

The Present and Future of Image and Video Generation Technologies

KAIST 소프트웨어대학원/KTAI 콜로кви엄

May 2025

성민혁 | KAIST 전산학부



Minhyuk Sung

- Associate Professor at KAIST
School of Computing
Affiliated with the Graduate School of AI
and with the Graduate School of Metaverse
2021 – Present
- Research Scientist at Adobe Research
2019 – 2020
- Ph.D. at Stanford University
Computer Science Department
2019

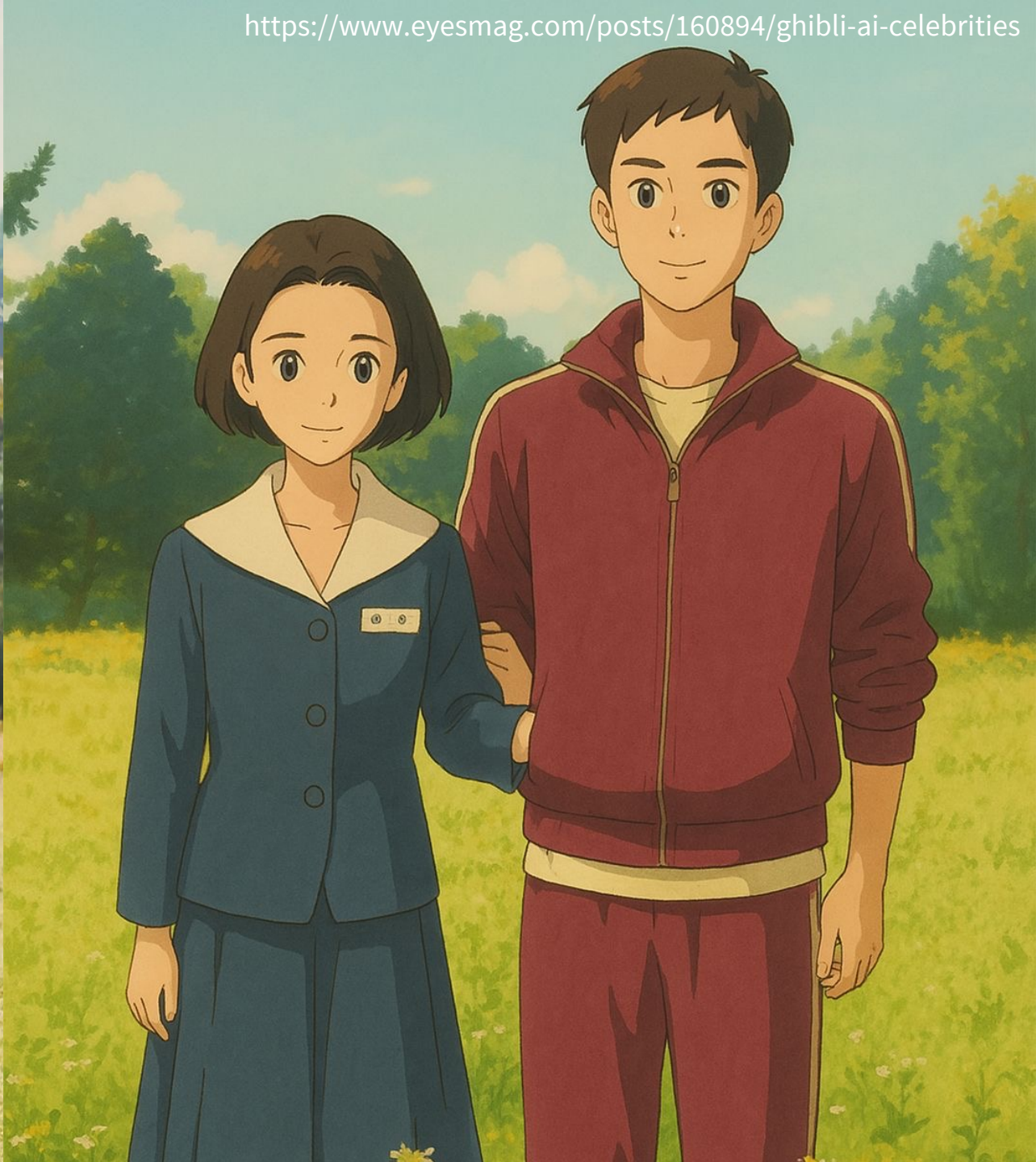
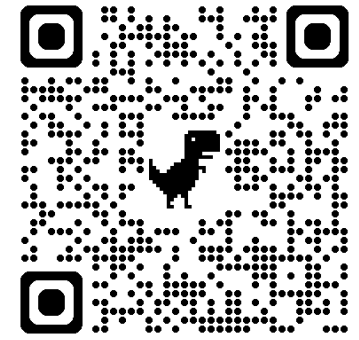




Image Generation Example

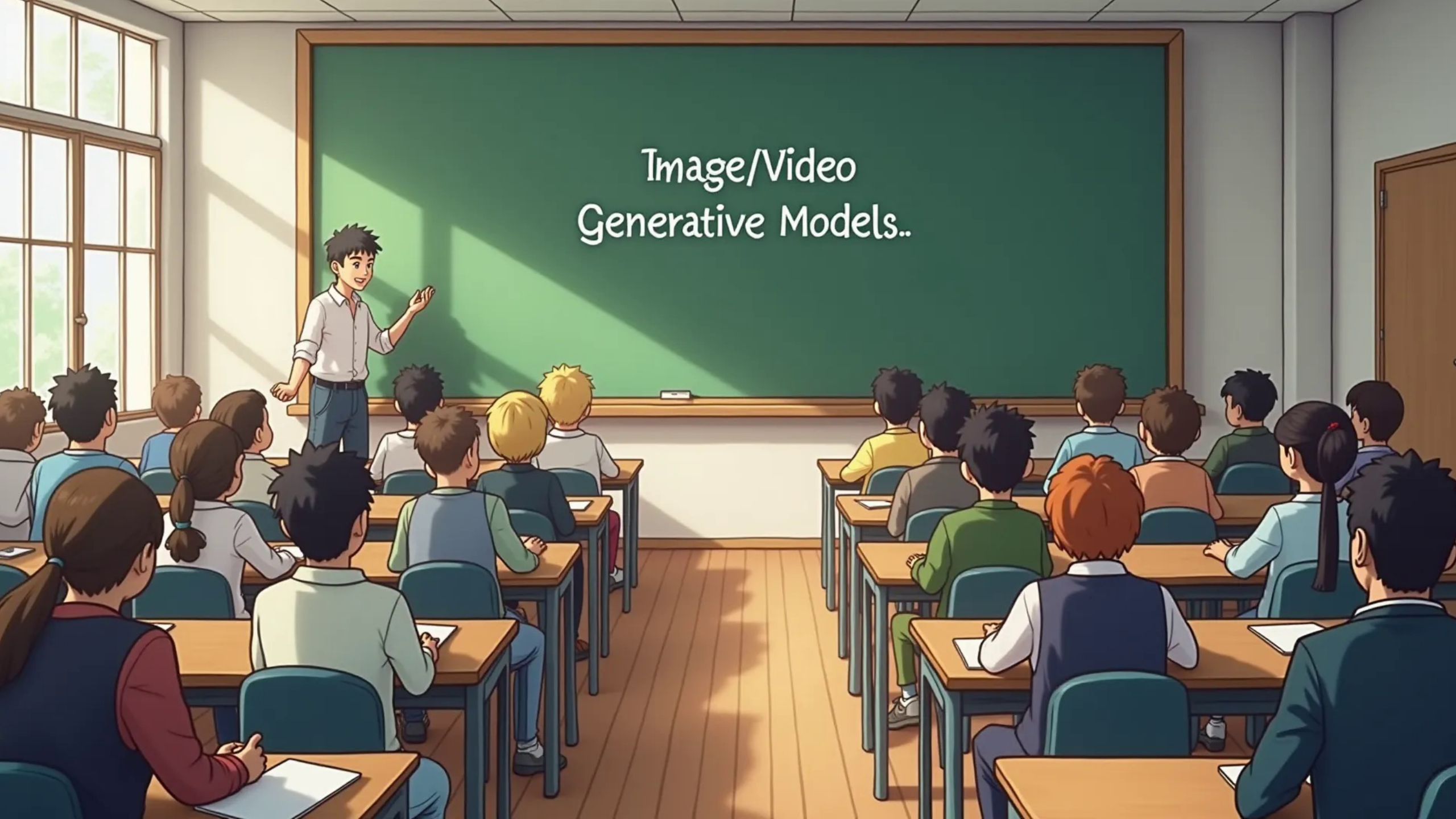


<https://huggingface.co/spaces/black-forest-labs/FLUX.1-dev>

<https://huggingface.co/spaces/black-forest-labs/FLUX.1-schnell>

A classroom for “Image/Video Generative Models” with a large audience. On the board, the text reads: “Image/Video Generative Models.” Create the image in Ghibli style.

Image/Video Generative Models..





February 2022



April 2022



July 2022



November 2022



March 2023



March 2023 (new version)



June 2023



December 2023

Midjourney | Prompt: "a hyper-realistic image of Harry Potter."

Analysis of Recent Image Generative Models

An Empirical Study of GPT-4o Image Generation Capabilities

**Sixiang Chen^{1*}, Jinbin Bai^{2*}, Zhuoran Zhao^{1*}, Tian Ye^{1*}, Qingyu Shi³, Donghao Zhou⁴, Wenhao Chai⁵,
Xin Lin⁶, Jianzong Wu³, Chao Tang³, Shilin Xu³, Tao Zhang⁶, Haobo Yuan⁶, Yikang Zhou⁶,
Wei Chow², Linfeng Li², Xiangtai Li^{3†}, Lei Zhu^{1,7†}, Lu Qi^{6†}**

¹The Hong Kong University of Science and Technology (GZ) ²National University of Singapore

³Peking University ⁴The Chinese University of Hong Kong ⁵University of Washington ⁶Wuhan University

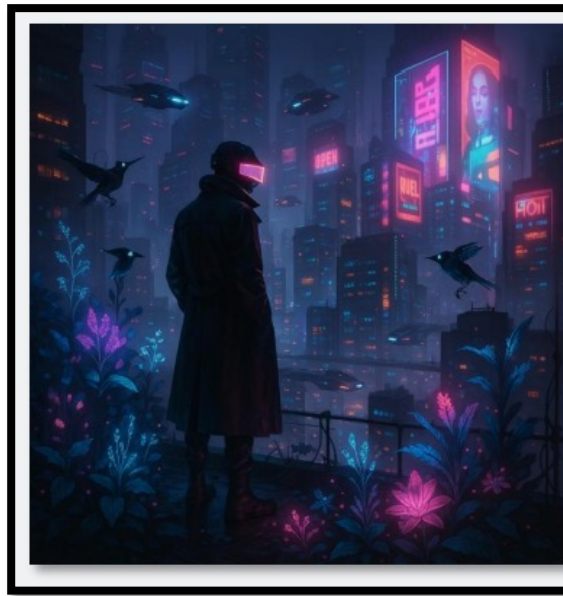
⁷The Hong Kong University of Science and Technology

Precise Text Alignment

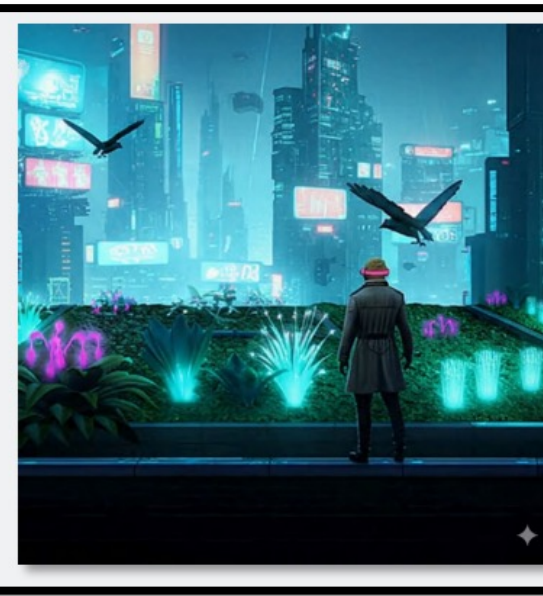
SD 2.1



GPT 4o



Gemini 2.0 Flash



FLUX



“On the rooftop of a skyscraper in a bustling cyberpunk city, a figure in a trench coat and neon-lit visor stands amidst a garden of bio-luminescent plants, overlooking the maze of **flying cars** and **towering holograms**. Robotic birds flit among the foliage, digital billboards flash advertisements in the distance.”

Object Count and Composition

SD 2.1



GPT 4o



Gemini 2.0 Flash



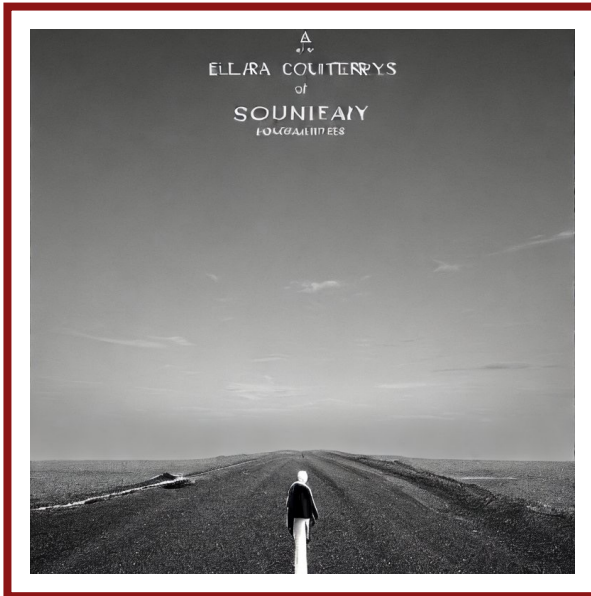
Midjourney v6.1



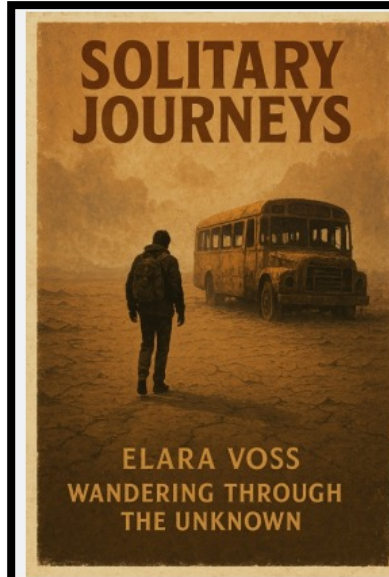
“Three differently colored apples (yellow, green, red from left to right) with a Coca-Cola bottle placed behind the middle apple.”

Text Rendering

SD 2.1



GPT 4o



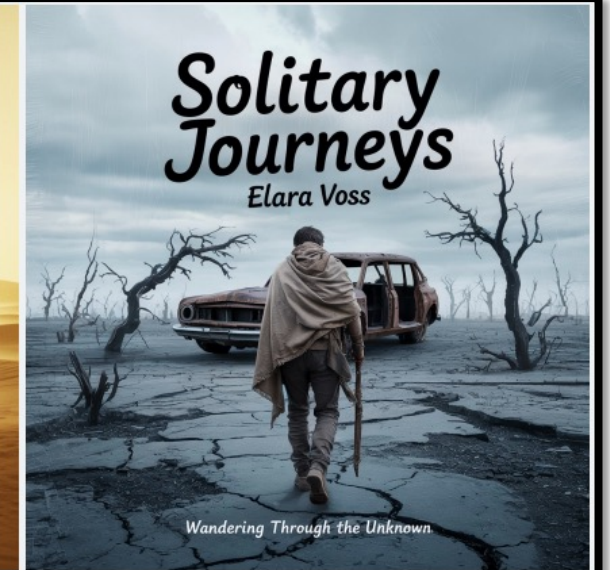
POSTA



Gemini 2.0 Flash

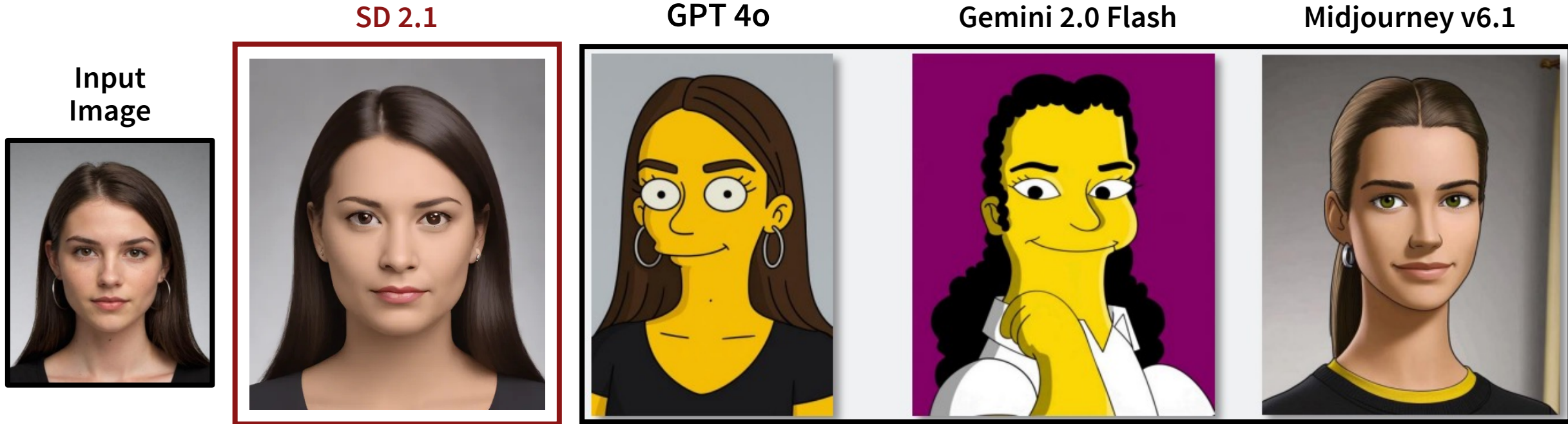


Ideogram 3.0



“Create a poster with the theme of a **Journey of Solitude**. The background should depict a lone figure walking toward an unusable form of transportation. The scene should evoke a sense of being lost, helplessness, and desolation, capturing the emotional weight of losing oneself in a barren, unforgiving landscape. Title: Solitary Journeys Subtitle: **Elara Voss** Information: **WANDERING THROUGH THE UNKNOWN.**”

Stylization



“Generate the Simpsons style of this picture.”





Google DeepMind Veo 2

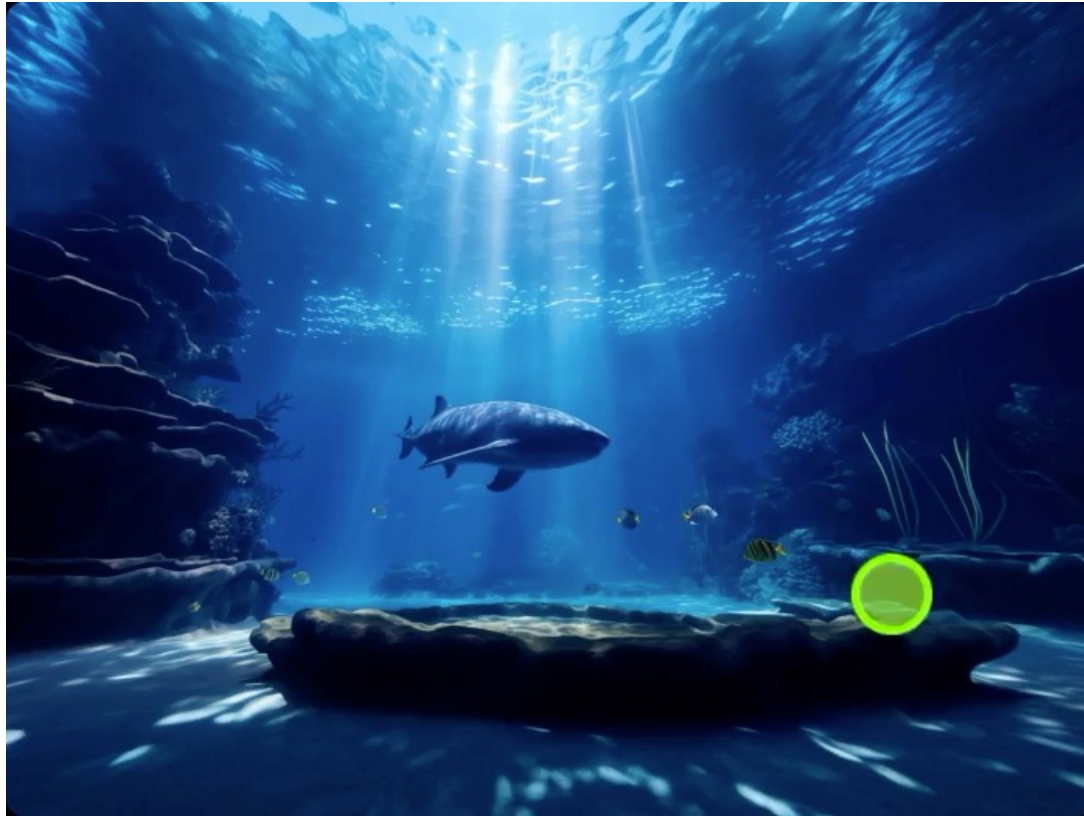
Prompt: A cinematic, high-action tracking shot follows an incredibly cute dachshund wearing swimming goggles as it leaps into a crystal-clear pool. The camera plunges underwater with the dog, capturing the joyful moment of submersion and the ensuing flurry of paddling with adorable little paws. Sunlight filters through the water, illuminating the dachshund's sleek, wet fur and highlighting the determined expression on its face. The shot is filled with the vibrant blues and greens of the pool water, creating a dynamic and visually stunning sequence that captures the pure joy and energy of the swimming dachshund.



KlingAI

Prompt: Rocks, spheres flying around a silver mirrored moon, digital art style, Clement's Archer, Felicia Simon, Vray, desert-wave, Daniel Anshan.

Control



Motion Brush



Frame Interpolation



Jason Allen won the digital-art competition at the Colorado State Fair last year for his piece "Théâtre D'opéra Spatial" that he created using the AI software Midjourney. Recently, the US Copyright Office refused to grant him a copyright for his piece, writing, "We have decided that we cannot register this copyright claim because the deposit does not contain any human authorship." He plans to appeal.



AIFF 2025

Finalists Announced

[View Finalists](#)

June 5th in New York at Alice Tully Hall | June 12th in LA at The Broad Stage Theater

Presenting Partners

TRIBECA
FESTIVAL

IMAX

CapCut

the gotham

.monks

NVIDIA

VIVA
TECHNOLOGY

GENERAL
ATLANTIC

Goldman
Sachs

runway

FACE

AI FILM FESTIVAL

https://www.instagram.com/runwayapp/reel/DJZb3sIKYha/?locale=ja_JP&hl=en

runway



[Home](#) [News](#) [Sport](#) [Business](#) [Innovation](#) [Culture](#) [Arts](#) [Travel](#) [Earth](#) [Audio](#) [Video](#) [Live](#)

AI was enemy No. 1 during Hollywood strikes. Now it's in Oscar-winning films

31 March 2025

Share  Save 

Regan Morris
BBC News, Los Angeles



Getty Images

Image/Video Generative Models



Image/Video Generative Models

Top Image Generative AI Companies (2025)		
The image generative AI sector is highly competitive, with several companies leading in technology, market share, and innovation. Below are the top companies shaping the field in 2025, based on market share, funding, and industry influence:		
Company	Notable Product/Model	Key Details & Market Position
Black Forest Labs	Flux family	Emerged as a market leader, holding close to 40% of image generation usage share in 2025 ⁹ ² . Known for advanced models and rapid ascent.
Google	Imagen3	Holds nearly 30% of usage share in 2025. Significant investment and rapid growth in image generation ⁹ ⁵ .
Stability AI	Stable Diffusion	Early leader in open-source image generation, still a major player with broad adoption ² ⁹ .
Midjourney	Midjourney	Popular for artistic and creative image generation, widely used among designers and artists ² .
Adobe	Firefly, Photoshop AI	Dominates creative professional market, integrating AI into Creative Cloud suite ⁶ .
Meta	AI image tools on Facebook/Instagram	Leverages massive user base, integrates AI image generation into social media platforms ⁶ ⁴ .
OpenAI	DALL-E series	Pioneered text-to-image generation, remains influential with DALL-E 3 and ongoing model improvements ⁷ .
Luma AI	Dream Machine	Focuses on 3D and high-fidelity image/video generation, growing rapidly in creative sectors ² .
Leonardo AI	Alchemy Refiner, Image Generator	Offers advanced content creation tools, gaining traction in creative industries ² .
Runway	Gen-1, Gen-2, Gen-3, Gen-4	Known for generative video and image tools, widely used in media and entertainment ⁷ .

Top Video Generative AI Companies in 2025		
The video generative AI landscape in 2025 is led by several innovative companies offering advanced tools for content creation, editing, and automation. The following companies are consistently ranked among the top in the industry:		
Company	Notable Product/Platform	Key Features and Highlights
Synthesia	Synthesia Studio	Studio-quality AI avatar videos, 230+ avatars, 140+ languages, widely used in enterprise and training ¹ ⁵ .
Runway	Gen-3, Gen-4, Runway ML	Next-gen AI video editing and generation, popular for creative and professional applications ¹ ⁵ ⁶ ⁹ .
Pika	Pika Art	AI-powered video editing and generation from captions and images, rapid user growth ¹ ⁴ ⁹ ⁵ ⁷ .
InVideo	InVideo AI	Full-length video creation, AI scriptwriting, media selection, accessible for creators ⁴ ⁹ ⁵ .
Kling	Kling AI Video Generator	High-quality motion video generation, gaining attention for realism and innovation ¹ ⁶ ⁹ .
Google	Veo (Veo 2)	Advanced generative video with strong physics and realism, limited free access ¹ ⁶ ⁹ .
OpenAI	Sora	Text-to-video generation, visually compelling results, ongoing development ¹ ⁴ ⁹ ⁶ ⁹ .
Luma Labs	Dream Machine	Focus on 3D and realistic video generation, paid-only access, strong realism ¹ ⁶ ⁹ .

About

Midjourney is an independent research lab exploring new mediums of thought and expanding the imaginative powers of the human species.

We are a small self-funded team focused on design, human infrastructure, and AI. We have 11 full-time staff and an incredible set of advisors.

Executives

David Holz

*Previously: Founder Leap Motion, Researcher
at NASA, Max Planck*

Advisors

Jim Keller

*Lead Silicon at Apple, AMD, Tesla, Intel,
Coauthor X86-64, CTO Tenstorrent*

Philip Rosedale

Founder of Second Life, CTO RealNetworks

Nat Friedman

CEO Github, Chairman of GNOME Foundation

Bill Warner

*Founder of Avid Technology, inventor of
nonlinear video editing*

Briefing

AI Startup Midjourney Expects \$200 Million in Revenue



By Kate Clark

2 years ago

Source: The Information

Most Popular

The Briefing

Why the Charter-Cox
\$34.5 Billion Cable Deal
Matters

May 16, 2025

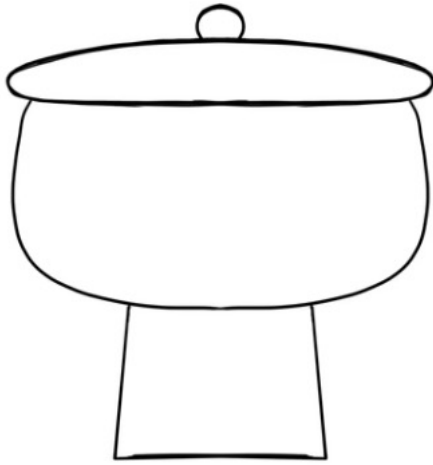
Midjourney, a generative artificial intelligence startup that charges users between \$10 and \$120 for a monthly subscription, is on pace to surpass \$200 million in revenue this year, according to The Information's [reporting](#). The company, which is built atop the gaming chat app Discord, is profitable and has not raised any money from venture capital firms since it was founded two years ago.

Notably, Midjourney shares profits with its employees and gives a slice of its revenues to Discord, which has helped the AI startup scale.

Midjourney founder David Holz said in an interview with The Informaiton that his goal is to build a company "kind of like Craigslist," the classifieds company that was famously bootstrapped. Holz added that Midjourney is "this weird thing that no one knows how to compete with that just sort of stands alone."

Adobe

How it Works



1. Sketch

Take a picture of your sketch, or draw within Vizcom's studio.



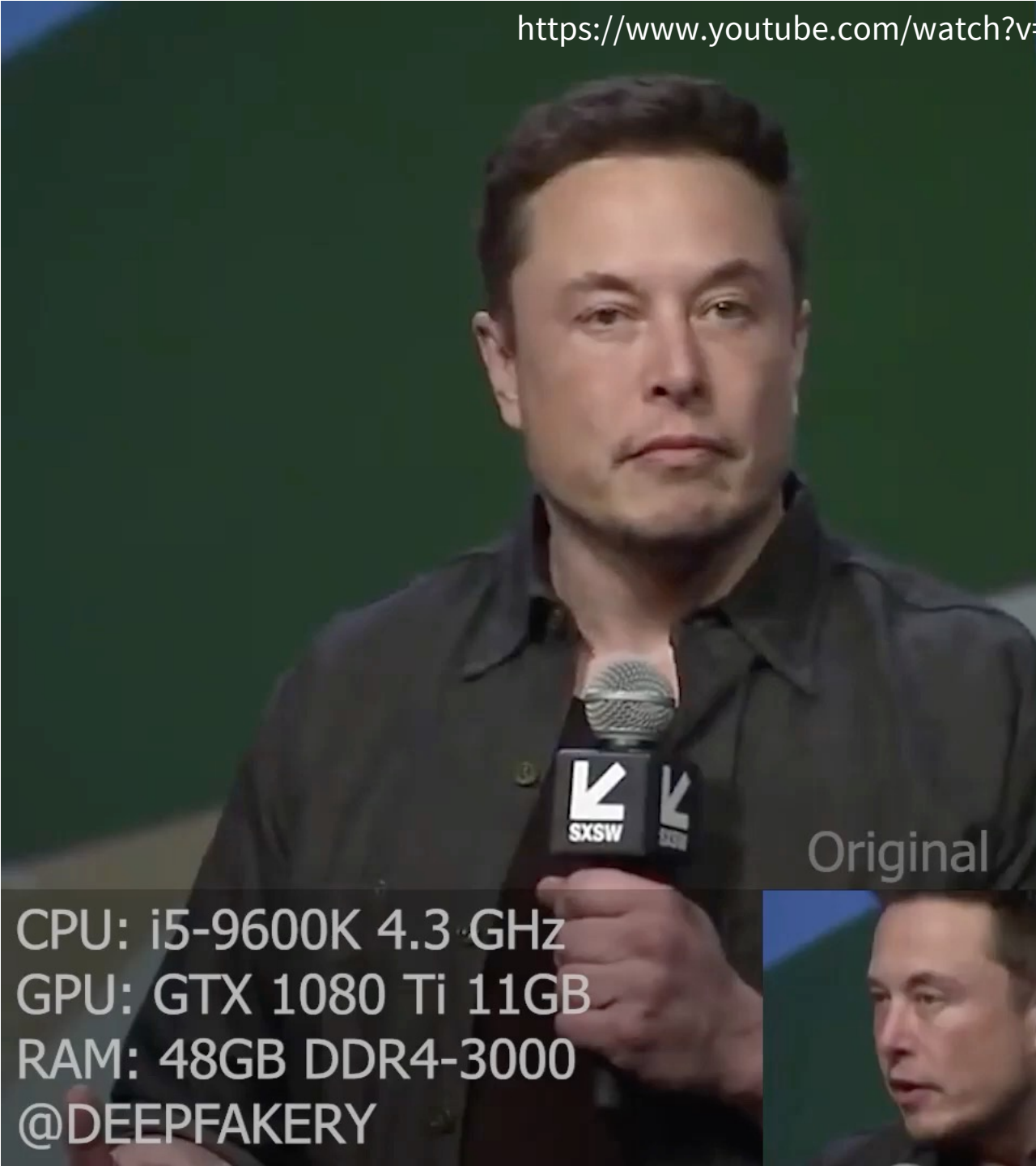
2. Render

Prompt and render your design with one of the default styles, or with your own Palette.



