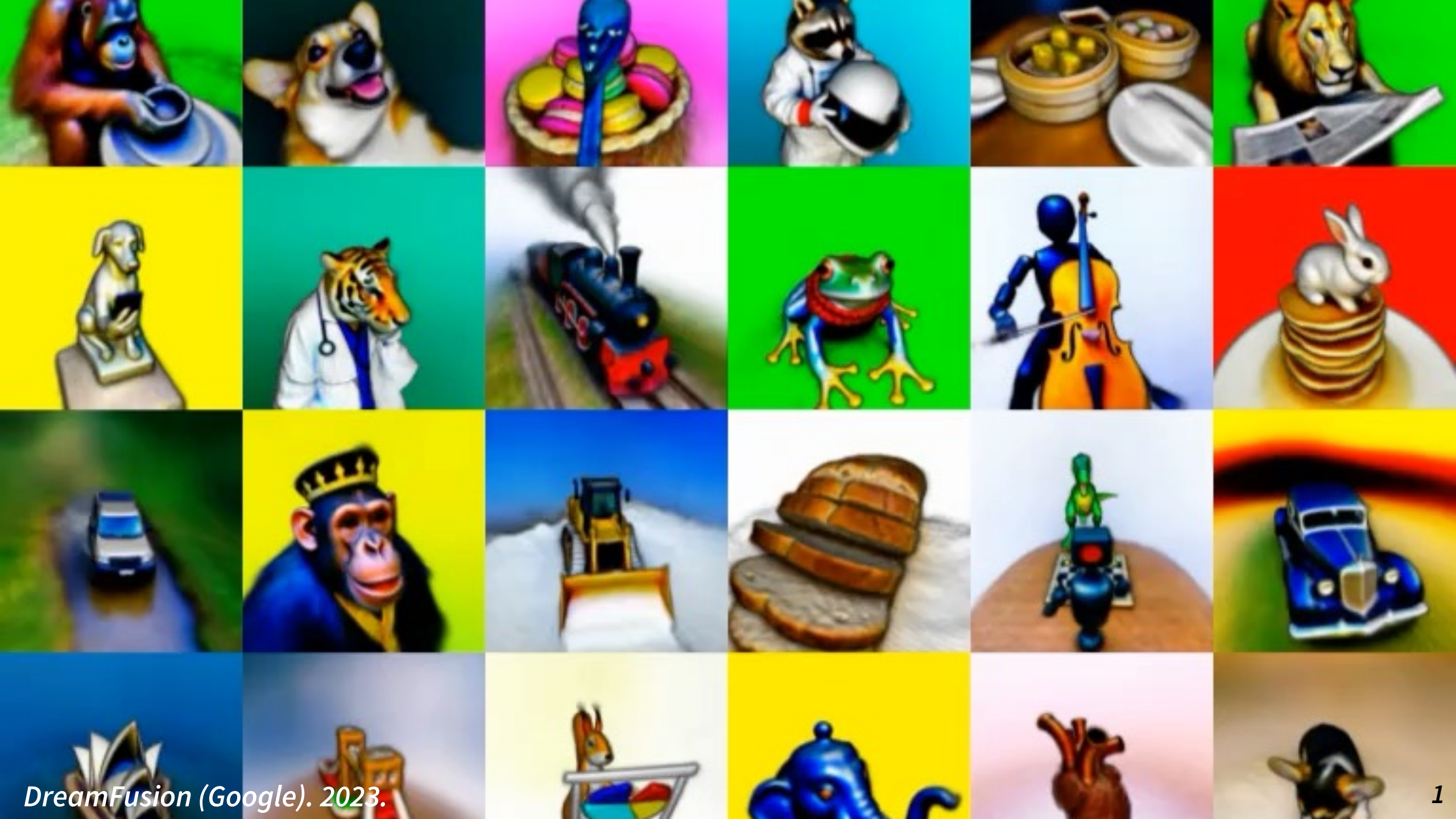


*Diffusion Models for
Visual Content Generation*

3D Generation

Presenter: Minhyuk Sung

Eurographics 2024 Tutorial



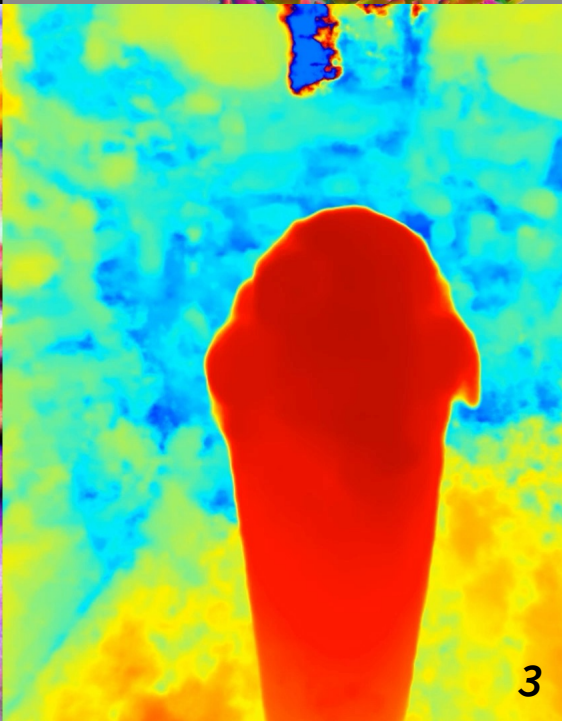
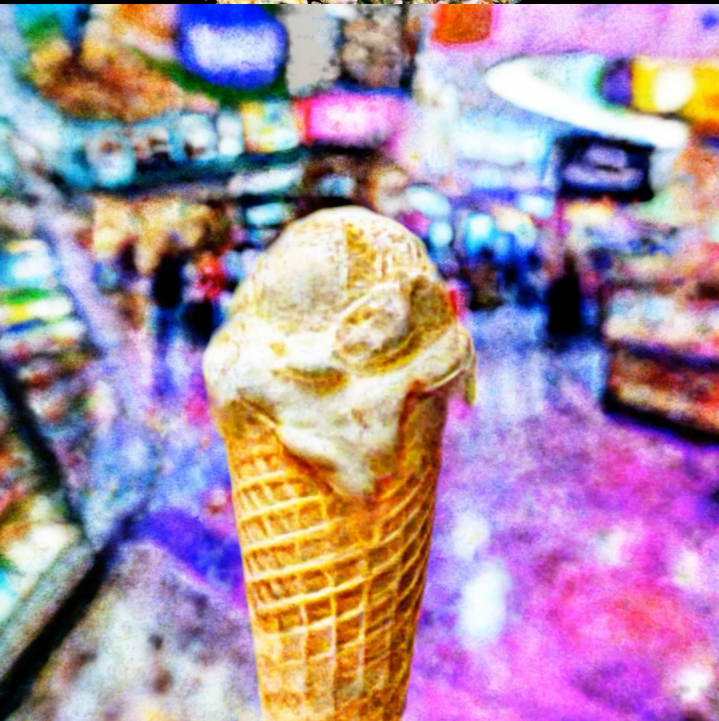
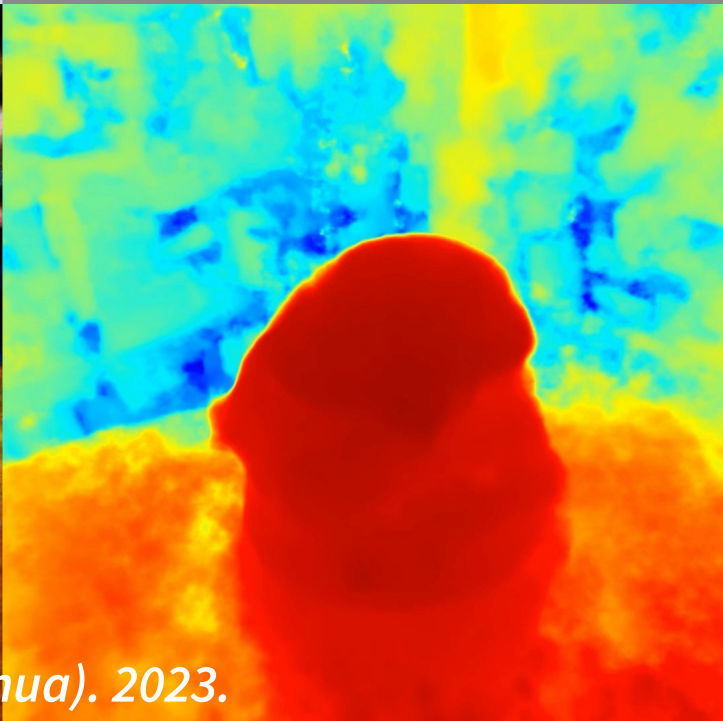
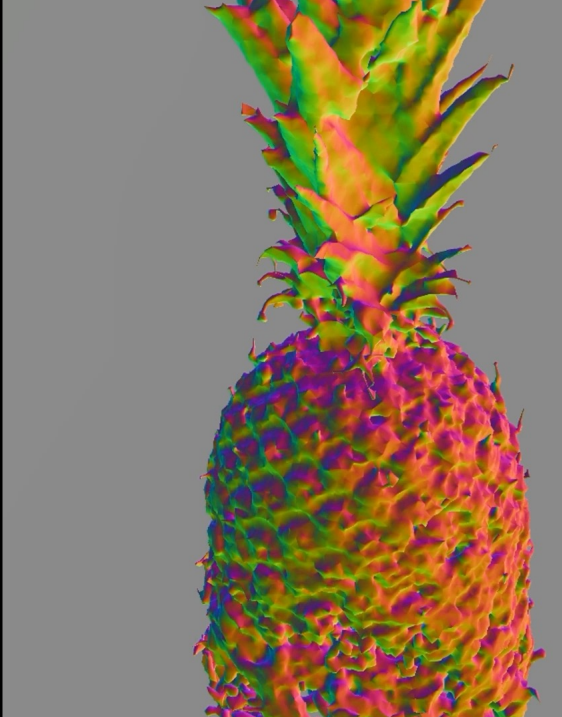
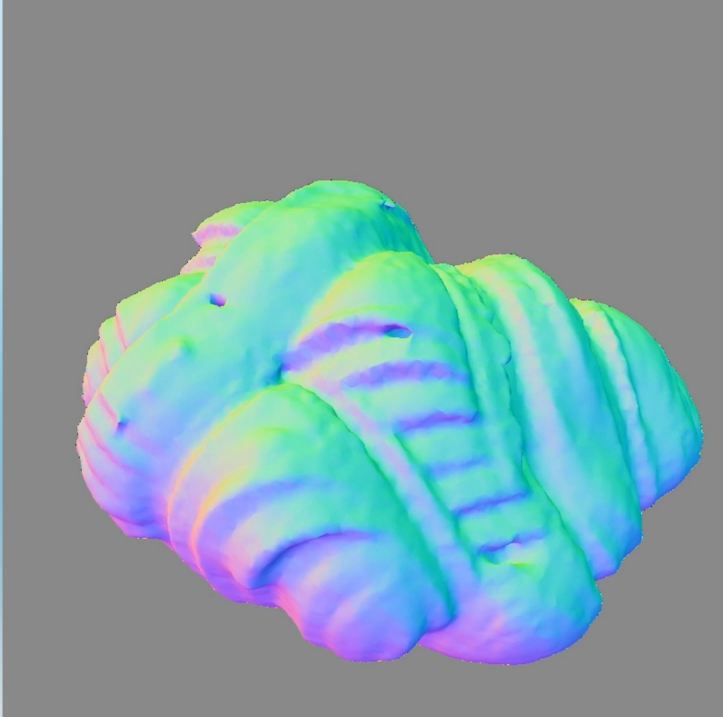


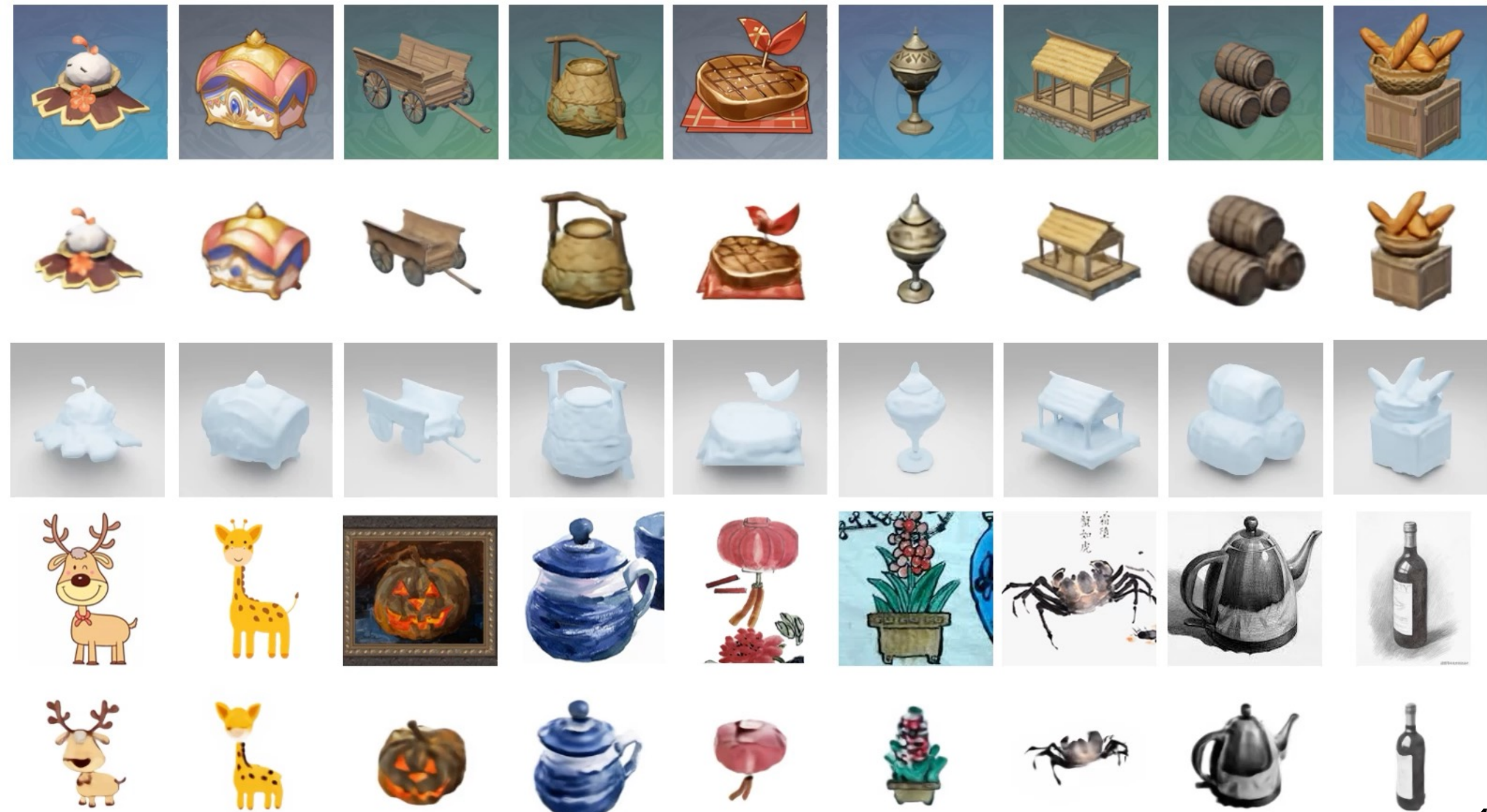
High-Resolution Text-to-3D Content Creation

