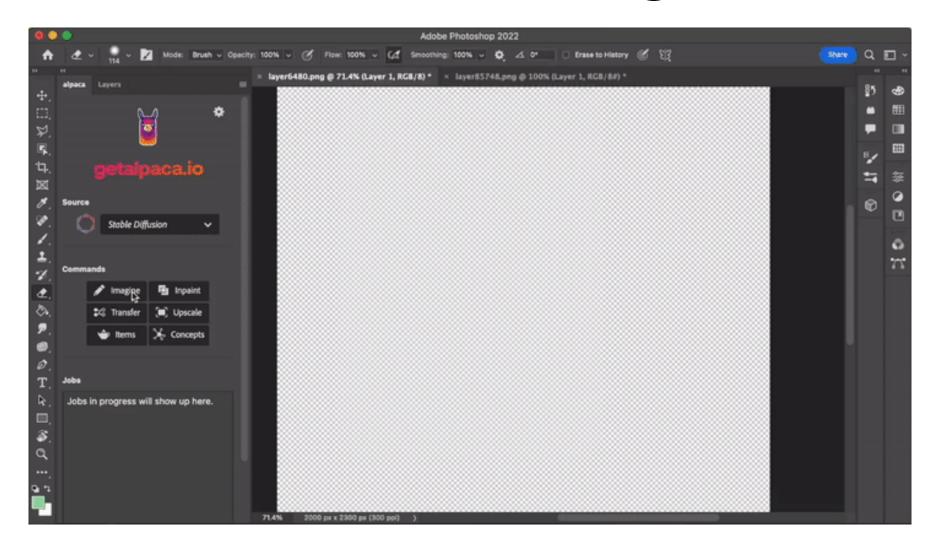
https://twitter.com/LighthiserScott/status/1567355079228887041?t=kXXCAVtuO5lJCGcro3Ma3A&s=19

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

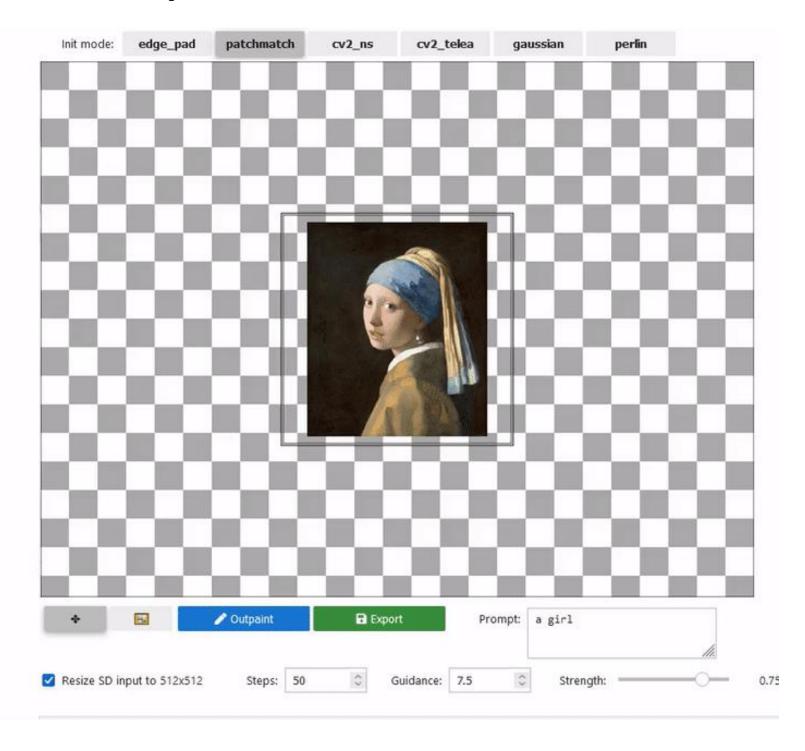
Prompt Marketplace (promptbase.com)



Uls / Plug-Ins for Photoshop, GIMP etc



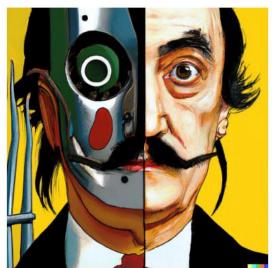
https://twitter.com/wbuchw/status/1563162131024920576



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

What if you have 1,000+ GPUs/TPUs

DALL-E 2, Imagen









a shiba inu wearing a beret and black turtleneck

a close up of a handpalm with leaves growing from it







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

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







Sprouts in the shape of text 'Imagen' coming out of a A photo of a Shiba Inu dog with a backpack riding a A high contrast portrait of a very happy fuzzy panda fairytale book.