3. 3D Model

Generate a 3d model of your rendering to view your idea from a new angle. Export the model to view in AR or to 3d print.



Original

CPU: i5-9600K 4.3 GHz
GPU: GTX 1080 Ti 11GB
RAM: 48GB DDR4-3000
@DEEPPFAKERY



Deepfake

DeepFaceLab 2.0
Build: NVIDIA 6/27/2020
Model: Quick96
Iterations: 1,000,000



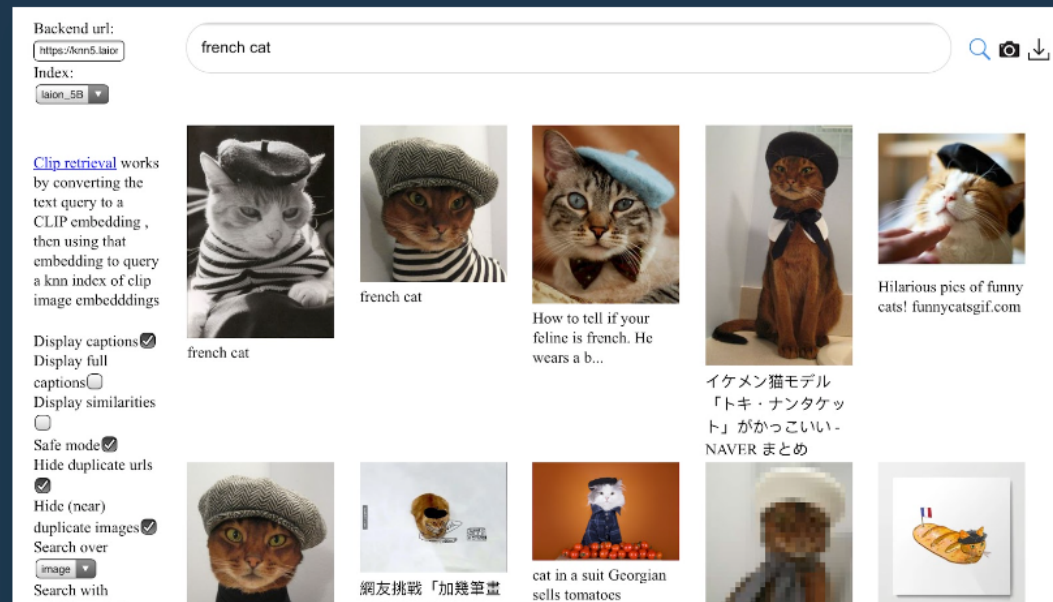
CREATED BY
**THE DOR
BROTHERS**

LAION-5B: A NEW ERA OF OPEN LARGE-SCALE MULTI-MODAL DATASETS

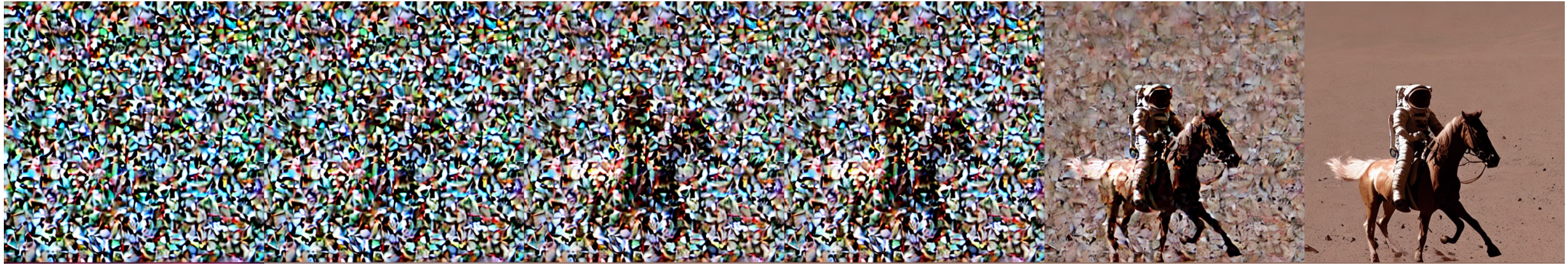
by: Romain Beaumont, 31 Mar, 2022

We present a dataset of **5.85 billion** CLIP-filtered image-text pairs, 14x bigger than LAION-400M, previously the biggest openly accessible image-text dataset in the world - see also our [NeurIPS2022 paper](#)

Authors: Christoph Schuhmann, Richard Vencu, Romain Beaumont, Theo Coombes, Cade Gordon, Aarush Katta, Robert Kaczmarczyk, Jenia Jitsev



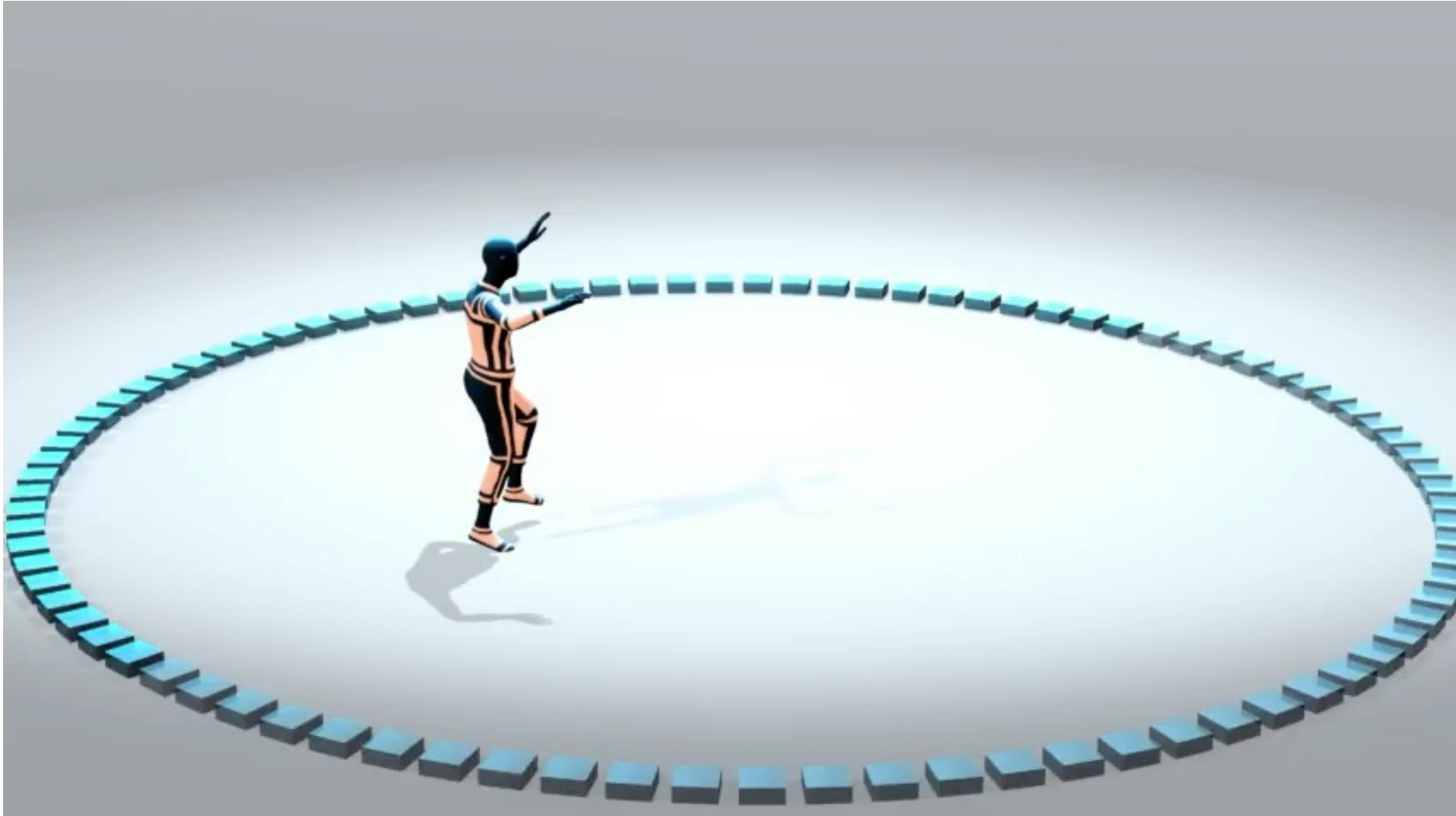
Diffusion Models



Audio Generation

Text input: A traditional Irish fiddle playing a lively reel.
Up Next: The sound of a light saber

Motion Generation



3D Generation



Molecule Generation



The image is a blue banner for the Nobel Prize in Chemistry 2024. At the top left is a gold Nobel medal. To its right, the text reads "NOBELPRISET I KEMI 2024" and "THE NOBEL PRIZE IN CHEMISTRY 2024". At the top right is the logo of the Kungl. Vetenskaps-Akademien (The Royal Swedish Academy of Sciences). Below the header, three portraits of the laureates are shown. Each portrait has a small vertical credit on its left: "Photo: University of Washington" for David Baker, "Photo: The Royal Society" for Demis Hassabis, and "Photo: BBSA Foundation" for John M. Jumper. Below each portrait is the laureate's name and affiliation. At the bottom, the award reasons are listed in Swedish and English. The text "#NobelPrize" is at the bottom left, and "THE NOBEL PRIZE" is at the bottom right.

NOBELPRISET I KEMI 2024
THE NOBEL PRIZE IN CHEMISTRY 2024

KUNGL. VETENSKAPS- AKADEMIEN
THE ROYAL SWEDISH ACADEMY OF SCIENCES

David Baker
University of Washington
USA

Demis Hassabis
Google DeepMind
United Kingdom

John M. Jumper
Google DeepMind
United Kingdom

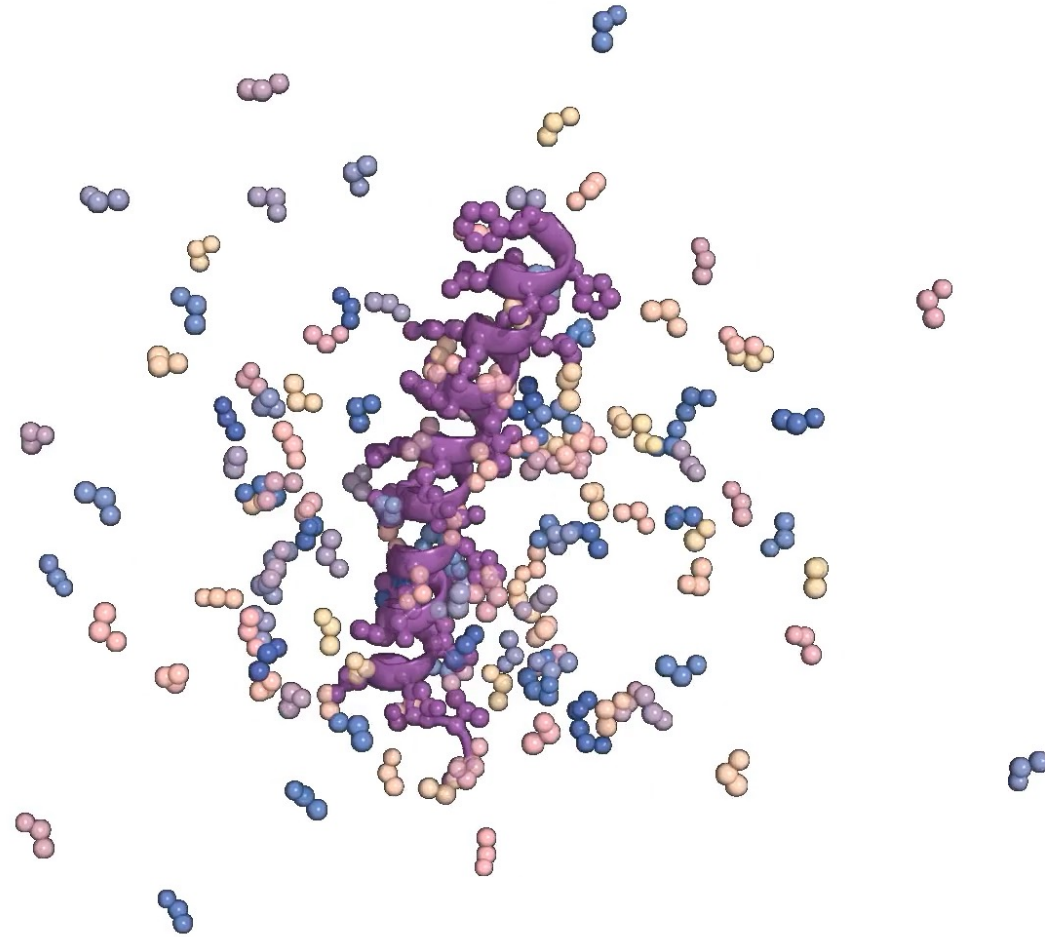
"för datorbaserad proteindesign"
"for computational protein design"

"för proteinstrukturprediktion"
"for protein structure prediction"

#NobelPrize

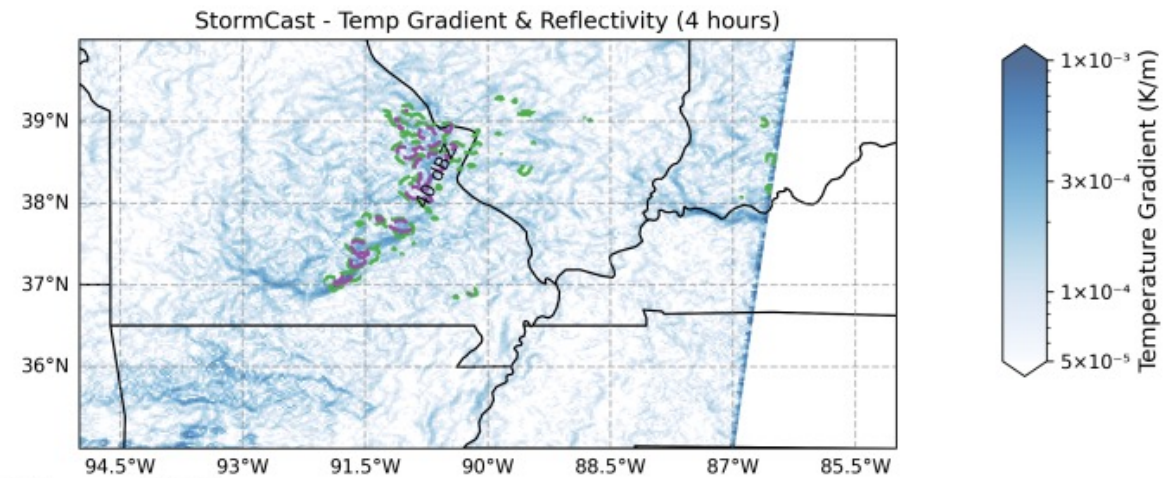
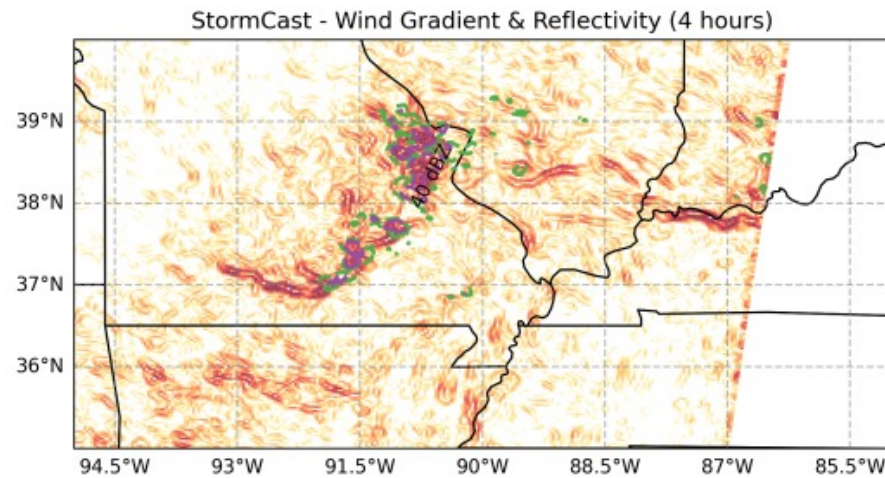
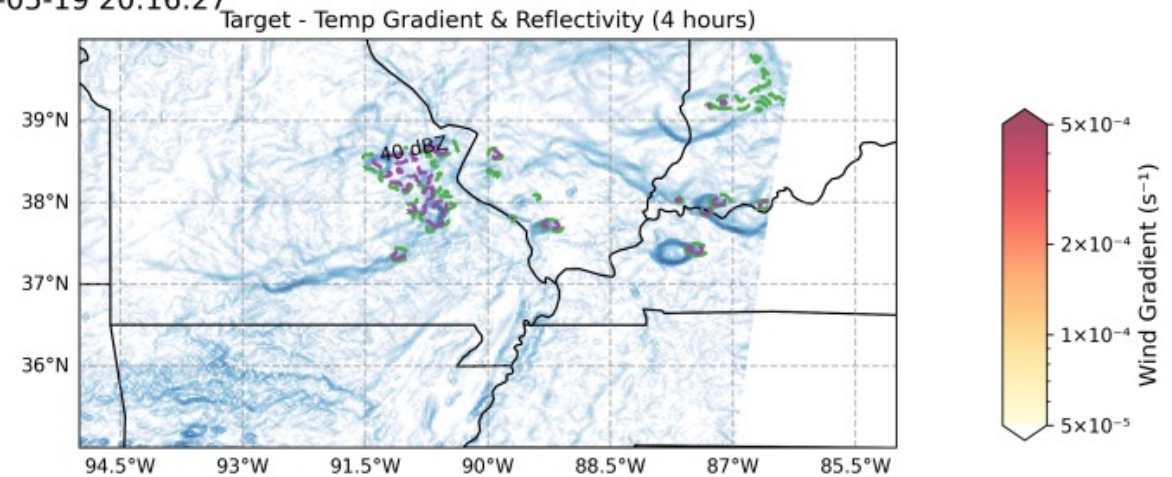
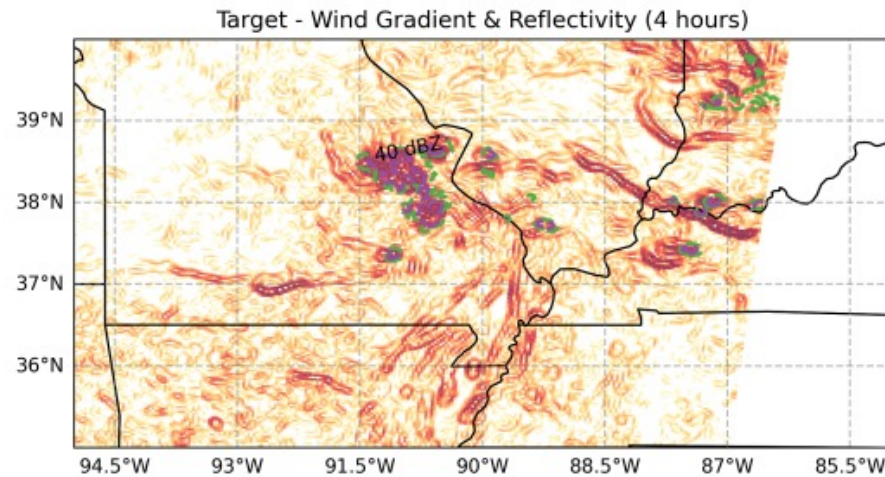
THE NOBEL PRIZE

Molecule Generation



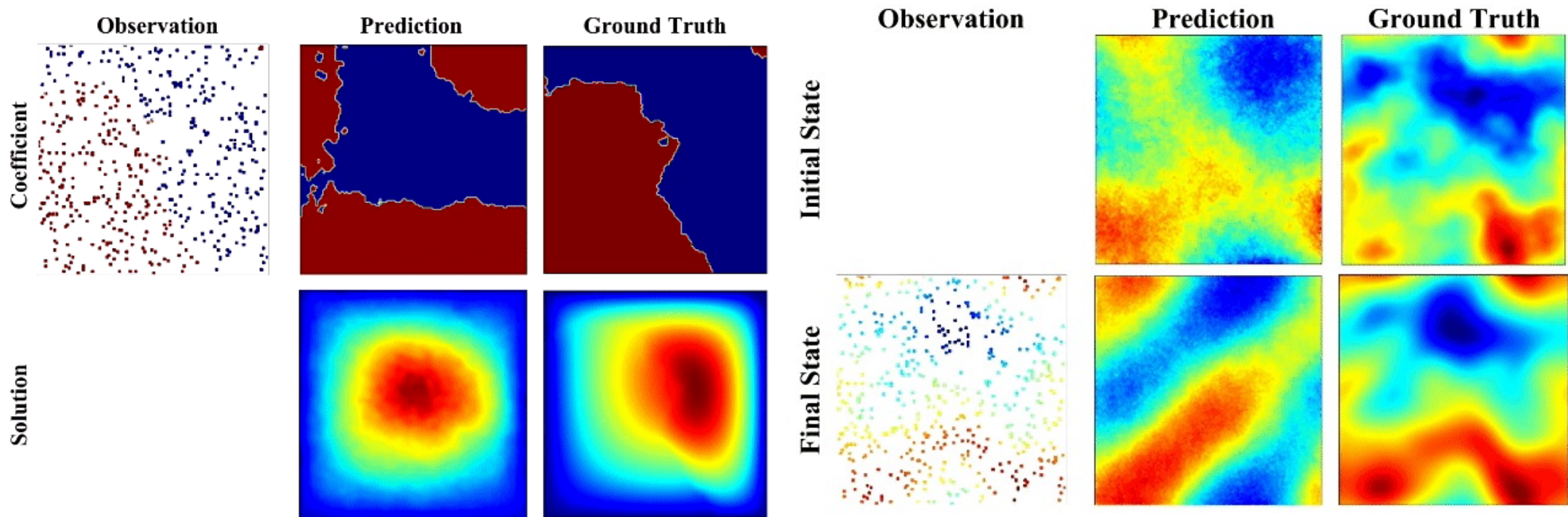
Weather Forecasting

2022-05-19 20:16:27



--- 40 dBZ --- 50 dBZ

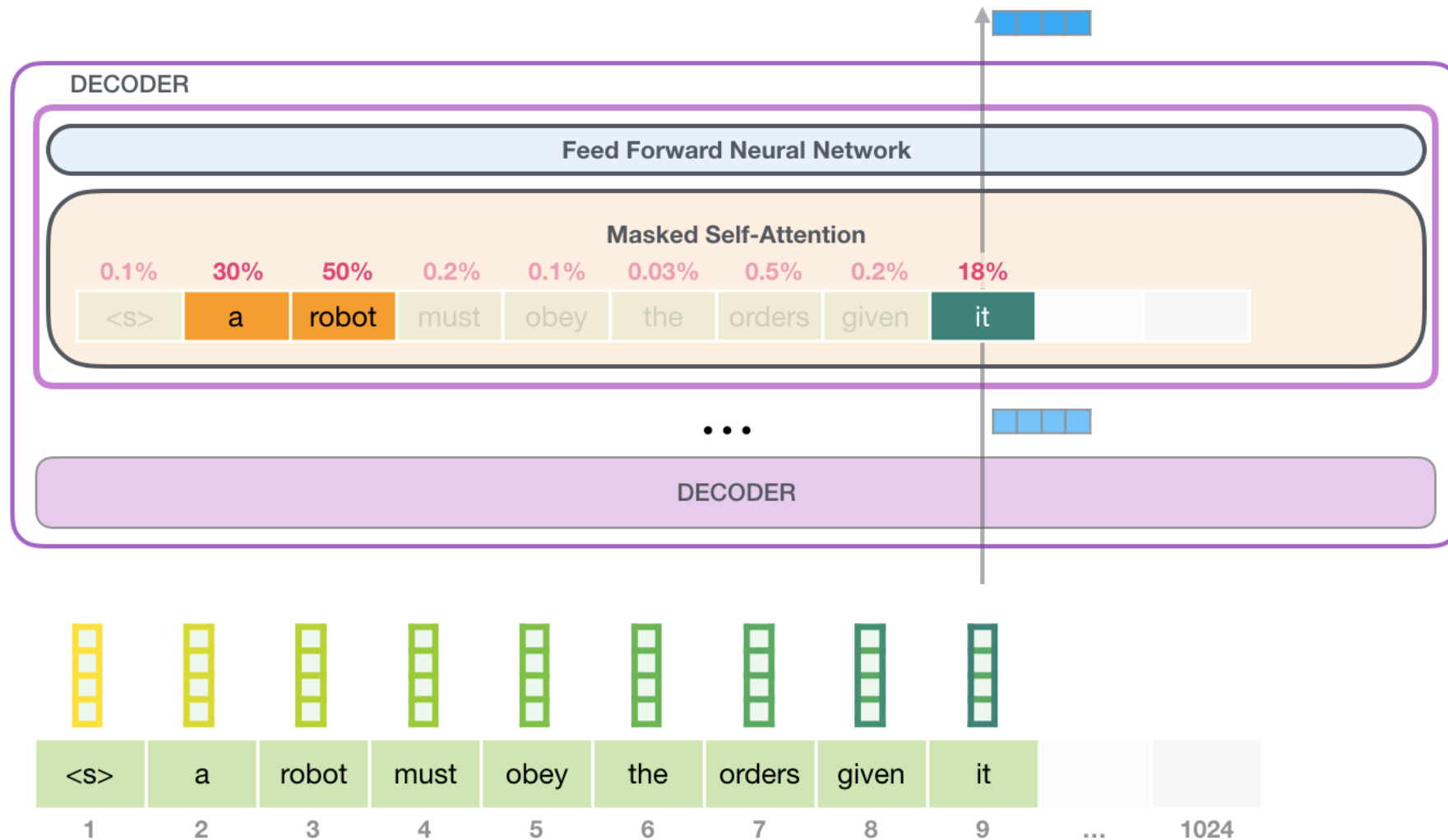
AI for Science – PDE Solving



How About Language Models?

Large Language Models (LLMs), such as GPT-3 and GPT-4, utilize a process called tokenization. Tokenization involves breaking down text into smaller units, known as tokens, which the model can process and understand. These tokens can range from individual characters to entire words or even larger chunks, depending on the model. For GPT-3 and GPT-4, a Byte Pair Encoding (BPE) tokenizer is used. BPE is a subword tokenization technique that allows the model to dynamically build a vocabulary during training, efficiently representing common words and word fragments. Although the core tokenization process remains similar across different versions of these models, the specific implementation can vary based on the model's architecture and training objectives.

Autoregressive Models



Tokenization










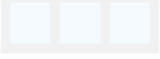

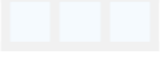

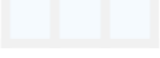

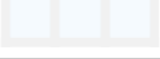





In computer science, we refer to human languages, like English and Mandarin, as "natural" languages. In contrast, languages designed to interact with computers, like Assembly and LISP, are called "machine" languages, following strict syntactic rules that leave little room for interpretation. While computers excel at processing their own highly structured languages, they struggle with the messiness of human language.

```
[637, 7595, 11222, 11, 581, 6716, 316, 5396, 22772, 11, 1299, 7725, 326, 133467, 11, 472, 392, 45497, 1, 22772, 13, 730, 20663, 11, 22772, 6884, 316, 20255, 483, 26971, 11, 1299, 15594, 326, 451, 96929, 11, 553, 4358, 392, 48082, 1, 22772, 11, 3992, 12035, 20449, 37247, 9607, 484, 7668, 3389, 3435, 395, 39300, 13, 8406, 26971, 19383, 540, 12323, 1043, 2316, 8916, 47557, 22772, 11, 1023, 24797, 483, 290, 13017, 1811, 328, 5396, 6439, 13]
```

Try yourself: <https://platform.openai.com/tokenizer>

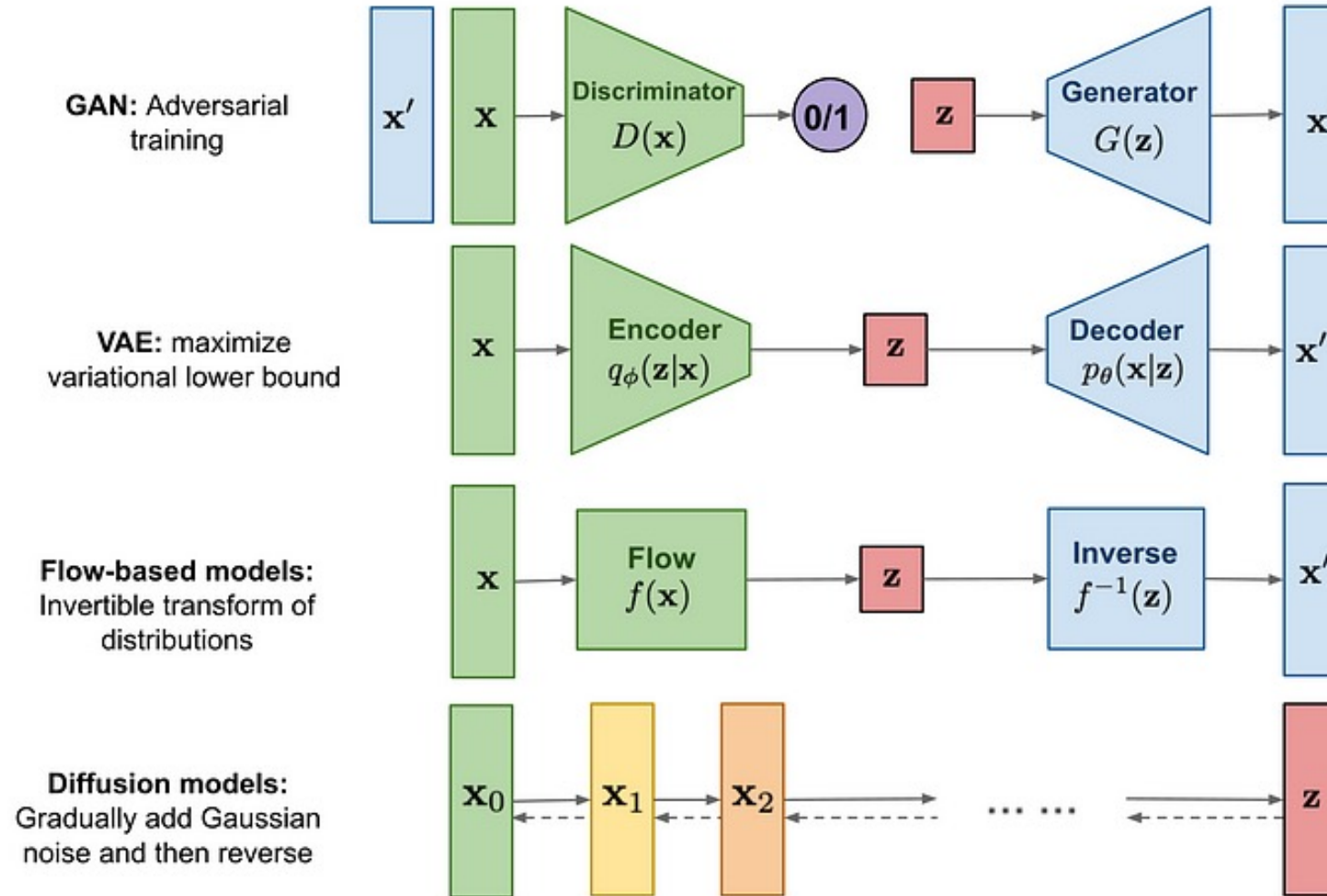
Learning Categorical Distributions

Word	Value vector	Score	Value X Score
<S>		0.001	
a		0.3	
robot		0.5	
must		0.002	
obey		0.001	
the		0.0003	
orders		0.005	
given		0.002	
it		0.19	
		Sum:	

Autoressive Models vs. Diffusion Models

- Autoregressive models are suitable for sequential discrete data.
- Diffusion models are suitable for continuous data in high-dimensional spaces.