Chen-Hsuan Lin* Jun Gao* Luming Tang* Towaki Takikawa* Xiaohui Zeng*
Xun Huang Karsten Kreis Sanja Fidler# Ming-Yu Liu# Tsung-Yi Lin

*# : equal contributions

NVIDIA Corporation



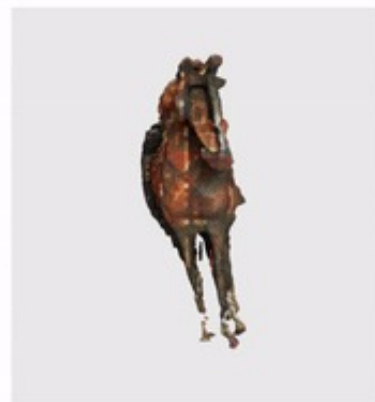
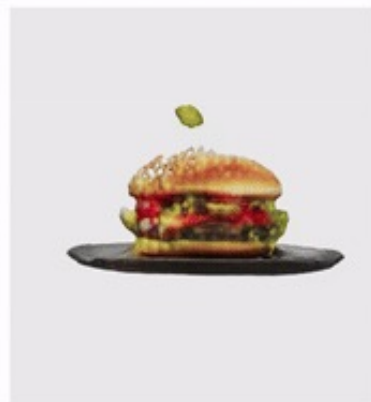
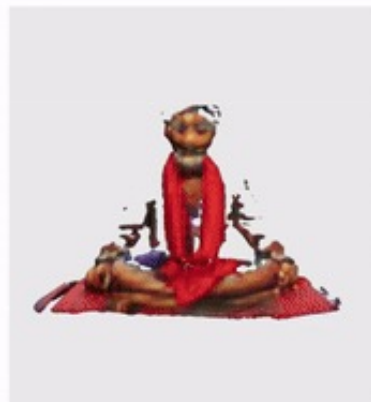


Liu et al., SyncDreamer: Generating Multiview-consistent Images from a Single-view Image, arXiv 2023.

Zero123-XL



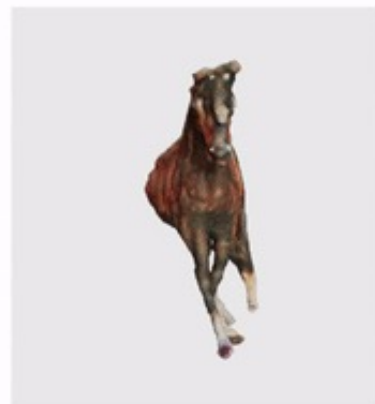
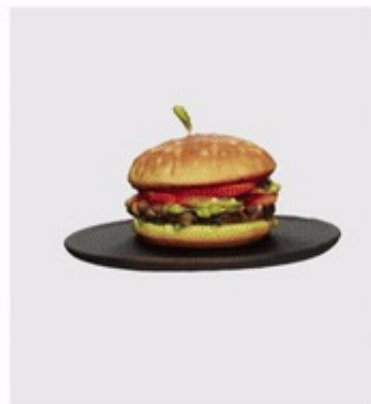
OpenLRM



Stable Zero123



TripoSr (ours)



3D Diffusion Models

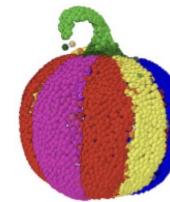
3D Diffusion Models

Point clouds

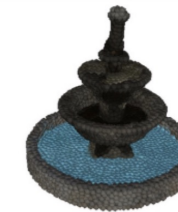
- *ShapeGF (Cai et al., 2020)*
- *DPM (Luo and Hu, 2021)*
- *LION (Nichol et al., 2022)*



"a corgi wearing a red santa hat"



"a multicolored rainbow pumpkin"



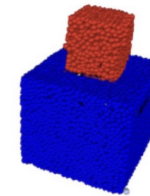
"an elaborate fountain"



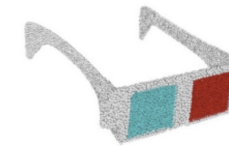
"a traffic cone"



"a vase of purple flowers"



"a small red cube is sitting on top of a large blue cube. red on top, blue on bottom"



"a pair of 3d glasses, left lens is red right is blue"



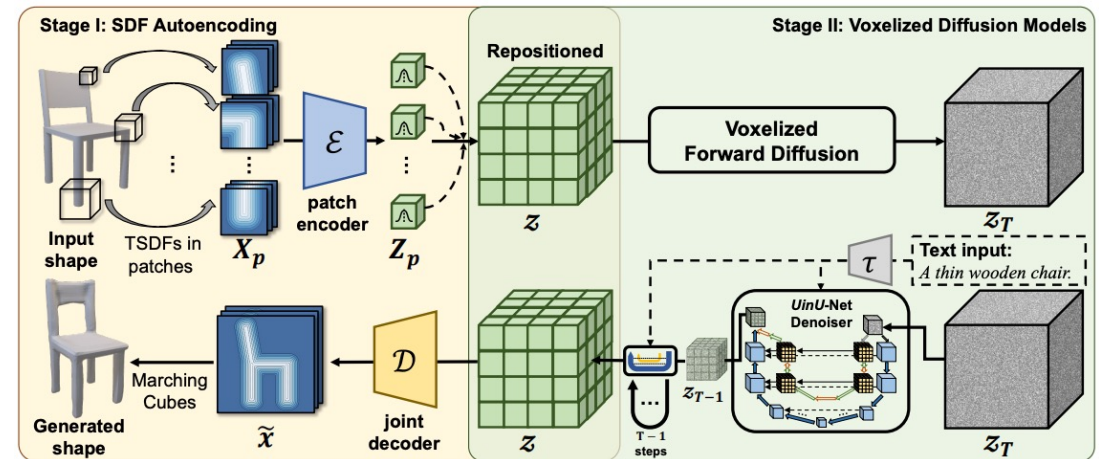
"an avocado chair, a chair imitating an avocado"

LION, Nichol et al. 2022.

3D Diffusion Models

Voxel representation

- *PVD (Zhou et al., 2022)*
- *Diffusion-SDF (Li et al., 2022)*

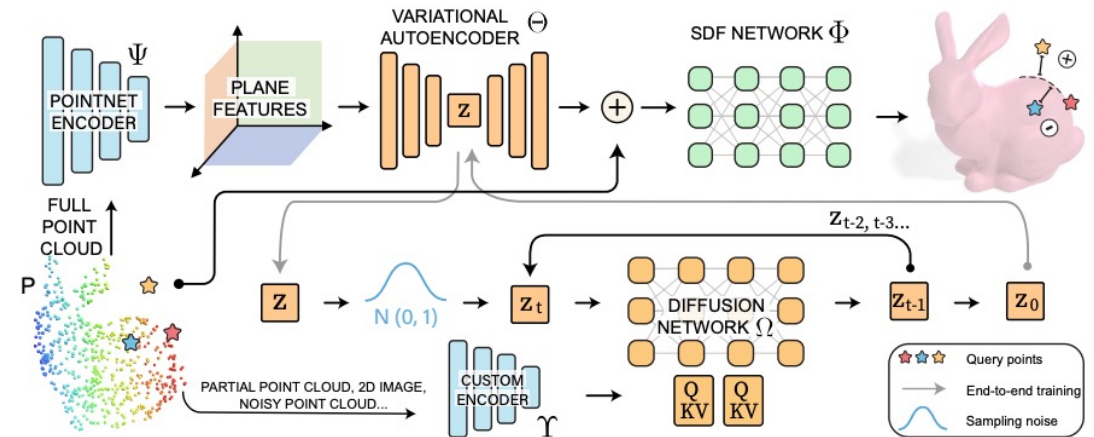


Diffusion-SDF, Li et al., 2022.

3D Diffusion Models

Latent representation

- *SDFusion*
(Cheng et al., 2022)
- *DiffusionSDF* (w/o hyphen,
Chou et al., 2023)



DiffusionSDF, Chou et al., 2023

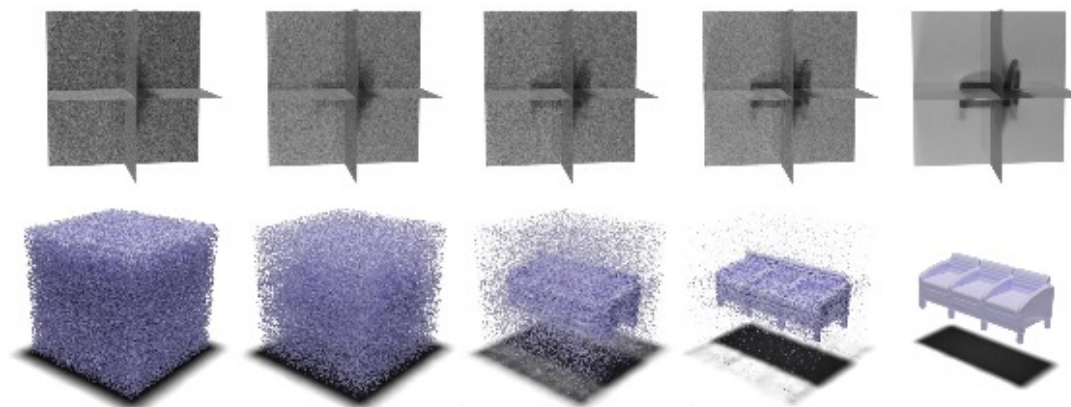
3D Diffusion Models

Triplane representation

- *NFD (Shue et al., 2022)*

Diffusion in the spectral domain:

- *NeuralWavelet (Hui et al., 2022)*



Shape Generation



Shape Interpolation

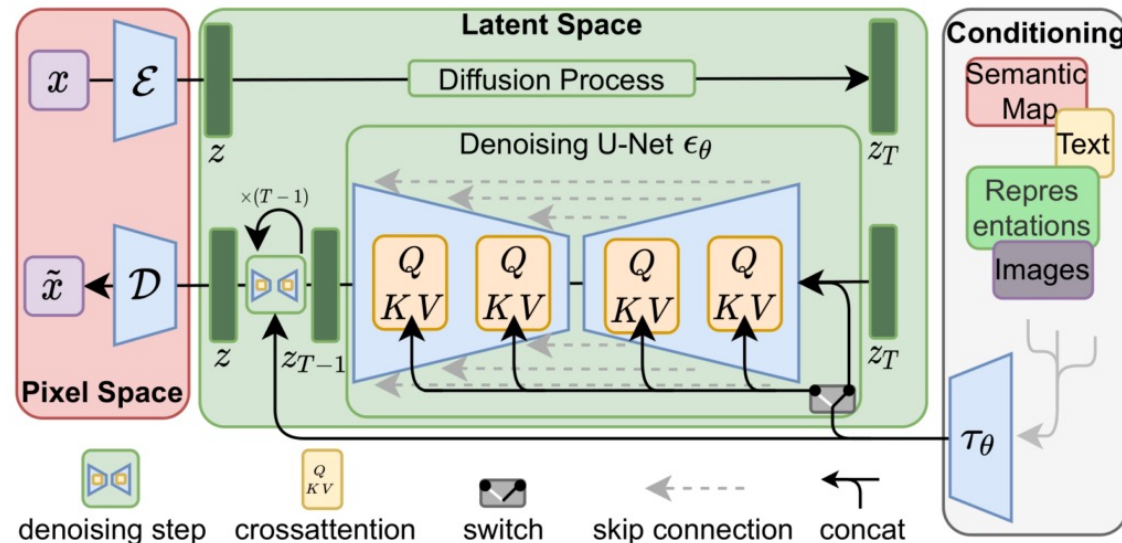
NFD, Shue et al., 2022.

Diffusion w/ Different Representations

- *Implicit representation*
- *Explicit representation*

Diffusion w/ Different Representations

- **Implicit** representation (i.e., **latent** features)
 - E.g., Latent diffusion (Rombach et al., 2022)
 - (+) **Best quality** of the generated data.
 - (-) **Requires retraining** for each conditional generation setup.



Latent Diffusion, Rombach et al., 2022.

Diffusion w/ Different Representations

- **Explicit** representation (E.g., **pixels** in images)
 - E.g.: The original DDPM model (Ho et al., 2020) for images.
 - (–) **Suboptimal performance** due to the high dimensionality.
 - (–) Cannot change the **resolution** of the data.
 - (+) Can be directly leveraged in **conditional generation** setups in a **zero-shot** manner.