A photo of a Shiba Inu dog with a backpack riding a 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 Butter- 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.

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?











Benchmarking Spatial Relationships in Text-to-Image Generation



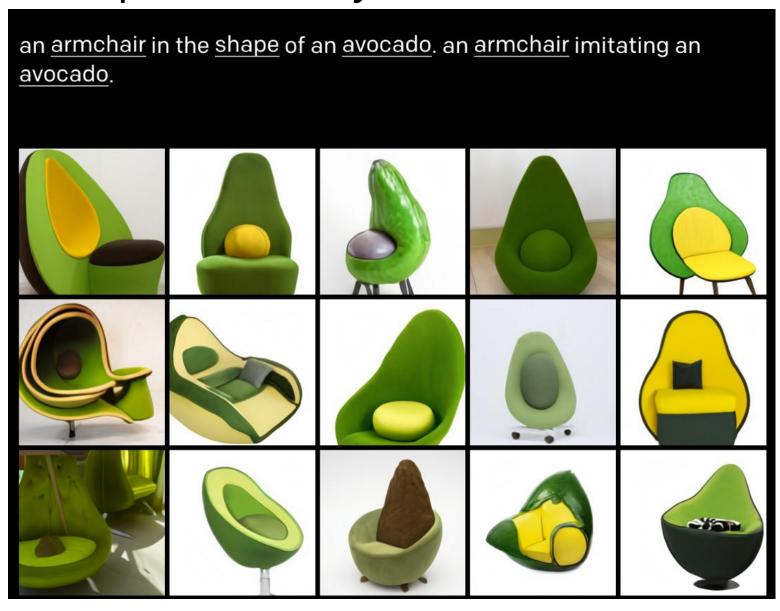
Figure 1: We benchmark T2I models on their competency with generating appropriate spatial relationships in their visual renderings. Although text inputs may explicitly mention these spatial relationships, T2I models lack such spatial understanding.

visort2i.github.io

https://github.com/microsoft/VISOR

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

Attribute-Level Compositionality Compositionality

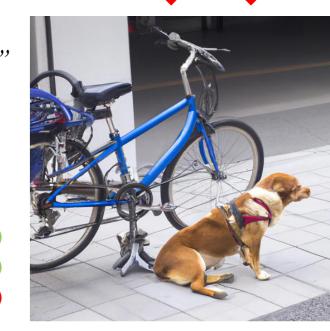


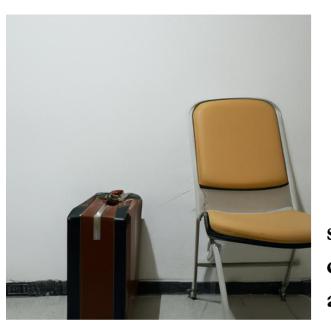
Object-Level / Spatial



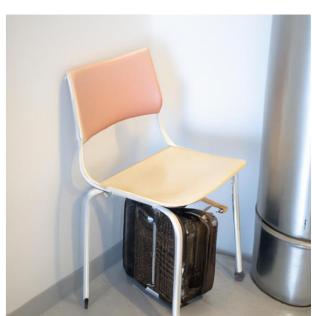












Follow-up (Method to Improve Spatial Reasoning in T2I)



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

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

⁵University of Maryland, Baltimore County



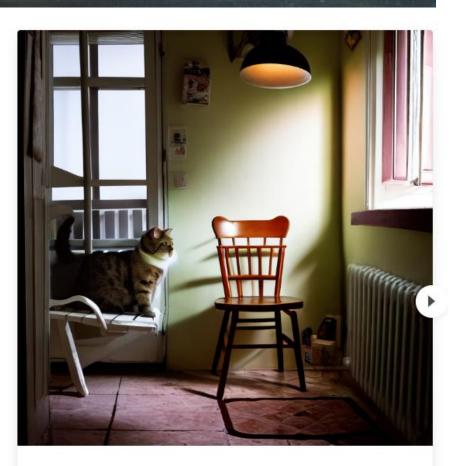
A giraffe to the right of a truck.



A hair drier to the right of a wine glass.