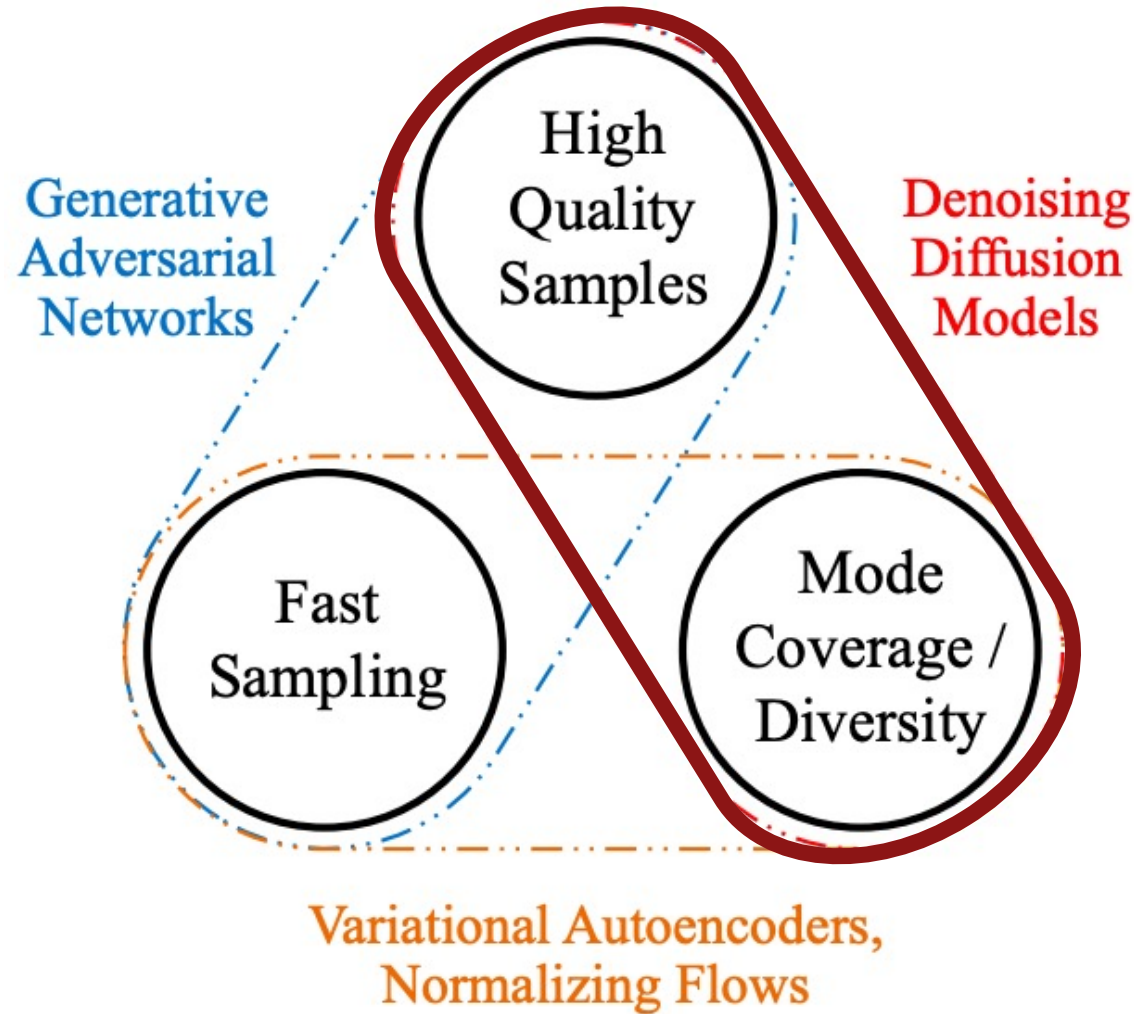
Generative Models



StyleGAN2



Generative Models – Comparison



Higher Quality & Diversity

Diffusion Models Beat GANs on Image Synthesis

Prafulla Dhariwal*
OpenAI
prafulla@openai.com

Alex Nichol*
OpenAI
alex@openai.com

Abstract

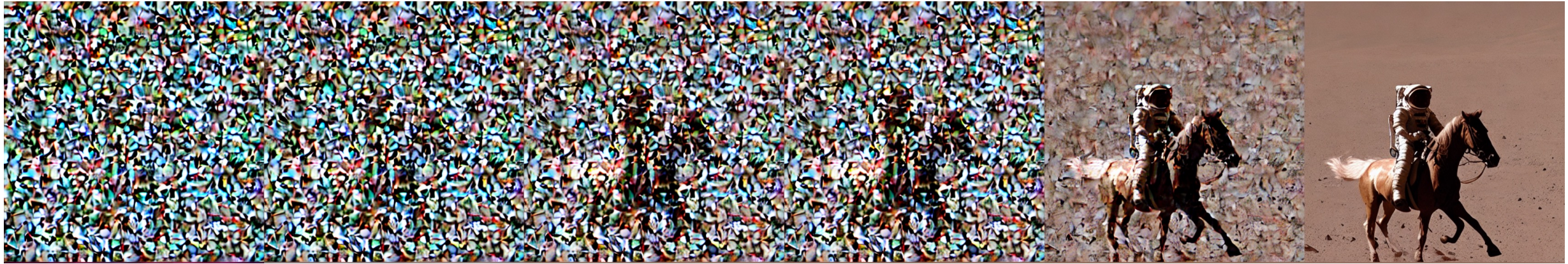
We show that diffusion models can achieve image sample quality superior to the current state-of-the-art generative models. We achieve this on unconditional image synthesis by finding a better architecture through a series of ablations. For conditional image synthesis, we further improve sample quality with classifier guidance: a simple, compute-efficient method for trading off diversity for fidelity using gradients from a classifier. We achieve an FID of 2.97 on ImageNet 128×128 , 4.59 on ImageNet 256×256 , and 7.72 on ImageNet 512×512 , and we match BigGAN-deep even with as few as 25 forward passes per sample, all while maintaining better coverage of the distribution. Finally, we find that classifier guidance combines well with upsampling diffusion models, further improving FID to 3.94 on ImageNet 256×256 and 3.85 on ImageNet 512×512 . We release our code at <https://github.com/openai/guided-diffusion>.

Diffusion Models

- (+) High quality
- (+) Diversity
- (—) Slow

Diffusion Models

The generative process of a diffusion model is an **iterative denoising** process.





CS492(D): Diffusion Models and Their Applications

Minhyuk Sung, KAIST, Fall 2024

Table of contents

[Time & Location](#)

[Description](#)

[Course Staff](#)

[Prerequisites](#)

[Grading](#)

[Paper List](#)

[Useful Resources](#)

[Important Dates](#)

[Schedule](#)

[AI Coding Assistant Tool Policy](#)



a painting of a truck



Let's Consider a Collection of Real Photos.



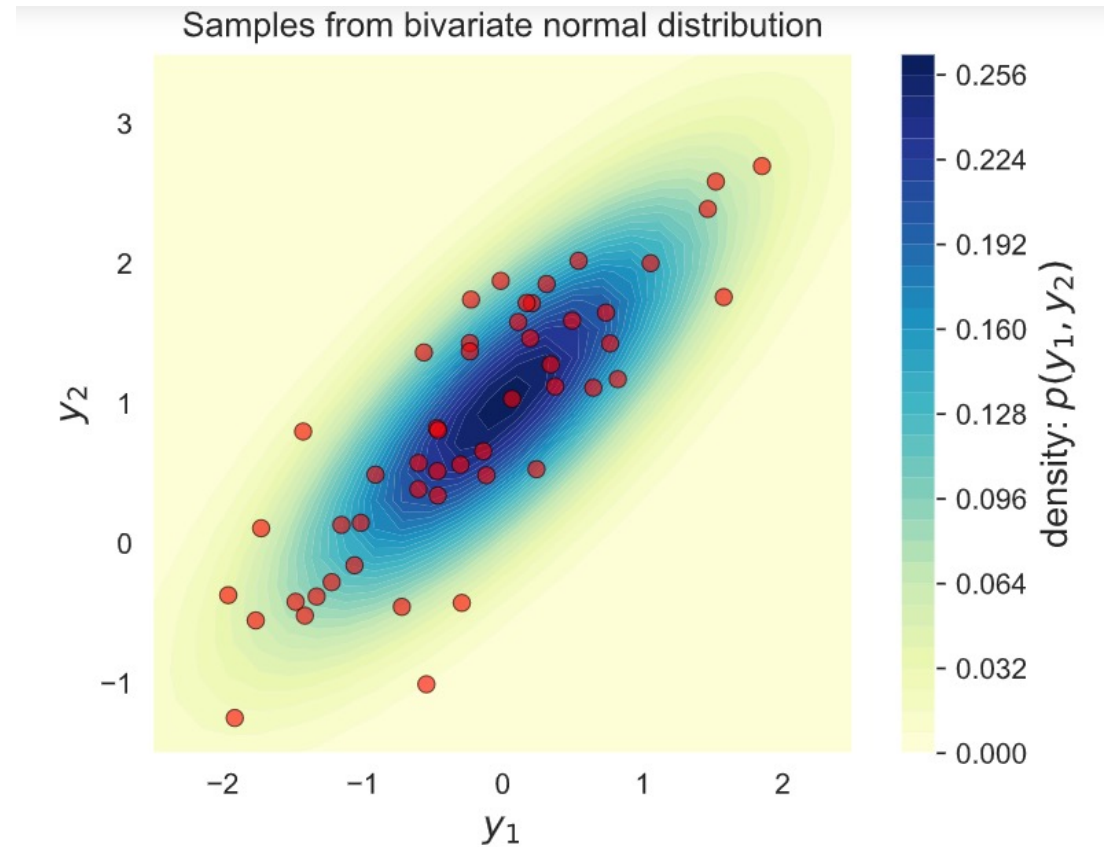
How to Generate a **New** Photo?



A Simpler Example

Let's consider a collection of
2D points.

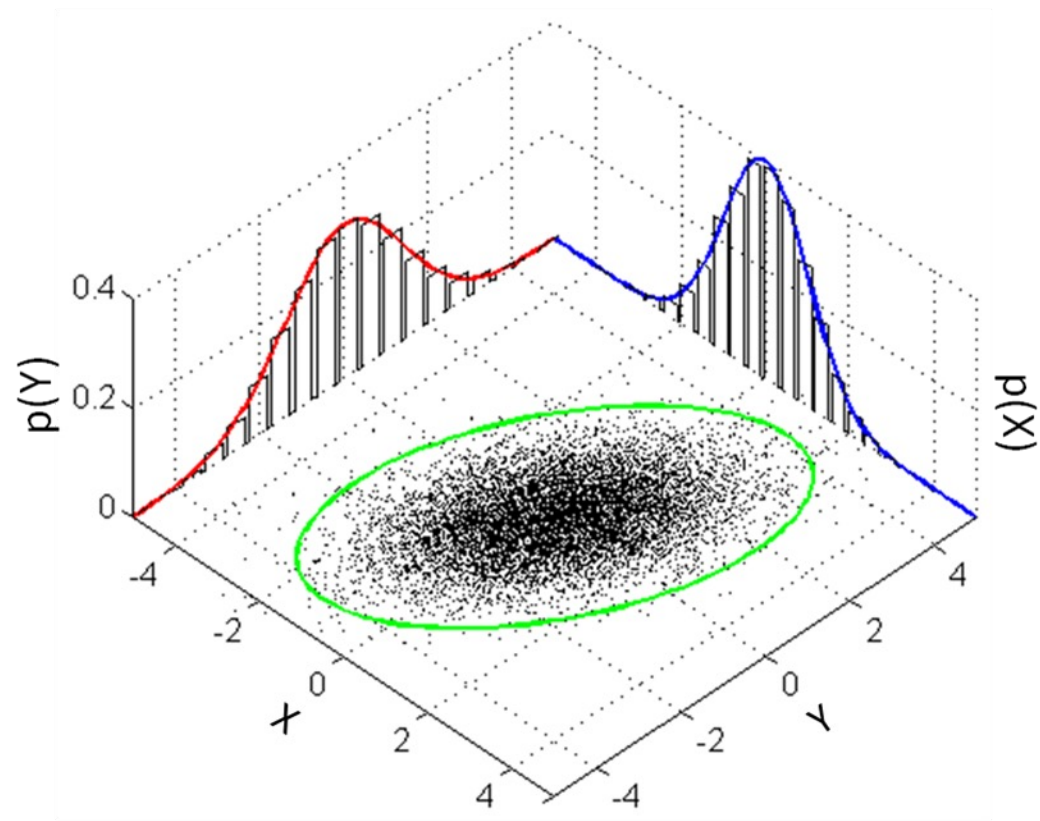
How can we sample
a **new 2D point**?



A Simpler Example

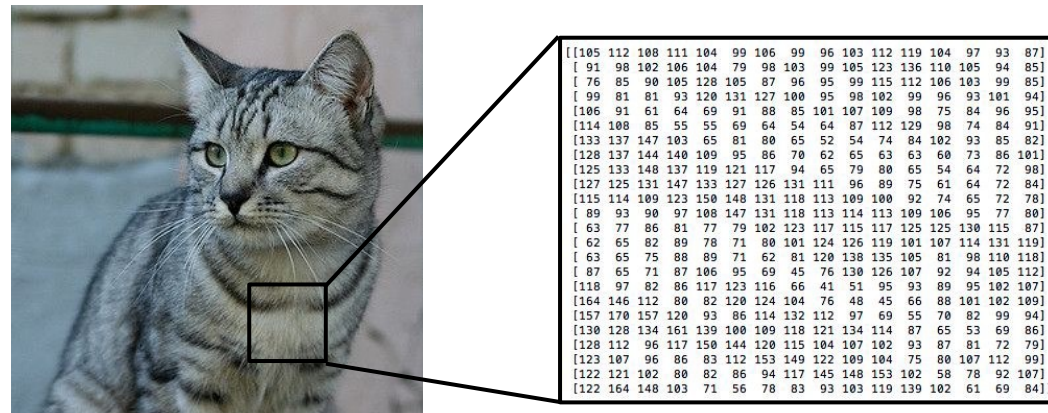
The 2D points are samples from a specific **probability distribution**.

If the distribution has a specific form (e.g., **Gaussian**), we can sample from it directly.



Statistical Perspective for Real Images

- Let's consider RGB images with a resolution of 256×256 .
- An image can be represented by a $256 \times 256 \times 3$ vector.
- This means that **an image is a point** in a $256 \times 256 \times 3$ -dimensional space.



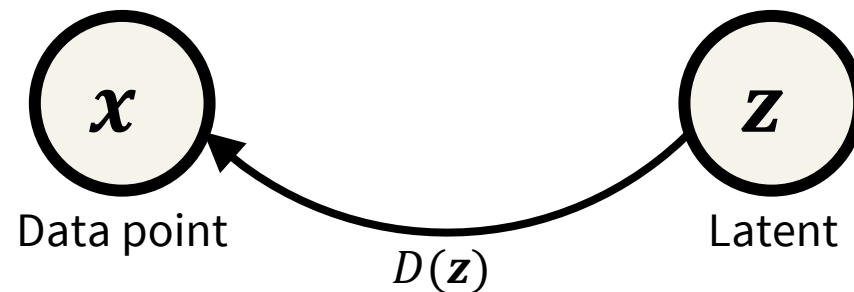
Statistical Perspective for Real Images

- Images are **sample points** in a high-dimensional space.
- Can we derive the PDF of the **data distribution** from the samples?
No, we don't know the distribution; we only have **samples**.

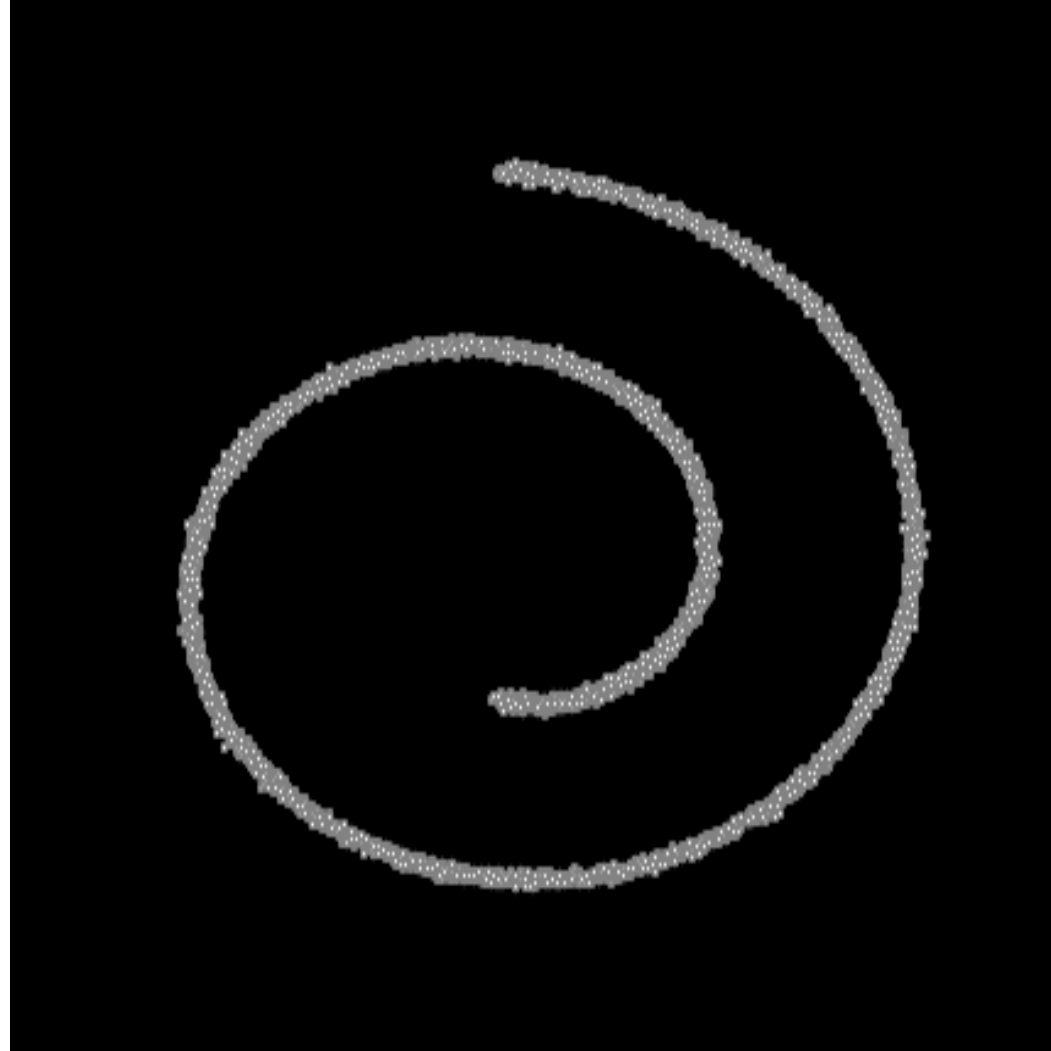


The Basic Idea

- **Map** a simple distribution which PDF is known $p(\mathbf{z})$ (e.g., a standard Gaussian distribution $\mathcal{N}(\mathbf{x}; \mathbf{0}, \mathbf{I})$) to the **data distribution** $p(\mathbf{x})$.
- Sample from $p(\mathbf{z})$ and map it to a data point.

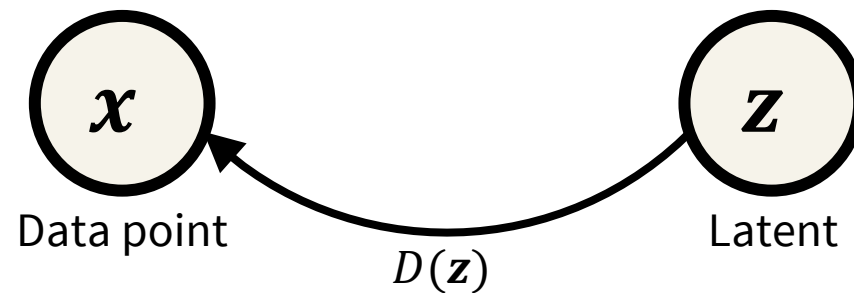


The Basic Idea



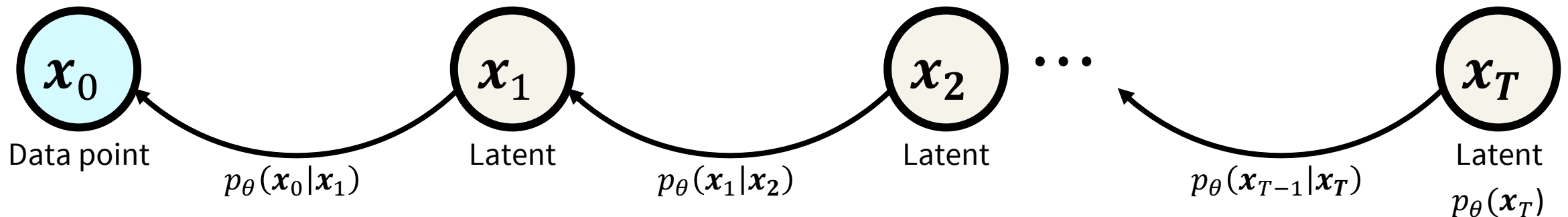
GANs / VAEs

- GANs and VAEs map a latent distribution to the data distribution **directly**.



Diffusion Models

- Diffusion models map a latent distribution to the data distribution in a **sequential, iterative manner**.



Denoising Process

- Diffusion models map a latent distribution to the data distribution in a **sequential, iterative manner**.
- This iterative process can be seen as a **denoising** process.

\mathbf{x}_t



Refinement Process

- At each step, the **expected final output** can be estimated easily.
- When viewed from this perspective, the denoising process can also be interpreted as a **refinement** process.

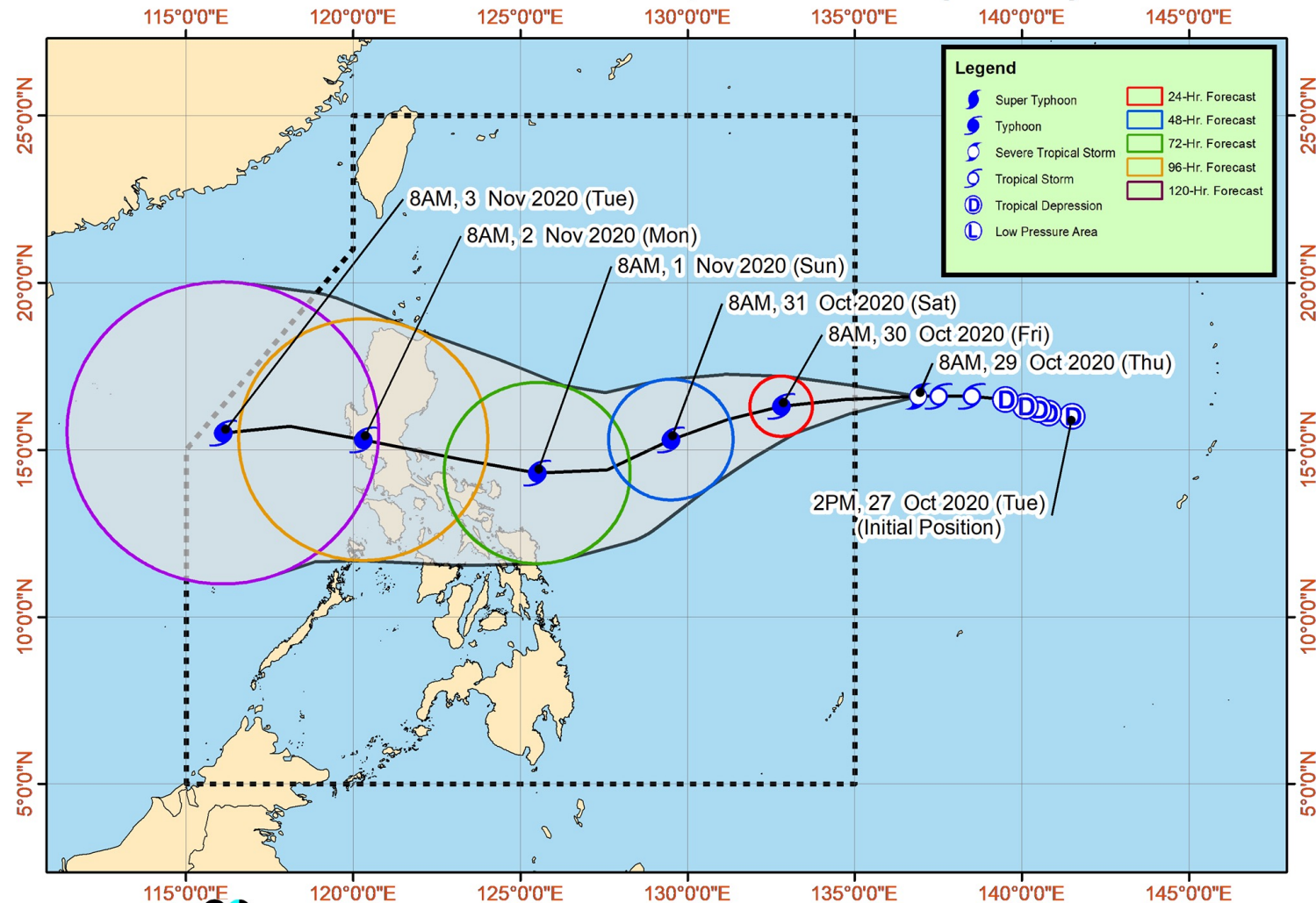


Diffusion Models



Path Prediction

Track of Severe Tropical Storm {GONI}





"Begin with the end in mind."
Stephen R. Covey

A journey without a destination is just wandering

“Begin with the end in mind”

What's Your Desired Future?

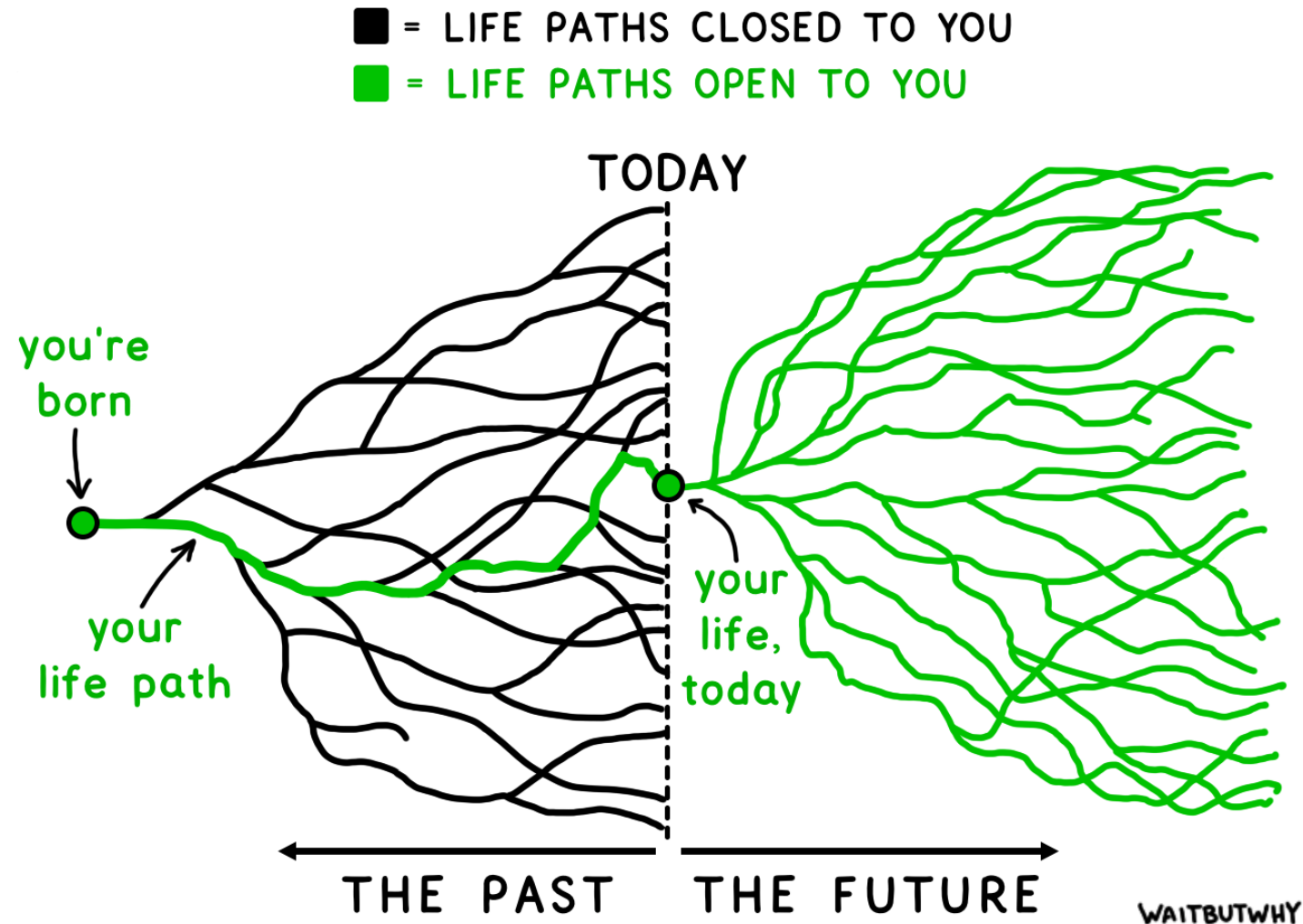
Begin With the End in Mind means to look out across the coming hours, days, weeks, months, and years with a clear vision of your desired destination—and then to proactively pursue that future.

Even if your vision for the future changes, having an end in mind for a meeting, project, or a lifetime gives you more clarity about where to start and how to adapt.

There's no better time than today to decide where you want to be tomorrow.

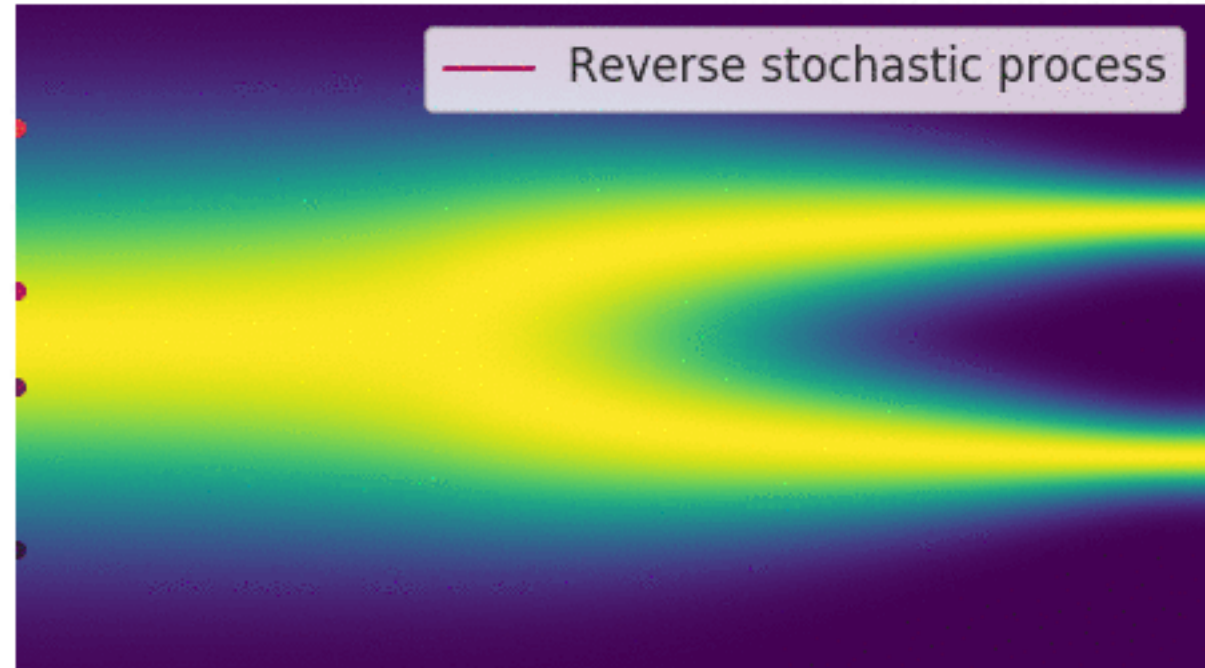


Life Paths



Stochastic Differential Equations

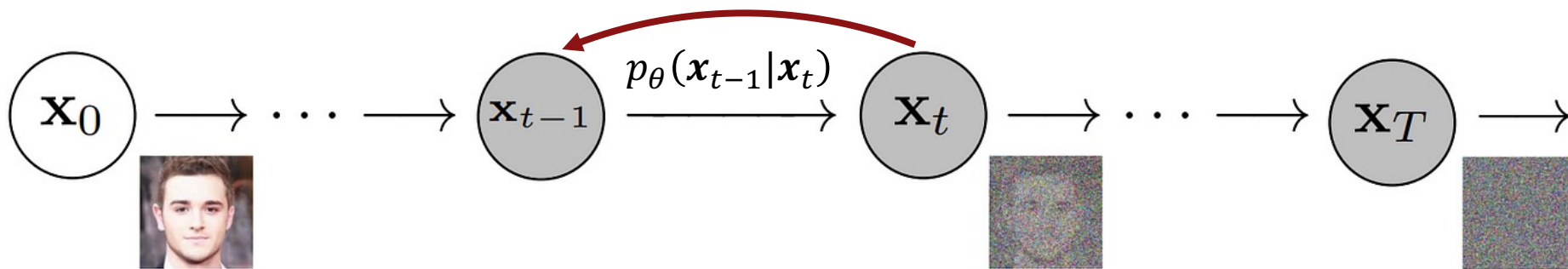
In a continuous-time domain, the mappings are formulated as **Stochastic Differential Equations** (SDEs).



