

SIGGRAPH 2024 Course: Generative Models for Visual Content Editing and Creation

ANYI RAO, Stanford University, USA
YUANBO XIANGLI, Cornell University, USA
YUWEI GUO, Chinese University of Hong Kong, China
MIA TANG, Stanford University, USA
CHENLIN MENG, Stanford University, USA
MANEESH AGRAWALA, Stanford University, USA



Fig. 1. Synthetic images generated with ControlNet [Zhang et al. 2023]

Authors' addresses: Anyi Rao, Stanford University, Palo Alto, CA, USA, anyirao@stanford.edu; Yuanbo Xiangli, Cornell University, New York, USA; Yuwei Guo, Chinese University of Hong Kong, Hong Kong, China; Mia Tang, Stanford University, Palo Alto, CA, USA; Chenlin Meng, Stanford University, Palo Alto, CA, USA; Maneesh Agrawala, Stanford University, Palo Alto, CA, USA.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish,

Interest in generative models is surging in academia and industry, with their impressive capabilities and creativity outputs. Crucially, these models are also providing users with a growing degree of control over the generation process via texts or visuals prompts. Concretely, large-scale text-to-image foundation models like Stable Diffusion [Rombach et al. 2021],

to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
© 2024 Association for Computing Machinery.
0730-0301/2024/8-ART \$15.00
<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

SDXL [Podell et al. 2023], eDiff-I [Balaji et al. 2022], DALL-E 3 [Betker et al. 2023]; text-to-video foundation models like Imagen Video [Ho et al. 2022] and Make-a-video [Singer et al. 2022], Sora [OpenAI 2024] have boosted the growth of visual content editing and generation. Representatively, works such as AnimateDiff [Guo et al. 2023], ControlNet [Zhang et al. 2023] democratized video creation with diverse user-defined conditions, and have become practical tools for graphic designs and personalized media. There has also been a revolution in 3D asset generation in terms of fidelity and efficiency. Harvesting the powerful prior of 2D diffusion models, works such as DreamFusion [Poole et al. 2022], Magic3D [Lin et al. 2023], Zero123 [Liu et al. 2023], Wonder3D [Long et al. 2023] were enabled high-quality text-and image-to-3D object generation, with plausible geometry and physical properties to support their usage in gaming and simulation tasks. At the meantime, the emergence of high-quality large-scale 3D data [Deitke et al. 2023a,b; Yu et al. 2023] also empowered direct generative model training in 3D space [Hong et al. 2023; Xu et al. 2023]. Inspired by the success of 3D asset generation, scene-level 3D synthesis also gained increasing interest. Work such as GeNVS [Chan et al. 2023], ReconFusion [Wu et al. 2023] also benefit from 2D diffusion priors to achieve high-quality novel view synthesis. Another branch of work, such as AssetField [Xiangli et al. 2023], BlockPlanner [Xu et al. 2021] regard scenes as a composition of 3D assets guided by layouts, that can be generatively modeled in a data-driven manner whilst guarantee user controllability.

This course covers the advances in generative models over the last few years, with a slight shift towards the controllability and creativity tasks enabled by generative models. We will first go over the fundamental machine learning and deep learning techniques relevant to generative models. Next, we will showcase recent representative work in controllable image, video and 3D content generation and compositional representation learning. Finally, we will conclude with a discussion on the future application of this technology, societal impact and open research problems. After the course, the attendees will learn basic knowledge about diffusion models and how such models can be applied to different applications.

P.S. Website: <https://cveu.github.io/event/sig2024.html>; Twitter: https://twitter.com/cveu_workshop

CCS Concepts: • **Information systems** → Multimedia content creation.

Additional Key Words and Phrases: Generative Models, Creativity Support

ACM Reference Format:

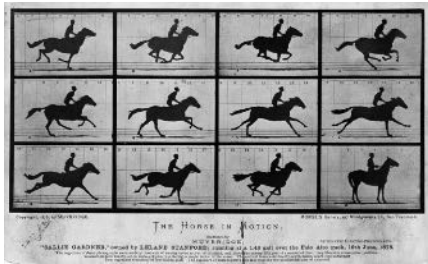
Anyi Rao, Yuanbo Xiangli, Yuwei Guo, Mia Tang, Chenlin Meng, and Maneesh Agrawala. 2024. SIGGRAPH 2024 Course: Generative Models for Visual Content Editing and Creation. *ACM Trans. Graph.* 1, 1 (August 2024), 2 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

REFERENCES

- Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, et al. 2022. ediffi: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324* (2022).
- James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. 2023. Improving image generation with better captions. *Computer Science*. <https://cdn.openai.com/papers/dall-e-3.pdf> 2, 3 (2023), 8.
- Eric R. Chan, Koki Nagano, Matthew A. Chan, Alexander W. Bergman, Jeong Joon Park, Axel Levy, Miika Aittala, Shalini De Mello, Tero Karras, and Gordon Wetzstein. 2023. GeNVS: Generative Novel View Synthesis with 3D-Aware Diffusion Models. In *arXiv*.
- Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan Fan, Christian Laforte, Vikram Voleti, Samir Yitzhak Gadre, Eli VanderBilt, Aniruddha Kembhavi, Carl Vondrick, Georgia Gkioxari, Kiana Ehsani, Ludwig Schmidt, and Ali Farhadi. 2023a. Objaverse-XL: A Universe of 10M+ 3D Objects. *arXiv preprint arXiv:2307.05663* (2023).
- Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. 2023b. Objaverse: A Universe of Annotated 3D Objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13142–13153.
- Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh Agrawala, Dahua Lin, and Bo Dai. 2023. AnimateDiff: Animate Your Personalized Text-to-Image Diffusion Models without Specific Tuning.
- Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. 2022. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303* (2022).
- Yicong Hong, Kai Zhang, Jiuxiang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan Sunkavalli, Trung Bui, and Hao Tan. 2023. LRM: Large Reconstruction Model for Single Image to 3D. *arXiv:2311.04400* [cs.CV]
- Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. 2023. Magic3D: High-Resolution Text-to-3D Content Creation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. 2023. Zero-1-to-3: Zero-shot One Image to 3D Object. *arXiv:2303.11328* [cs.CV]
- Xiaoxiao Long, Yuan-Chen Guo, Cheng Lin, Yuan Liu, Zhiyang Dou, Lingjie Liu, Yuexin Ma, Song-Hai Zhang, Marc Habermann, Christian Theobalt, et al. 2023. Wonder3D: Single Image to 3D using Cross-Domain Diffusion. *arXiv preprint arXiv:2310.15008* (2023).
- OpenAI. 2024. Video generation models as world simulators. <https://openai.com/research/video-generation-models-as-world-simulators>.
- Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. 2023. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952* (2023).
- Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. 2022. DreamFusion: Text-to-3D using 2D Diffusion. *arXiv* (2022).
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2021. High-Resolution Image Synthesis with Latent Diffusion Models. *arXiv:2112.10752* [cs.CV]
- Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. 2022. Make-a-video: Text-to-video generation without text-video data. *arXiv preprint arXiv:2209.14792* (2022).
- Rundi Wu, Ben Mildenhall, Philipp Henzler, Keunhong Park, Ruiqi Gao, Daniel Watson, Pratul P. Srinivasan, Dor Verbin, Jonathan T. Barron, Ben Poole, and Aleksander Holynski. 2023. ReconFusion: 3D Reconstruction with Diffusion Priors. *arXiv* (2023).
- Yuanbo Xiangli, Linning Xu, Xingang Pan, Nanxuan Zhao, Bo Dai, and Dahua Lin. 2023. AssetField: Assets Mining and Reconfiguration in Ground Feature Plane Representation. *arXiv:2303.13953* [cs.CV]
- Linning Xu, Yuanbo Xiangli, Anyi Rao, Nanxuan Zhao, Bo Dai, Ziwei Liu, and Dahua Lin. 2021. BlockPlanner: City Block Generation With Vectorized Graph Representation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5077–5086.
- Yinghao Xu, Hao Tan, Fujun Luan, Sai Bi, Peng Wang, Jiahao Li, Zifan Shi, Kalyan Sunkavalli, Gordon Wetzstein, Zexiang Xu, and Kai Zhang. 2023. DMV3D: Denoising Multi-View Diffusion using 3D Large Reconstruction Model. *arXiv:2311.09217* [cs.CV]
- Xianggang Yu, Mutian Xu, Yidan Zhang, Haolin Liu, Chongjie Ye, Yushuang Wu, Zizheng Yan, Tianyou Liang, Guanying Chen, Shuguang Cui, and Xiaoguang Han. 2023. MVImgNet: A Large-scale Dataset of Multi-view Images. In *CVPR*.
- Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. 2023. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 3836–3847.

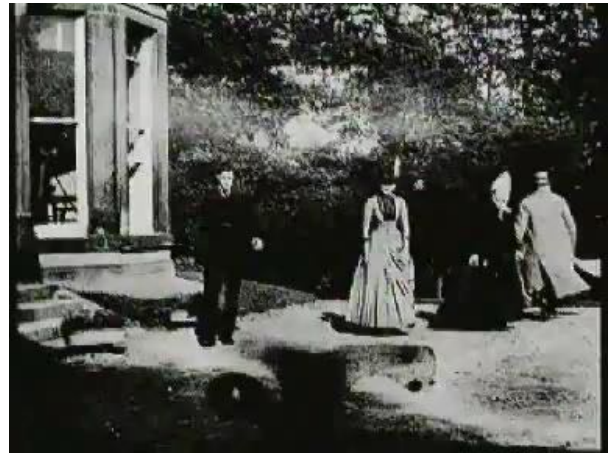
The Birth of Videos

The Horse in Motion (1878)



The first motion picture ever made
Eadweard Muybridge

Roundhay Garden Scene (1888)



The first film with 20 frames
Louis Le Prince

<https://www.thevintagenews.com/2016/06/27/46591-2/> [https://headsup.scout24.com/what-was-the-first-movie-ever-made/#:~:text=Roundhay%20Garden%20Scene%20\(1888\),it%20is%20technically%20a%20movie](https://headsup.scout24.com/what-was-the-first-movie-ever-made/#:~:text=Roundhay%20Garden%20Scene%20(1888),it%20is%20technically%20a%20movie)

Video and Its Origins in Magic

The Vanishing Lady (1897)



Alter Time and Space through Editing
George Melies

Un Homme De Tete (1898)



The Father of Visual Effects
George Melies

Creative Video and Its Origins in Magic



@kassupalen – TikTok 2020



@zachking – TikTok 2019

Rao, Gaba, et al, Organizing ICCV23, ECCV22, ICCV21 Creative Video Editing and Understanding Workshop

Text to Video Generation: SORA



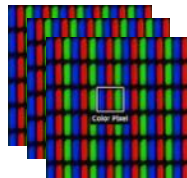
SORA: Prompt: The story of a robot's life in a cyberpunk setting.

How Visual Content is Created?

Visual Content from Pixels

$rgb(w, h) t$

appearance, width, height



Visual Content from a Camera Navigating in the 3D Environment

$$(x, y, z, \alpha, \beta, \gamma, f) t$$

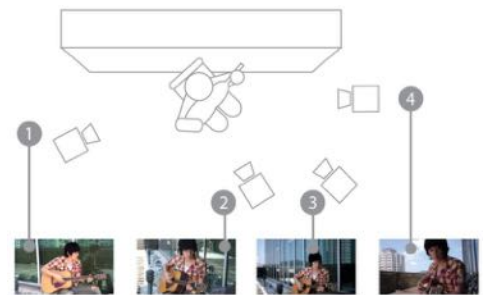
Position, Angle, Focal Length



Visual Content from Multi-view Editing

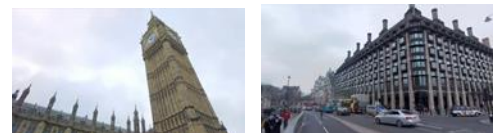
$$i, t$$

camera index, time



$$(\alpha, \beta, f) t$$

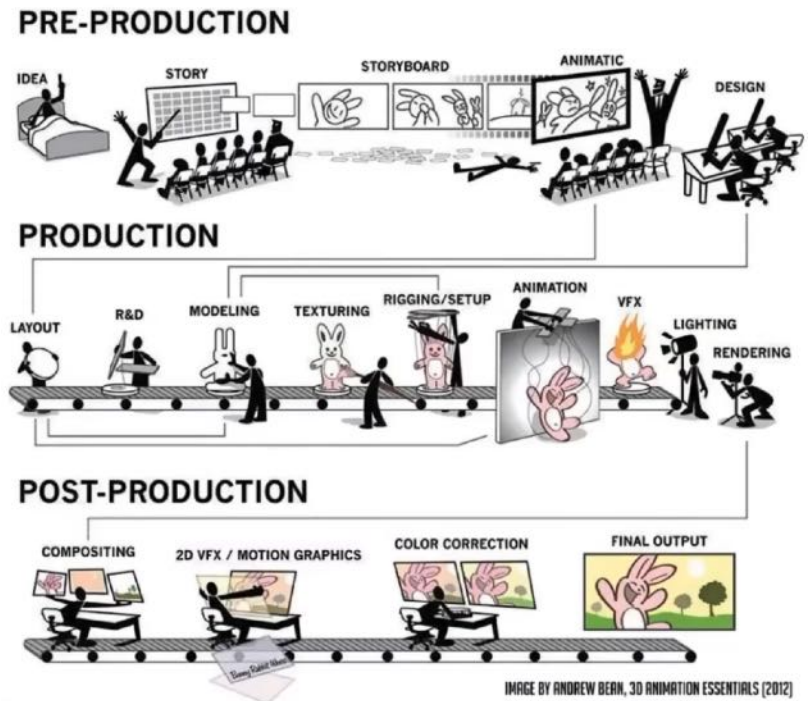
horizontal/vertical angle, focal, time



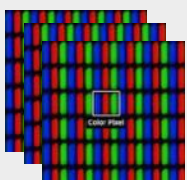
Visual Content from Professional Pipeline SIGGRAPH 2024 DENVER+ 28 JUL — 1 AUG

A Comprehensive workflow that combines

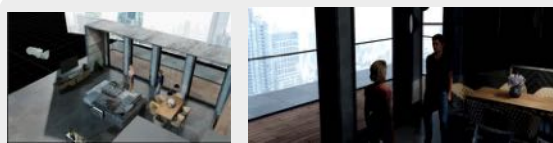
- People's efforts
- Natural language processing
- Computer graphics
- Computer vision
- Animation
- VFX
- Artificial intelligence
- More....



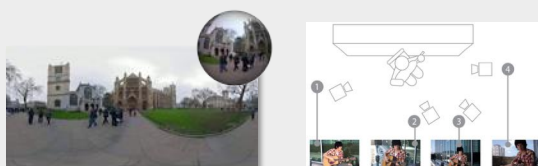
How Visual Content is Created?



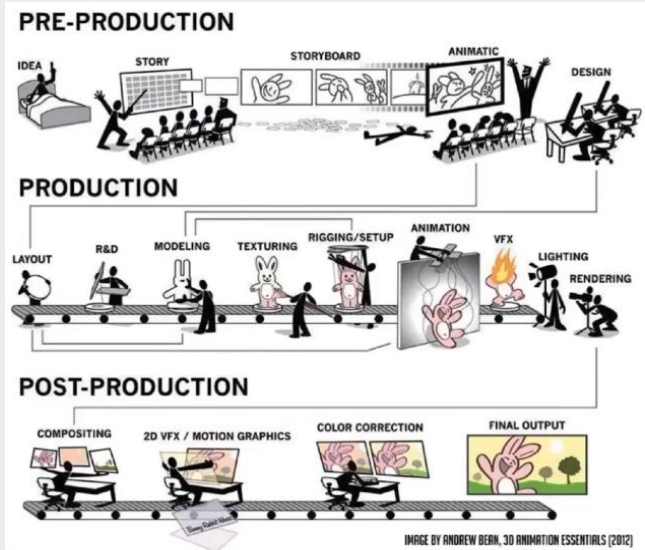
Pixels



Camera in 3D



Editing from Multi view Footage



Professional Video Pipeline

🙄 More e.g., Remixing, Augmenting, Editing

Introduction to Generative Models

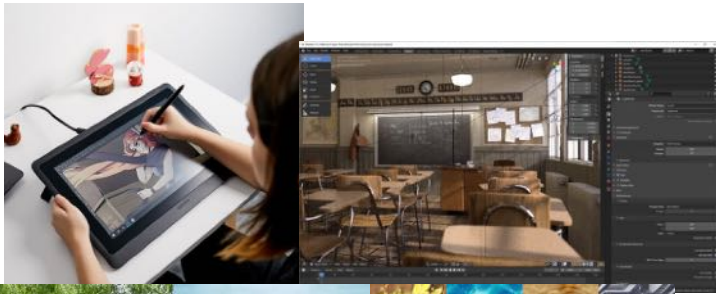
1

Agenda

- Introduction to Diffusion
- Conditional generation and guidance
- Implementation Architectures

2

To make a beautiful synthetic image...



Past



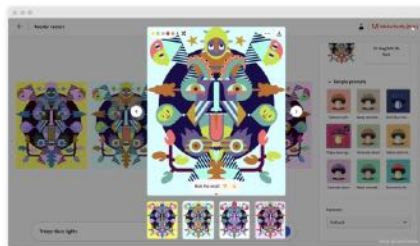
Now!

A cute corgi sitting on a beach sipping on a glass of lemonade. 4K, photorealistic.



3

Generative AI Applications



Art & Design



content Generation



Entertainment

4

The Landscape of Generative AI



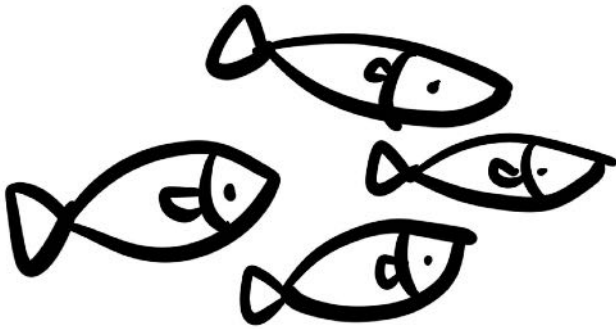
How do we generate new data?



our goal:

Generate fish that
looks and behaves like it
belongs to this river

How do we generate new data?

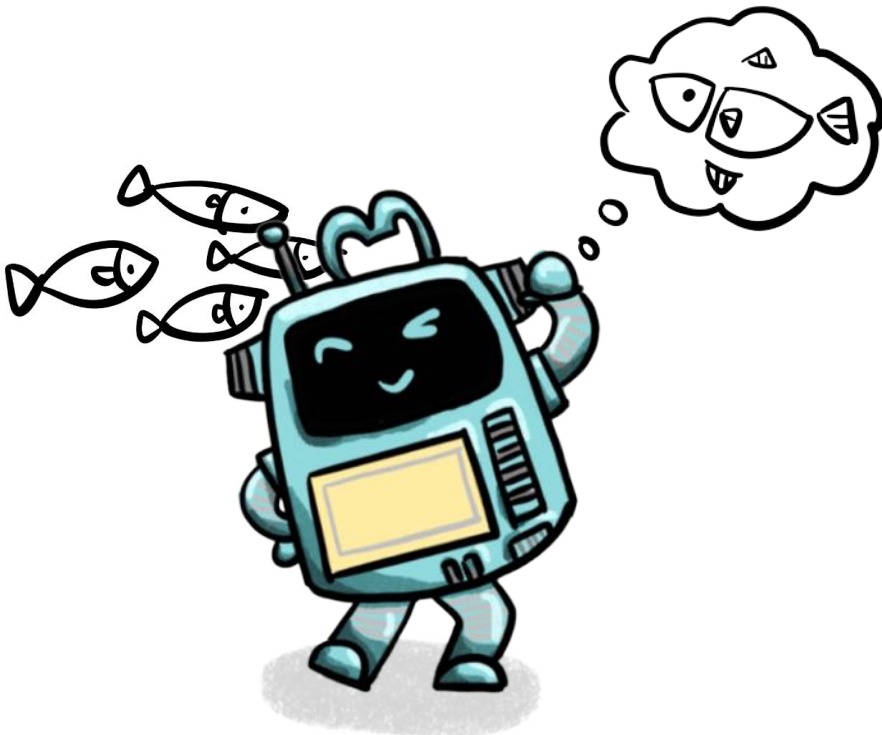


Step 1:

capture ^{fish from the river}
a lot of

7

How do we generate new data?



Step 2:

Train a neural network
to learn the fish
distribution by analyzing
the captured fish!

8

How do we generate new data?

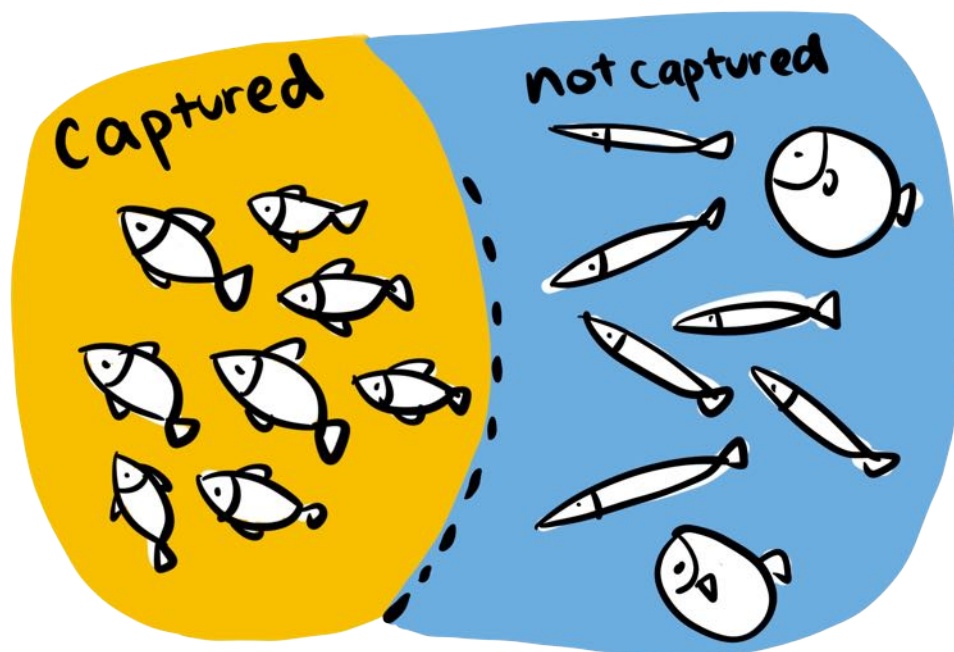
Step 3:



- Use the trained neural network to generate new fish.
- Ensure the generated fish is good: have characteristics learned from the dataset.

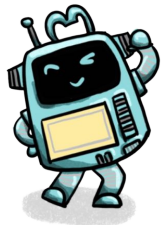
9

Are we modeling the actual distribution?

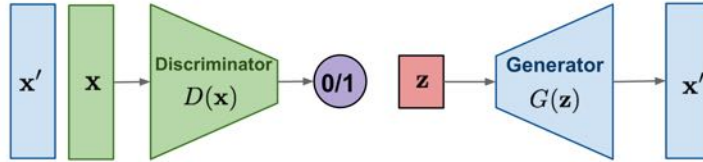


10

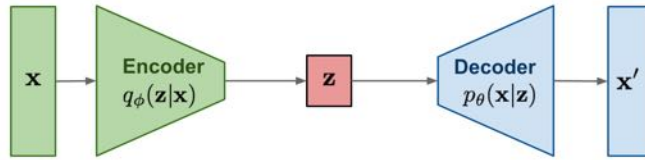
Models



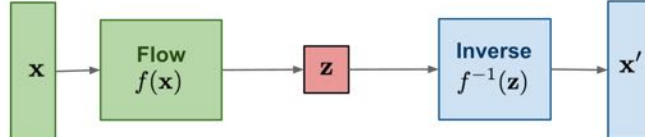
GAN: Adversarial training



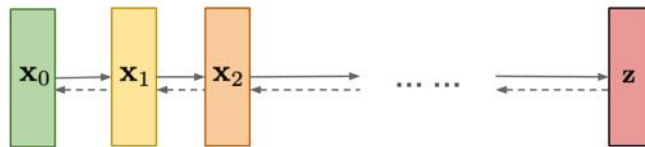
VAE: maximize variational lower bound



Flow-based models: Invertible transform of distributions

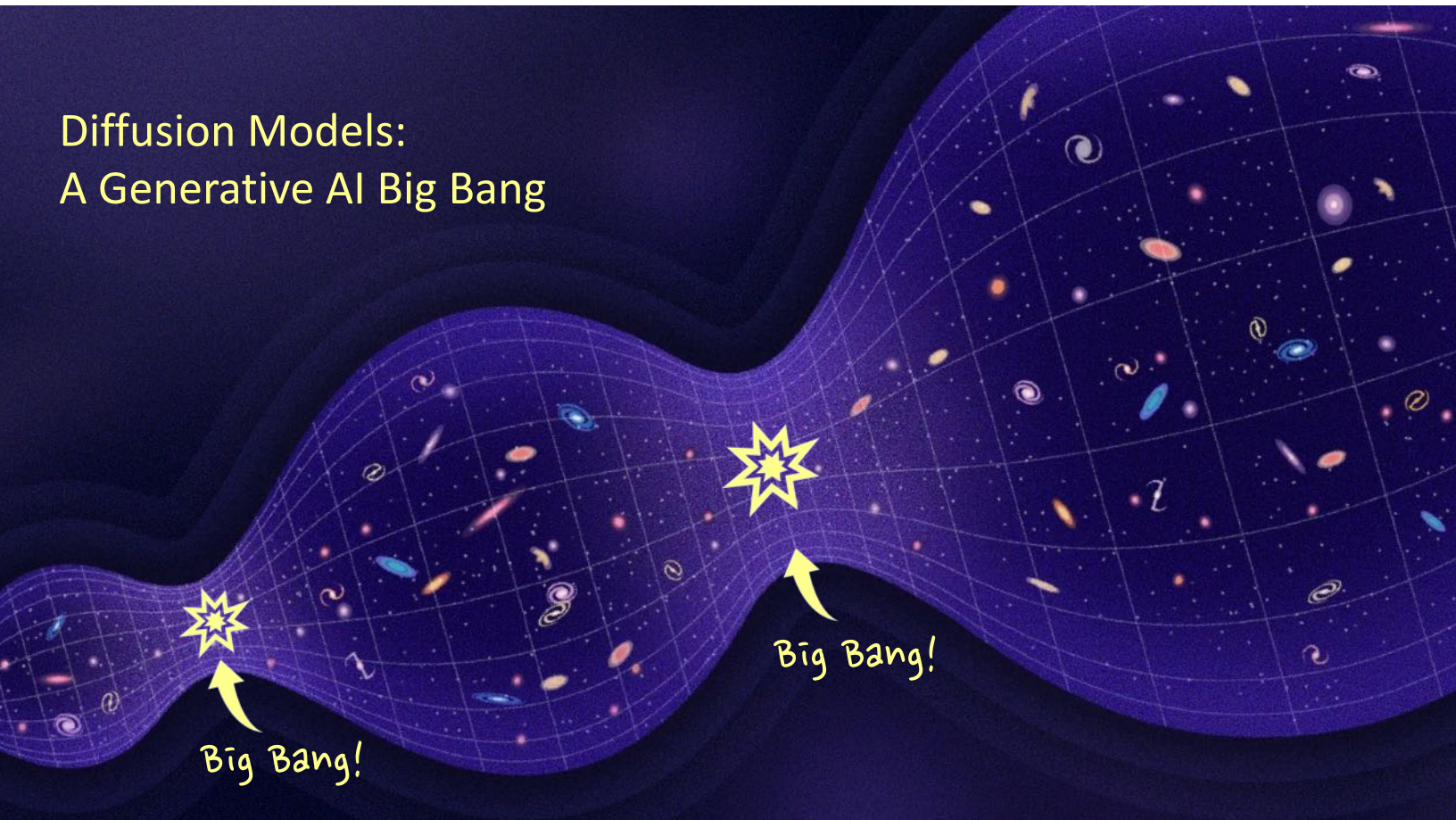


Diffusion models: Gradually add Gaussian noise and then reverse



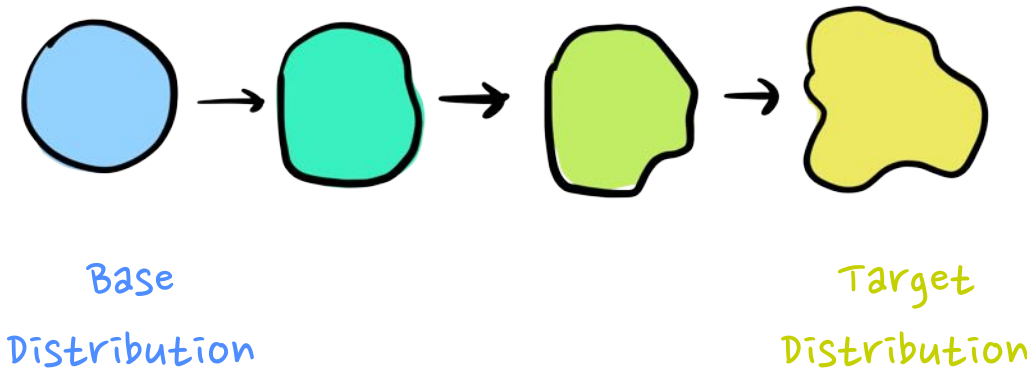
<https://lilianweng.github.io/posts/2021-07-11-diffusion-models/>

Diffusion Models: A Generative AI Big Bang



Diffusion models

- Main idea: iteratively convert a base distribution to the target distribution via Markov chain



Denoising Diffusion Probabilistic Models

Jonathan Ho
UC Berkeley
jonathanho@berkeley.edu

Ajay Jain
UC Berkeley
ajayj@berkeley.edu


Pieter Abbeel
UC Berkeley
pabbeel@cs.berkeley.edu

Abstract

We present high-quality image synthesis results using diffusion probabilistic models, a class of latent variable models inspired by considerations from nonequilibrium thermodynamics. Our best results are obtained by training on a weighted variational bound designed according to a novel connection between diffusion probabilistic models and denoising score matching with Langevin dynamics, and our models naturally admit a progressive lossy decompression scheme that can be interpreted as a generalization of autoregressive decoding. On the unconditional CIFAR-10 dataset, we obtain an Inception score of 9.46 and a state-of-the-art FID score of 3.17. On 256x256 LSUN, we obtain sample quality similar to ProgressiveGAN. Our implementation is available at <https://github.com/hojonathanho/diffusion>.

1 Introduction

Deep generative models of all kinds have recently exhibited high quality samples in a wide variety of data modalities. Generative adversarial networks (GANs), autoregressive models, flows, and variational autoencoders (VAEs) have synthesized striking image and audio samples [14, 27, 3, 38, 38, 25, 10, 32, 44, 57, 26, 33, 45], and there have been remarkable advances in energy-based modeling and score matching that have produced images comparable to those of GANs [11, 55].



Basics of diffusion models

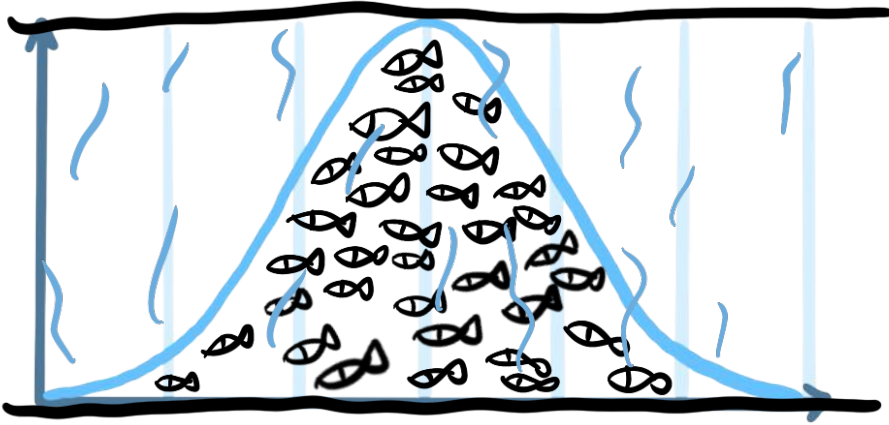
Forward
Diffusion
Process

Reverse
Diffusion
Process

Training
&
Sampling

Refresh on distributions

Gaussian Distribution



Defined by:

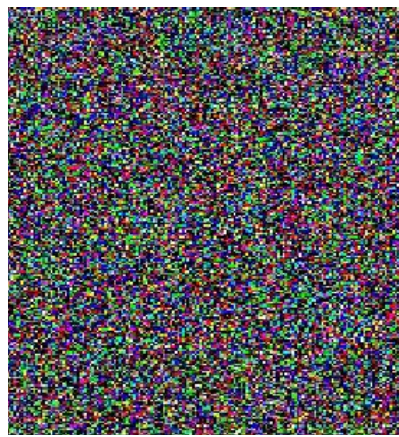
- Mean
- Std. deviation
 $= \sqrt{\text{variance}}$

15

Gaussian noise



+



=



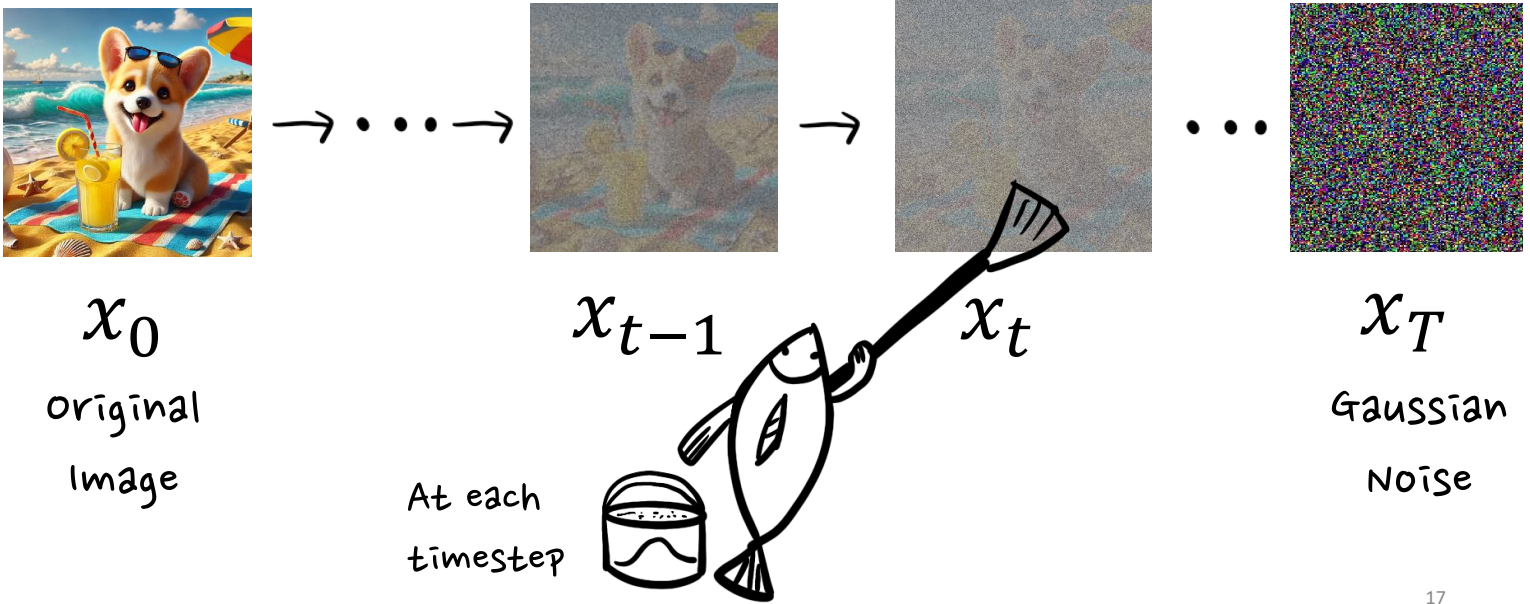
original
Image

Gaussian
Noise

Noised
Image

16

Forward diffusion process



Forward diffusion process

$$q(x_t | x_0) = \mathcal{N}(x_t; \underbrace{\sqrt{1 - \beta_t}}_{\text{mean}}, \underbrace{\beta_t I}_{\text{var.}})$$

variance schedule

Apply forward process one by one?