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



A cat sitting <u>on</u> a chair with a lamp to the <u>right</u> and a window <u>above</u>, 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

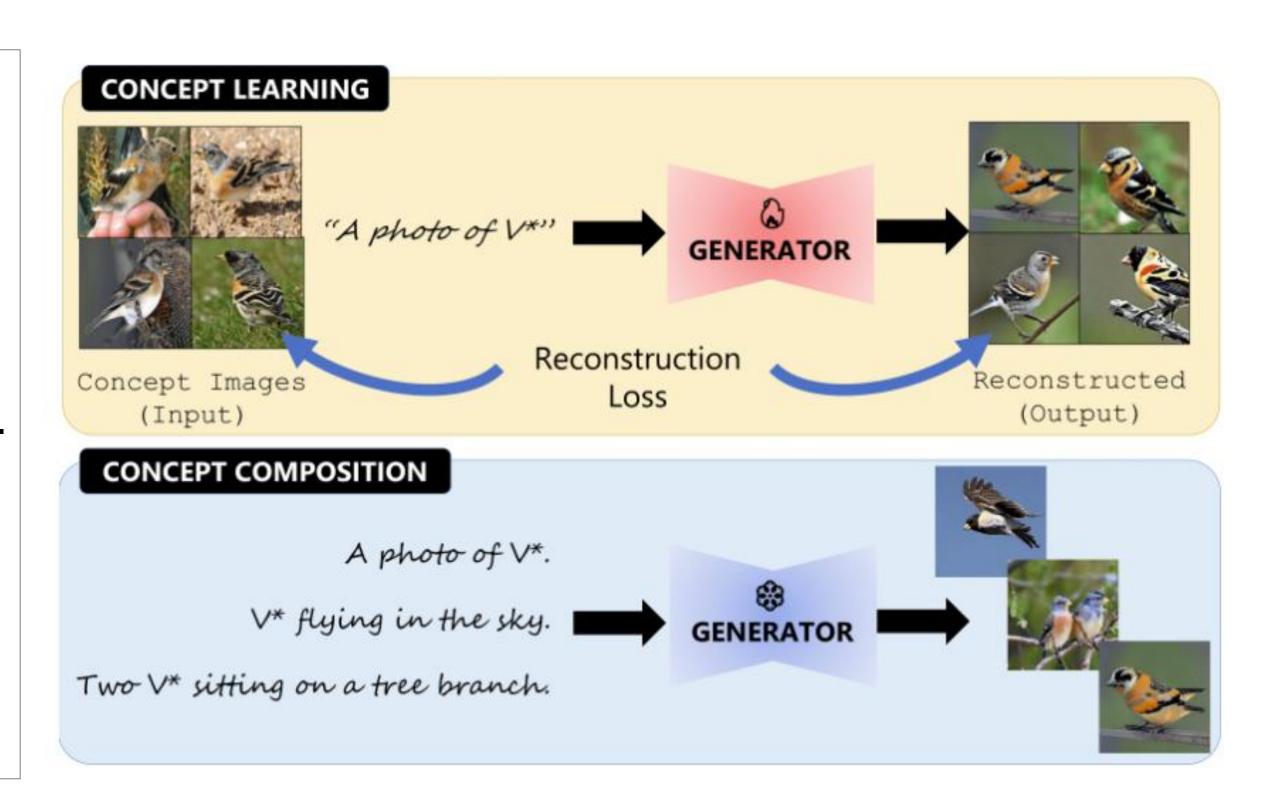
Arizona State University
 University of Maryland Baltimore County