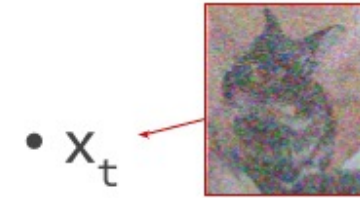
Neural Network as a Mean Predictor

At each step, a neural network predicts the **mean** of the distribution for the next denoised data point:

$$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mu_{\theta}(\mathbf{x}_t, t), \tilde{\sigma}_t^2 \mathbf{I})$$



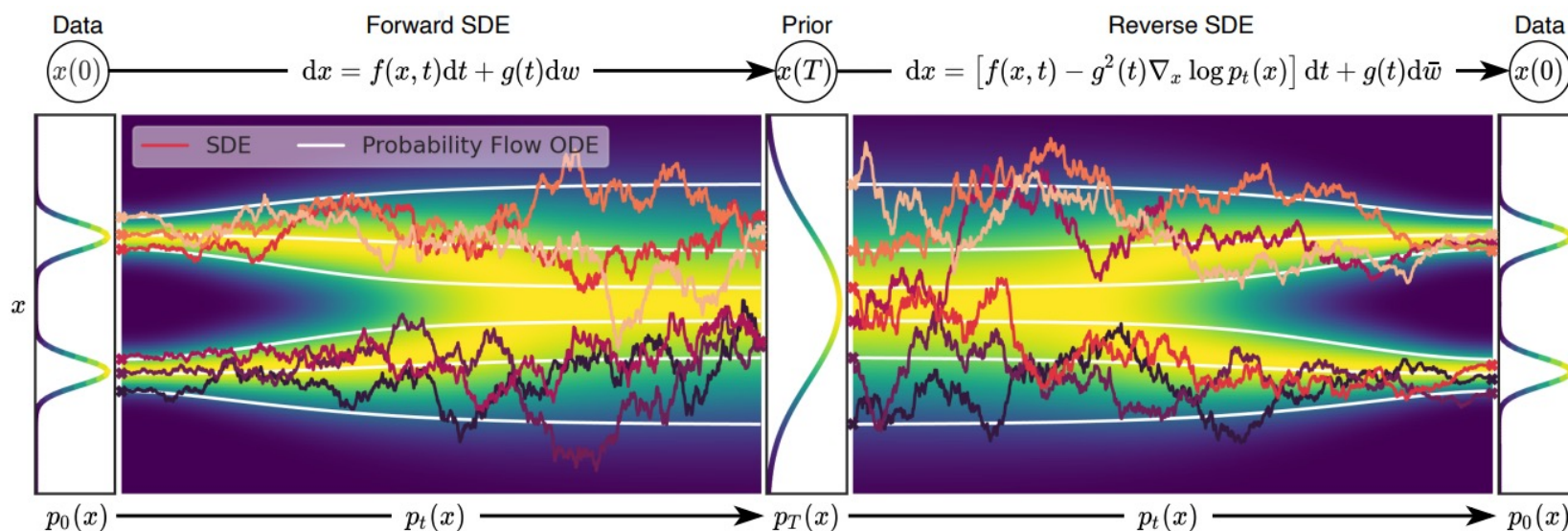
Stochastic Denoising Process



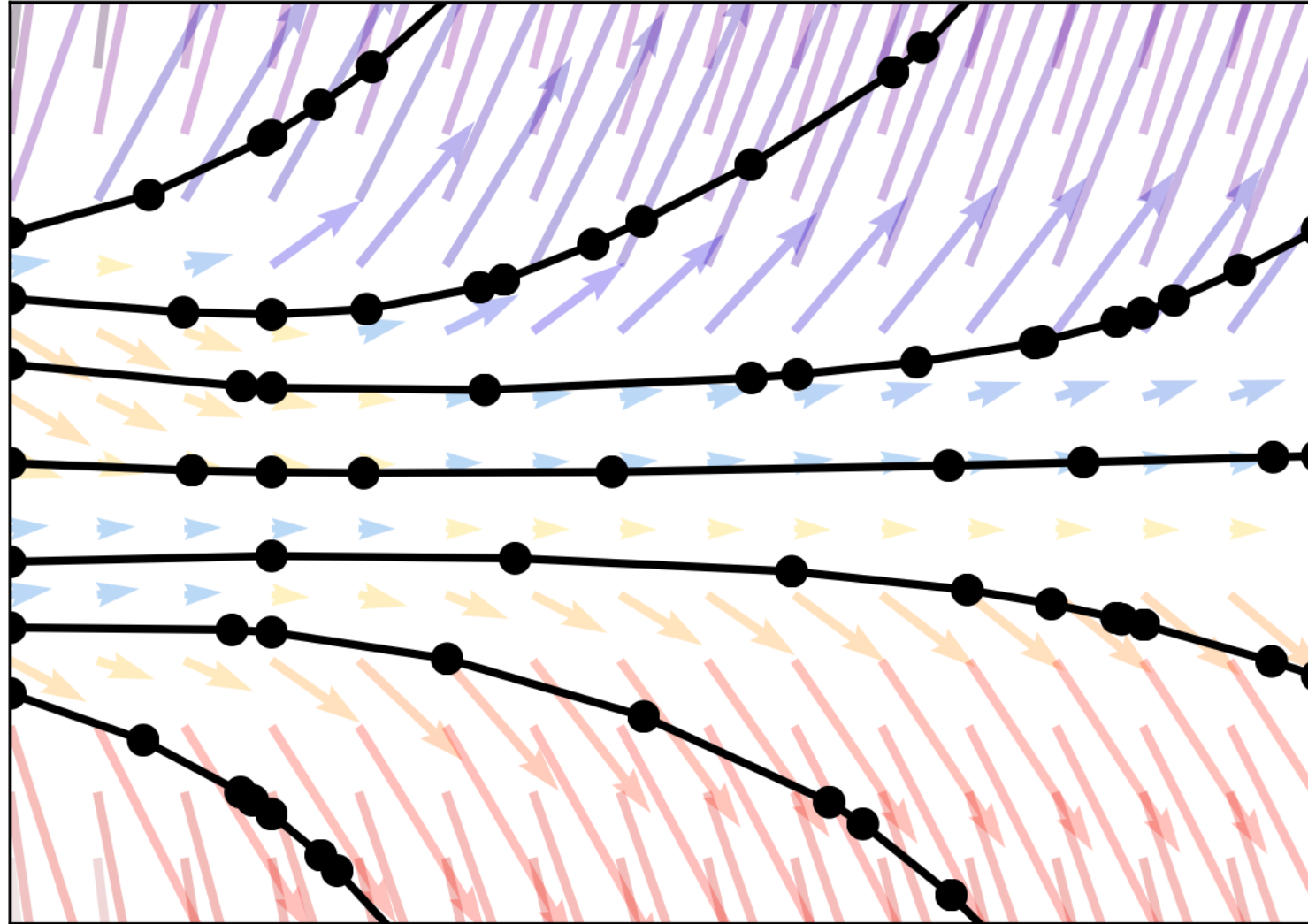
Probability Flow ODE

A SDE has a corresponding ODE (**no stochasticity**) that has same marginal distribution at every time step.

- Red curves: SDE
- White curves: ODE

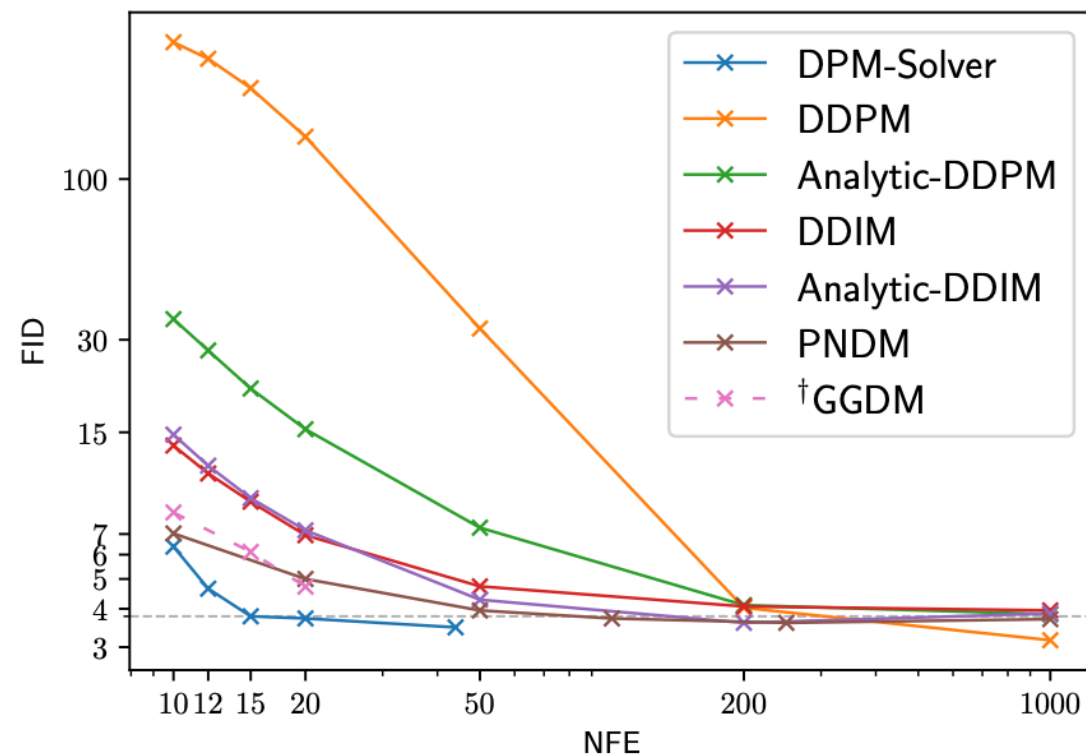


ODE Solvers



ODE Solvers

The ODE can be solved much faster using a different ODE solver at test time, without requiring any fine-tuning of the network.

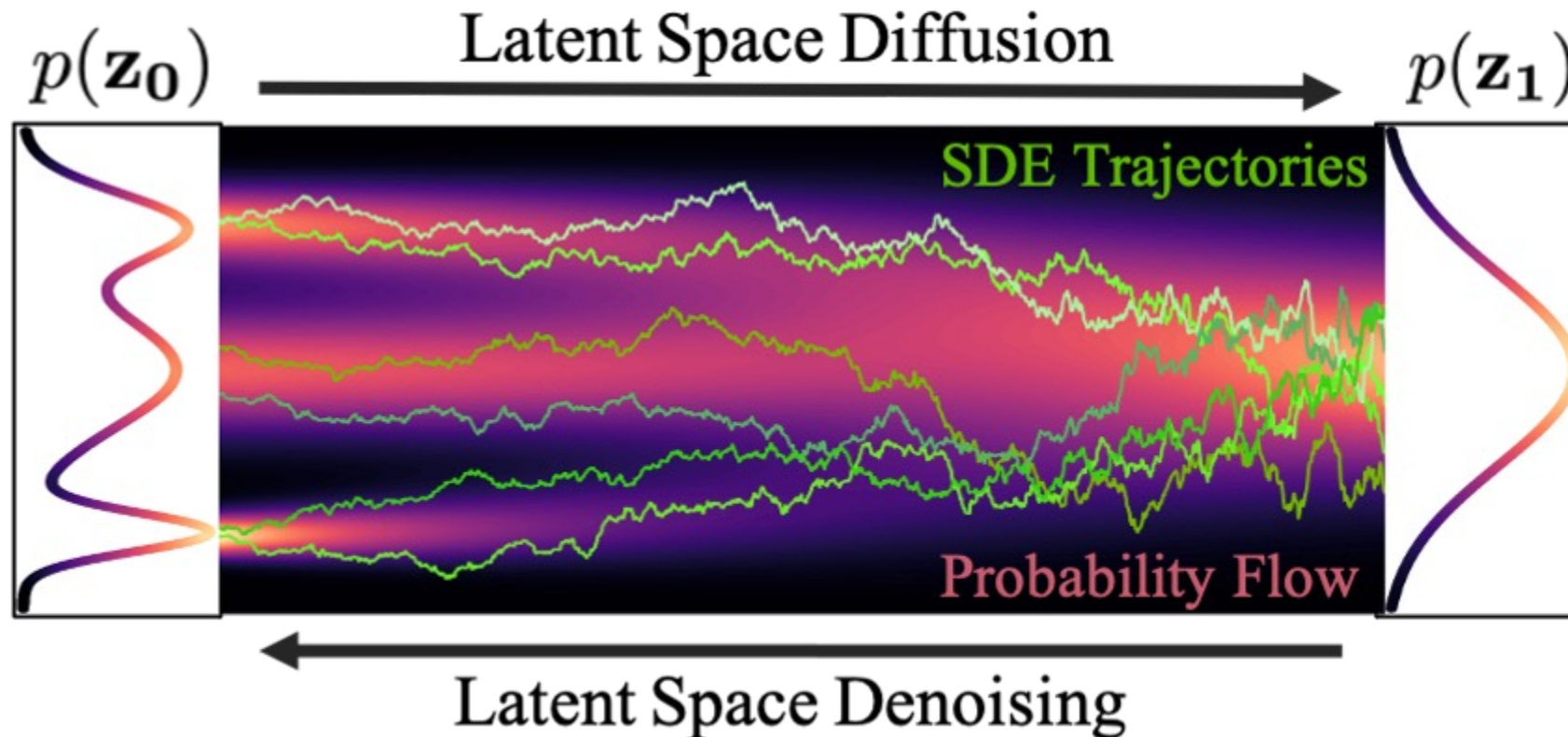


Diffusion Models

- (+) High quality
 - (+) Diversity
 - (—) Slow
-
- (+)(Relatively) easy to implement and train
 - (+) Easy to convert a conditional model
 - (+) Easy to personalize
 - (+) Easy to align to the given reward model
 - (+) Easy to distill knowledge
 - ...

Guided Generation

The denoising process can be guided in multiple different ways.



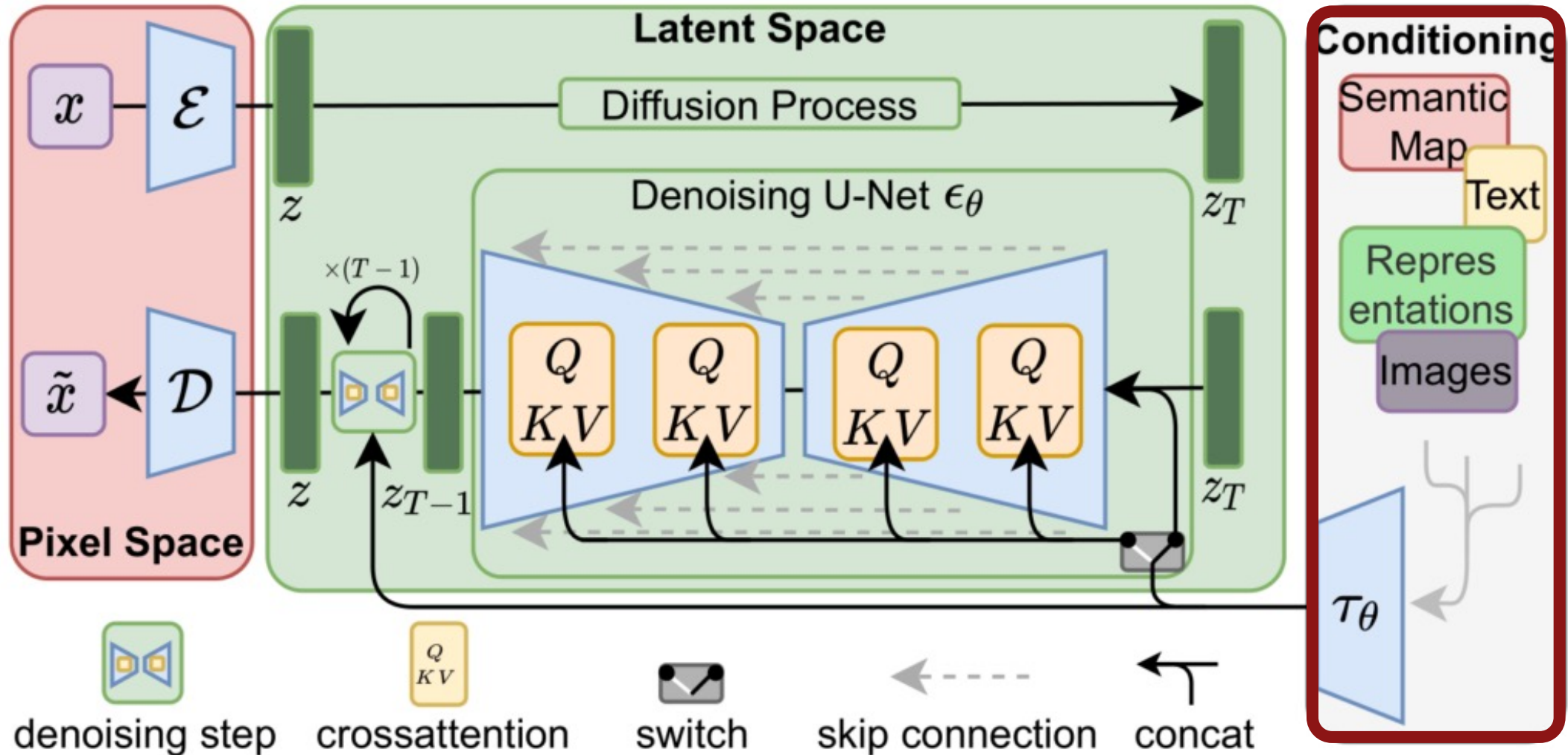
Guided Generation

Three Categories:

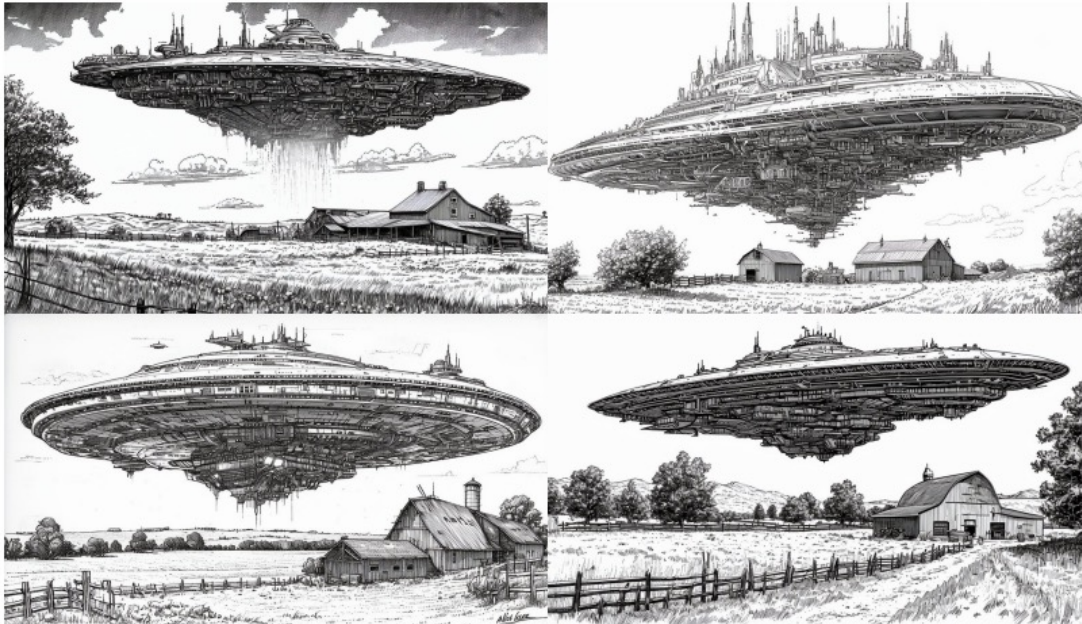
1. Manipulation in Neural Network Prediction
2. Guidance Based on the Expected Final Output
3. Particle Filtering

1. Manipulation in Neural Network Prediction

Network Architecture



Example: Text-to-Image (T2I)



detailed pen and ink drawing of a massive complex alien space ship above a farm in the middle of nowhere.



photo of a bear wearing a suit and tophat in a river in the middle of a forest holding a sign that says "I cant bear it".

Example: Text-to-Image (T2I)



tilt shift aerial photo of a cute city made of sushi on a wooden table in the evening.



dark high contrast render of a psychedelic tree of life illuminating dust in a mystical cave.

Example: Text-to-Image (T2I)



an anthropomorphic fractal person behind the counter at a fractal themed restaurant.



beautiful oil painting of a steamboat in a river in the afternoon. On the side of the river is a large brick building with a sign on top that says SD3.

Example: Text-to-Image (T2I)



an anthropomorphic pink donut with a mustache and cowboy hat standing by a log cabin in a forest with an old 1970s orange truck in the driveway

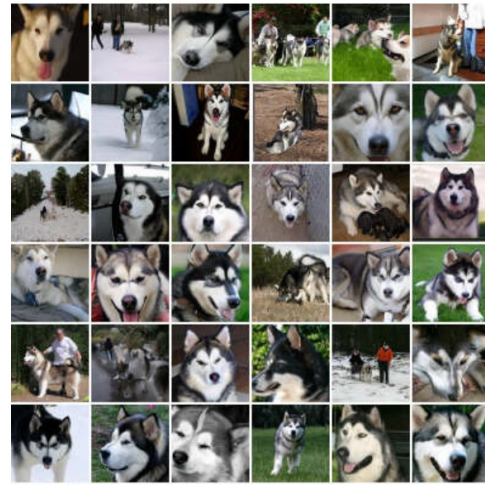


fox sitting in front of a computer in a messy room at night. On the screen is a 3d modeling program with a line render of a zebra.

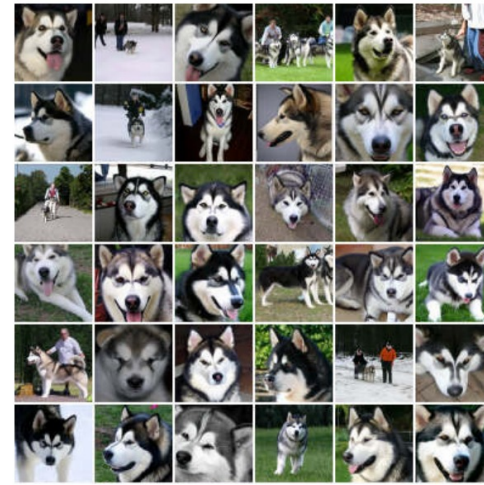
Classifier-Free Guidance (CFG)



(a) $w = 0$



(b) $w = 1$



(c) $w = 2$



(d) $w = 4$



Negative Prompt

Prompt: Portrait of zimby anton fadeev cyborg propaganda poster

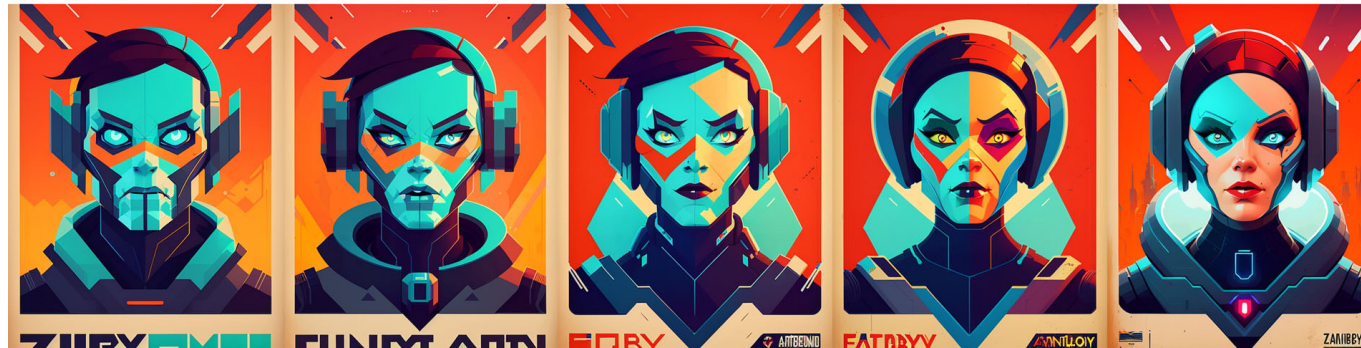
[NPW] Weight: 0.0

[NPW] Weight: 0.25

[NPW] Weight: 0.5

[NPW] Weight: 0.75

[NPW] Weight: 1.0



[-] Male

[NPW] Weight: 0.0

[NPW] Weight: 0.25

[NPW] Weight: 0.5

[NPW] Weight: 0.75

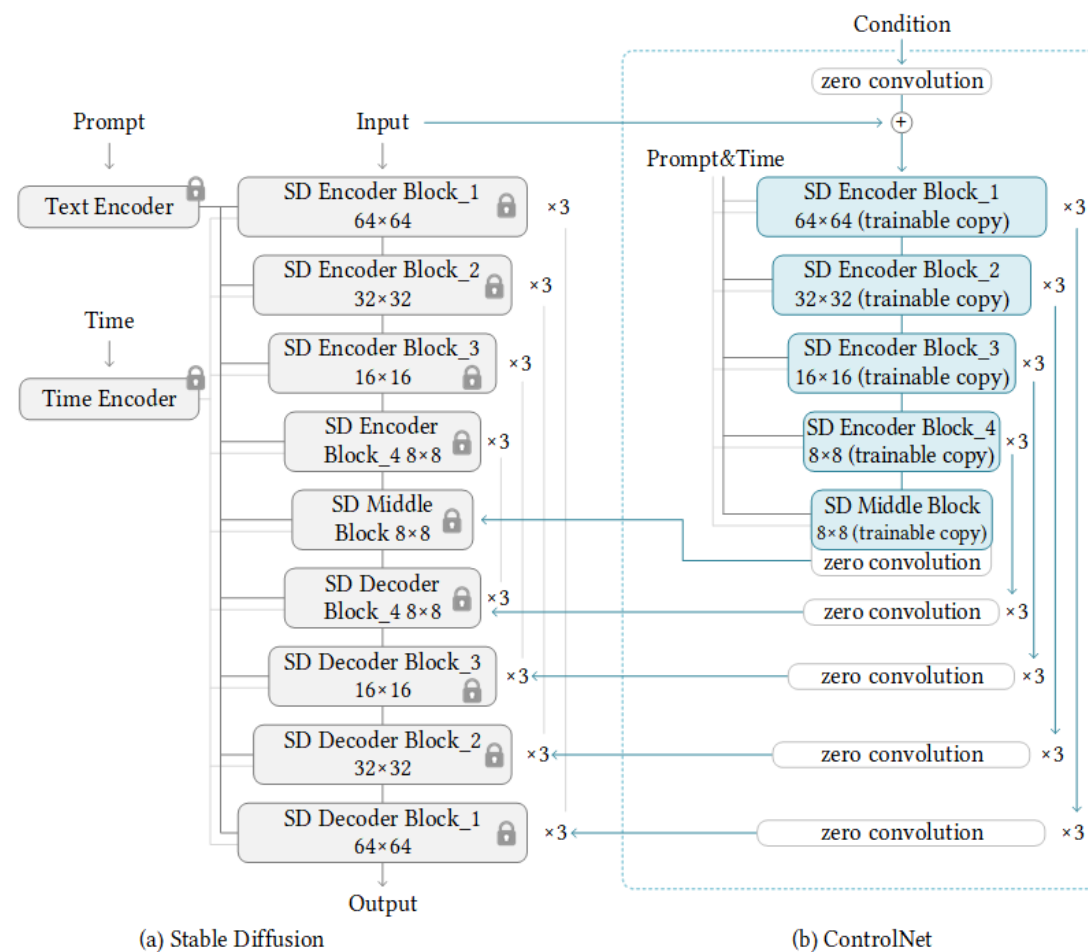
[NPW] Weight: 1.0



[+] Female

ControlNet

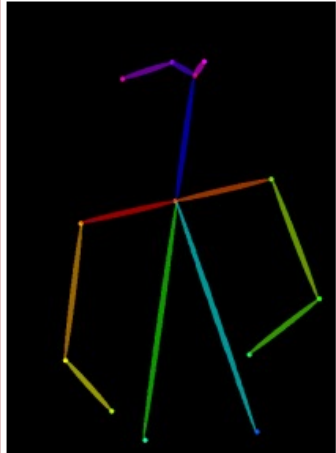
Can we **convert** a pretrained unconditional image diffusion model into an **image-conditioned** generative model using a relatively much **smaller set of input-output pairs** ($\sim 50k$)?



