Diffusion Model Track

- Methods and Application in Image Synthesis

Sakura 2024/11/17

bili_sakura@zju.edu.cn



Timeline of novel works on diffusion (image) model

> OpenAl (Tim Brooks, Yang Song)

iGPT (Chen et al., 2020) → ADM (Dhariwal & Nichol, 2021) → DALL-E (Ramesh et al., 2021) → GLIDE (Nichol et al., 2022) → DALL-E 2 (Ramesh et al., 2022) → DALL-E 3 (Betker et al., 2023) → Consistency Models (Lu & Song, 2024)

➤ Google (Jonathan Ho)

Classifier-Free Guidance (Ho et al., 2021) → Prompt-to-Prompt (Hertz et al., 2022) → Imagen (Saharia et al., 2022) → Null-Text Inversion (Mokady et al., 2023) → Imagen (Imagen 3 Team, 2024)

> Stanford (Ermon Group)

Score-based Generative Models (Song & Ermon, 2019) > Stochastic Differential Equations (Song et al., 2020) >

- → DDIM (Song et al., 2021) → SDEdit (Meng et al., 2021) → Distillation (Meng et al., 2023) →
- → DPO-Diffusion(Wallace et al., 2024)
- > NVIDIA (Jiaming Song)

K-Diffusion (Karras et al., 2022) → VideoLDM (Blattmann et al., 2023) → DiffiT (Hatamizadeh et al., 2024)

Stability AI (Robin Rombach)

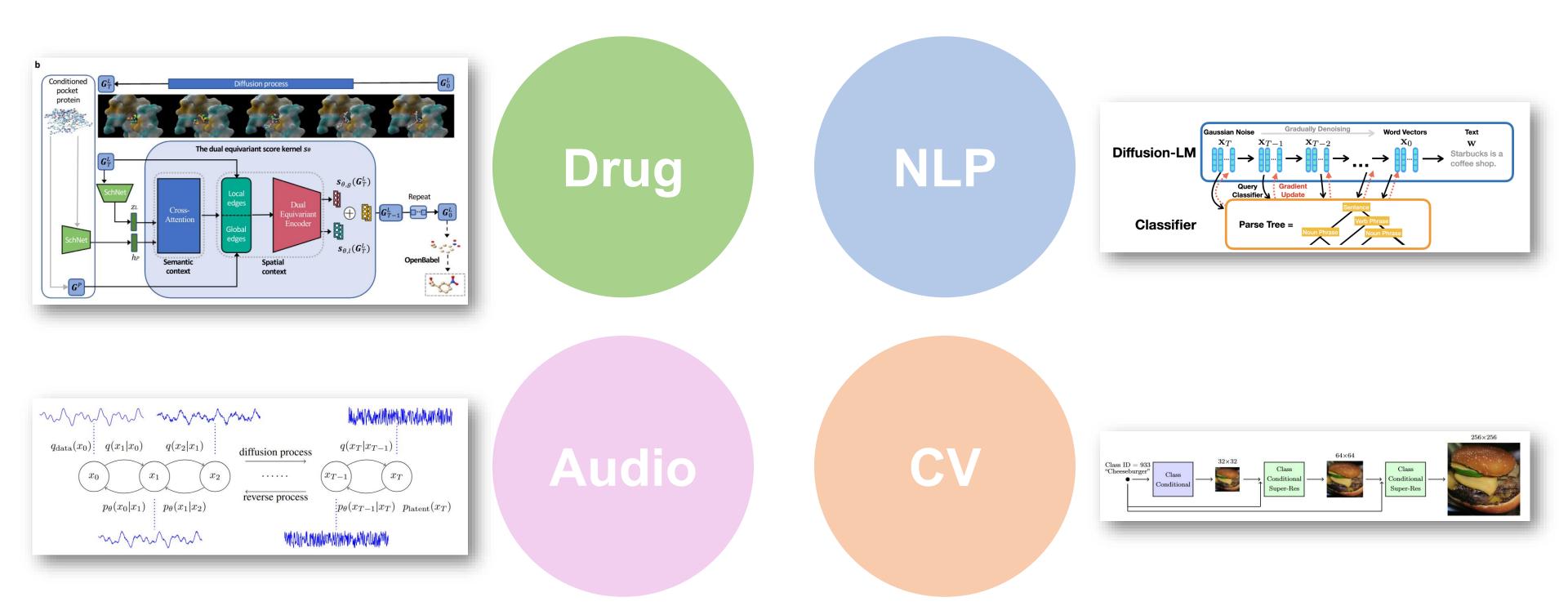
SD (Rombach et al., 2022) → SDXL (Podell et al., 2023) → SDXL Turbo (Nichol et al., 2024) → SD3(Esser et al., 2024)

Unclassified

DDPM (Ho et al., 2020); ControlNet (Zhang et al., 2023); InstructPix2Pix (Brooks et al., 2023);

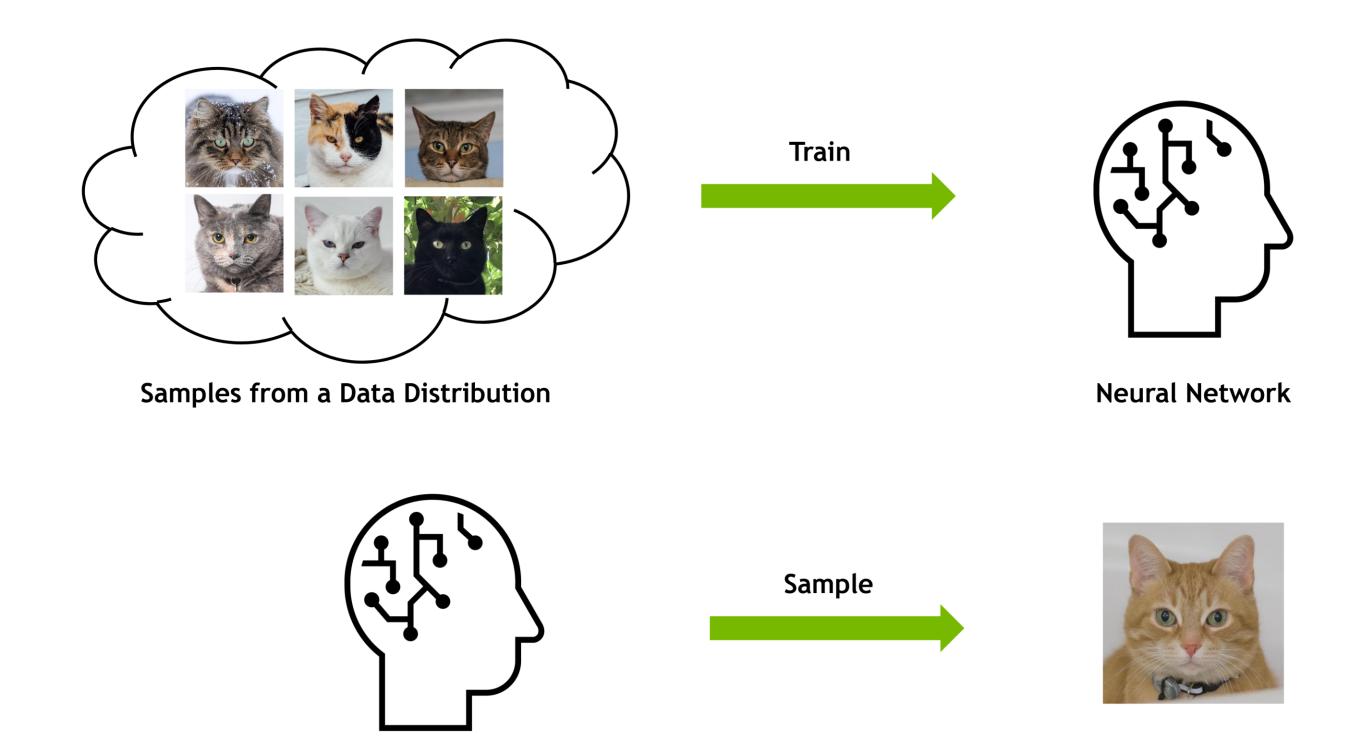
© Sakura, 2024. All rights reserved.

