

# Embodied Vision

Generative Modeling

Tsung-Wei Ke

Spring 2026



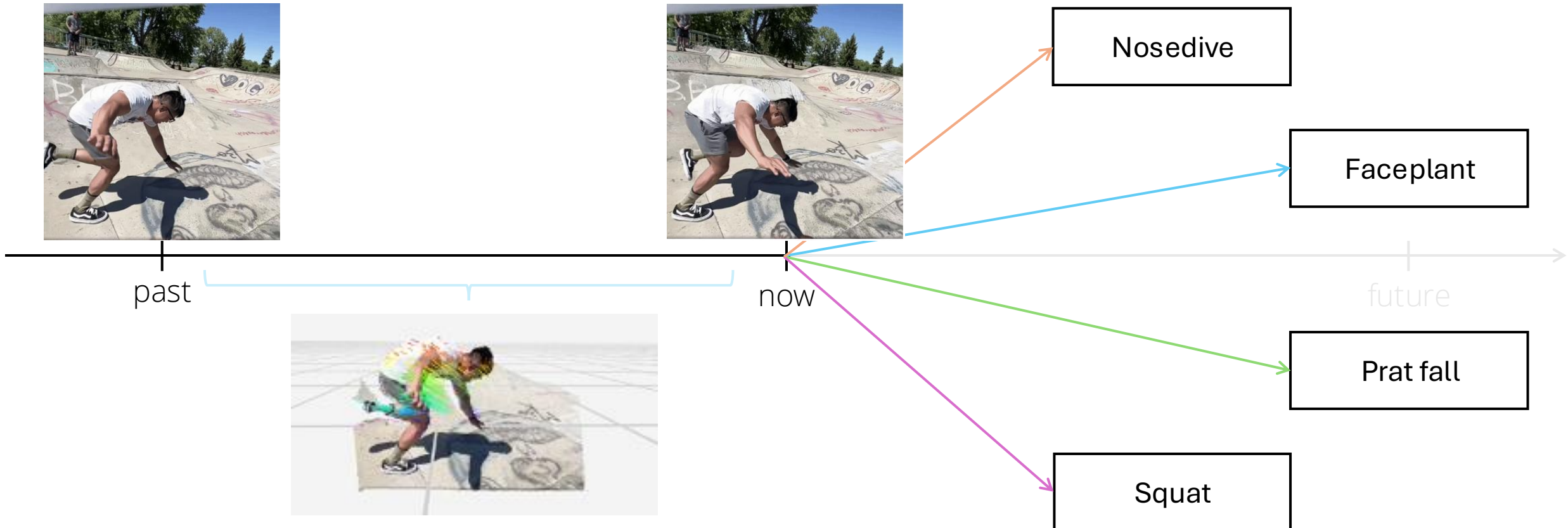
# Content

- Generative Modeling
  - Autoregressive Models
  - Variational AutoEncoders
  - Denoising Diffusion Probabilistic Models
  - Flow Matching Models
- 4D Generative Modeling
  - 3D Geometry Generation
  - 3D Generation with Multi-view Video Generation
  - 4D Geometry Generation
  - 4D Generation with Multi-view Video Generation

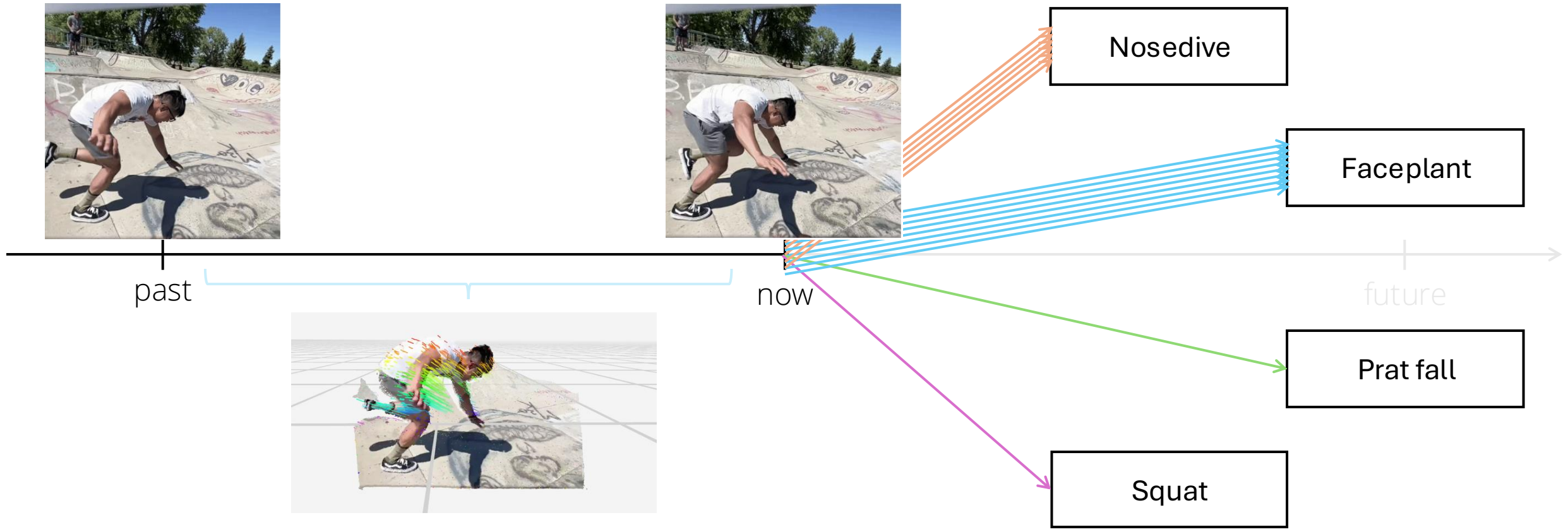
# We Inferred the Past Motion of an Observed Object... How About the Future Motion of the Object?



# Unfortunately, the Future is Not Unique. We May Have Multiple Plausible Futures

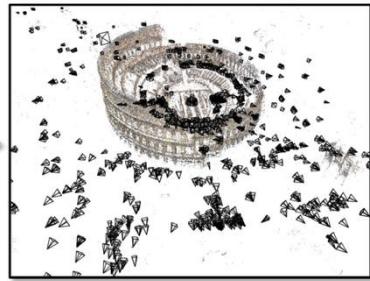


# Can We Use Generative Models to Predict Futures?



# Can We Use Generative Models to Predict Futures?

- Previously, we obtained 3D from 2D images; We obtained 4D from 2D videos.



- Can we generate 3D/4D world without 2D data?

# Can We Use Generative Models to Generate a 3D World?

- Previously, we obtained 3D from 2D images.



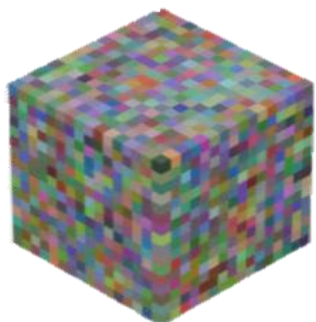
- How to generate a 3D world?
  - ~~Generate 3d assets directly → The 3D generative is not competent yet, as the amount of available 3D data is still small~~
  - Generate multi-view 2D images and lift to 3D → Yes! 2D generation is already great, since we have a huge amount of 2D images available on the Internet.

# Content

- Generative Modeling
  - Autoregressive Models
  - Variational AutoEncoders
  - Denoising Diffusion Probabilistic Models
  - Flow Matching Models
- 4D Generative Modeling
  - 3D Geometry Generation
  - 3D Generation with Multi-view Video Generation
  - 4D Geometry Generation
  - 4D Generation with Multi-view Video Generation

# Generate 3D Assets with Diffusion Models

- Idea: Generate 3D assets with diffusion models
- The availability and scalability of training data?
- The choice of 3D representations
  - How to model shape and visual appearance?
  - Do we need a joint representations?



# The availability and scalability of training data



Source	# Objects
IKEA [32]	219
GSO [17]	1K
EGAD [41]	2K
OmniObject3D [63]	6K
PhotoShape [46]	5K
ABO [13]	8K
Thingi10K [67]	10K
3d-Future [19]	10K
ShapeNet [9]	51K
Objaverse 1.0 [14]	800K
<b>Objaverse-XL</b>	<b>10.2M</b>

# Representations

- Representations of shapes
  - Voxel grids
  - Point clouds
  - Meshes
  - 3D GS
  - ...
- Representations of visual appearance
  - Texture maps
  - RGB
  - Spherical harmonics
  - ...
- What are better design choices?



# Structured 3D Latents for Scalable and Versatile 3D Generation

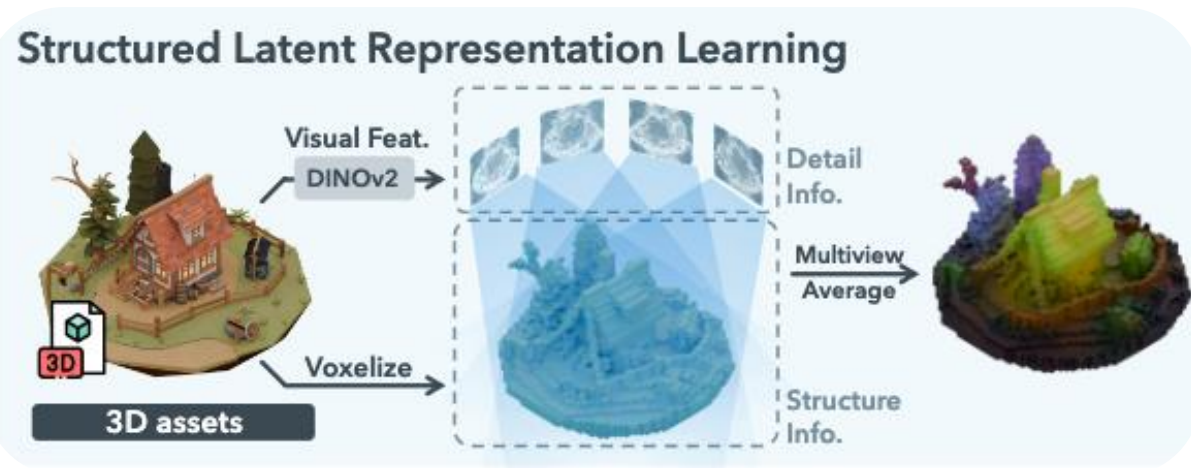
Jianfeng Xiang<sup>1,3</sup> Zelong Lv<sup>2,3</sup> Sicheng Xu<sup>3</sup> Yu Deng<sup>3</sup> Ruicheng Wang<sup>2,3</sup>  
Bowen Zhang<sup>2,3</sup> Dong Chen<sup>3</sup> Xin Tong<sup>3</sup> Jiaolong Yang<sup>3</sup>

<sup>1</sup>Tsinghua University    <sup>2</sup>USTC    <sup>3</sup>Microsoft Research

CVPR 2025 Highlight

# Structured Latents

- A unified 3D latent representations for high-quality 3D generation via **active voxels**
  - Represent shape by active voxel positions  $\{p_i\}_{i=1}^L$ , which are sparse voxels intersecting the object's surface
  - Represent appearance by image features  $\{f_i\}_{i=1}^L$ , which extracted from densely rendered views and attached onto active voxels



# Issue: Image Features are Overly High Dimensional

- DINOv2 is used to extract image features, which has 2048 dimensions. Generating high-dimensional data is computationally expensive
- Idea: compress image features  $\{f_i\}_{i=1}^L$  to low-dimensional latents  $\{z_i\}_{i=1}^L$  with VAE by reconstructing (1) 3D GS, (2) Meshes and (3) Radiance fields

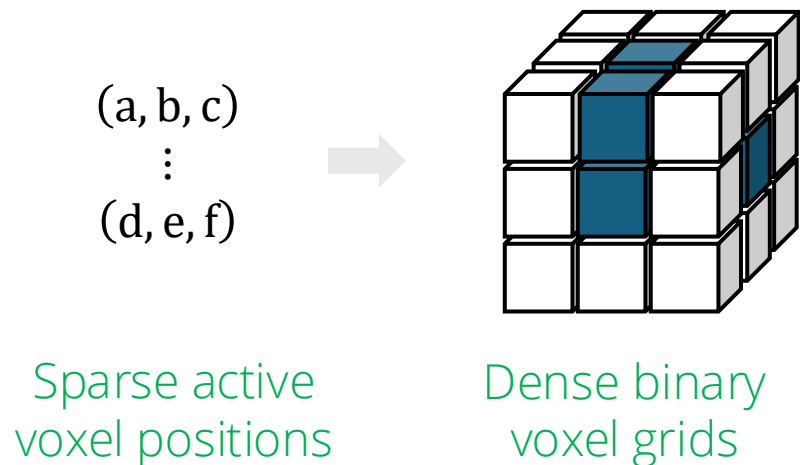


# Generate Unified 3D Latent Representations

- A unified 3D latent representation  $\{(p_i, z_i)\}_{i=1}^L$  for high-quality 3D generation via **active voxels**, with active voxel positions  $\{p_i\}_{i=1}^L$  and visual latents  $\{z_i\}_{i=1}^L$
- A two-stage generation pipeline: first generate sparse structures  $\{p_i\}_{i=1}^L$  then generate visual representations  $\{z_i\}_{i=1}^L$
- Consider flow matching as the generative model

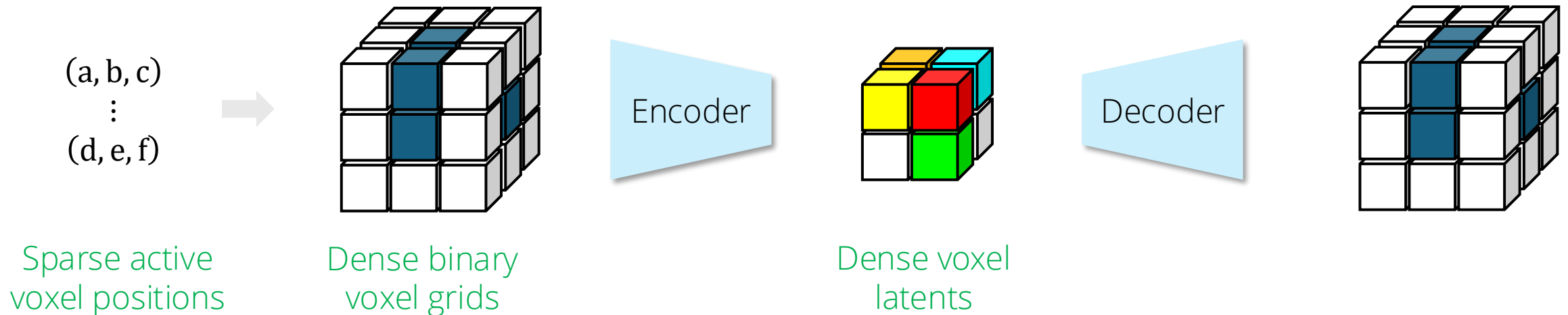
# How to Generate Sparse Structures $\{p_i\}_{i=1}^L$

- Challenge 1: Active voxel positions  $\{p_i\}_{i=1}^L$  are “set” representations, while generating such unordered data with diffusion models is challenging
- Solution 1: Reorganize active voxels as binary voxel grids



# How to Generate Sparse Structures $\{p_i\}_{i=1}^L$

- Challenge 1: Active voxel positions  $\{p_i\}_{i=1}^L$  are “set” representations, while generating such unordered data with diffusion models is challenging
- Solution 1: Reorganize active voxels as binary voxel grids
- Challenge 2: Voxel grids are costly to generate
- Solution 2: Compress voxel grids into low-dim voxel latents with VAE and learn diffusion models to generate these latents