ControlNet



Input Canny edge



Input human pose



Default



“masterpiece of fairy tale, giant deer, golden antlers”



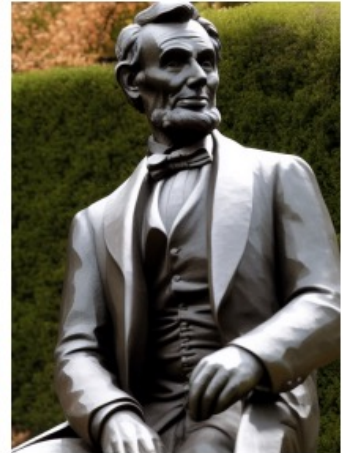
“..., quaint city Galic”



Default

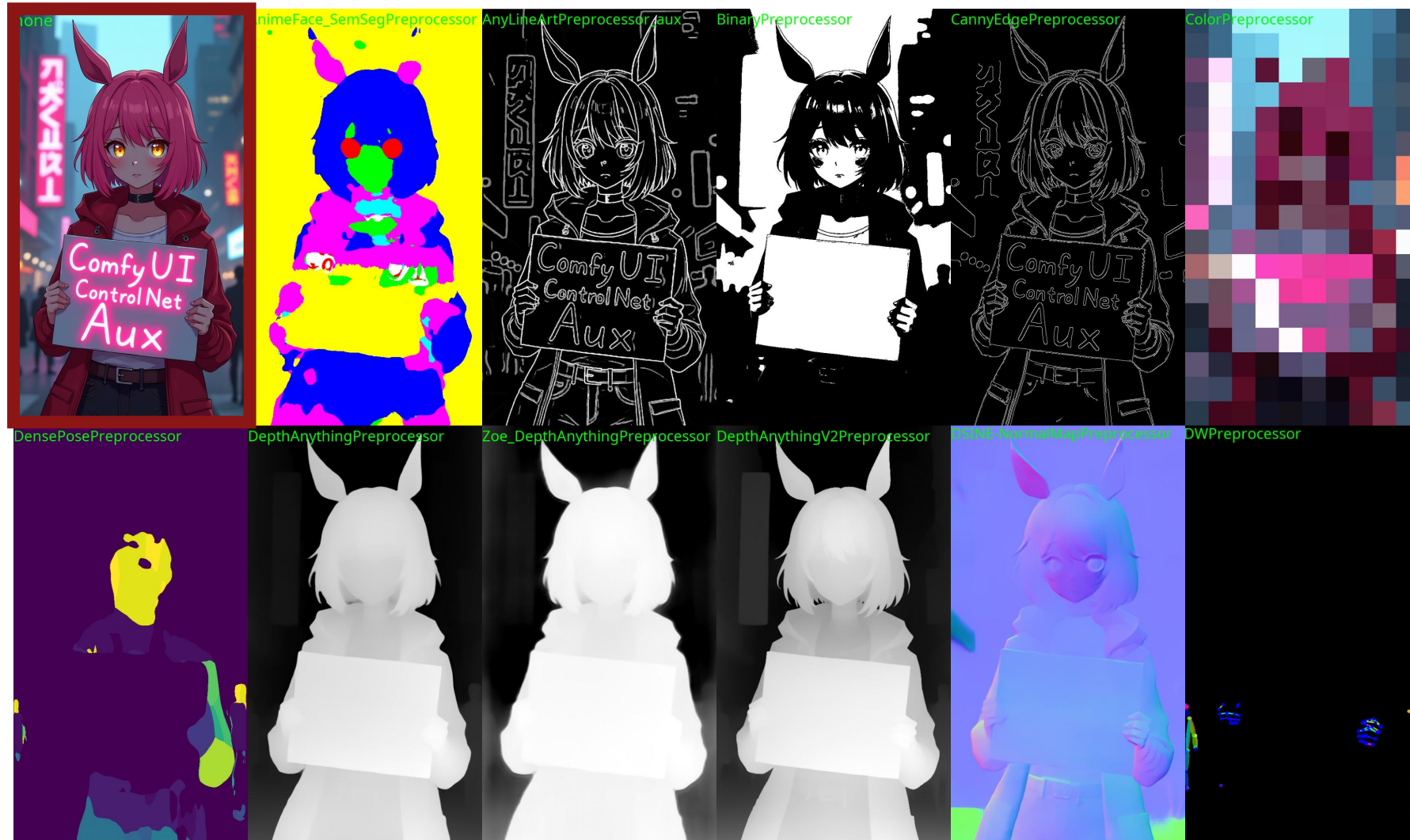


“chef in kitchen”

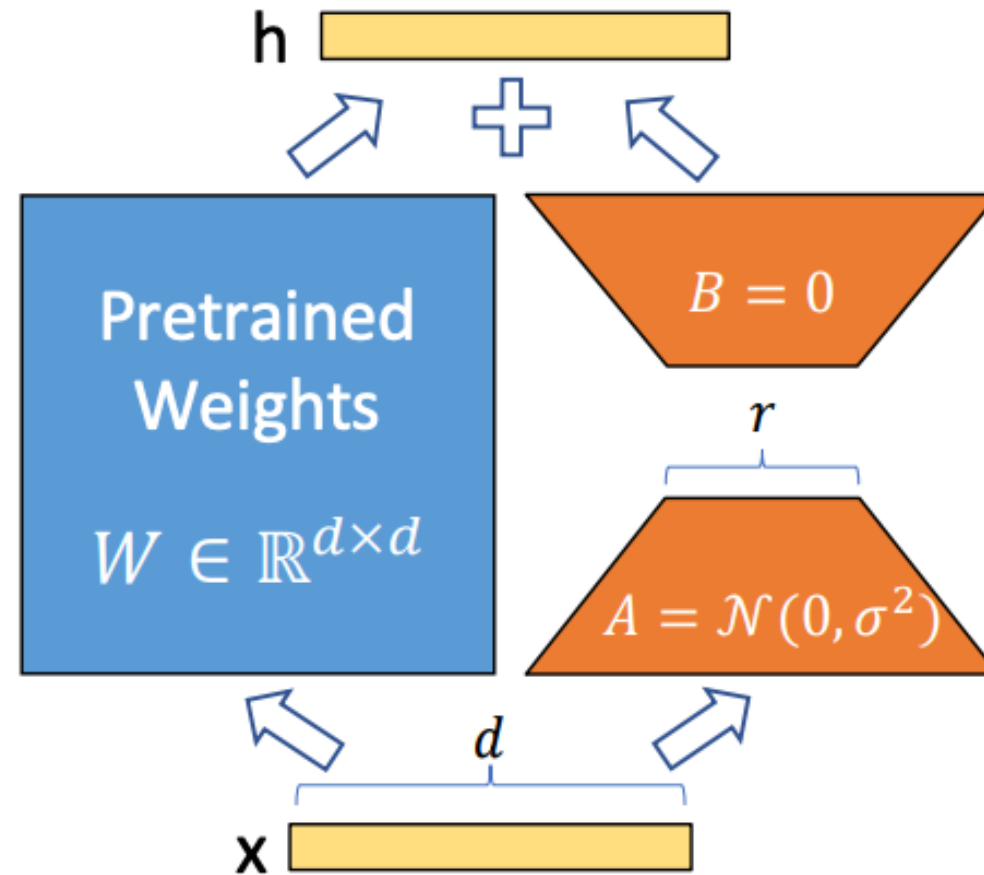


“Lincoln statue”

ControlNet



Low-Rank Adaptation (LoRA)



LoRA-Based Personalization



LoRA

LoRA-Based Personalization

Input images



A [V] backpack in the Grand Canyon



A [V] backpack with the night sky



A [V] backpack in the city of Versailles



A wet [V] backpack in water



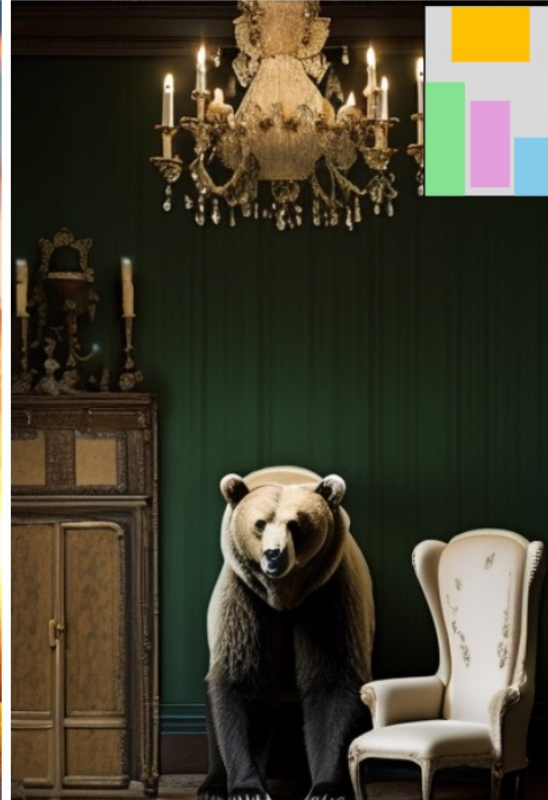
A [V] backpack in Boston

DreamBooth

2D Spatial Grounding



"A car and a castle in a beautiful landscape and a balloon in the sunset sky."



"A bear is in an antique living room with a chair and a chiffoier, while a chandelier hanging from the ceiling."

"A balloon and fireworks are in the night sky, and below is a beautiful forest with a car and a chair and a waterfall."



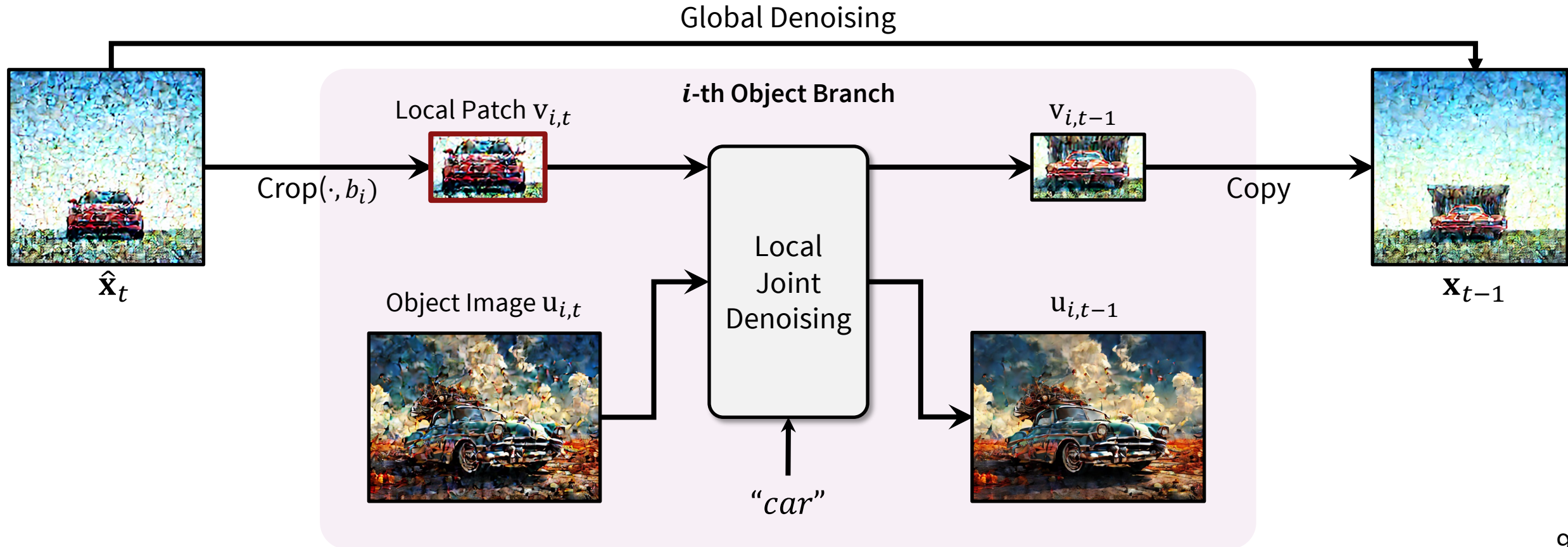
"A banana and an apple are beneath a book and a flower is lying on the book in a room"



"Aurora lights up the sky and a horse and a house are on the grassy meadow with a mountain in the background."

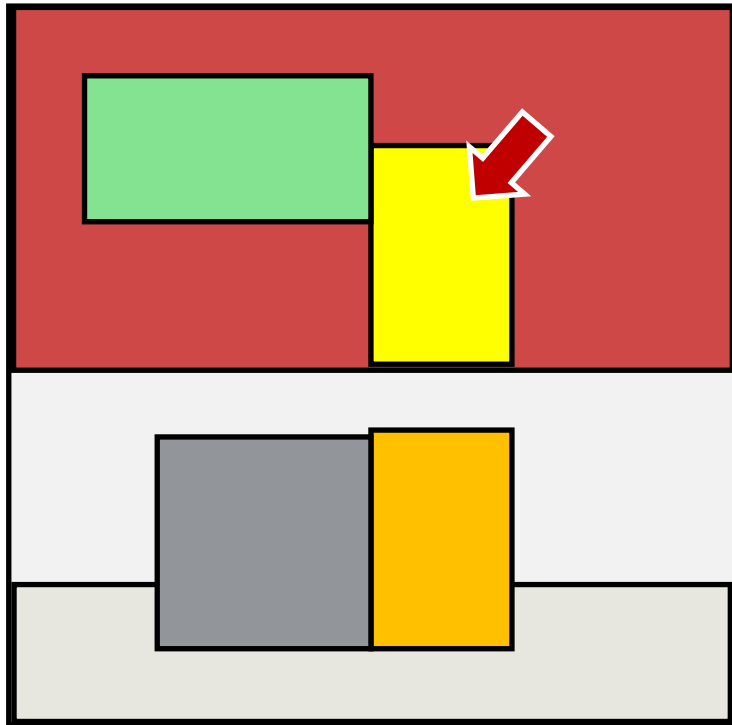
Joint Diffusion with DiT

Generate an object image for each bounding box and jointly denoise it with the corresponding region of the main image.

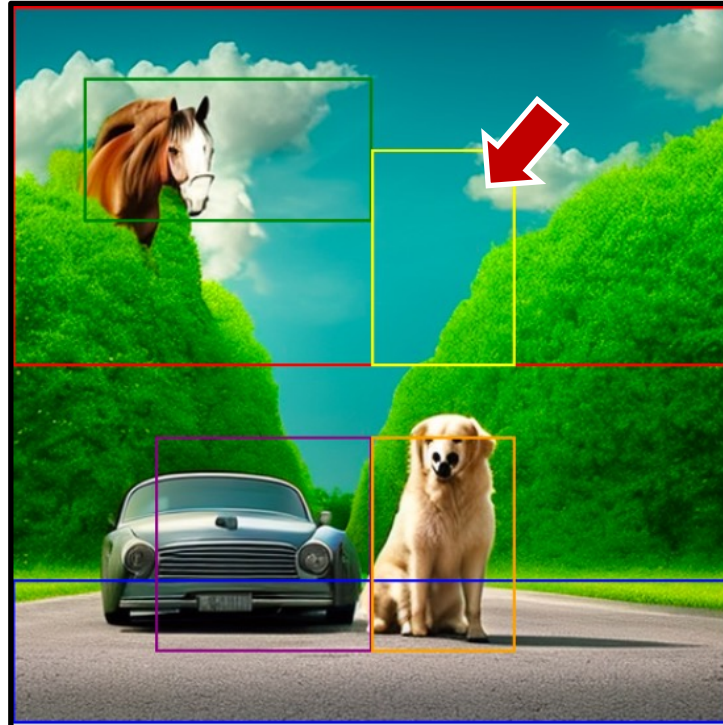


GrounDiT Results

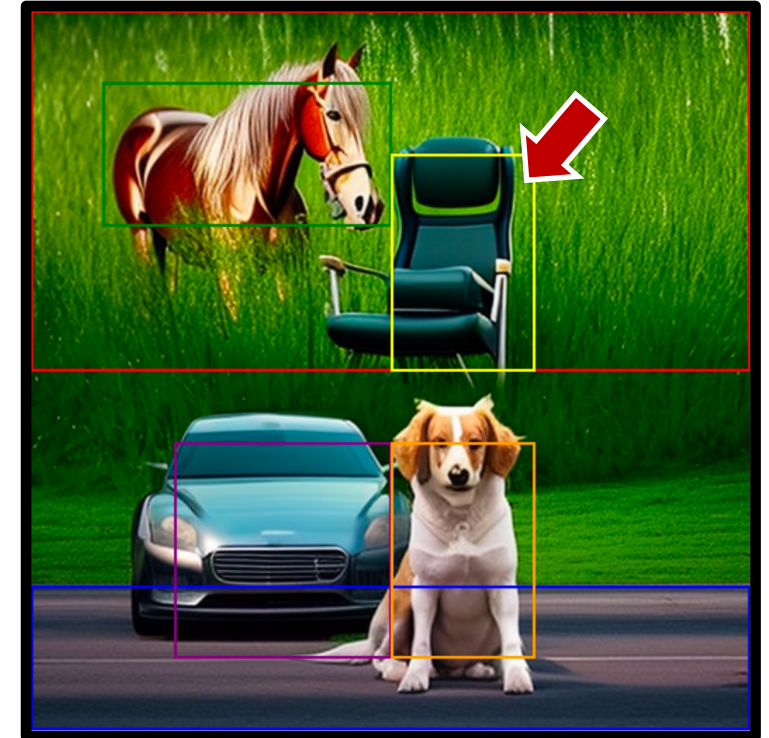
Layout



R&B [Xiao et al., ICLR 2024]



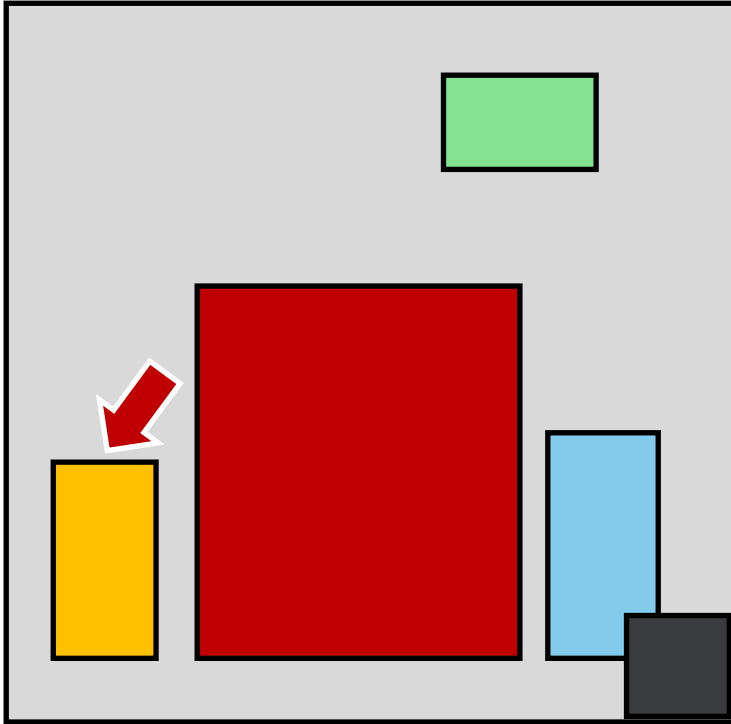
Ours



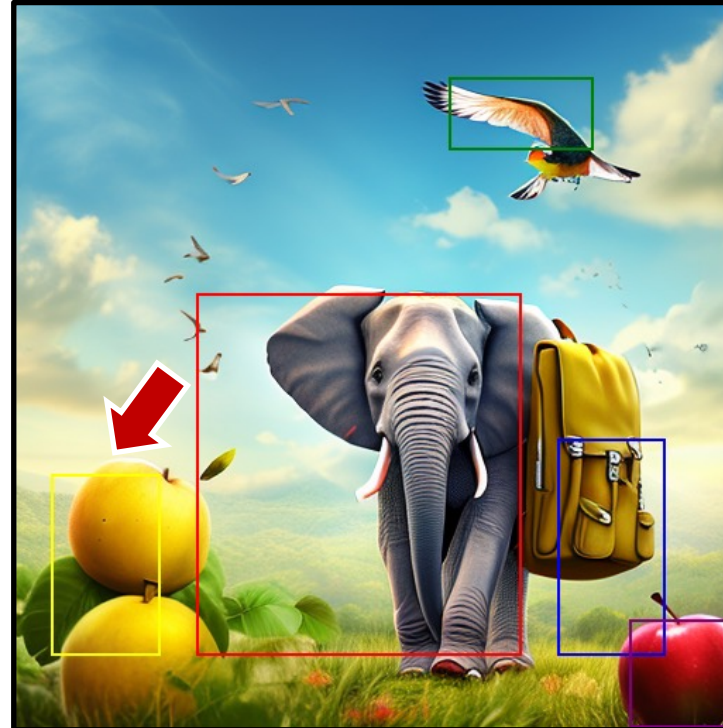
*“A car and a dog on the road while horse and
a chair is on the grass.”*

GrounDiT Results

Layout



R&B [Xiao et al., ICLR 2024]



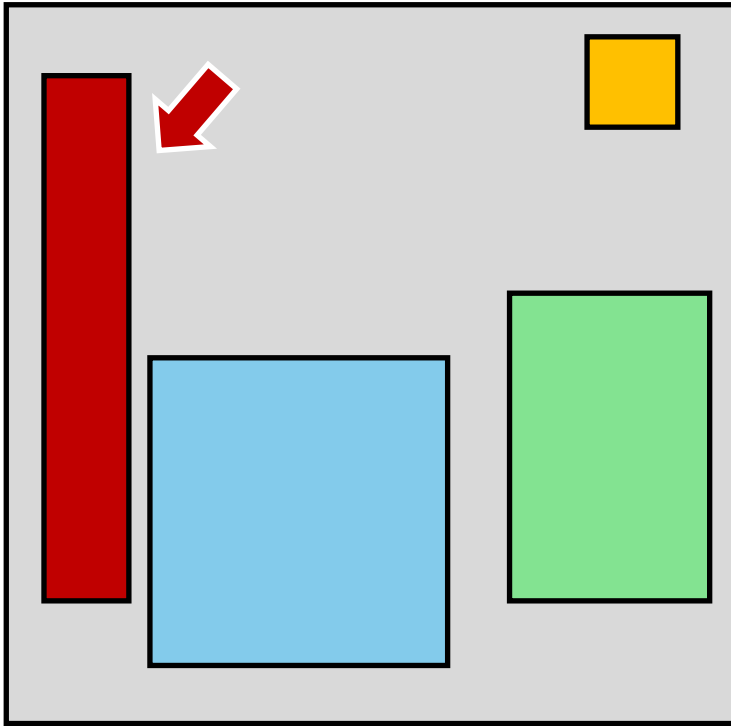
Ours



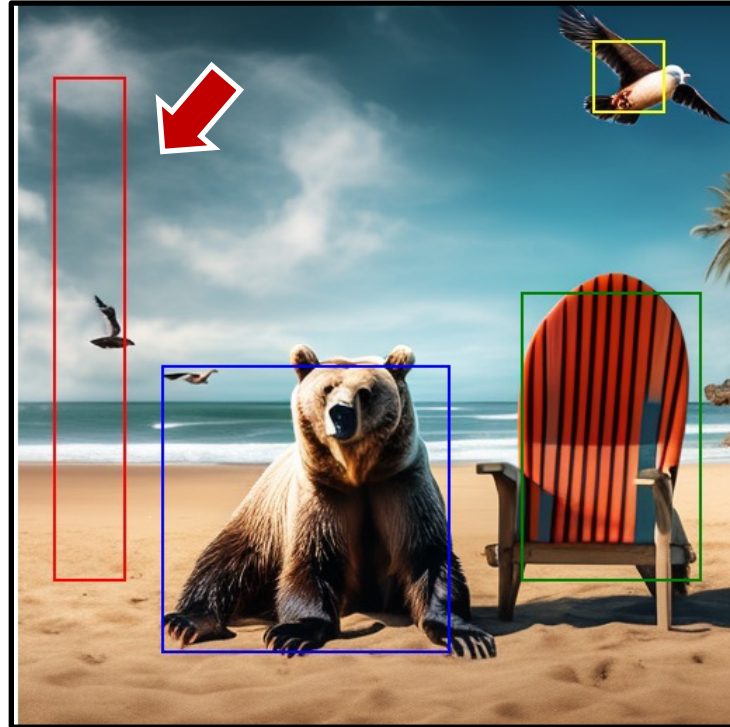
“A banana and an apple and an elephant and a backpack in the meadow with bird flying in the sky.”

GrounDiT Results

Layout



R&B [Xiao et al., ICLR 2024]



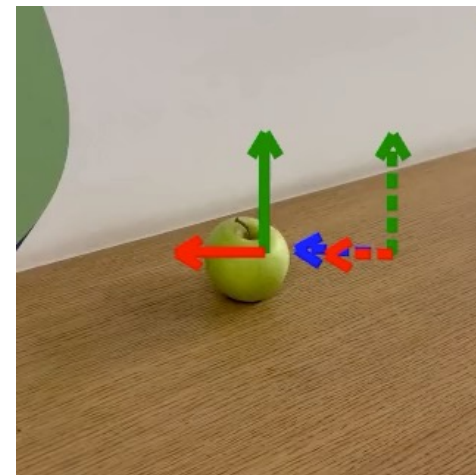
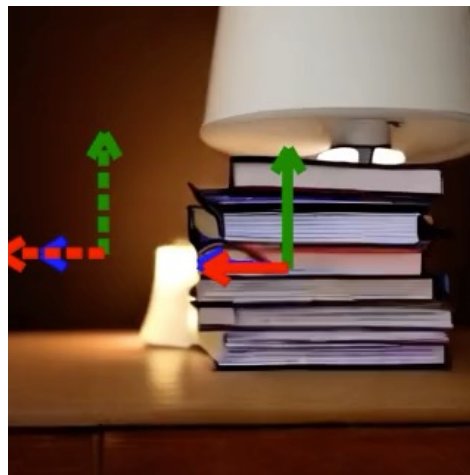
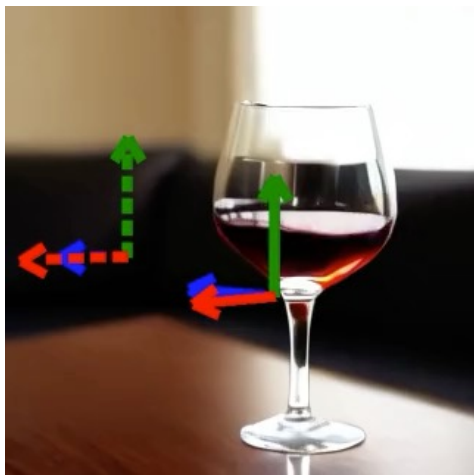
Ours



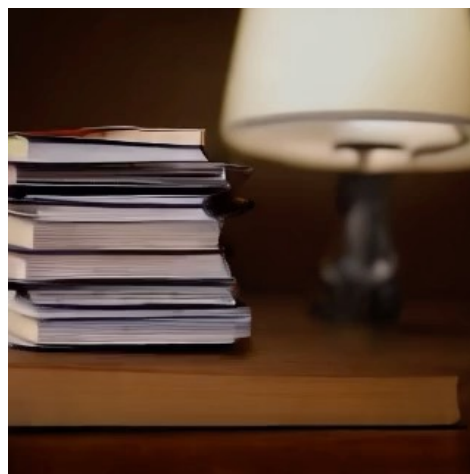
*“A bear sitting between a surfboard and a chair
with a bird flying in the sky.”*

Compositional Video Editing

Input

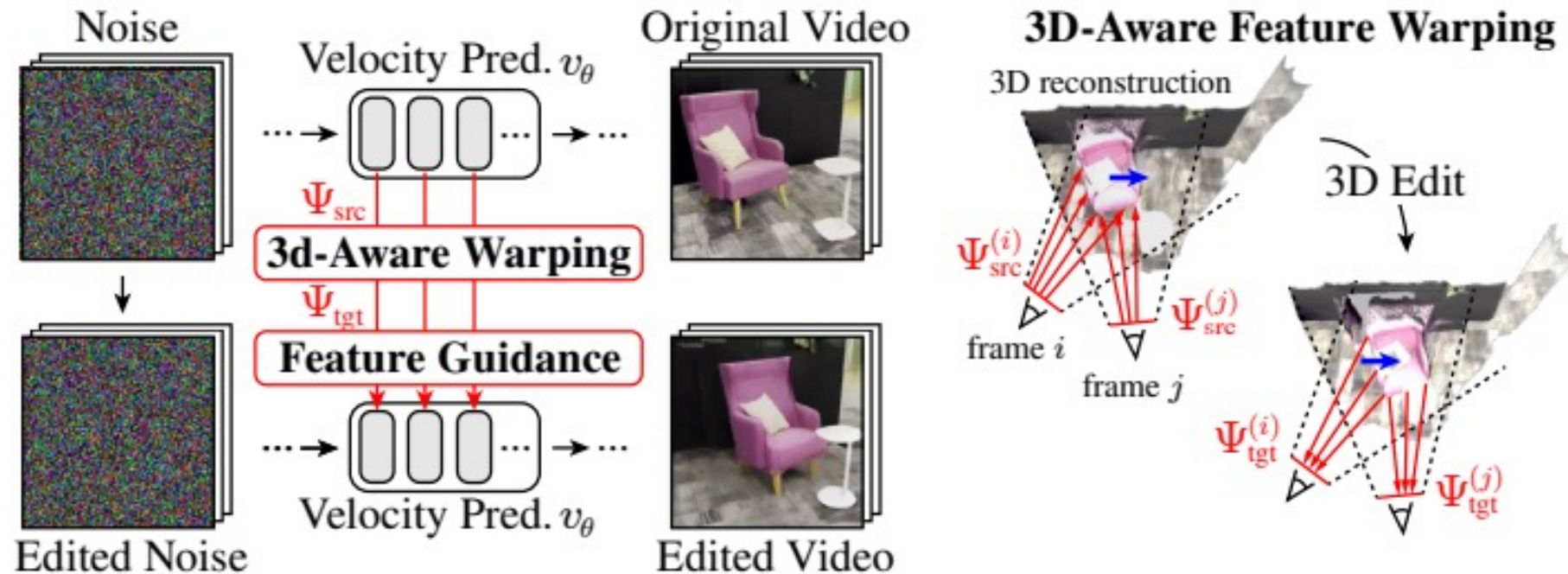


Edited Video



VideoHandles

Key Idea: Warp the **attention maps** of the score prediction network based on **3D reconstruction** and the given **3D transformation**.