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

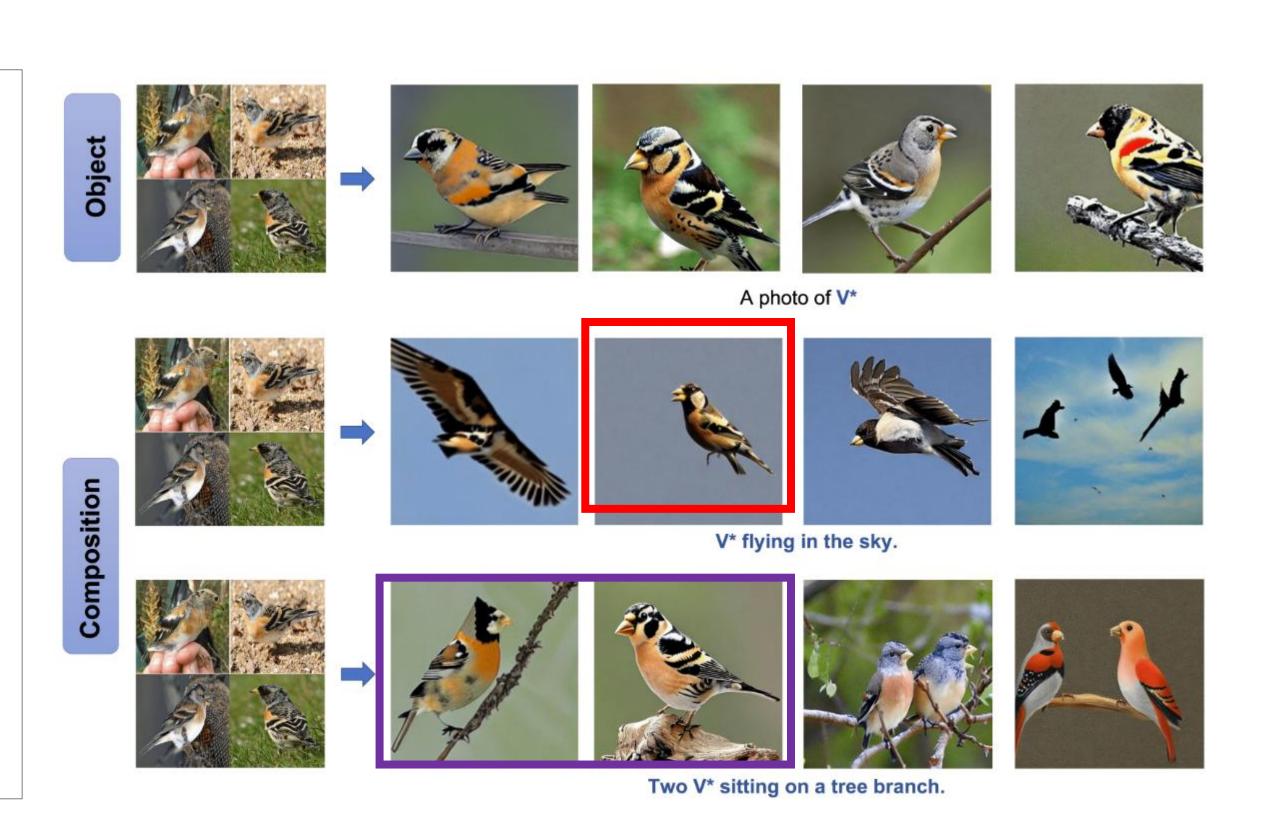


ConceptBed

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 ...



Dataset

Evaluation Metric:
 "Concept Confidence Deviation"

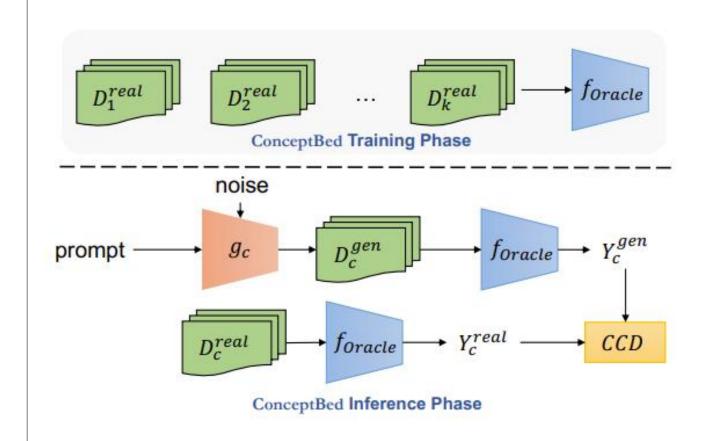


Figure 3: The outline of **CONCEPTBED** evaluation framework.

Domain

Art Painting Cartoon Photo Sketch

4 concepts with 7 classes each

ConceptBed Dataset

Objects

DogsAircraftCarFruitCatsBoatClockBirdMonkeysBottleFish...