Broad Applications of Diffusion Models



Deep Generative Learning

Learning to generate data



The Landscape of Deep Generative Learning

Bayesian Networks

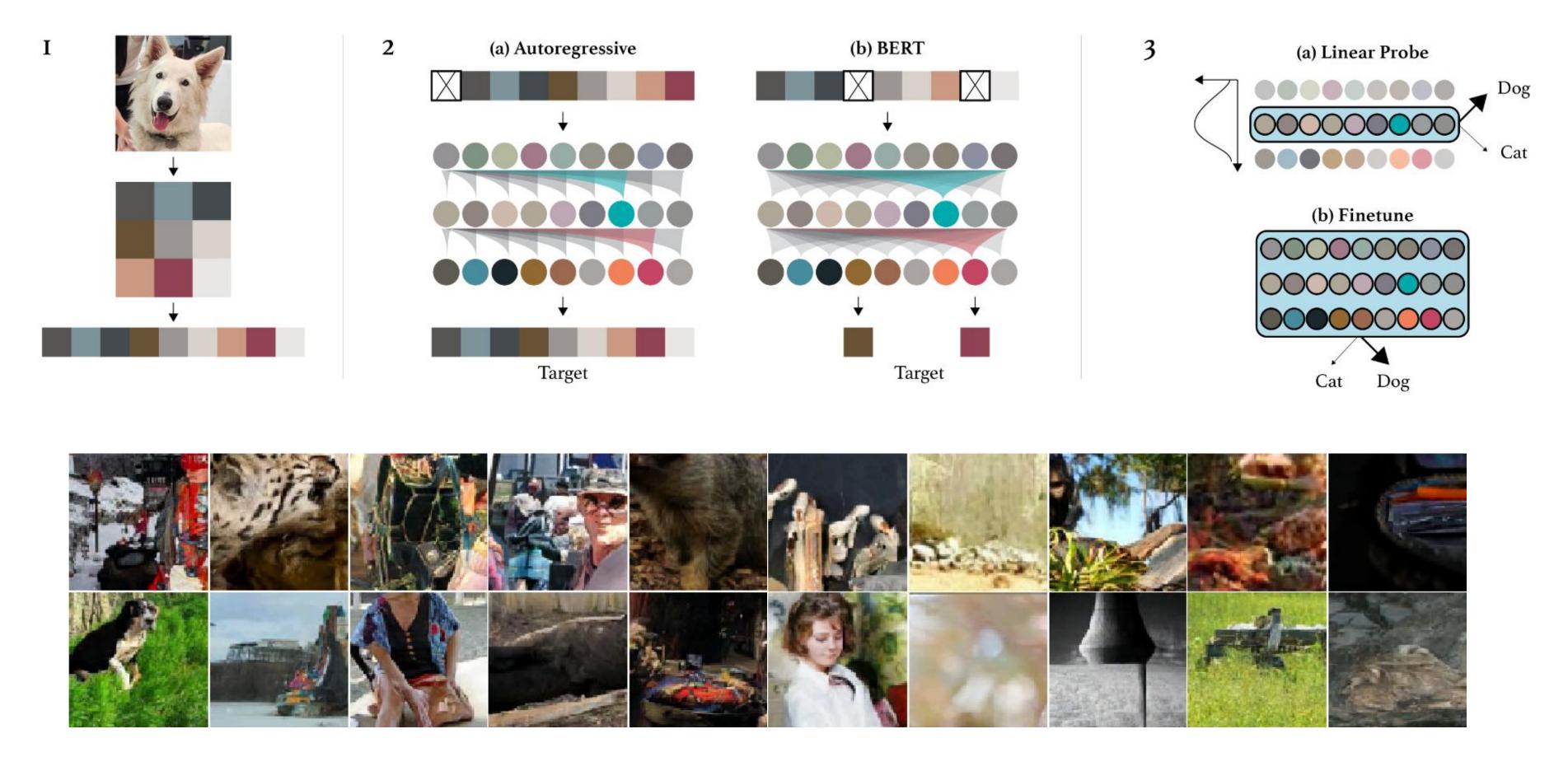
Variational Autoencoders Normalizing Flows