VideoHandles Results

Input

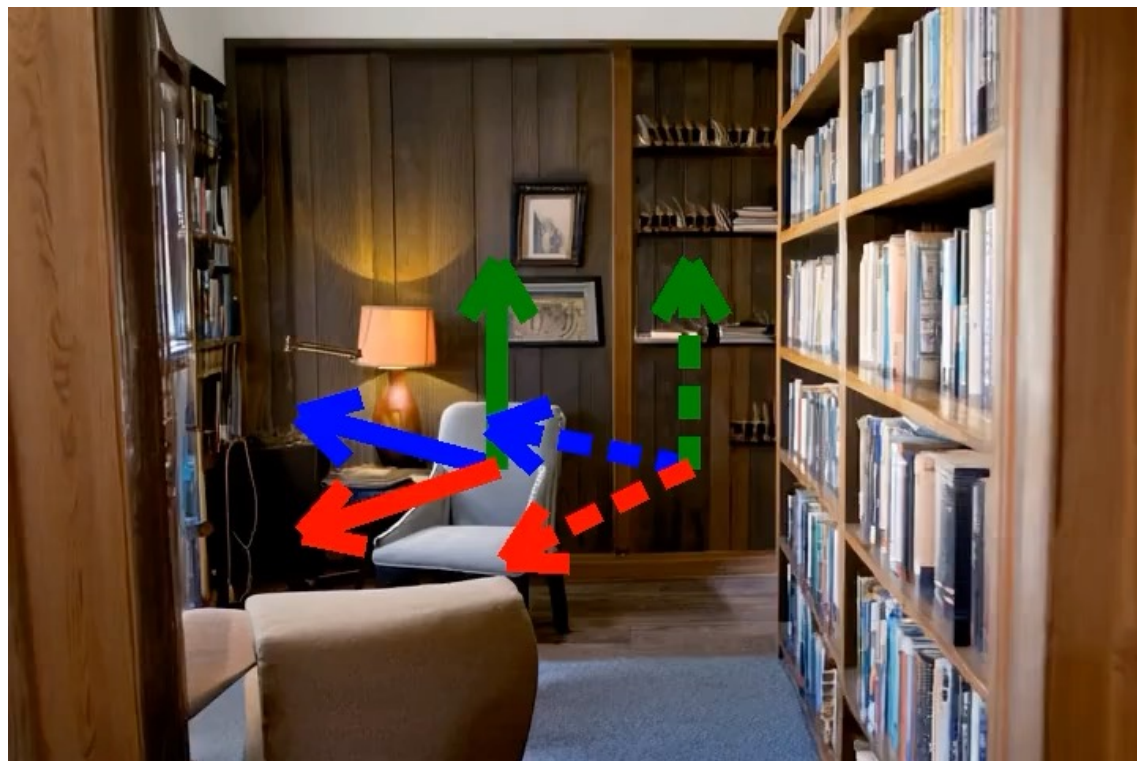


Output



VideoHandles Results

Input



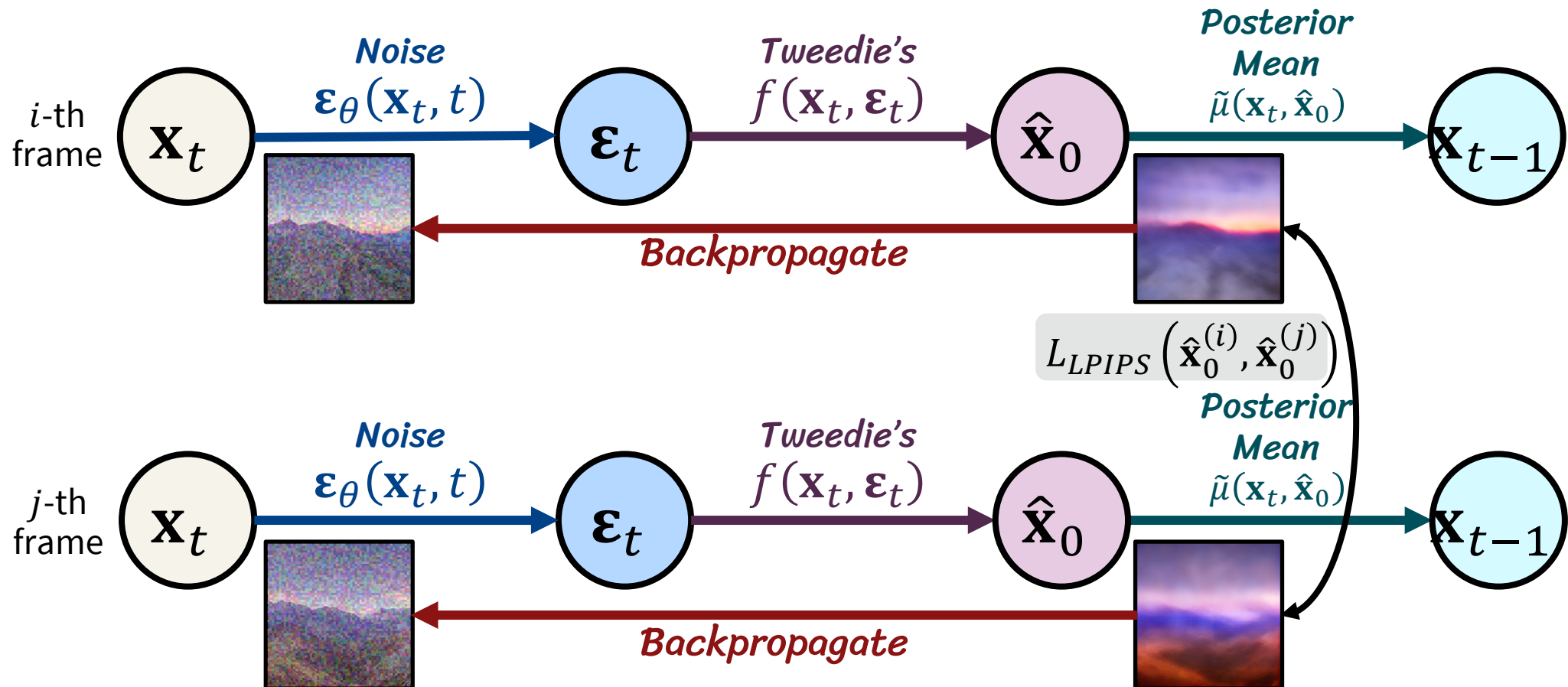
Output



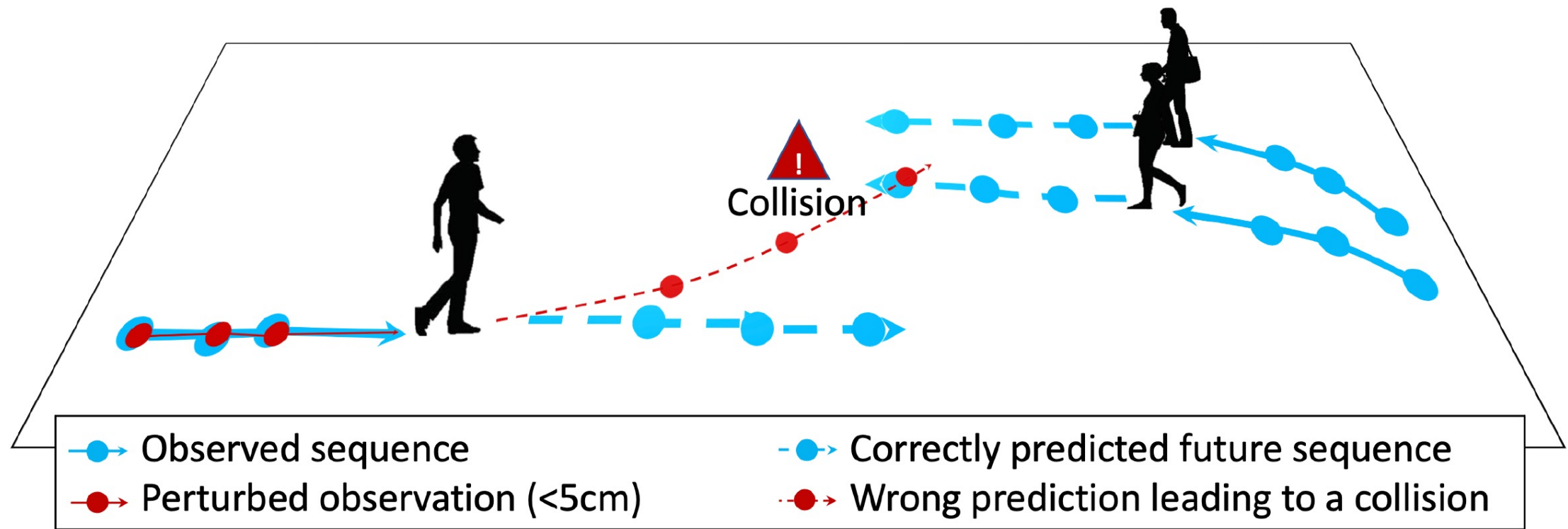
2. Guidance Based on the Expected Final Output

Guidance via Lookahead

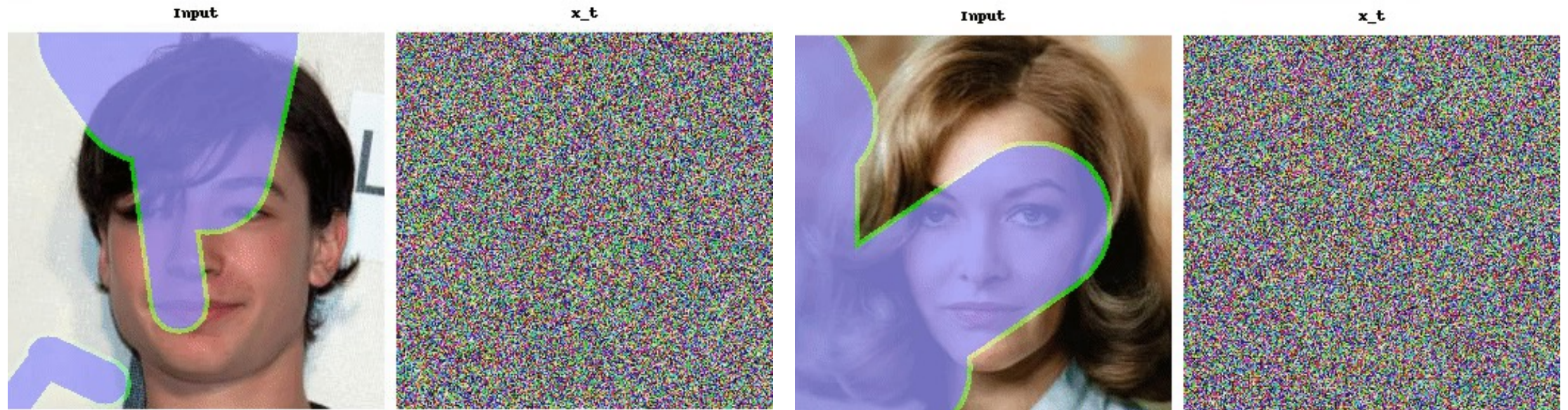
Perform backpropagation from the expected final output to the noisy data point.



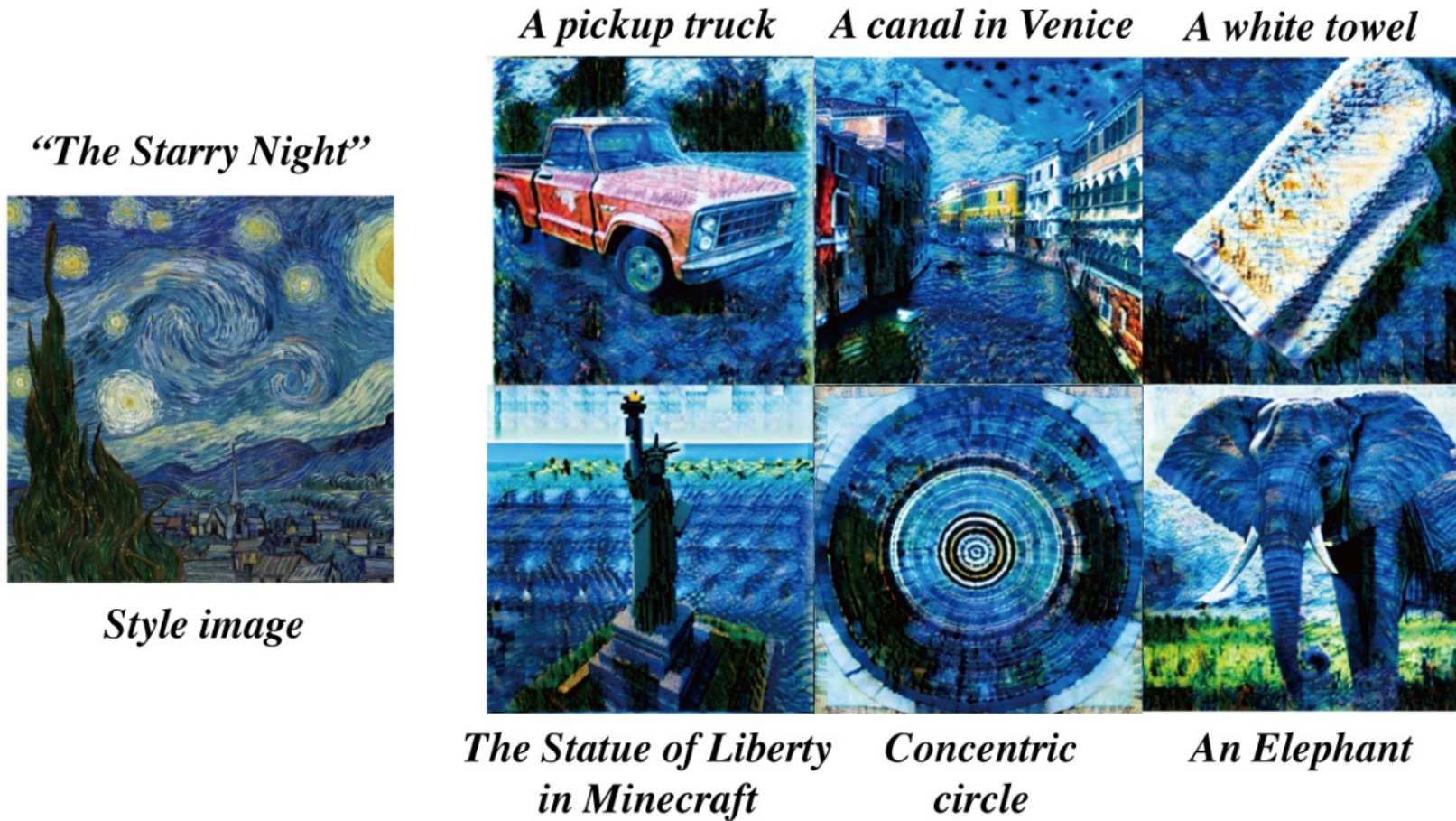
Path Correction via Lookahead



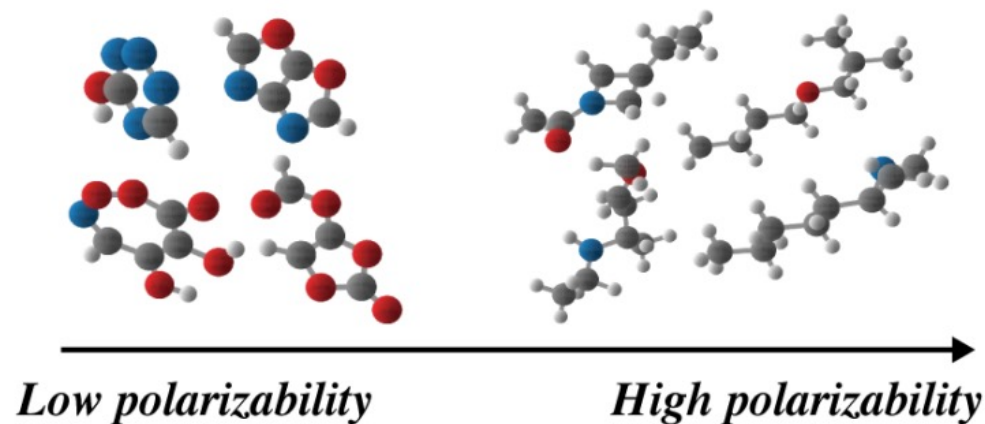
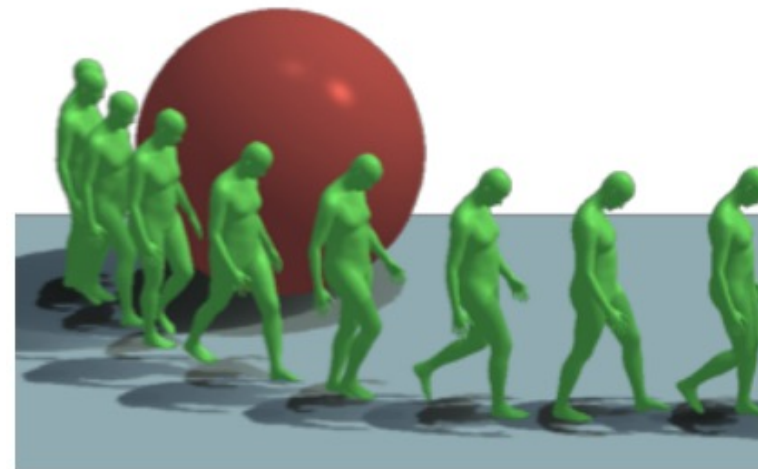
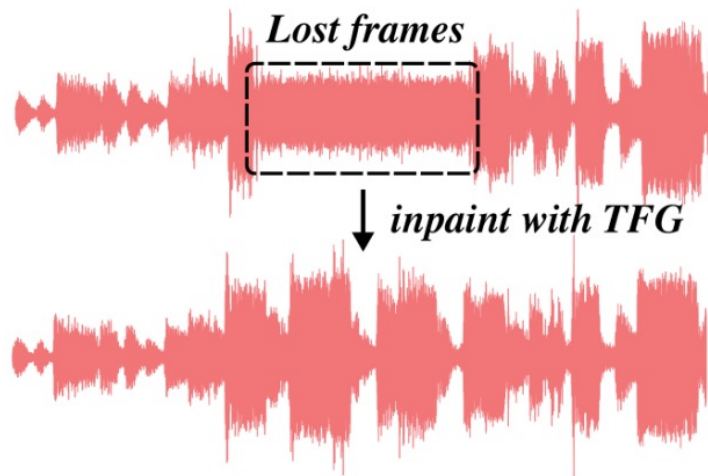
Application: Inpainting



Application: Style Transfer

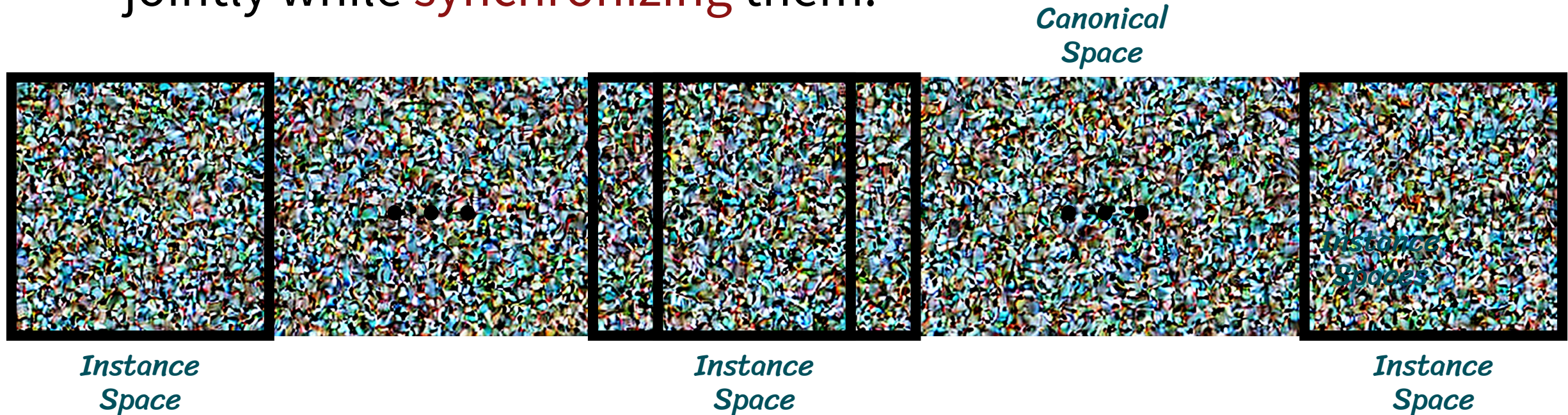


Application in Audio/Motion/Molecule



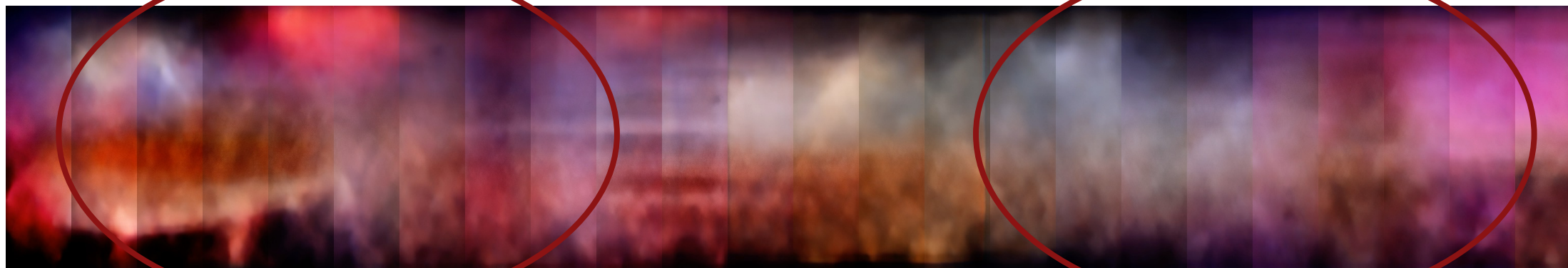
Diffusion Synchronization

- Given a **large image space**, define multiple **instance spaces** with overlaps.
- Perform the generative processes for all instance frames jointly while **synchronizing** them.

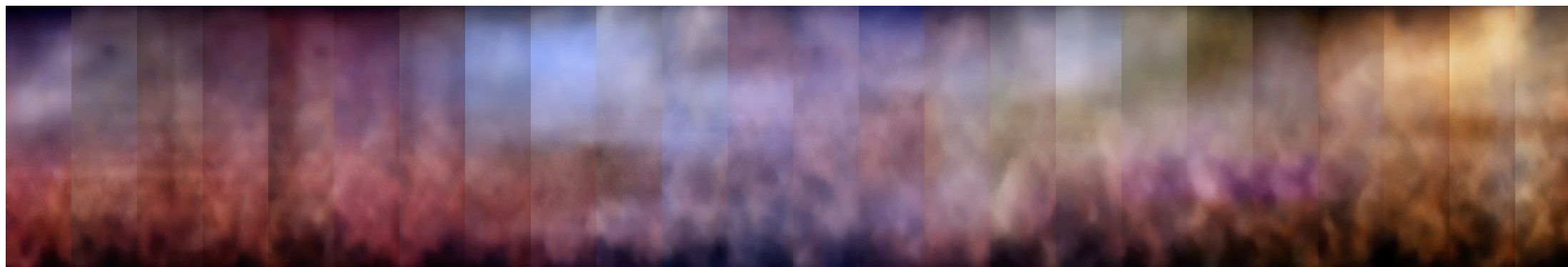


Application: Wide Image Generation

MultiDiffusion (Bar-Tal et al.)



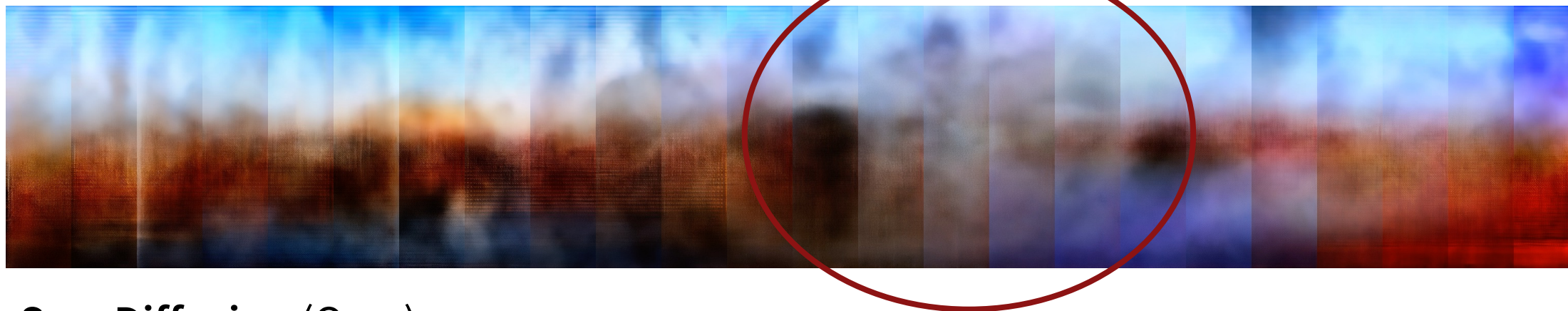
SyncDiffusion (Ours)



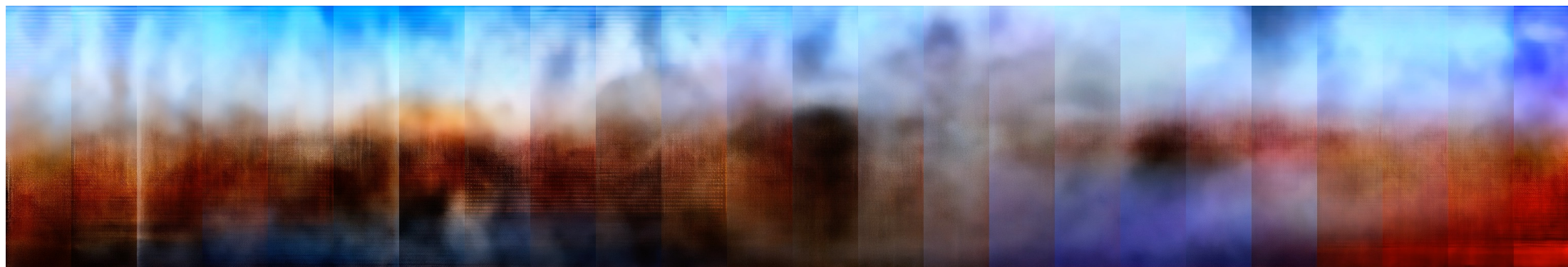
“A photo of a rock concert”

Application: Wide Image Generation

MultiDiffusion (Bar-Tal et al.)



SyncDiffusion (Ours)



“Skyline of New York City”

Application: Wide Image Generation

MultiDiffusion (Bar-Tal et al.)



SyncDiffusion (Ours)



“Silhouette wallpaper of a dreamy scene with shooting stars”

More results are available on <https://syncdiffusion.github.io/>.

Application: Wide Image Generation

MultiDiffusion (Bar-Tal et al.)



SyncDiffusion (Ours)



“An illustration of a beach in La La Land style”

More results are available on <https://syncdiffusion.github.io/>.

Application: Wide Image Generation

MultiDiffusion (Bar-Tal et al.)



SyncDiffusion (Ours)



“A cinematic view of a castle in the sunset”

More results are available on <https://syncdiffusion.github.io/>.

Diffusion Synchronization

The synchronization can be conducted in any canonical space (e.g., panorama space, texture space, etc).

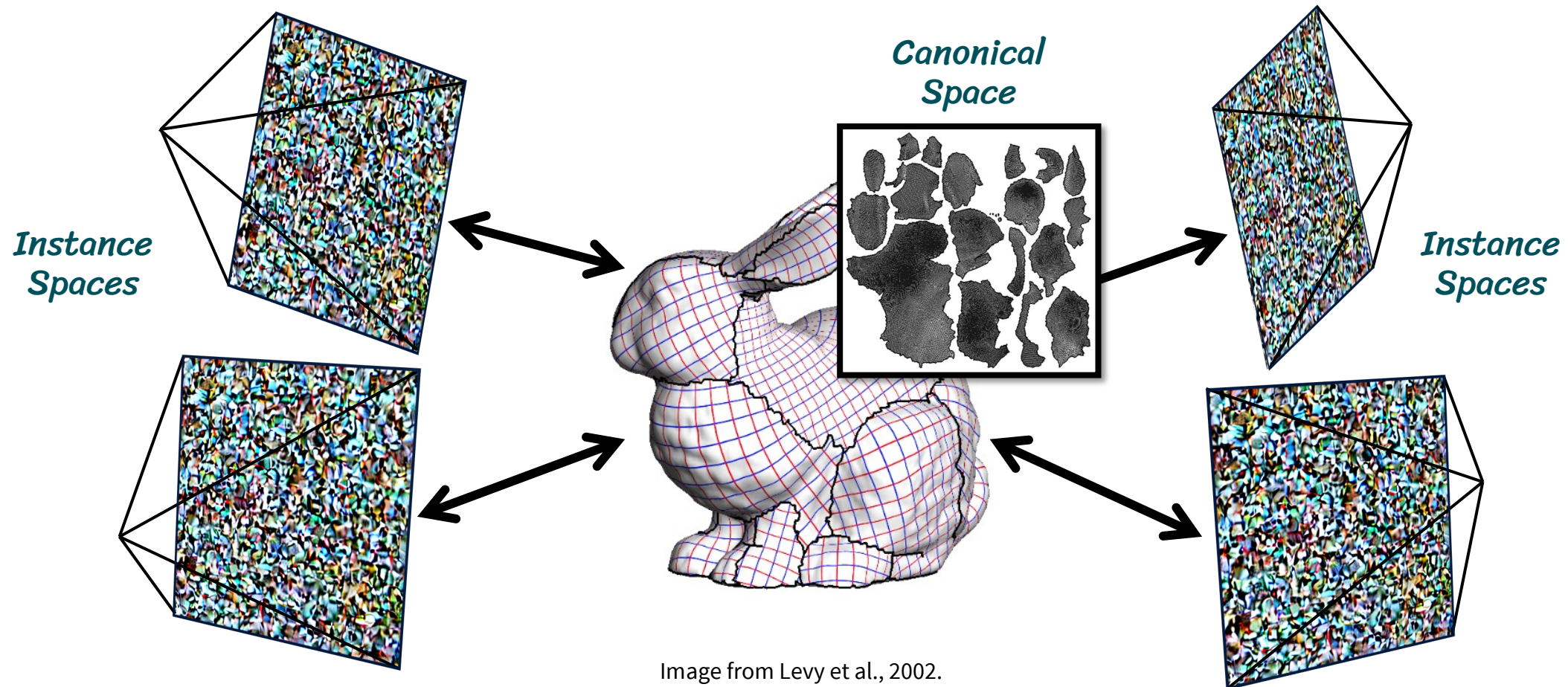


Image from Levy et al., 2002.

Diffusion Synchronization

1. **Unproject** Tweedie estimate of instance spaces into the canonical space.

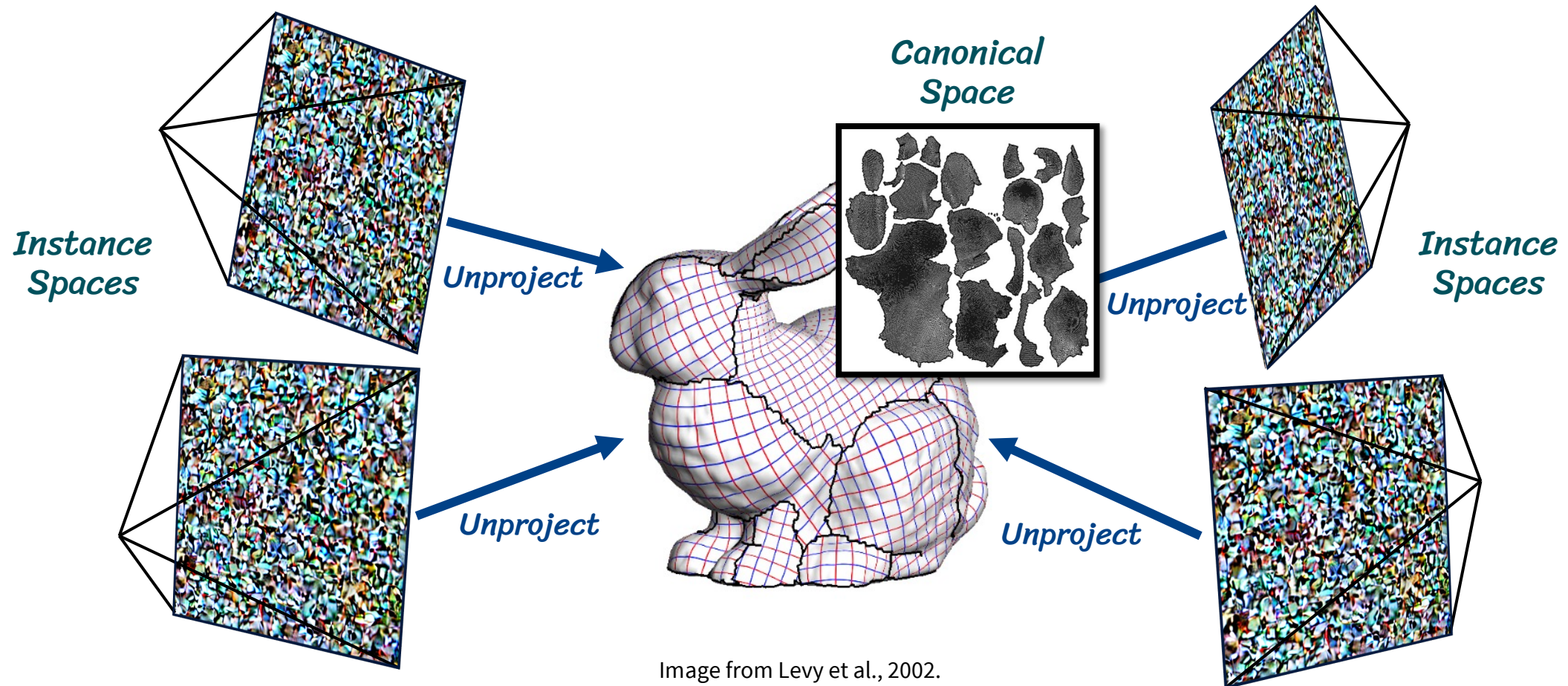
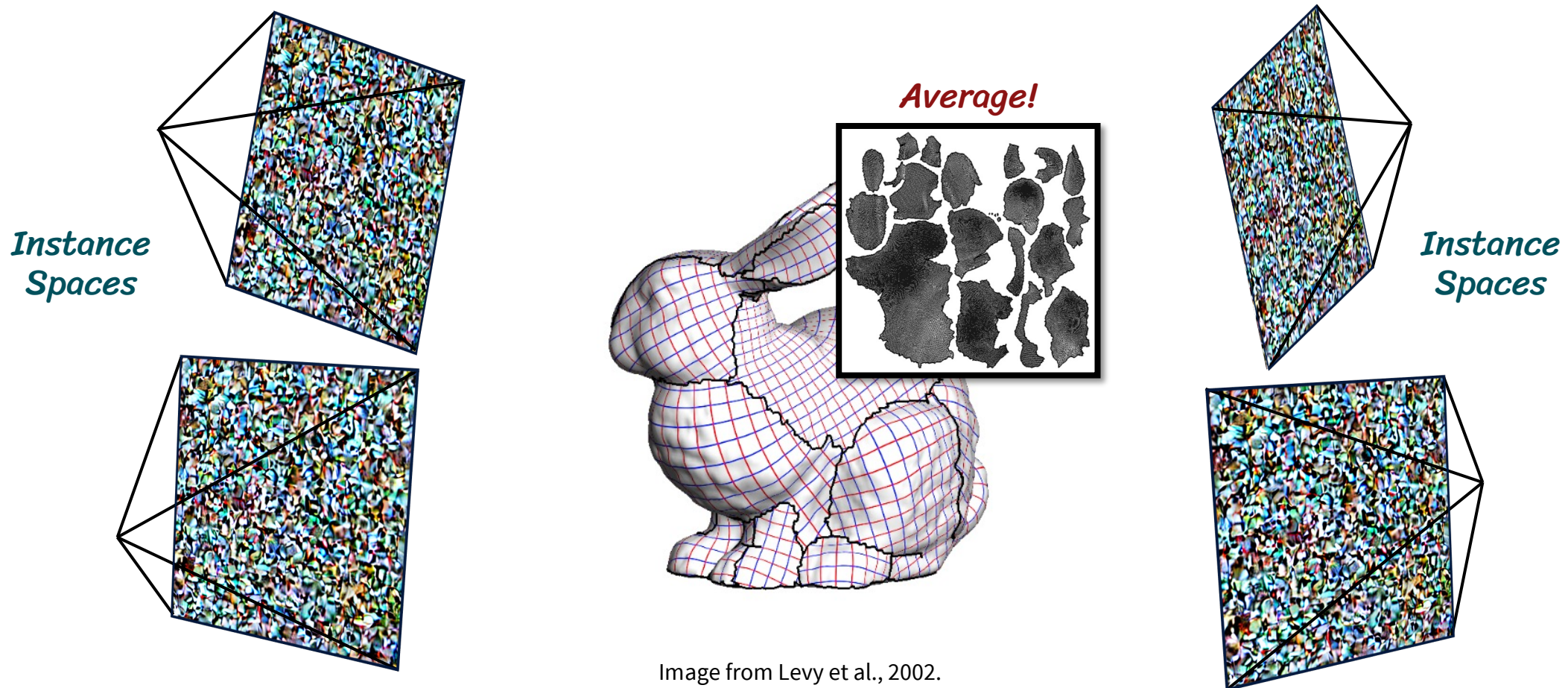


Image from Levy et al., 2002.