Too much to store in my memory! or disk!



$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1})$$

19

Reparameterization trick

Forward process $q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t\mathbf{I})$

Rewriting Def. $x_t = \sqrt{1 - \beta_t}x_{t-1} + \sqrt{\beta_t}\epsilon_t, \quad \epsilon_t \sim \mathcal{N}(0, \mathbf{I})$

Define variables $\alpha_t = 1 - \beta_t \quad \bar{\alpha}_t = \prod_{i=1}^t \alpha_i$

Gaussian Recap We can merge gaussians with different variances. $\mathcal{N}(0, \sigma_1^2\mathbf{I}), \mathcal{N}(0, \sigma_2^2\mathbf{I}) \rightarrow \mathcal{N}(0, (\sigma_1^2 + \sigma_2^2)\mathbf{I})$

Plug in

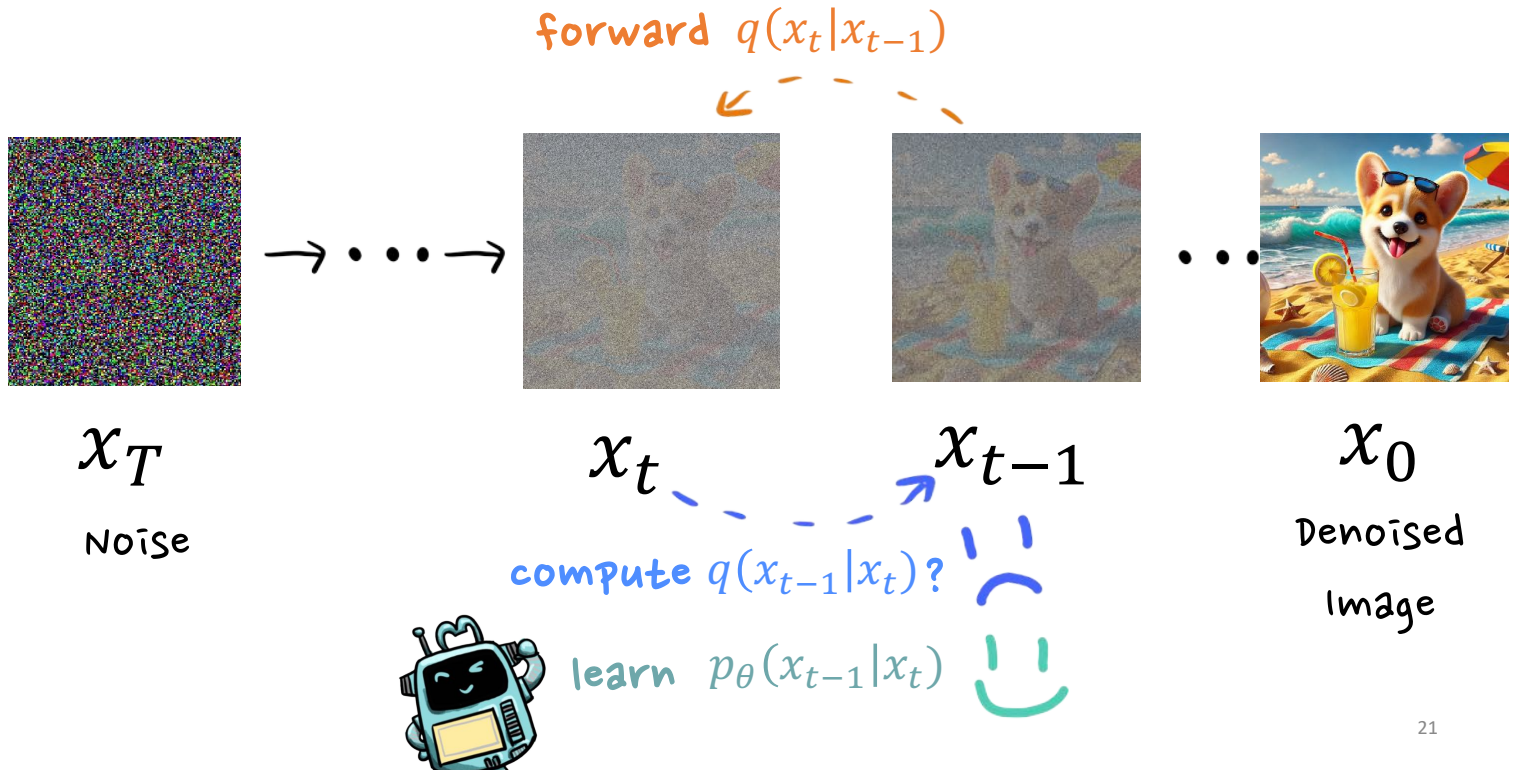
definitions!

$$\begin{aligned} \mathbf{x}_t &= \sqrt{\alpha_t}\mathbf{x}_{t-1} + \sqrt{1 - \alpha_t}\epsilon_{t-1} \\ &= \sqrt{\alpha_t\alpha_{t-1}}\mathbf{x}_{t-2} + \sqrt{1 - \alpha_t\alpha_{t-1}}\bar{\epsilon}_{t-2} \\ &= \dots \\ &= \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon \end{aligned}$$

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$$

20

Reverse diffusion process



21

Training diffusion models

Model: $\epsilon_\theta \left(\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon \right)$

Loss: $MSE \left[\epsilon - \epsilon_\theta \left(\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon \right) \right]$

Algorithm 1 Training

- 1: repeat
 - 2: $x_0 \sim q(x_0)$
 - 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
 - 4: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 5: Take gradient descent step on $\nabla_\theta \left\| \epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t) \right\|^2$
 - 6: until converged
-

Loss calculation

DDPM Training Loop

22

Algorithm 2 Sampling

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

DDPM Sampling Loop

23

Forward
Diffusion
Process

Reverse
Diffusion
Process

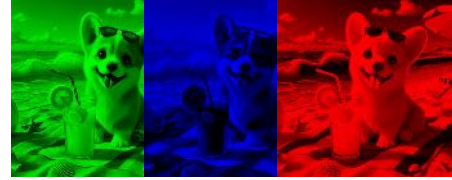
Training
&
Sampling

24

Pixels are expensive!



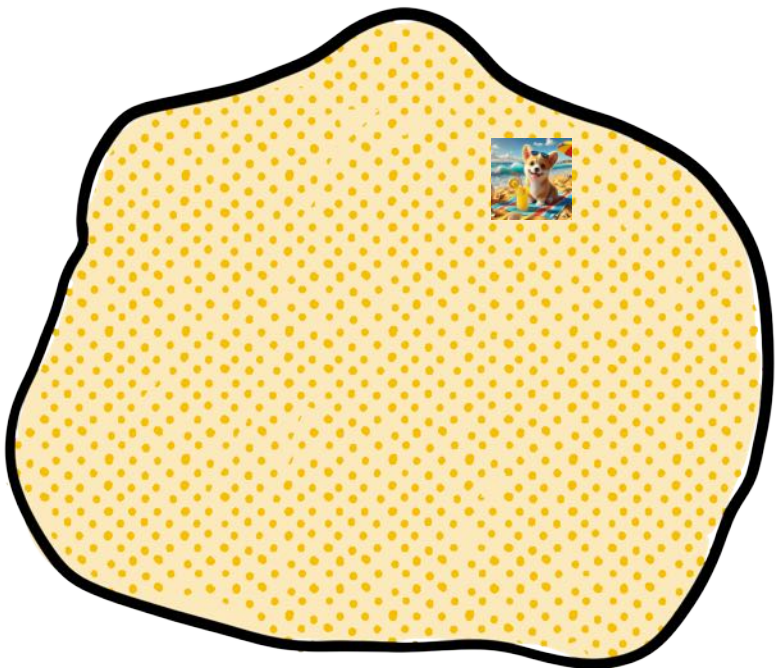
High Quality Image



Need to store a lot of data!

25

Latent diffusion



Pixel Space

VAE
→



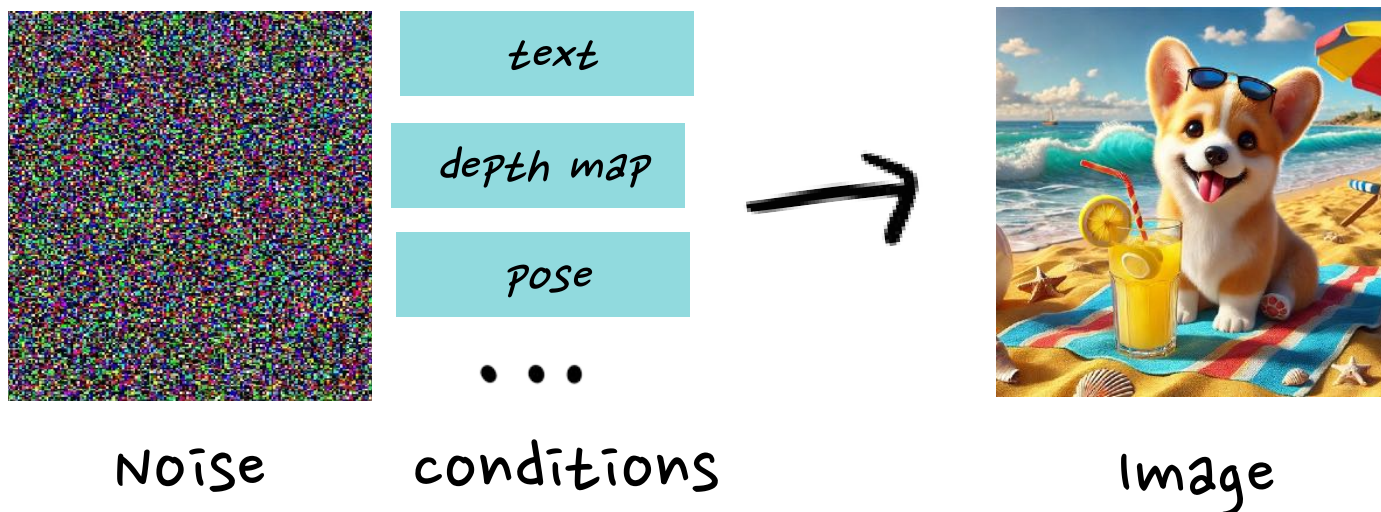
Latent Space

26

Conditional generation and guidance

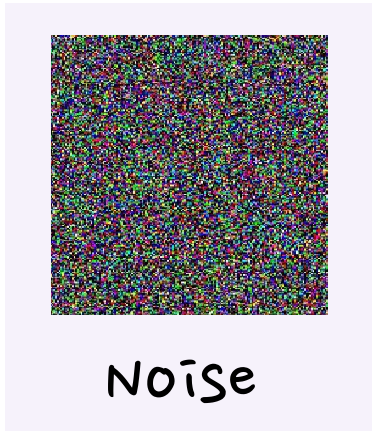
27

So far, we've seen noise to image



28

Score function

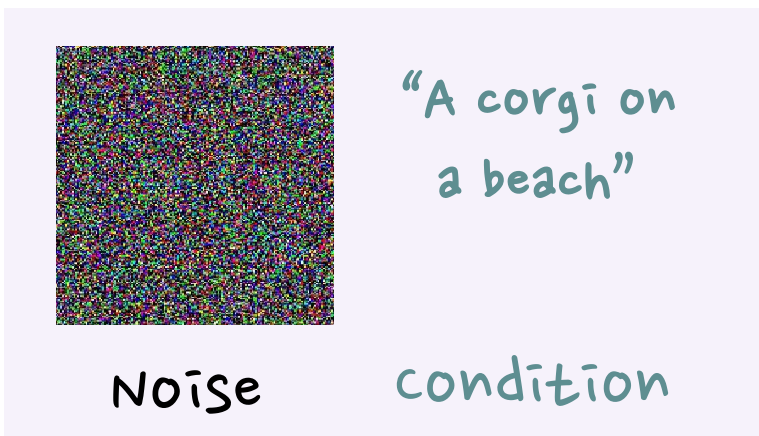


Image

$$p(x) \approx \nabla_x \log p(x)$$

29

Conditional diffusion models



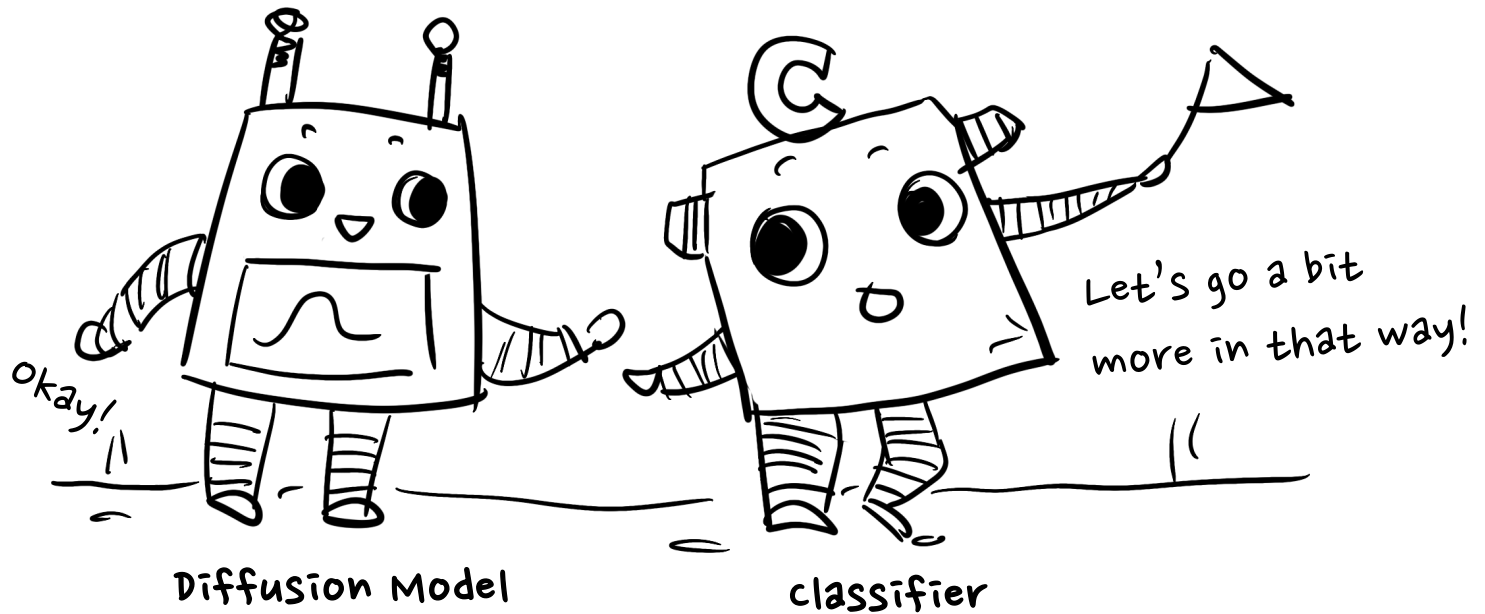
Image

$$p(x|y) \approx \nabla_x \log p(x|y) = \nabla_x \log p(y|x) + \nabla_x \log p(x)$$

conditioning term

unconditional
score function³⁰

Classifier guidance



31

Classifier guidance

Guidance Strength

$$\nabla_x \log p(x|y) = w \nabla_x \log p(y|x) + \nabla_x \log p(x)$$



$w = 1.0$



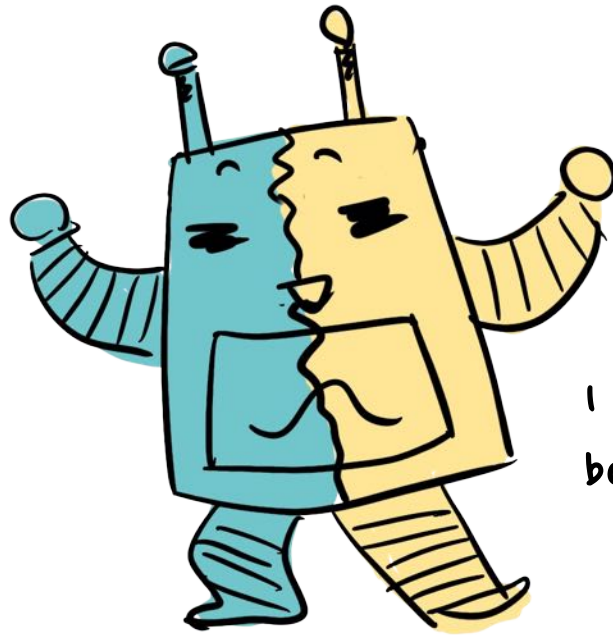
$w = 10.0$

32

Classifier-Free Guidance (CFG)

conditioning Dropout:

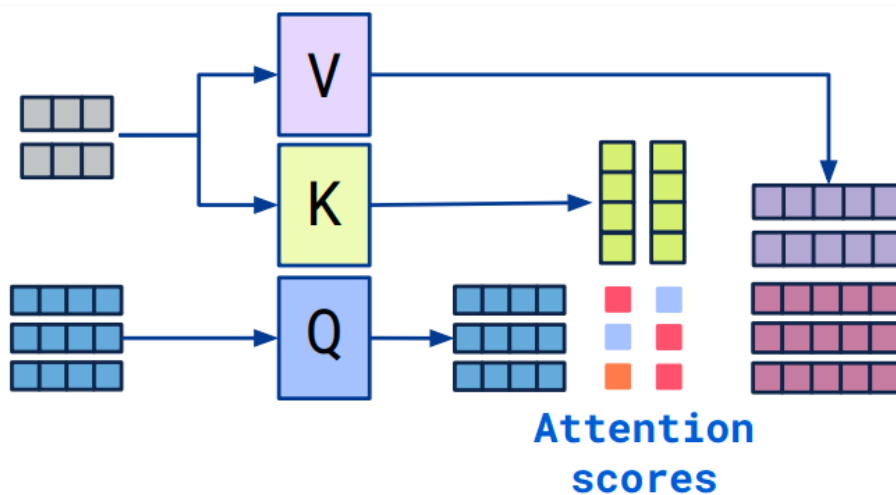
10 to 20 percentage of the time, the conditioning information is removed.



I can do both!

33

How do we feed conditioning signals?



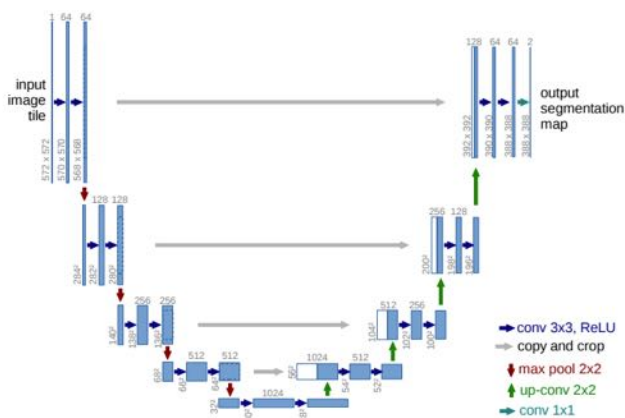
cross-attention layers!

34

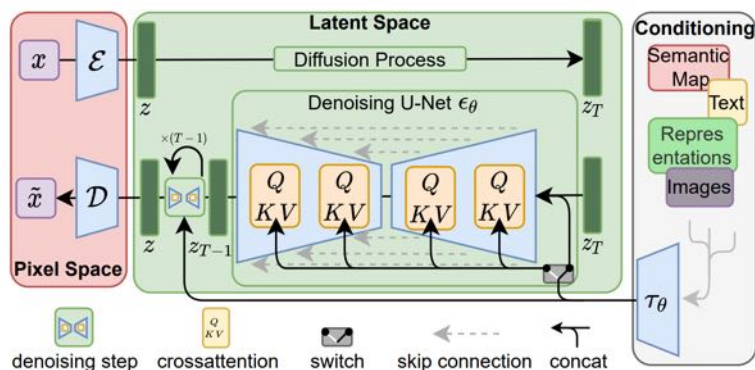
Implementation Architectures

35

U-Net Architecture



U-Net architecture



U-Net based
diffusion architecture

36

U-Net Architecture



Imagen.



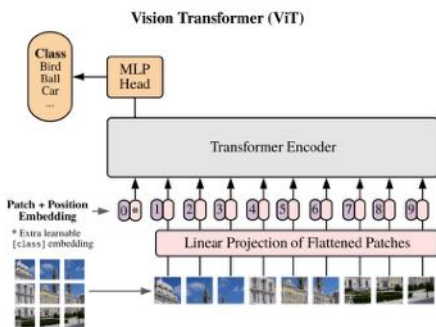
Stable Diffusion



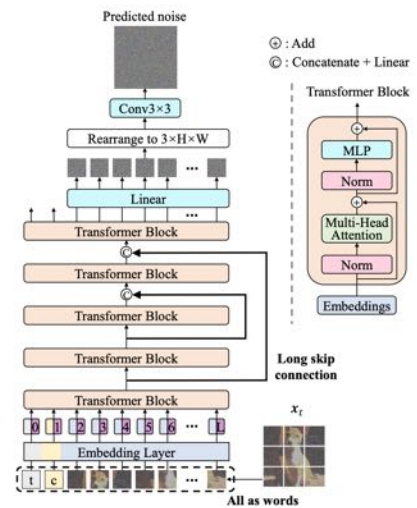
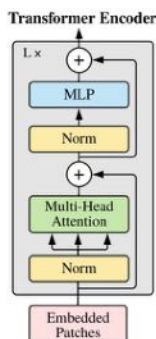
eDiff-1

Saharia et al., "Photorealistic text-to-image diffusion models with deep language understanding", NeurIPS 2022
 Rombach et al., "High-Resolution Image Synthesis with Latent Diffusion Models", CVPR 2022
 Balaji et al., "ediffi: Text-to-image diffusion models with an ensemble of expert denoisers", arXiv 2022

Transformer Architecture



vision transformer.



Transformer based diffusion model

Dosovitskiy et al., "An image is worth 16x16 words: Transformers for image recognition at scale", ICLR 2021
 Bao et al., "All are Worth Words: a ViT Backbone for Score-based Diffusion Models", arXiv 2022

Transformer Architecture



Scalable Diffusion Models
with Transformers



One Transformer Fits All
Distributions in Multi-Modal
Diffusion at Scale



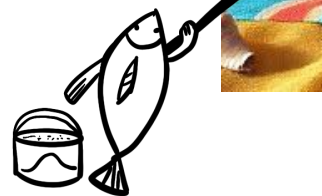
Simple diffusion:
End-to-end diffusion for
high resolution images

Peebles and Xie, "[Scalable Diffusion Models with Transformers](#)", arXiv 2022
Bao et al., "[One Transformer Fits All Distributions in Multi-Modal Diffusion at Scale](#)", arXiv 2023
Hoogeboom et al., "[simple diffusion: End-to-end diffusion for high resolution images](#)", arXiv 2023

39

Summary

- Basics of diffusion models
 - Forward & reverse diffusion process
 - Sampling and training
 - Latent diffusion
- conditional generation
 - Classifier guidance
 - Classifier-free guidance (CFG)
 - Adding condition using cross-attention
- Implementation architectures
 - U-net
 - Vision Transformers



Generated with DALL-E

40

Start from Text-to-Image Large Models

DALL-E 3



Stable Diffusion

Stable Diffusion XL



1

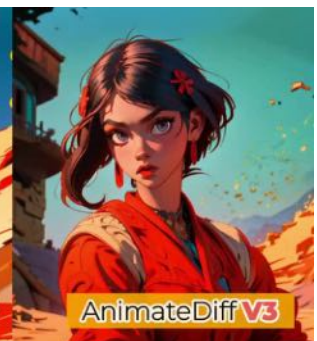
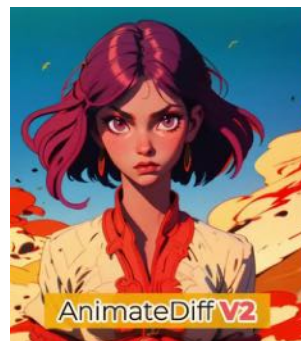
Start from Text-to-Video Large Models



Gen-2: The Next Step Forward for Generative AI

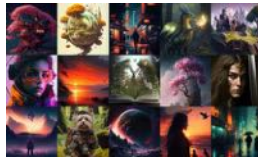
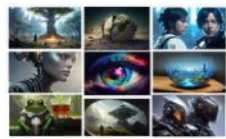
Google Research

LUMIERE

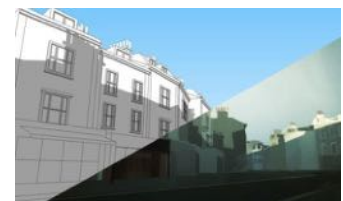
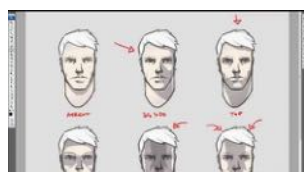
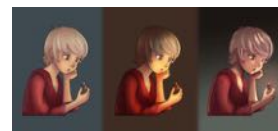
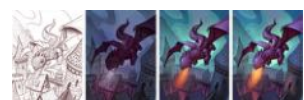
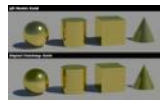
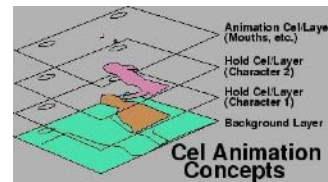


2

Generating High-quality Images ...



But visual creation is more than just generating beautiful images ...



More Control Other Than Texts?

stanford memorial church with neon signage in the style of bladerunner



Iteration 1

stanford memorial church and main quad with palm trees in the style of bladerunner



Iteration 3

nighttime rain stanford memorial church and main quad with palm trees, **night market food stalls and neon signs** in the style of bladerunner



Iteration 8

nighttime rain stanford memorial church and main quad with palm trees, night market food stalls and neon signs **like downtown tokyo**

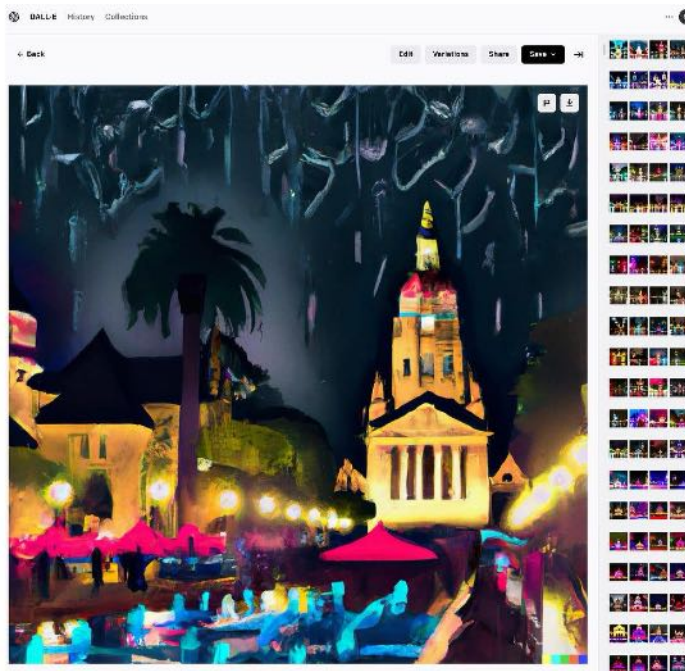


Iteration 17

More Control Other Than Texts?

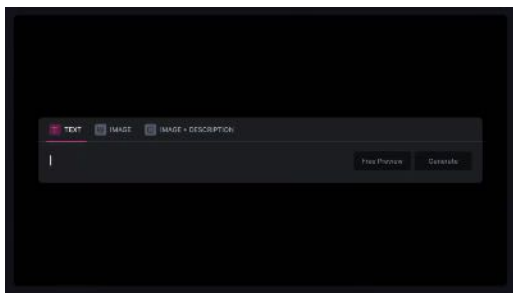
nighttime rain stanford memorial church and main quad with palm trees, night market **japadog** food stalls and neon signs, **neo tokyo bladerunner style film still illustration**

Iteration 21

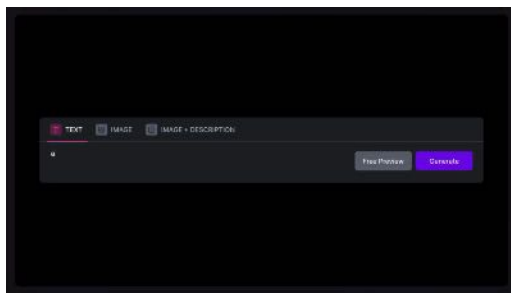


Lots of trial-and-error!

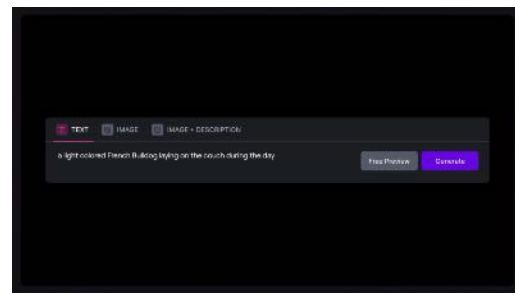
Text Control is Limited in Creation



Iteration 1



Iteration 5



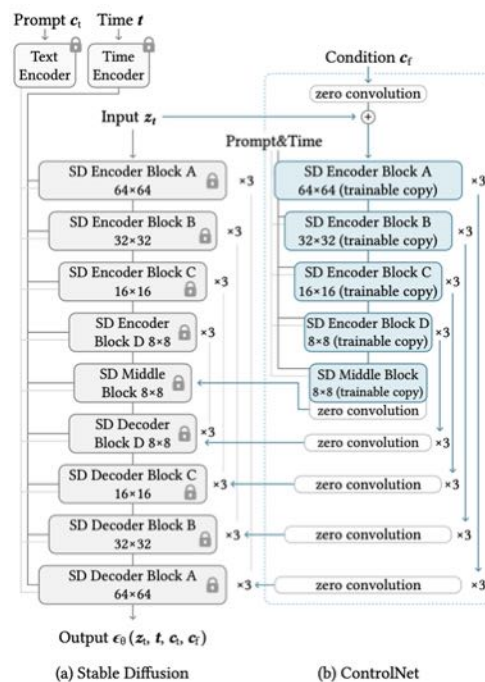
Iteration 17

Text does not match user's **mental representation**, which leads to lots of **trial-and-error!**

7

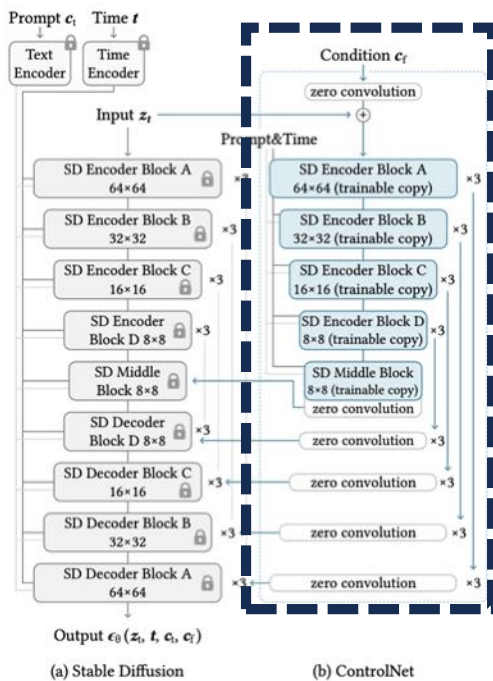
<https://magrawala.substack.com/p/unpredictable-black-boxes-are-terrible>

ControlNet



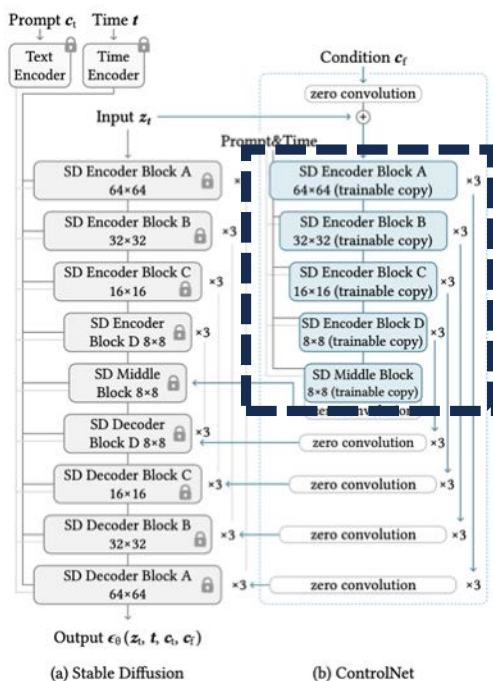
8

Architecture of ControlNet



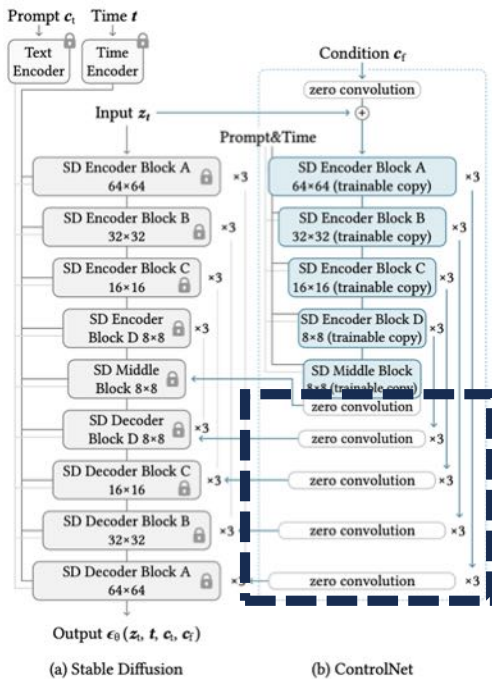
- Using **external model** to process control signals.
- **Re-using pretrained weights** as the backbone of control model.
- Connecting with **zero-initialized layers** to reduce initial noise.

Architecture of ControlNet



- Using **external model** to process control signals.
- **Re-using pretrained weights** as the backbone of control model.
- Connecting with **zero-initialized layers** to reduce initial noise.

Architecture of ControlNet



- Using **external model** to process control signals.
- **Re-using pretrained weights** as the backbone of control model.
- Connecting with **zero-initialized layers** to reduce initial noise.

External model to process control signals

Finetuning
diffusion model
weights

v.s.