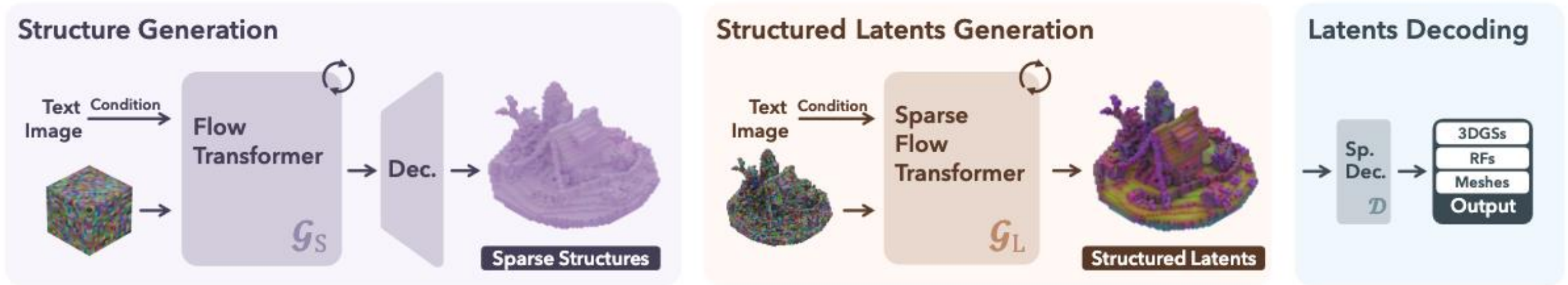


# How to Generate Visual Representations $\{z_i\}_{i=1}^L$

- Challenge 1: Visual representations  $\{z_i\}_{i=1}^L$  are “set” representations, while generating such unordered data with diffusion models is challenging
- Solution 1: Reorganize active voxels as  $N \times N \times N \times d$  voxel grids ( $d$  denotes the dimension of  $z_i$ )
- Challenge 2: The  $N \times N \times N \times d$  voxel grids is overly high-dimensional
- Solution 2: Compress the  $N \times N \times N \times d$  voxel grids to lower-dimensional voxel grids with sparse 3D convolution

# Generate Unified 3D Latent Representations

## 3D Assets Generation



# Results: Text to 3D Asset



*A rustic log cabin with a stone chimney and a wooden porch.*



*Portable transistor radio, dark cover, speaker grille, brand logo on front.*



*Futuristic red toy blaster with transparent magazine.*



*Blocky, orange and teal robot with articulated limbs.*

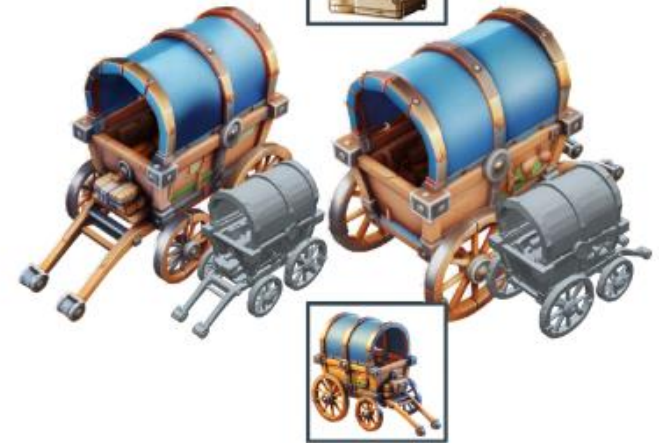


*Metallic dog-like robot with articulated legs and futuristic design elements.*



*Yellow and black bulldozer with movable front blade.*

# Results: Image to 3D Asset



# Results: 3D Asset Editing



*Cozy log cabin*



*Mossy stone bricks*



*Sci-fi inspired*



*Toy blocks*

Input



Remove



- Arm

Add



+ Weapon

Replace



Leg → Track

# Results: Composable 3D World

Wooden and iron chest.

Wooden crate.

Round dark wooden table.

Phonograph.



# SAM 3D: 3Dfy Anything in Images

**SAM 3D Team, Xingyu Chen\***, **Fu-Jen Chu\***, **Pierre Gleize\***, **Kevin J Liang\***, **Alexander Sax\***, **Hao Tang\***, **Weiyao Wang\***, **Michelle Guo**, **Thibaut Hardin**, **Xiang Li<sup>◦</sup>**, **Aohan Lin**, **Jiawei Liu**, **Ziqi Ma<sup>◦</sup>**, **Anushka Sagar**, **Bowen Song<sup>◦</sup>**, **Xiaodong Wang**, **Jianing Yang<sup>◦</sup>**, **Bowen Zhang<sup>◦</sup>**, **Piotr Dollár<sup>†</sup>**, **Georgia Gkioxari<sup>†</sup>**, **Matt Feiszli<sup>†§</sup>**, **Jitendra Malik<sup>†§</sup>**

Meta Superintelligence Labs

\*Core Contributor (Alphabetical, Equal Contribution), <sup>◦</sup>Intern, <sup>†</sup>Project Lead, <sup>§</sup>Equal Contribution

We present SAM 3D, a generative model for visually grounded 3D object reconstruction, predicting geometry, texture, and layout from a single image. SAM 3D excels in natural images, where occlusion and scene clutter are common and visual recognition cues from context play a larger role. We achieve this with a human- and model-in-the-loop pipeline for annotating object shape, texture, and pose, providing visually grounded 3D reconstruction data at unprecedented scale. We learn from this data in a modern, multi-stage training framework that combines synthetic pretraining with real-world alignment, breaking the 3D “data barrier”. We obtain significant gains over recent work, with at least a 5 : 1 win rate in human preference tests on real-world objects and scenes. We will release our code and model weights, an online demo, and a new challenging benchmark for in-the-wild 3D object reconstruction.

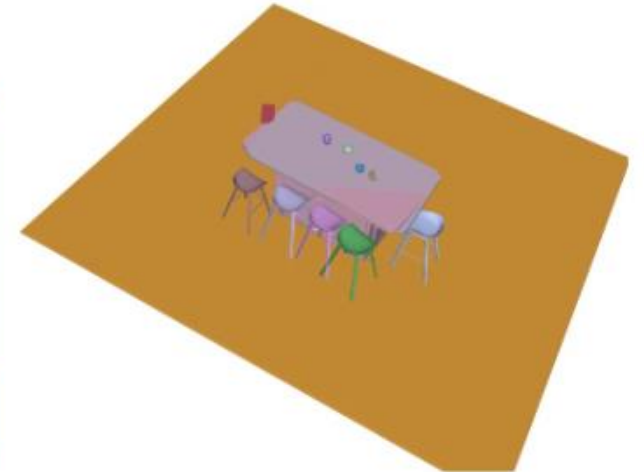
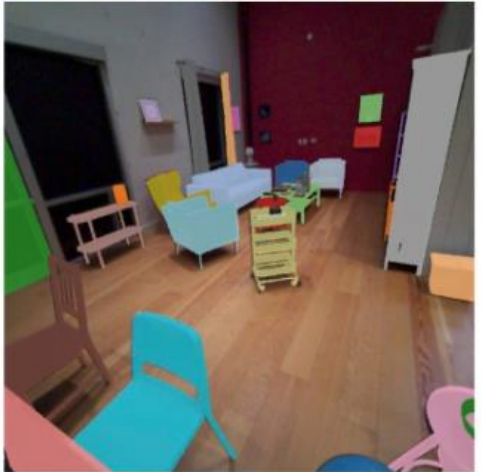
**Demo:** <https://www.aidemos.meta.com/segment-anything/editor/convert-image-to-3d>

**Code:** <https://github.com/facebookresearch/sam-3d-objects>

**Website:** <https://ai.meta.com/sam3d>



# Results: Reconstruct Digital Twins from RGB Images



# However, 3D Asset Generation Has Scale Issues



10M+ object meshes

Is it possible to generate 3D from image generation?  
Will it be better than 3D asset generation?

Backend url:  Index:

Search: french cat

Clip retrieval works by converting the text query to a CLIP embedding, then using that embedding to query a kmn index of clip image embeddings

Display captions   
Display full captions   
Display similarities   
Safe mode   
Hide duplicate urls   
Hide (near) duplicate images   
Search over   
Search with multilingual clip

french cat

french cat

How to tell if your feline is french. He wears a b...

Hilarious pics of funny cats! funnycatsgif.com

イケメン猫モデル「トキ・ナンタケット」がかっこいい - NAVER まとめ

hipster cat

cat in a suit Georgian sells tomatoes

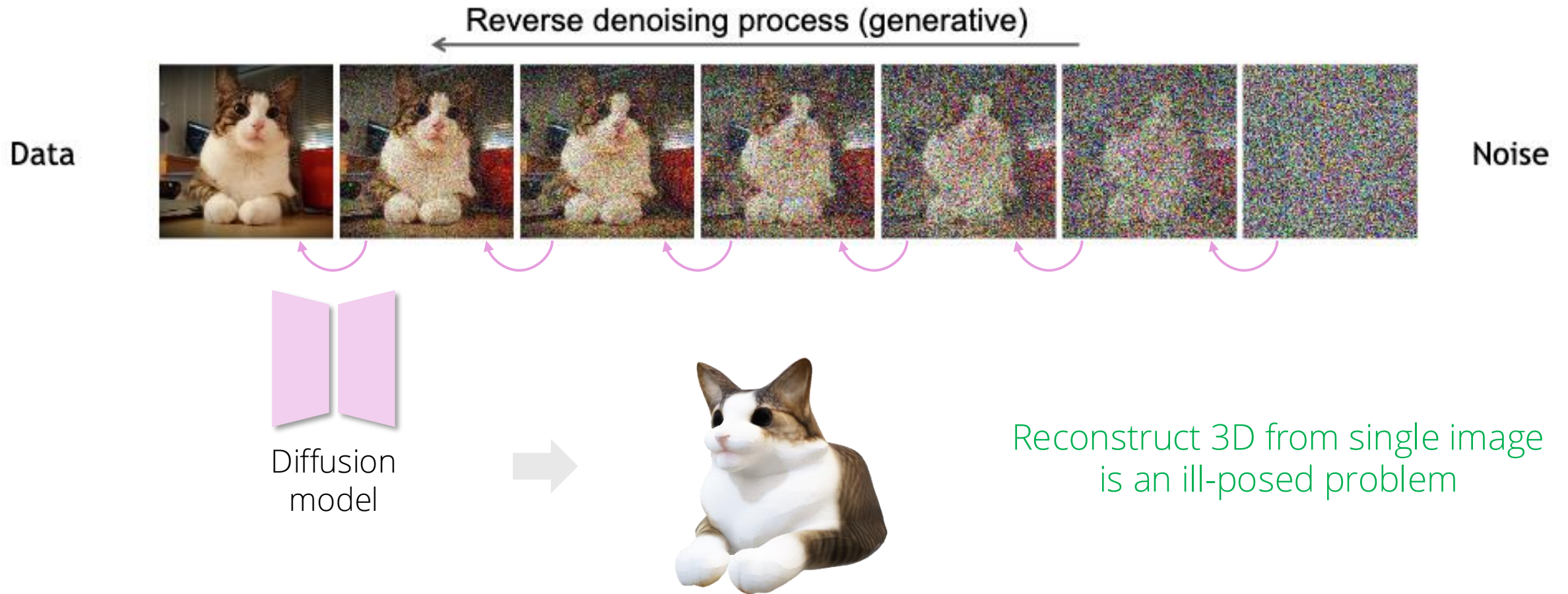
French Bread Cat Loaf Metal Print

5B+ text-image pairs

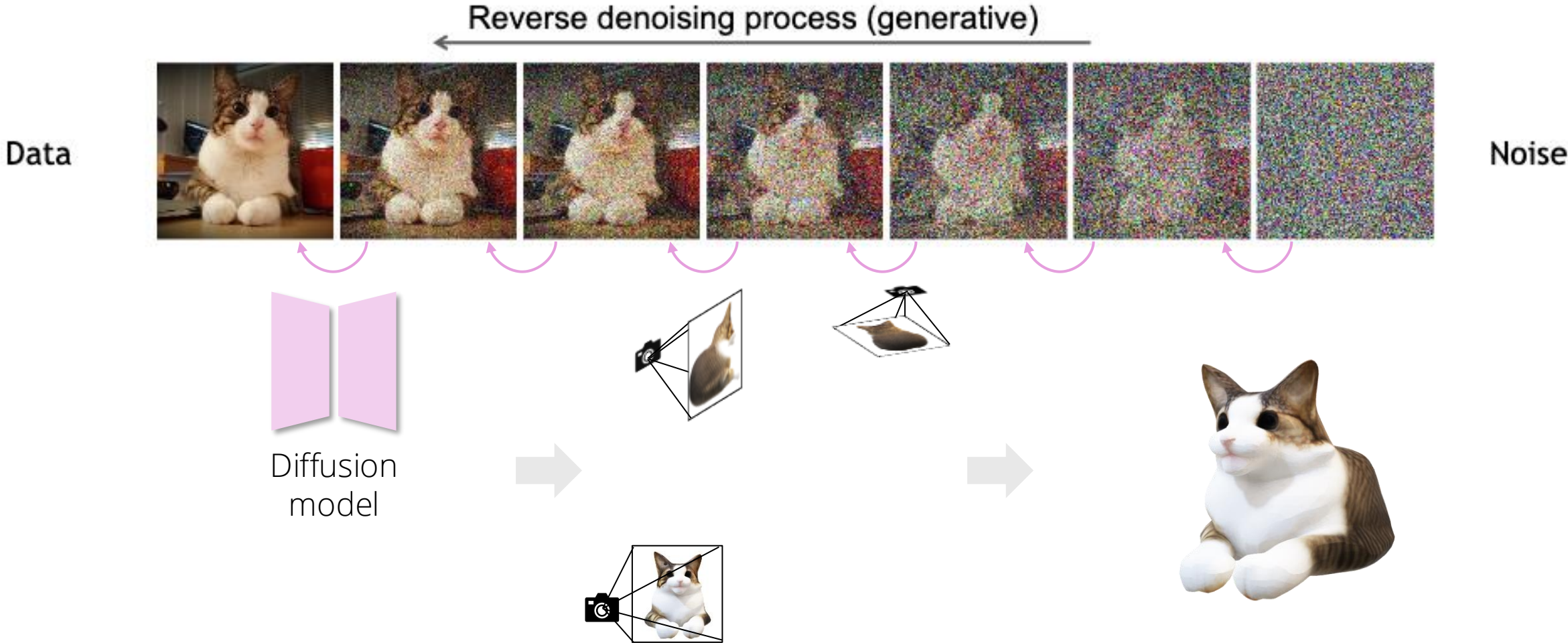
# Content

- Generative Modeling
  - Autoregressive Models
  - Variational AutoEncoders
  - Denoising Diffusion Probabilistic Models
  - Flow Matching Models
- 4D Generative Modeling
  - 3D Geometry Generation
  - 3D Generation with Multi-view Video Generation
  - 4D Geometry Generation
  - 4D Generation with Multi-view Video Generation

# Can We Obtain 3D Object from 2D Image Generation?



# Can We Generate Multiple Views of the Input Image from Image Diffusion Model?



# Let's Revisit The Learning Objective of Diffusion Models

- $p_\theta(x_0) := \int p_\theta(x_{0:T}) dx_{0:T}$
- $p_\theta(x_{0:T}) := p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)$
- Diffusion models minimize the variational bound (ELBO) on negative log likelihood

$$\begin{aligned} \mathbb{E}[-\log p_\theta(x_0)] &\leq \mathbb{E}_q \left[ -\log \frac{p_\theta(x_{0:T})}{q(x_{1:T}|x_0)} \right] = \mathbb{E}_q \left[ -\log p(x_T) - \sum_{t \geq 1} \log \frac{p_\theta(x_{t-1}|x_t)}{q(x_t|x_{t-1})} \right] \\ &= \mathbb{E}_q \left[ \underbrace{D_{KL}(q(x_T|x_0) \| p_\theta(x_T))}_{L_T} + \sum_{t \geq 1} \underbrace{D_{KL}(q(x_{t-1}|x_t, x_0) \| p_\theta(x_{t-1}|x_t))}_{L_{t-1}} - \underbrace{\log p_\theta(x_0|x_1)}_{L_0} \right] \end{aligned}$$