Restricted
Boltzmann Machines

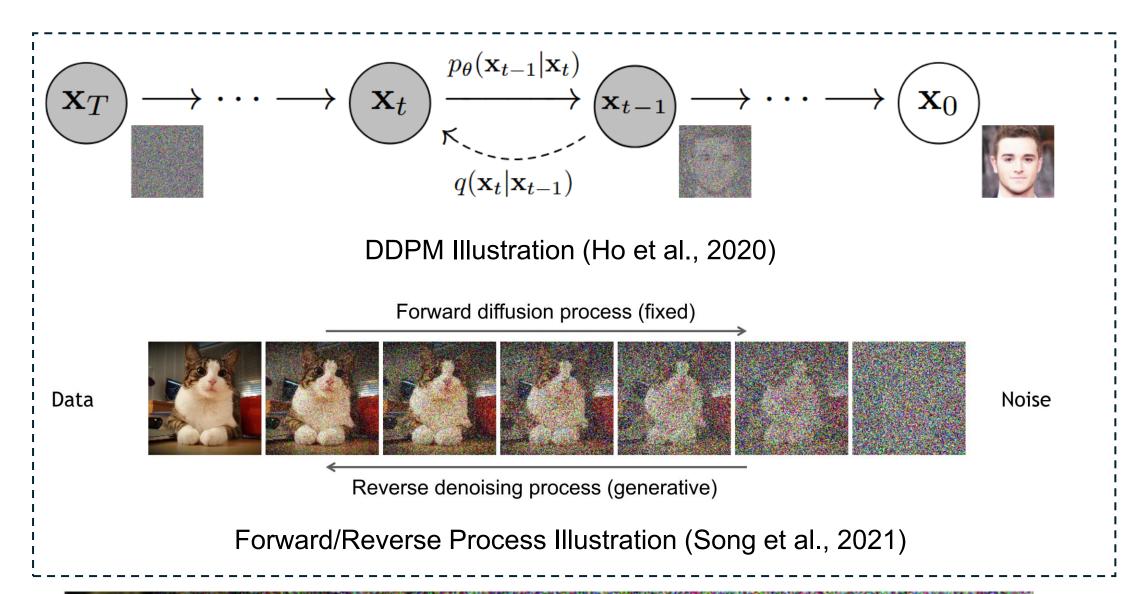
Energy-based Models

Generative Adversarial Networks Autoregressive Models

Denoising Diffusion Models



Class-unconditional samples from iGPT-L trained on input images of resolution 96×96 (Chen et al., 2020).





Instead of linear schedule (top) in DDPM, Improved DDPM uses cosine (bottom) schedules. The cosine schedule adds noise more slowly.

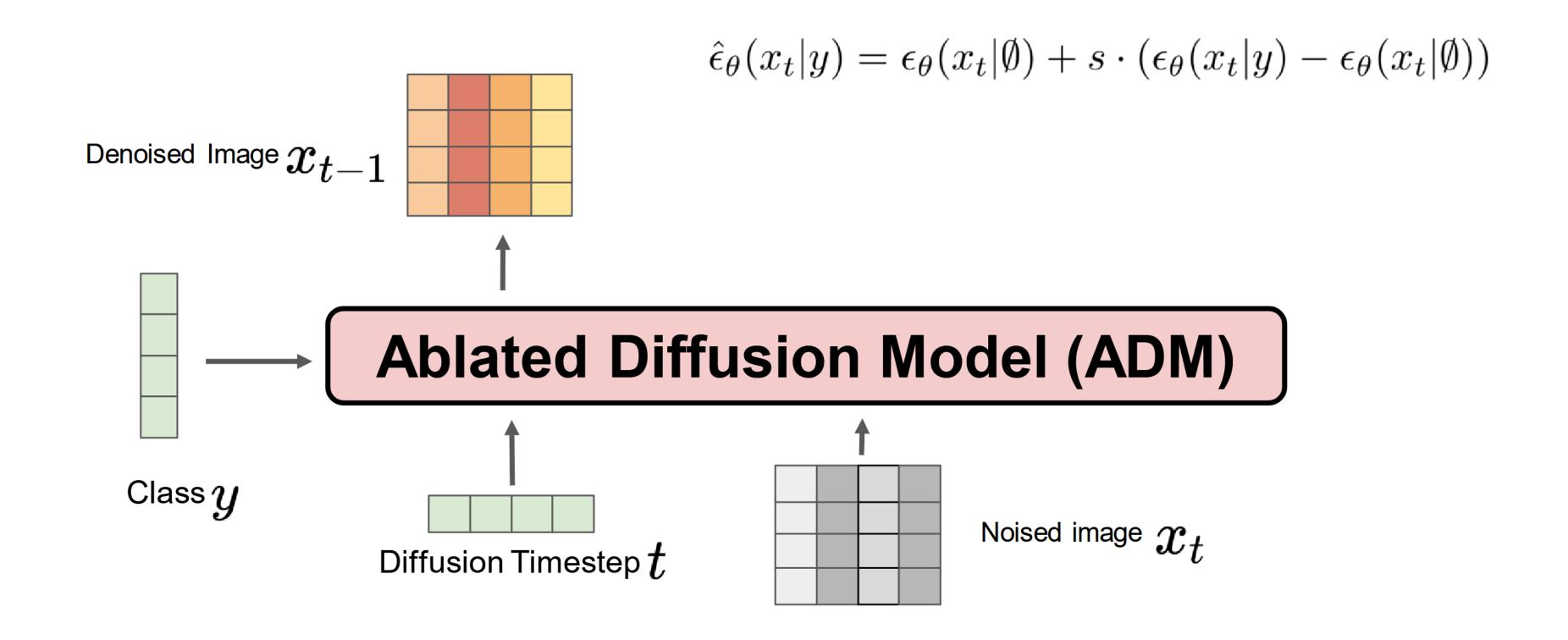
It	ters	T	Schedule	Objective	NLL	FID
	00K 00K	1K 4K	linear linear	$L_{ m simple} \ L_{ m simple}$	3.99 3.77	32.5 31.3
20 20	00K 00K 00K 00K	4K 4K 4K 4K	linear cosine cosine cosine	$L_{ m hybrid} \ L_{ m simple} \ L_{ m hybrid} \ L_{ m vlb}$	3.66 3.68 3.62 3.57	32.2 27.0 28.0 56.7
	.5M .5M	4K 4K	cosine cosine	$L_{ m hybrid} \ L_{ m vlb}$	3.57 3.53	19.2 40.1

Ablating schedule and objective on ImageNet 64 × 64.