Training
external control
model

- Compossible control (multiple conditions)
- Minimal influence to the base model (the base model can be changed)
- Reduced overfitting risk (training with small dataset becomes easier)

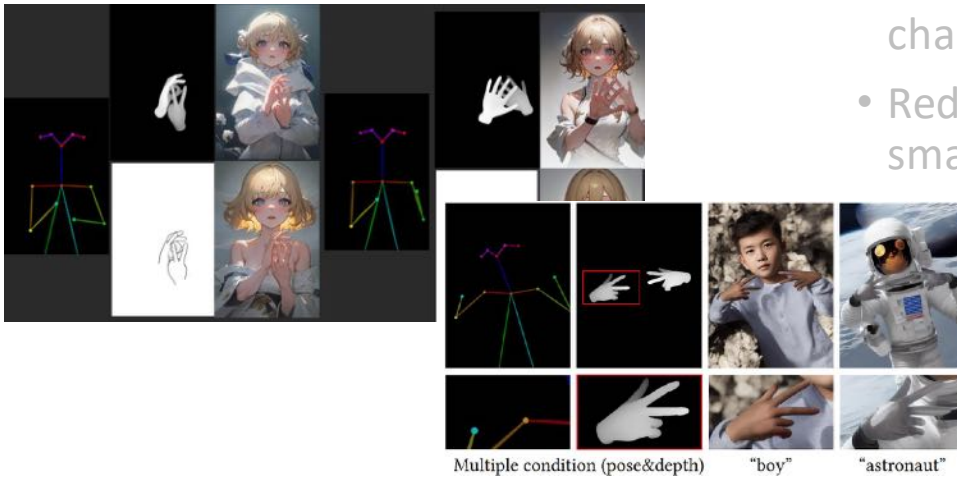
External model to process control signals SIGGRAPH 2024 DENVER+ 28 JUL - 1 AUG

Finetuning
diffusion model
weights

v.s.

Training
external control
model

- Compossible control (multiple conditions)
- Minimal influence to the base model (the base model can be changed)
- Reduce overfitting (training with small dataset becomes easier)



13

External model to process control signals SIGGRAPH 2024 DENVER+ 28 JUL - 1 AUG

Finetuning
diffusion model
weights

v.s.

Training
external control
model

- Compossible control (multiple conditions)
- Minimal influence to the base model (the base model can be changed)
- Reduce overfitting (training with small dataset becomes easier)



"house"

SD 1.5

Comic Diffusion

Protogen 3.4

14

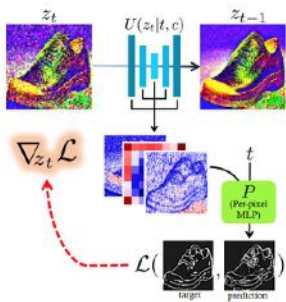


without ControlNet
(using Stability's "official" method to add the channels to input layer, same as their depth-to-image structure)

SD + ControlNet

- Compossible control (multiple conditions)
- Minimal influence to the base model (the base model can be changed)
- Reduce overfitting (training with small dataset becomes easier)

Reusing pretrained backbone



Some insights from some previous works ...

In the paper "Sketch-Guided Text-to-Image Diffusion Models" (from 2022 November), Voynov *et.al.* discussed that one of the major challenge of "sketches" guided diffusion is **the difficult alignment of complex scenes with mixed and ambiguous semantics.**



Figure 14. **Failure cases.** The quality of the results may drop for different initialization, and on complex scenes with mixed and ambiguous semantics.

This motivates us to find a **stronger backbone** to solve the semantic alignment and understanding problem ...



By the way, this is the result from ControlNet 1.1.

Reusing pretrained backbone

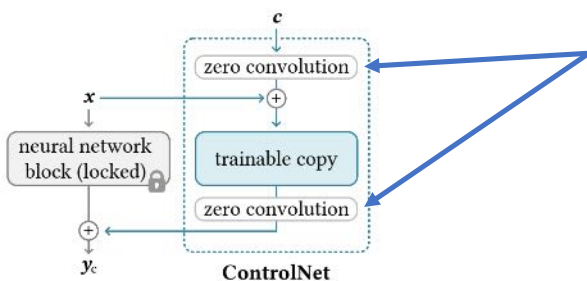
The ability to “guess” contexts without accurate prompts ...

The diagram shows two architectures: (a) Stable Diffusion and (b) ControlNet. Both share a common backbone of SD Encoder and Decoder blocks. In ControlNet, the encoder and decoder blocks are replaced by trainable copies, and zero-convolution layers are added at various stages to receive conditioning information. Below the diagram, an 'Input' sketch of a vase is shown alongside a 'high-quality and extremely detailed image' of a golden vase. To the right, a screenshot of the 'Control Stable Diffusion with Scribble Maps' interface shows a grid of generated images, including a backpack, a basket of fruit, and a red car, demonstrating the model's ability to generate detailed content from simple prompts.

17

Using zero-initialized layers

The ability to “guess” contexts without accurate prompts ...



Zero-initialized connection layers

- Reduce initial harmful noise
- Protect the trainable copy

$$y_c = \mathcal{F}(x; \Theta) + \mathcal{Z}(\mathcal{F}(x + \mathcal{Z}(c; \Theta_{z1}); \Theta_c); \Theta_{z2}),$$

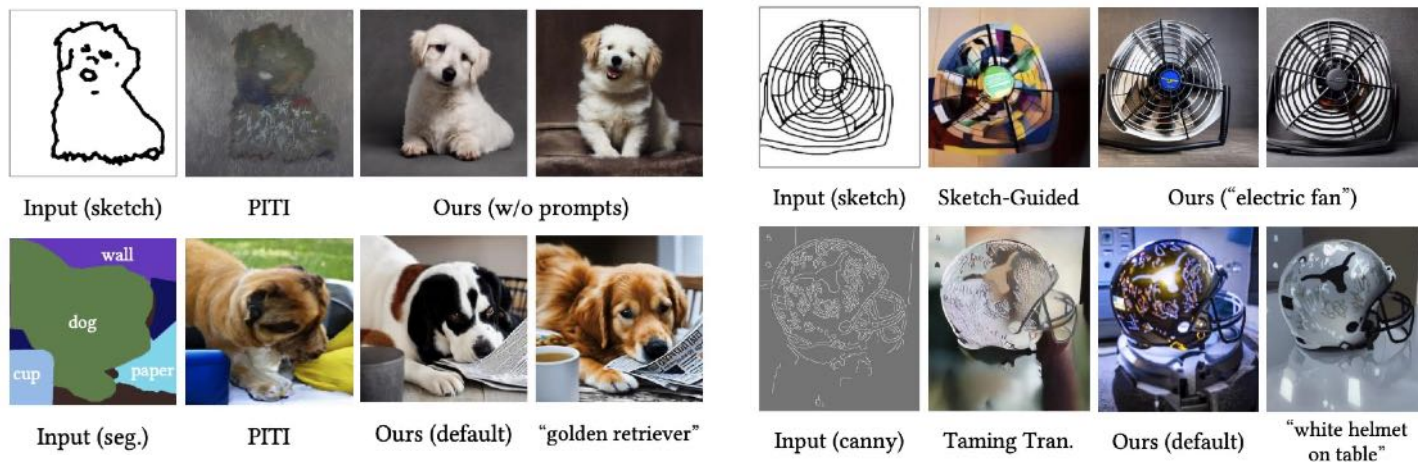
In the first training step, $y_c = y$.

18



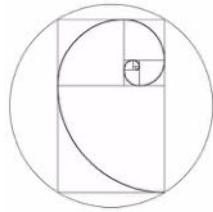
All experiments are conducted with Stable Diffusion 1.5

Comparisons



Unleash Human Creativity

 **Reddit**
<https://www.reddit.com/StableDiffusion/comments/>
SDBattle: Week 5 - ControlNet Cross Walk Challenge! Use ...
Mar 20, 2023 — Welcome back to the weekly Stable Diffusion Battle Challenge! Excited to see what you all make! Join us for more battles over at r/SDBattles. If ...



 **Reddit**
<https://www.reddit.com/StableDiffusion/comments/>
SDBattle: Week 3 - ControlNet Fibonacci Challenge! Use ...
Mar 6, 2023 — I think I'll pass this battle and just see the results and if they disprove my statement. If I would go into photoshop and add some gradients or ...



 **Reddit**
<https://www.reddit.com/StableDiffusion/comments/>
SDBattle: Week 8 - ControlNet The Thinker Challenge! Use ...
Apr 11, 2023 — In this week's SD Battle I hand made a pose based on The Thinker statue since ControlNet was having a hard time generating one.

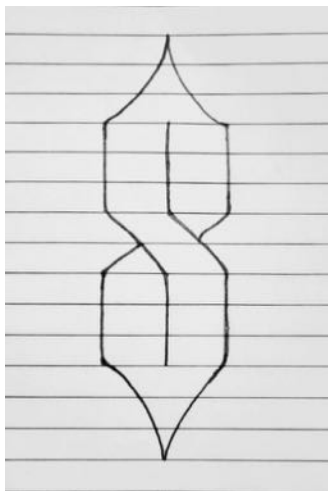


 **Reddit**
<https://www.reddit.com/StableDiffusion/comments/>
SDBattle: Week 6 - ControlNet Dog Paw Challenge! Use ...
Mar 27, 2023 — Join us for more thrilling battles over at r/SDBattles, where the action never stops and the stakes are always high! Don't miss this chance to ...

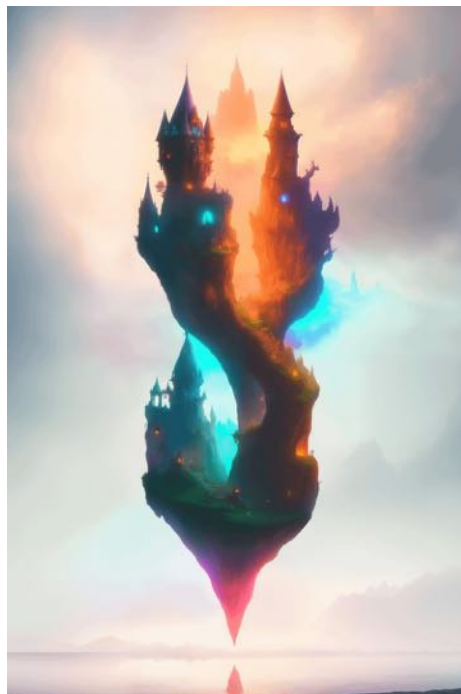
Input

Results Conditioned on the Canny Map from Input

Unleash Human Creativity



Input

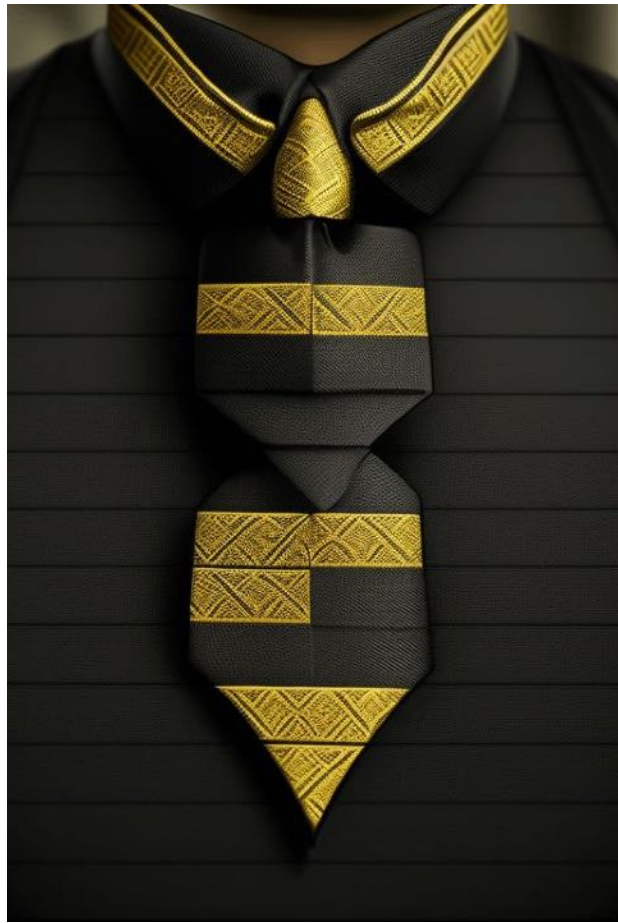




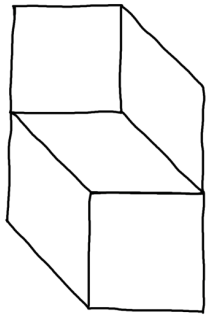
 **SIGGRAPH 2024**
DENVER+ 28 JUL — 1 AUG



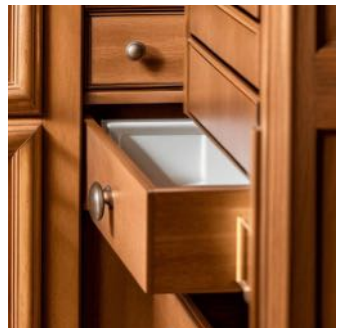
 **SIGGRAPH 2024**
DENVER+ 28 JUL — 1 AUG



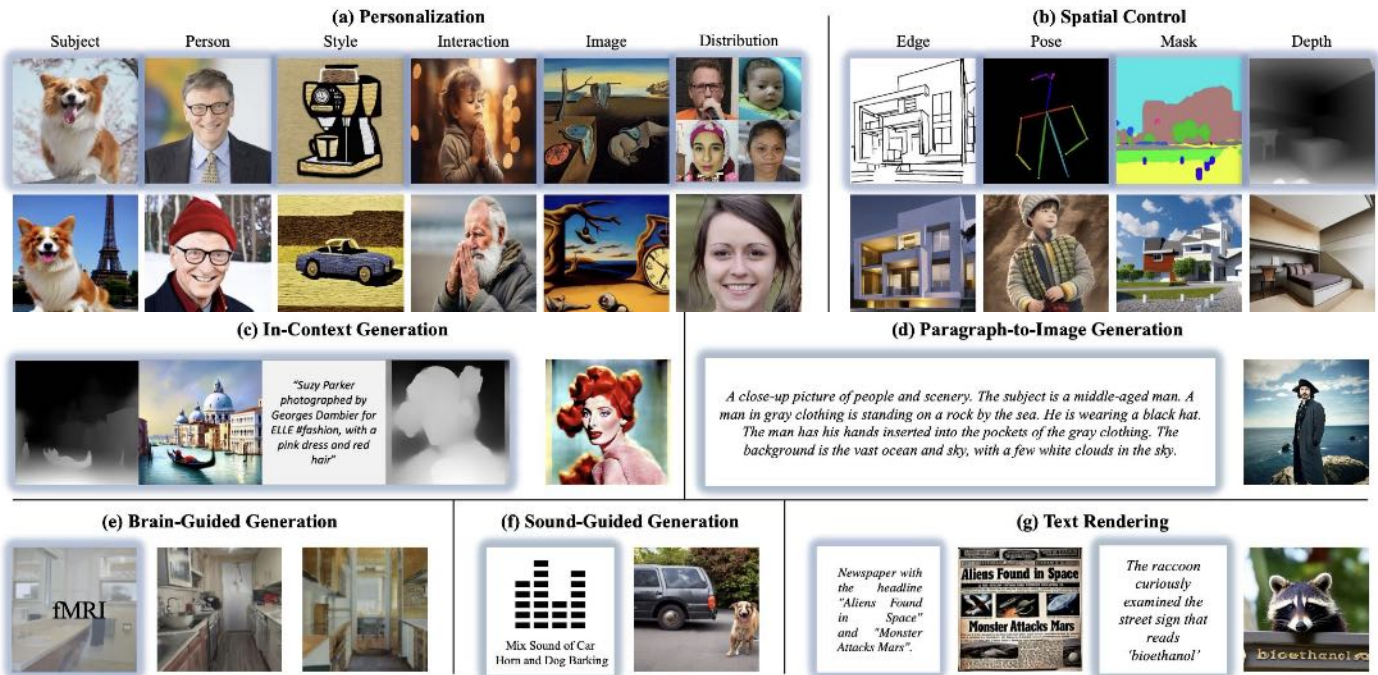
Unleash Human Creativity



Input



Controllable Generation in General

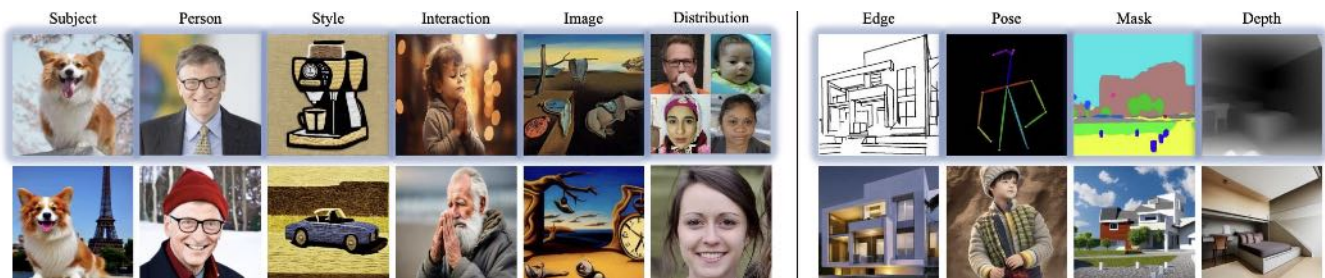


27

Cao, Pu, et al. "Controllable generation with text-to-image diffusion models: A survey." *arXiv preprint arXiv:2403.04279* (2024).

Take Away

- Text control is limited
- Better control leads to higher quality



28

Extend Image Diffusion Models for Videos

1. Data Format



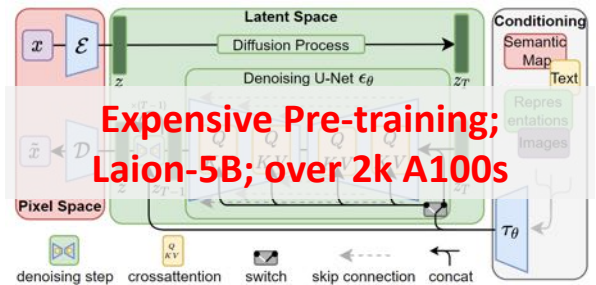
image



video

2. Dataset availability

LAION (image, 5B) vs. WebVid (video, 10M)



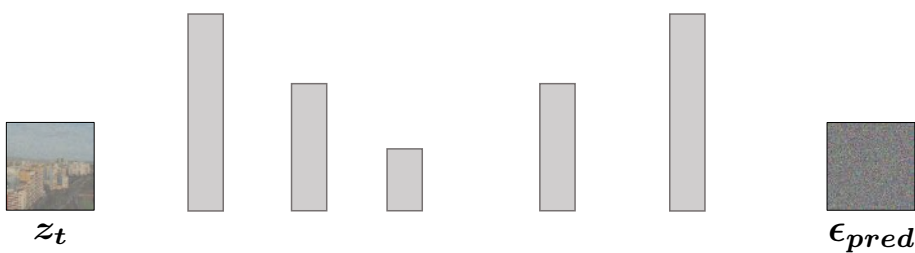
Extend Image Diffusion Models for Videos

Goal: leveraging powerful T2I prior knowledge

Reasons:

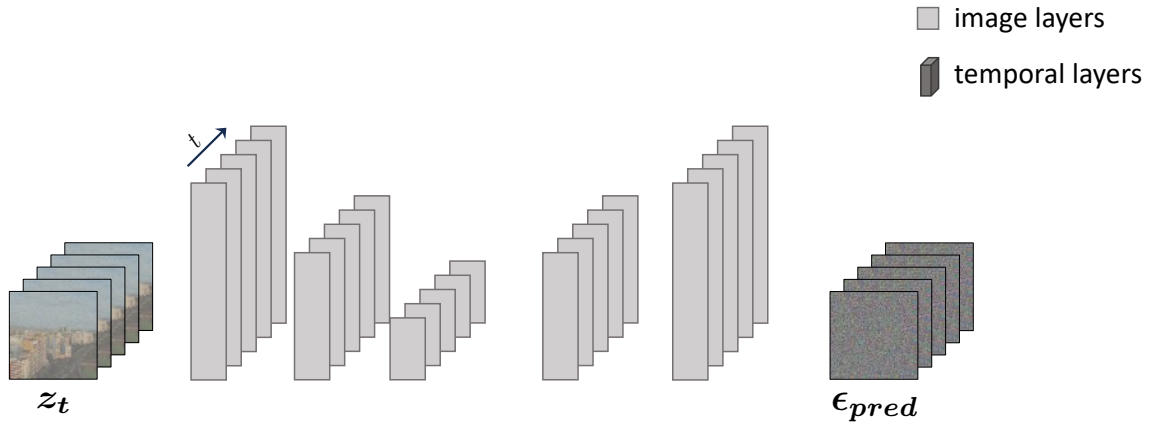
- (1) better initialization than from scratch
- (2) dataset scale (5B vs. 10M)

- image layers
- temporal layers



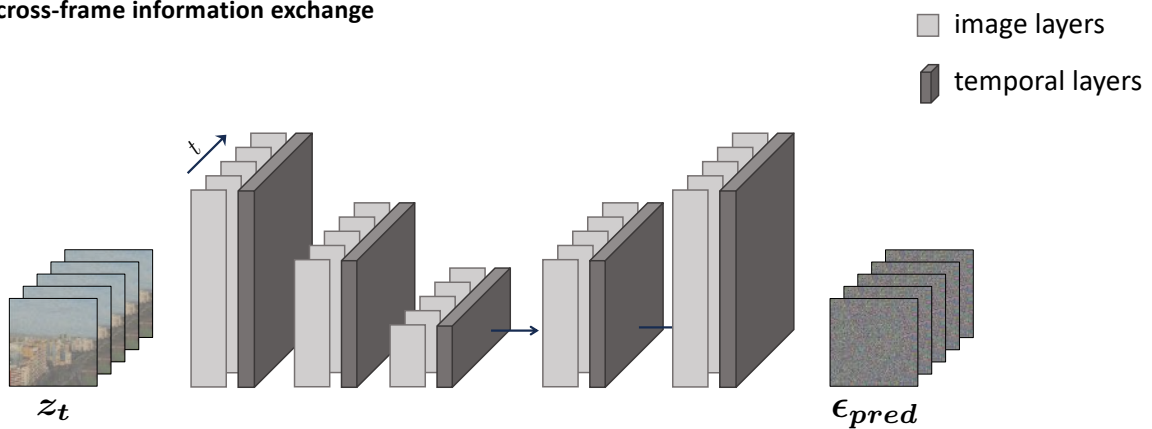
Extend Image Diffusion Models for Videos

1. Repeat the image generator along the time axis (e.g., 16/24 frames)



Extend Image Diffusion Models for Videos

1. Repeat the image generator along the time axis (e.g., 16/24 frames)
2. Enable cross-frame information exchange



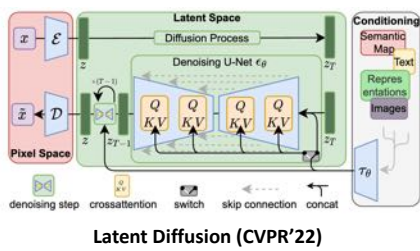
AnimateDiff:

Animate Your Personalized Text-to-Image Diffusion Model without Specific Tuning

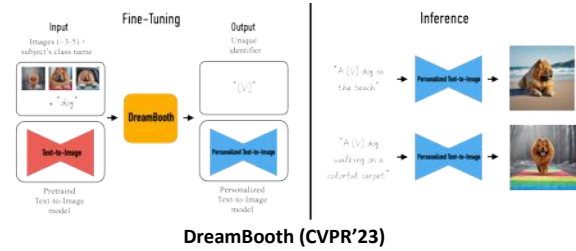


AnimateDiff: Repurpose image diffusion model for video generation

Image Generation Foundation Models, e.g.,



Model Personalization Methods, e.g.,



High-quality Personalized Models on HuggingFace and CivitAI

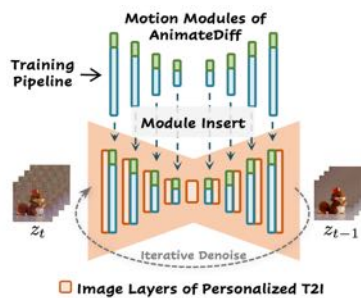


AnimateDiff: Method

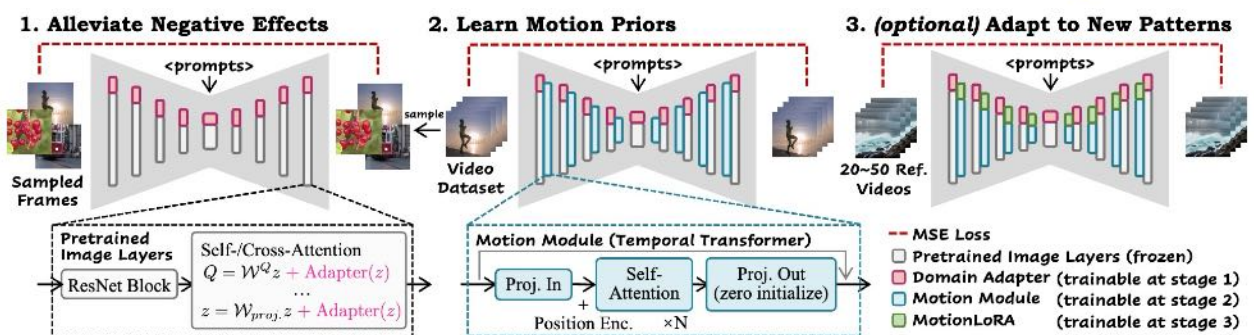
Our Goals: Animating personalized text-to-image diffusion models with a motion module, which

- Preserves original models' visual quality, → 1st stage, Domain Adapter
- Learns transferable motion priors from real-world videos, and → 2nd stage, Motion Module
- Efficiently adapts to specific motion patterns. → 3rd stage, MotionLoRA

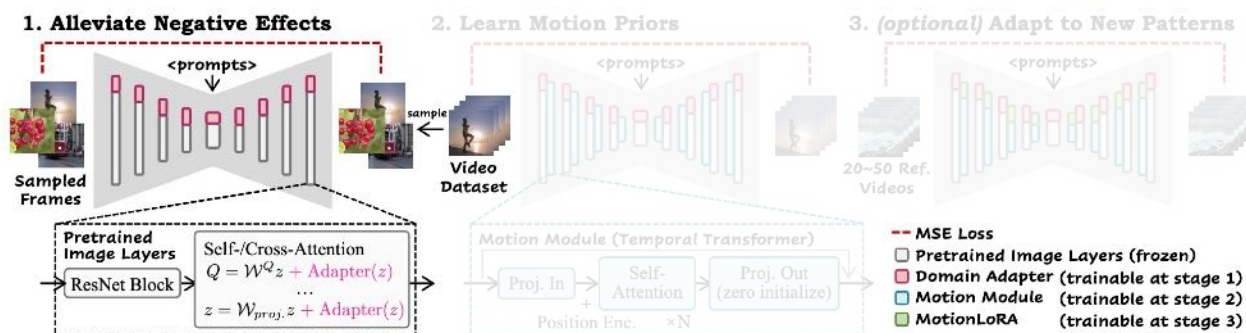
At inference, we directly insert the pre-trained motion module without needing specific tuning.



AnimateDiff: Method



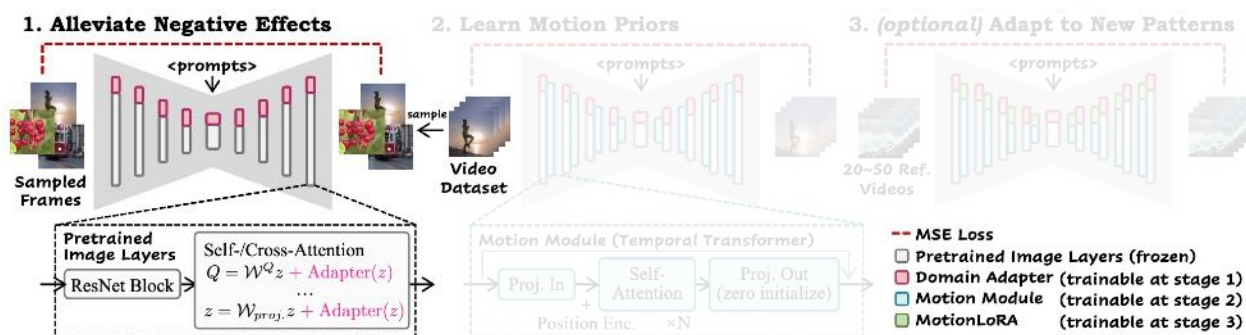
AnimateDiff: Method



Training Domain Adapter (1st stage): Alleviate Negative Effects from Training Data

- Video datasets' lower quality: watermarks, motion blurs, and compression artifacts
- Solution: learning such visual patterns with domain adapter and removing it at inference

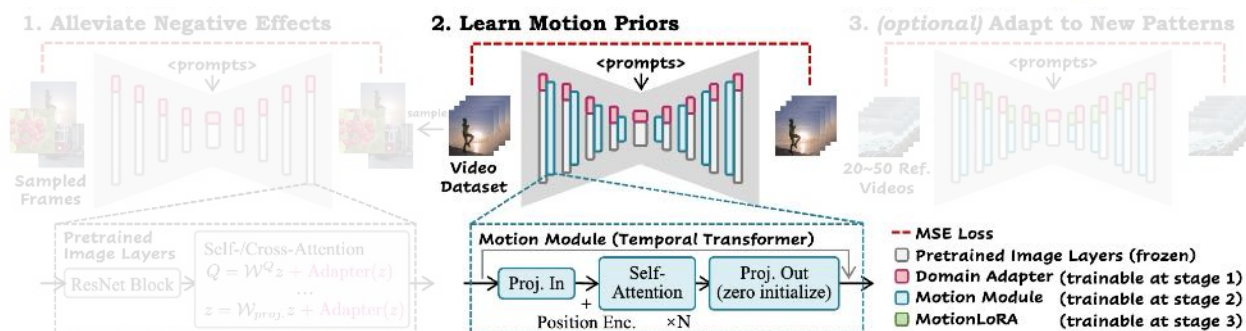
AnimateDiff: Method



Ablation Study: a lower domain adapter's effect helps preserve the original model's visual quality



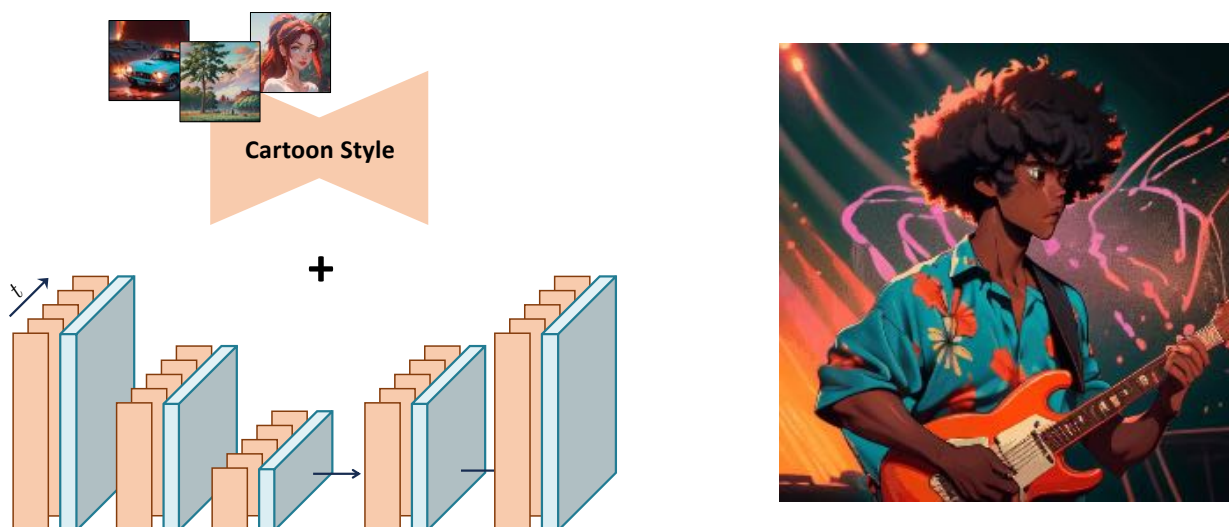
AnimateDiff: Method



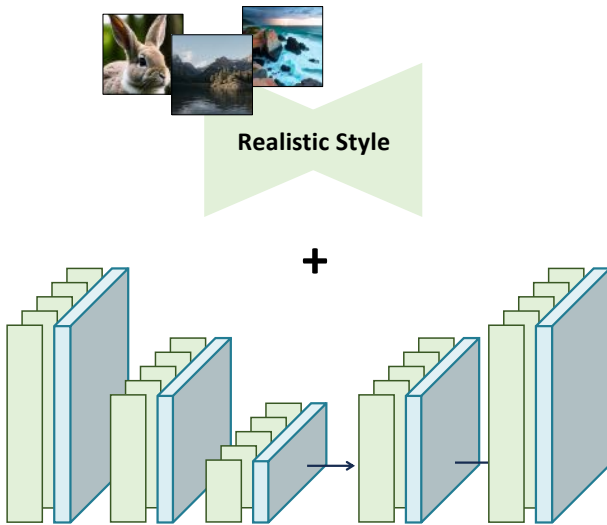
Training Motion Module (2nd stage): learning general motion priors from real-world videos

- **Model inflation:** from 2D image to 3D video
- **Temporal self-attention + position embeddings:** modeling cross-frame interactions
- Motion modules are inserted between frozen 2D image layers

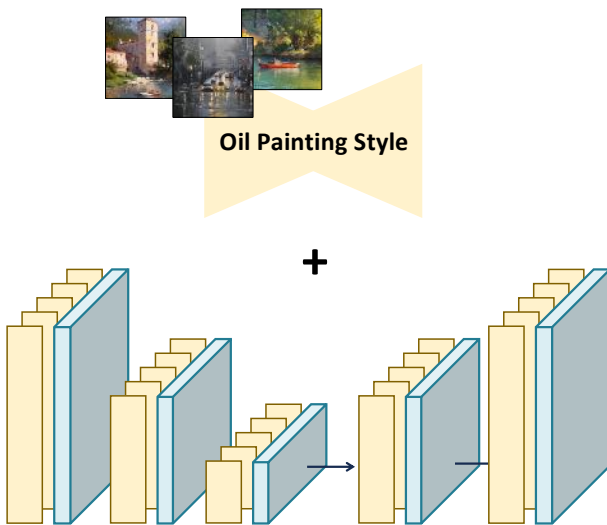
AnimateDiff: Method



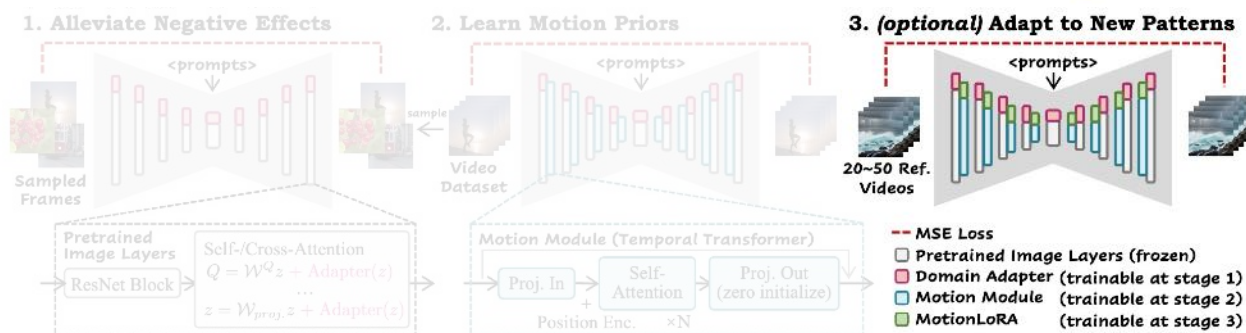
AnimateDiff: Method



AnimateDiff: Method



AnimateDiff: Method



Training MotionLoRA (3rd stage, optional): adapting to specific motion patterns

- Motion patterns like zooming and rolling are common in productions
- **Solution:** training additional LoRA adapter upon motion module's pre-trained weights, with few numbers of reference videos