# The Learning Objective of Diffusion Models

- $\tilde{\mu}_t(x_t, x_0) := \frac{\sqrt{\bar{\alpha}_{t-1}\beta_t}}{1-\bar{\alpha}_t} x_0 + \frac{\sqrt{\bar{\alpha}_t(1-\bar{\alpha}_{t-1})}}{1-\bar{\alpha}_t} x_t$
- $x_t(x_0, \epsilon) = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1-\bar{\alpha}_t} \epsilon$
- To maximize the likelihood  $p_\theta(x_0)$ , we need to minimize the negative log likelihood of intermediate diffusion samples

$$L_{t-1} = \mathbb{E}_q [D_{KL}(q(x_{t-1}|x_t, x_0) \| p_\theta(x_{t-1}|x_t))] = \mathbb{E}_q \left[ \frac{1}{2\sigma_t^2} \|\tilde{\mu}_t(x_t, x_0) - \mu_\theta(x_t, t)\|^2 \right] + C$$

Do you know how to show this?  
Check the Eq 8 and 9 of the paper

This results from how we  
parameterize our diffusion model

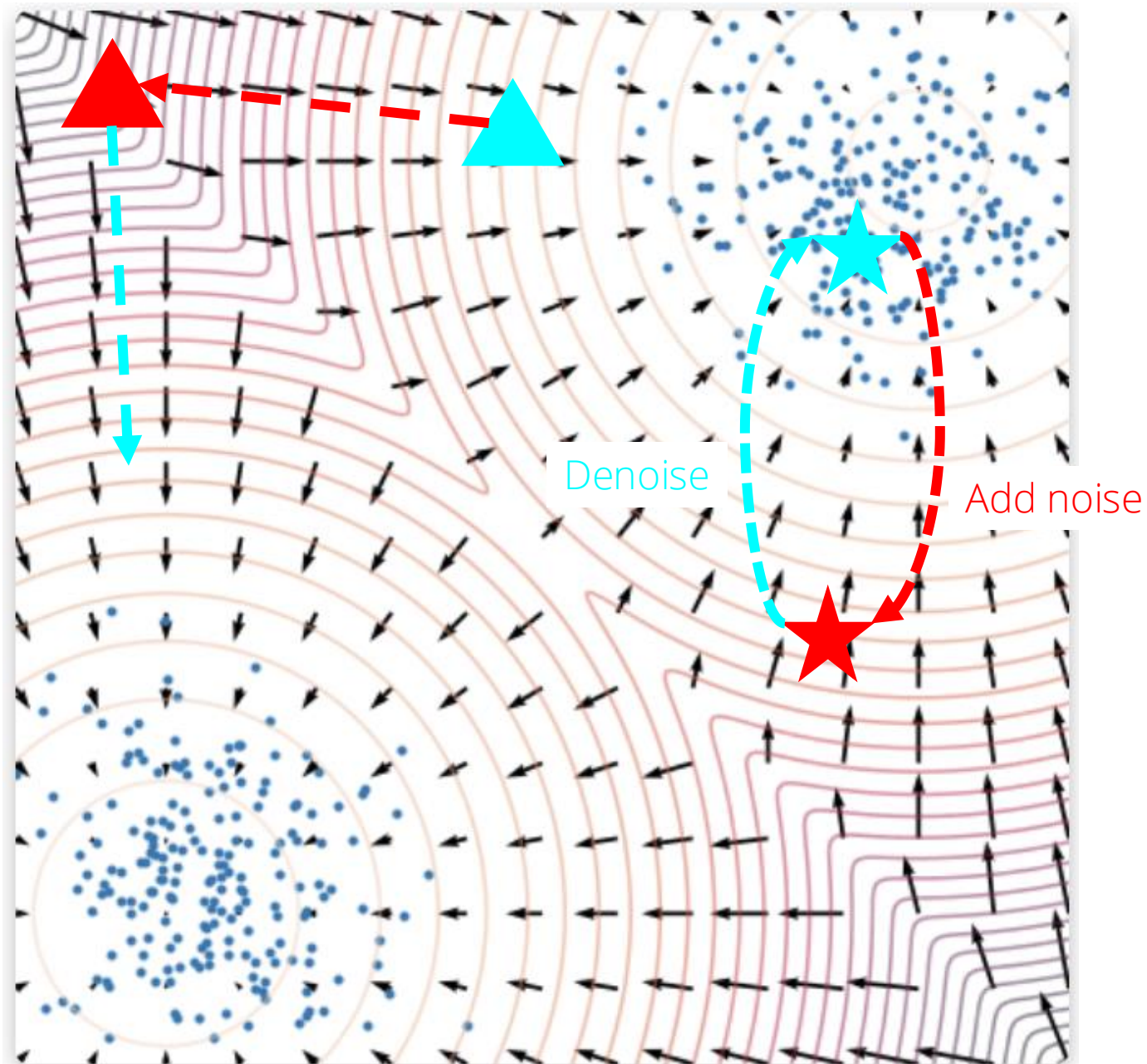
$$L_{t-1} \propto \mathbb{E}_{x_0, \epsilon} \left[ \left\| \frac{1}{\sqrt{\alpha_t}} \left( x_t(x_0, \epsilon) - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon \right) - \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(x_t, t) \right) \right\|^2 \right]$$

# We Can Optimize the Generated Sample by Minimizing the Intermediate Diffusion Loss

- $x_t(x_0, \epsilon) = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon$
- To maximize the likelihood  $p_\theta(x_0)$ , we need to minimize the negative log likelihood of intermediate diffusion samples

$$L_{t-1} \propto \mathbb{E}_{x_0, \epsilon} [\|\epsilon - \epsilon_\theta(x_t, t)\|^2] = \mathbb{E}_{x_0, \epsilon} [\|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)\|^2]$$

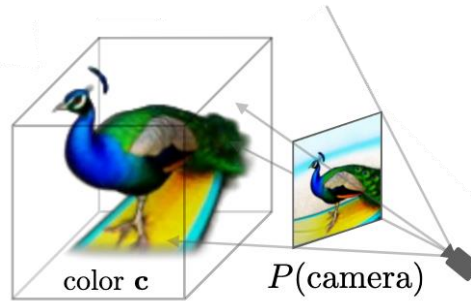
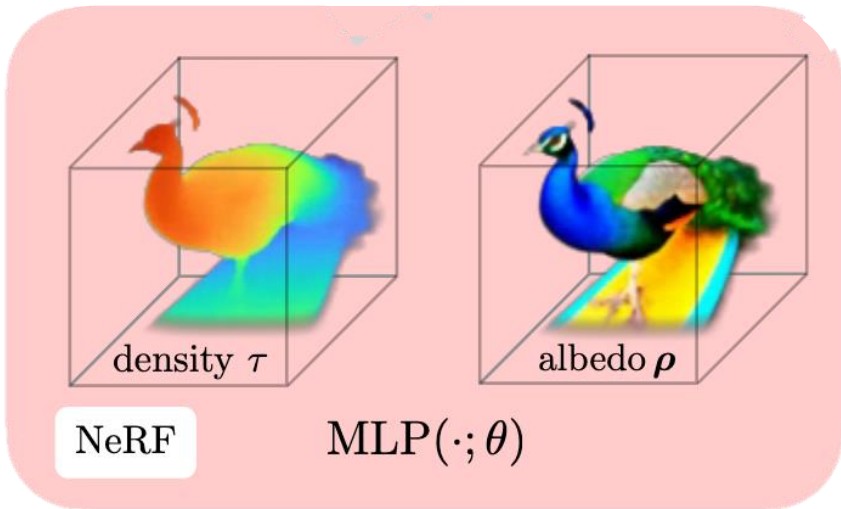
- But here, instead of optimizing the diffusion model parameterized by  $\theta$ , we aim to optimize  $x_0$ . In other words, if  $x_0$  is sampled from the same distribution (of the training data), the diffusion loss would be small.



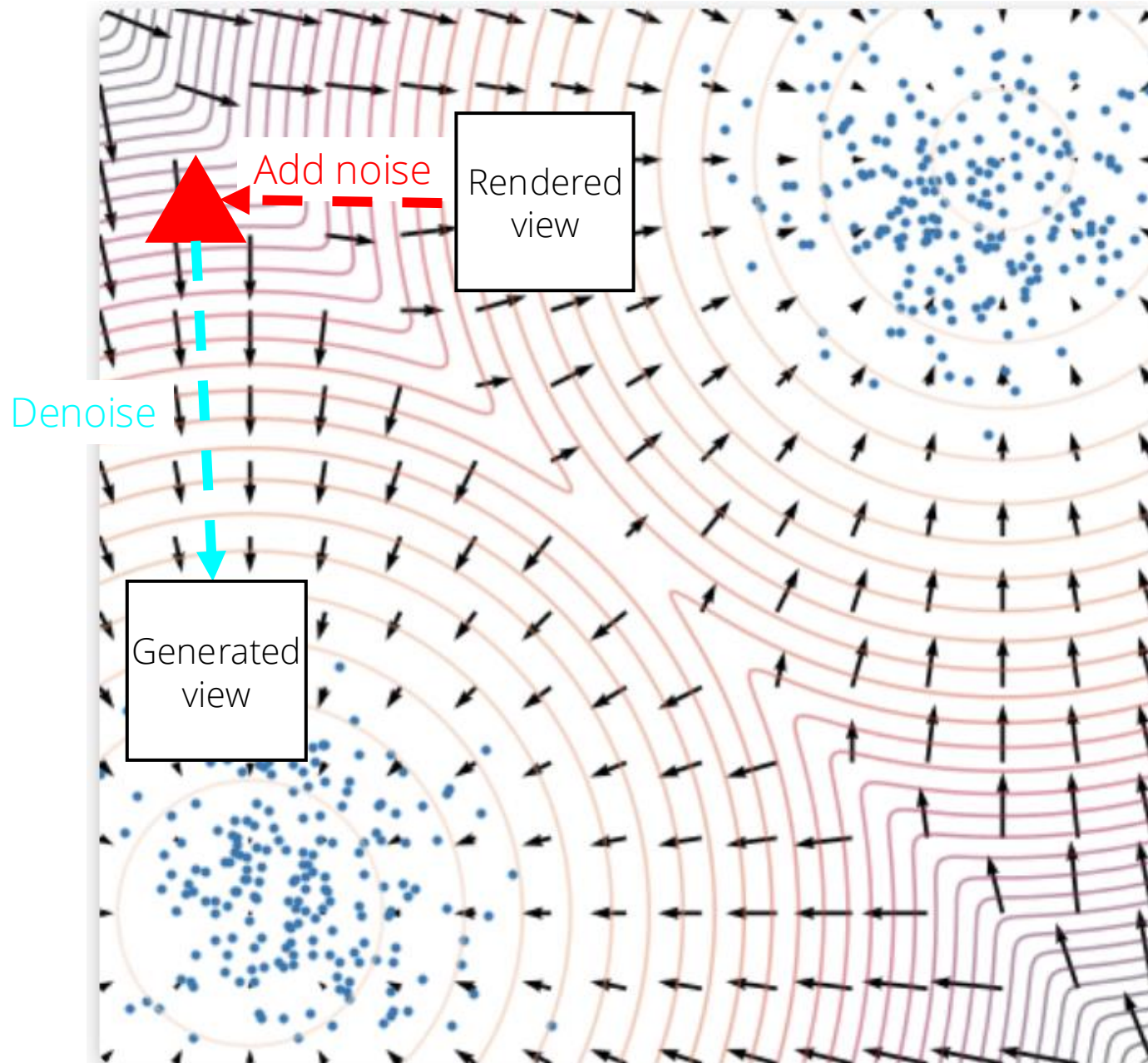
Observation: Noising good  $x_0$  is more easily to denoise back to clean samples than bad  $x_0$ .

Idea: optimizing denoising loss to search good  $x_0$

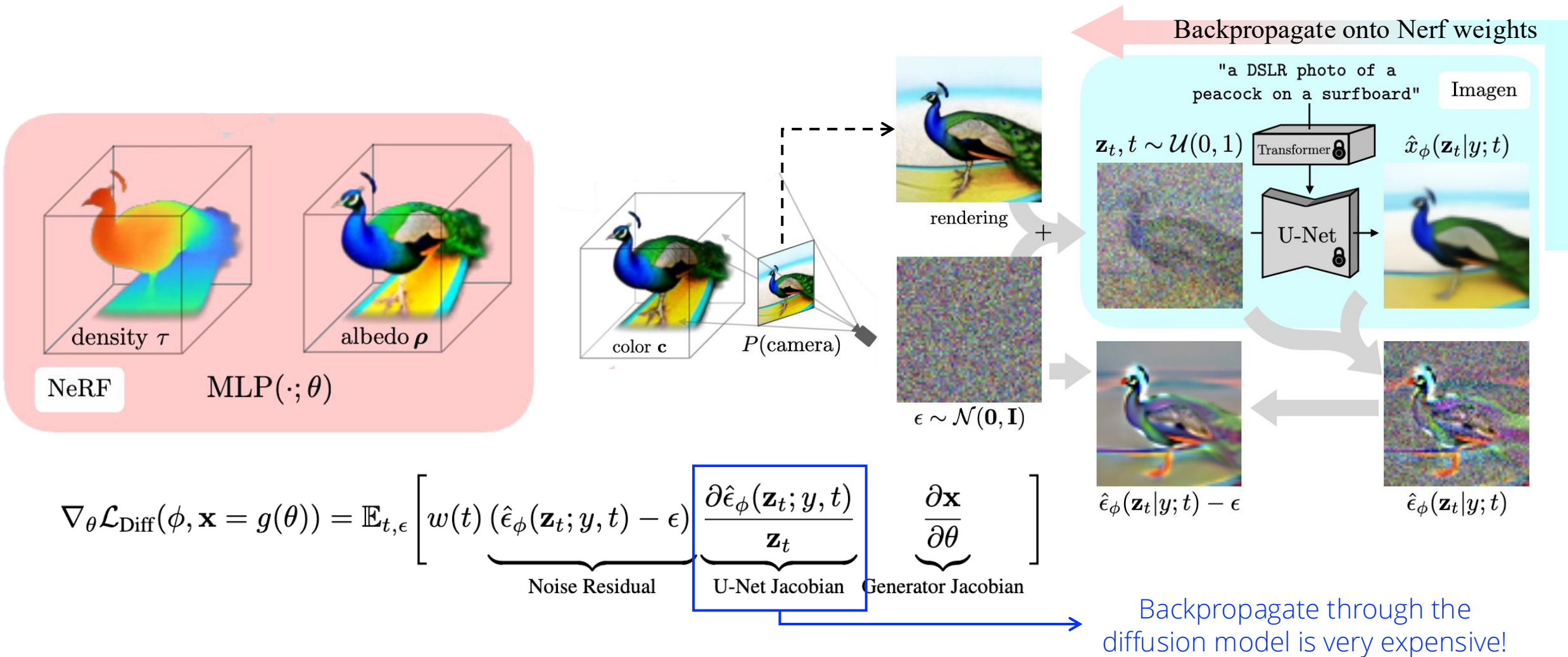
# 3D-from-2D Reconstruction Recipe: Render a View from Differentiable 3D Model



- **Goal:** A 3D volume to be optimized, which can be voxel grids, NeRF, or 3DGS
- **Idea:** optimize the differentiable 3D volume with the prior knowledge in 2D generative models
- **Supervision signal:** rendering loss between rendered and generated images