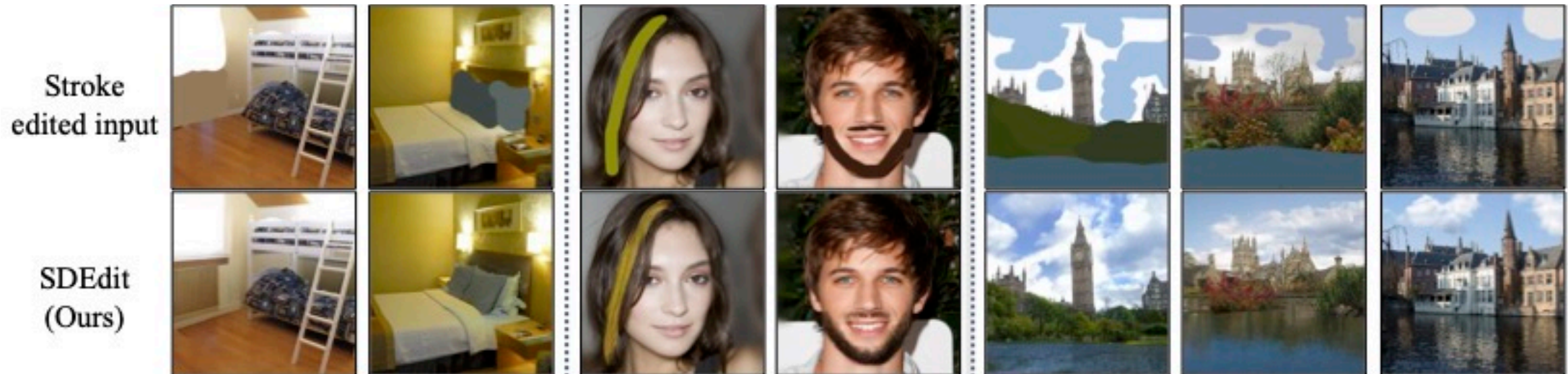
SDEdit [Meng et al., 2022]

Image **editing** using a pretrained pixel-space diffusion model.

$$\mathbf{x}^{(0)} = m \odot \mathbf{x}_a^{(0)} + (1 - m) \odot \mathbf{x}_b^{(0)}$$

$$\mathbf{x}^{(t)} = \text{denoise}(\mathbf{x}^{(0)}, t)$$

$$\mathbf{x}'^{(0)} = \text{add_noise}(\mathbf{x}^{(t)}, t)$$



RePaint [Lugmayr et al., 2022]

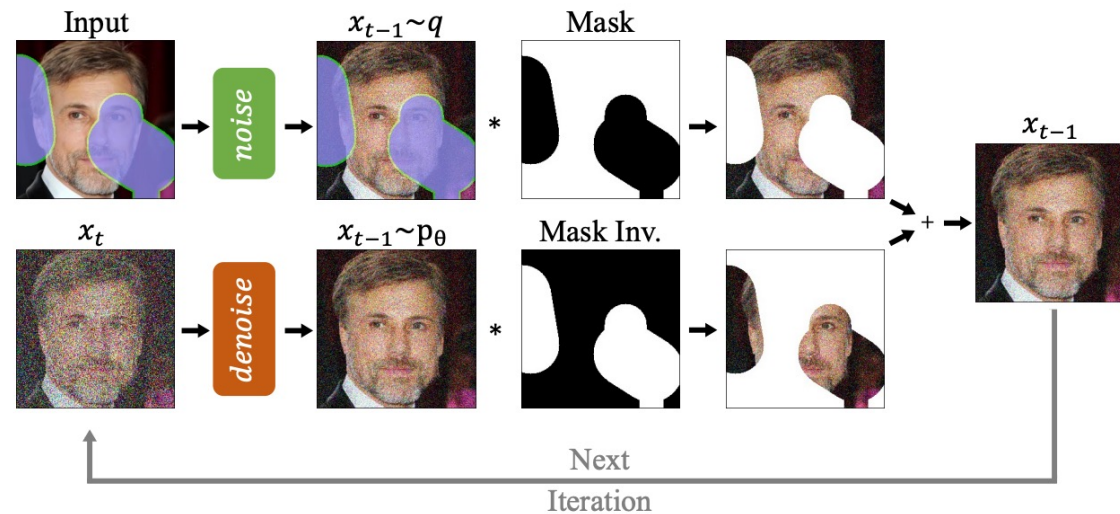
Image *inpainting* using a pretrained pixel-space diffusion model.

$$\mathbf{x}_f^{(t-1)} = \text{denoise}(\mathbf{x}^{(t)}, 1)$$

$$\mathbf{x}_b^{(t-1)} = \text{add_noise}(\mathbf{x}^{(0)}, t)$$

$$\mathbf{x}^{(t-1)} = m \odot \mathbf{x}_f^{(t-1)} + (1 - m) \odot \mathbf{x}_b^{(t-1)}$$

Repeat for $t = T, \dots, 1$.



Hybrid Representation

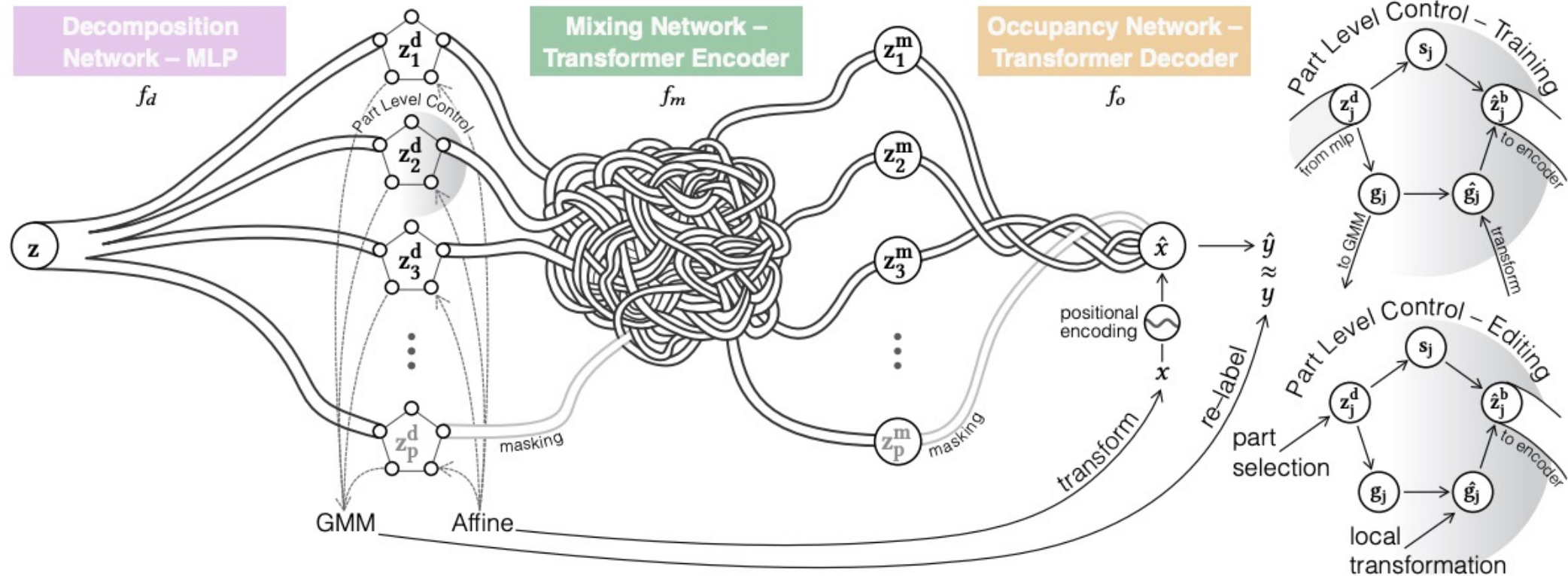
Leverages a novel *hybrid* representation describing

- *global* part-level structure *explicitly*, and
- *local* geometry *implicitly*.



SPAGHETTI [Hertz et al., 2022]

The part decomposition and the explicit/implicit representations are learned in an **unsupervised** way.



Part-Level Representation

For each part of an object learned in an *unsupervised* way,

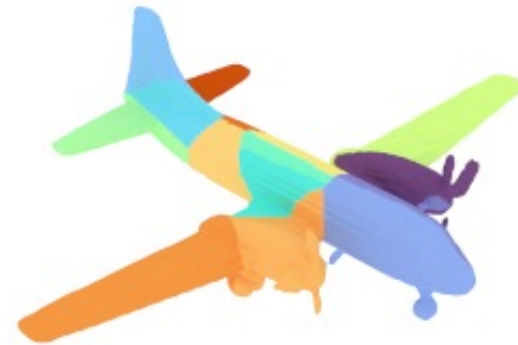
- *Explicit* parameters of *Gaussian blubs* indicate position, scale, and rotation.



Part-Level Representation

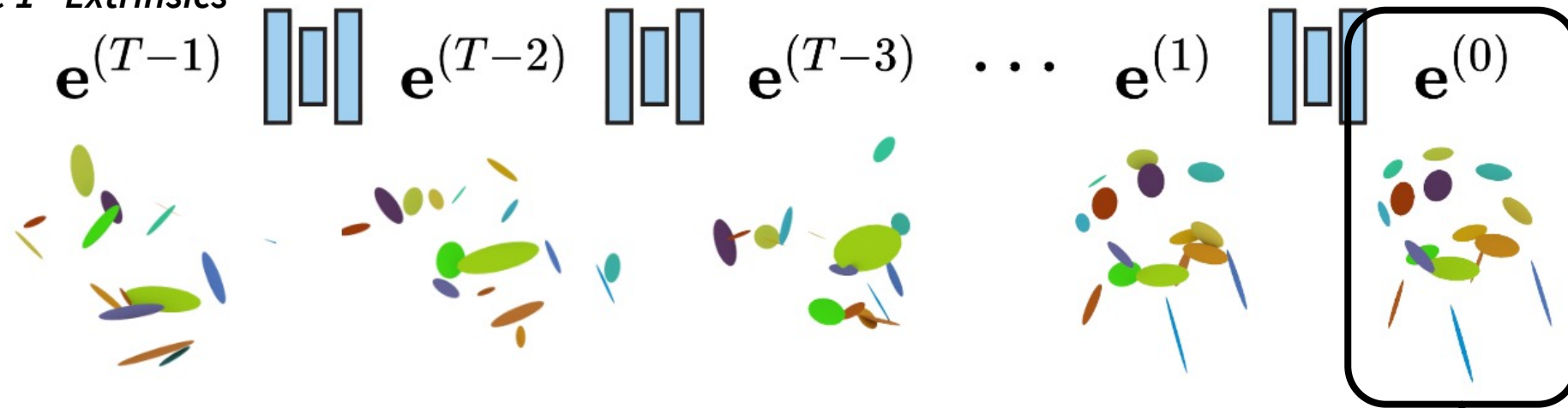
For each part of an object learned in an *unsupervised* way,

- *Implicit* latent feature is decoded into an *occupancy* function.

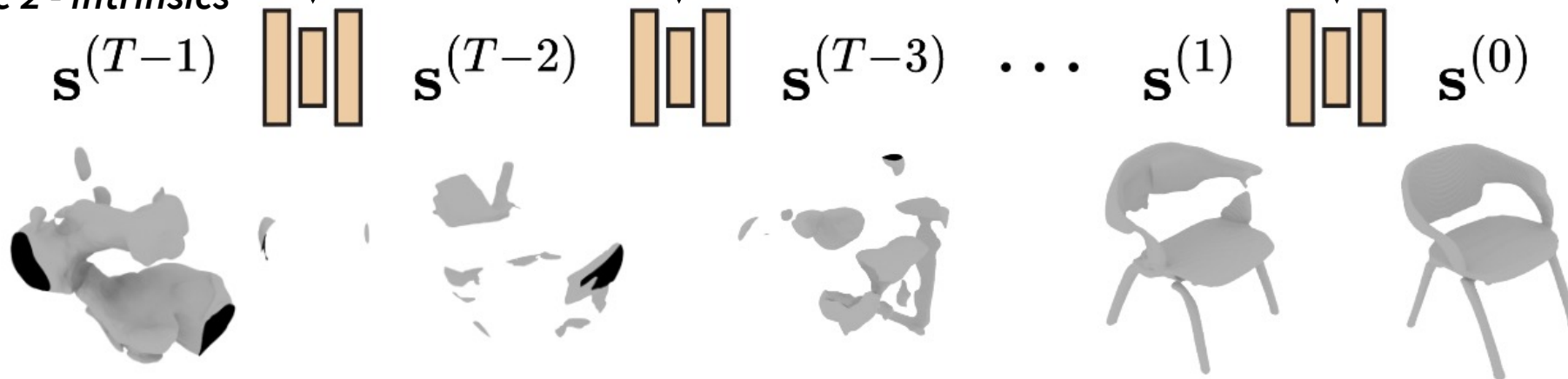


Two-Phase Cascaded Diffusion

Phase 1 - Extrinsics



Phase 2 - Intrinsic



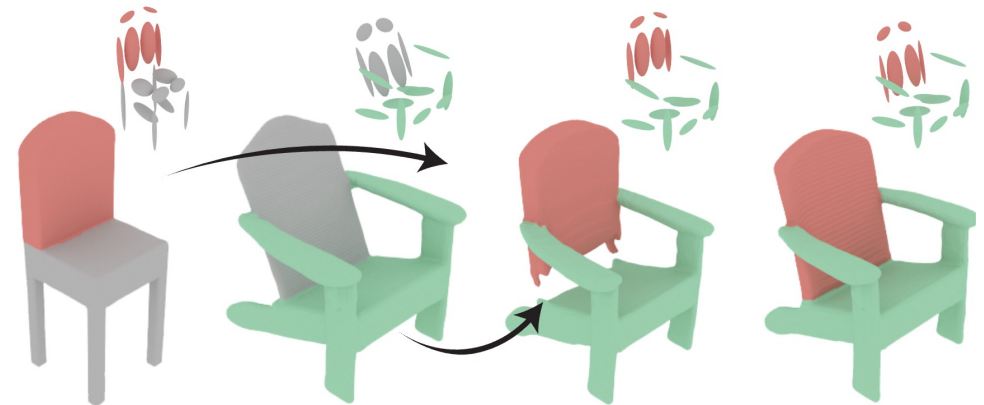
Diversity of 3D Shapes



Applications



Part Completion



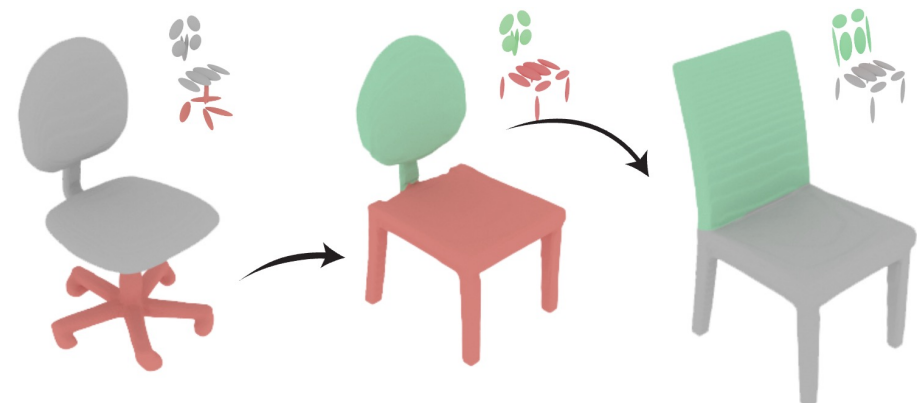
Part Mixing



"Chair has **round arms** and **wheels**."

"Its the one with **gaps** in the **back**."