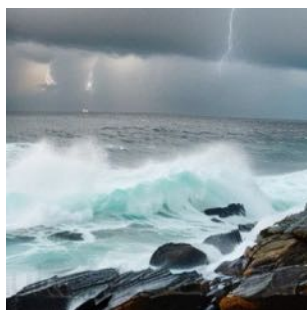
AnimateDiff: Method



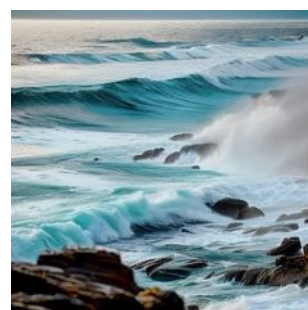
zoom in



rolling



zoom out + rolling



right + up

AnimateDiff: Experiments

Training

- Dataset: WebVid-10M
- Pre-trained text-to-image model: Stable Diffusion V1.5

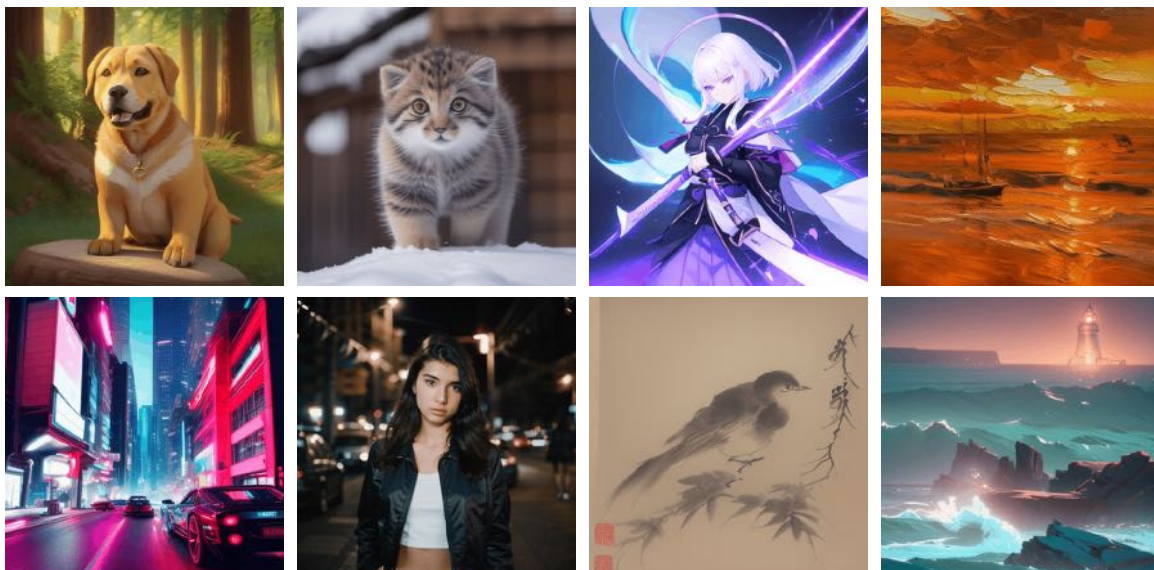
Evaluations

- Diverse model collected from the community

Model Name	Domain	Type
ToonYou	2D Cartoon	T2I Base Model
MeinaMix	2D Anime	T2I Base Model
Lyriel	Stylistic	T2I Base Model
RCNZ Cartoon 3d	3D Cartoon	T2I Base Model
epiC Realism	Realistic	T2I Base Model
Realistic Vision	Realistic	T2I Base Model
Oil painting	Stylistic	LoRA
MoXin	Stylistic	LoRA
TUSUN	Concept	LoRA

AnimateDiff: Experiments

Qualitative Results: on eight different community model



AnimateDiff: Experiments

Quantitative Evaluation

- Our method is preferred by user study and CLIP metrics in text/domain fidelity and temporal smoothness

Method	User Study (↑)			CLIP Metric (↑)		
	Text.	Domain.	Smooth.	Text.	Domain.	Smooth.
Text2Video-Zero	1.620	2.620	1.560	32.04	84.84	96.57
Tune-a-Video	2.180	1.100	1.615	35.98	80.68	97.42
Ours	2.210	2.280	2.825	31.39	87.29	98.00

AnimateDiff: Experiments

Compatibility with Text-to-Image Models' Adapter

- AnimateDiff can be directly used with pre-trained T2I adapters, e.g., ControlNet, for controllable generations
- Depth-guided generation with ControlNet-depth



Generating Higher Spatial/Temporal Resolution

Cascaded pipeline

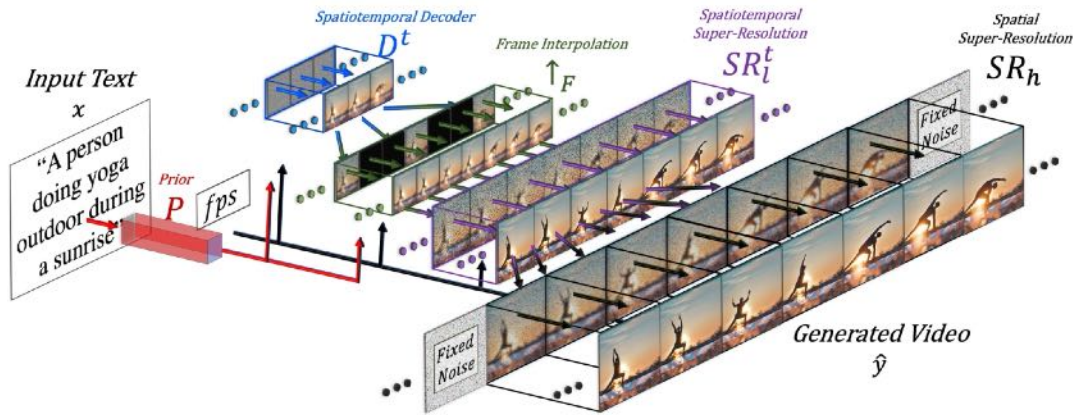


Image credit: Make-A-Video

Generating Higher Spatial/Temporal Resolution

Spatial-Temporal Architecture

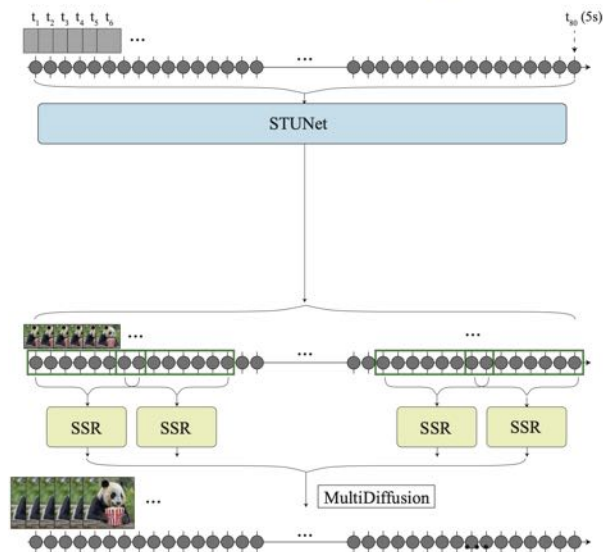
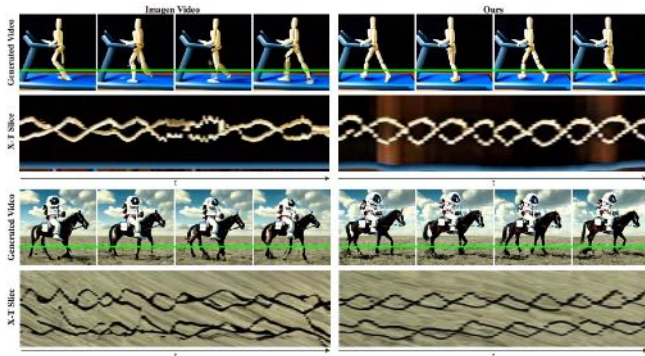
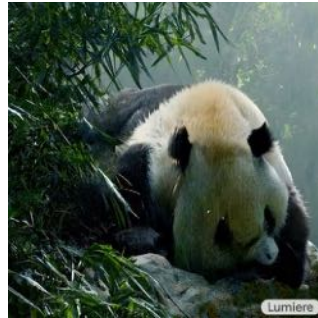


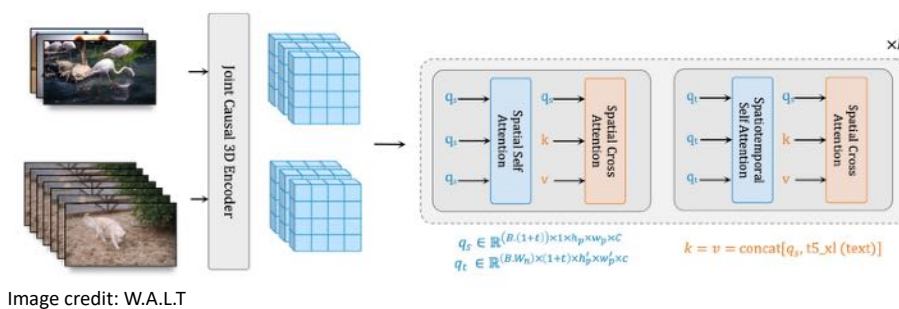
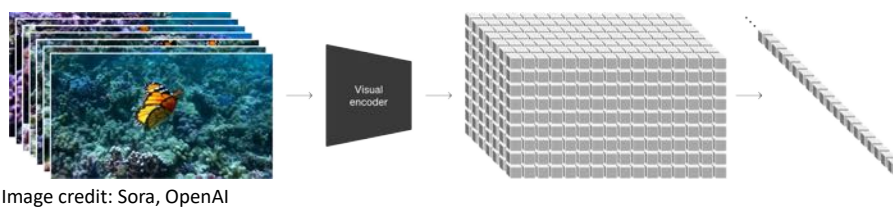
Image credit: Lumiere

Generating Higher Spatial/Temporal Resolution

Spatial-Temporal Architecture

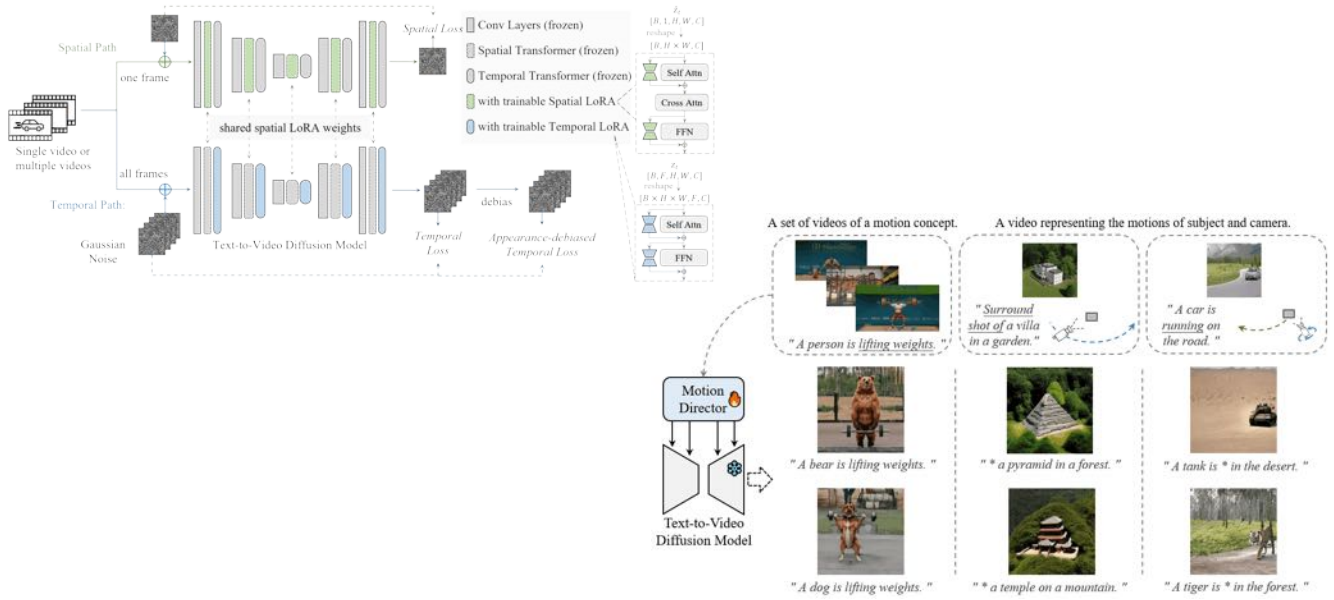


Transformer-based Approaches



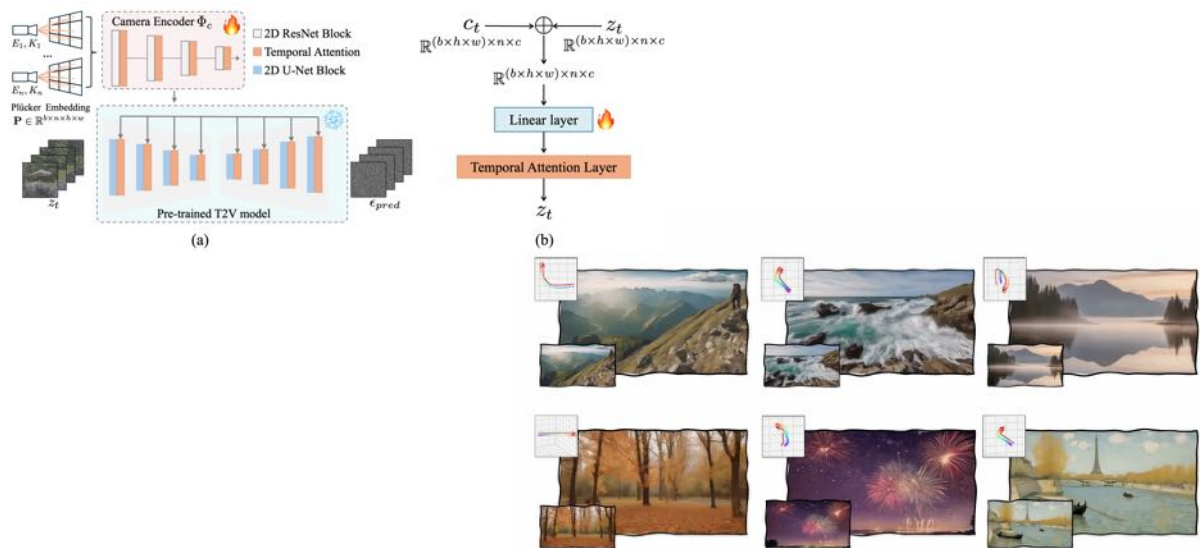
Controllable Generation

Motion: MotionDirector: Motion Customization of Text-to-Video Diffusion Models



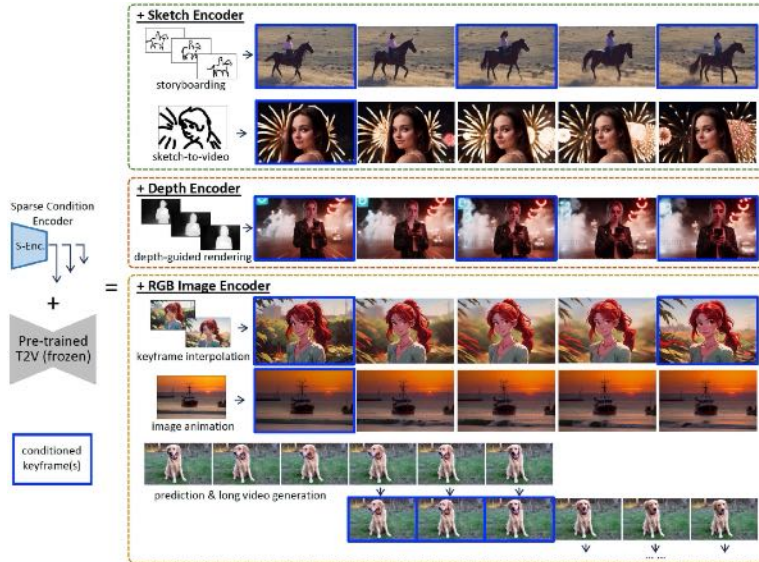
Controllable Generation

Camera: CameraCtrl: Enabling Camera Control for Video Diffusion Models



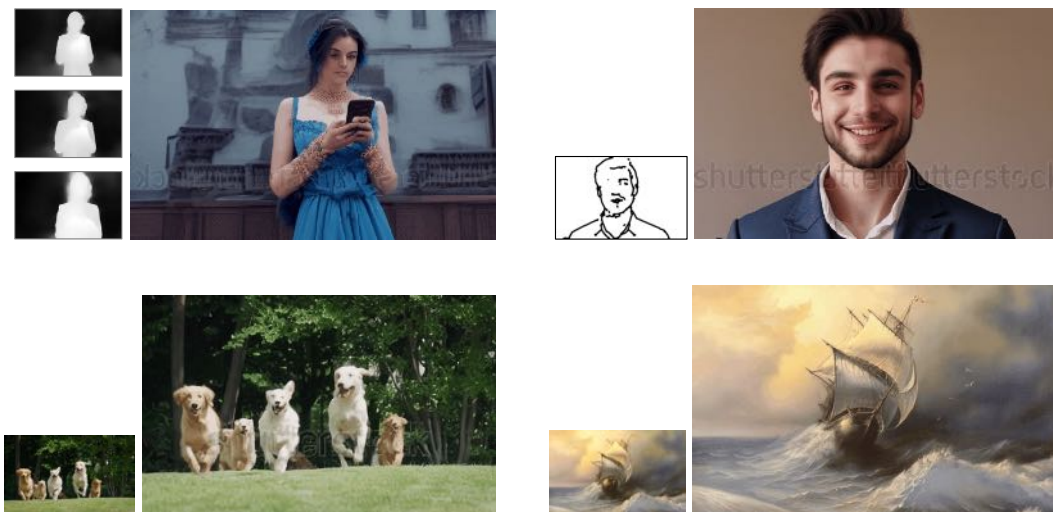
Controllable Generation

Layout & Pixel: SparseCtrl: Adding Sparse Controls to Text-to-Video Diffusion Models



Controllable Generation

Layout & Pixel: SparseCtrl: Adding Sparse Controls to Text-to-Video Diffusion Models



Diffusion Based Video Editing

Global Stylization: Structure and Content-Guided Video Synthesis with Diffusion Models (GEN-1)

Decouple the structure and appearance via depth maps

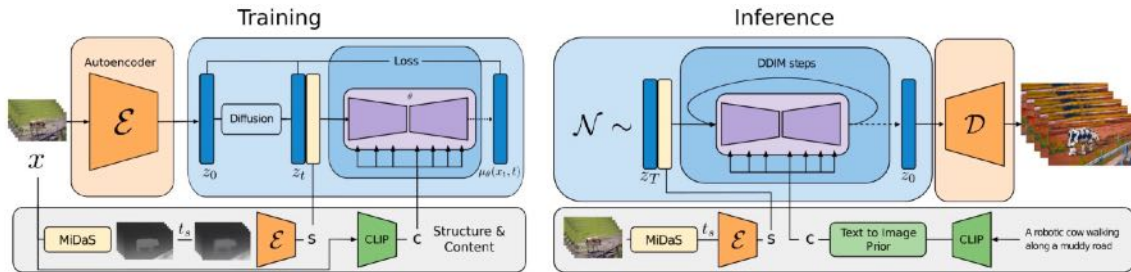
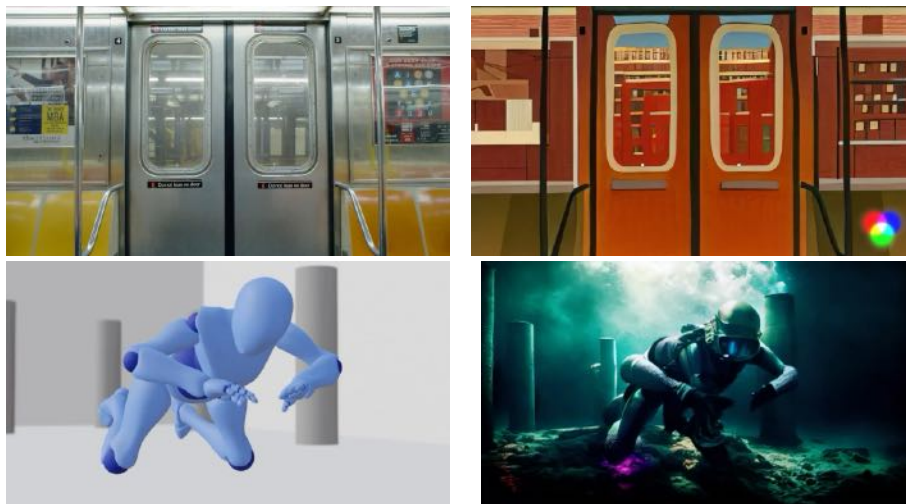


Image credit: GEN-1

Diffusion Based Video Editing

Global Stylization: Structure and Content-Guided Video Synthesis with Diffusion Models (GEN-1)

Decouple the structure and appearance via depth maps



Diffusion Based Video Editing

Local: Fate/Zero: Fusing Attentions for Zero-shot Text-based Video Editing

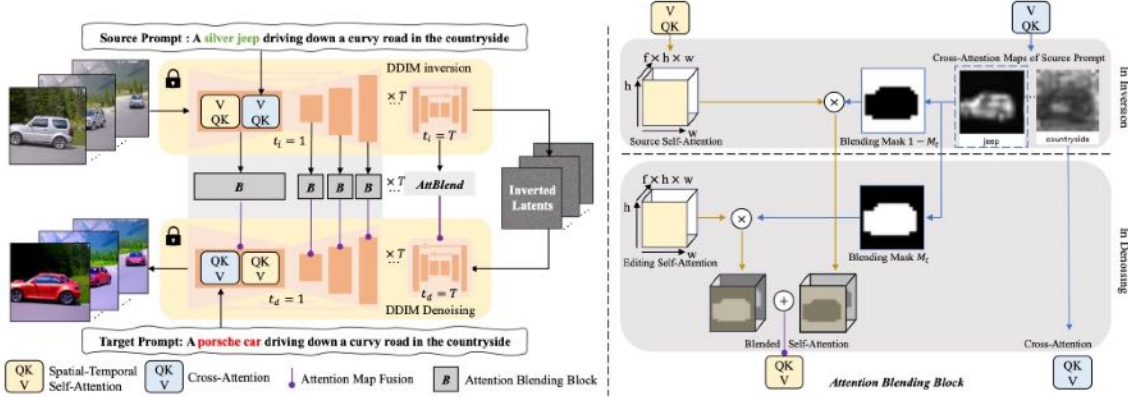


Image credit: FateZero

Diffusion Based Video Editing

Local: Fate/Zero: Fusing Attentions for Zero-shot Text-based Video Editing

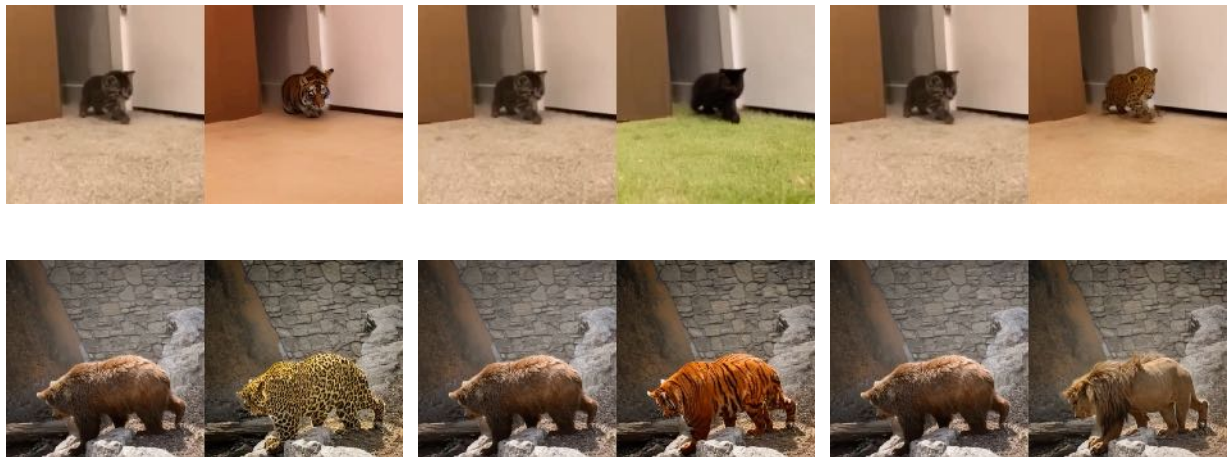
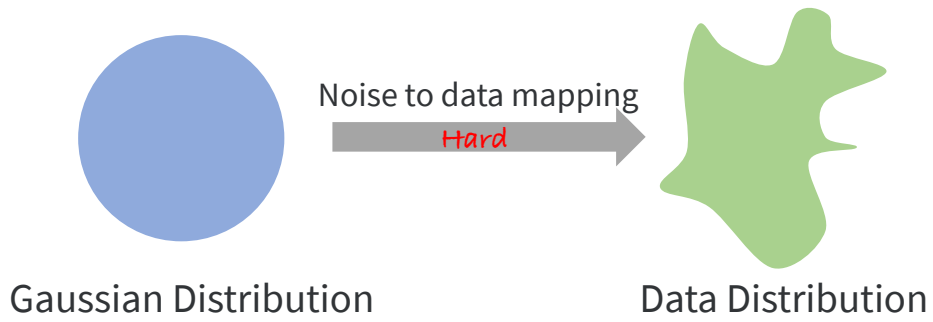


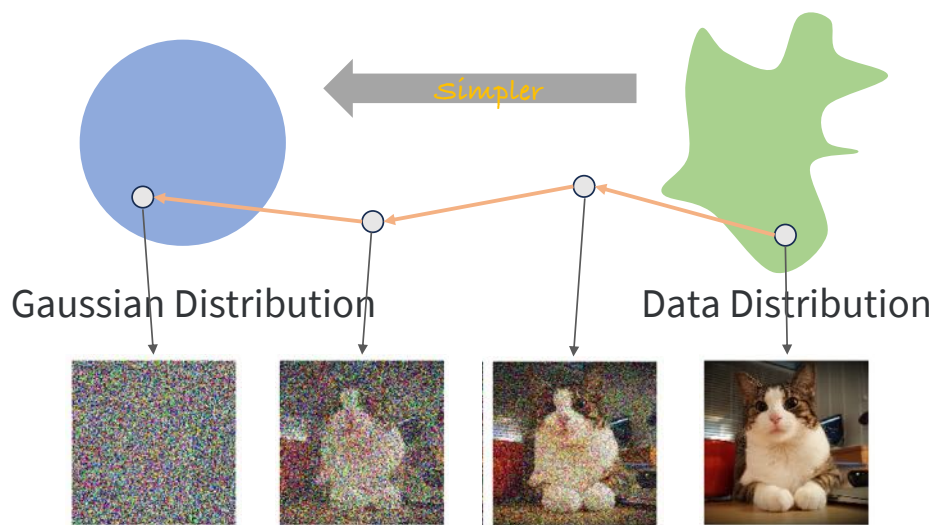
Image credit: FateZero

Recap Diffusion



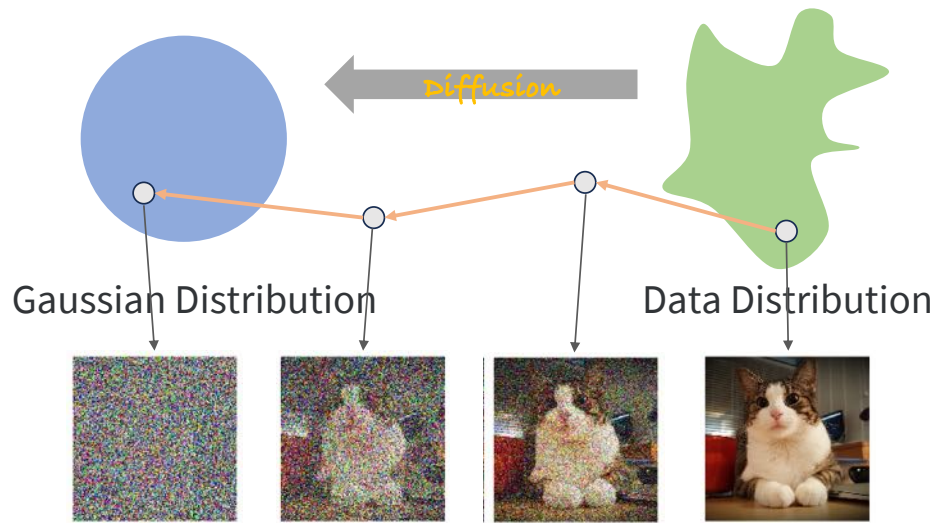
1

Recap Diffusion

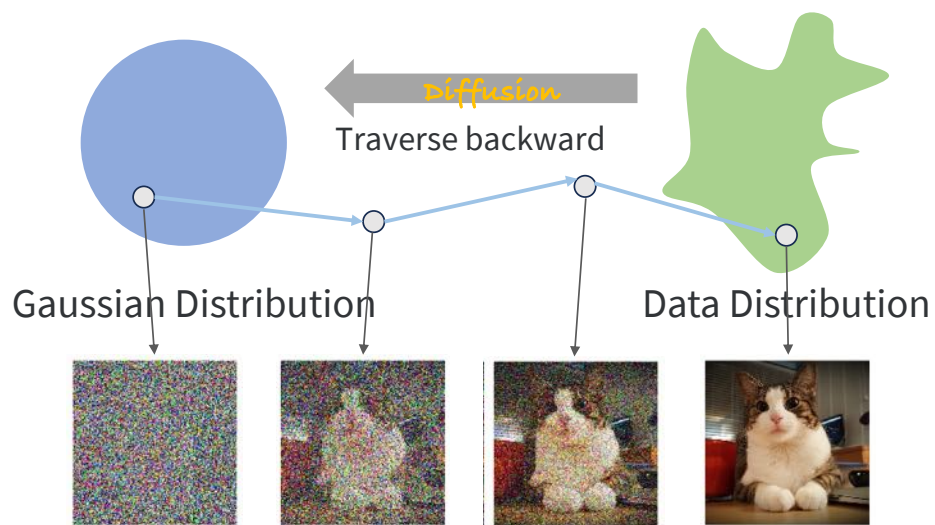


2

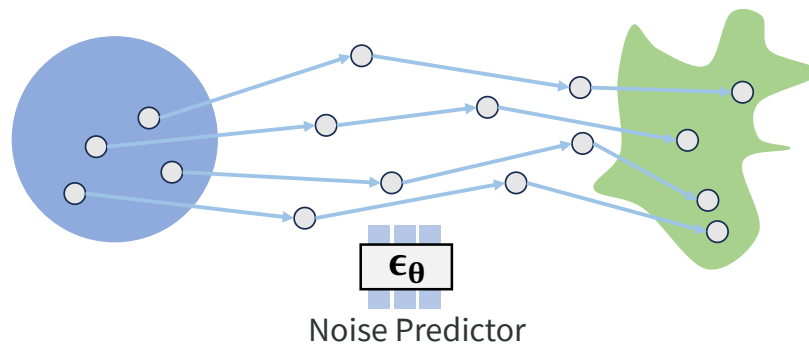
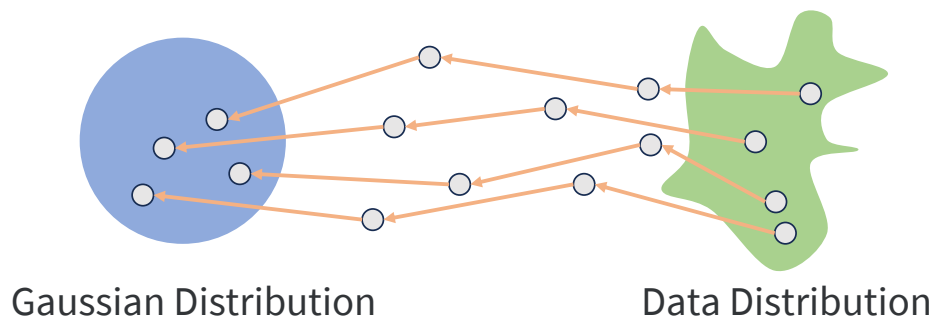
Recap Diffusion



Recap Diffusion

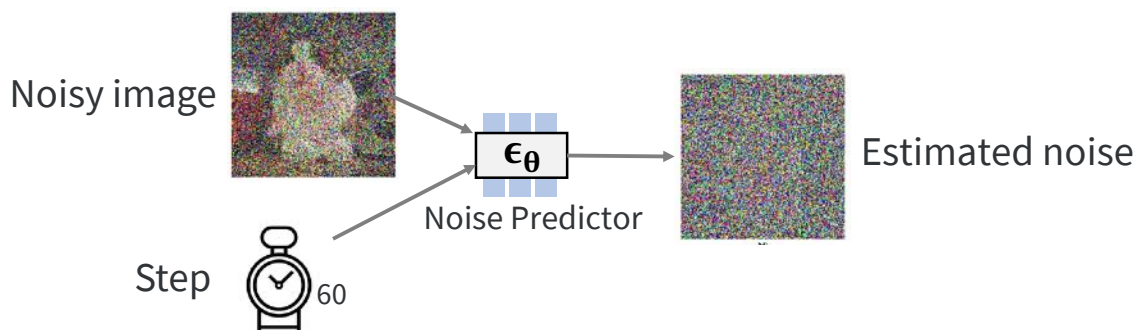
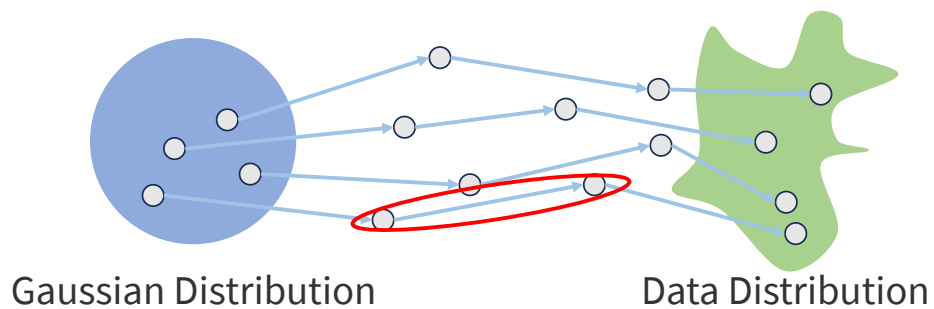