15 High-level Concepts, 80 Low-level concepts

Total: 284 unique concepts

ConceptBed Compositions

Atrributes

Yellow Wing Orange Beak Brown Eyes

• • •

200 concepts with 112 attributes

Attribute

V* with red feathers V*'s ears are up A red V* eating

3028 Prompts

Relation

V* has fish in the mouth V* is sitting in a bucket A horse toy beside the V*

2891 Prompts

Action

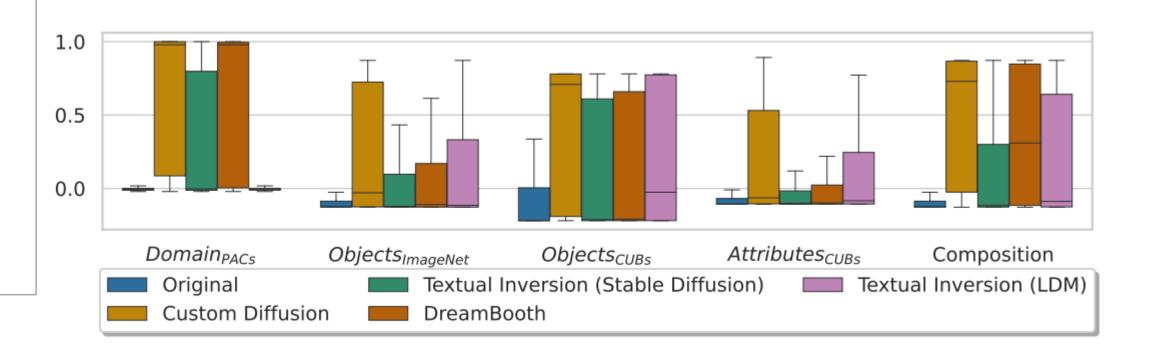
V* is licking herself
V* is running at the beach
A red V* eating

1267 Prompts

Counting

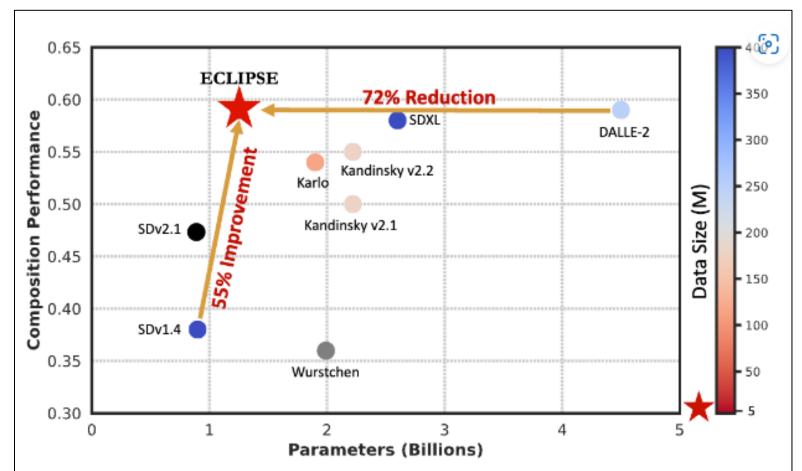
Two V* sitting together A photo of three V* Two V* with a cat

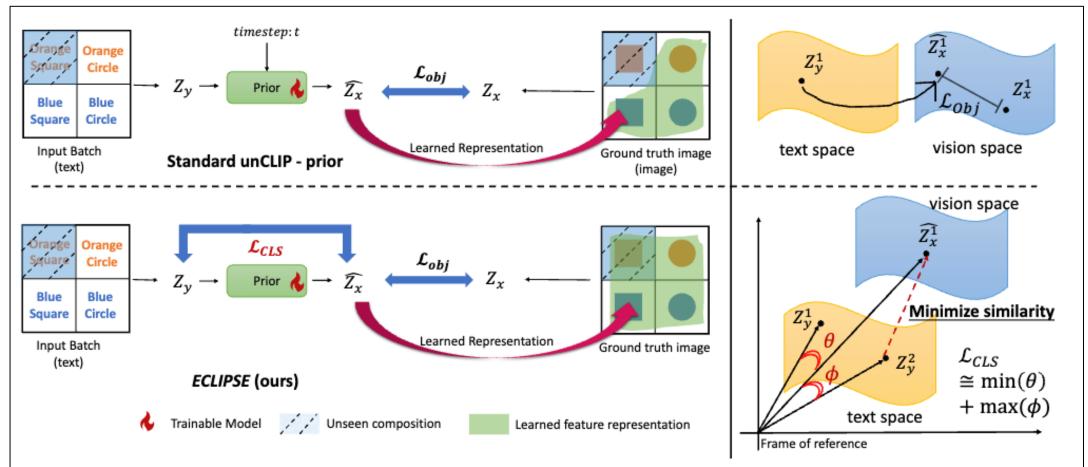
203 Prompts



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.