Text-to-3D Generation



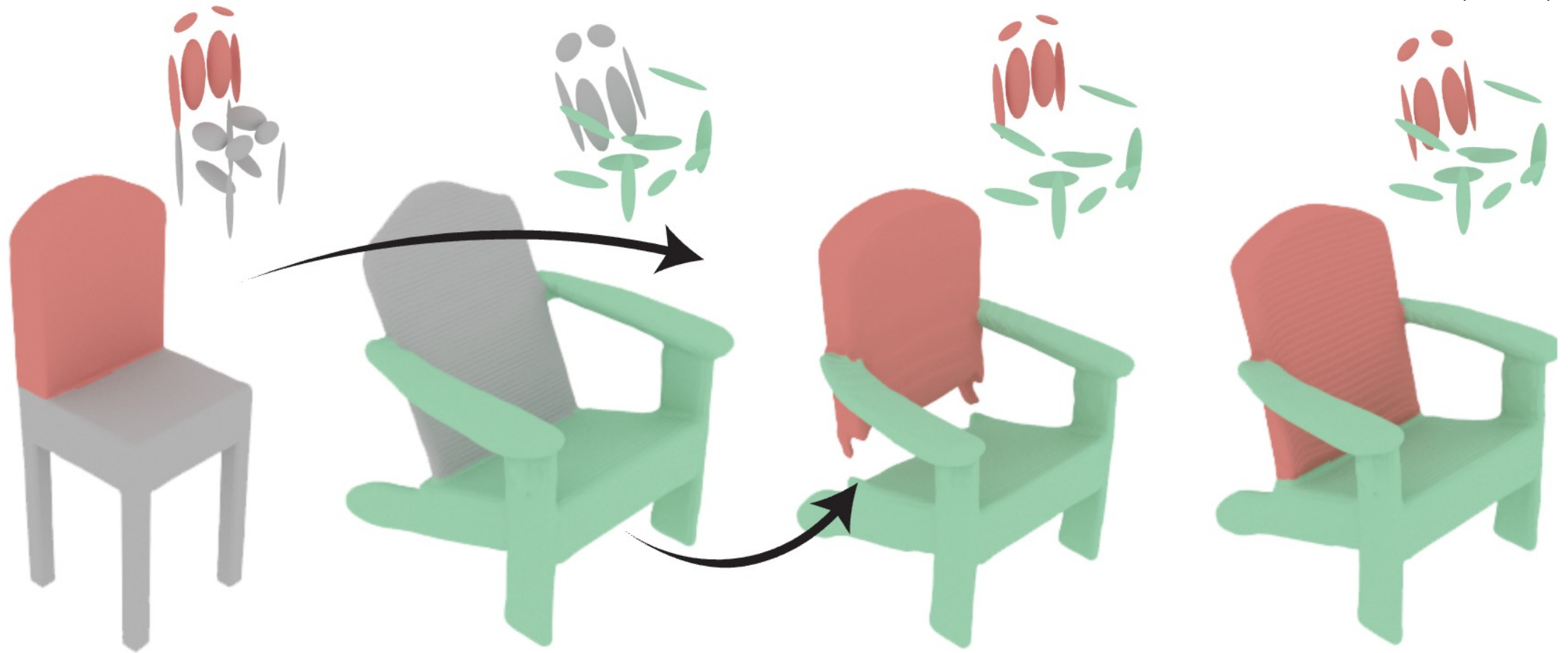
"A chair with **four legs**"

"**rectangle back** chair"

Text-Guided Part Editing

Part Mixing

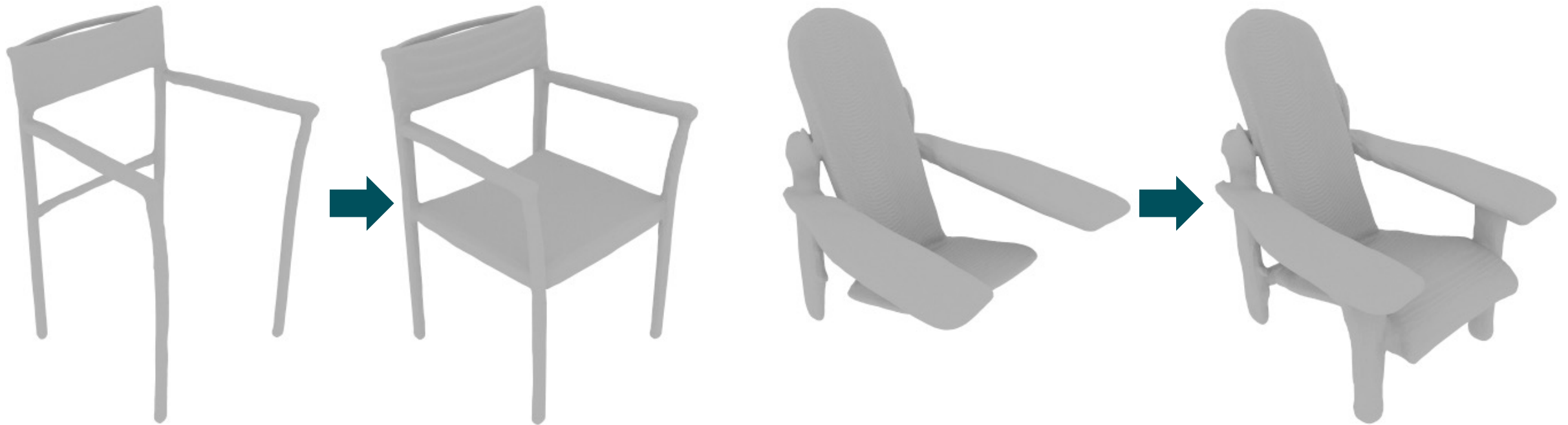
$$\mathbf{x}^{(0)} = m \odot \mathbf{x}_a^{(0)} + (1 - m) \odot \mathbf{x}_b^{(0)}$$
$$\mathbf{x}^{(t)} = \text{denoise}(\mathbf{x}^{(0)}, t)$$
$$\mathbf{x}'^{(0)} = \text{add_noise}(\mathbf{x}^{(t)}, t)$$



Part Completion

$$\begin{aligned} \mathbf{x}_f^{(t-1)} &= \text{denoise}(\mathbf{x}^{(t)}, 1) \\ \mathbf{x}_b^{(t-1)} &= \text{add_noise}(\mathbf{x}^{(0)}, t) \\ \mathbf{x}^{(t-1)} &= m \odot \mathbf{x}_f^{(t-1)} + (1 - m) \odot \mathbf{x}_b^{(t-1)} \end{aligned}$$

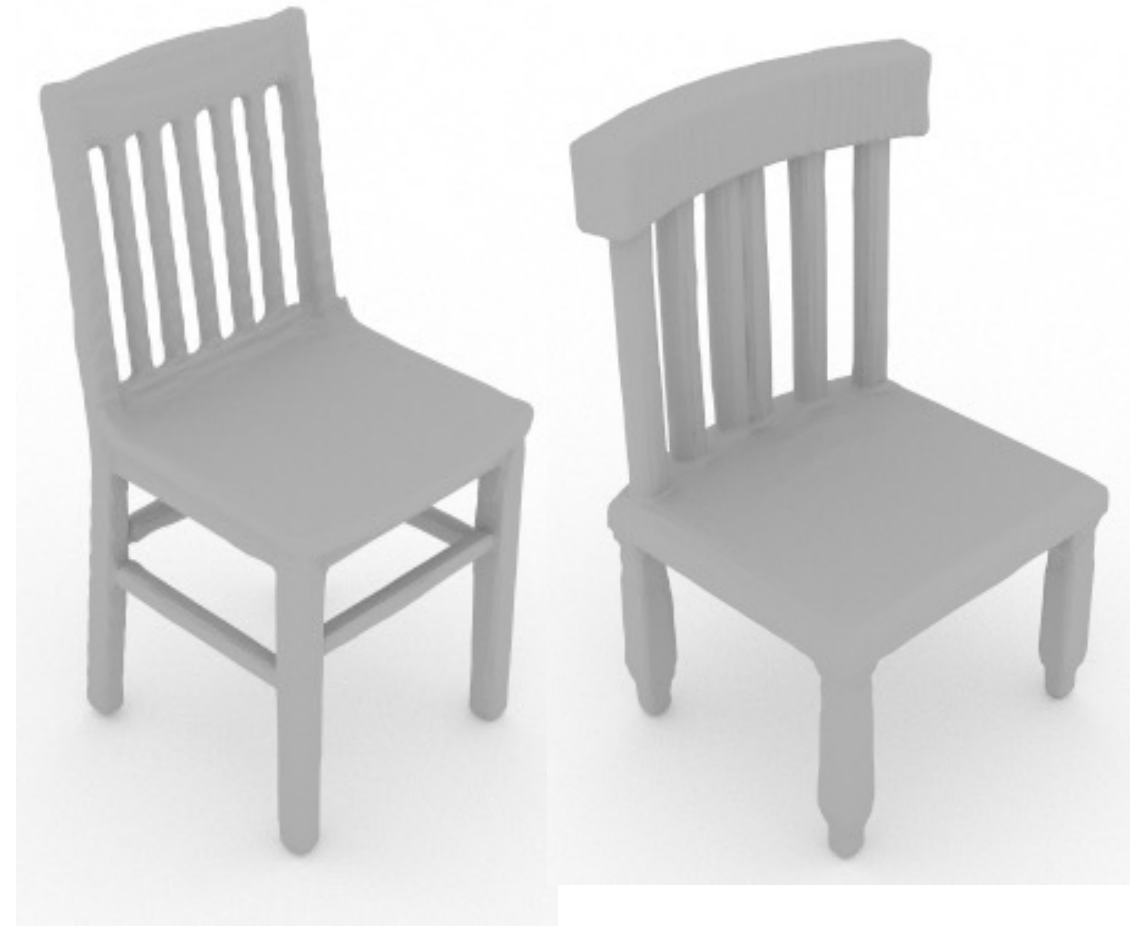
Repeat for $t = T, \dots, 1$.



Text-to-3D Generation

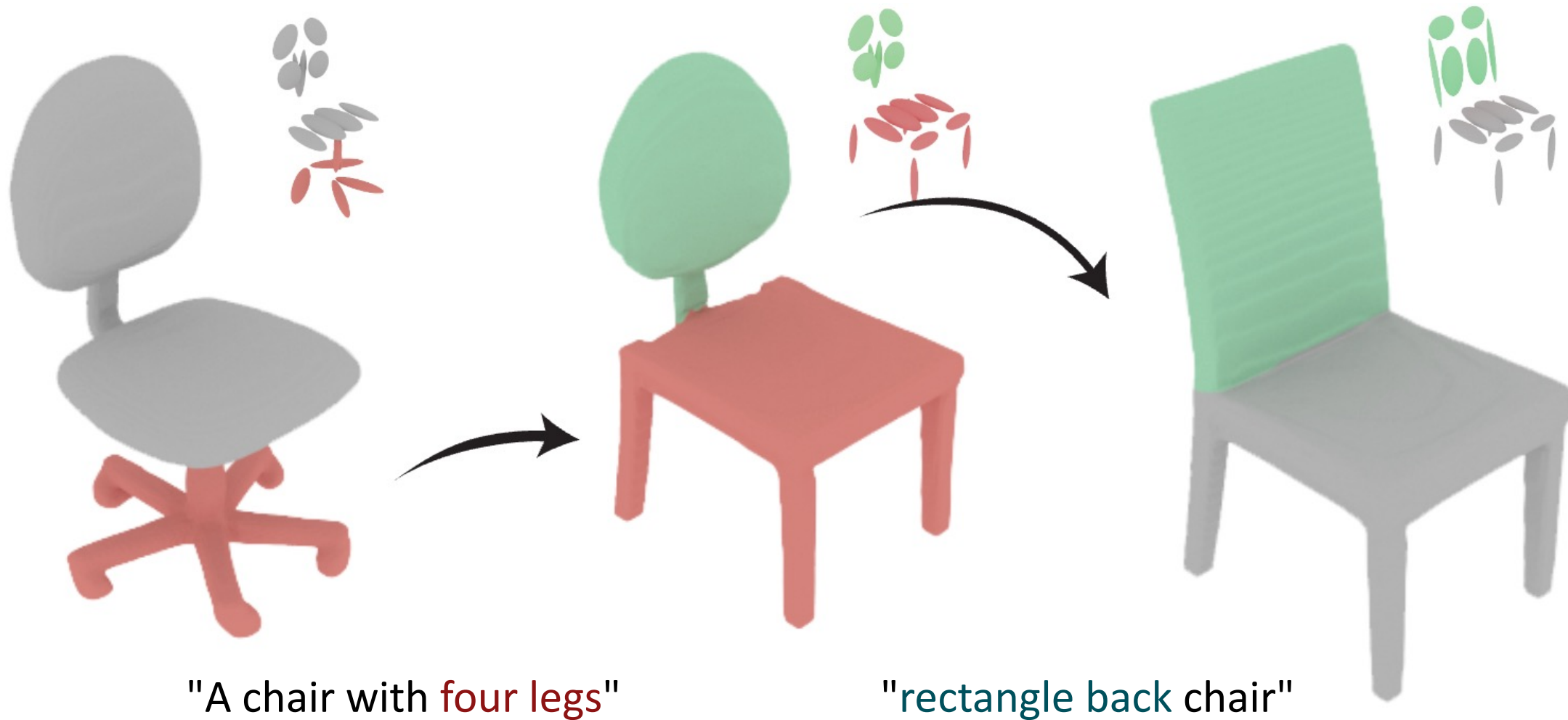


"Chair has **round arms** and **wheels**."



"Its the one with **gaps** in the **back**."

Text-Guided Part Editing



Limitation?

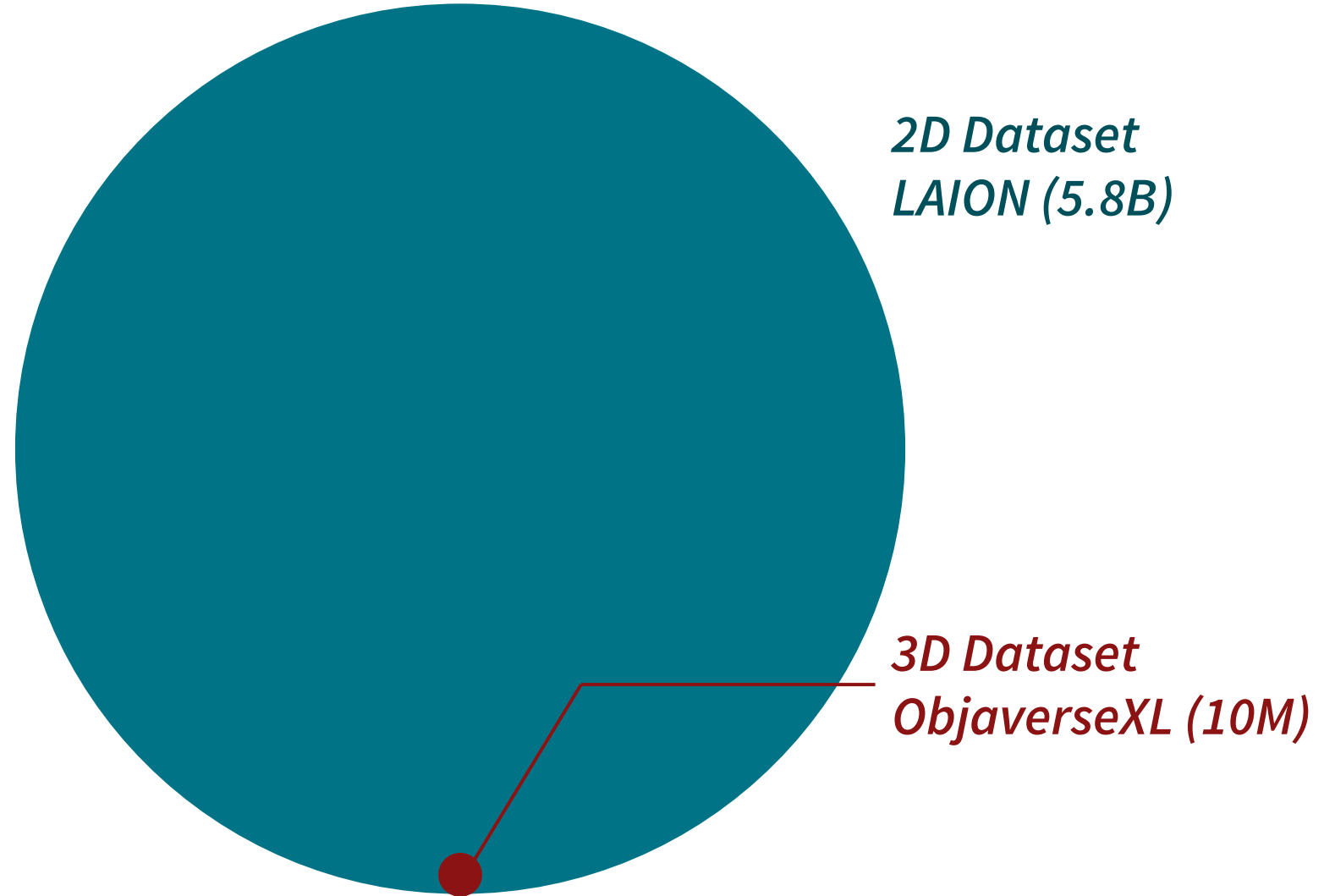


Challenge: Lack of Large-Scale 3D Dataset



Objaverse-XL
Allen Institute

Data Deficiency



Diversity of Imaginable 3D Shapes



“frog wearing a sweater”