Recap Diffusion



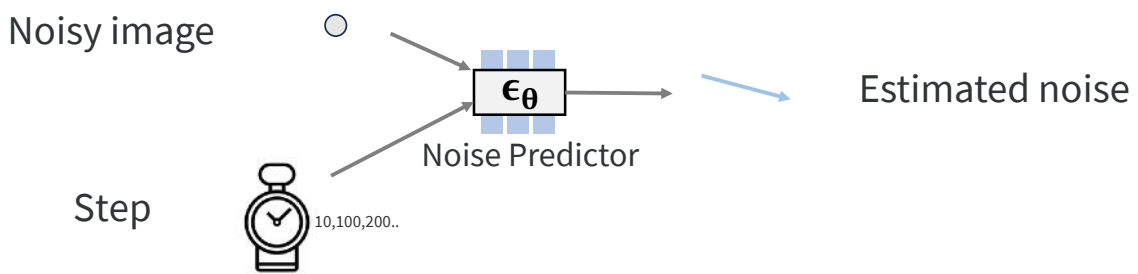
5

Recap Diffusion



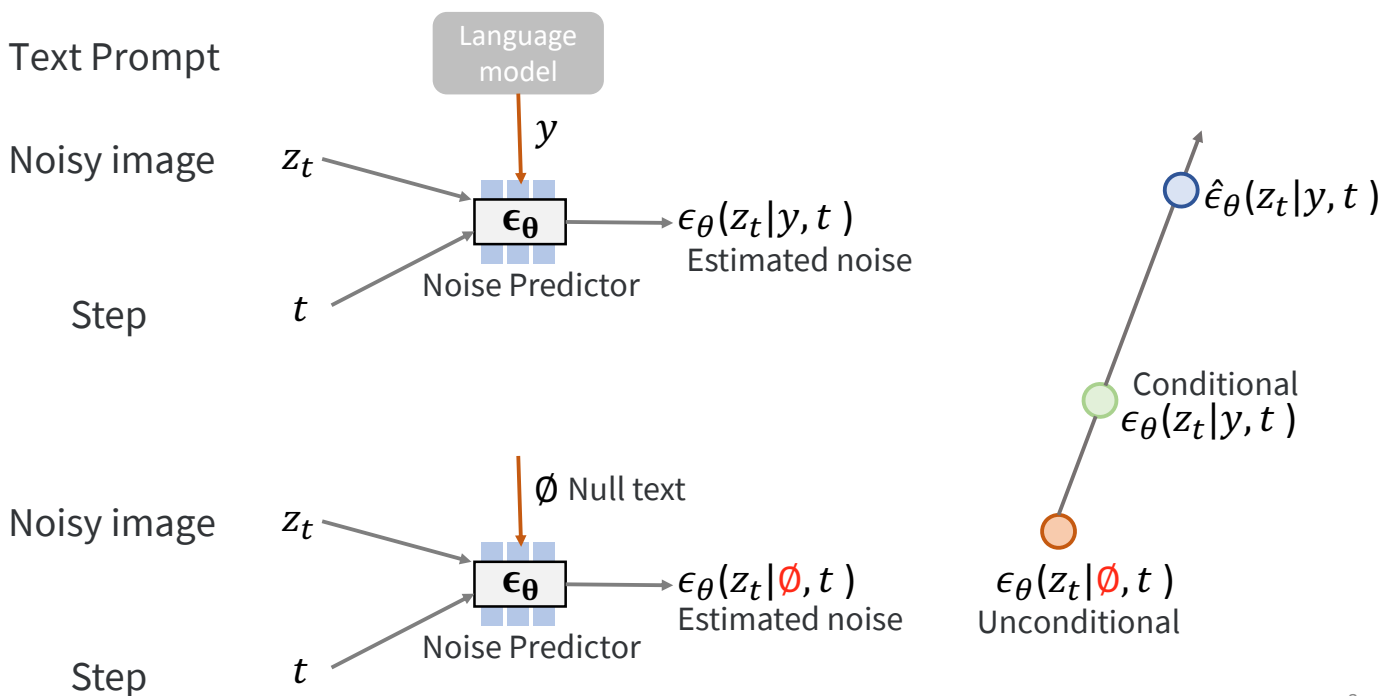
6

Recap Diffusion



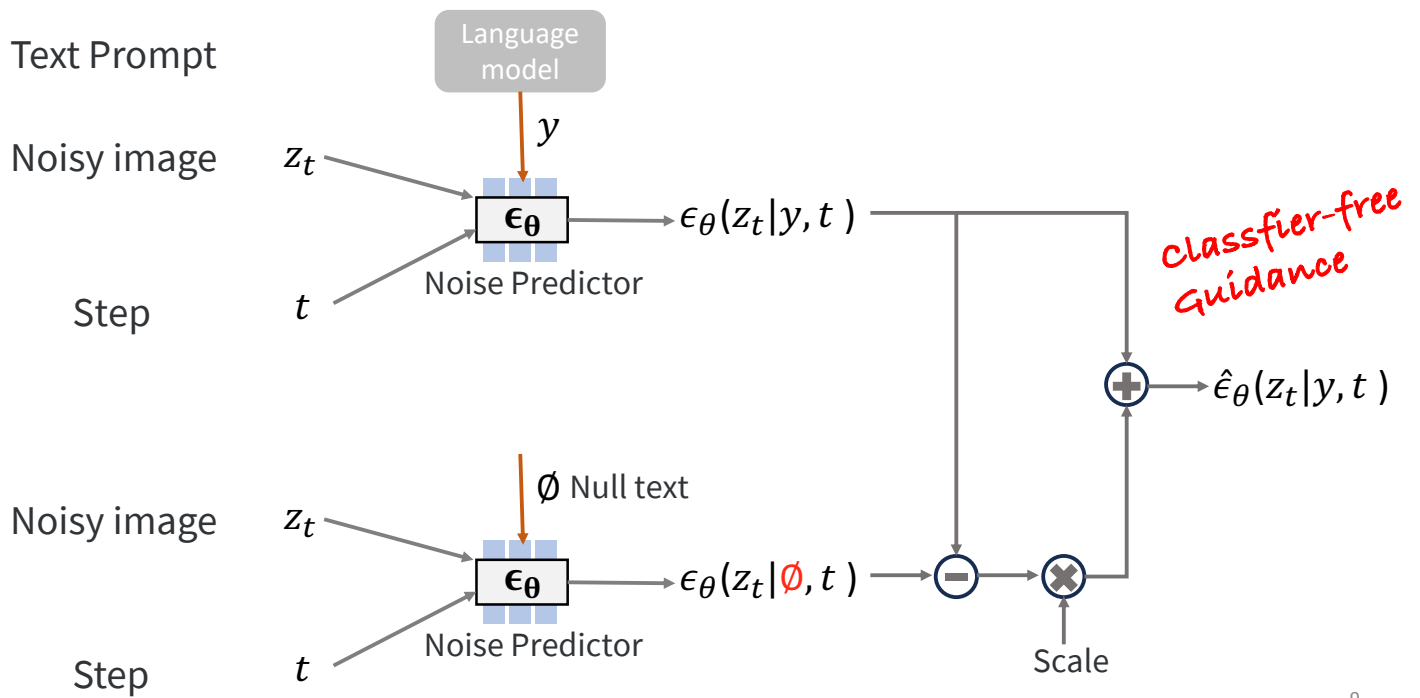
7

Text Conditioned Diffusion



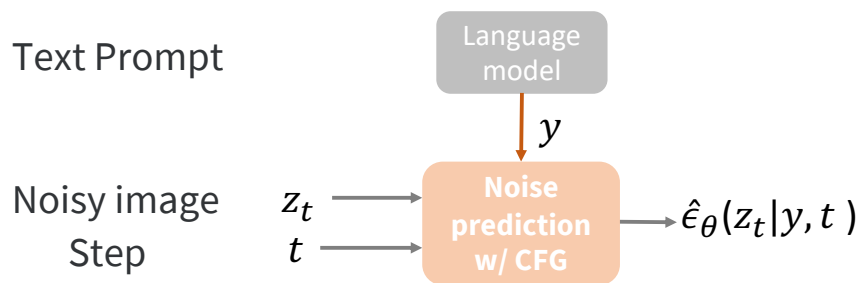
8

Text Conditioned Diffusion

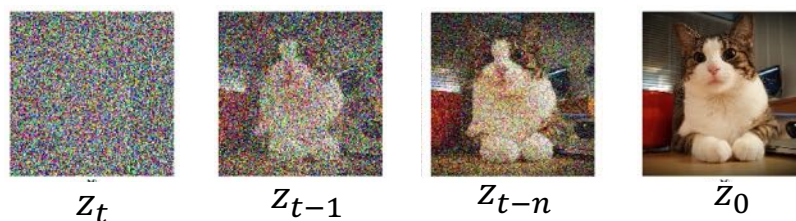


9

2D Diffusion to 3D

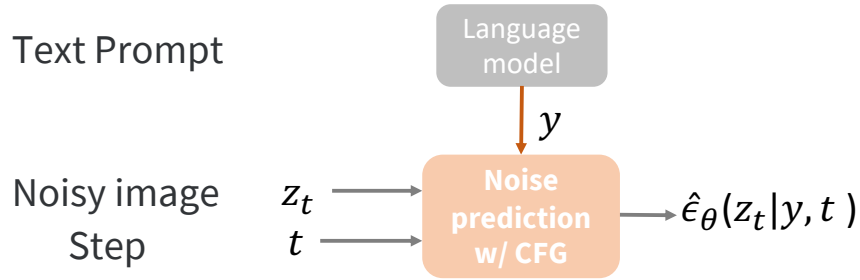


update samples in pixel space ----- 2D images



10

2D Diffusion to 3D



update samples in pixel space ----- 2D images

How about 3D?



Parameter
 $x = g(\theta)$ *update in parameter space?*
 Rendered Image Differentiable Rendering

3D Content Generation *Empowered by 2D Diffusion Priors*

$$L_{diff}(\phi, x) = \mathbb{E}_{t \sim U(0,1), \epsilon \sim \mathcal{N}(0,I)} [w(t) \|\epsilon_\phi(z_t | y, t) - \epsilon\|_2^2]$$

$$z_t = a_t x + \sigma_t \epsilon$$

Training a diffusion model:
 $\phi^* = \operatorname{argmin}_\phi L_{diff}(\phi, x)$

With a trained diffusion model:
 $x^* = \operatorname{argmin}_x L_{diff}(\phi, x)$



Parameter
 $x = g(\theta)$
 Rendered Image Differentiable Rendering

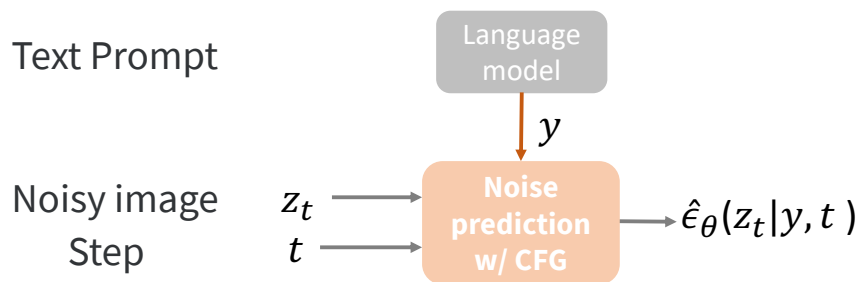
3D Content Generation *Empowered by 2D Diffusion Priors*

$$L_{diff}(\phi, g(\theta)) = \mathbb{E}_{t \sim \mathcal{U}(0,1), \epsilon \sim \mathcal{N}(0,I)} [w(t) \|\epsilon_\phi(z_t|y, t) - \epsilon\|_2^2]$$

$$\nabla_\theta L_{diff}(\phi, g(\theta)) = \mathbb{E}_{t, \epsilon} [w(t) (\hat{\epsilon}_\phi(z_t|y, t) - \epsilon) \frac{\partial \hat{\epsilon}_\phi(z_t|y, t)}{\partial z_t} \frac{\partial x}{\partial \theta}]$$

Noise Residual
U-Net Jacobian
Generator Jacobian

update 3D representation w/ gradient descent



13

3D Content Generation *Empowered by 2D Diffusion Priors*

$$L_{diff}(\phi, g(\theta)) = \mathbb{E}_{t \sim \mathcal{U}(0,1), \epsilon \sim \mathcal{N}(0,I)} [w(t) \|\epsilon_\phi(z_t|y, t) - \epsilon\|_2^2]$$

$$\nabla_\theta L_{diff}(\phi, g(\theta)) = \mathbb{E}_{t, \epsilon} [w(t) (\hat{\epsilon}_\phi(z_t|y, t) - \epsilon) \frac{\partial \hat{\epsilon}_\phi(z_t|y, t)}{\partial z_t} \frac{\partial x}{\partial \theta}]$$

Noise Residual
~~U-Net Jacobian~~
~~Generator Jacobian~~

$$\nabla_\theta L_{SDS}(\phi, g(\theta)) = \mathbb{E}_{t, \epsilon} [w(t) (\hat{\epsilon}_\phi(z_t|y, t) - \epsilon) \frac{\partial x}{\partial \theta}]$$



DreamFusion: Text-to-3D using 2D Diffusion

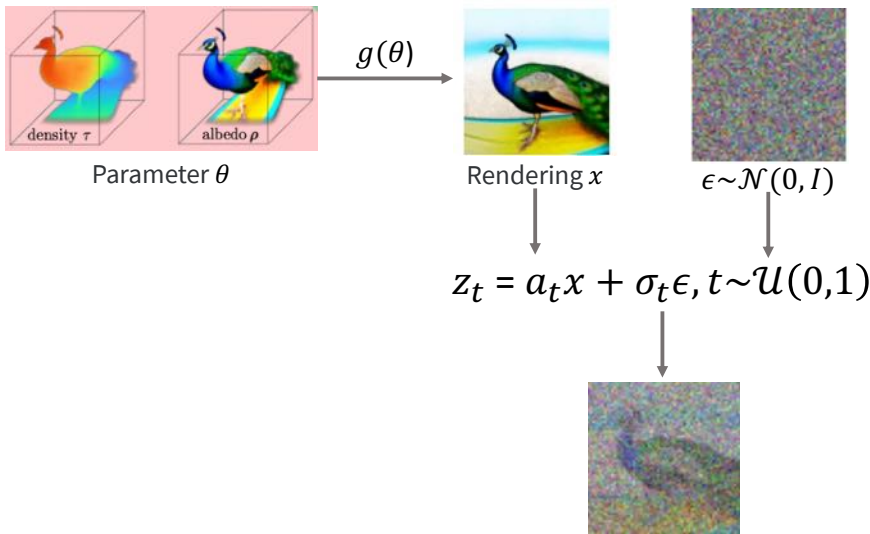
Ben Poole (Google Research), Ajay Jain (UC Berkeley), Jonathan T. Barron (Google Research), Ben Mildenhall (Google Research)

Score Distillation Sampling

14

3D Content Generation *Empowered by 2D Diffusion Priors*

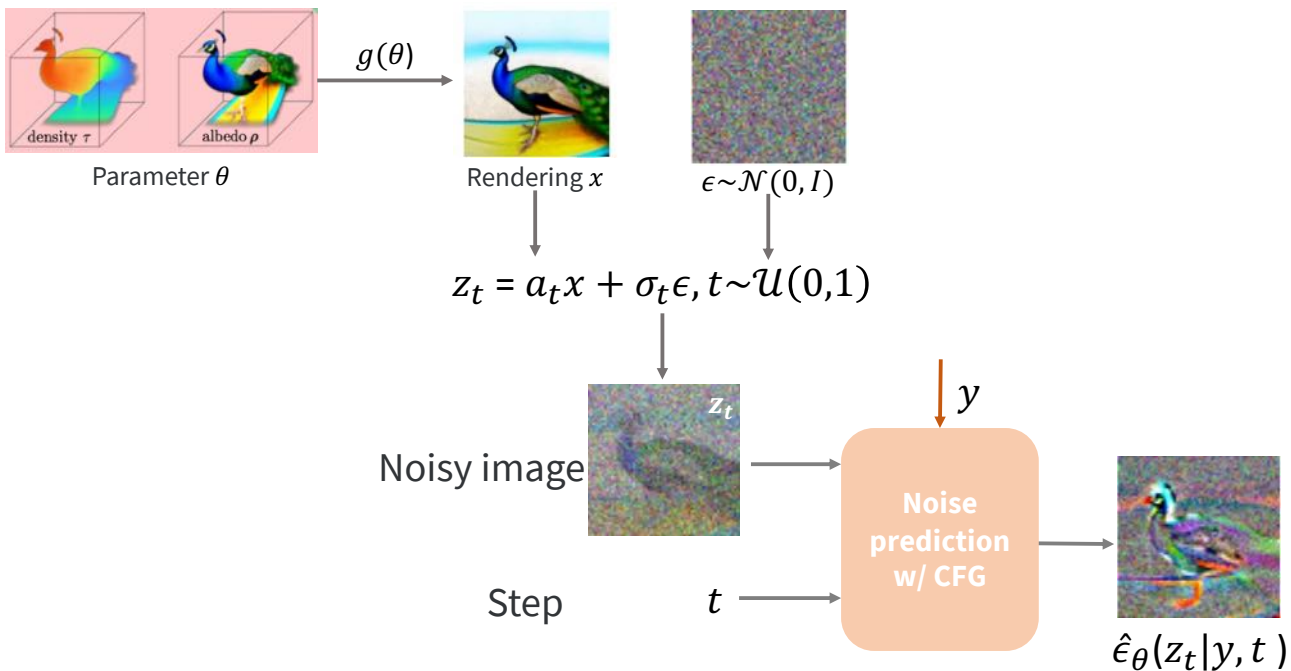
Score Distillation Sampling (SDS) Loss



15

3D Content Generation *Empowered by 2D Diffusion Priors*

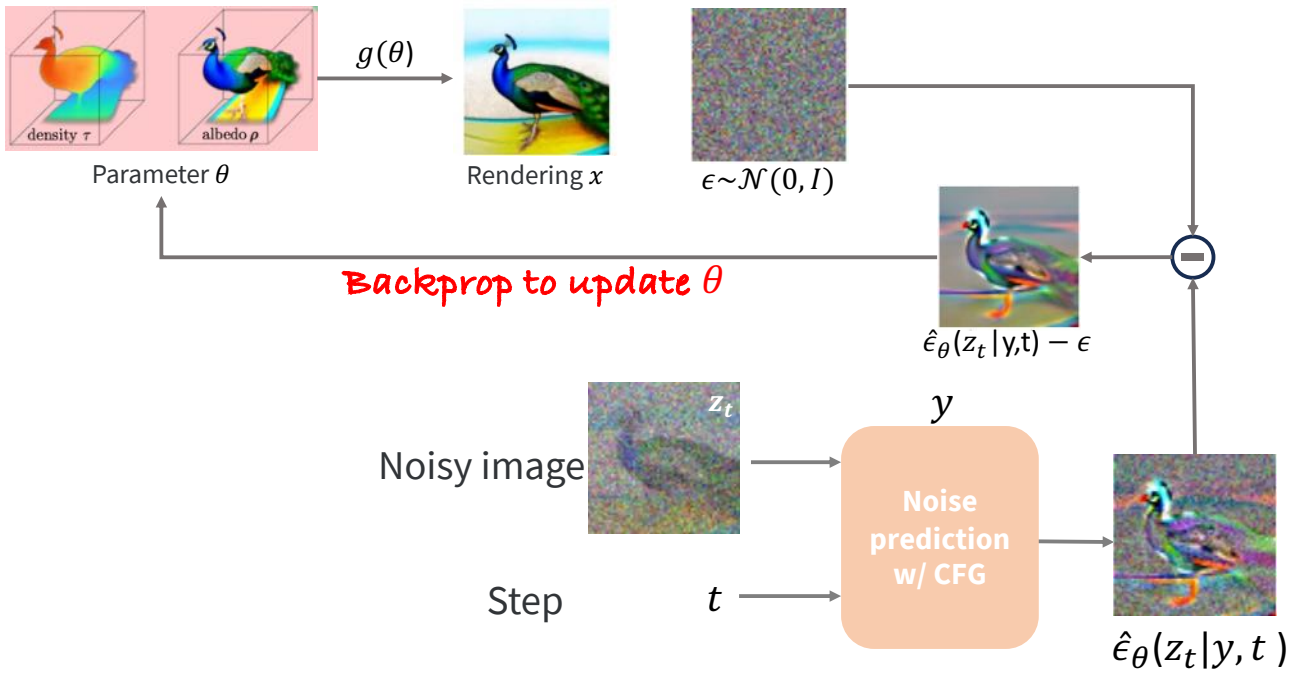
Score Distillation Sampling (SDS) Loss



16

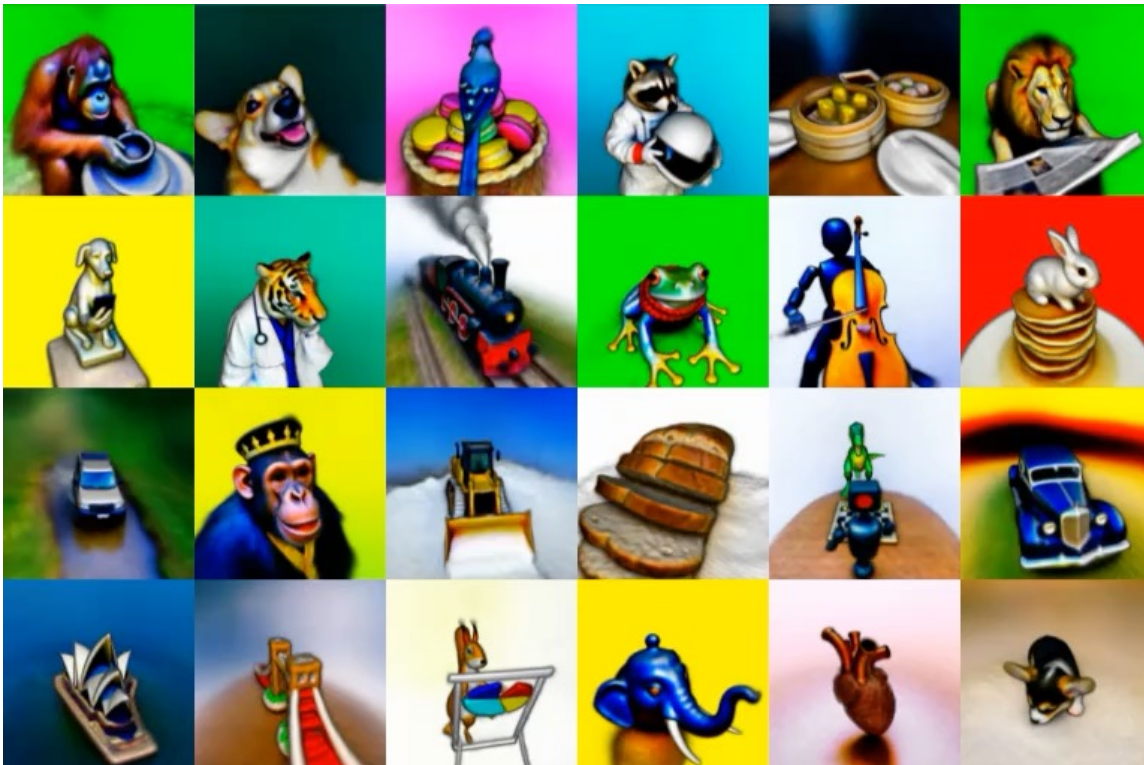
3D Content Generation *Empowered by 2D Diffusion Priors*

Score Distillation Sampling (SDS) Loss



17

3D Content Generation *Empowered by 2D Diffusion Priors*



18

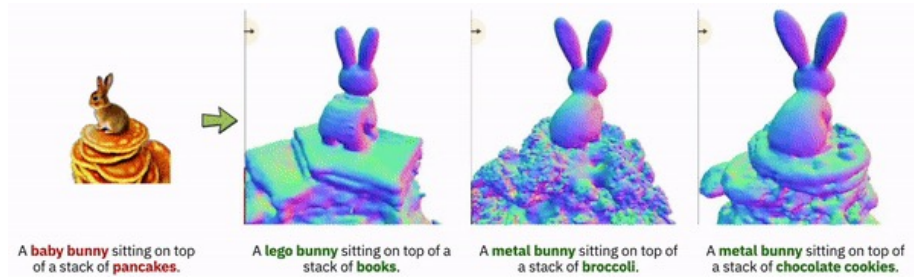
3D Content Generation *Empowered by 2D Diffusion Priors*

DreamFusion: Text-to-3D using 2D Diffusion



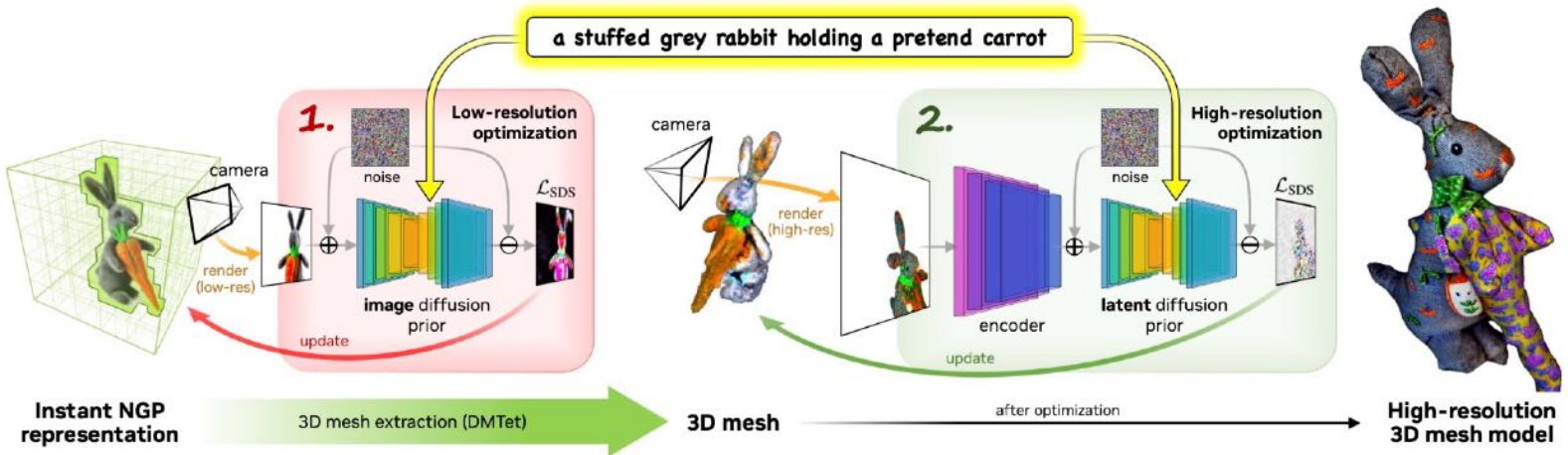
↓ Higher resolution

Magic3D: High-Resolution Text-to-3D Content Creation



3D Content Generation *Empowered by 2D Diffusion Priors*

Magic3D: High-Resolution Text-to-3D Content Creation



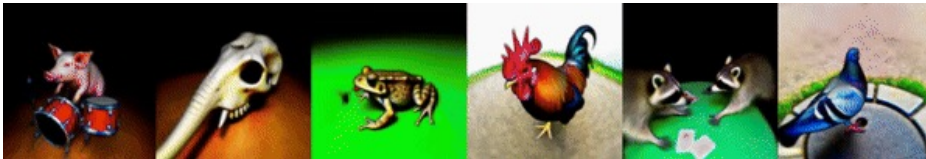
3D Content Generation *Empowered by 2D Diffusion Priors*



21

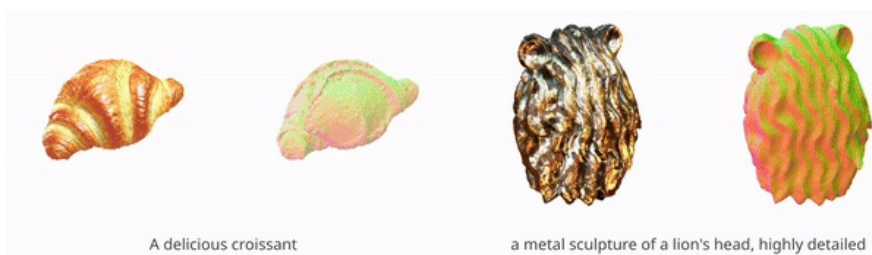
3D Content Generation *Empowered by 2D Diffusion Priors*

DreamFusion: Text-to-3D using 2D Diffusion



↓ **Richer appearance**

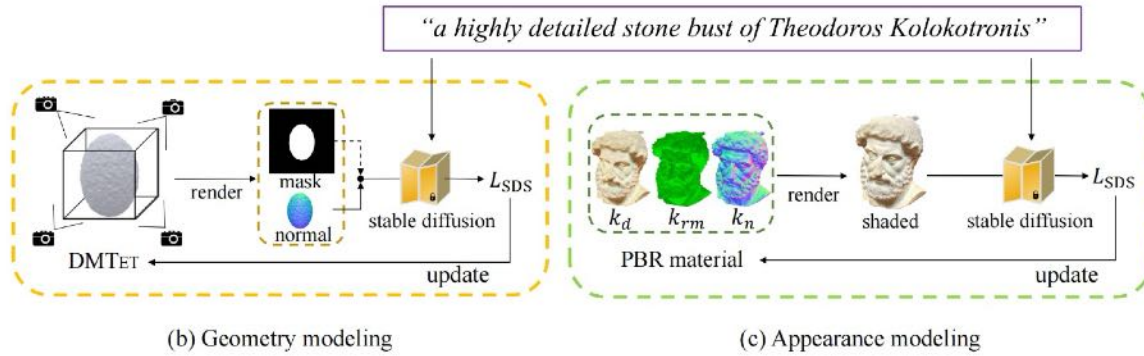
Fantasia3D: Disentangling Geometry and Appearance for High-quality Text-to-3D Content Creation



22

3D Content Generation *Empowered by 2D Diffusion Priors*

Fantasia3D: Disentangling Geometry and Appearance for High-quality Text-to-3D Content Creation



23

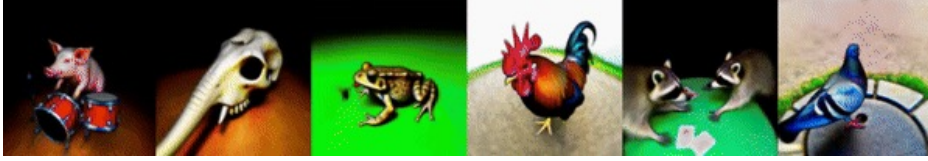
3D Content Generation *Empowered by 2D Diffusion Priors*



24

3D Content Generation *Empowered by 2D Diffusion Priors*

DreamFusion: Text-to-3D using 2D Diffusion



↓ *Single image to 3D*

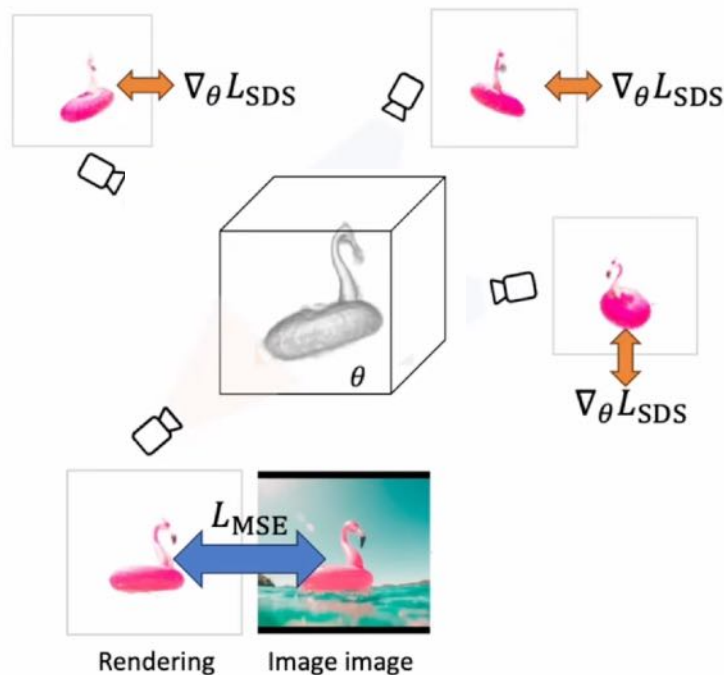
RealFusion: 360° Reconstruction of Any Object from a Single Image



25

3D Content Generation *Empowered by 2D Diffusion Priors*

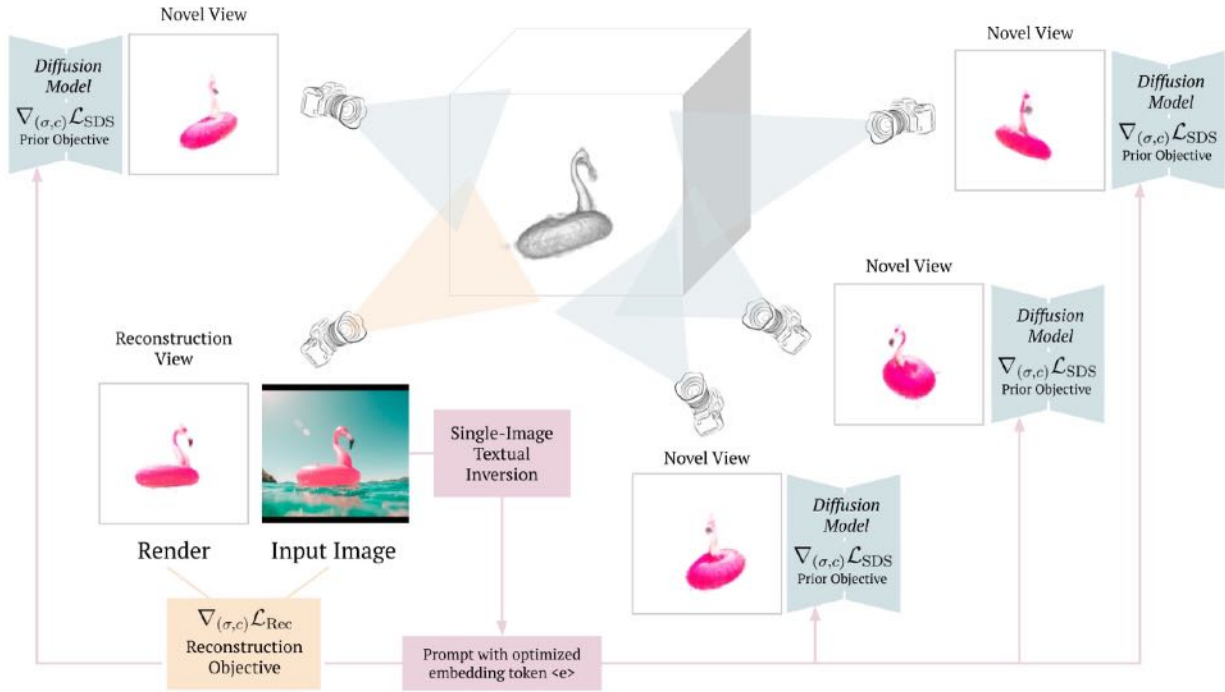
RealFusion: 360° Reconstruction of Any Object from a Single Image