Class-conditional ImageNet 64×64 samples generated using 250 sampling steps from Lhybrid model (FID 2.92) which shows high diversity

Conditioned Diffusion Model



Recall that when sampling from a conditional model $p(\mathbf{x}|\mathbf{y})$ we basically need an estimation of $abla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y})$

Using Bayes' rule, we have:

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) = \nabla_{\mathbf{x}_t} \log \frac{p(\mathbf{y} | \mathbf{x}_t) p(\mathbf{x}_t)}{p(\mathbf{y})} = \nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$$

Classifier gradient

Introduce a guidance scale (w):

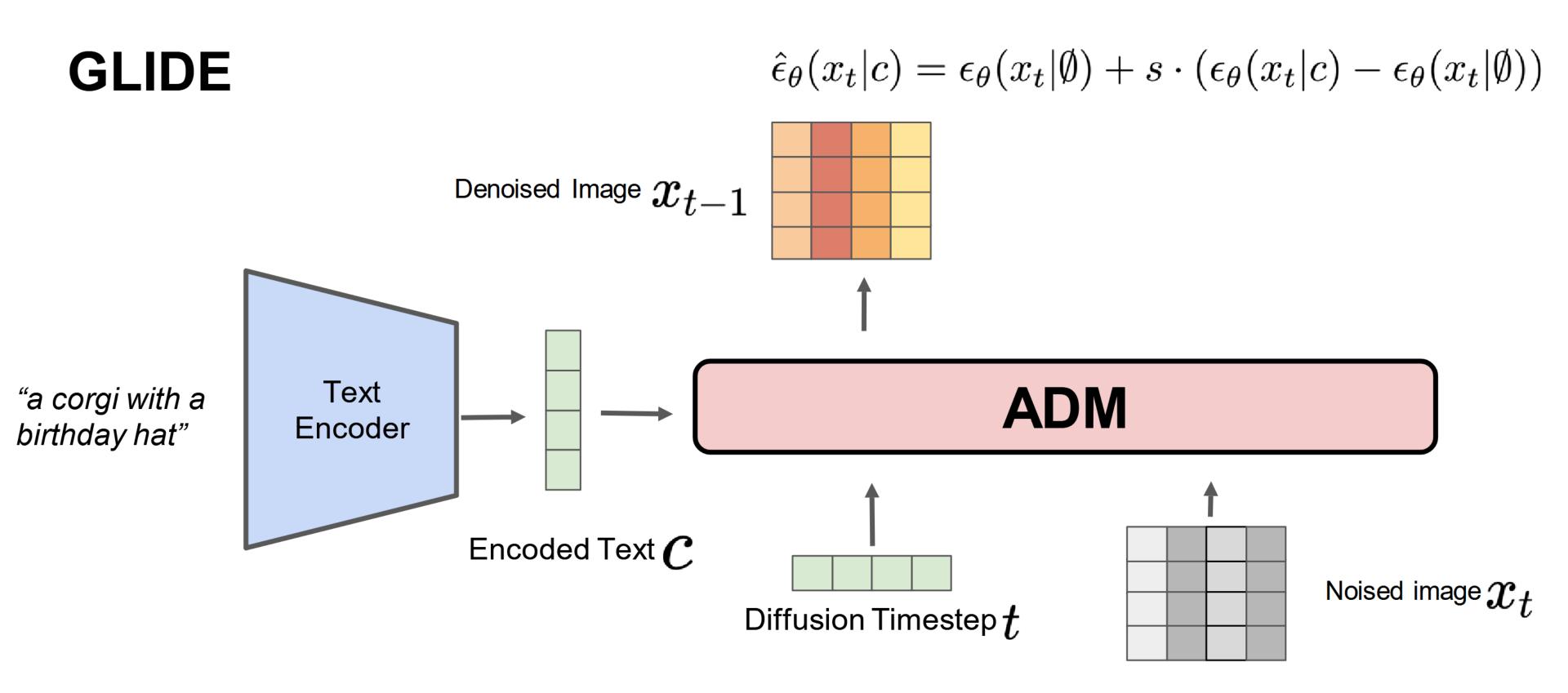
$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}) \approx \mathbf{w} \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$$