“eggshell broken in two with an adorable chick standing next to it”



“ghost eating a hamburger”



“a pig wearing a backpack”

*We have a small-scale 3D dataset
but images on an internet scale.*

*Images are projections of 3D
from specific angles.*



3D Generation

*How to generate a 3D object
from a collection of 2D images?*

3D Reconstruction

*How to reconstruct a 3D object
from a collection of 2D images
of a specific object?*

NeRF



3D Reconstruction

- *Input: A set of **images** with **camera poses**.*
- *Output: A representation of the **3D object**.*





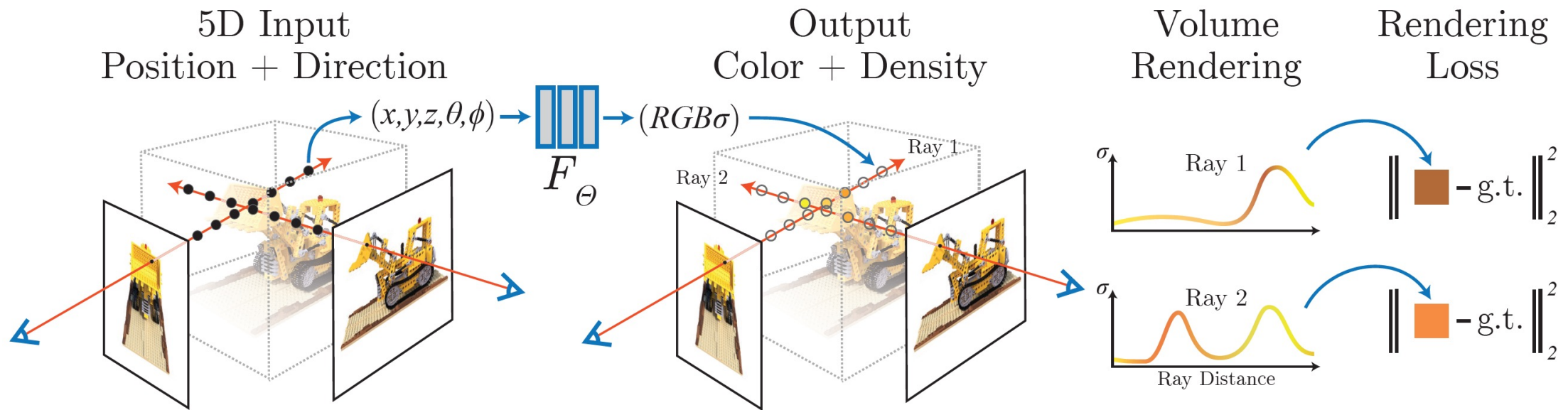
▼ Metrics

59.97 (16.67 ms)

Gaussian Splatting

NeRF Optimization

1. **Render** a NeRF representation into a specific view.
2. Compute the **difference** with the **given image**.
3. Update the NeRF using **gradient descent**.



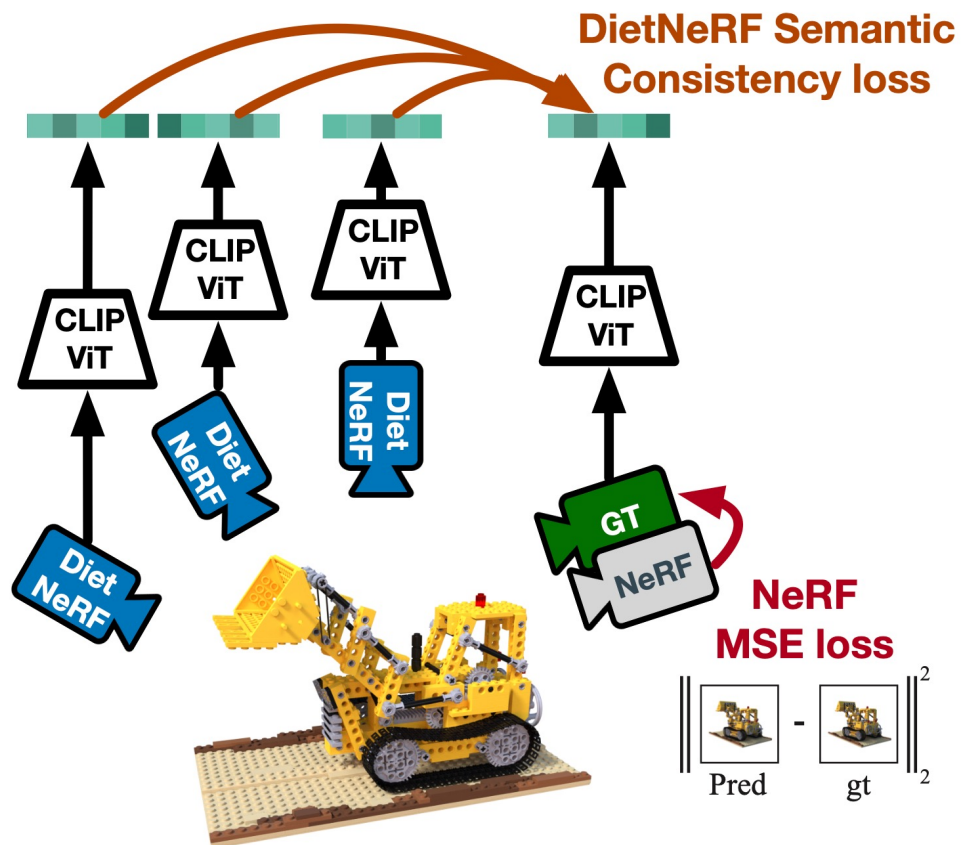
*Can we perform NeRF reconstruction
with a few images?*

Then, we need additional information!

Few-Shot NeRF

DietNeRF [Jain et al., ICCV 2021]

Use *priors* learned by **CLIP** [Radford et al., 2021], a *text-image model*.



“a bulldozer is a bulldozer from any perspective”

CLIP [Radford et al., 2021]

CLIP takes a text-image pair as input and assesses the *alignment* between the text and the image.

“A beautiful painting of a dog riding a dolphin” +  = CLIP 70.7%

STEP 500

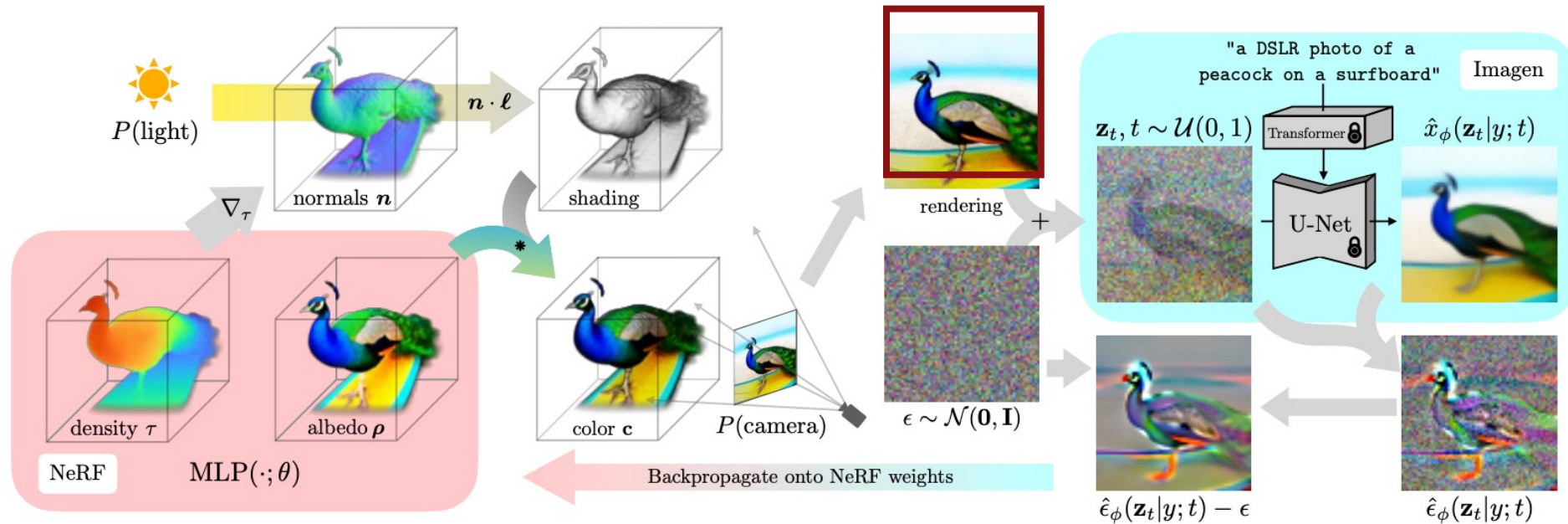
The image shows a text prompt on the left, a plus sign, a painting of a dog riding a dolphin in the center, an equals sign, and the CLIP score of 70.7% on the right. The painting is a colorful, detailed work of art depicting a dog sitting on the back of a dolphin in a body of water, with a town and a church in the background. The text is in a monospaced font. The CLIP score is in a bold, sans-serif font. The entire scene is enclosed in a light blue box with a black border.

Knowledge Distillation in 3D Generation

1. *Render* the NeRF representation into a specific view

2. Compute the *alignment* to the given text.

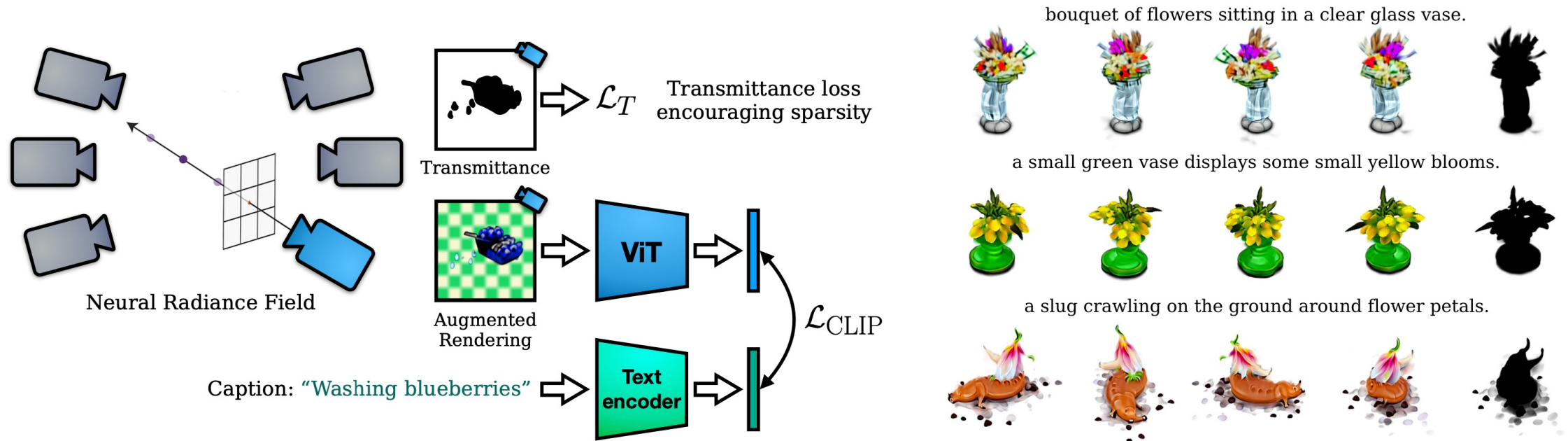
3. Update the NeRF using the *gradient descent*.



Extrem Case: Zero-Shot NeRF

DreamFields [Jain et al., CVPR 2022]

Given *a text prompt* but *no images*, generate a 3D shape by maximizing similarity between a *rendered image* and the input prompt in the CLIP embedding space.



DreamFields [Jain et al., CVPR 2022]

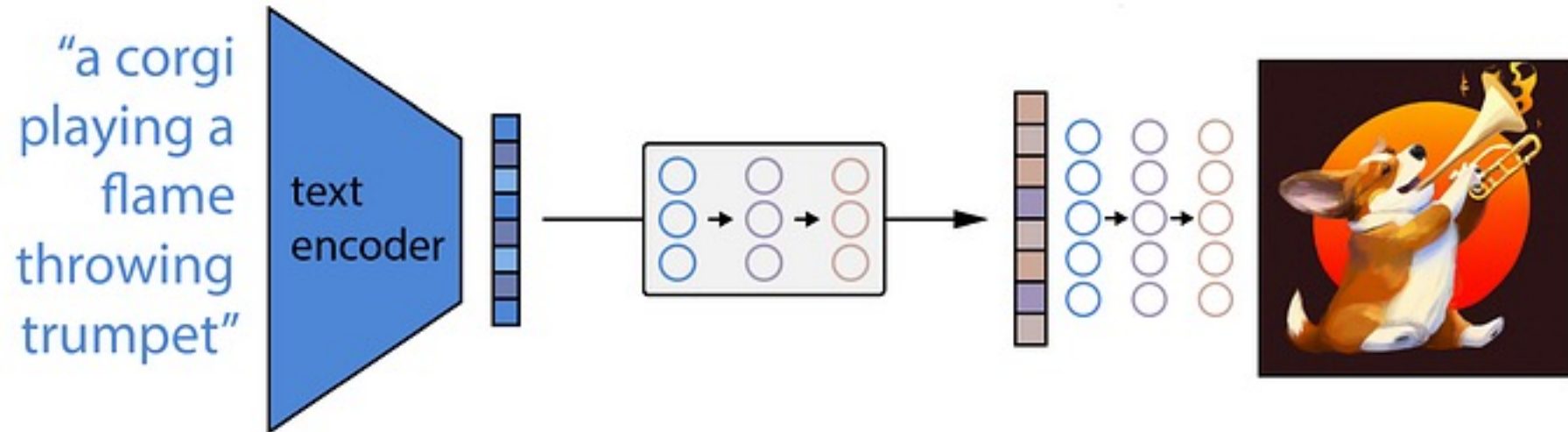


The first 3D generative model
using 2D priors!

an archair in the shape of a _____.
an archair imitating a _____.

a teapot in the shape of a _____.
a teapot imitating a _____.