26

3D Content Generation *Empowered by 2D Diffusion Priors*

RealFusion: 360° Reconstruction of Any Object from a Single Image



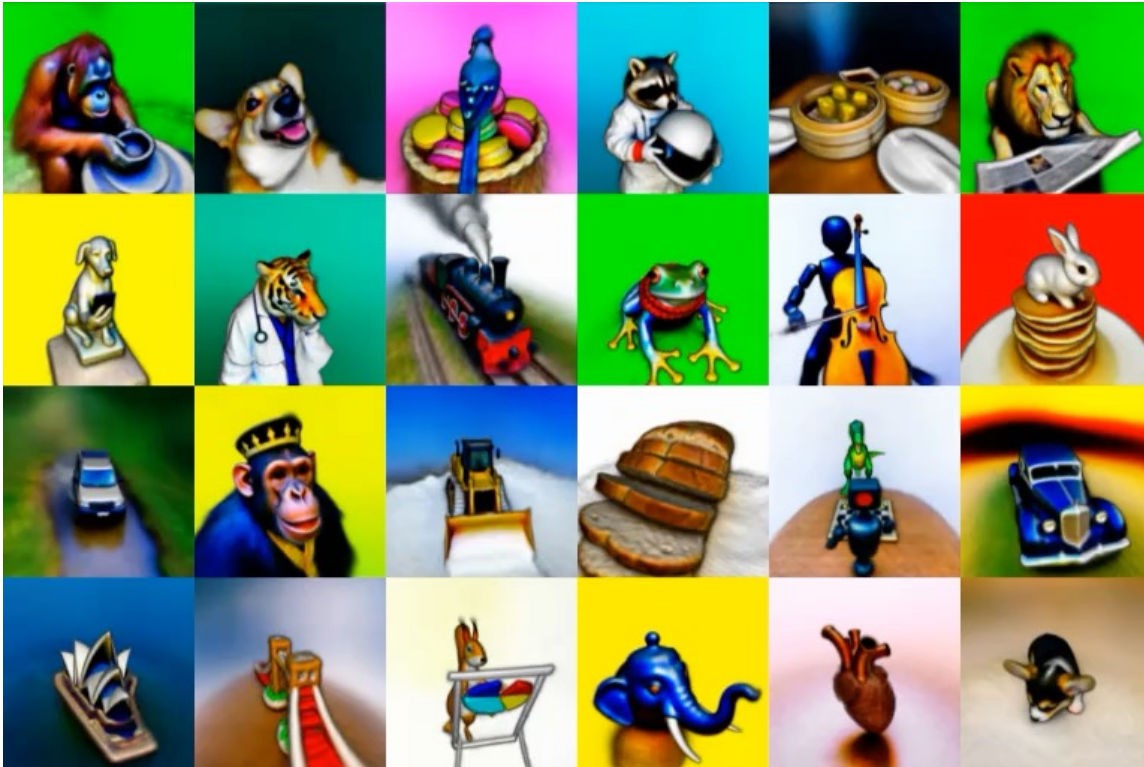
27

3D Content Generation *Empowered by 2D Diffusion Priors*



8

3D Content Generation *Empowered by 2D Diffusion Priors*



29

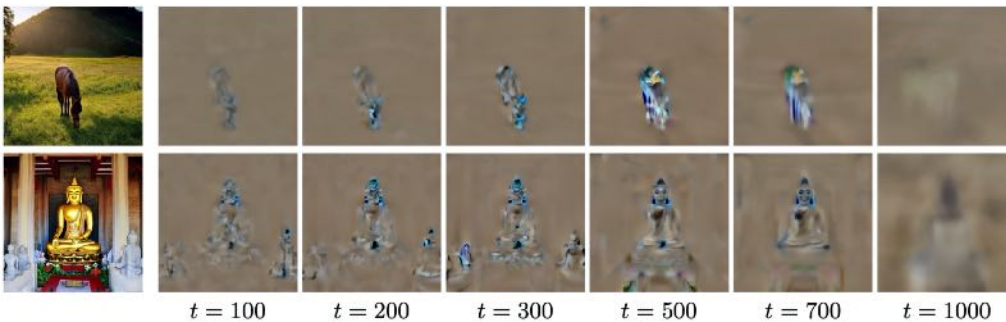
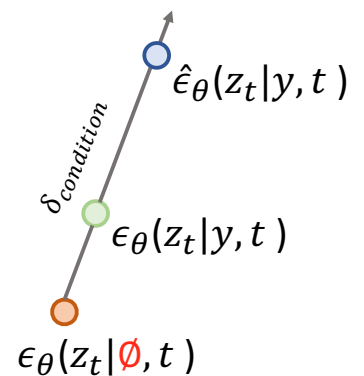
3D Content Generation *Empowered by 2D Diffusion Priors*

Photorealistic appearance?

$$\nabla_{\theta} L_{SDS}(\phi, g(\theta)) = \mathbb{E}_{t, \epsilon} [w(t) (\hat{\epsilon}_{\theta}(z_t|y, t) - \epsilon) \frac{\partial x}{\partial \theta}]$$

Estimated noise
with classifier free guidance

$$\hat{\epsilon}_{\theta}(z_t|y, t) = \epsilon_{\phi}(z_t|\emptyset, t) + s(\underbrace{\epsilon_{\phi}(z_t|y, t) - \epsilon_{\phi}(z_t|\emptyset, t)}_{\delta_{condition}})$$



aligned with the condition;
uncorrelated with the added noise ϵ

30

3D Content Generation Empowered by 2D Diffusion Priors

Photorealistic appearance?

$$\nabla_{\theta} L_{SDS}(\phi, g(\theta)) = \mathbb{E}_{t, \epsilon} [w(t) (\hat{\epsilon}_{\theta}(z_t | y, t) - \epsilon) \frac{\partial x}{\partial \theta}]$$

Estimated noise
with classifier free guidance

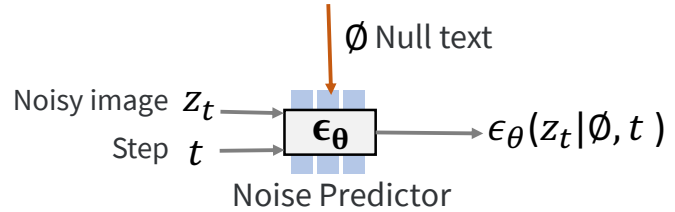
$$\hat{\epsilon}_{\theta}(z_t | y, t) = \epsilon_{\phi}(z_t | \emptyset, t) + s \delta_{condition}$$

$$z_t = a_t x + \sigma_t \epsilon$$

Training: Real image x

SDS: Rendered image $x = g(\theta)$

Domain difference



31

3D Content Generation Empowered by 2D Diffusion Priors

Photorealistic appearance?

$$\nabla_{\theta} L_{SDS}(\phi, g(\theta)) = \mathbb{E}_{t, \epsilon} [w(t) (\hat{\epsilon}_{\theta}(z_t | y, t) - \epsilon) \frac{\partial x}{\partial \theta}]$$

Estimated noise
with classifier free guidance

$$\hat{\epsilon}_{\theta}(z_t | y, t) = \epsilon_{\phi}(z_t | \emptyset, t) + s \delta_{condition}$$

$$\epsilon_{\phi}(z_t | \emptyset, t) = \delta_{domain} + \delta_{denoising}$$



x_{ID}



$\delta_{denoising}$



x_{OOD}



δ_{domain}



$x_{OOD} + \delta_{domain}$

32

3D Content Generation Empowered by 2D Diffusion Priors

Photorealistic appearance?

$$\nabla_{\theta} L_{SDS}(\phi, g(\theta)) = w(t)(\delta_{domain} + s\delta_{condition} + \delta_{denoising} - \epsilon) \frac{\partial x}{\partial \theta}$$

Domain-correction
Align w/ text
Diff b/w predicted noise & added noise

Training: Real image x
SDS: Rendered image $x = g(\theta)$
Domain difference



33

3D Content Generation Empowered by 2D Diffusion Priors

Photorealistic appearance?

$$\nabla_{\theta} L_{SDS}(\phi, g(\theta)) = w(t)(\delta_{domain} + s\delta_{condition} + \delta_{denoising} - \epsilon) \frac{\partial x}{\partial \theta}$$

Domain-correction
Align w/ text
Diff b/w predicted noise & added noise

Discard this one!

$$s\delta_{condition} = s(\underbrace{\epsilon_{\phi}(z_t|y, t)}_{\text{Conditional}} - \underbrace{\epsilon_{\phi}(z_t|\emptyset, t)}_{\text{Unconditional}})$$

$$\epsilon_{\phi}(z_t|\emptyset, t) = \delta_{domain} + \delta_{denoising} \quad \text{🤔 hard to separate}$$

Assumption: $\delta_{condition=p_{neg}} = -\delta_{domain}$



"unrealistic, blurry, low quality, out of focus, ugly, low contrast, dull, dark, low-resolution, gloomy"

34

3D Content Generation Empowered by 2D Diffusion Priors

Photorealistic appearance?

$$\nabla_{\theta} L_{SDS}(\phi, g(\theta)) = w(t)(\underbrace{\delta_{domain}}_{\text{Domain-correction}} + s\underbrace{\delta_{condition}}_{\text{Align w/ text}} + \underbrace{\delta_{denoising} - \epsilon}_{\text{Diff b/w predicted noise \& added noise}}) \frac{\partial x}{\partial \theta}$$

Disgard this one!

$$s\delta_{condition} = s(\underbrace{\epsilon_{\phi}(z_t|y, t)}_{\text{Conditional}} - \underbrace{\epsilon_{\phi}(z_t|\emptyset, t)}_{\text{Unconditional}})$$

$$\epsilon_{\phi}(z_t|\emptyset, t) = \delta_{domain} + \delta_{denoising} \quad \text{😬 hard to separate}$$

$$\epsilon_{\phi}(z_t|\emptyset, t) - \epsilon_{\phi}(z_t|y = p_{neg}, t) = \delta_{domain} + \cancel{\delta_{denoising}} - (\cancel{\delta_{domain}} + \cancel{\delta_{denoising}} + \cancel{\delta_{condition=p_{neg}}})$$

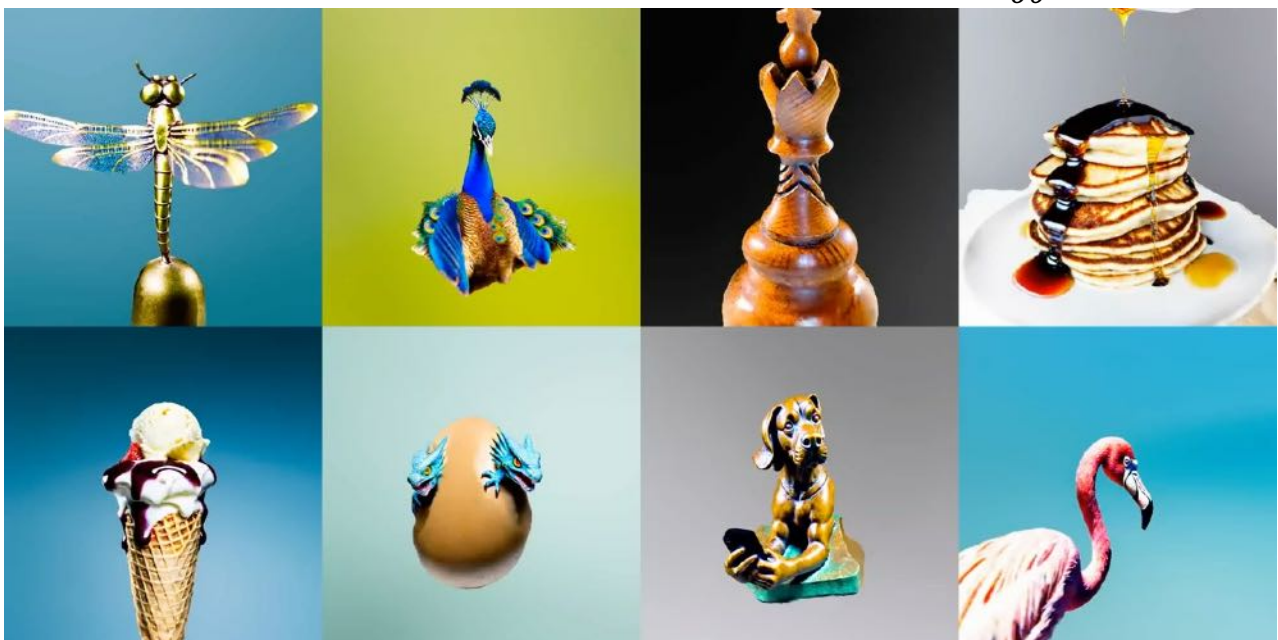
Assumption: $\delta_{condition=p_{neg}} = -\delta_{domain}$

35

3D Content Generation Empowered by 2D Diffusion Priors

NFSD: Noise Free Score Distillation

$$\nabla_{\theta} L_{NFSD}(\phi, g(\theta)) = w(t)(\delta_{domain} + s\delta_{condition}) \frac{\partial x}{\partial \theta}$$



36

3D Content Generation *Empowered by 2D Diffusion Priors*

ProlificDreamer: High-Fidelity and Diverse Text-to-3D Generation with Variational Score Distillation

$$\nabla_{\theta} L_{VSD}(\phi, g(\theta)) = w(t)(\delta_{domain} + s\delta_{condition}) \frac{\partial x}{\partial \theta}$$

$$\nabla_{\theta} L_{VSD}(\phi, g(\theta)) = w(t)(\hat{\epsilon}_{\theta}(z_t|y, t) - \epsilon_{LoRA}(z_t|y, t, c)) \frac{\partial x}{\partial \theta}$$

Conditional model LoRA model,
trained on the rendered images

$$\epsilon_{LoRA}(z_t|y, t, c) = \delta_{denoising} \quad \text{No domain difference this time!}$$

$\nabla_{\theta} L_{SDS}(\phi, g(\theta)) = \mathbb{E}_{t,\epsilon} [w(t)(\hat{\epsilon}_{\theta}(z_t|y, t) - \epsilon) \frac{\partial x}{\partial \theta}]$
Estimated noise with classifier free guidance

$\hat{\epsilon}_{\theta}(z_t|y, t) = \epsilon_{\phi}(z_t|\emptyset, t) + s\delta_{condition}$
 $\epsilon_{\phi}(z_t|\emptyset, t) = \delta_{domain} + \delta_{denoising}$

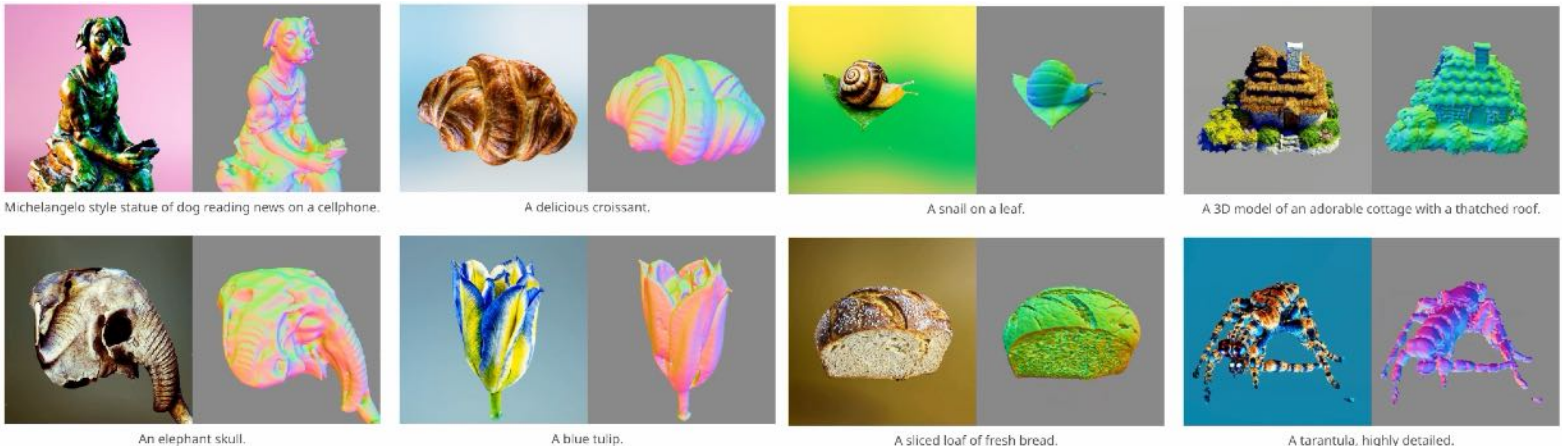
x_D $\delta_{denoising}$ x_{OOD} δ_{domain} $x_{OOD} + \delta_{domain}$

$$\hat{\epsilon}_{\theta}(z_t|y, t) - \epsilon_{LoRA}(z_t|y, t, c) = \delta_{domain} + \cancel{\delta_{denoising}} + s\delta_{condition} - \cancel{\delta_{denoising}}$$

3D Content Generation *Empowered by 2D Diffusion Priors*

ProlificDreamer: High-Fidelity and Diverse Text-to-3D Generation with Variational Score Distillation

$$\nabla_{\theta} L_{VSD}(\phi, g(\theta)) = w(t)(\delta_{domain} + s\delta_{condition}) \frac{\partial x}{\partial \theta}$$



3D Content Generation *Empowered by 2D Diffusion Priors*

Janus (multi-face) problem



3D Content Generation *Empowered by 2D Diffusion Priors*

Janus (multi-face) problem

Dalle-2

Stable Diffusion v2

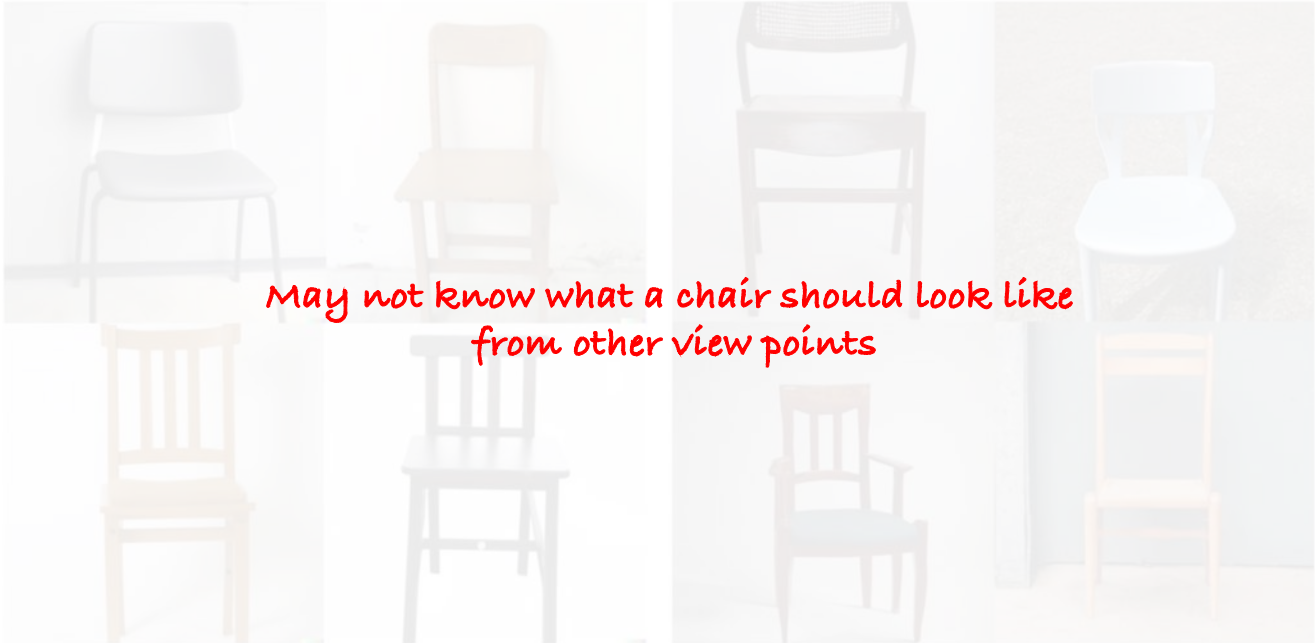


3D Content Generation *Empowered by 2D Diffusion Priors*

Janus (multi-face) problem

Dalle-2

Stable Diffusion v2



41

3D Content Generation *Learn from 3D data*

Janus (multi-face) problem

How to use 3D data?



A Universe of Annotated 3D Objects

Objaverse 1.0 is a Massive Dataset with 800K+ Annotated 3D Objects

Objaverse-XL

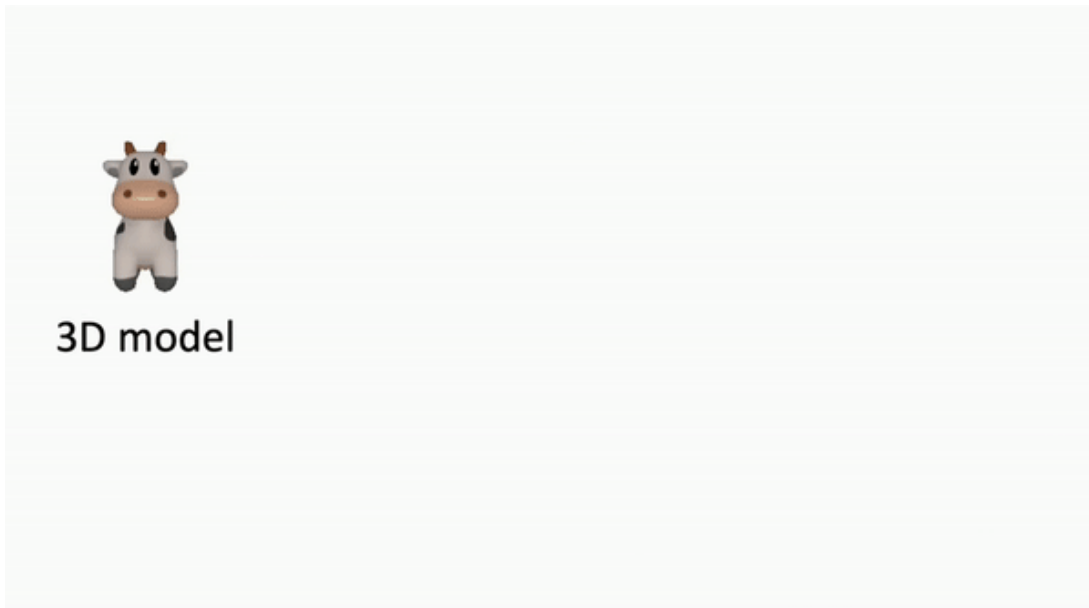
A Universe of 10M+ 3D Objects

 Image dataset: 5B

42

3D Content Generation *Learn from 3D data*

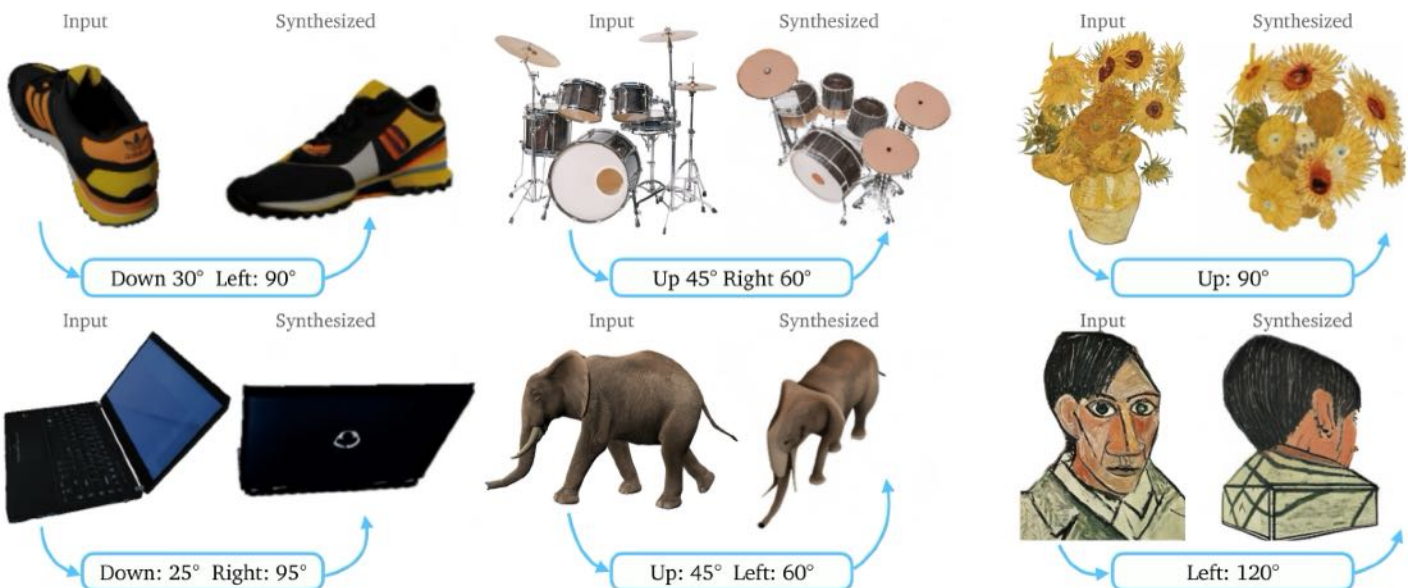
Use 3D models to get multi-view images



43

3D Content Generation *Learn from 3D data*

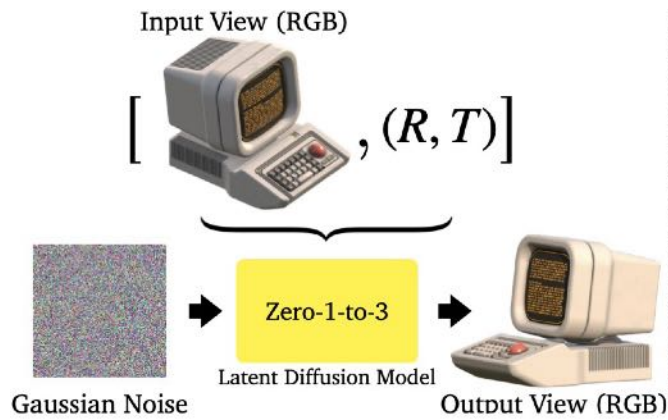
Zero-1-to-3: Zero-shot One Image to 3D Object



44

3D Content Generation *Learn from 3D data*

Zero-1-to-3: Zero-shot One Image to 3D Object

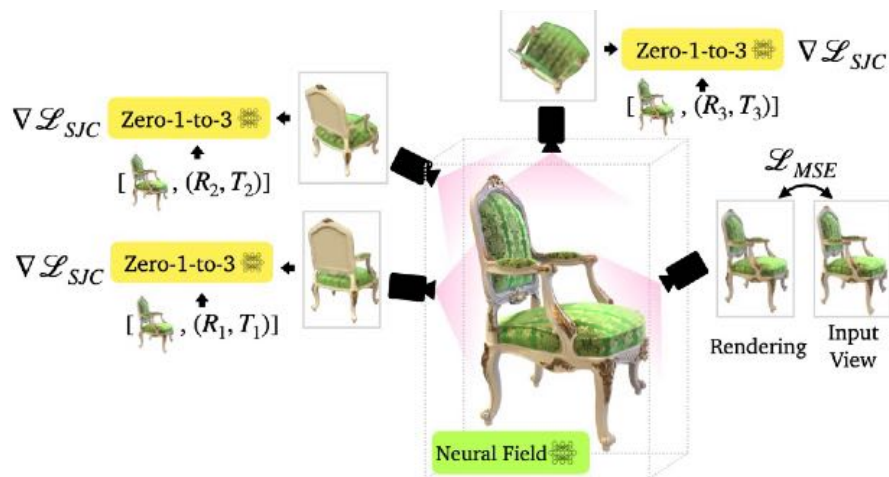


Idea: Finetune a 2D diffusion model to generate novel views (i.e. condition on camera pose)

45

3D Content Generation *Learn from 3D data*

Zero-1-to-3: Zero-shot One Image to 3D Object



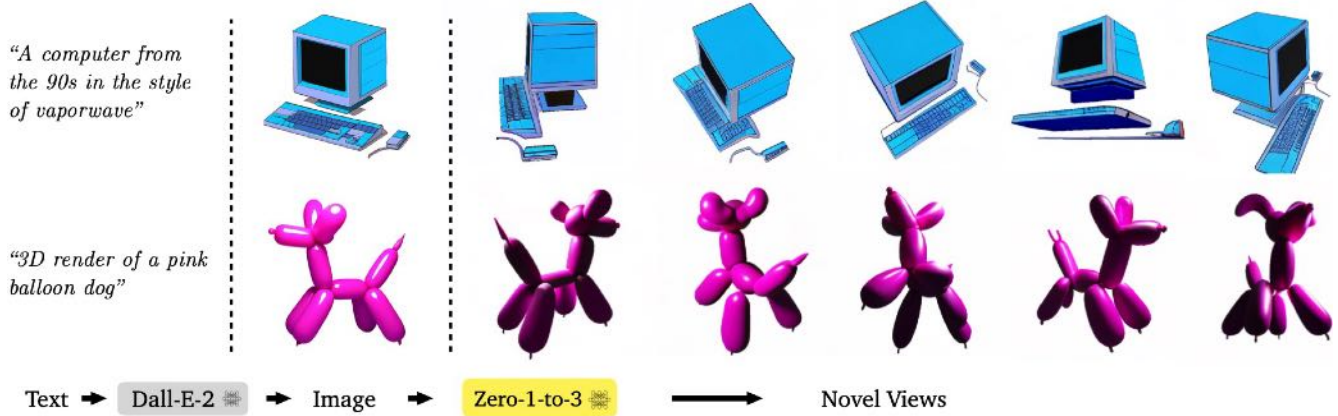
The trained model can be used for (single-view) 3D reconstruction

46

3D Content Generation *Learn from 3D data*

Zero-1-to-3: Zero-shot One Image to 3D Object

Text2img2NVS



47

3D Content Generation *Learn from 3D data*

Zero-1-to-3: Zero-shot One Image to 3D Object

Single view 3D reconstruction



48

3D Content Generation *Learn from 3D data*

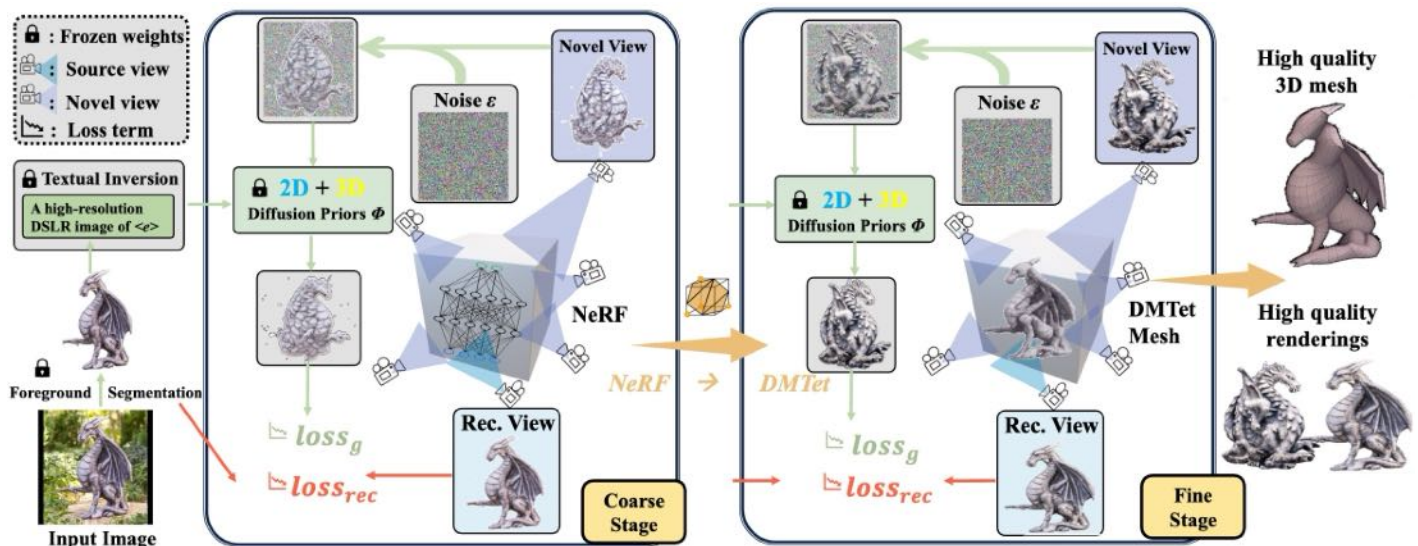
But 3D prior only is blurry, 2D prior only lacks geometry...