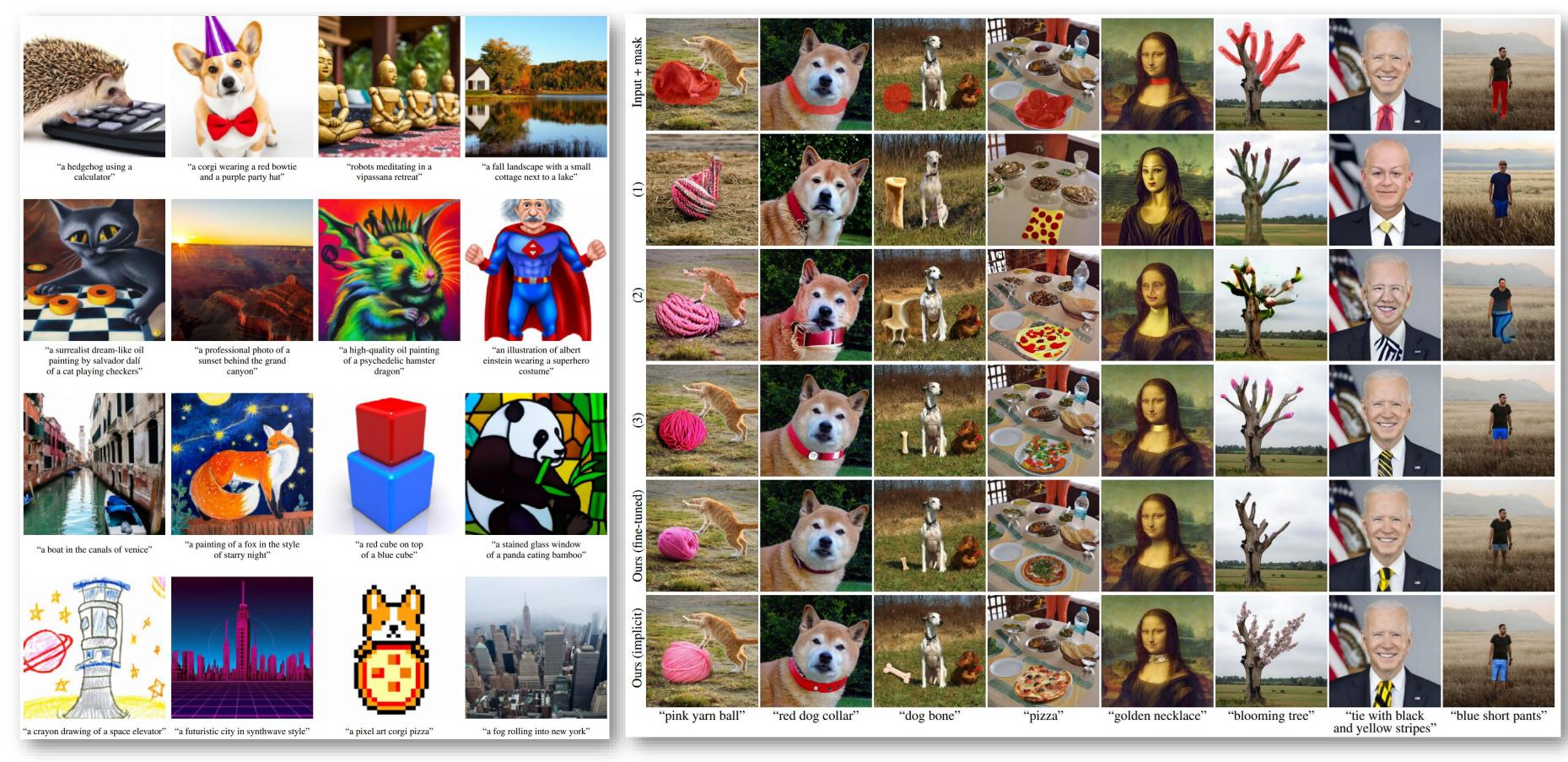




Samples from an unconditional diffusion model with classifier guidance to condition on the class "Pembroke Welsh corgi". Using classifier scale 1.0 (left; FID: 33.0) does not produce convincing samples in this class, whereas classifier scale 10.0 (right; FID: 12.0) produces much more class-consistent images.

Score model

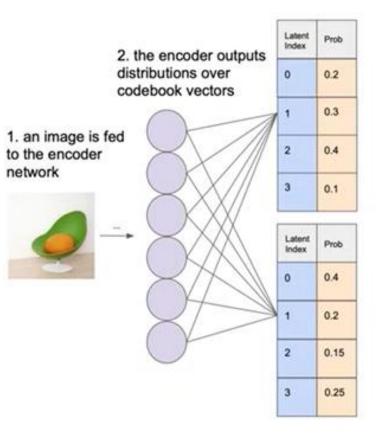




Photorealistic Image Samples.

Comparison of Image Inpainting Quality on Real Images.

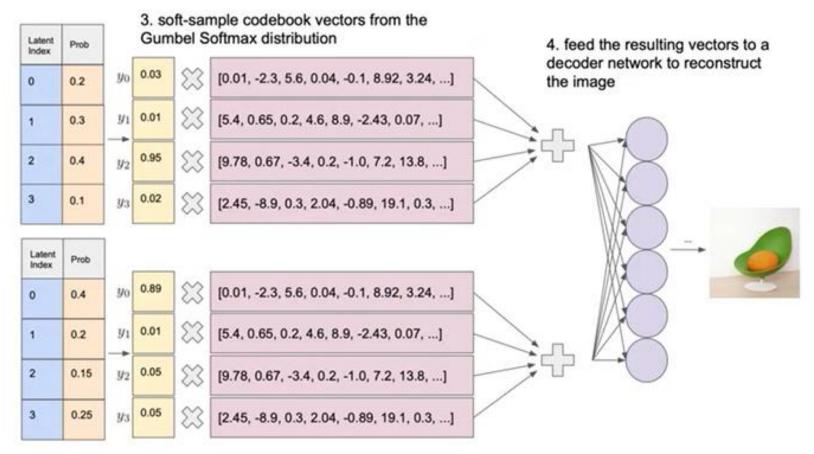
Stage 1



Latent Index	Codebook Vector
0	[0.01, -2.3, 5.6, 0.04, -0.1, 8.92, 3.24,]
1	[5.4, 0.65, 0.2, 4.6, 8.9, -2.43, 0.07,]
2	[9.78, 0.67, -3.4, 0.2, -1.0, 7.2, 13.8,]
3	[2.45, -8.9, 0.3, 2.04, -0.89, 19.1, 0.3,]

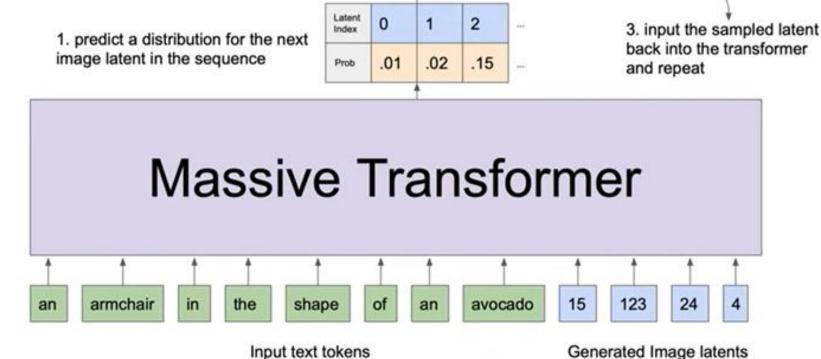
2. sample a latent

from this distribution



Stage 2





5. out pops an image of an avocado chair!

4. feed codebook vectors into the dVAE decoder

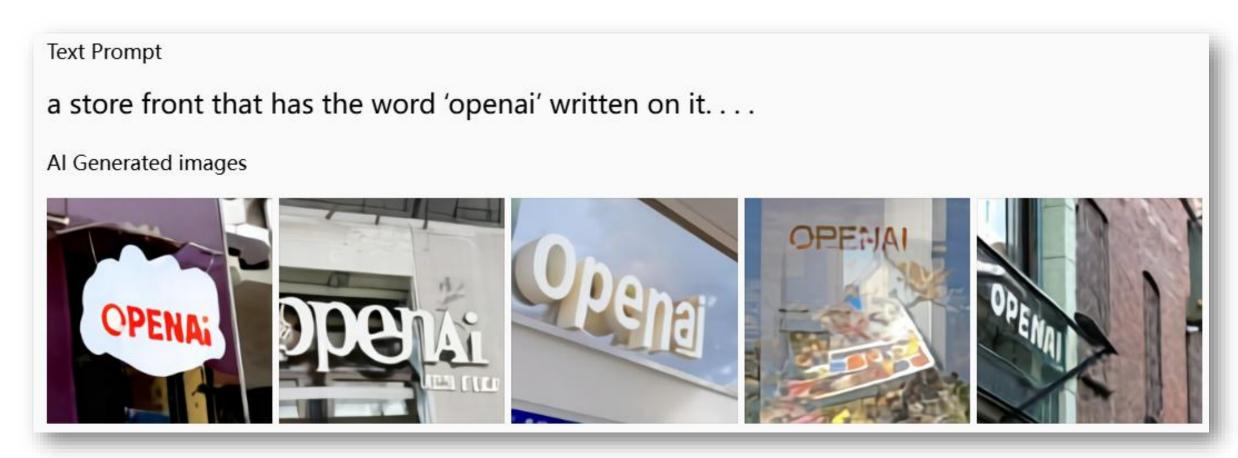
look up vectors in dVAE codebook

15 123 24 4 2 6 ...