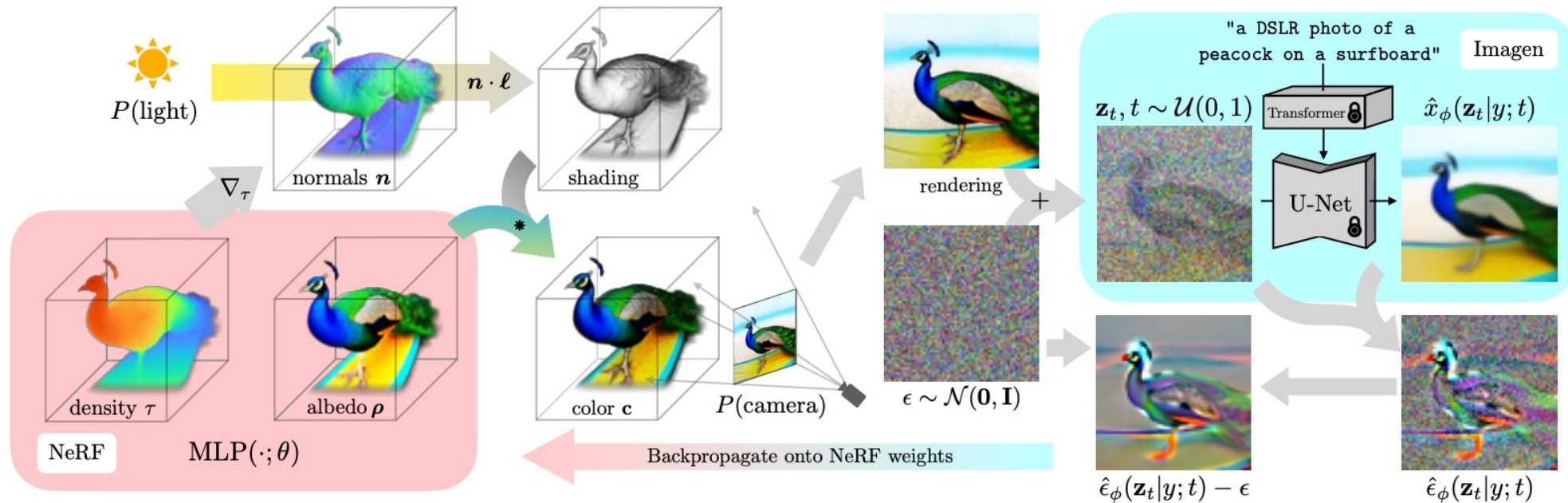
Can we use an *image diffusion model* instead of CLIP?



Score Distillation Sampling (SDS)

DreamFusion [Poole et al., ICLR 2023]

Proposed the idea of **Score Distillation Sampling (SDS)**, leveraging a pretrained diffusion model to **measure the plausibility** of rendered images.



Score Distillation Sampling (SDS)

- How can we utilize a pretrained diffusion model to *measure the plausibility* of rendered images?

- Review the *loss* function:

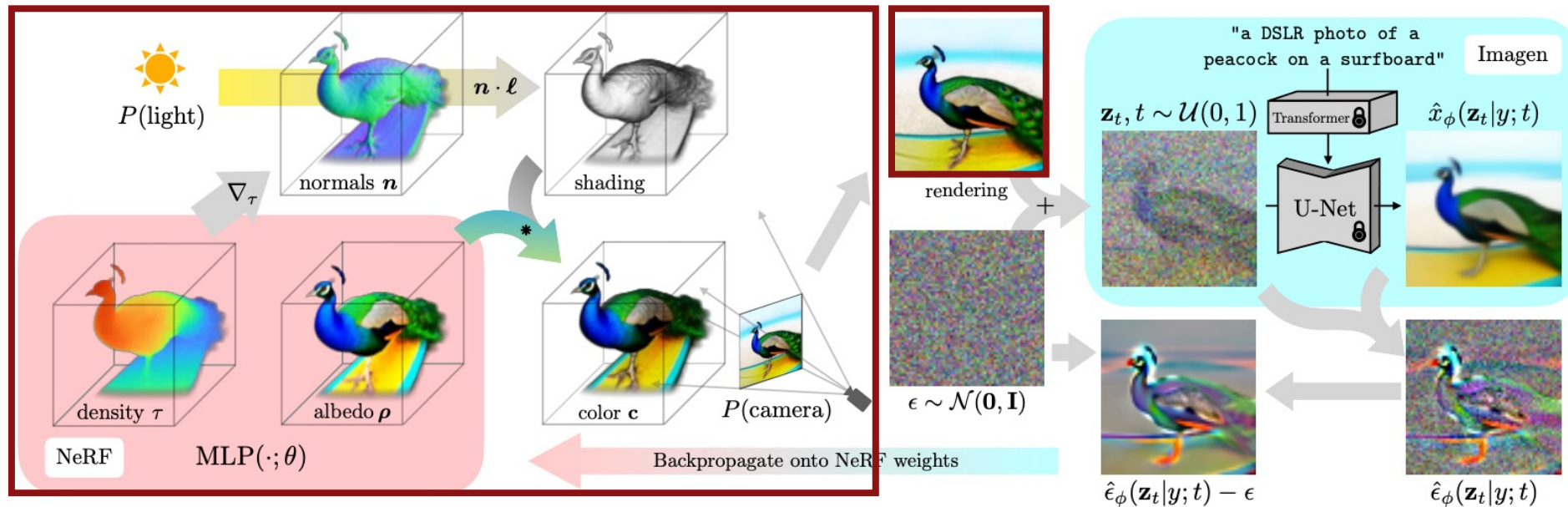
$$\mathcal{L} = \mathbb{E}_{t \sim [1, T], \mathbf{x}_0, \boldsymbol{\epsilon}_t} \left[\left\| \boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_\theta \left(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_t, t \right) \right\|^2 \right]$$

If the training of the diffusion model has converged, the loss for *real* data \mathbf{x}_0 will be close to zero.

Score Distillation Sampling (SDS)

1. **Render** the NeRF representation into a specific view.

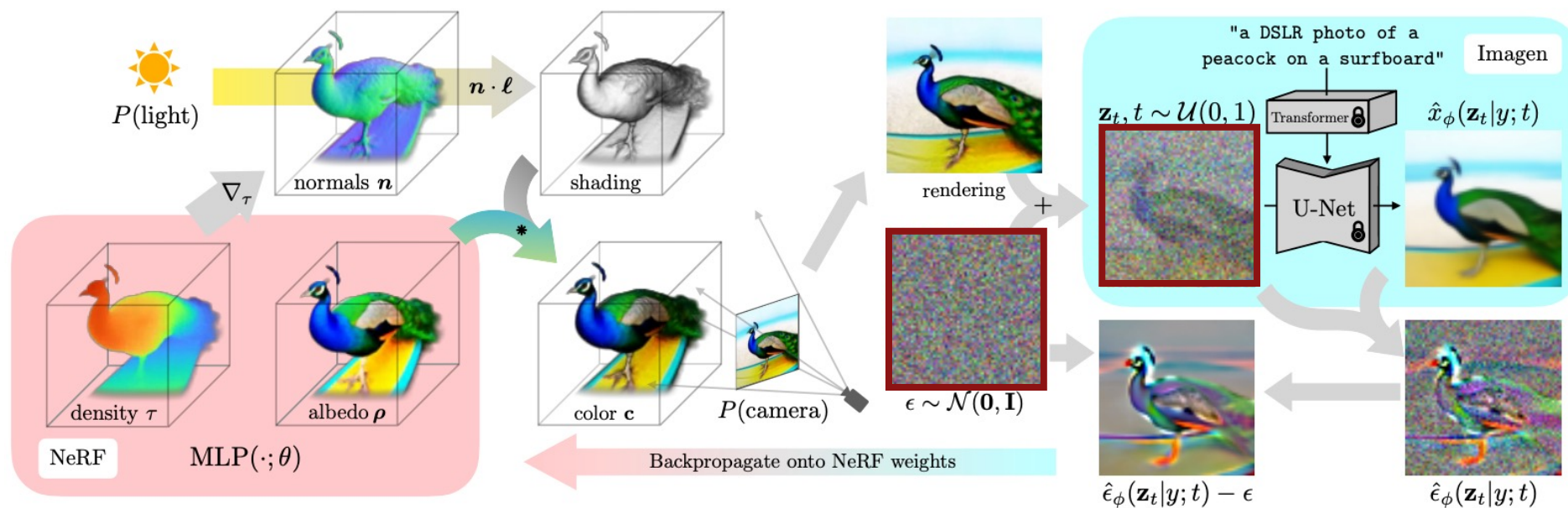
Let ϕ denote the NeRF parameter, and $\mathbf{x}_0 = g(\phi)$ denote the rendered image.



Score Distillation Sampling (SDS)

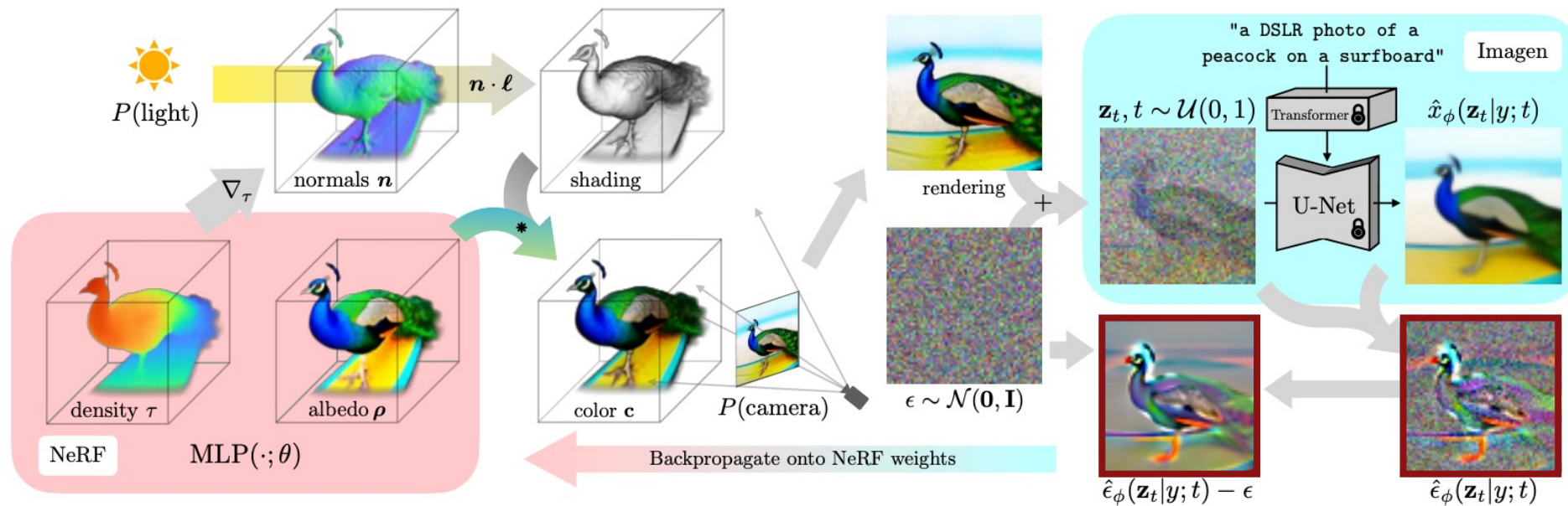
2. Add **noise** to the rendered image \mathbf{x}_0 :

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_t.$$



Score Distillation Sampling (SDS)

3. Perform **gradient descent** on \mathcal{L} with respect to the NeRF parameters ϕ .



DreamFusion Results



“frog wearing a sweater”



“eggshell broken in two
with an adorable chick
standing next to it”



“ghost eating a hamburger”



“a pig wearing a backpack”

Why SDS Instead of Reverse Diffusion?

*This is a scenario where the images are **parameterized differently from how they were represented during the training of the diffusion model.***

- *Training: Per-pixel colors.*
- *Inference: NeRF rendering.*

Example: Vector Images / Sketches

The same idea but with a **different parameterization** of images.



A blue poison-dart frog sitting on a water lily. [...]



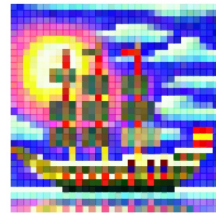
the Great Pyramid. [...]



A rabbit cutting grass with a lawnmower. [...]



Translation. [...]



a pirate ship landing on the moon. [...]



A snail on a leaf. [...]



Yeti taking a selfie. [...]



a hedgehog. [...]



The space between infinity. [...]



A realistic photograph of a cat. [...]



A watercolor painting of a cat. [...]



A Japanese woodblock print of one cat. [...]

Example: Mesh Editing

The same idea but with a **different parameterization of images.**

Source



ARAP

[Sorkine and Alexa, 2007]



APAP

[Yoo et al., 2024]



Stable-DreamFusion

Stable-Dreamfusion

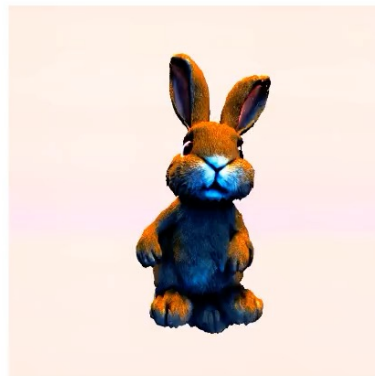
A pytorch implementation of the text-to-3D model Dreamfusion, powered by the [Stable Diffusion](#) text-to-2D model.

ADVERTISEMENT: Please check out [threestudio](#) for recent improvements and better implementation in 3D content generation!

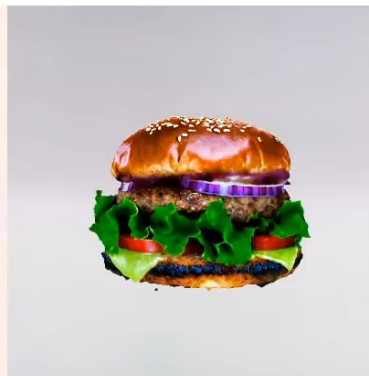
NEWS (2023.6.12):

- Support of [Perp-Neg](#) to alleviate multi-head problem in Text-to-3D.
- Support of Perp-Neg for both [Stable Diffusion](#) and [DeepFloyd-IF](#).