49

3D Content Generation *Learn from 3D data*

Magic123: One Image to High-Quality 3D Object Generation Using Both 2D and 3D Diffusion Priors

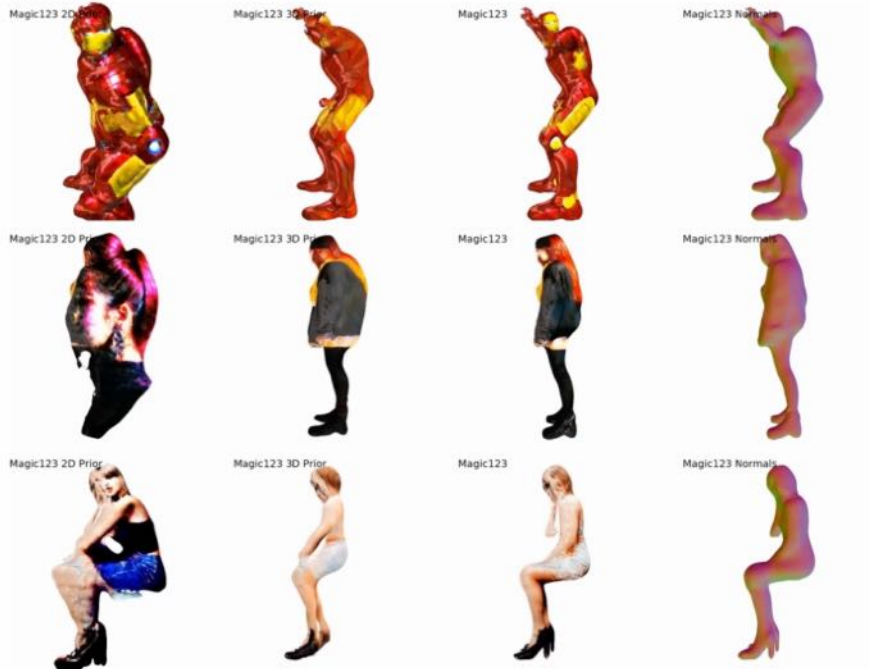


Idea: Combine both 2D and 3D diffusion priors

50

3D Content Generation *Learn from 3D data*

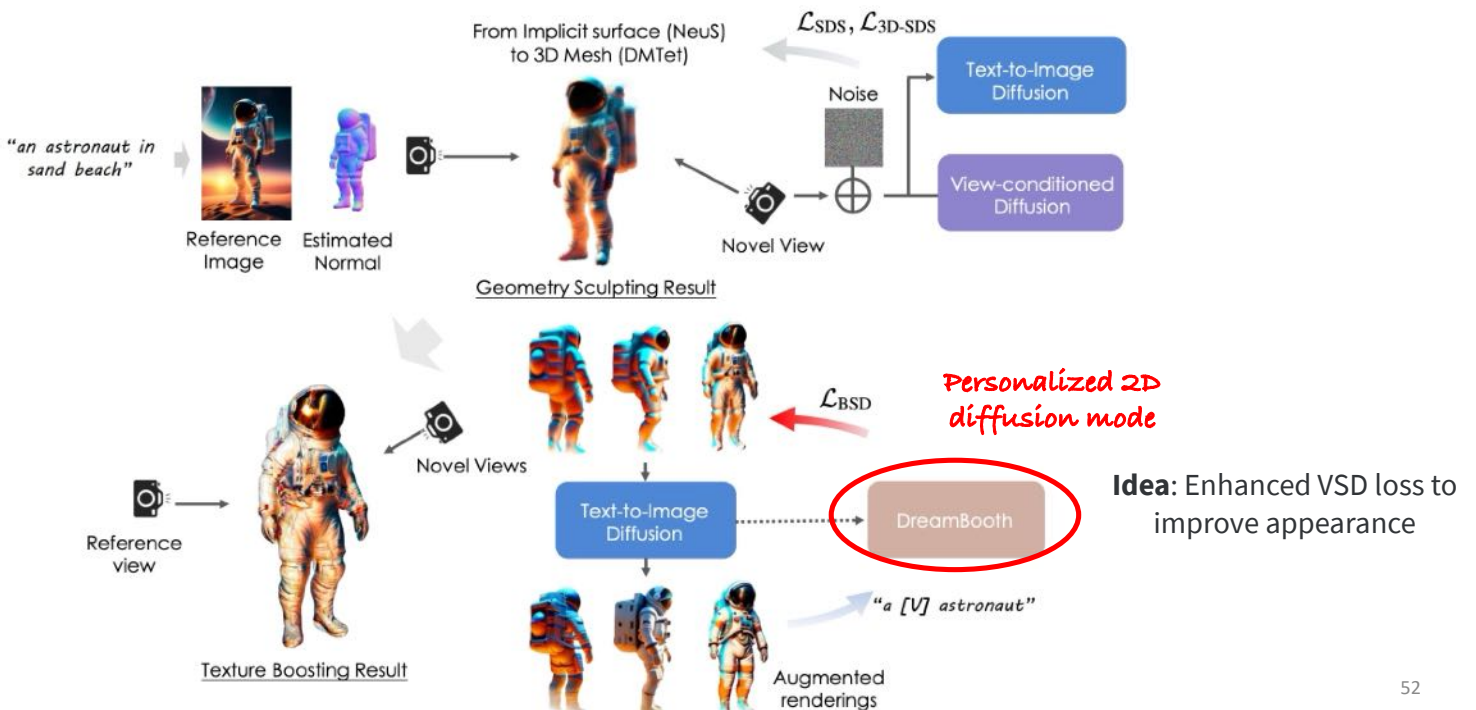
Magic123: One Image to High-Quality 3D Object Generation Using Both 2D and 3D Diffusion Priors



51

3D Content Generation *Learn from 3D data*

DreamCraft3D: Hierarchical 3D Generation with Bootstrapped Diffusion Prior



52

3D Content Generation *Learn from 3D data*

DreamCraft3D: Hierarchical 3D Generation with Bootstrapped Diffusion Prior



53

3D Content Generation *Learn from 3D data*

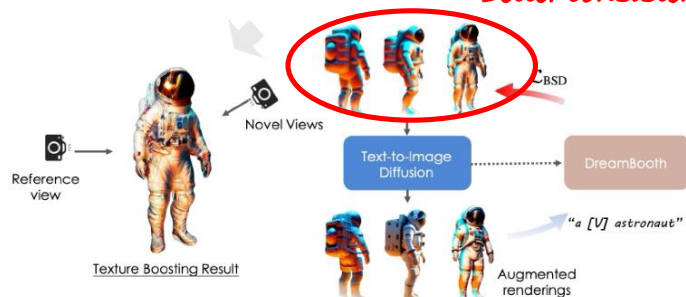
DreamCraft3D: Hierarchical 3D Generation with Bootstrapped Diffusion Prior

$$\text{Recall: } \nabla_{\theta} L_{\text{VSD}}(\phi, g(\theta)) = w(t)(\hat{\epsilon}_{\theta}(z_t|y, t, z) - \epsilon_{\text{LoRA}}(z_t|y, t, c, z)) \frac{\partial x}{\partial \theta}$$

$$\nabla_{\theta} L_{\text{BSD}}(\phi, g(\theta)) = w(t)(\epsilon_{\text{DreamBooth}}(z_t|y, t, r_{t'}(z), v) - \epsilon_{\text{LoRA}}(z_t|y, t, z, v)) \frac{\partial x}{\partial \theta}$$

$r_{t'}(z)$: "Augmented" image renderings
Restore a heavily noised image
Reveal high-frequency details but sacrifice fidelity

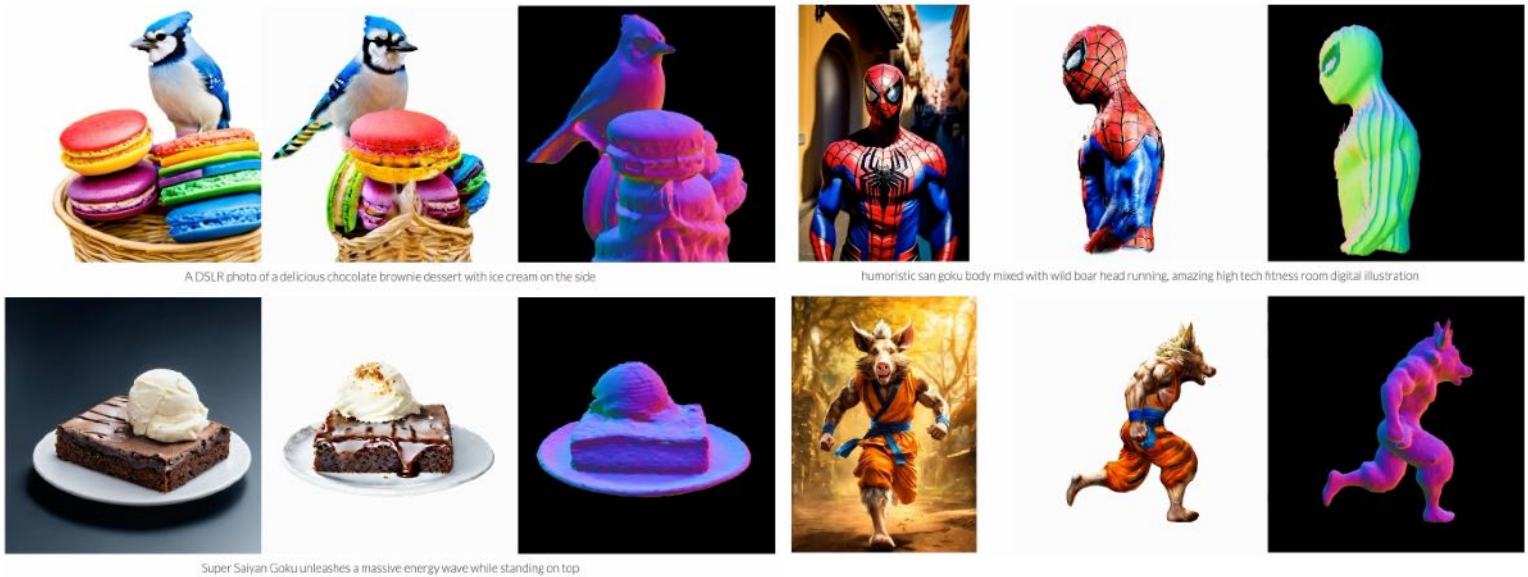
Better quality
Identity preserved
Better consistency



54

3D Content Generation *Learn from 3D data*

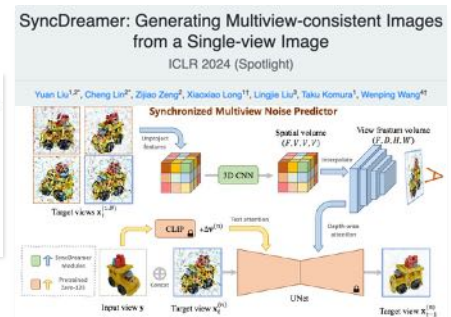
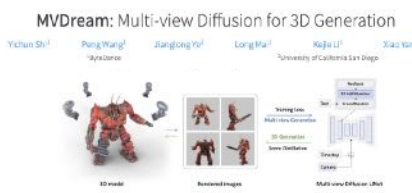
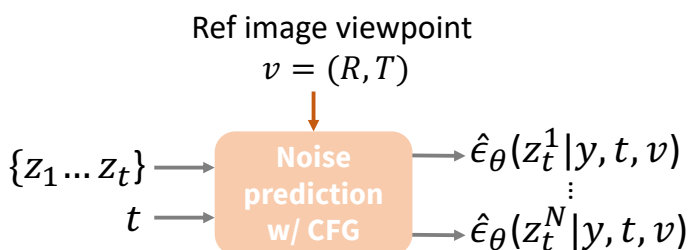
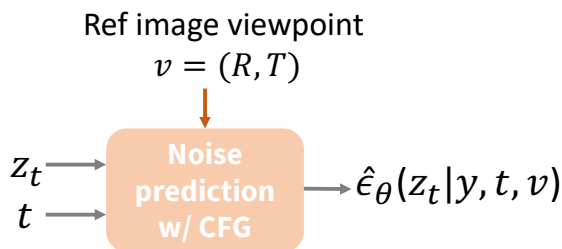
DreamCraft3D: Hierarchical 3D Generation with Bootstrapped Diffusion Prior



55

3D Content Generation *Learn from 3D data*

Single viewpoint prediction to multi viewpoint prediction

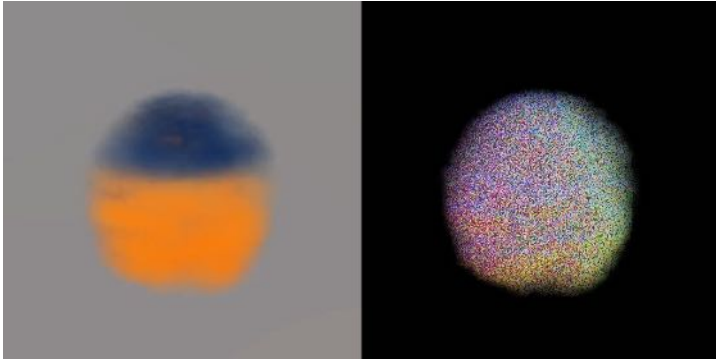


56

3D Content Generation *Learn from 3D data*

Single viewpoint prediction to multi viewpoint prediction

MVDream



SyncDreamer



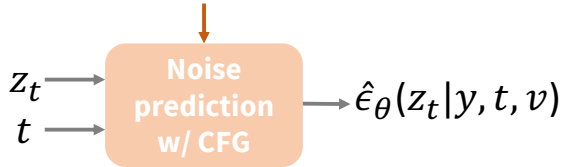
57

3D Content Generation *Learn from 3D data*

RGB + Geometry

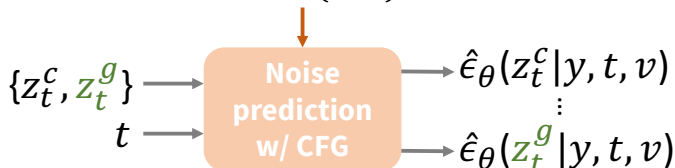
Ref image viewpoint

$$v = (R, T)$$



Ref image viewpoint

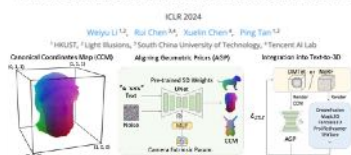
$$v = (R, T)$$



Zero-1-to-3: Zero-shot One Image to 3D Object



SweetDreamer: Aligning Geometric Priors in 2D Diffusion for Consistent Text-to-3D



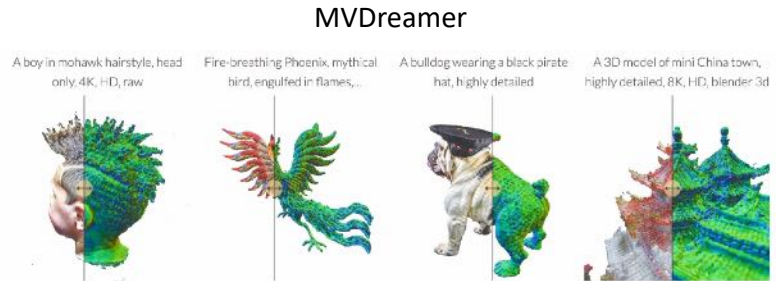
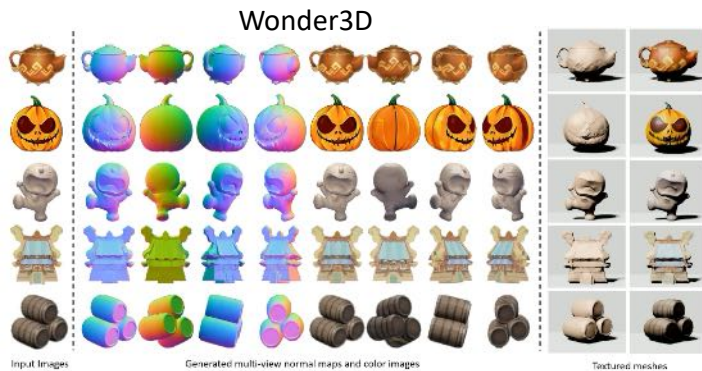
Wonder3D: Single Image to 3D using Cross-Domain Diffusion CVPR 2024 Highlight



58

3D Content Generation *Learn from 3D data*

RGB + Geometry



59

3D Content Generation *Learn from 3D data*

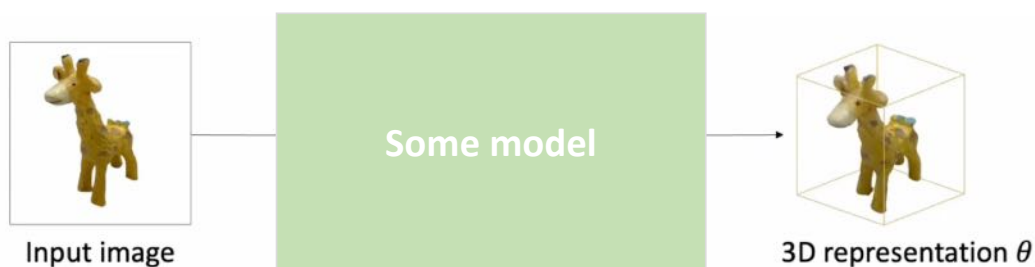
More direct approaches?

Optimize 3D representations

- 1) ref view match the input image
- 2) Novel views are photorealistic and view-consistent

But is time consuming

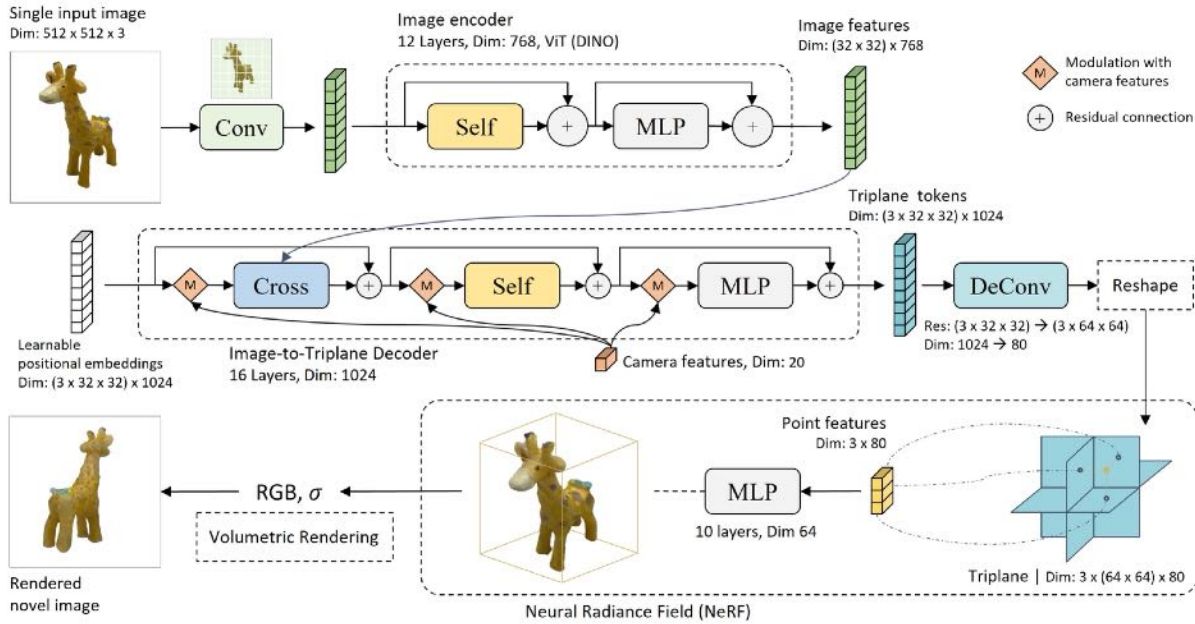
A direct inference approach?



60

3D Content Generation *Learn from 3D data*

LRM: LARGE RECONSTRUCTION MODEL FOR SINGLE IMAGE TO 3D



61

3D Content Generation *Learn from 3D data*

LRM: LARGE RECONSTRUCTION MODEL FOR SINGLE IMAGE TO 3D

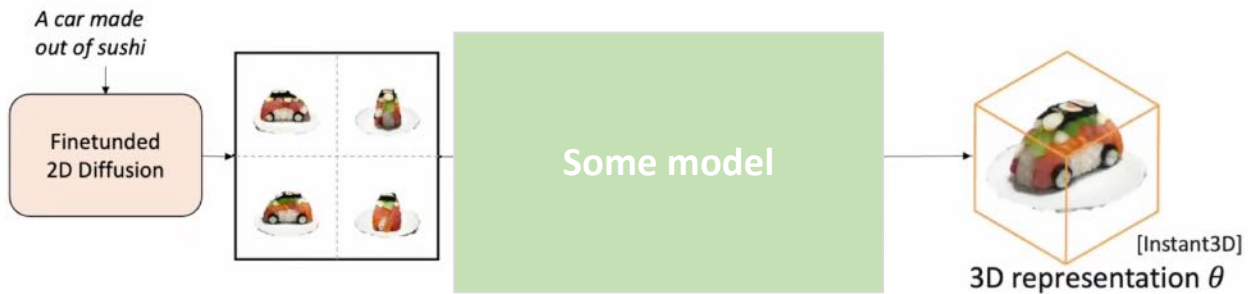


3D Content Generation *Learn from 3D data*

LRM: LARGE RECONSTRUCTION MODEL FOR SINGLE IMAGE TO 3D



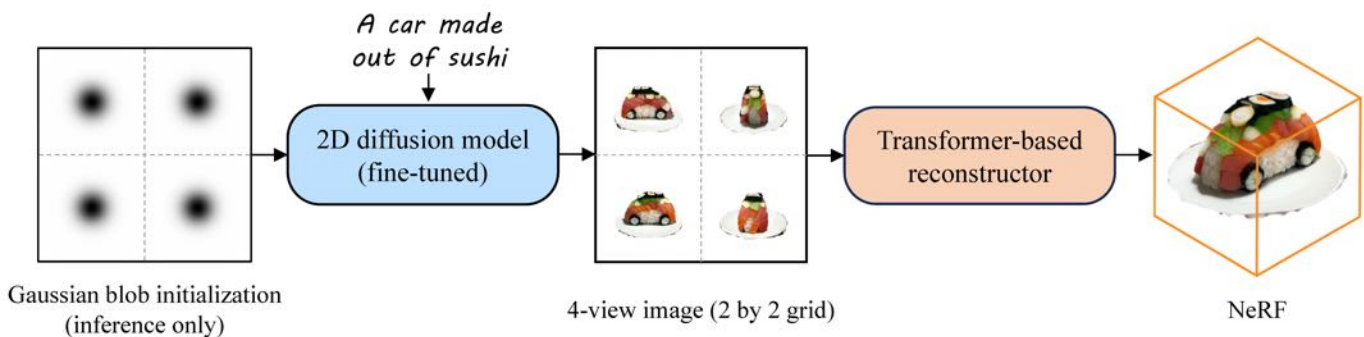
INSTANT3D: FAST TEXT-TO-3D WITH SPARSE-VIEW GENERATION AND LARGE RECONSTRUCTION MODEL



63

3D Content Generation *Learn from 3D data*

INSTANT3D: FAST TEXT-TO-3D WITH SPARSE-VIEW GENERATION AND LARGE RECONSTRUCTION MODEL

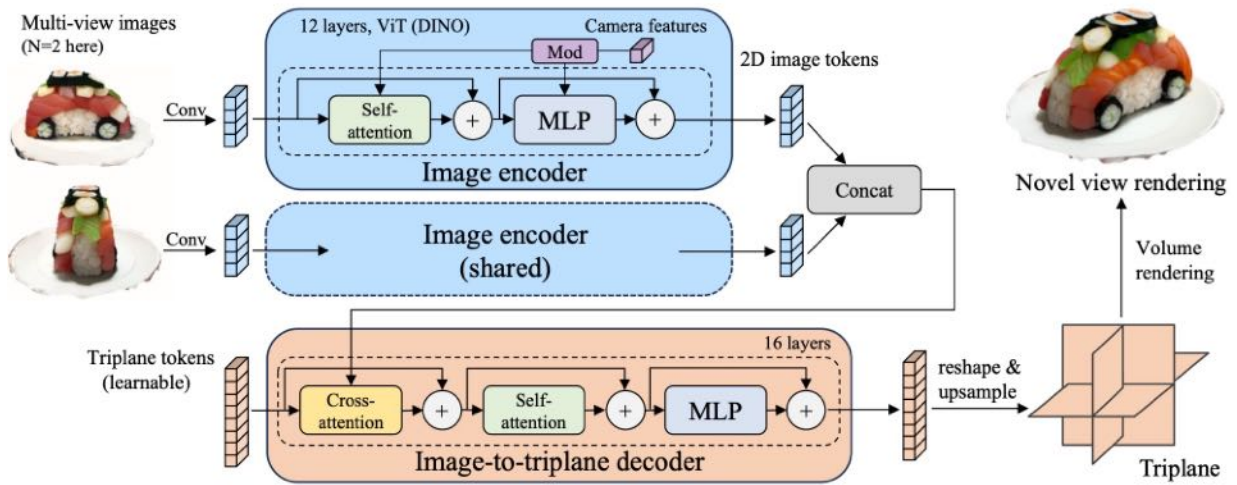


Idea: multi-view 2D diffusion + sparse view reconstruction

64

3D Content Generation *Learn from 3D data*

INSTANT3D: FAST TEXT-TO-3D WITH SPARSE-VIEW GENERATION AND LARGE RECONSTRUCTION MODEL



3D Content Generation *Learn from 3D data*

INSTANT3D: FAST TEXT-TO-3D WITH SPARSE-VIEW GENERATION AND LARGE RECONSTRUCTION MODEL

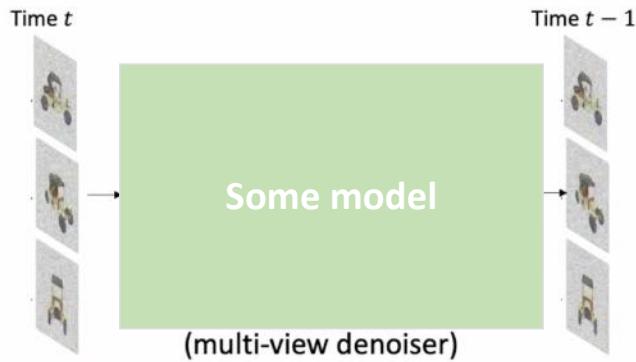


3D Content Generation *Learn from 3D data*

INSTANT3D: FAST TEXT-TO-3D WITH SPARSE-VIEW GENERATION AND LARGE RECONSTRUCTION MODEL



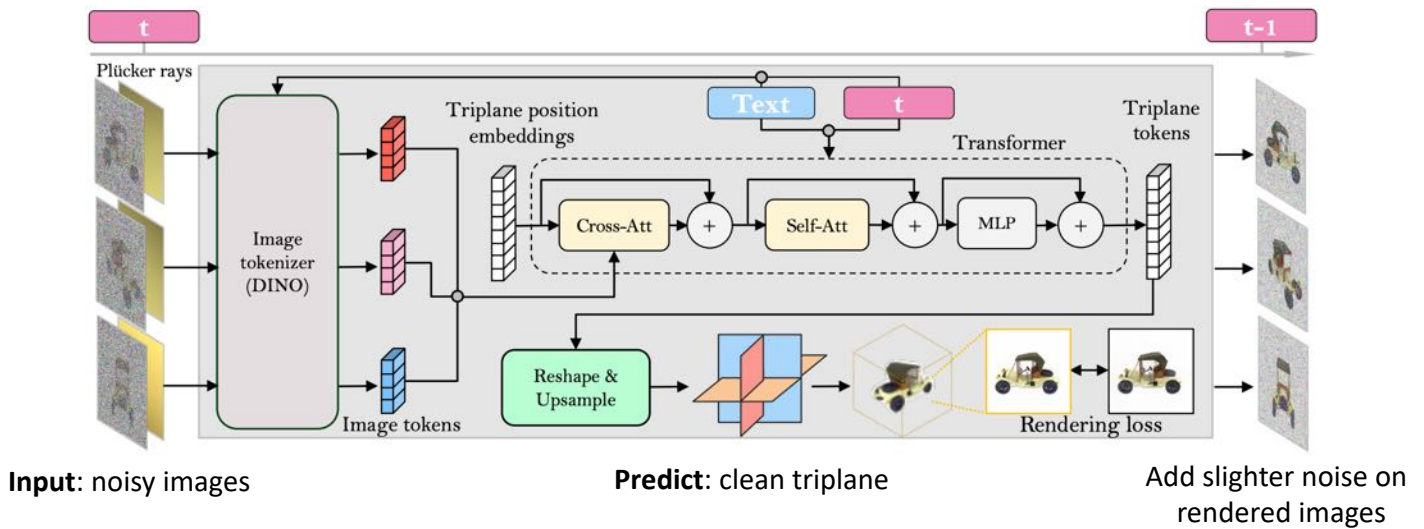
DMV3D:DENOISING MULTI-VIEW DIFFUSION USING 3D LARGE RECONSTRUCTION MODEL



67

3D Content Generation *Learn from 3D data*

DMV3D:DENOISING MULTI-VIEW DIFFUSION USING 3D LARGE RECONSTRUCTION MODEL



68

3D Content Generation *Learn from 3D data*

DMV3D:DENOISING MULTI-VIEW DIFFUSION USING 3D LARGE RECONSTRUCTION MODEL



69

Conclusion

2D priors with Score Distillation Sampling

- Higher resolution
- Richer appearance
- Single-view to 3D
- Photorealistic appearance

3D priors

- View-conditioned diffusion
- Multi-view diffusion
- View-conditioned geometry + appearance diffusion

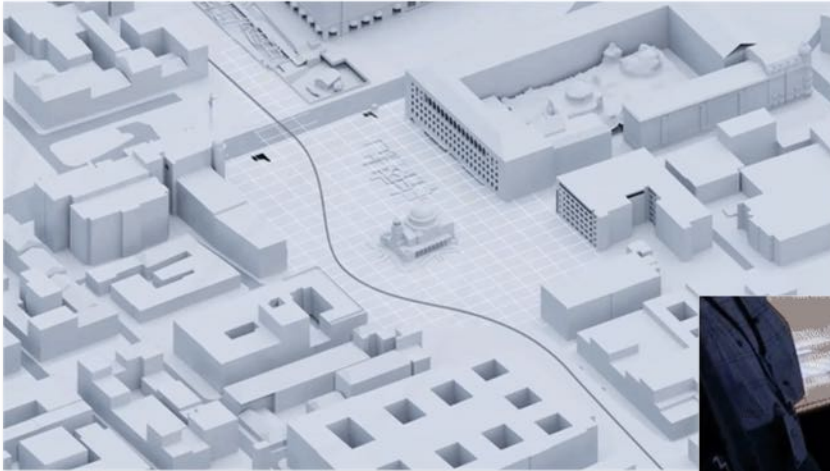
Feed-forward models (empowered by data + transformer)

- Single-image to 3D
- Multi-view to 3D
- Multi-view diffusion

Sometimes we want to manipulate existing scenes...

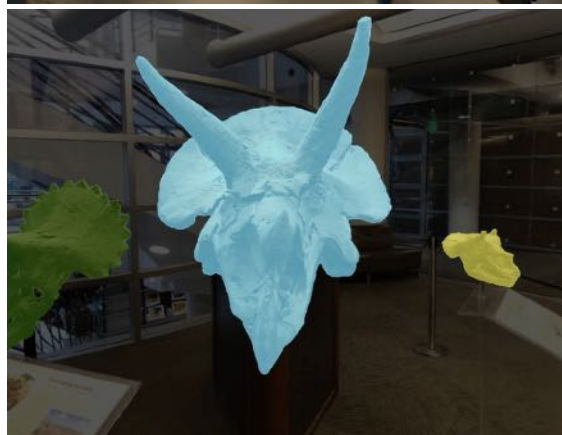
70

Scene Manipulation/Editing/Generation



Manipulate/Edit

Interpretable?



Manipulate/Edit

Interpretable?



Flexible?

Interactive editing

NSVF



Novel View Synthesis



Editable Scene Rendering

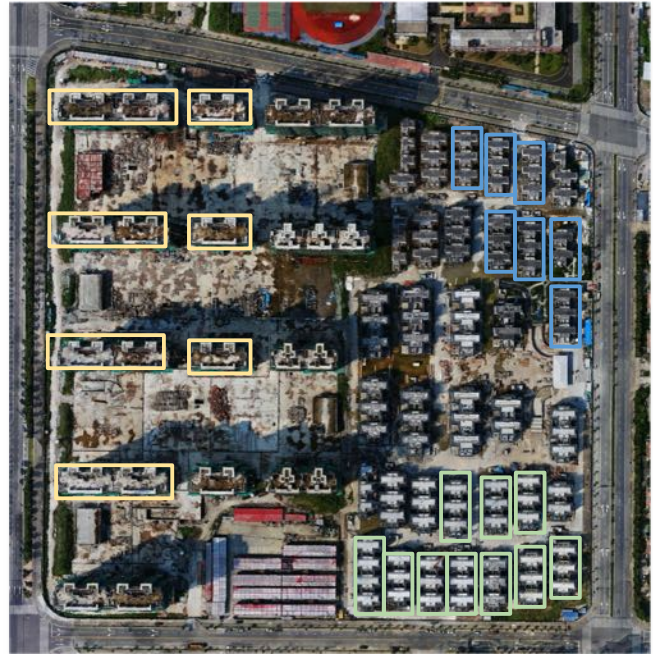
Liu, L., Gu, J., Lin, K.Z., Chua, T., & Theobalt, C. (2020). *Neural Sparse Voxel Fields*

Yang, B., Zhang, Y., Xu, Y., Li, Y., Zhou, H., Bao, H., Zhang, G., & Cui, Z. (2021). *Learning Object-Compositional Neural Radiance Field for Editable Scene Rendering.*

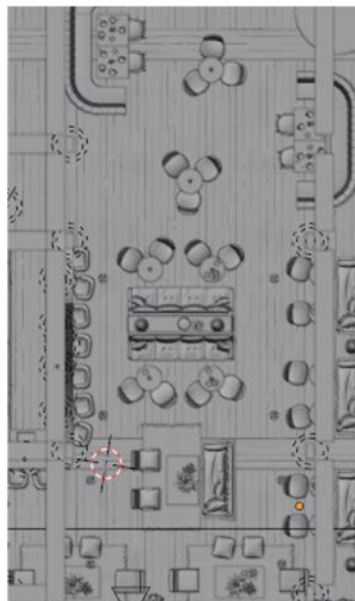
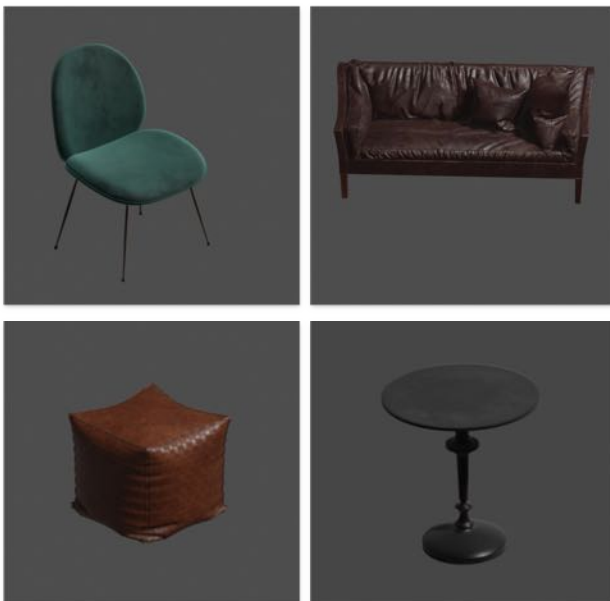
Scalable?