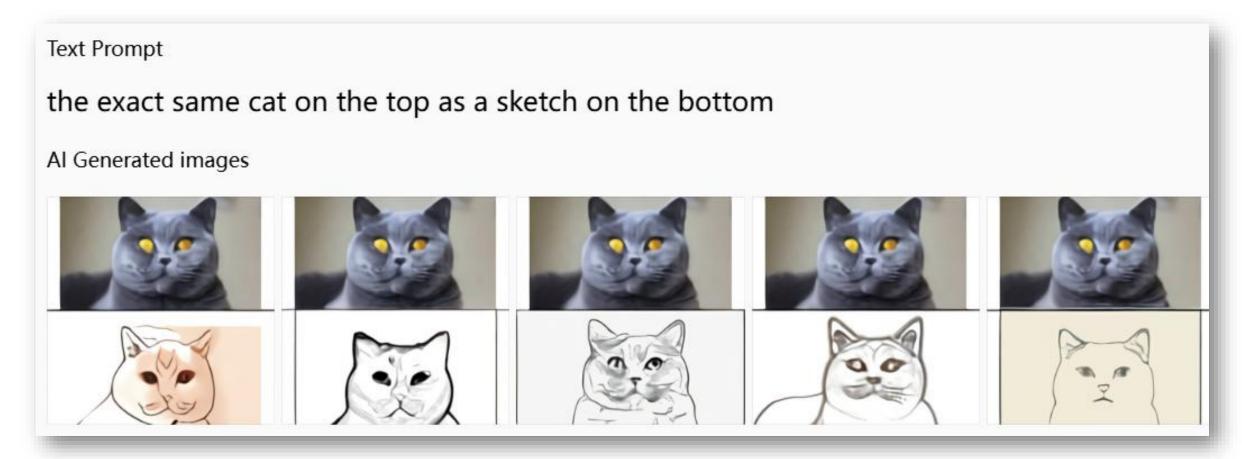
Ramesh et al. "Zero-Shot Text-to-Image Generation." ICML 2021