Text-to-3D



a rabbit, animated movie character, high detail 3d model



a DSLR photo of a delicious hamburger



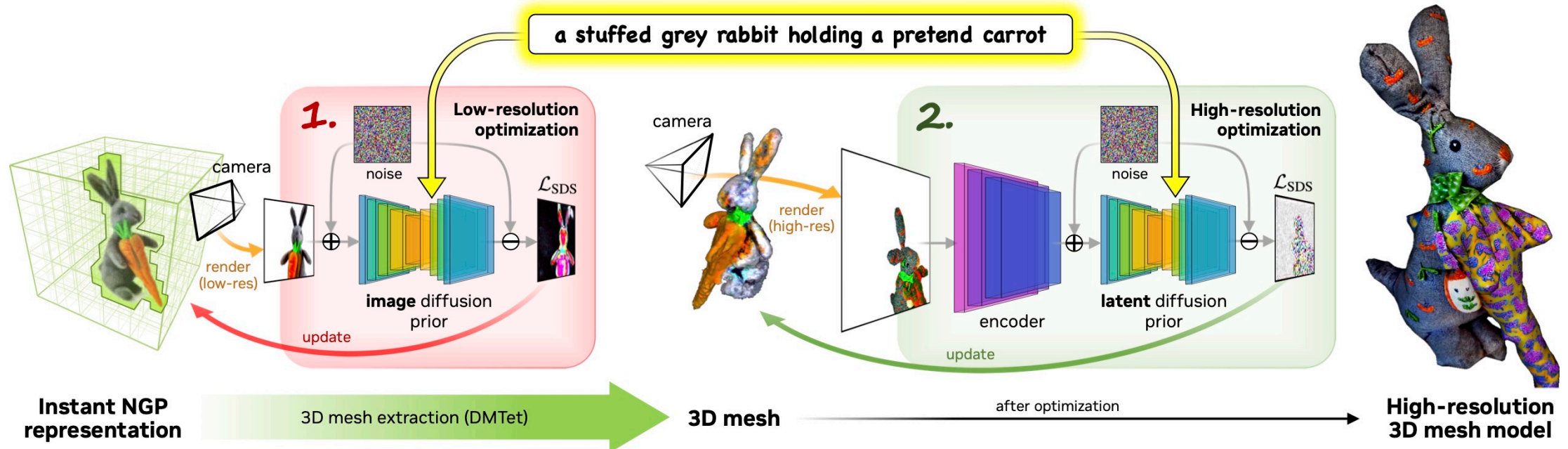
a highly detailed stone bust of Theodoros Kolokotronis



a small saguaro cactus planted in a clay pot

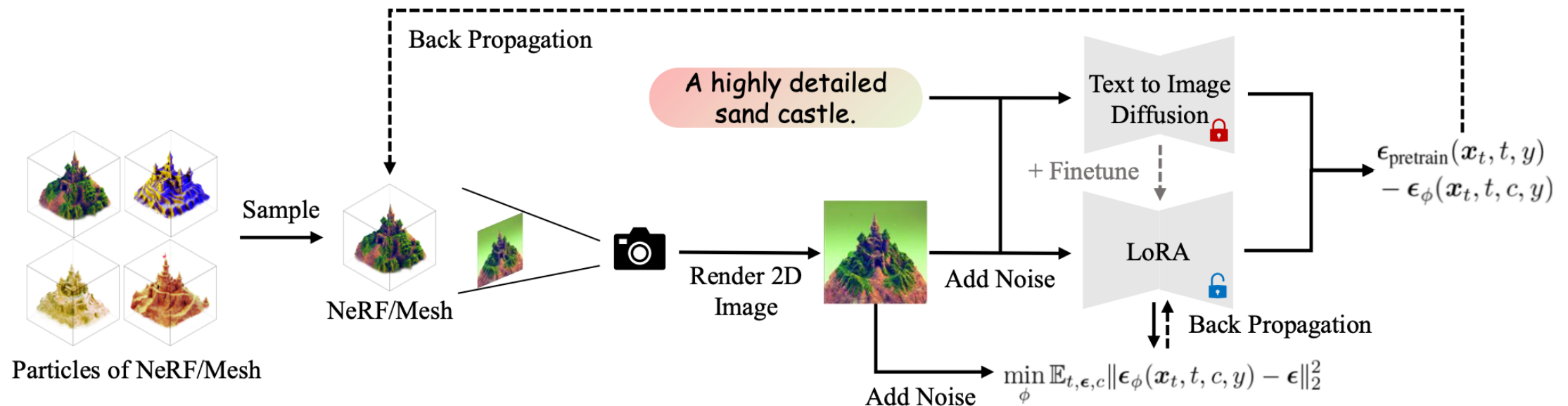
Magic3D [Lin et al., CVPR 2023]

- *Two-stage approach.*
- *Extract a coarse mesh in the first stage, and then texture the mesh in the second stage.*



ProlificDreamer [Wang et al., arXiv 2023]

- Minimize the SDS loss for the **multiple samples** of the NeRF parameters ϕ .
- Finetune the diffusion model with the **Low Rank Adaptation (LoRA)** technique.

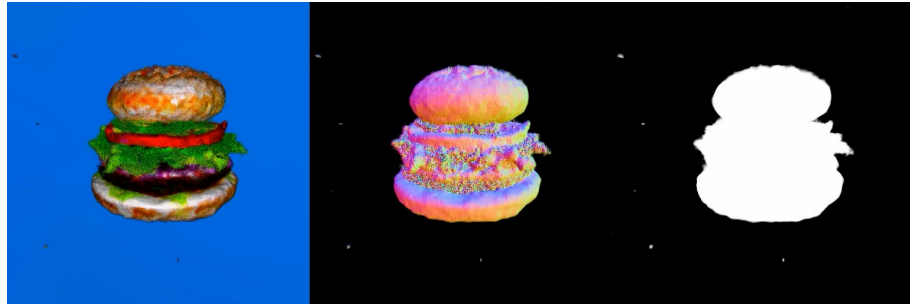


Limitation of SDS

It does not converge well without a **high CFG weight** (e.g., $w = 400$) and thus suffers from **model collapse**.

“a delicious hamburger”

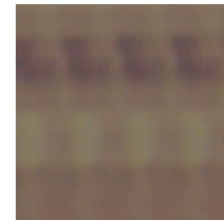
$w = 400$



$w = 7.5$



Credit: Jaihoon Kim



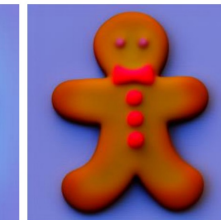
NeRF Initial.



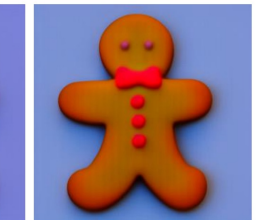
Seed=0



Seed=1



Seed=2

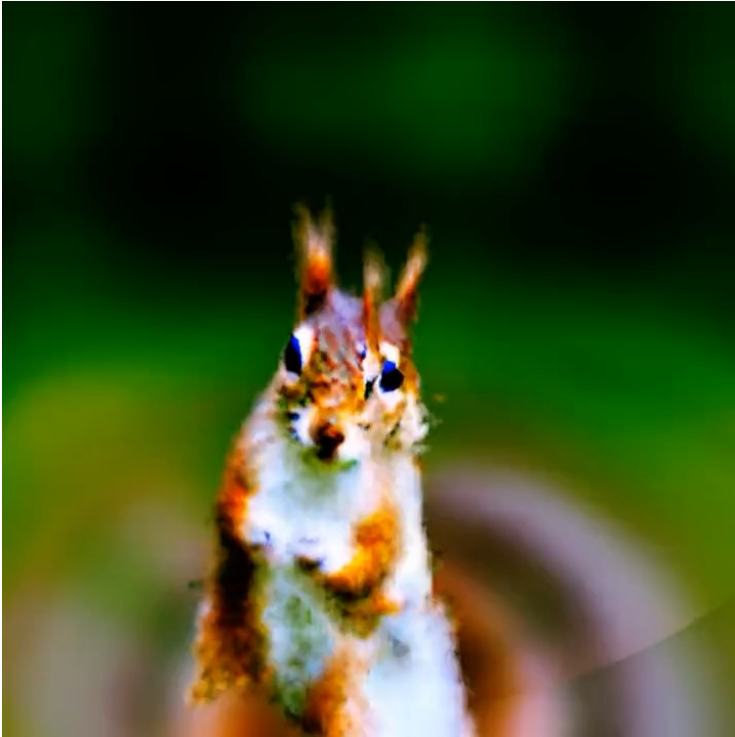


Seed=3

“gingerbread man”

Huang et al., *DreamTime: An Improved Optimization Strategy for Text-to-3D Content Creation*, arXiv 2023.

Limitation of 3D Generation from 2D Priors



Twitter @_akhaliq

Limitation of 3D Generation from 2D Priors

Supervision for **geometry** is still needed!



ProlificDreamer, Wang et al., 2023.

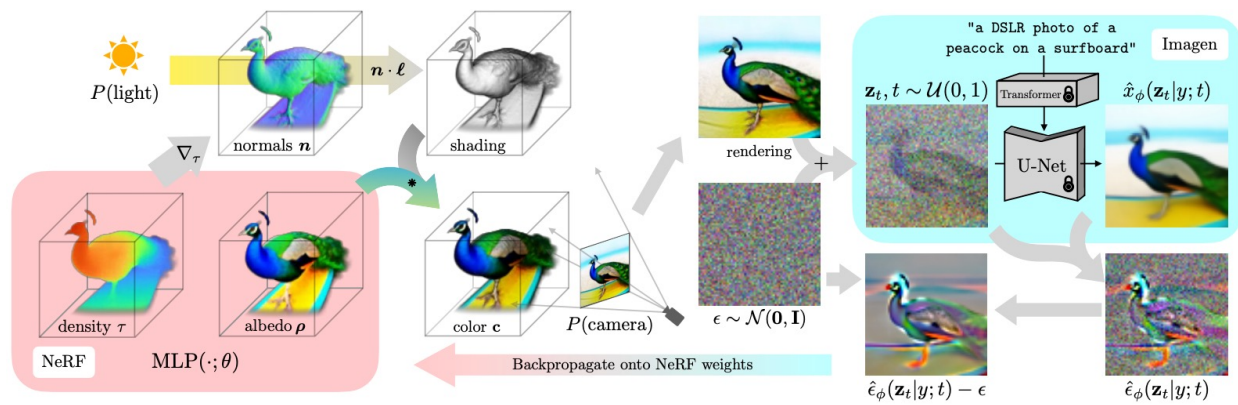


StableDreamFusion

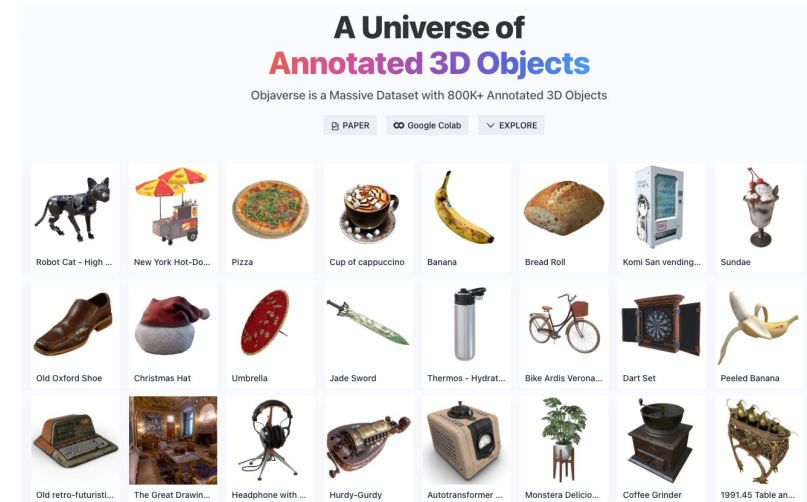
What's Next for 3D Generative Models?

1. Combining 3D Supervision

While *pretrained image generative models* will continue to be valuable for 3D generation, the key to producing realistic and solid 3D shapes would lie in utilizing *small-scale 3D priors*.



DreamFusion
Google

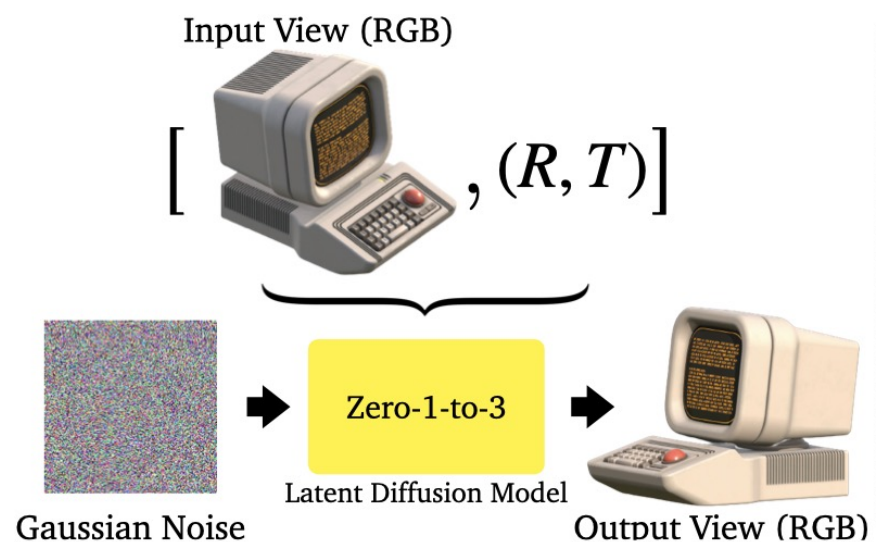


Objaverse
Allen Institute

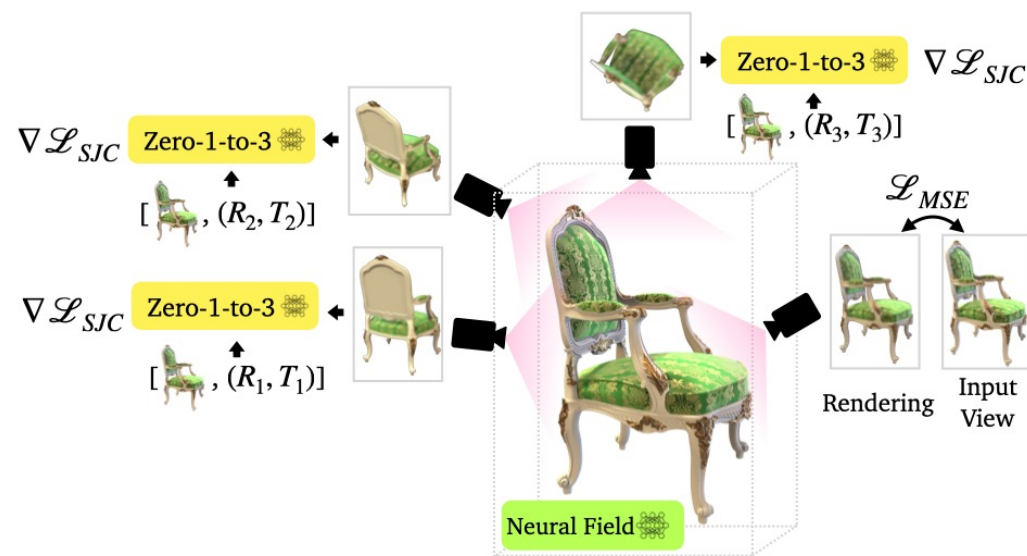
Zero-1-to-3 [Liu et al., 2023]

Novel view generation

An image diffusion model generating a novel view image conditioned by another view image and camera pose.