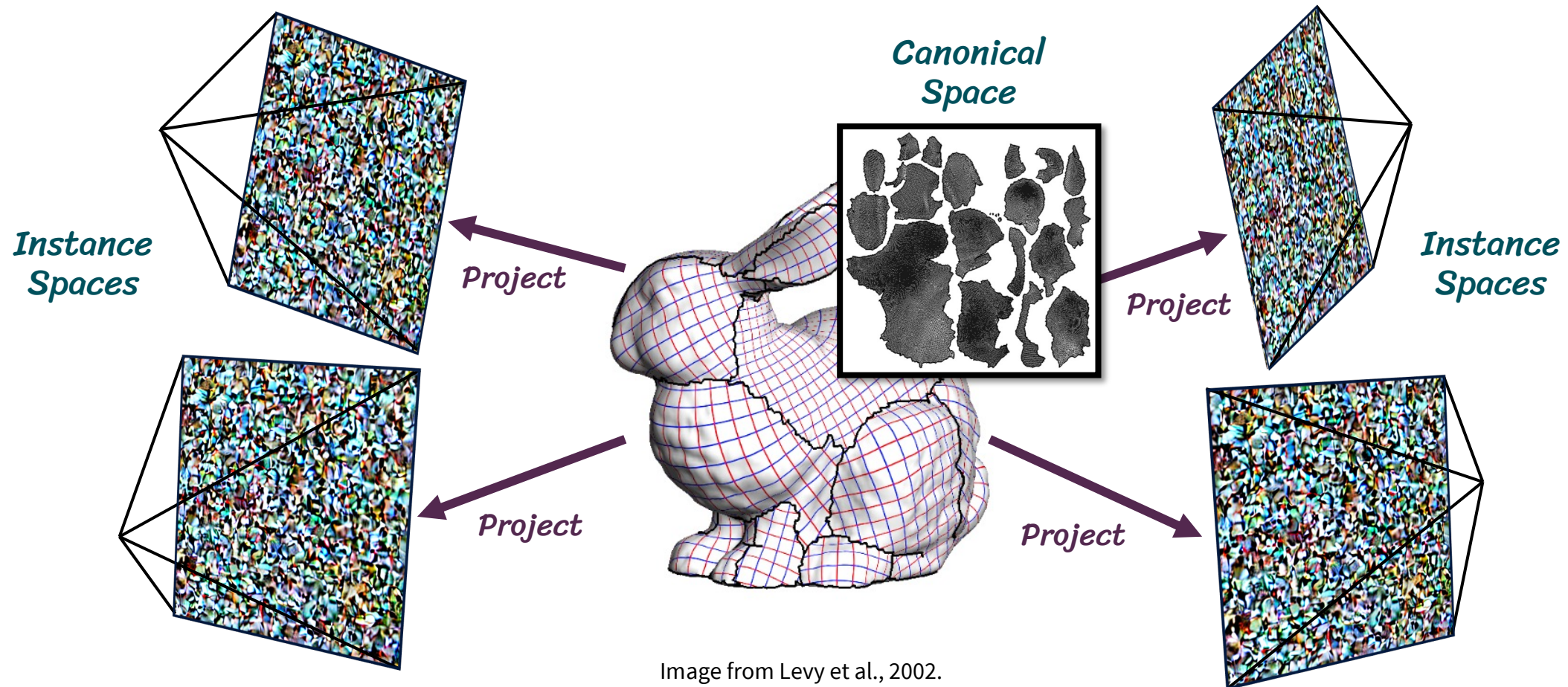
Content Synchronization

2. Compute the **loss**: the difference between each Tweedie estimate and their average in the canonical space.



Diffusion Synchronization

3. **Backpropagate** the loss to the noise **in** the instance spaces.



Application: 3D Texturing



"A hand carved wood turtle"



"A dumpster"



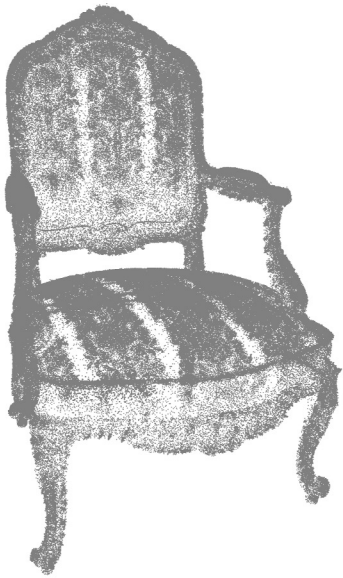
"A Chinese style lantern"



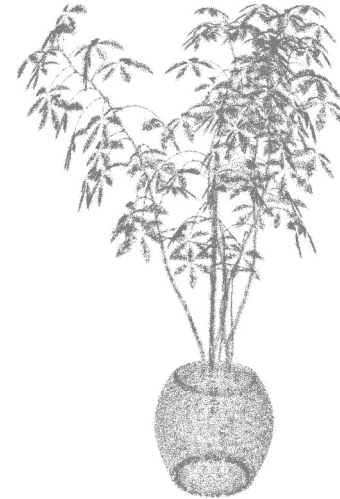
"A car with graffiti"

Application: 3D Texturing

Gaussian splat texture generation



"A majestic red chair"



"A photo of a tree with multicolored leaves"

Application: 3D Texturing

"A marble dresser"



"A wooden crate"



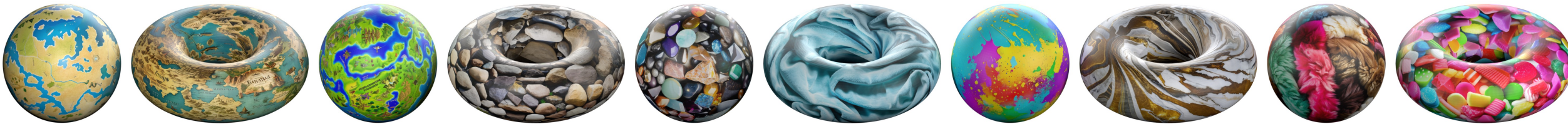
*"A mug with
spiral silver ornaments"*



*"A marble statue
of a wolf"*



"A cozy stone fireplace"



Application: 360° Panorama Generation

“Rocky desert landscape with towering saguaro cacti”



“Graffiti-covered alleyway with street art murals”



“Quirky steampunk workshop filled with gears and gadgets”



“Cozy neighborhood pub with outdoor seating”



“Abandoned factory with soft rays through dusty air”

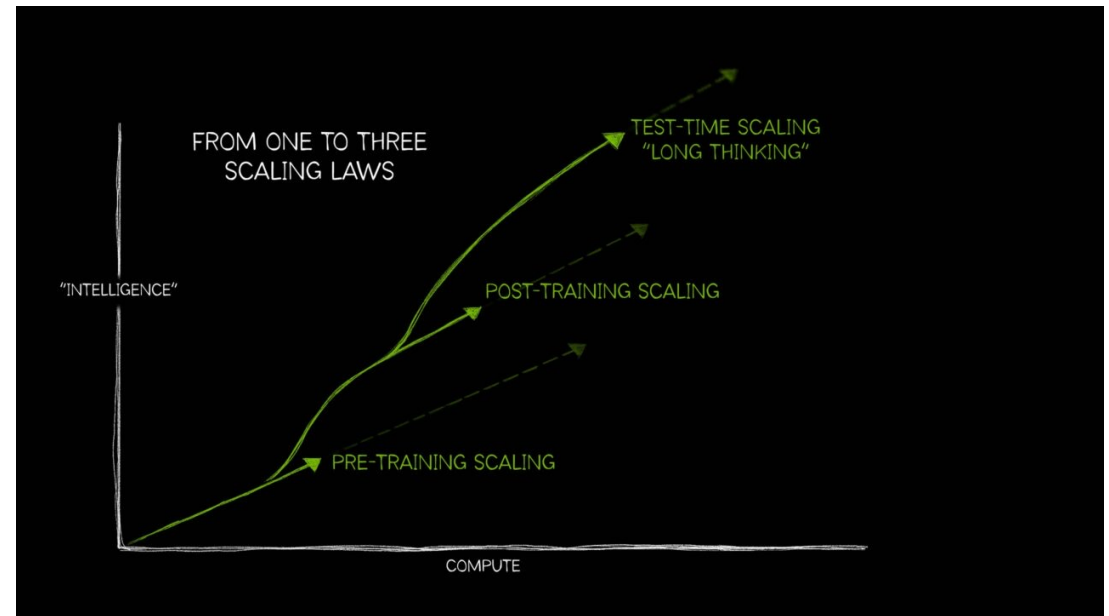
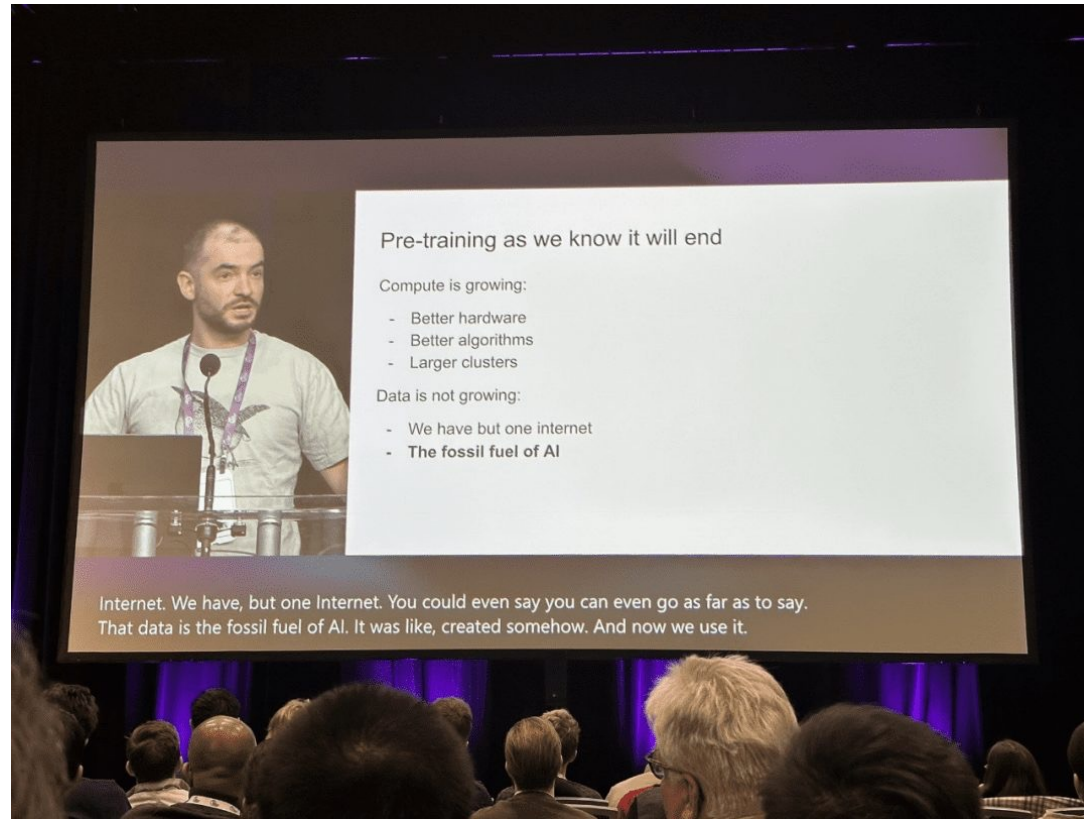


“Quaint canal lined with boats and cafes”



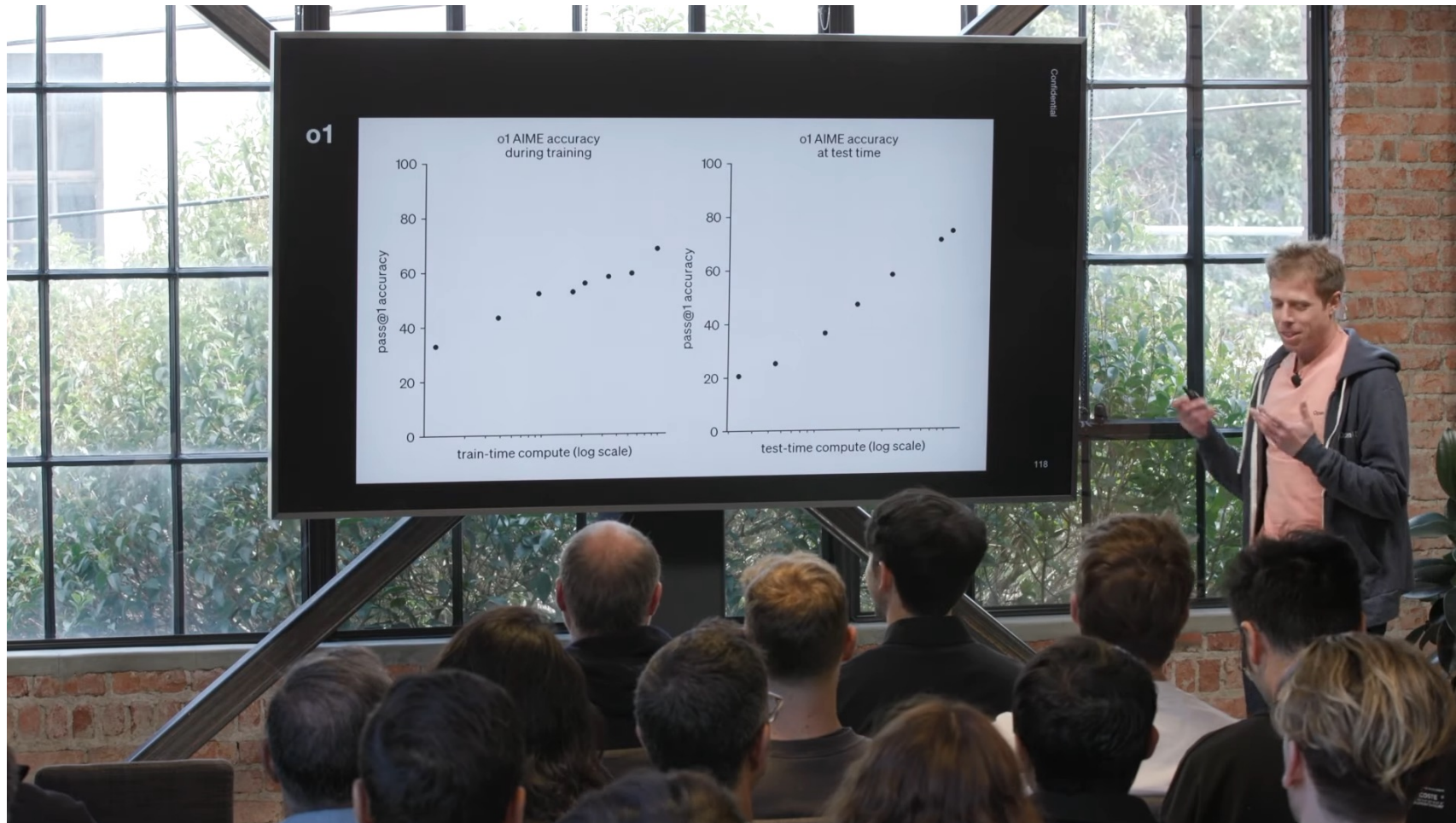
3. Particle Filtering

The End of the Pre-Training Era?

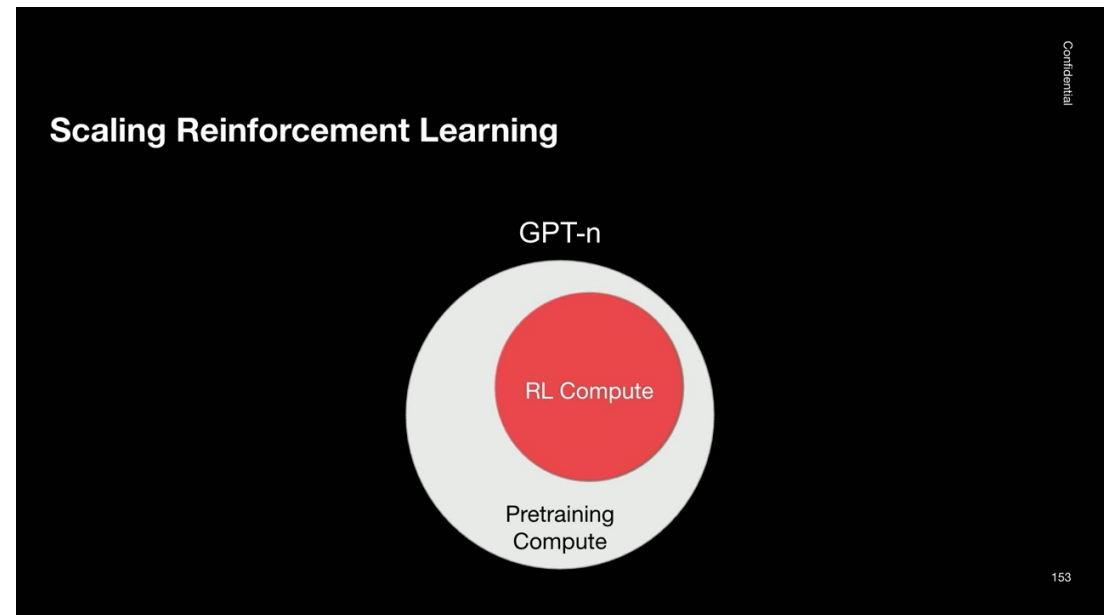
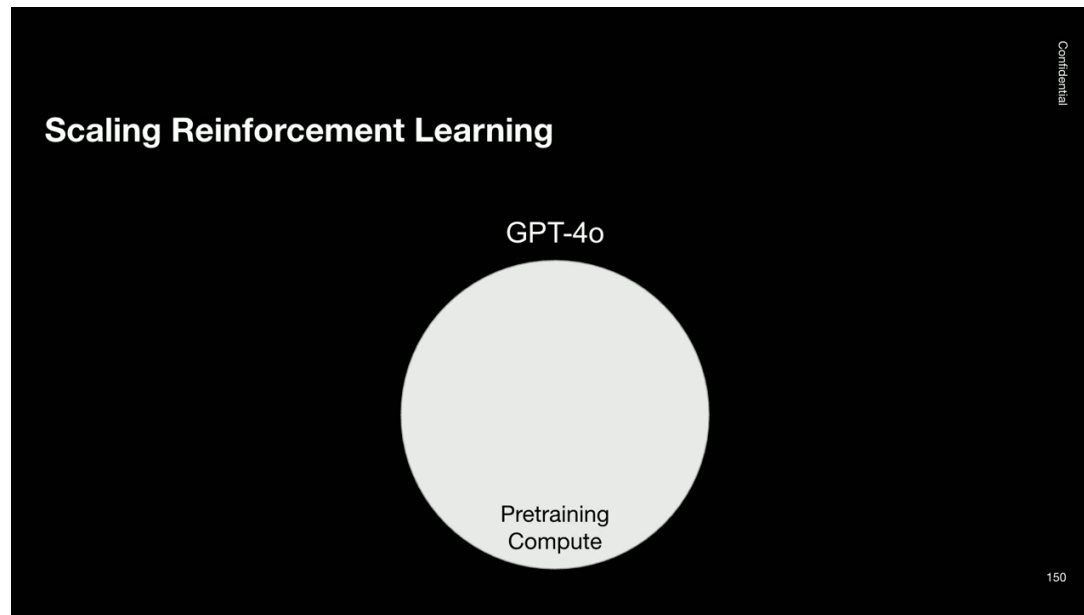


OpenAI's Dan Roberts Reasons

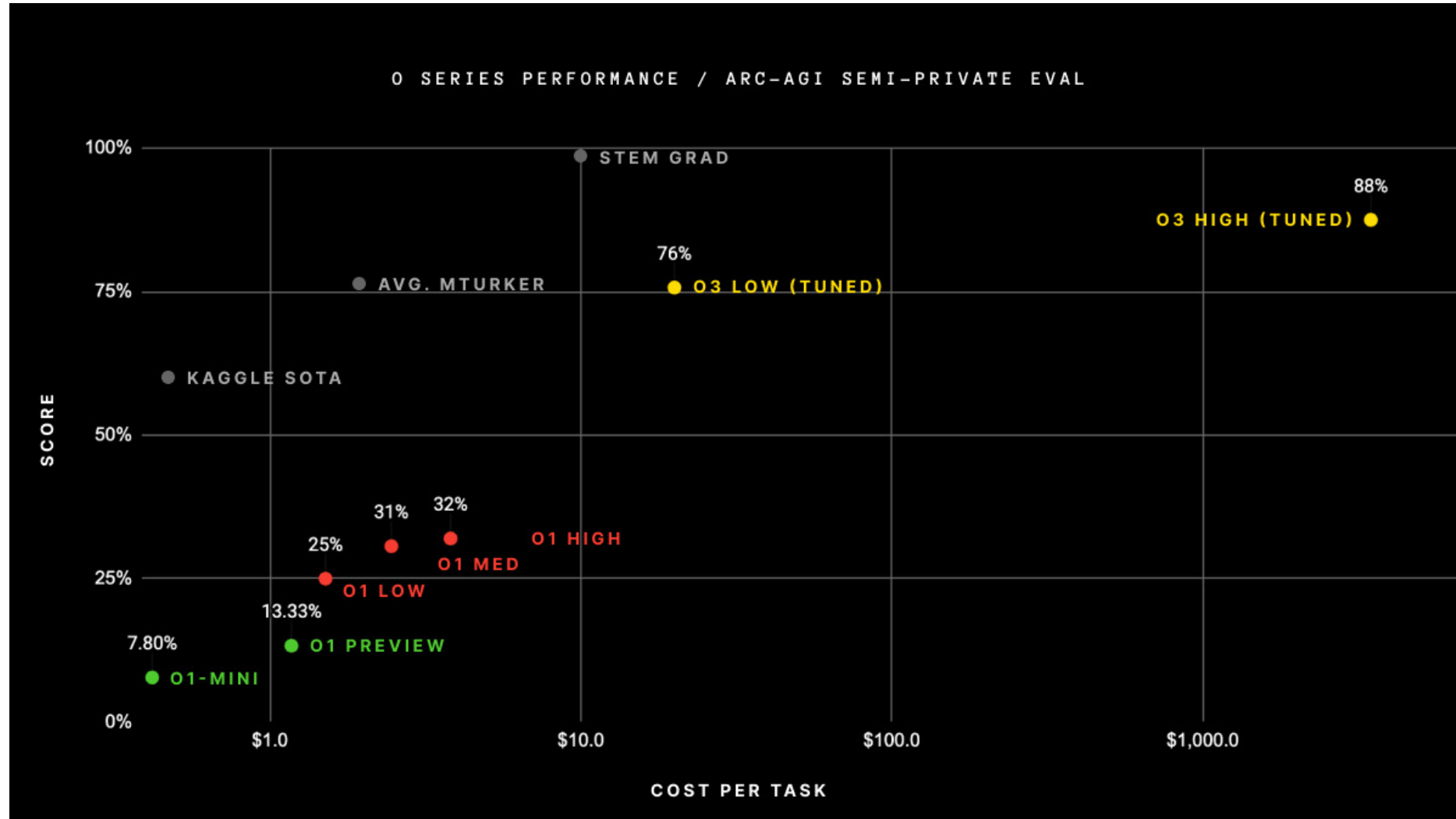
9 Years to AGI? OpenAI's Dan Roberts Reasons About Emulating Einstein





OpenAI's Dan Roberts Reasons




Inference-Time Scaling – GPT




Inference-Time Scaling – GPT




**o1-pro** Default 




Compare Try in Playground

REASONING

Higher

SPEED

Slowest

PRICE
\$150 • \$600
Input • Output

INPUT
  
Text, image

OUTPUT
  
Text

The o1 series of models are trained with reinforcement learning to think before they answer and perform complex reasoning. The o1-pro model uses more compute to think harder and provide consistently better answers.

o1-pro is available in the [Responses API only](#) to enable support for multiple tool calls and multi-turn model interactions before responding to API requests.

- ✦ 200,000 context window
- ↪ 100,000 max output tokens
- 📅 Oct 01, 2023 knowledge cutoff
- 💡 Reasoning token support

Pricing

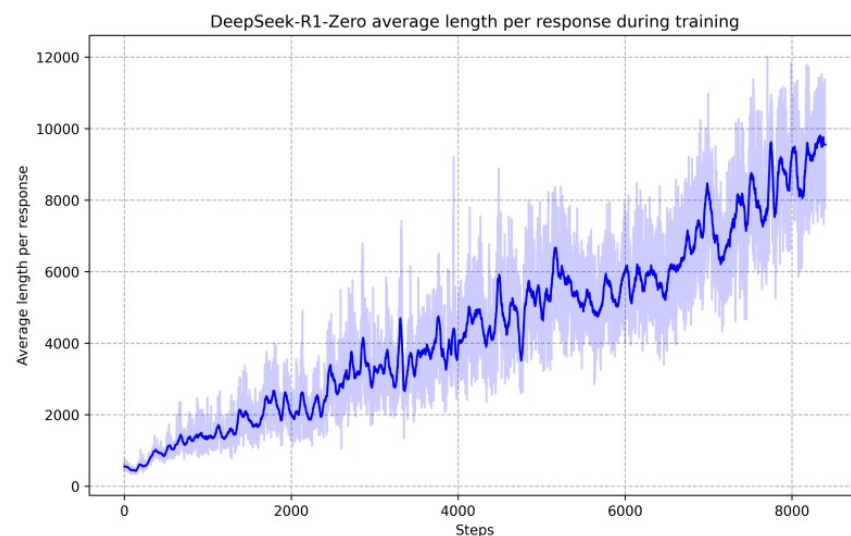
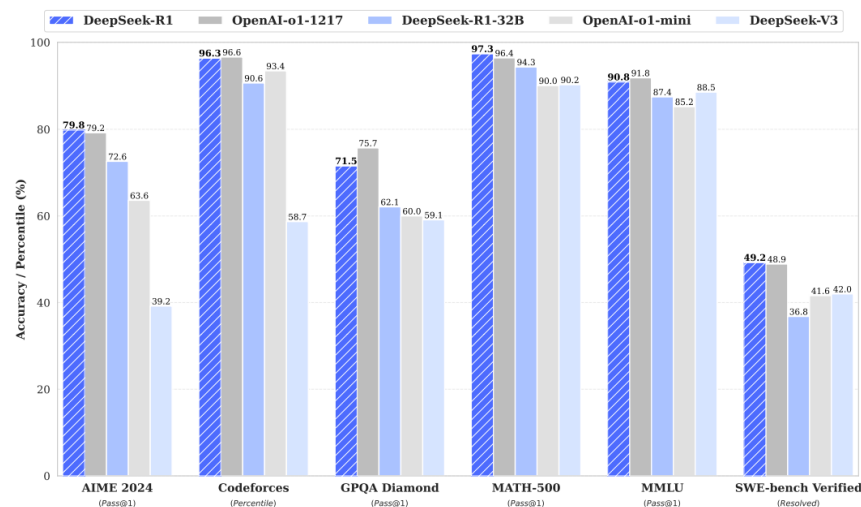
Pricing is based on the number of tokens used. For tool-specific models, like search and computer use, there's a fee per tool call. See details in the [pricing page](#).

Text tokens Per 1M tokens • Batch API price ☒

Input
\$150.00

Output
\$600.00

Inference-Time Scaling – DeepSeek



Question: If $a > 1$, then the sum of the real solutions of $\sqrt{a - \sqrt{a+x}} = x$ is equal to

Response: <think>

To solve the equation $\sqrt{a - \sqrt{a+x}} = x$, let's start by squaring both ...

$$(\sqrt{a - \sqrt{a+x}})^2 = x^2 \implies a - \sqrt{a+x} = x^2.$$

Rearrange to isolate the inner square root term:

$$(a - x^2)^2 = a + x \implies a^2 - 2ax^2 + (x^2)^2 = a + x \implies x^4 - 2ax^2 - x + (a^2 - a) = 0$$

...

Wait, wait. Wait. That's an aha moment I can flag here.

Let's reevaluate this step-by-step to identify if the correct sum can be ...

We started with the equation:

$$\sqrt{a - \sqrt{a+x}} = x$$

First, let's square both sides:

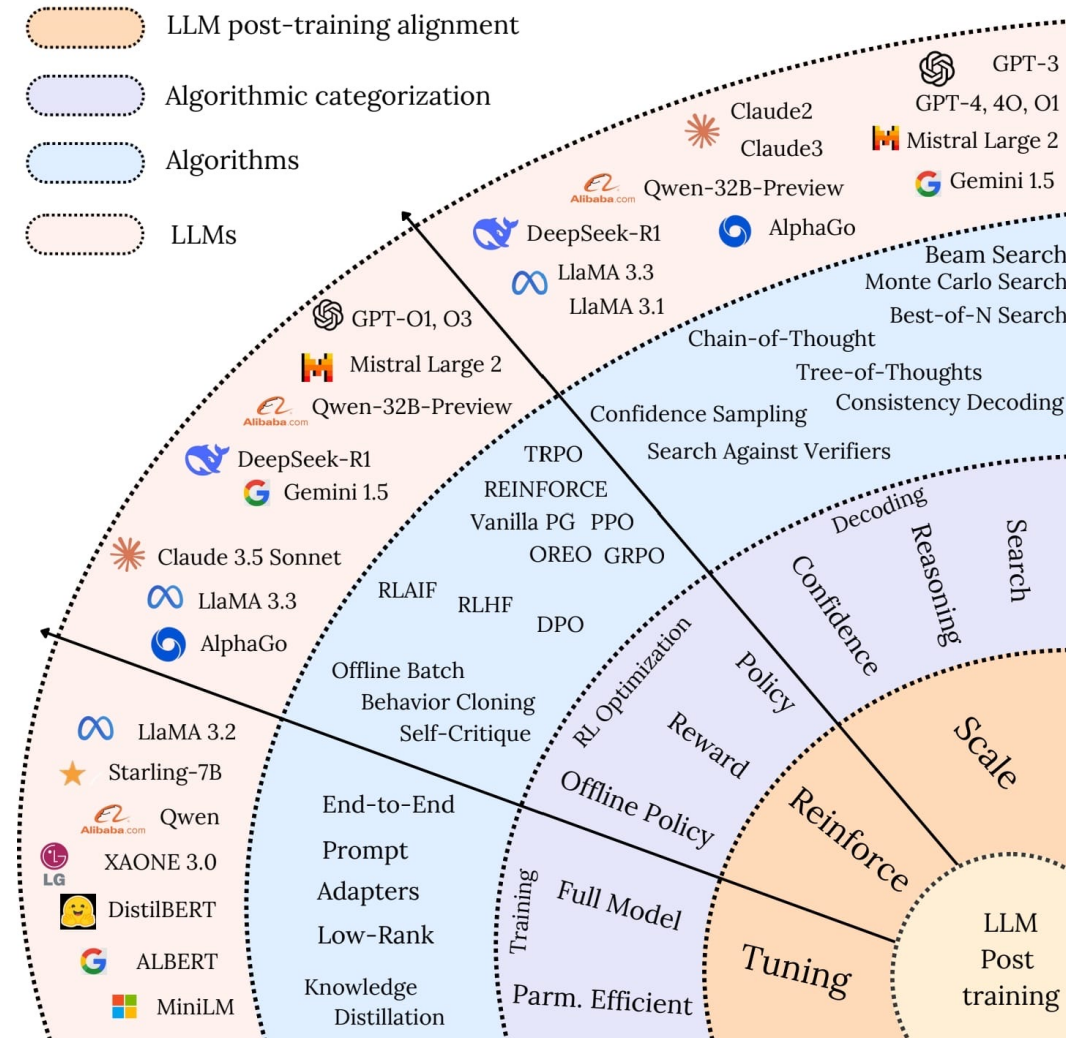
$$a - \sqrt{a+x} = x^2 \implies \sqrt{a+x} = a - x^2$$

Next, I could square both sides again, treating the equation: ...