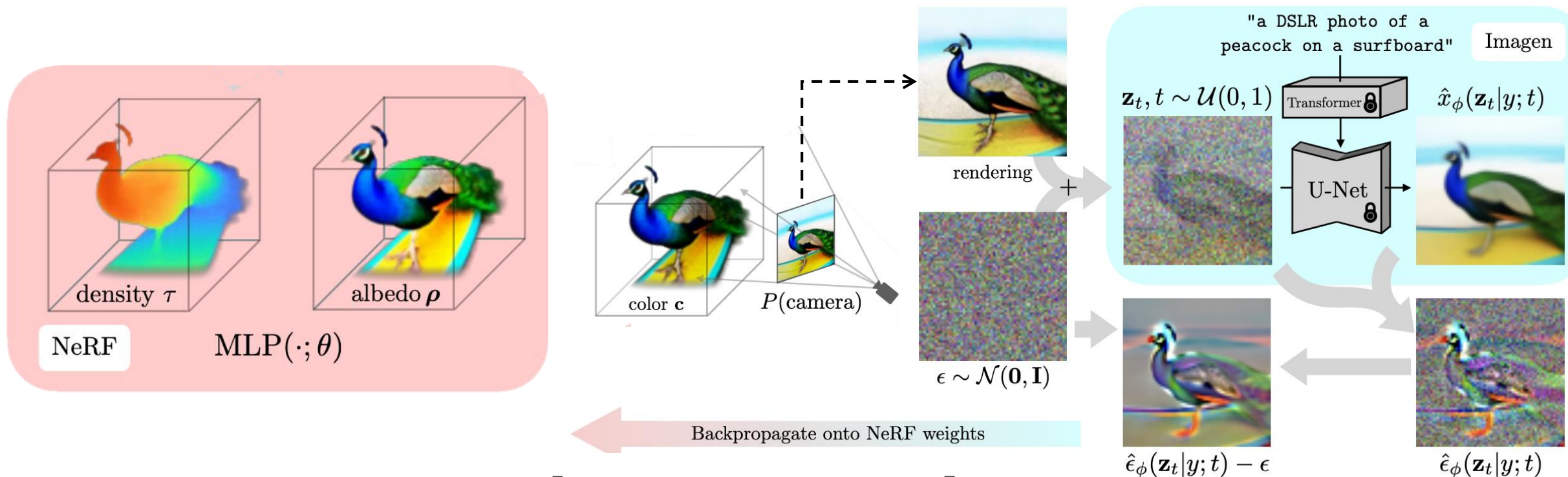


# 3D-from-2D Reconstruction Recipe: Render a View from Differentiable 3D Model



$$\nabla_{\theta} \mathcal{L}_{\text{Diff}}(\phi, \mathbf{x} = g(\theta)) = \mathbb{E}_{t, \epsilon} \left[ \underbrace{w(t)}_{\text{Noise Residual}} \underbrace{(\hat{\epsilon}_\phi(\mathbf{z}_t; y, t) - \epsilon)}_{\text{U-Net Jacobian}} \underbrace{\frac{\partial \hat{\epsilon}_\phi(\mathbf{z}_t; y, t)}{\partial \mathbf{z}_t}}_{\text{Generator Jacobian}} \right]$$

# 3D-from-2D Reconstruction Recipe



$$\nabla_{\theta} \mathcal{L}_{\text{SDS}}(\phi, \mathbf{x} = g(\theta)) \triangleq \mathbb{E}_{t, \epsilon} \left[ w(t) (\hat{\epsilon}_{\phi}(\mathbf{z}_t; \mathbf{y}, t) - \epsilon) \frac{\partial \mathbf{x}}{\partial \theta} \right]$$

Approximate the gradients by bypassing the whole diffusion model

# Results: Distillation of 3D from Image Diffusion

Problem: Without any 3D prior knowledge, generated views are inconsistent



an orangutan making a clay bowl on a throwing wheel\*



a raccoon astronaut holding his helmet†



a blue jay standing on a large basket of rainbow macarons\*



a corgi taking a selfie\*



a table with dim sum on it†



a lion reading the newspaper\*



Michelangelo style statue of dog reading news on a cellphone

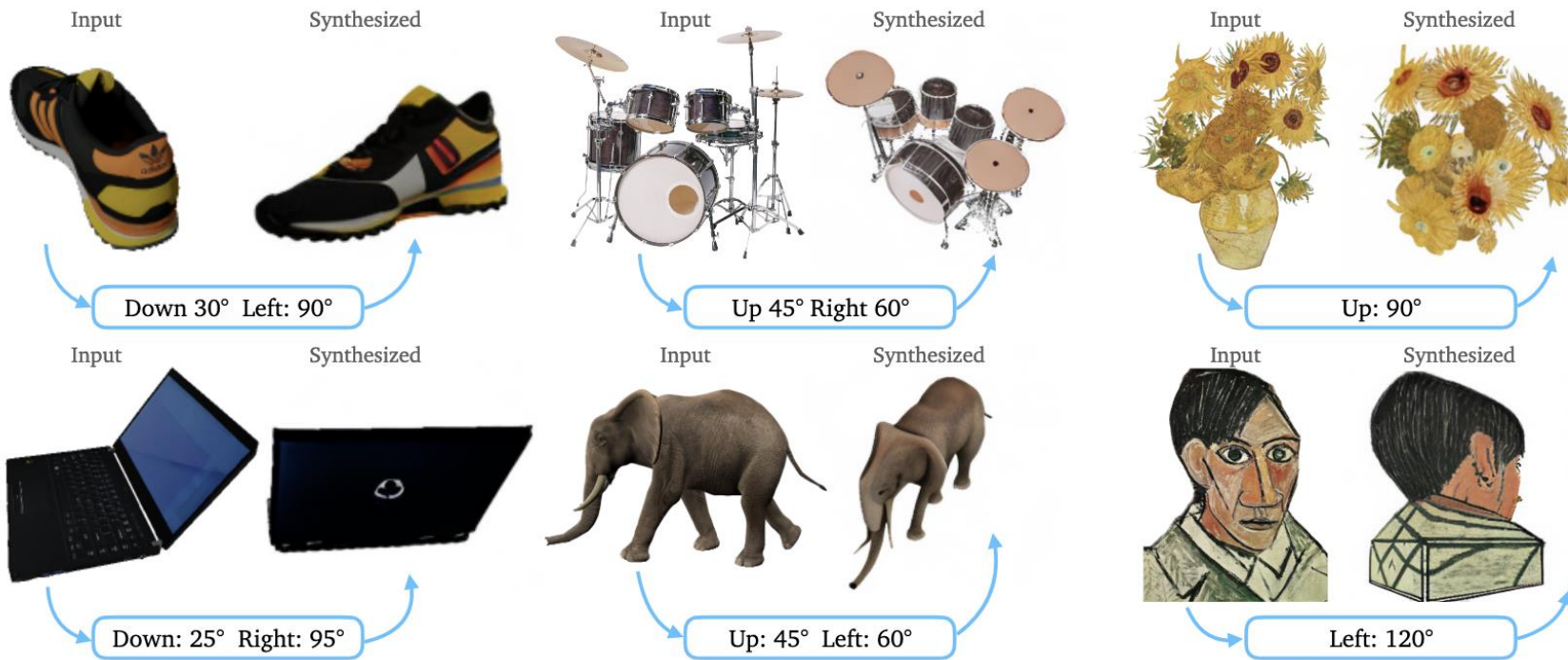


a tiger dressed as a doctor\*

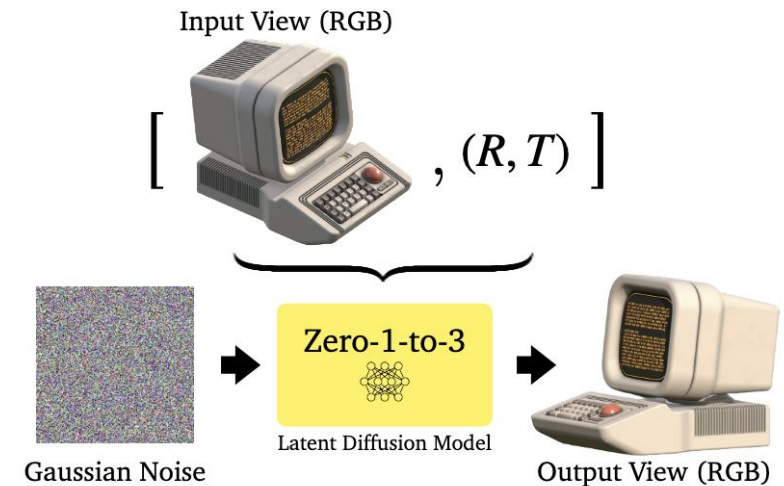


a steam engine train, high resolution\*

# Better Ideas: Fine-tune Image Diffusion Model for Multi-View Generation



Enable conditioning on camera angles



# Generate Training Data of Camera Poses and Rendered Views from 3D Asset Dataset



Source	# Objects
IKEA [32]	219
GSO [17]	1K
EGAD [41]	2K
OmniObject3D [63]	6K
PhotoShape [46]	5K
ABO [13]	8K
Thingi10K [67]	10K
3d-Future [19]	10K
ShapeNet [9]	51K
Objaverse 1.0 [14]	800K
<b>Objaverse-XL</b>	<b>10.2M</b>

# Results: More Consistent View Generation



# Inconsistency is still the Problem



**Multi-face Janus Problem**



**Content Drift Problem**

Figure 1: Typical multi-view consistency problems of 2D-lifting methods for 3D generation. Left: “A photo of a horse walking” where the horse has two faces. Right: “a DSLR photo of a plate of fried chicken and waffles with maple syrup on them” where the chicken gradually becomes a waffle.