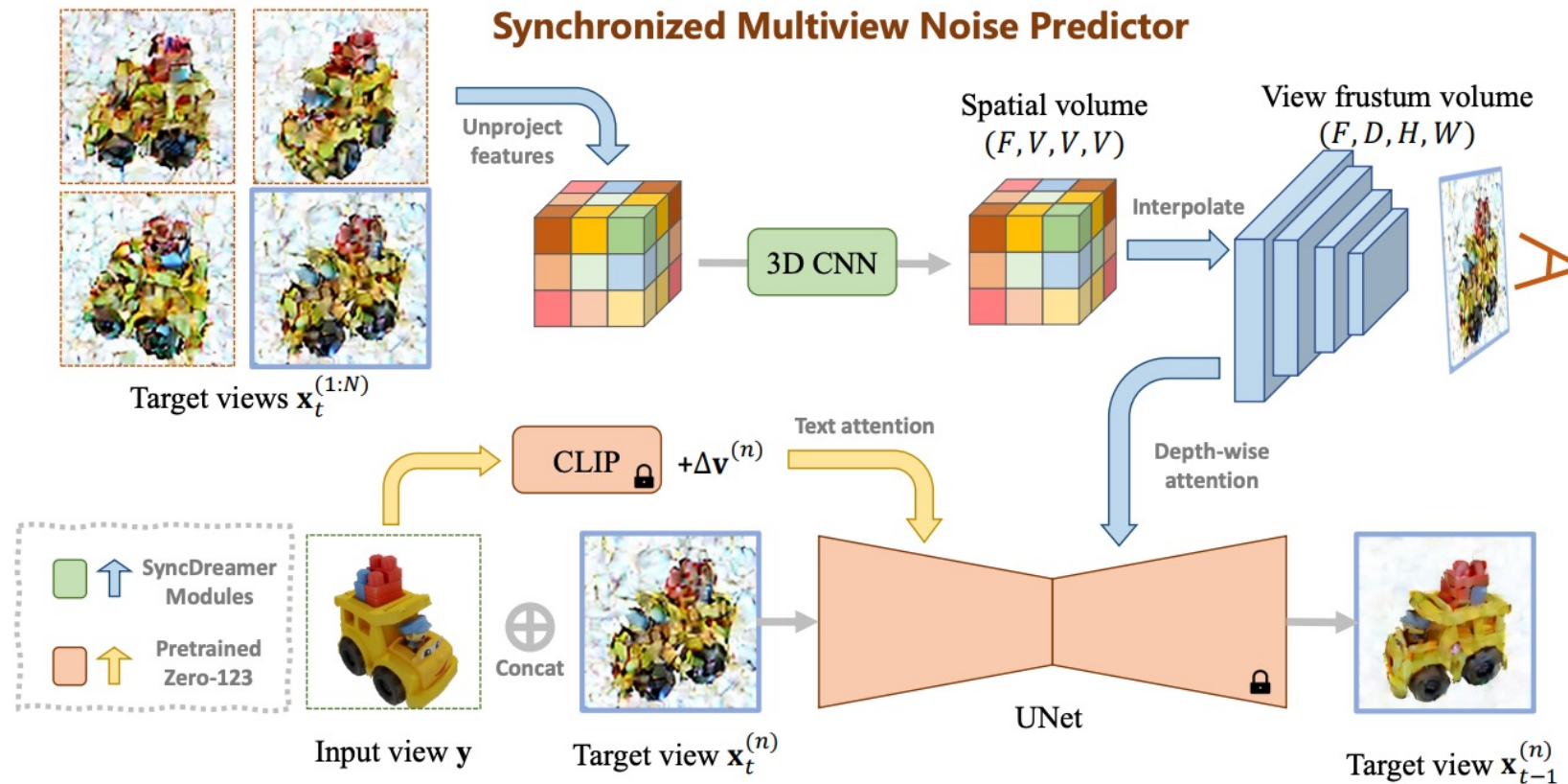
Novel View Synthesis



3D Reconstruction

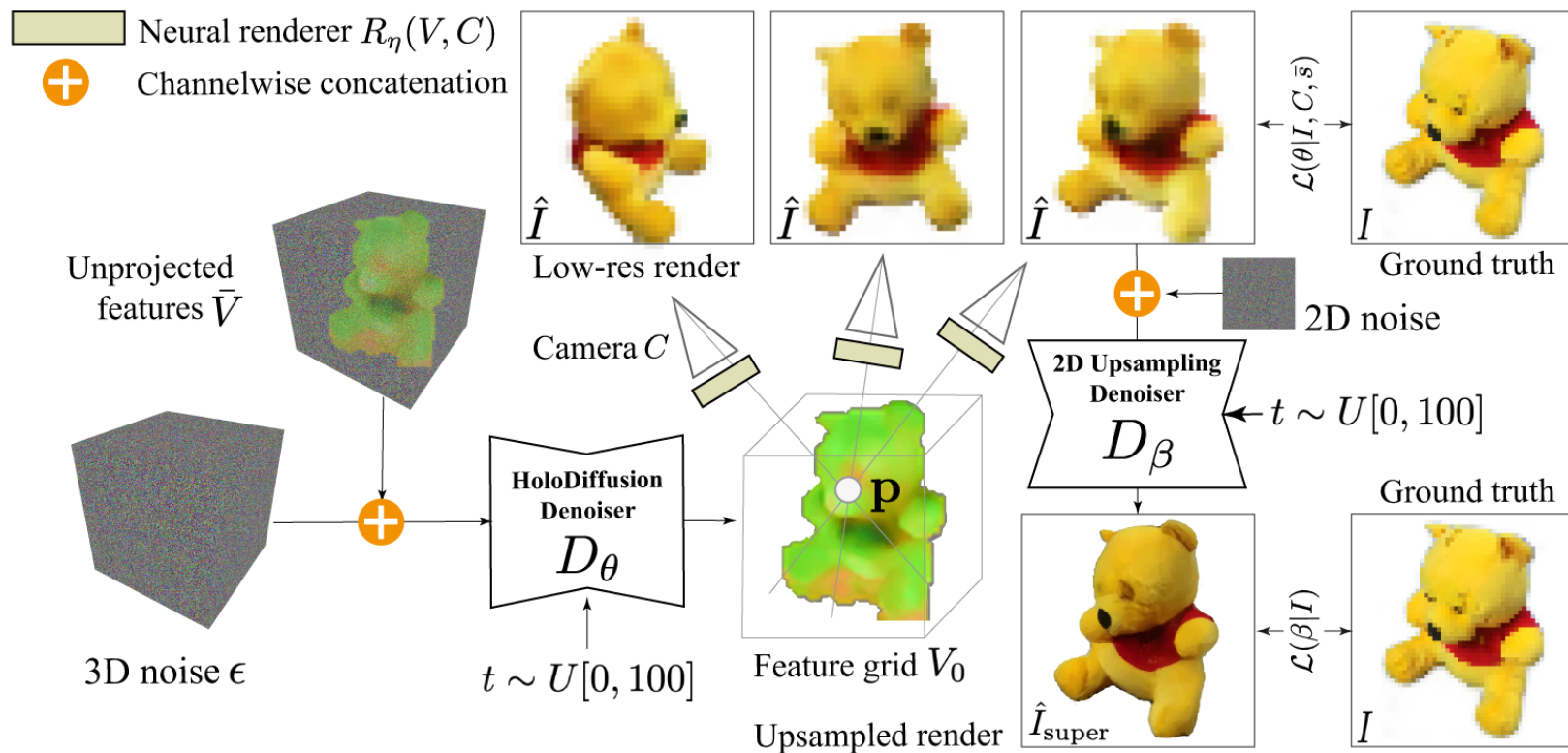
SyncDreamer [Liu et al., 2023]

Utilize Zero 1-to-3 to learn the joint probability distribution of multi-view images.



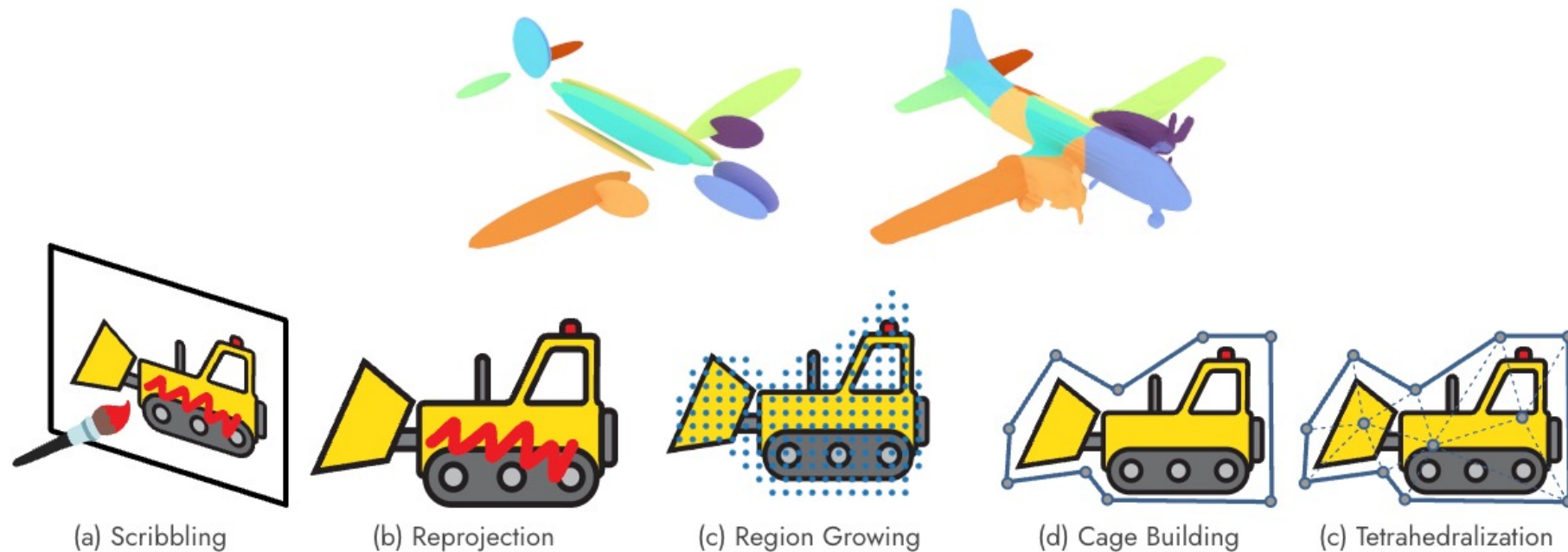
HoloFusion [Karnewar et al., 2023]

- Train a 3D diffusion model using multi-view images only.
- Can be extended to integrate 2D priors.



2. Generation → Editing

The focus of 3D generative models will shift towards creating *versatile models* capable of not only generating but *editing* and *manipulating 3D shapes*.



NeRFshop, Jambon et al., I3D 2023.

Delta Denoising Score [Hertz et al., 2023]

A new loss function for zero-shot *image* editing.



“Photo of sunset, winter, river.” → “Photo of sunset, winter, river with polar bears.”



“Amaryllis, North Window, oil on linen.” → “Cactus, North Window, oil on linen.”



“Forest painting, autumn trees.” → “Moon over forest painting, autumn trees.”

Posterior Distillation Sampling [Koo et al., 2024]

A new loss function for zero-shot **NeRF** editing.

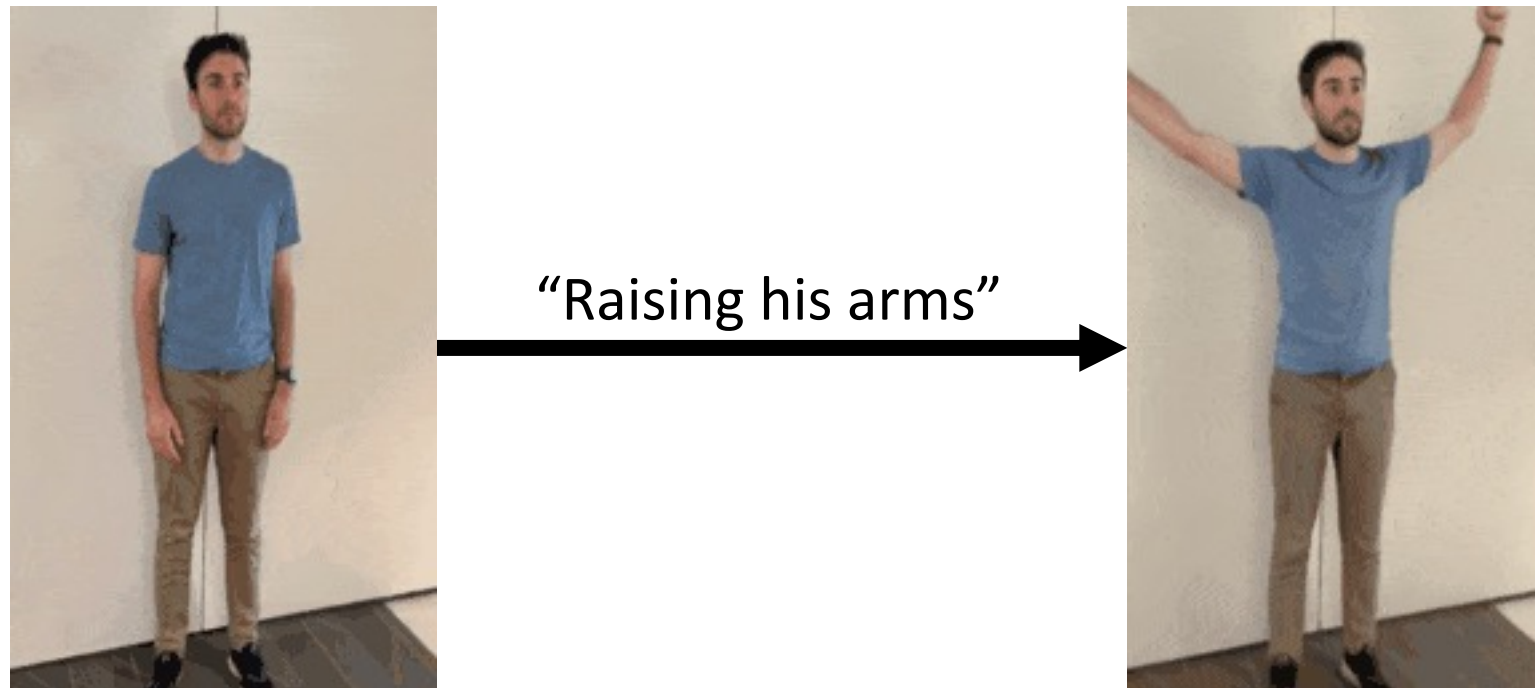


“Leonardo DiCaprio”



Posterior Distillation Sampling [Koo et al., 2024]

A new loss function for zero-shot **NeRF** editing.



3. Texture Generation

How can we generate *photorealistic texture* given a 3D mesh or Gaussian splats?