...

Table 3 | An interesting “aha moment” of an intermediate version of DeepSeek-R1-Zero. The model learns to rethink using an anthropomorphic tone. This is also an aha moment for us, allowing us to witness the power and beauty of reinforcement learning.

Inference-Time Scaling



Failure Cases of GPT-4o

GPT 4o



“Four drums, **seven** tomatoes, and five candles.”

Failure Cases of GPT-4o

GPT 4o



“In a room, all the chairs are occupied except one.”

Goal

GPT 4o



FLUX + Inference-Time Scaling



“Four drums, **seven** tomatoes, and five candles.”

Goal

GPT 4o



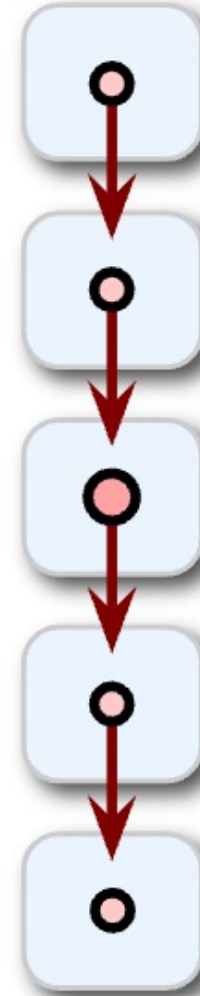
FLUX + Inference-Time Scaling



“In a room, all the chairs are **occupied** except one.”

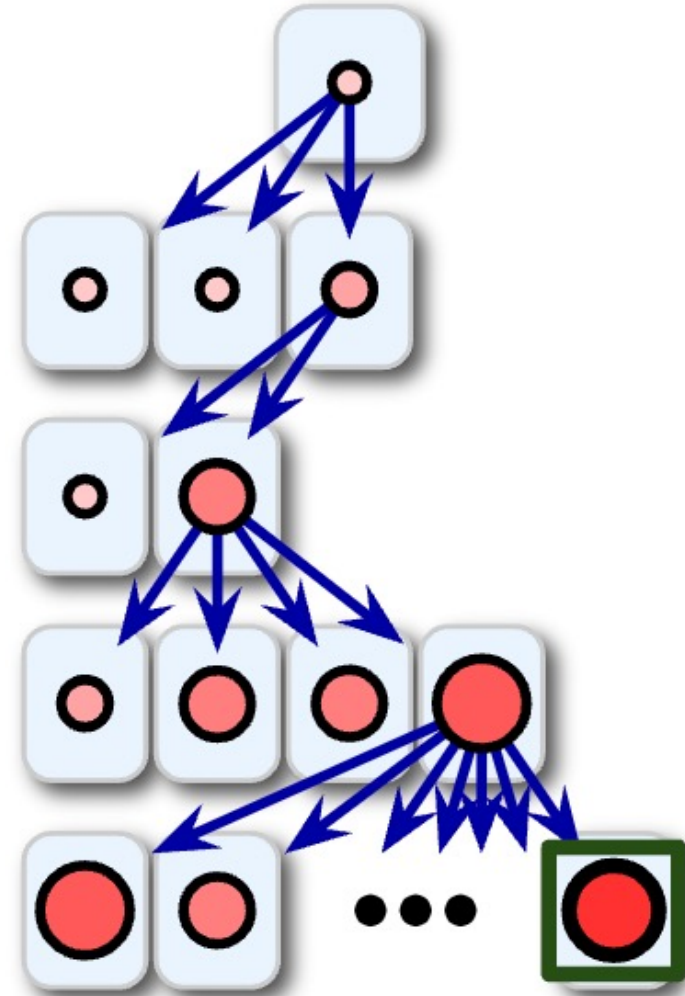
Particle Filtering

- The denoising process of a diffusion model can be viewed as a **Sequential Monte Carlo** process with a **single trajectory**.
- Only one particle is sampled at each step.

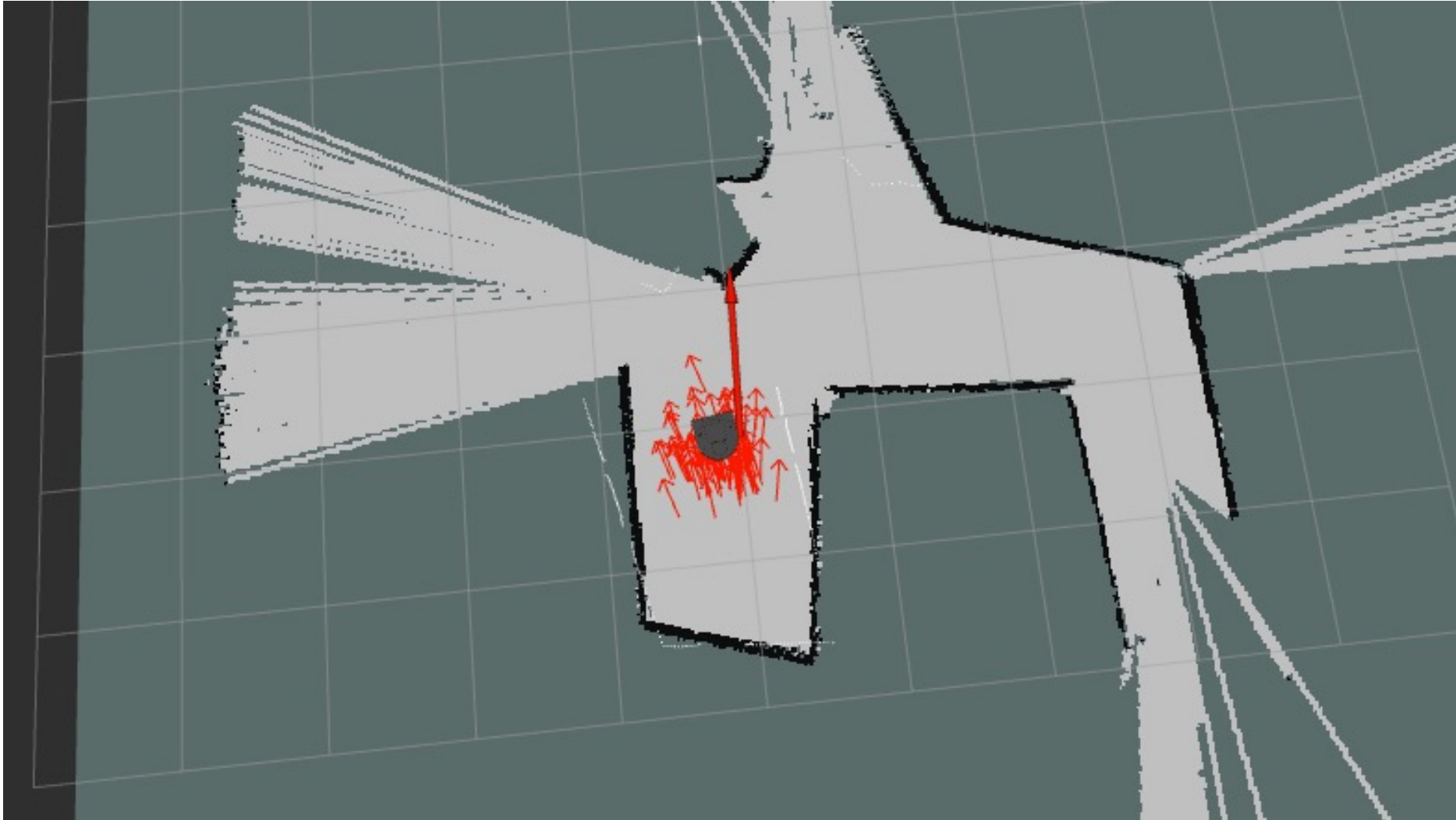


Particle Filtering

Could this be extended to
particle filtering with
multiple particles at each
denoising timestep?



Particle Filtering



Applications: Text Alignment

“Every painting in the gallery is framed and hung straight, except for one that is hanging crooked.”



FLUX



Ours

Applications: Text Alignment

*“The arcade machine is bigger than the television but **smaller** than the refrigerator.”*



FLUX



Ours

Applications: Text Alignment

*“A large suitcase is placed **beside** an open closet, with a folded jacket resting on top where a pair of shoes sit side by side in front of it.”*



FLUX



Ours

Applications: Text Alignment

(+) A nurse

(−) Stethoscope, hat, mask



FLUX



Ours

Applications: Object Count

2 cups, 3 paintings, 4 lamps, 4 bananas.



FLUX



Ours

Applications: Object Count

7 balloons, 4 bears, 4 swans.



Best of N



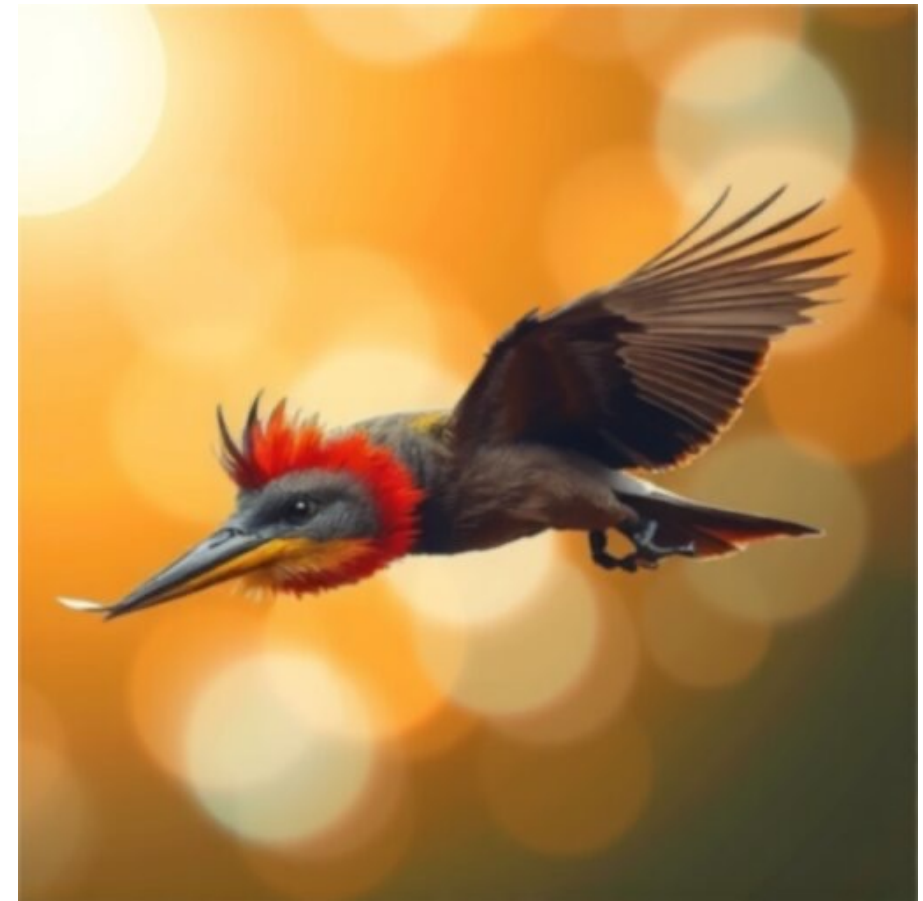
Ours

Applications: Aesthetic Image Generation

Bird



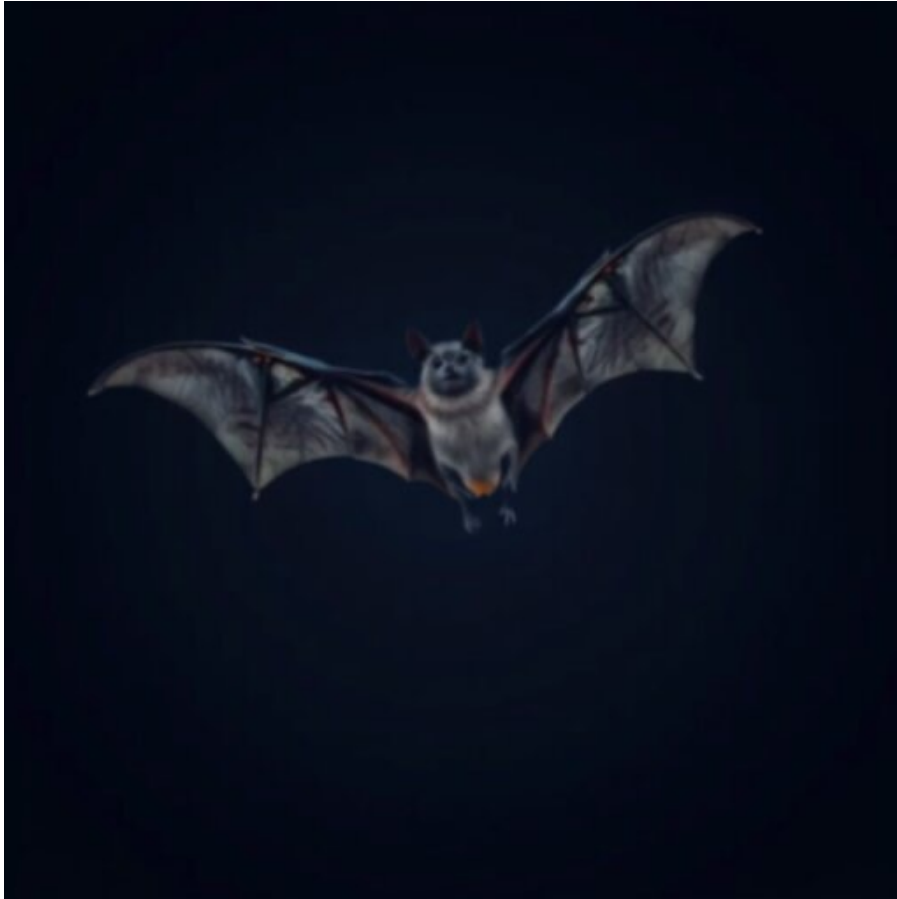
Best of N



Ours

Applications: Aesthetic Image Generation

Bat



Best of N



Ours

Application: Orientation Grounding

3D orientation grounding for multiple open-vocabulary objects in image generation.

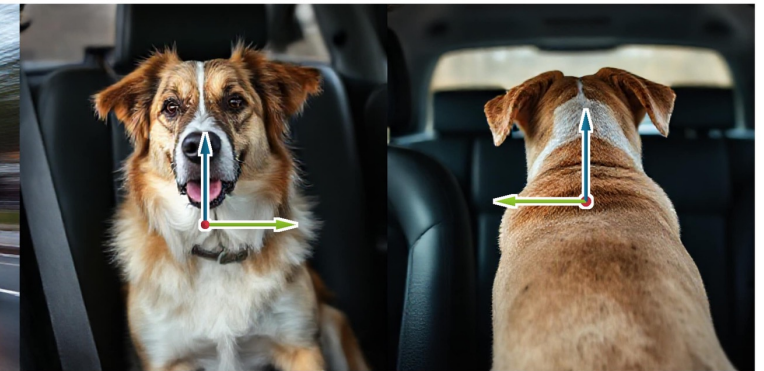
*“a **teddy bear** holding a stick next to a block toy.”*



*“a **bus** driving down the road really fast.”*



*“a **dog** that is sitting down in a backseat.”*



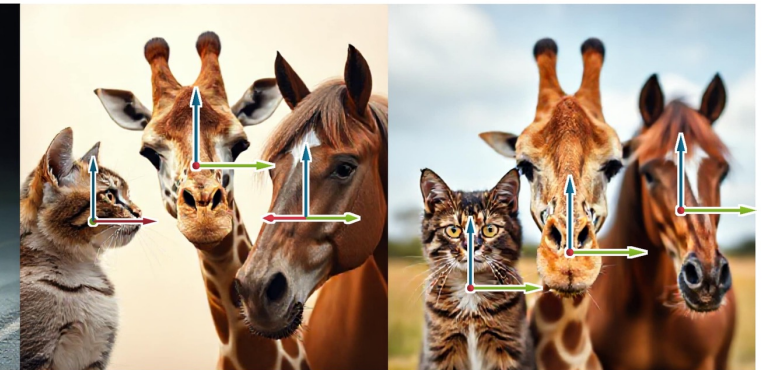
*“a **man**, and a woman.”*



*“a **motorcycle**, and a bear.”*

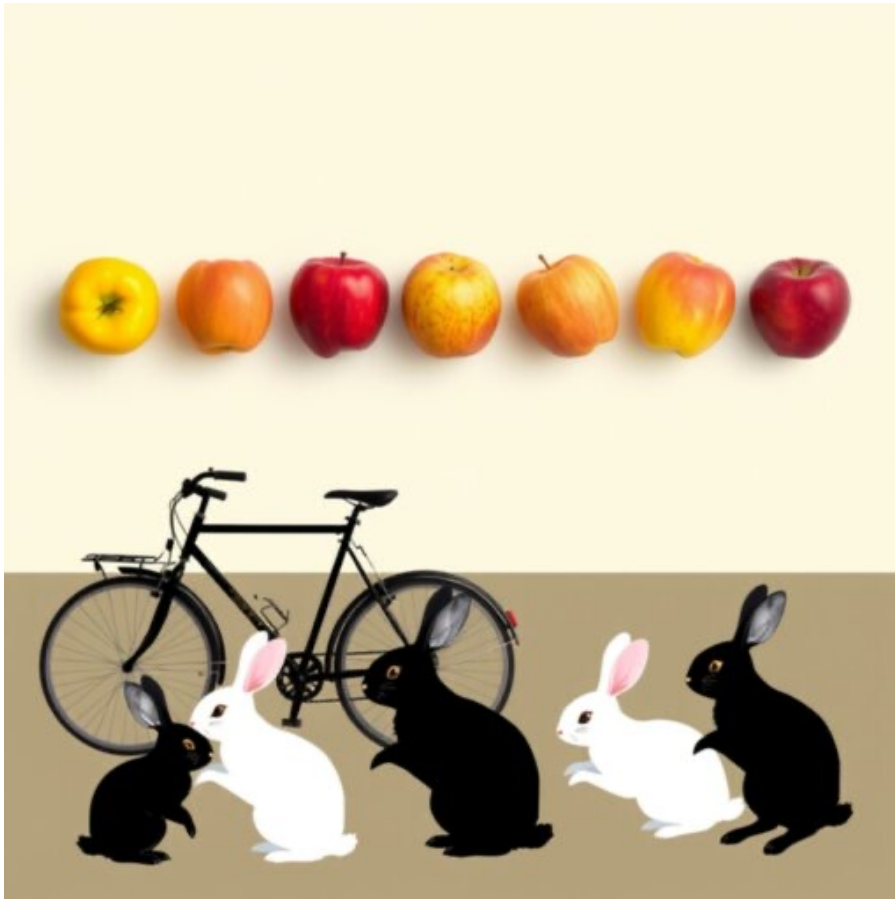


*“a **cat**, a giraffe, and a horse.”*



Applications: Object Count

8 apples, 3 bicycles, 5 rabbits.



Best of N



Ours

Applications: Object Count



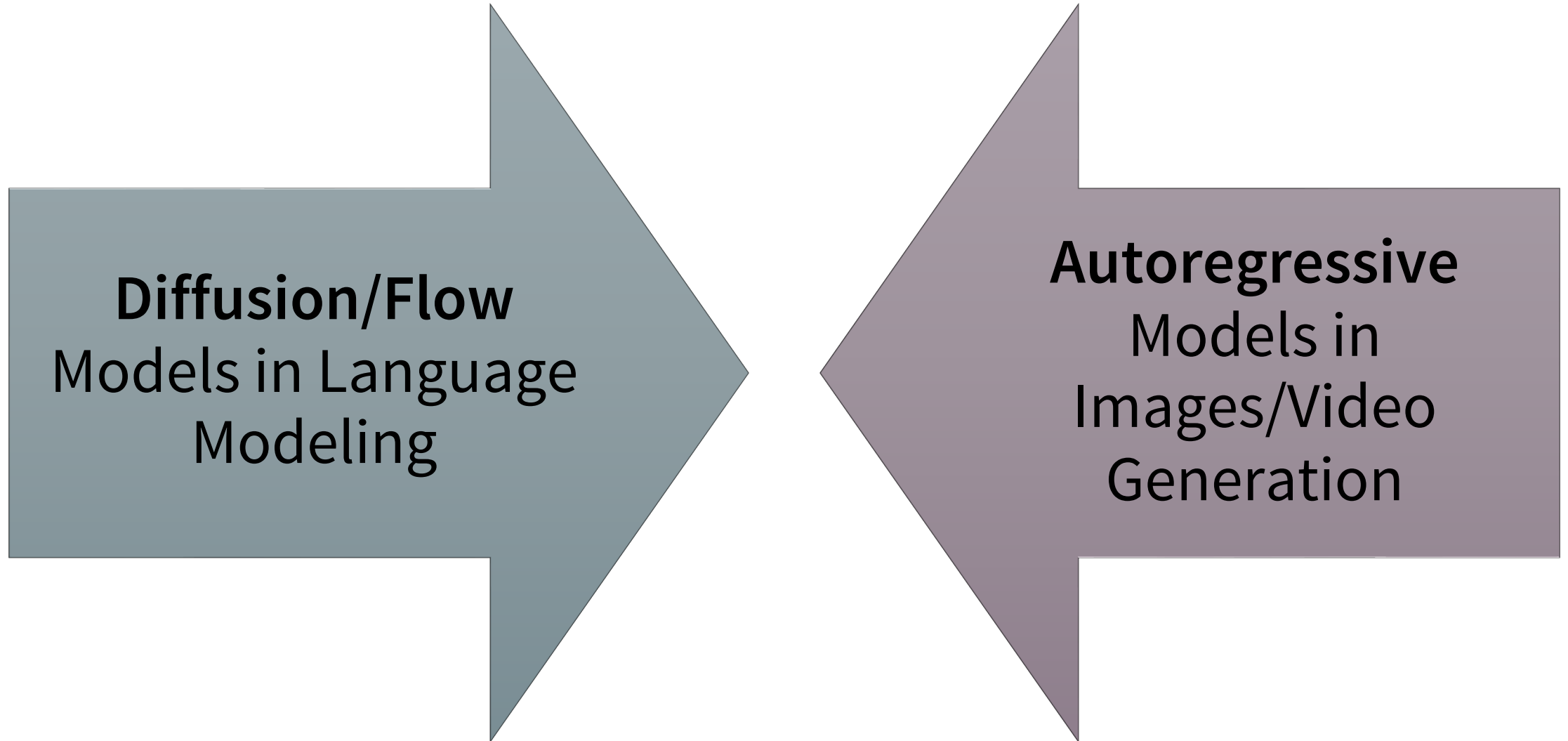
82 Blueberries



33 Blueberries

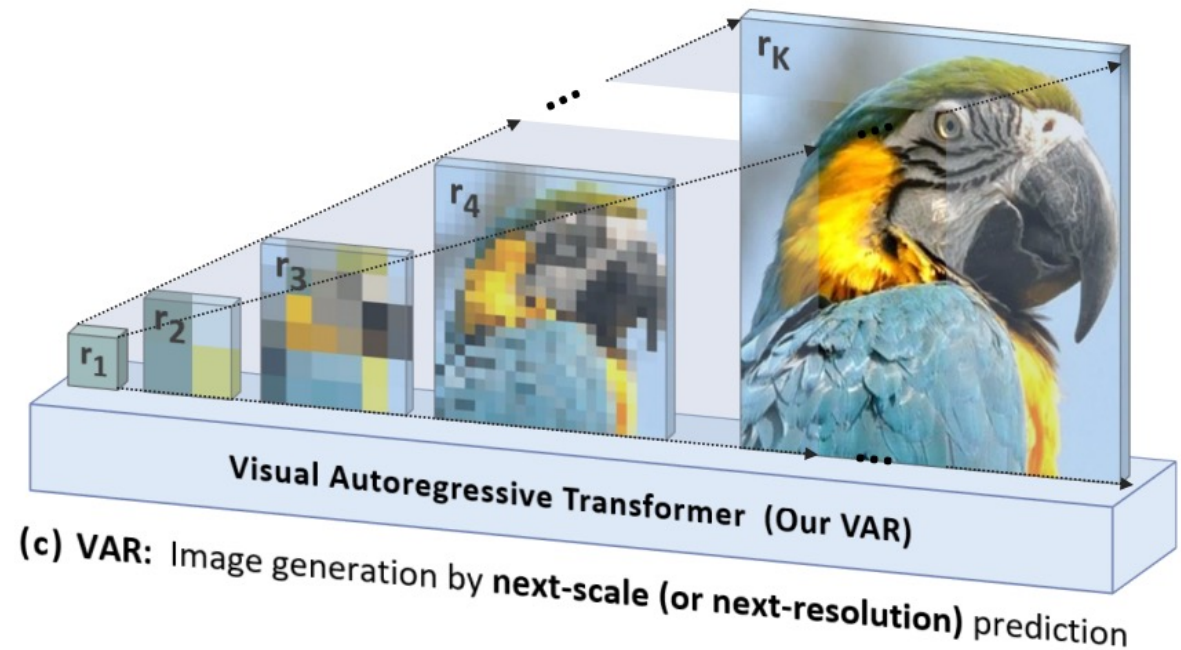
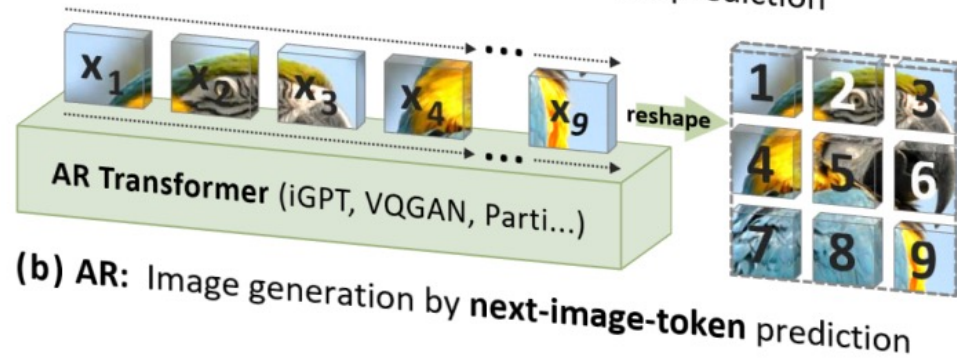
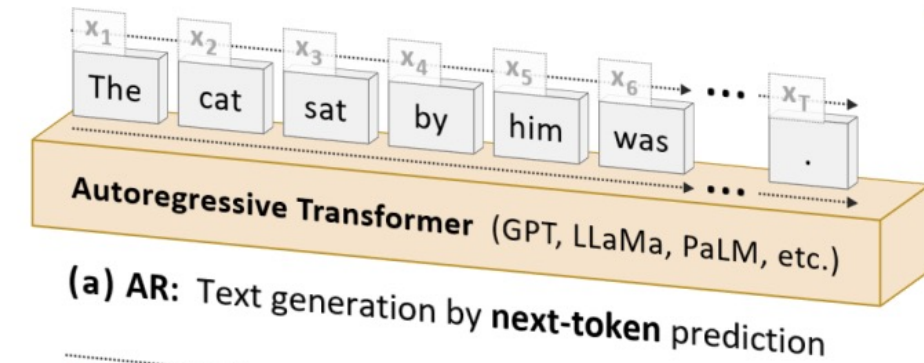
What's the Next?

Diffusion/Flow vs. Autoregressive



Visual Autoregressive Model

Three Different Autoregressive Generative Models



Visual Autoregressive Model

- (−) Lower image quality and lower generation speed.
- (−) Still a long way to go for efficient fine-tuning or training-free guided generation.
- (+) Could be better suited for integration with multimodal foundation models.

Large Language Diffusion Models



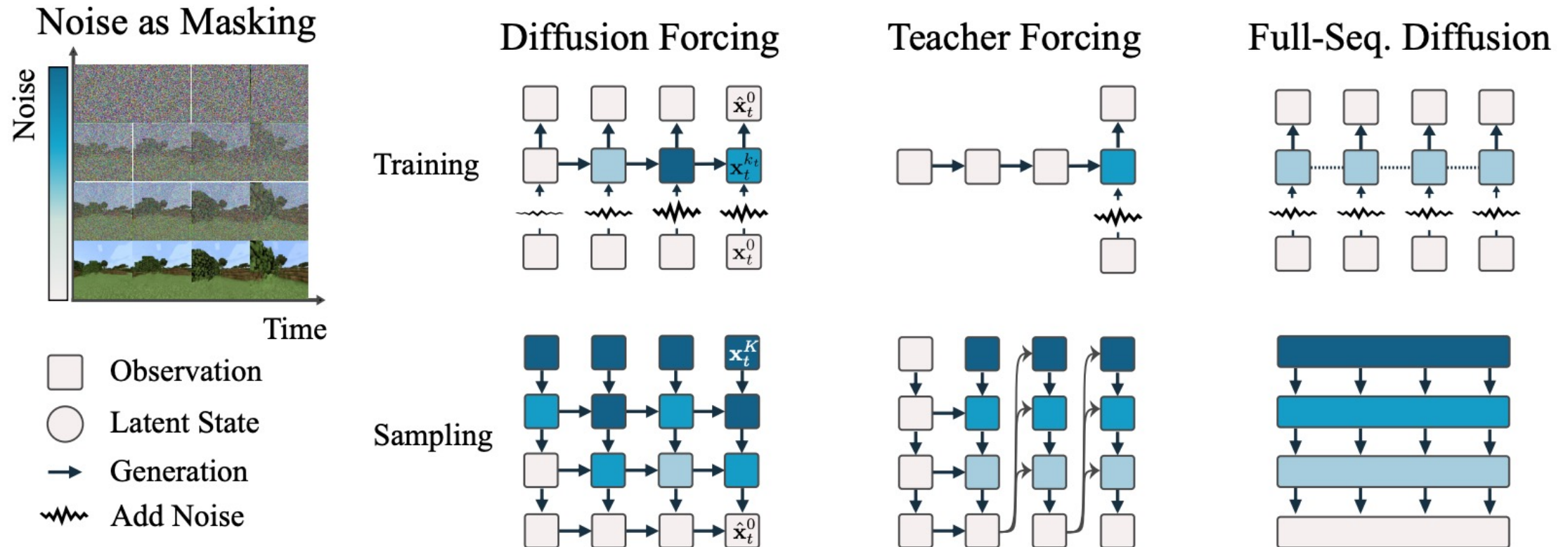
Mercury, the first commercial-grade diffusion LLM by Inception Labs

Diffusion/Flow vs. Autoregressive

Controllability will be one of the key factors
in determining the winner!

Hybrid: Autoregressive Diffusion Models

Diffusion/flow for each item;
autoregressive for the entire set, sequence, or graph.



Hybrid: Autoregressive Diffusion Models

- Video generation
- Other time-series data: audio, finance, weather, health, physics simulation, etc.
- Set/graph generation: multi-view, scene graph, CAD, molecule, etc.

*How can inference-time scaling be specialized for
autoregressive diffusion models?*

The Present and Future of Image and Video Generation Technologies



Minhyuk Sung
KAIST Visual AI Group
<https://visualai.kaist.ac.kr/>