# Idea: Concurrently Predict Multiple Views



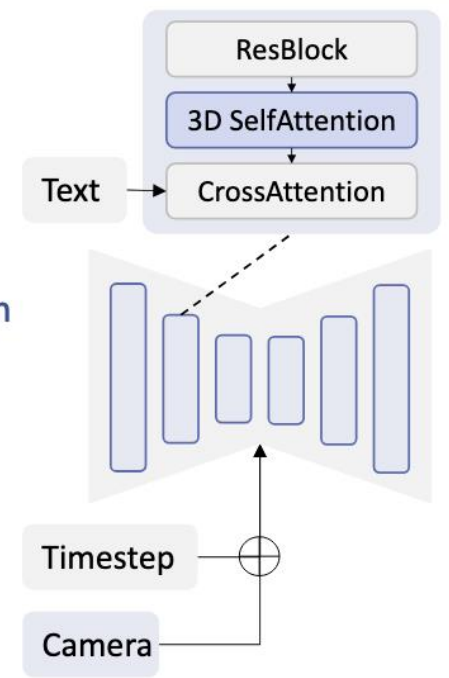
3D model



Rendered images

Training Loss  
Multi-view Generation

3D Generation  
Score Distillation



Multi-view Diffusion UNet

# More Consistent Multi-View Generation

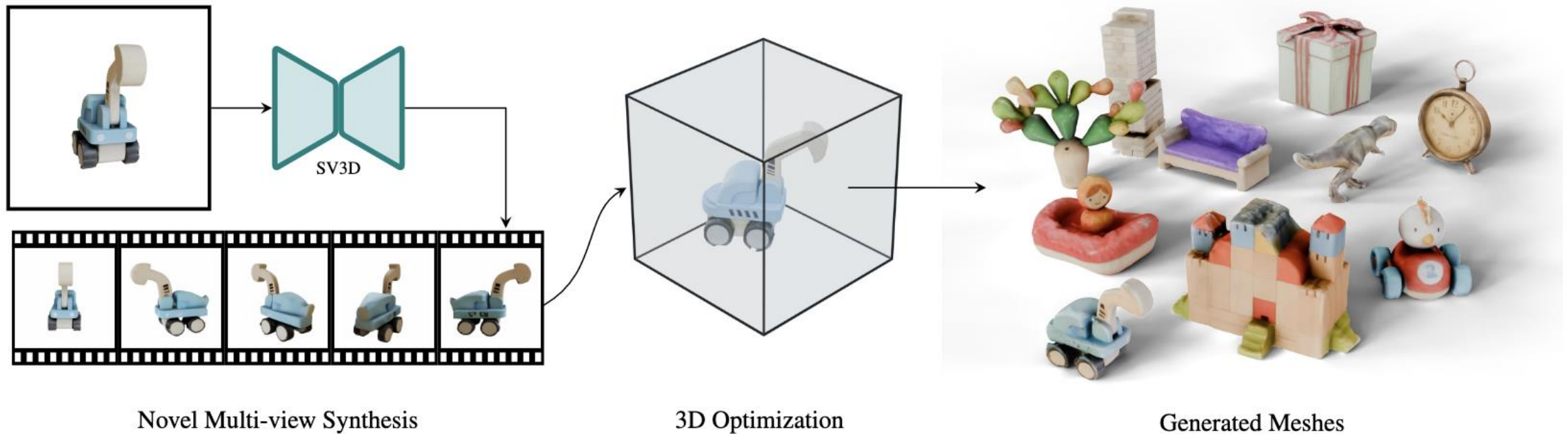


a bald eagle carved out of wood

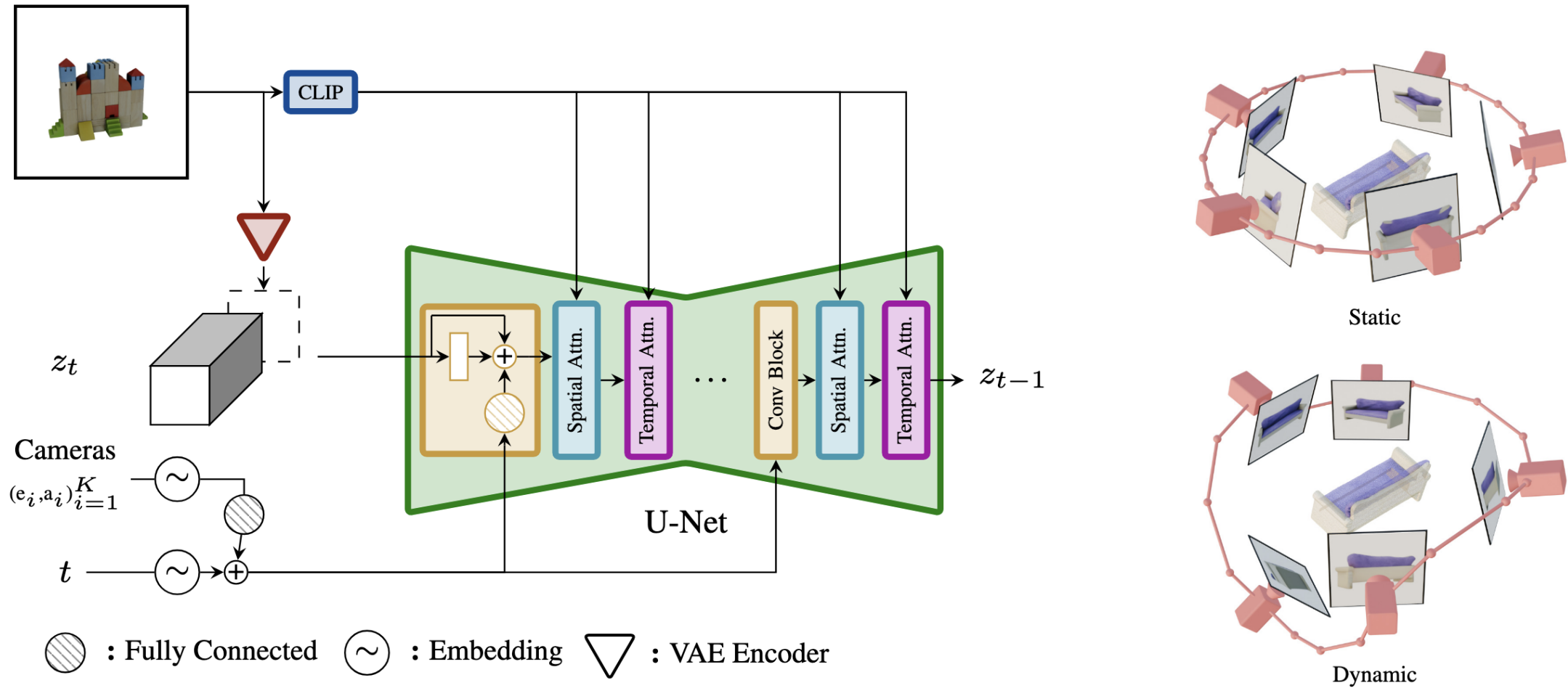


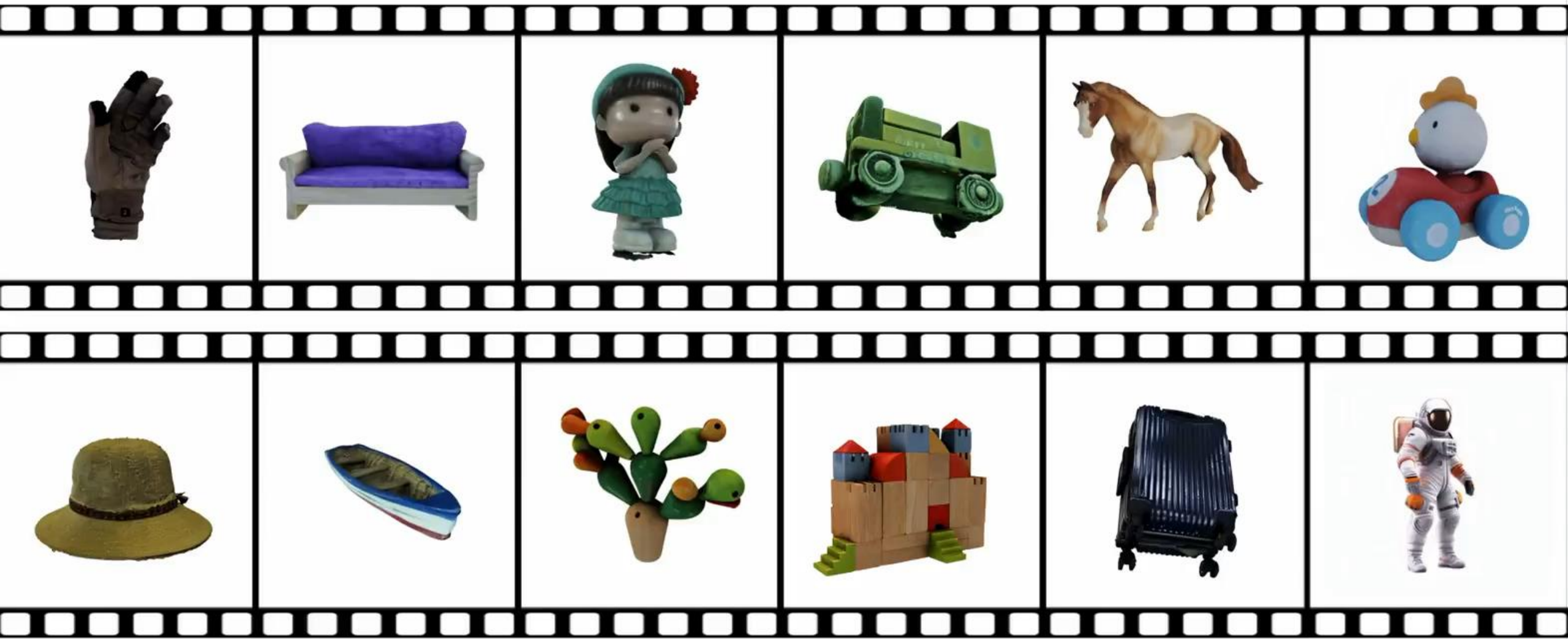
a bulldog wearing a black pirate hat

# A Better Idea: Concurrently Predict Videos of Multiple Views



# Camera-view Trajectory as Video

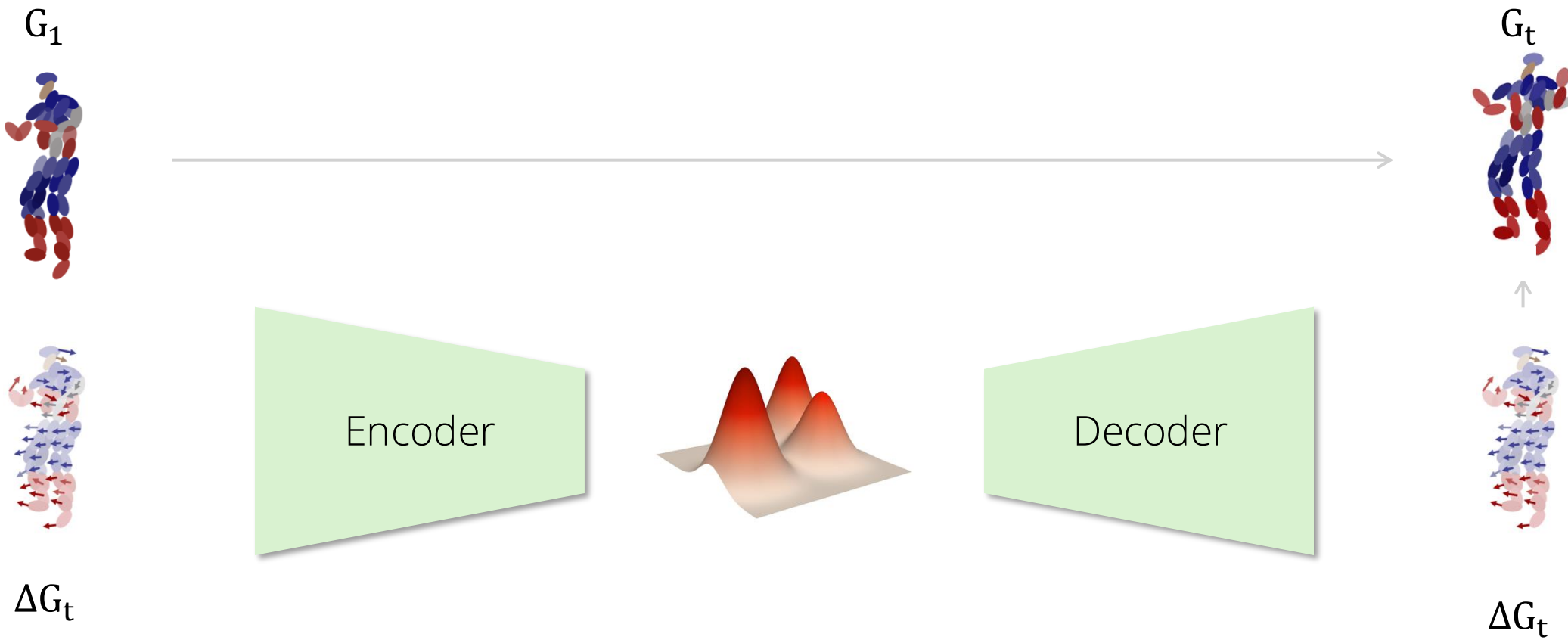




# Content

- Generative Modeling
  - Autoregressive Models
  - Variational AutoEncoders
  - Denoising Diffusion Probabilistic Models
  - Flow Matching Models
- 4D Generative Modeling
  - 3D Geometry Generation
  - 3D Generation with Multi-view Video Generation
  - 4D Geometry Generation
  - 4D Generation with Multi-view Video Generation

# We Have Covered in Previous Lecture. The Idea is to Decouple Shape and Motion Generation



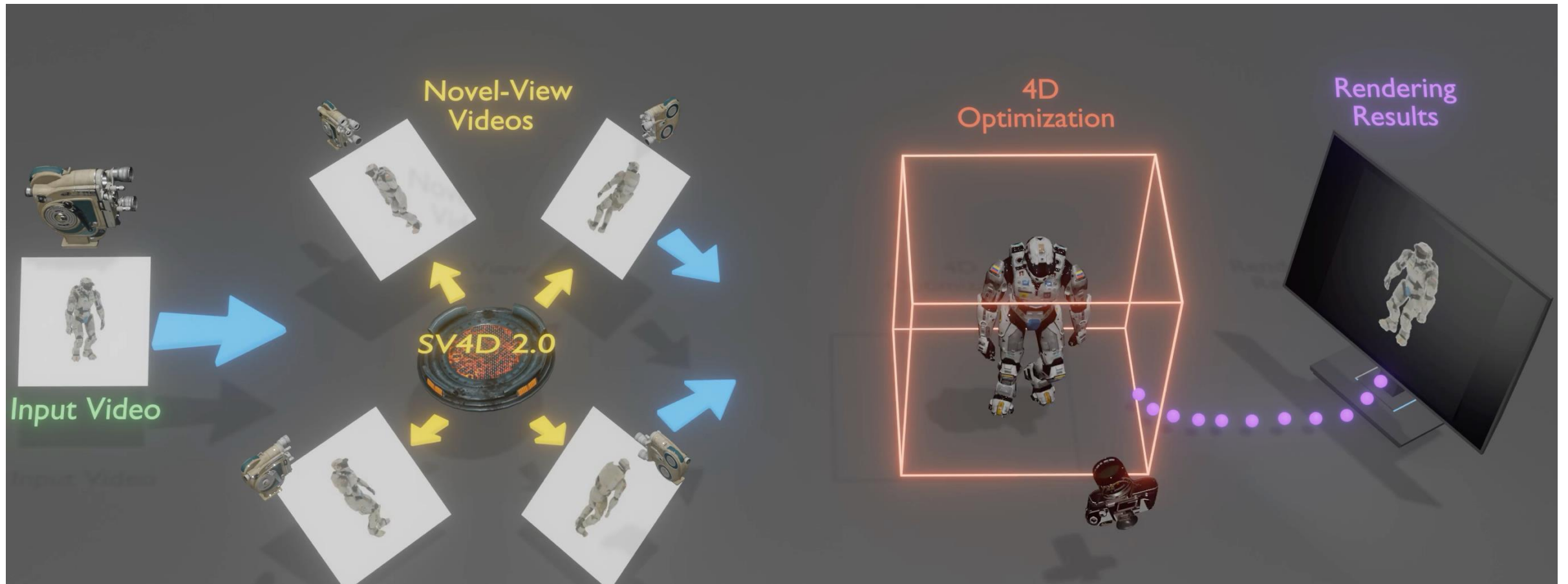
# Content

- Generative Modeling
  - Autoregressive Models
  - Variational AutoEncoders
  - Denoising Diffusion Probabilistic Models
  - Flow Matching Models
- 4D Generative Modeling
  - 3D Geometry Generation
  - 3D Generation with Multi-view Video Generation
  - 4D Geometry Generation
  - 4D Generation with Multi-view Video Generation

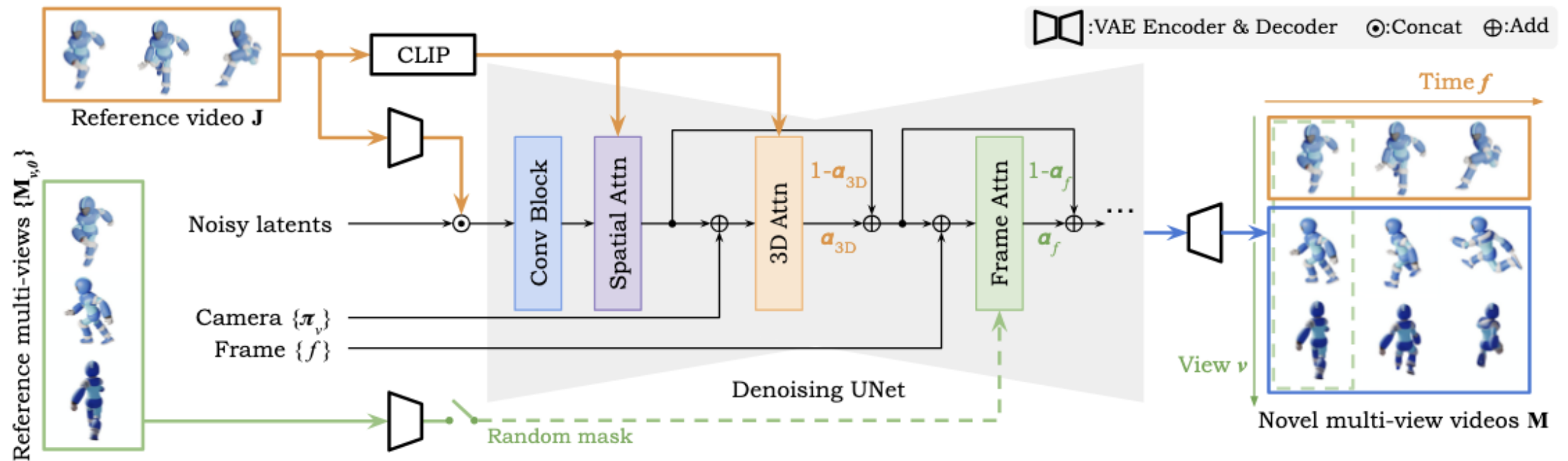
# How to Generate a 4D World with Video Diffusion?

- How to generate a 4D world?
  - Generate multi-view 2D videos and lift to 4D
  - Generate multi-view 2D images and lift to 3D, and generate motion of these 3D assets
  - Generate 3D videos

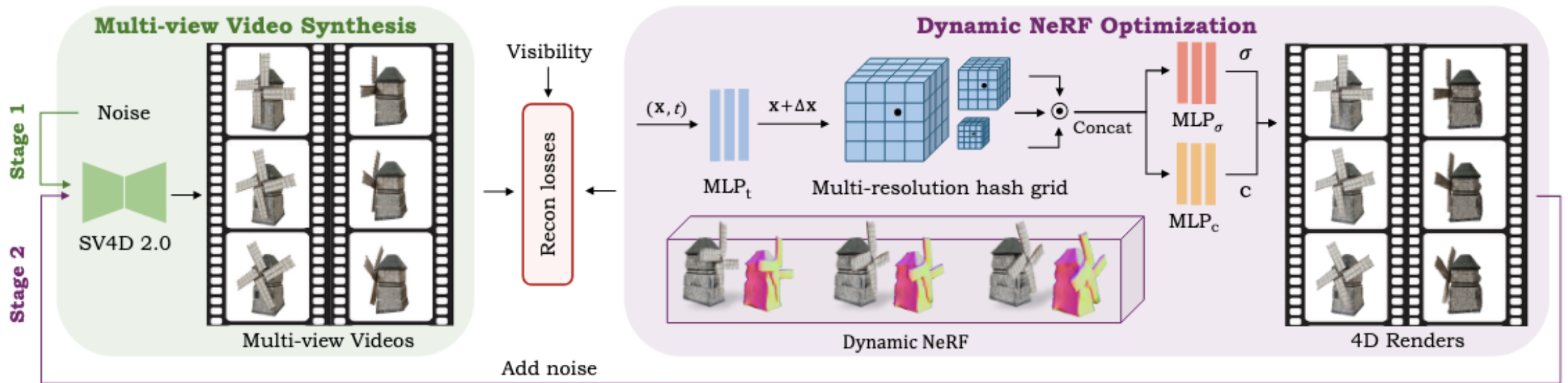
# Idea: Generate Multi-view Videos and Lift to 4D



# Generate Multi-view Videos and Lift to 4D



# Generate Multi-view Videos and Lift to 4D



# Generate Multi-view Videos and Lift to 4D



*Input Video*



*4D Optimization*



*Input Video*



*4D Optimization*



*Input Video*



*4D Optimization*



*Input Video*



*4D Optimization*

## CAT4D: Create Anything in 4D with Multi-View Video Diffusion Models

Rundi Wu<sup>1,2</sup> Ruiqi Gao<sup>1</sup> Ben Poole<sup>1</sup> Alex Trevithick<sup>1,3</sup>  
Changxi Zheng<sup>2</sup> Jonathan T. Barron<sup>1</sup> Aleksander Hołyński<sup>1</sup>