*This slide is derived https://www.youtube.com/watch?v=oENCNi4JxPY

Image Credit: https://openai.com/index/dall-e/

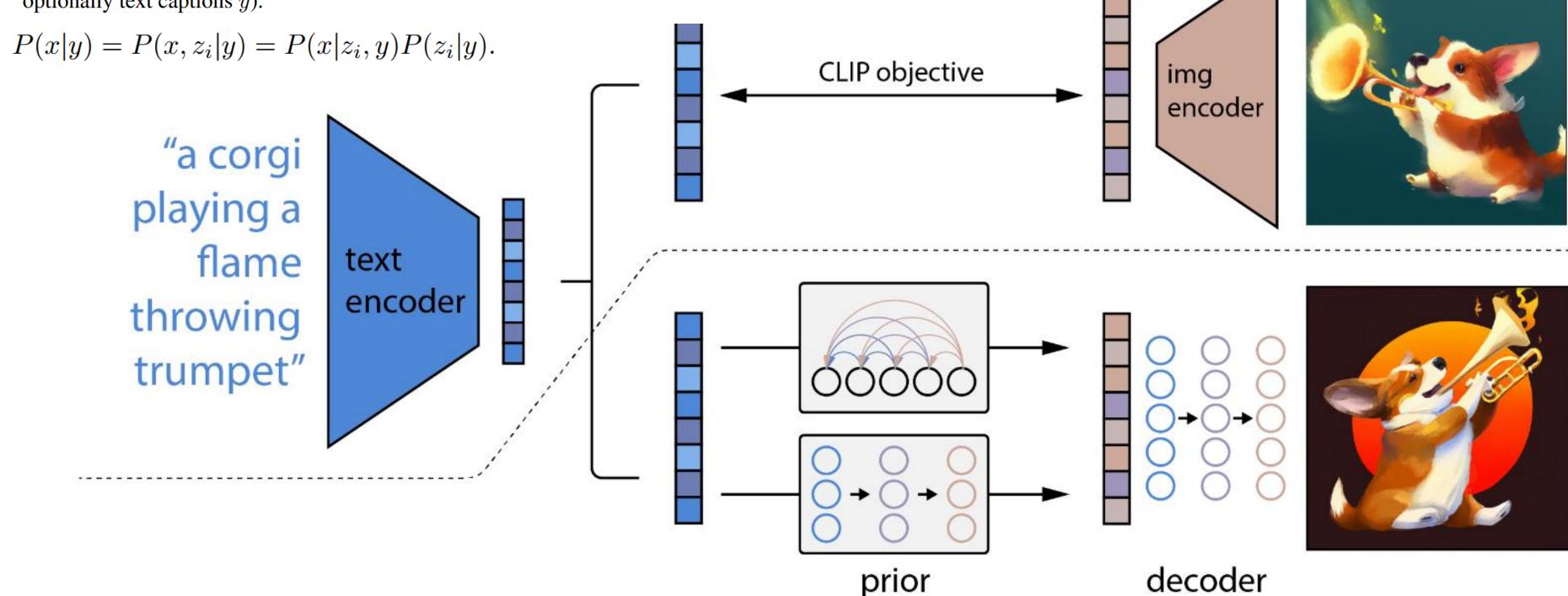


DALL-E Capabilities: Inferring contextual details

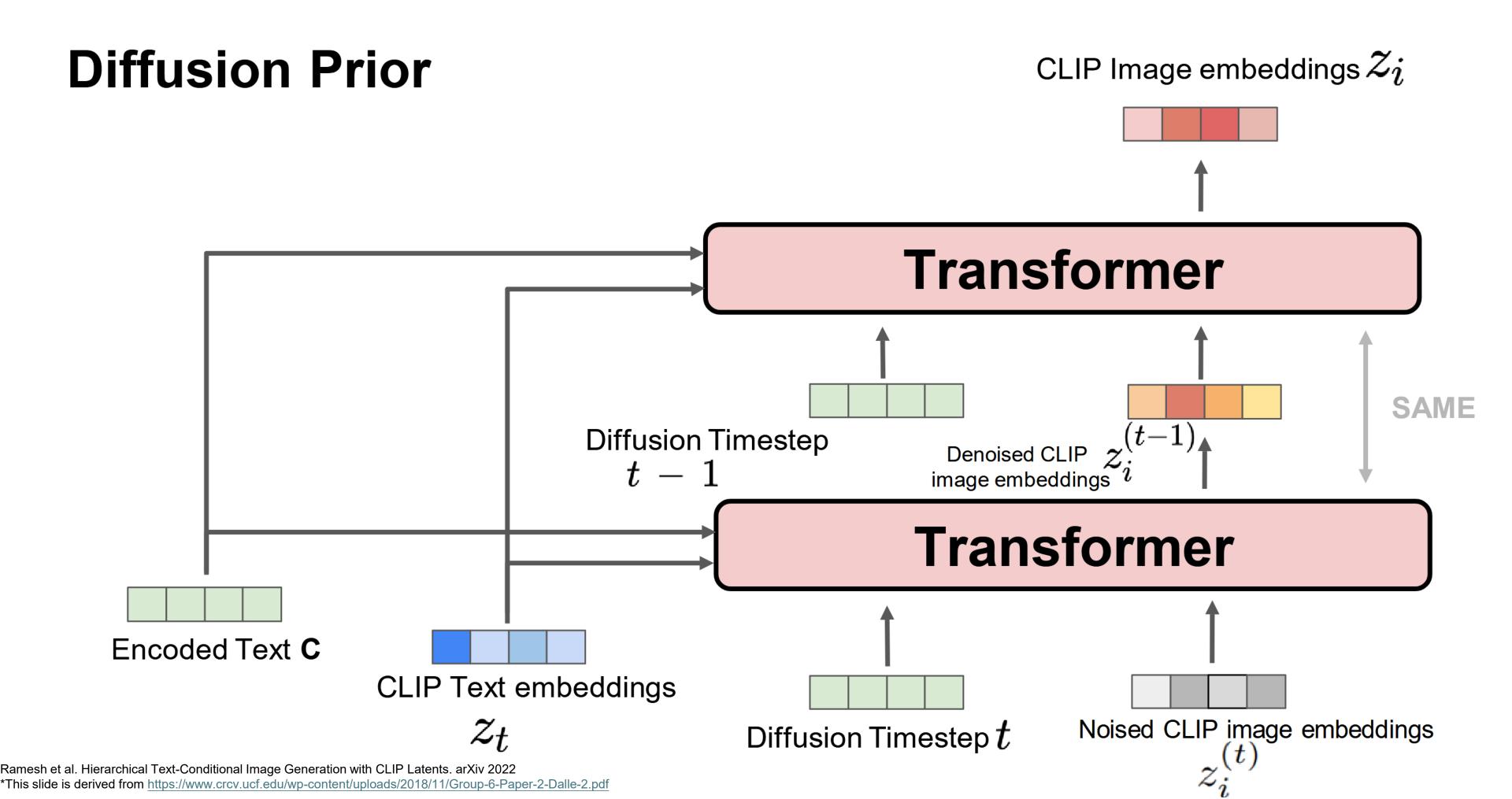


DALL-E Capabilities: Zero-shot visual reasoning (cherry-picked)

- A prior $P(z_i|y)$ that produces CLIP image embeddings z_i conditioned on captions y.
- A decoder $P(x|z_i, y)$ that produces images x conditioned on CLIP image embeddings z_i (and optionally text captions y).



Below the dotted line, we depict our text-to-image generation process: a CLIP text embedding is first fed to an autoregressive or **diffusion prior** to produce an image embedding, and then this embedding is used to condition a diffusion decoder which produces a final image.



Why the prior matters?

Condition decoder on captions alone



Condition decoder on Caption + text embedding impersonating image embeddings



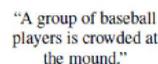
Prior + CLIP image embedding







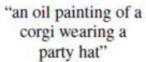








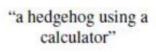




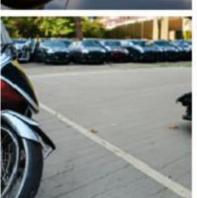




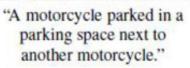
















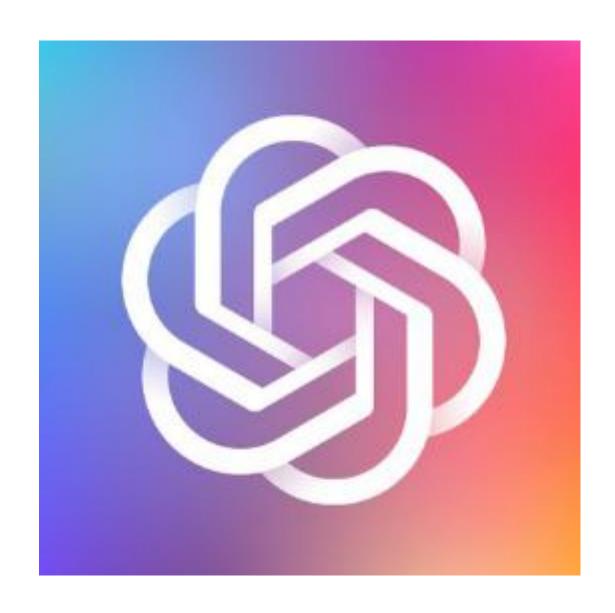


"This wire metal rack holds several pairs of shoes and sandals"



a photo of a landscape in winter \rightarrow a photo of a landscape in fall

Text diffs applied to images by **interpolating** between their CLIP image embeddings and a normalised difference of the CLIP text embeddings produced from the two descriptions.