Urban Fabric



Interior Design

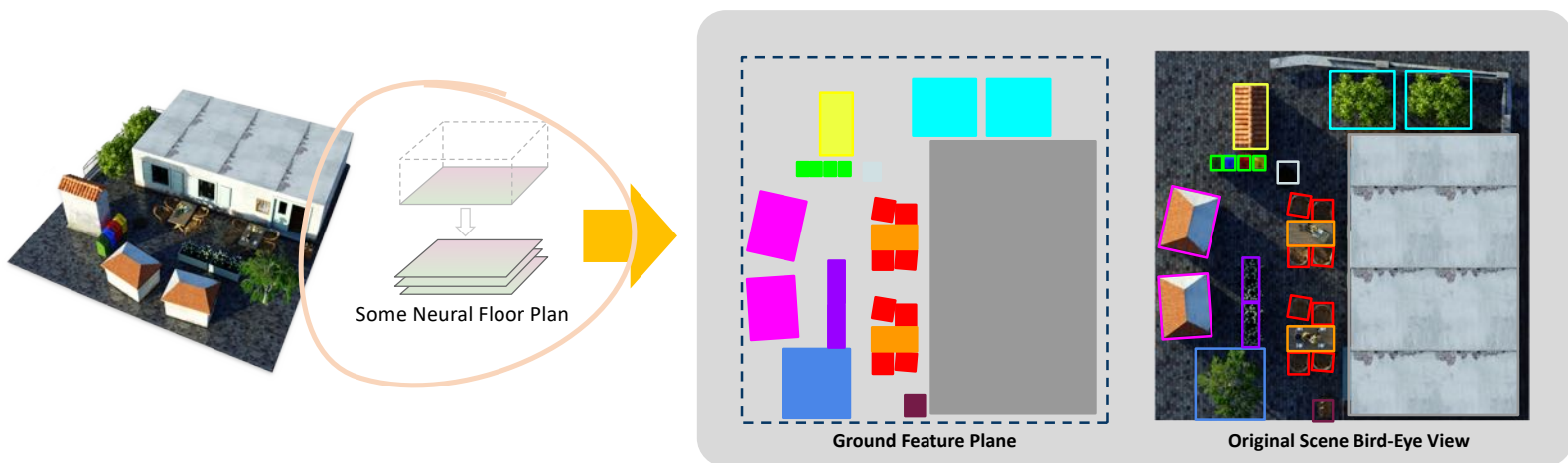


Extract Assets and Layout



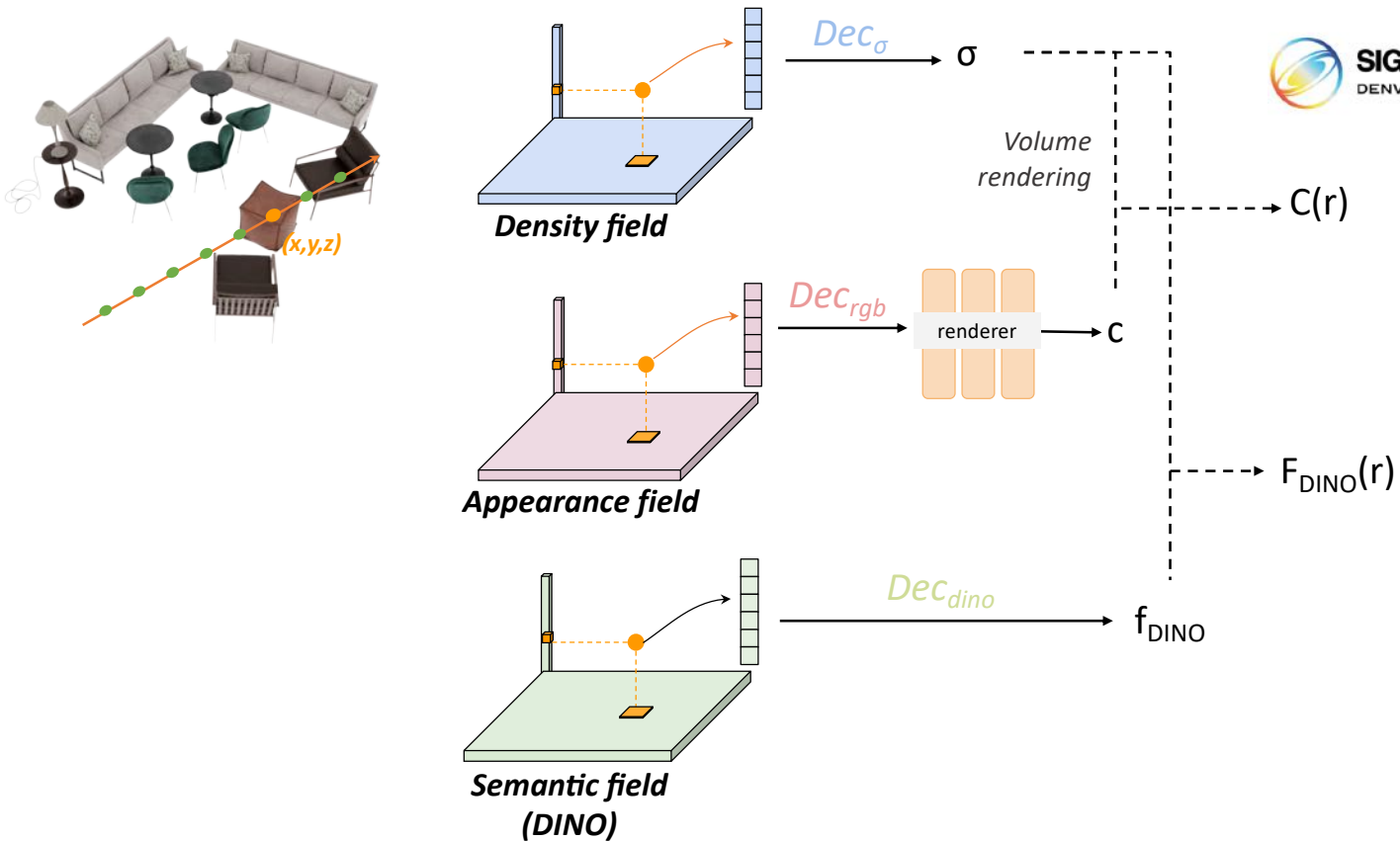
77

🤔 Can we find and categorize objects on this floor plan?

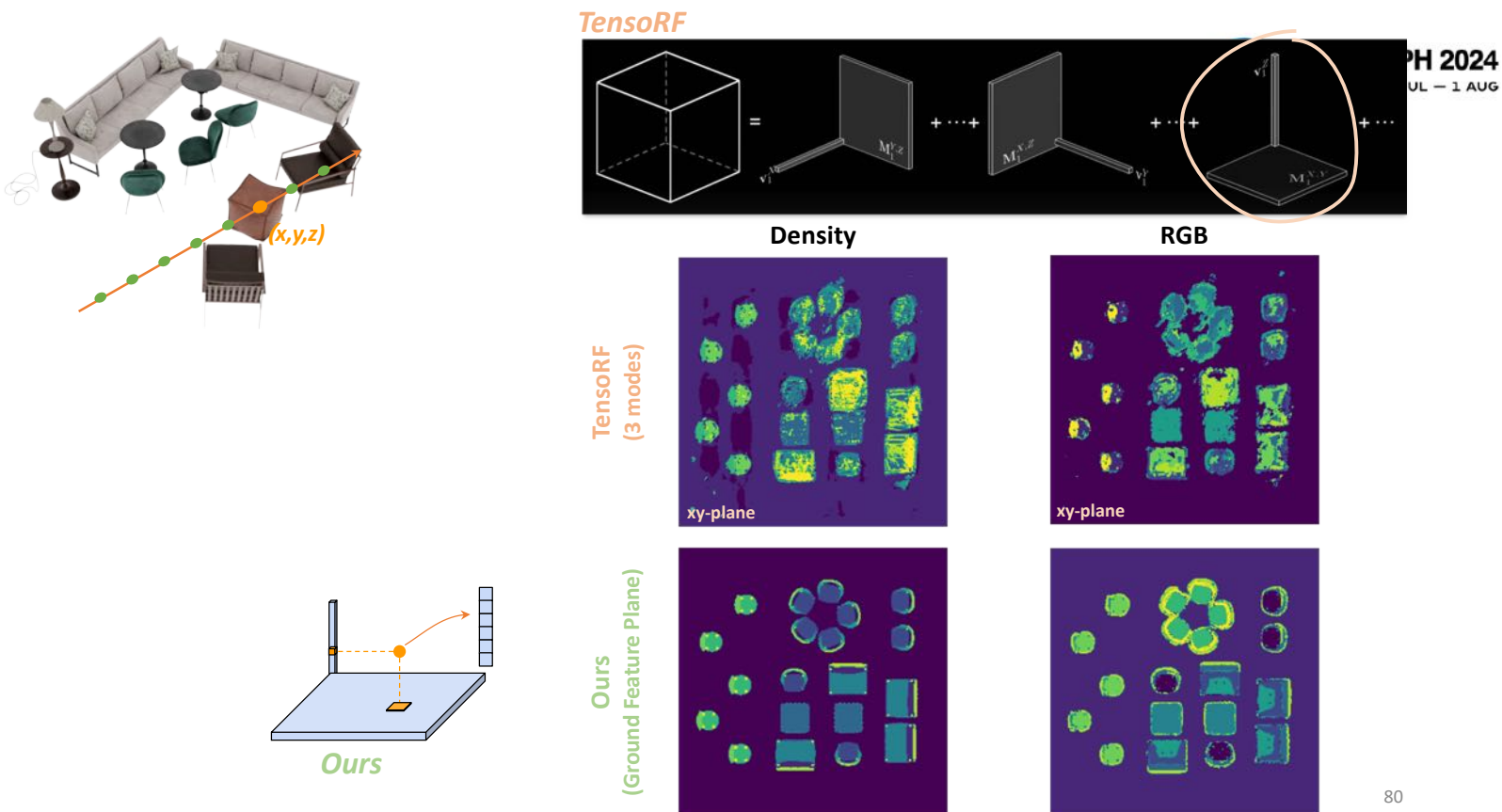


Clean Sharp Object clues

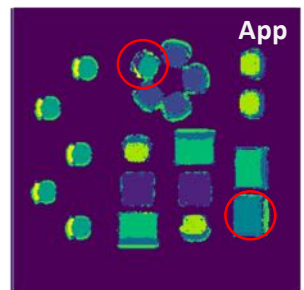
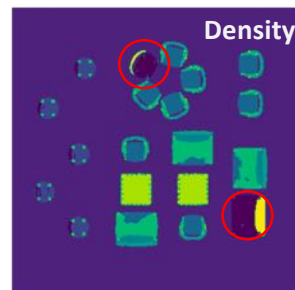
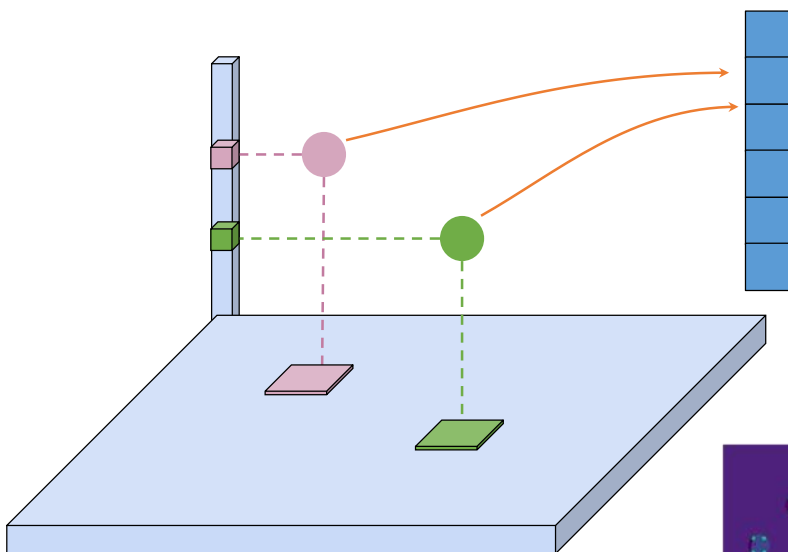
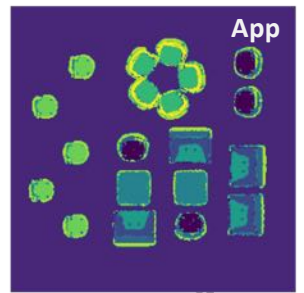
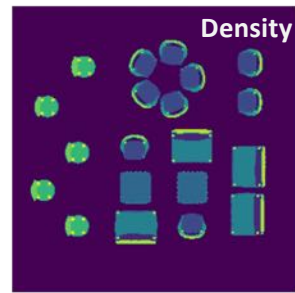
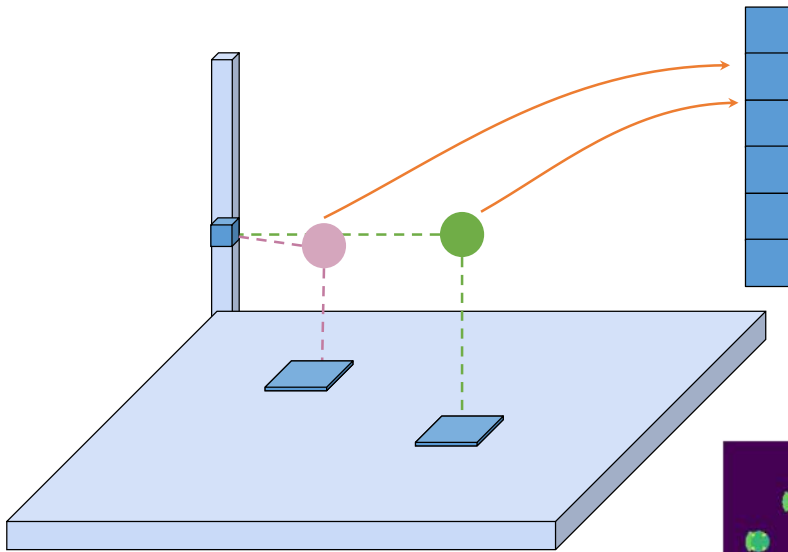
78



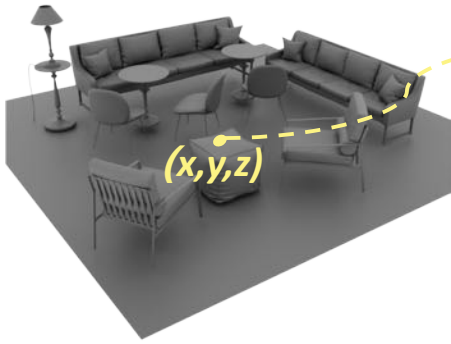
79



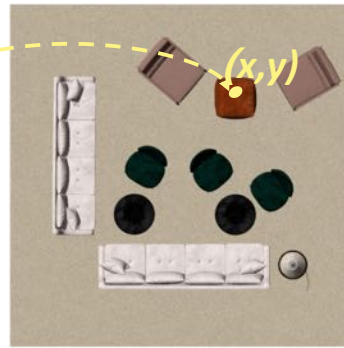
80



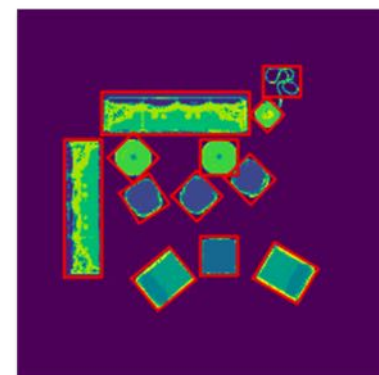
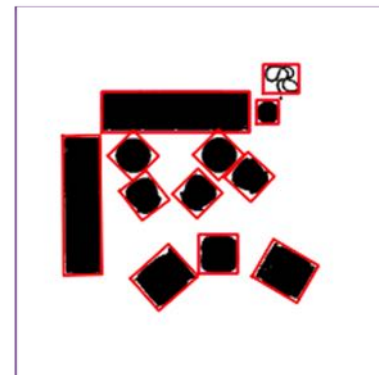
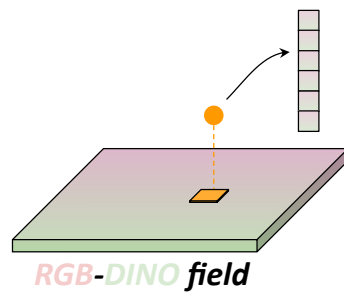
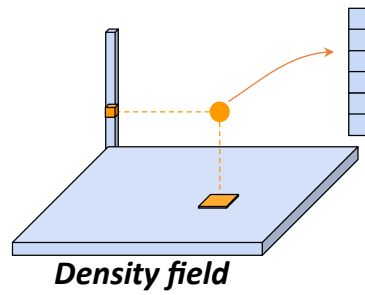
3D density feature

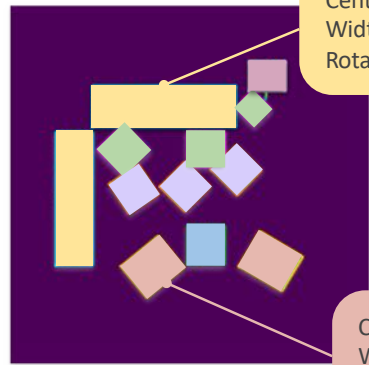
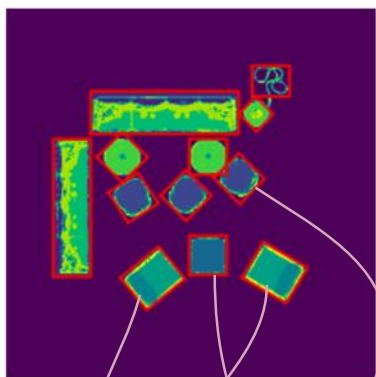
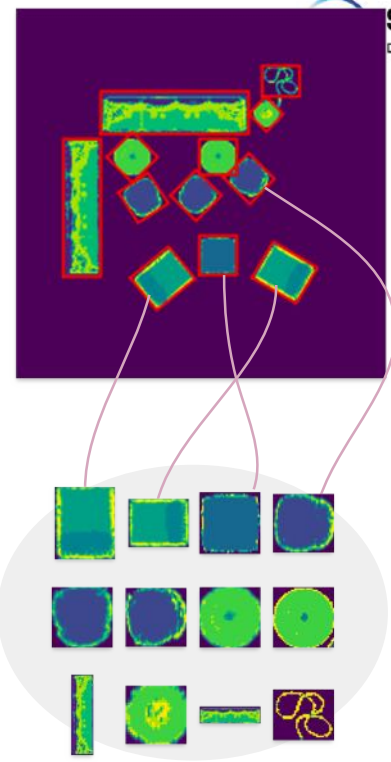
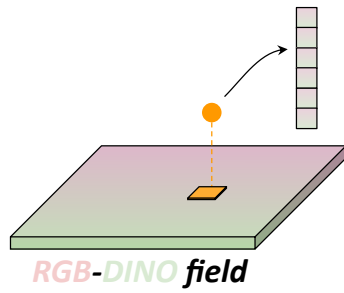


2D RGB-DINO plane feature



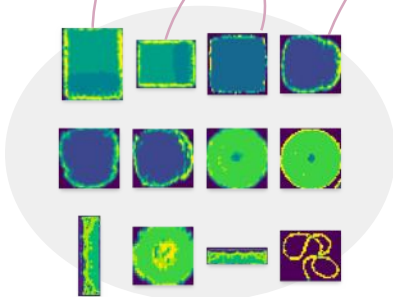
3D RGB-DINO feature



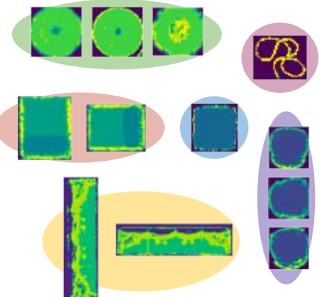


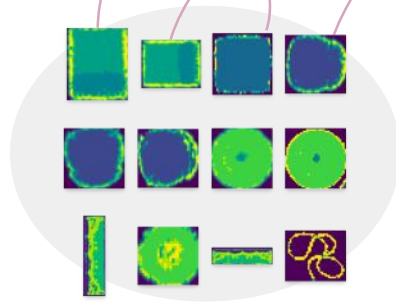
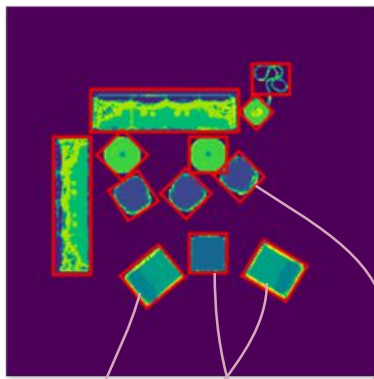
Center: (161,105)
Width/Height: (150, 43)
Rotation (from x-axis): 0.0

Center: (149,275)
Width/Height: (42, 50)
Rotation (from x-axis): 50.6

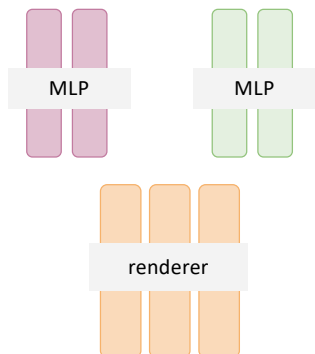
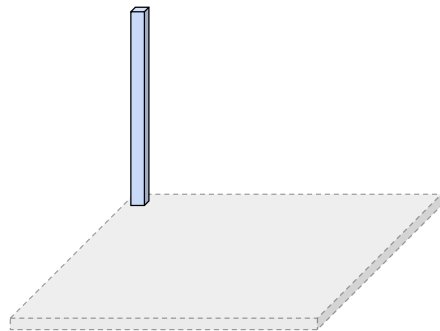
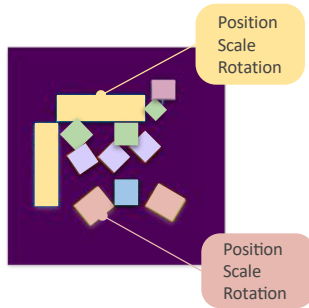
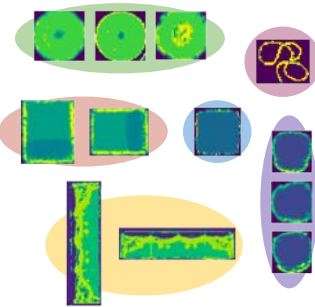


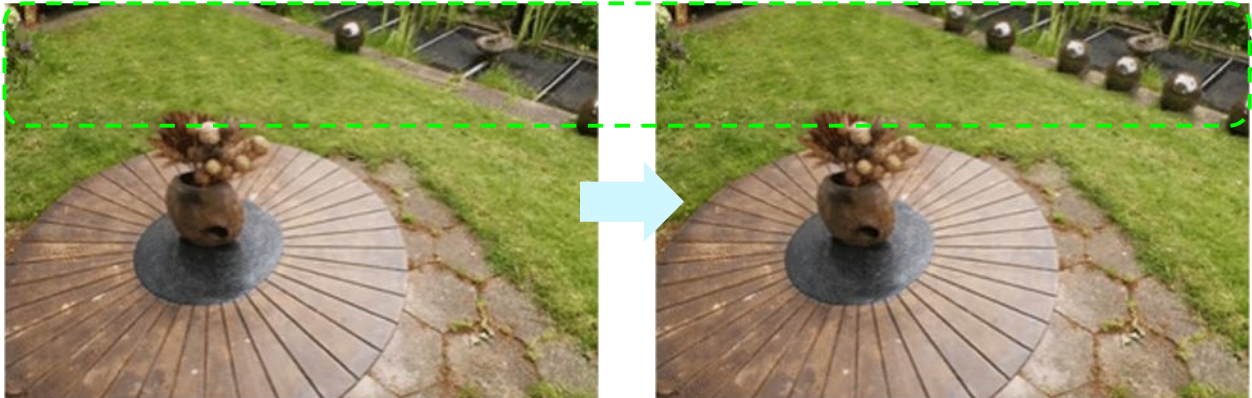
Clustering



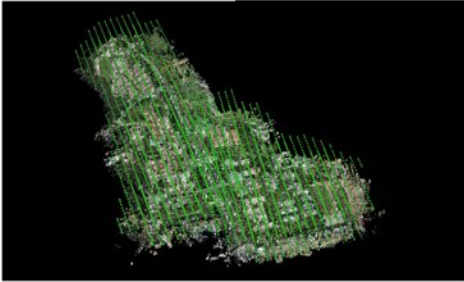
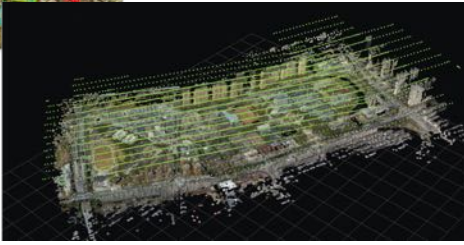


Clustering

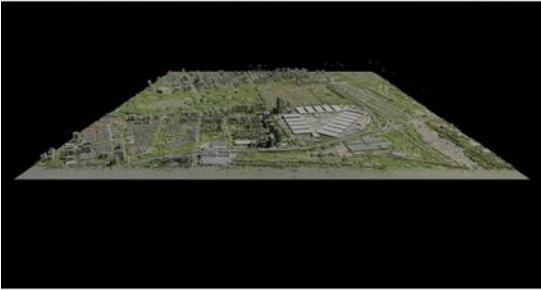
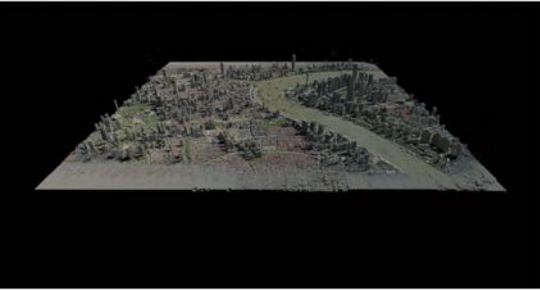




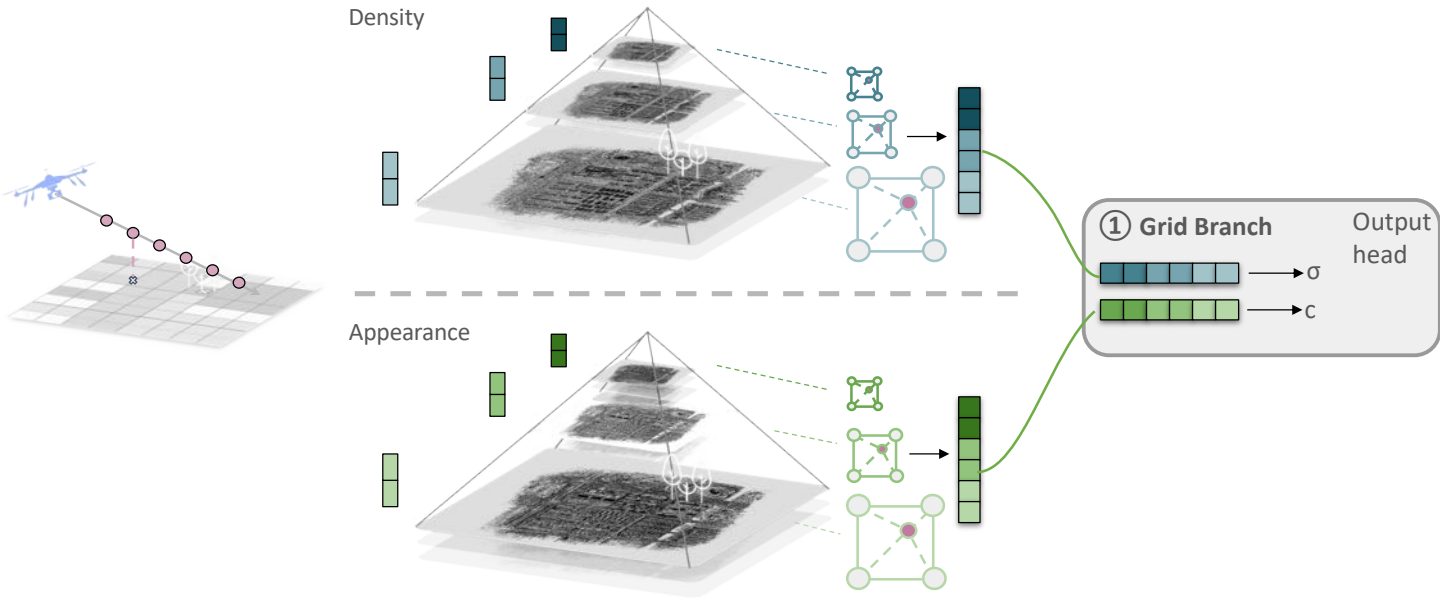
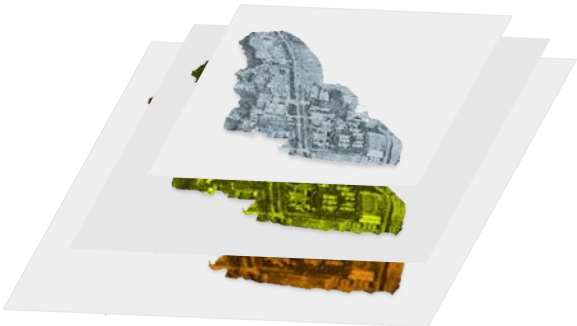
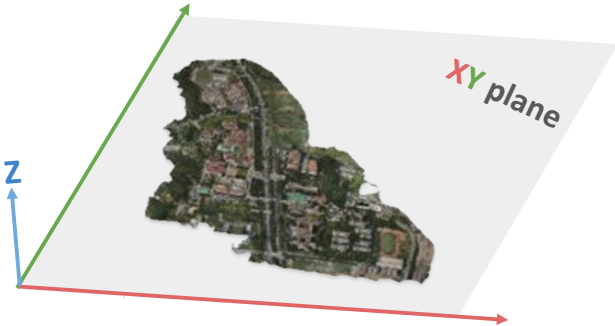
Oblique Photography

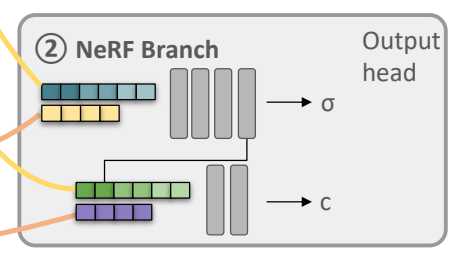
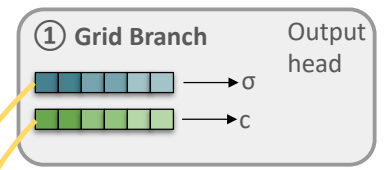
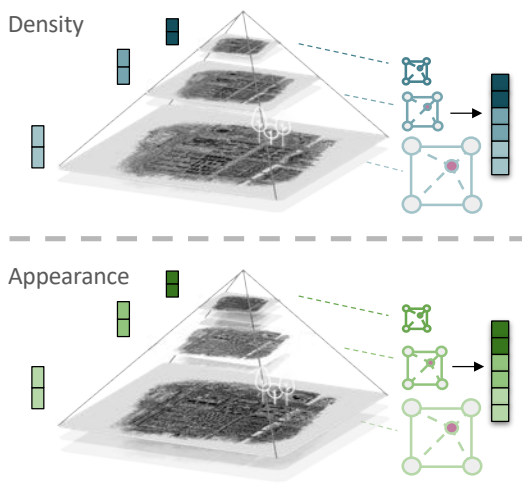
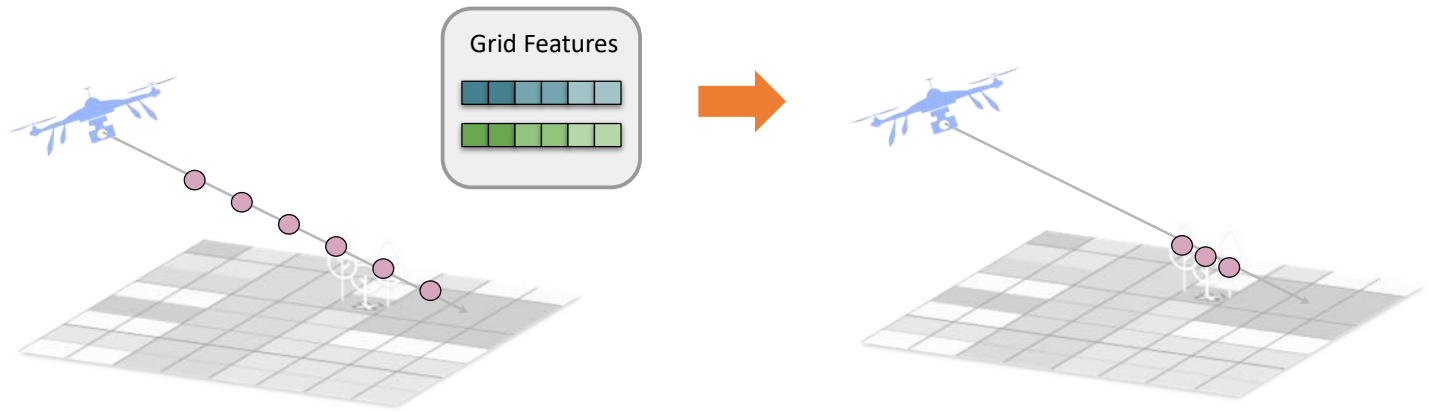


1km²~10k images
1 image~50 megapixels



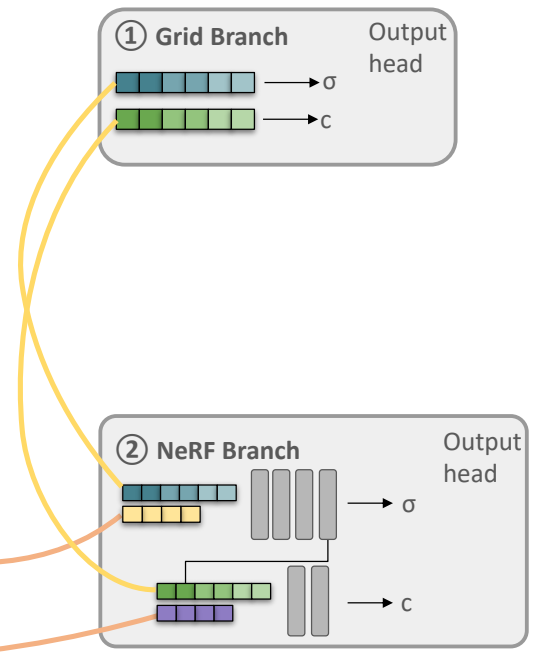
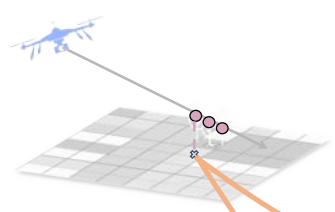
Plane-Axis Factorization

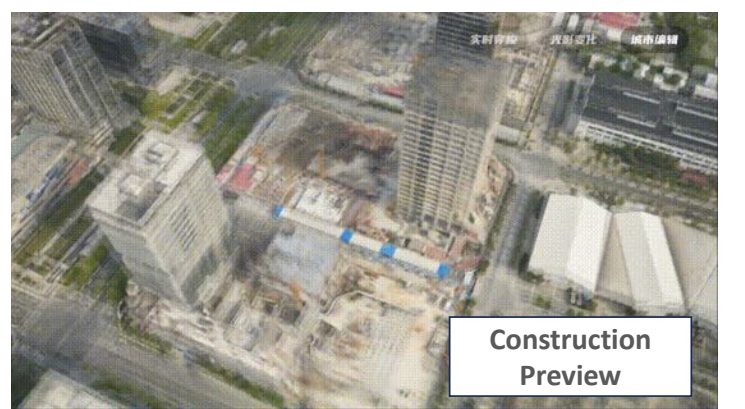
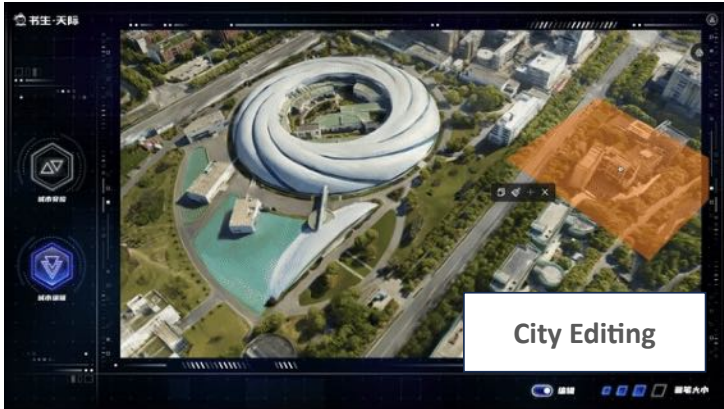
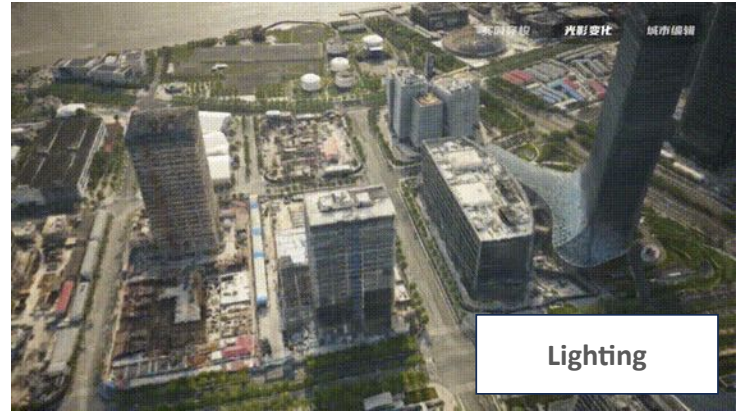




Fourier Embedding (position)
 $(\sin(x), \cos(x), \dots, \sin(2^{l-1}x), \cos(2^{l-1}x))$

Fourier Embedding (dir)
 $(\sin(x), \cos(x), \dots, \sin(2^{l-1}x), \cos(2^{l-1}x))$





Creative Scene Editing