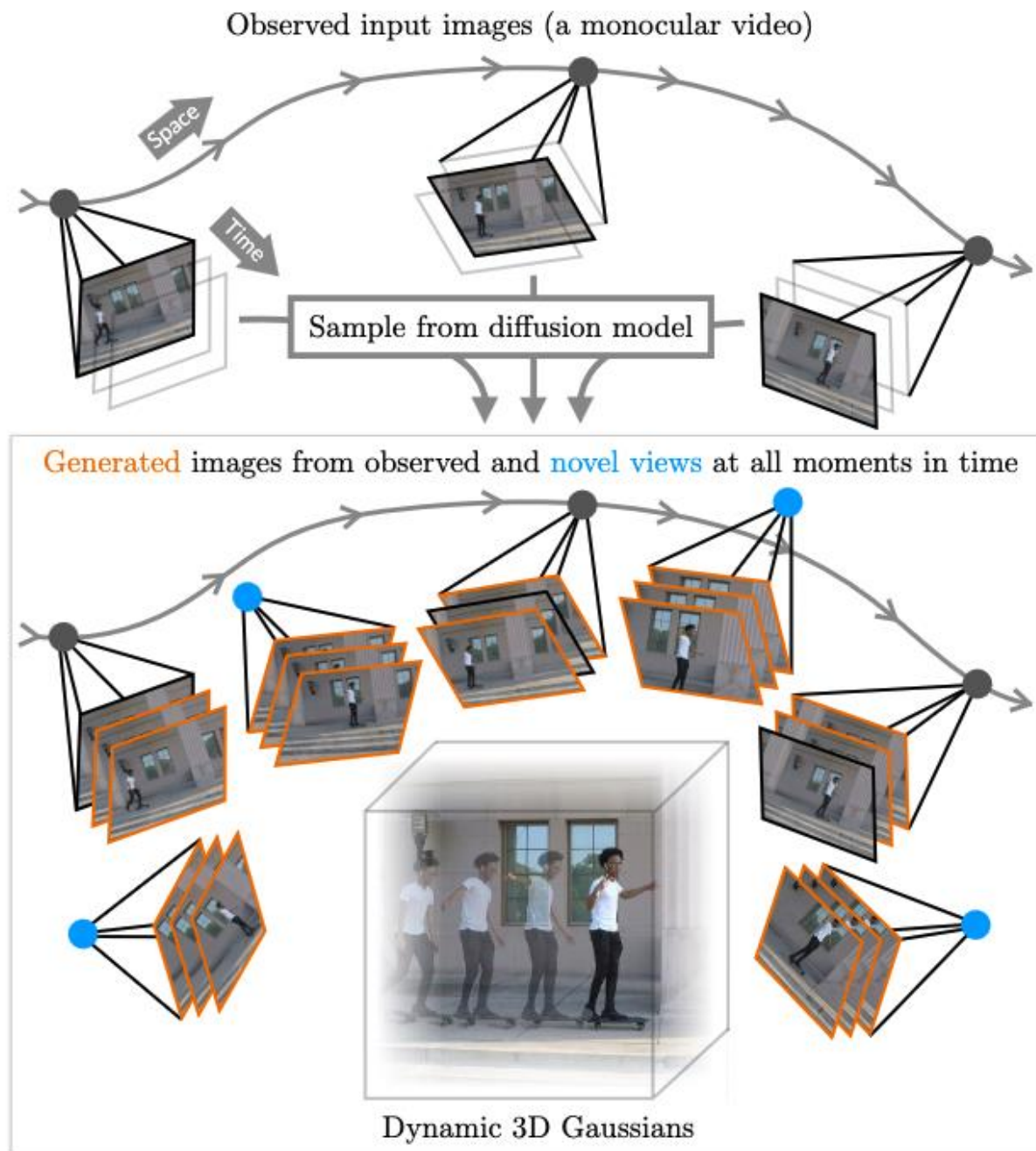
<sup>1</sup>Google DeepMind <sup>2</sup>Columbia University <sup>3</sup>UC San Diego



Sparse dynamic capture to 4D

Real video to 4D

Text-to-video-to-4D





Input video

Sample from multi-view  
video diffusion model



Generated multi-view videos

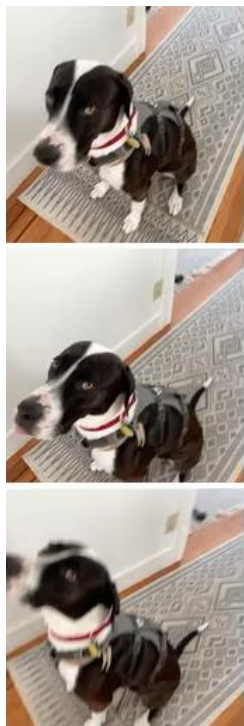
Optimize dynamic  
3D Gaussians



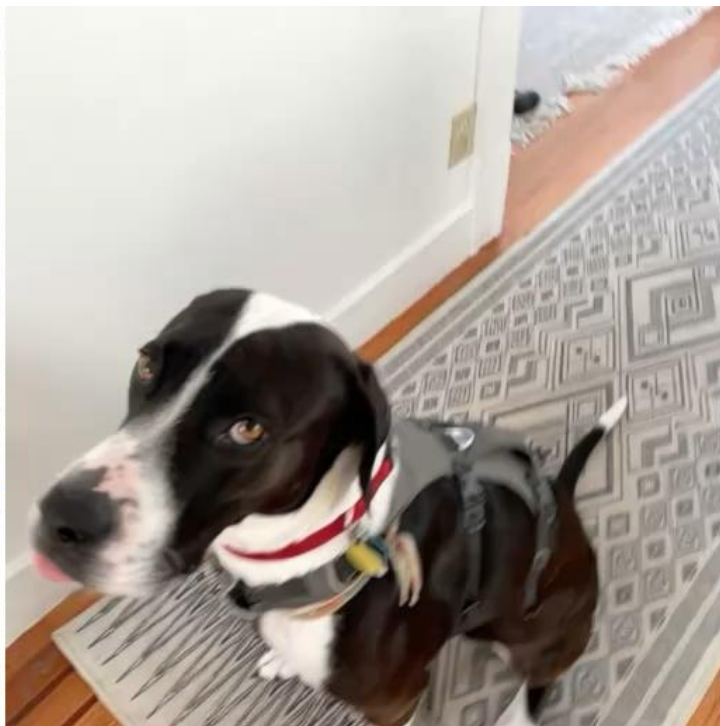
Dynamic 3D scene

# Results: 4D Reconstruction

input



Fixed View Varying Time



Varying View Fixed Time



Varying View Varying Time



# Results: 4D Reconstruction

input



Fixed View Varying Time



Varying View Fixed Time



Varying View Varying Time

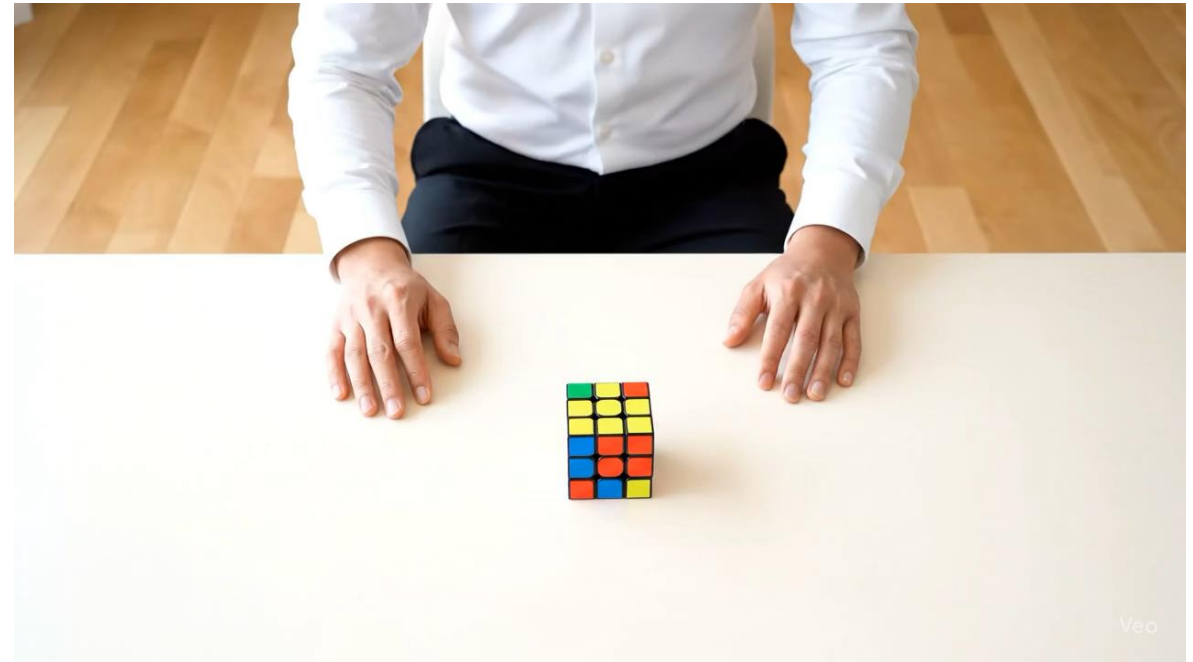


Have you seen the problem?

# Generated Videos are Often Physically Incorrect. We Need to Model the Motion Explicitly.

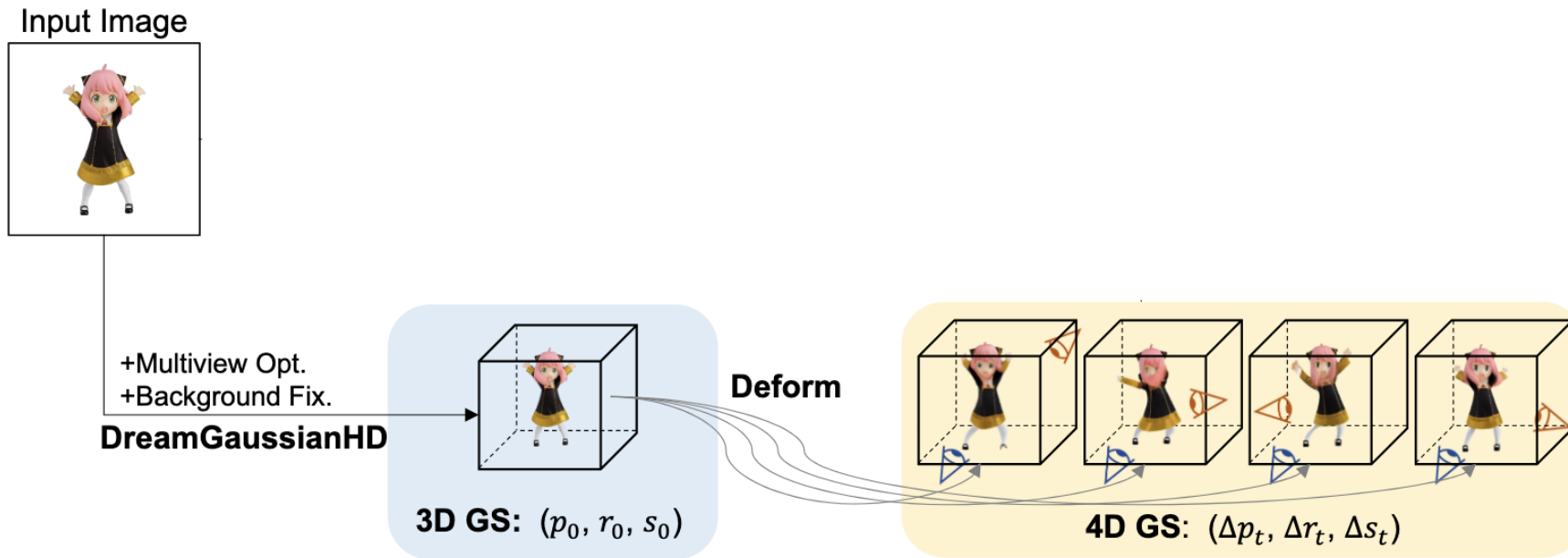


Video generated by Veo3

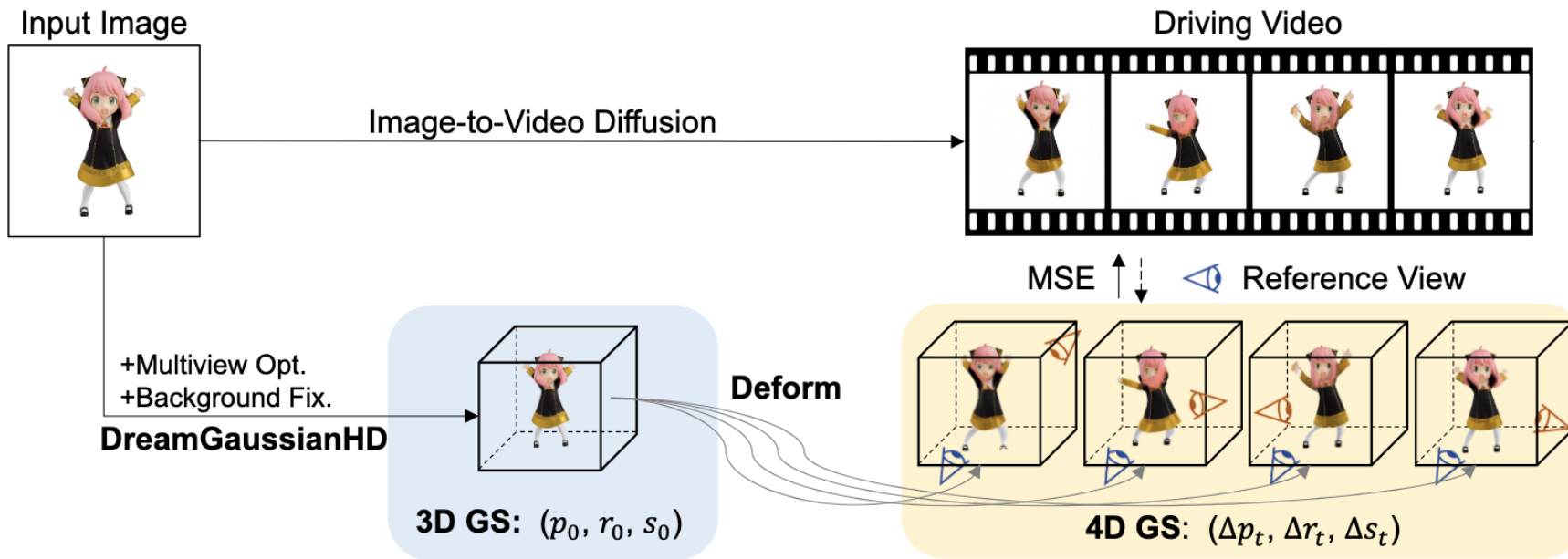


Video generated by Veo3

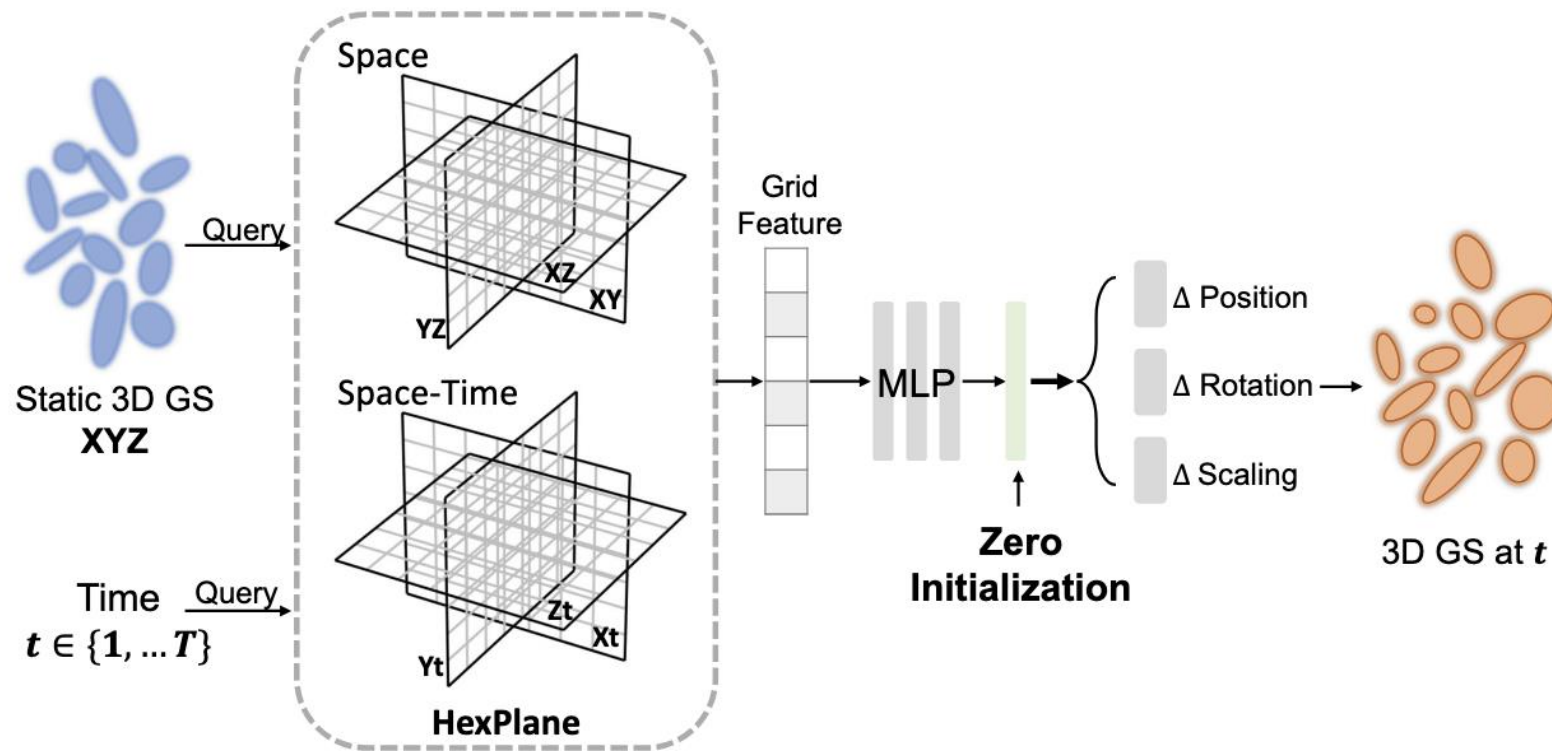
# Idea: Create 3D Assets by Generating Multi-view Images and Generate the 3D Motion