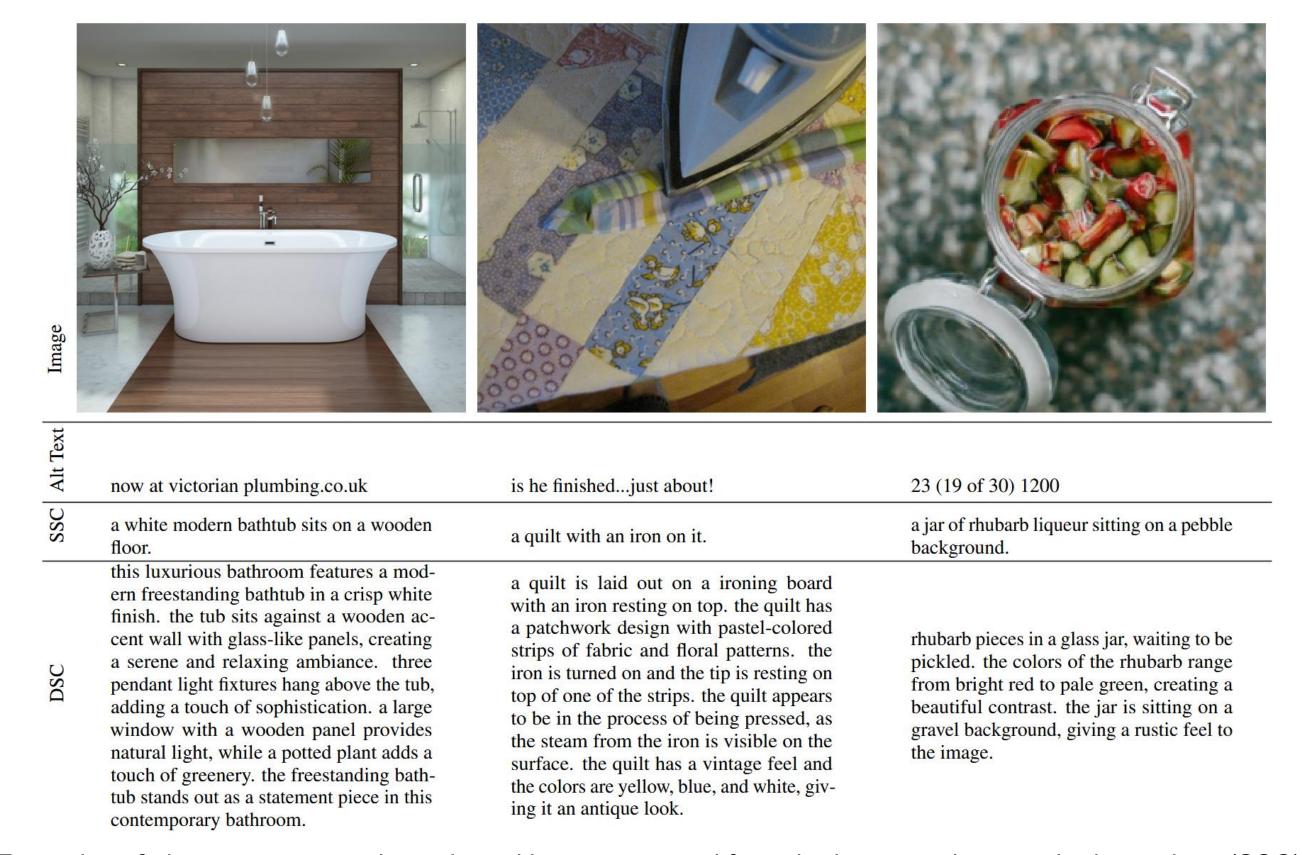


Variations of an input image by encoding with CLIP and then decoding with a diffusion model. The variations preserve both **semantic information** like the overlapping strokes in the logo, as well as **stylistic elements** like the color gradients in the logo, while varying the non-essential details.

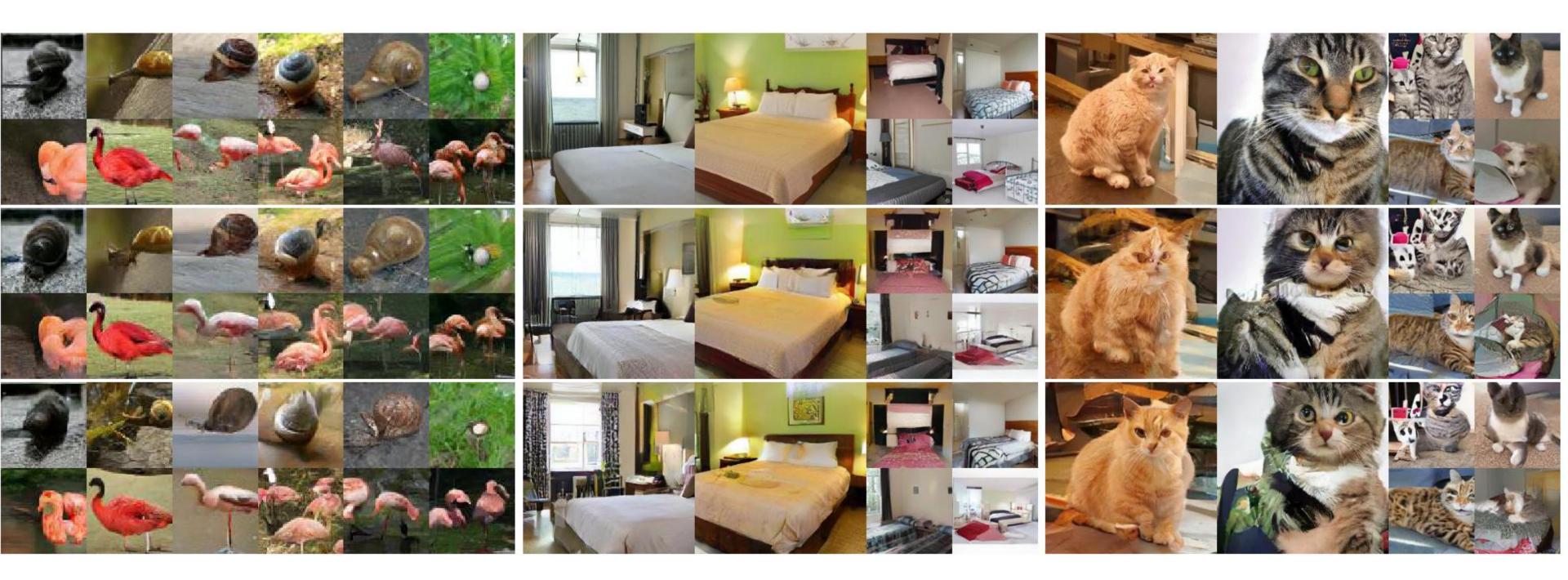


A illustration from a graphic novel. A bustling city street under the shine of a full moon. The sidewalks bustling with pedestrians enjoying the nightlife. At the corner stall, a young woman with fiery red hair, dressed in a signature velvet cloak, is haggling with the grumpy old vendor. the grumpy vendor, a tall, sophisticated man is wearing a sharp suit, sports a noteworthy moustache is animatedly conversing on his steampunk telephone.

Compelling results from DALL-E 3 with detailed text prompt.



Examples of alt-text accompanying selected images scraped from the internet, short synthetic captions (SSC), and descriptive synthetic captions (DSC).



Samples generated by EDM (top), **CT** + single-step generation (**middle**), and **CT** + 2-step generation (**Bottom**). All corresponding images are generated from the same initial noise.



(a) Left: The gray-scale image. Middle: Colorized images. Right: The ground-truth image.



(b) Left: The downsampled image (32×32). Middle: Full resolution images (256×256). Right: The ground-truth image (256×256).



(c) Left: A stroke input provided by users. Right: Stroke-guided image generation. Zero-shot image editing with a consistency model trained by consistency distillation on LSUN Bedroom 256× 256.