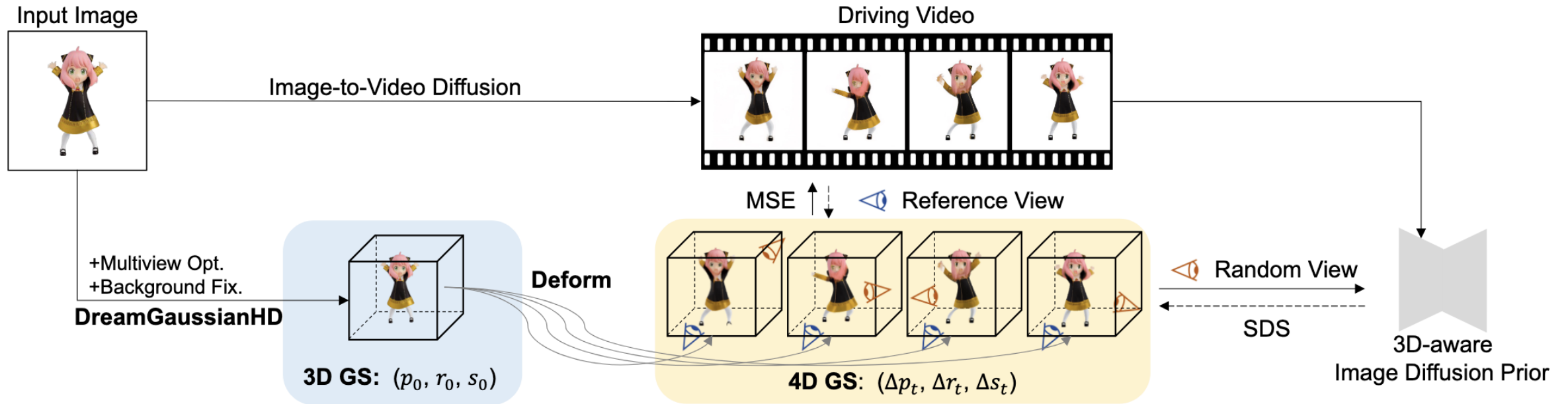
# Generate the Motion Field from a Reference Video



# Predict the Motion of 3D Gaussians with a Neural Network



# Create 3D Assets by Generating Multi-view Images and Generate the 3D Motion



Given a single-view image:

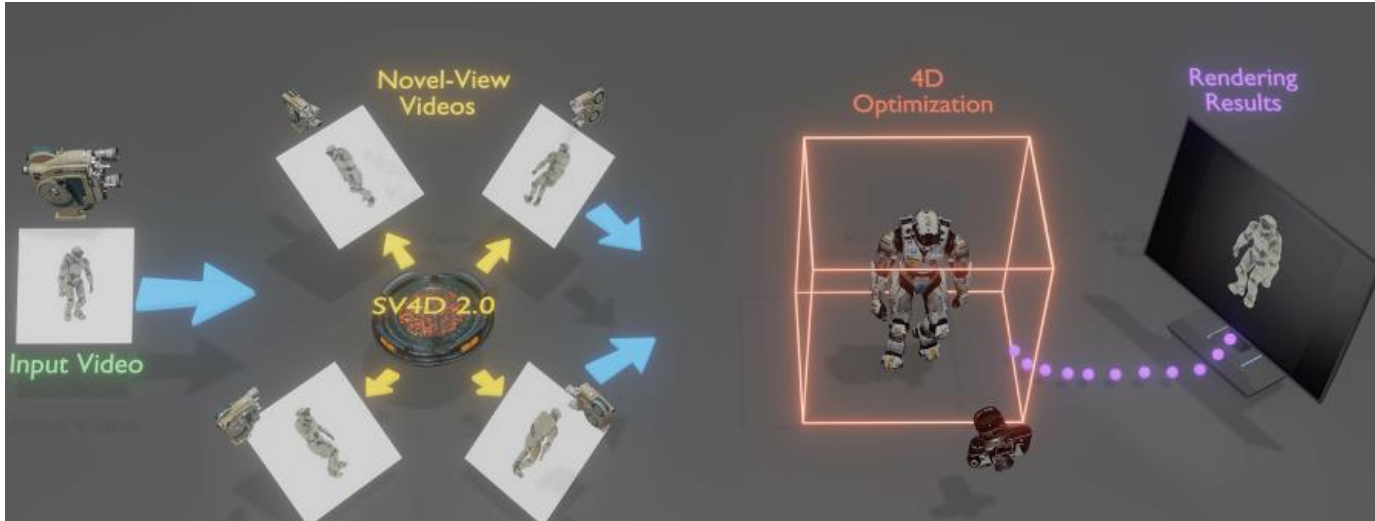
1. 3D-from-2D reconstruction
2. Flow prediction which deforms the 3D model at each frame
3. Image-to-Video prediction to guide the optimization of flow prediction

Once we have deformed 3D model at each frame:

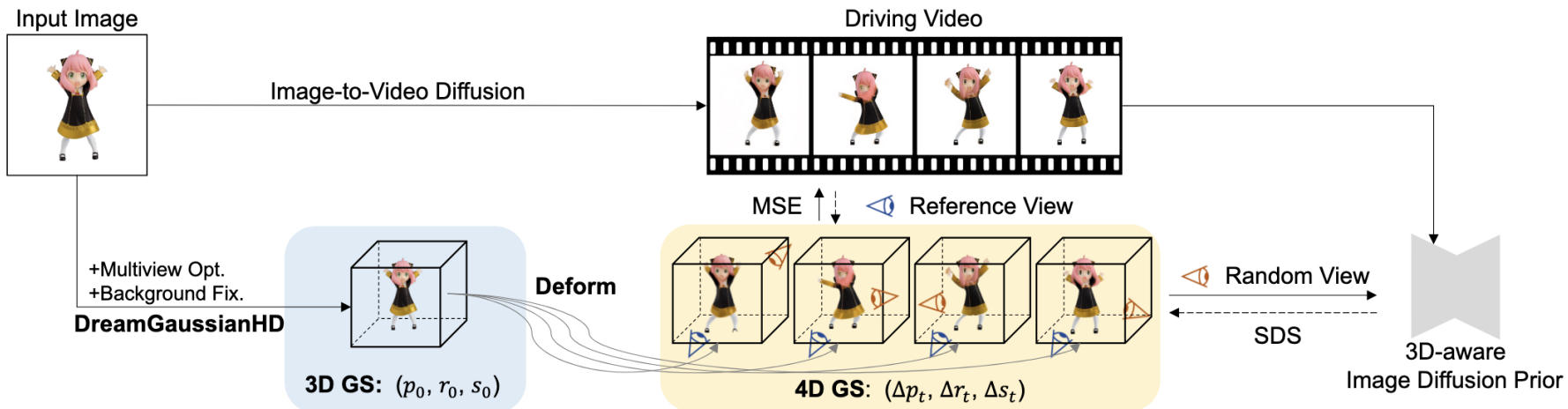
$$\mathcal{L}_{\text{Ref}} = \frac{1}{\mathcal{T}} \sum_{\tau=1}^{\mathcal{T}} \|f(\phi(S, \tau), o_{\text{Ref}}) - I_{\text{Ref}}^{\tau}\|_2^2,$$

$$\nabla_{\phi} \mathcal{L}_{\text{SDS}} = \mathbb{E}_{t, \tau, \epsilon, o} [(\epsilon_{\theta}(\hat{I}; t, I_{\text{Ref}}^{\tau}, o) - \epsilon) \frac{\partial I}{\partial \phi}], \text{ and } \hat{I} = f(\phi(S, \tau), o),$$

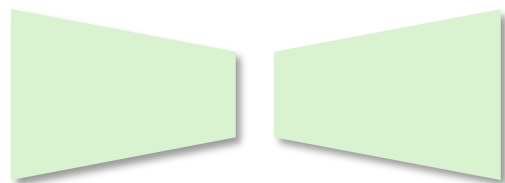
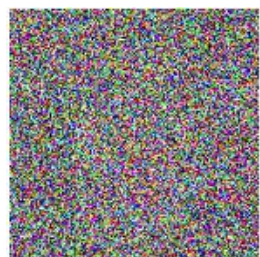
# What are the Problems?



- The distribution of motion is implicitly captured by the video generator
- The distribution of motion is coupled with the distribution of appearance
- Can we generate generate motion explicitly?

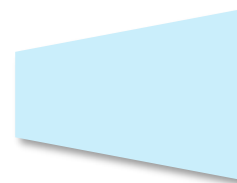
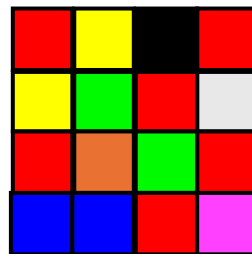


# Idea: Generate 3D Videos



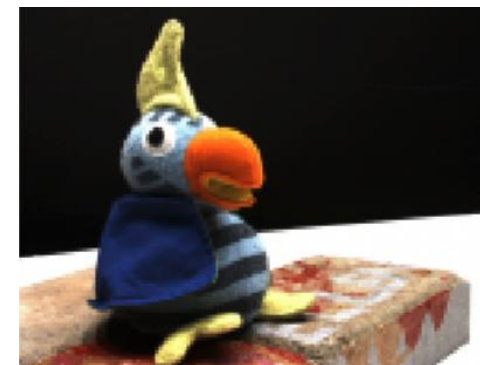
Generative model

Low-resolution features



Decoder

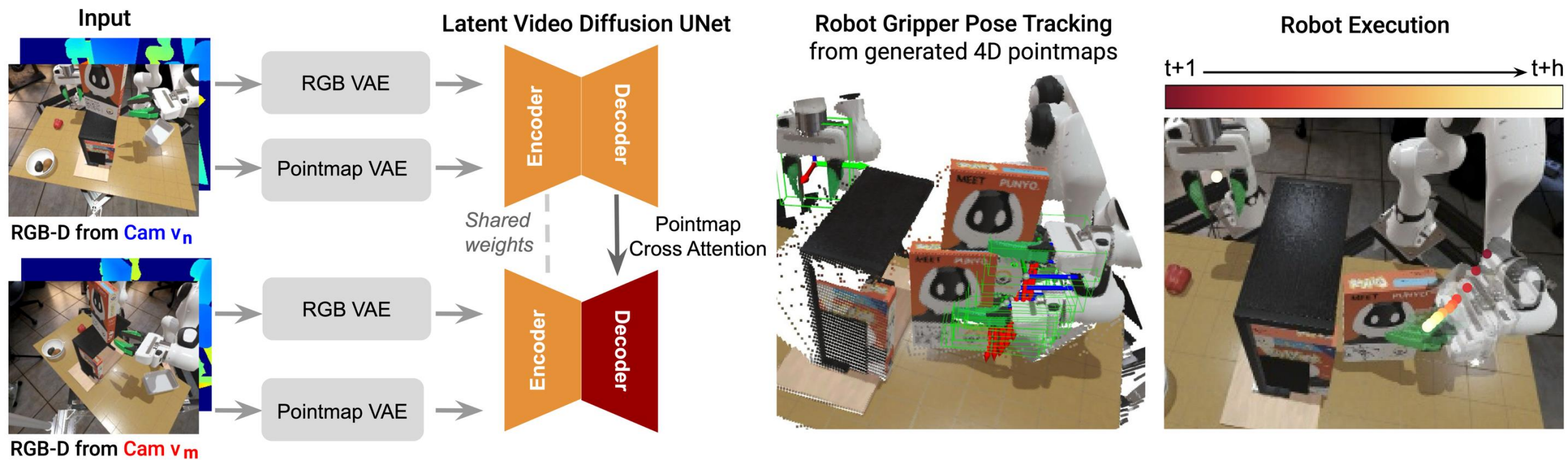
RGB image



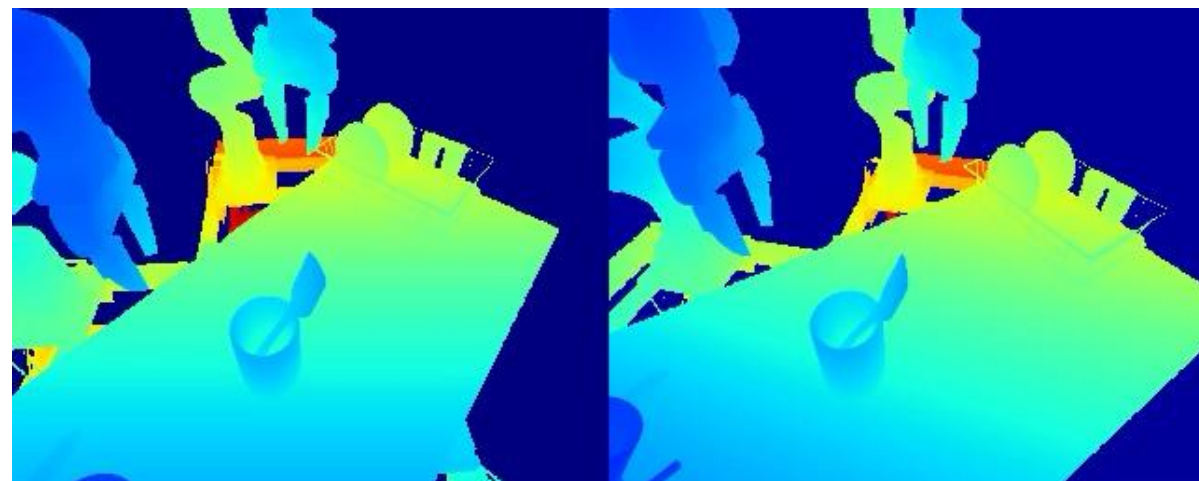
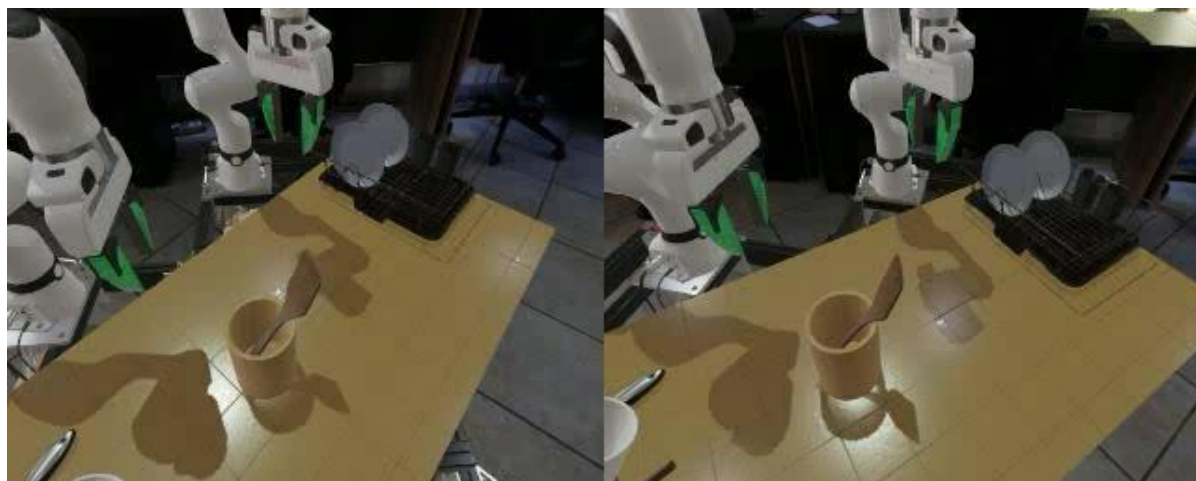
XYZ Point map



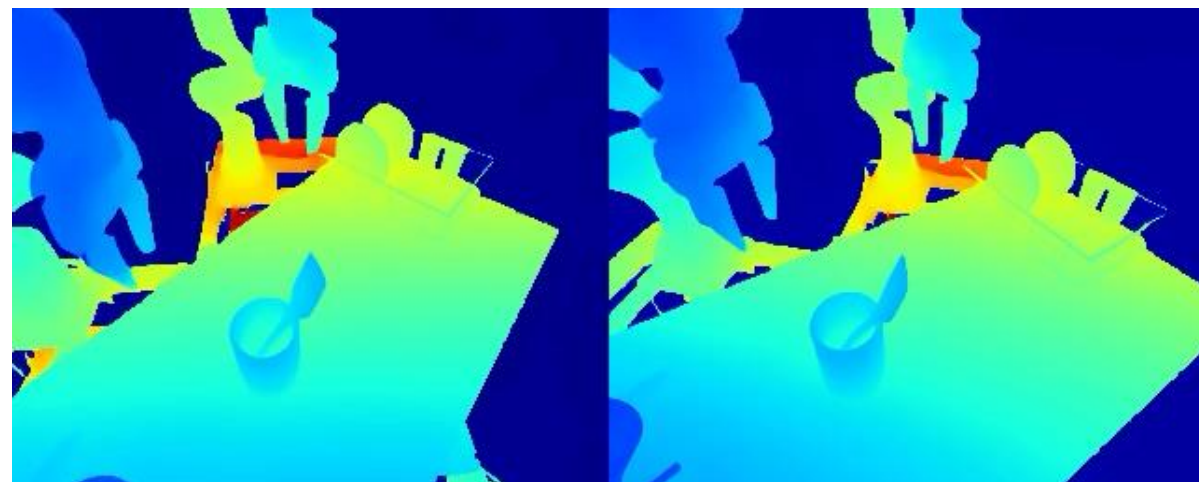
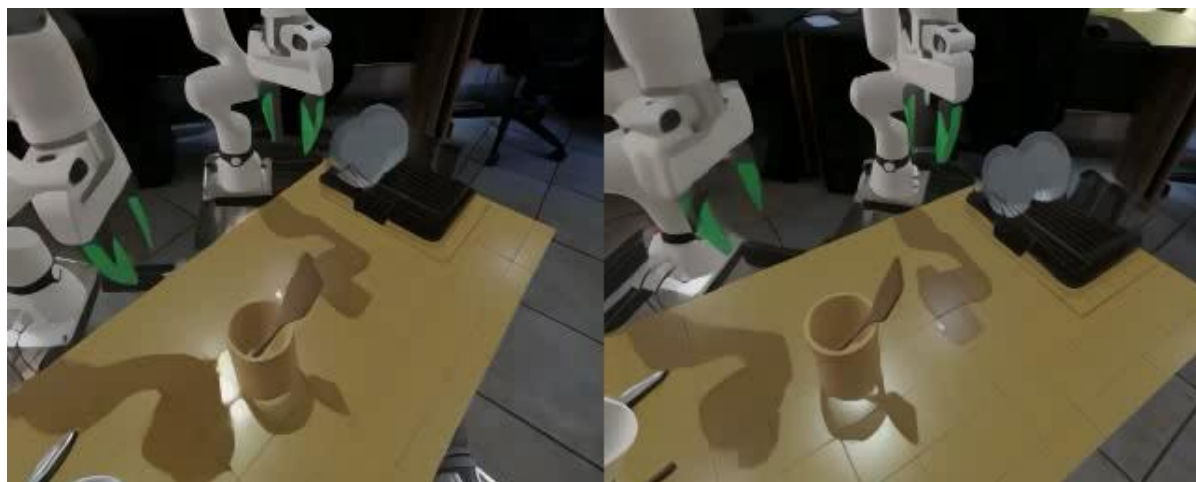
# 3D Video Generator



## Ground-truth



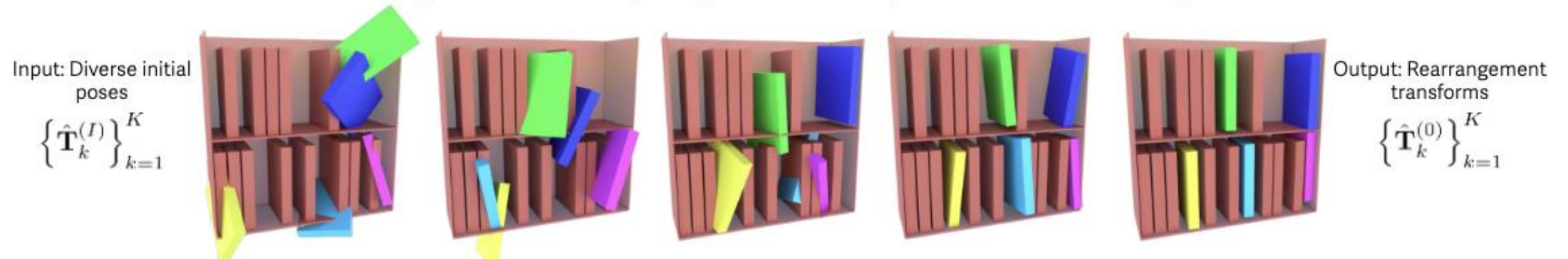
## Generation



What Other Motion Can We Generate  
with Diffusion Models?

# Diffusion Models that Generate Rigid-body Motion

(A) Evaluation: Starting from diverse initial poses, pose diffusion outputs a diverse set of rearrangement solutions



- Idea: Use a diffusion model to generate 6 DoF object poses for arranging objects in a scene
- Problem: 6 DoF object poses belongs to  $SE(3)$  group (a special set of  $3 \times 4$  matrices). However, the diffusion model we discussed is based on the Gaussian formulation. How to generalize the diffusion model to  $SE(3)$ ?

# Generalize Diffusion Models to $SE(3)$

- Idea 1: Formulate the noising and denoising process with  $SE(3)$

## **$SE(3)$ -DiffusionFields: Learning smooth cost functions for joint grasp and motion optimization through diffusion**

Julen Urain<sup>\*1</sup>, Niklas Funk<sup>\*1</sup>, Jan Peters<sup>1,2,3,4</sup>, Georgia Chalvatzaki<sup>1</sup>



**Fig. 1:** Pick and place task in which the robot has to pick a mug and move it to the target pose (in the shelves) without colliding. We exploit diffusion models for jointly optimizing both grasp and motion and show the successful trajectory from left to right.

# Generalize Diffusion Models to $SE(3)$

- Idea 2: Extend the definition of noise from a Gaussian noise to a random transform  $T_{\Delta}^{rand} \in SE(3)$ , and define the sample at diffusion step  $i$  as a transform  $T^{(i)} \in SE(3)$

$$x_{i-1} = \frac{1}{\sqrt{\alpha_i}} \left( x_i - \frac{1-\alpha_i}{\sqrt{1-\alpha_i}} \epsilon_{\theta}(x_i, i) \right) + \delta_i z$$

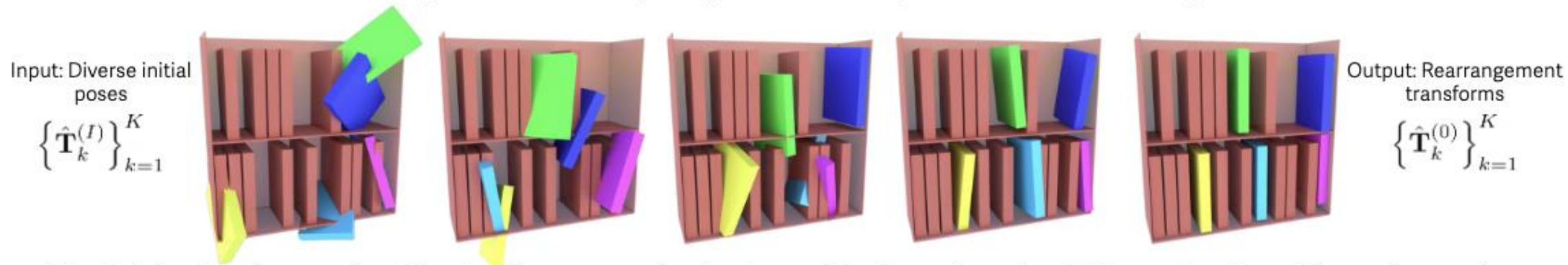


$$T^{(i-1)} = T_{\Delta}^{rand} T_{\Delta}^i T^{(i)}$$

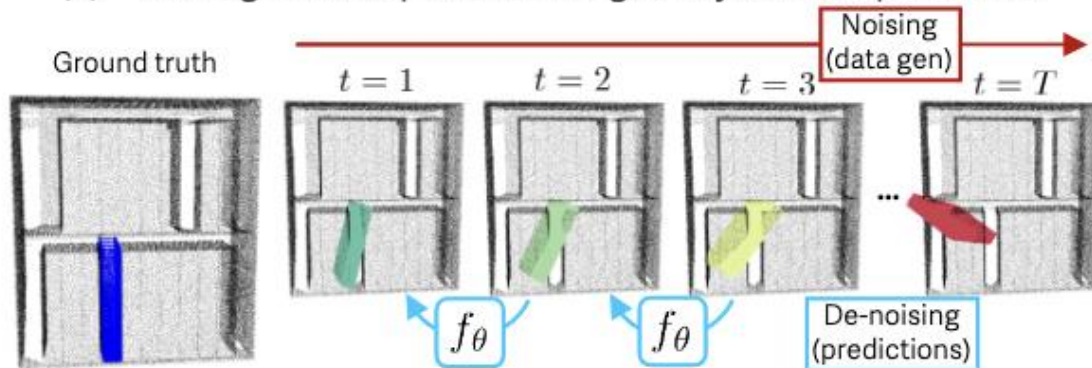
- In DDPM, we sample from a Gaussian noise and denoising it into a sample. Here, initial rotations are drawn from a uniform grid over  $SO(3)$  and translations are uniformly drawn from the bounding box of the scene.

# Diffusion Models that Generate Rigid-body Motion

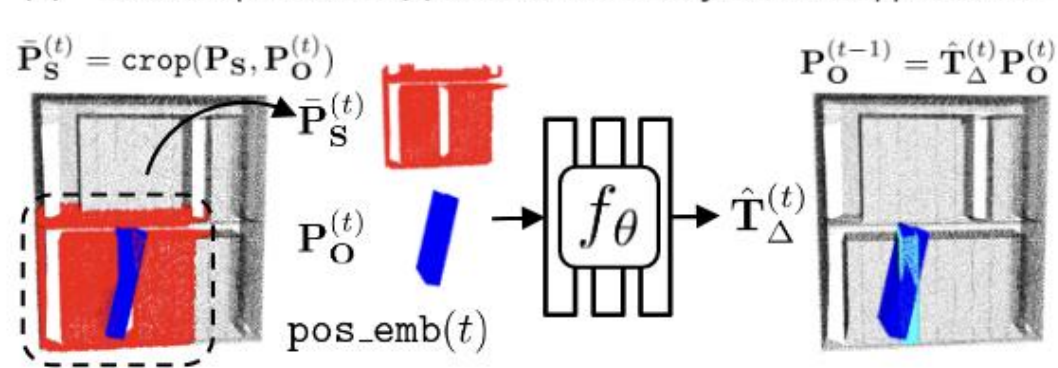
(A) Evaluation: Starting from diverse initial poses, pose diffusion outputs a diverse set of rearrangement solutions



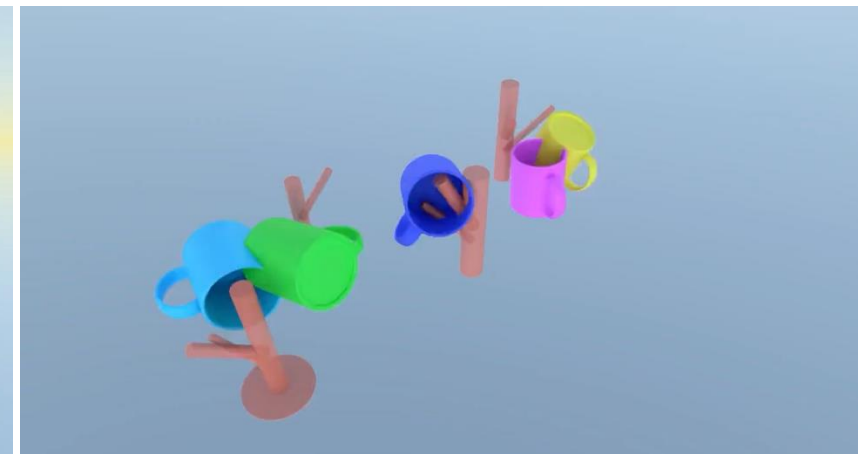
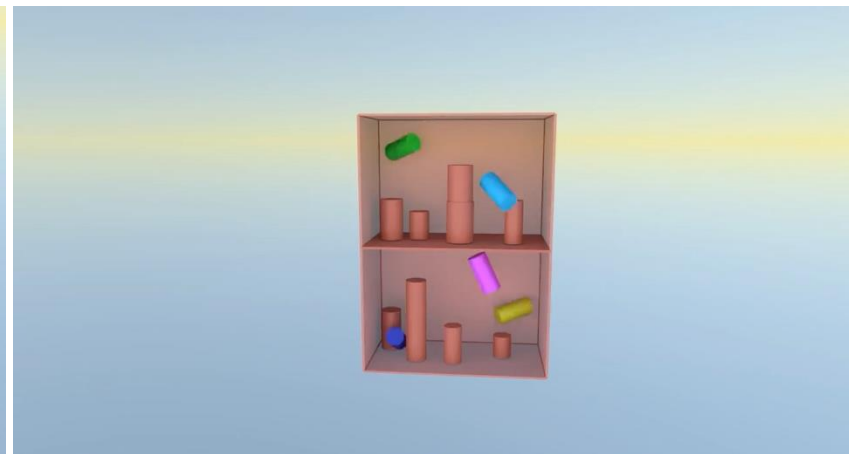
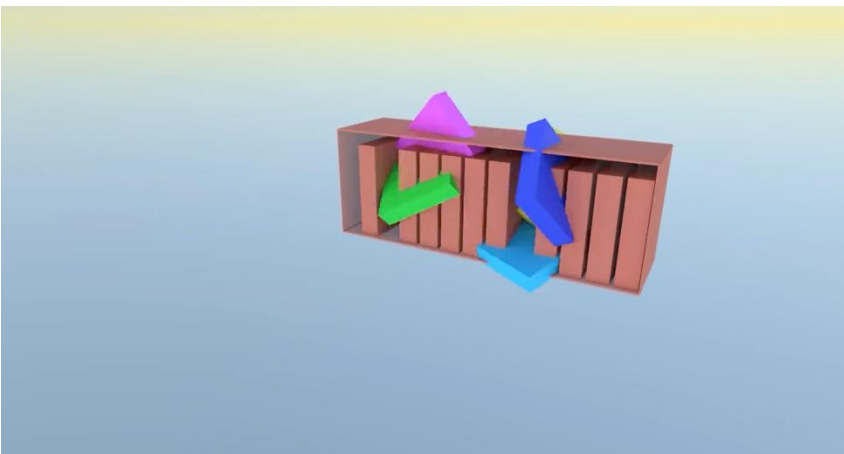
(B) Training: Iterative pose de-noising for object-scene point cloud



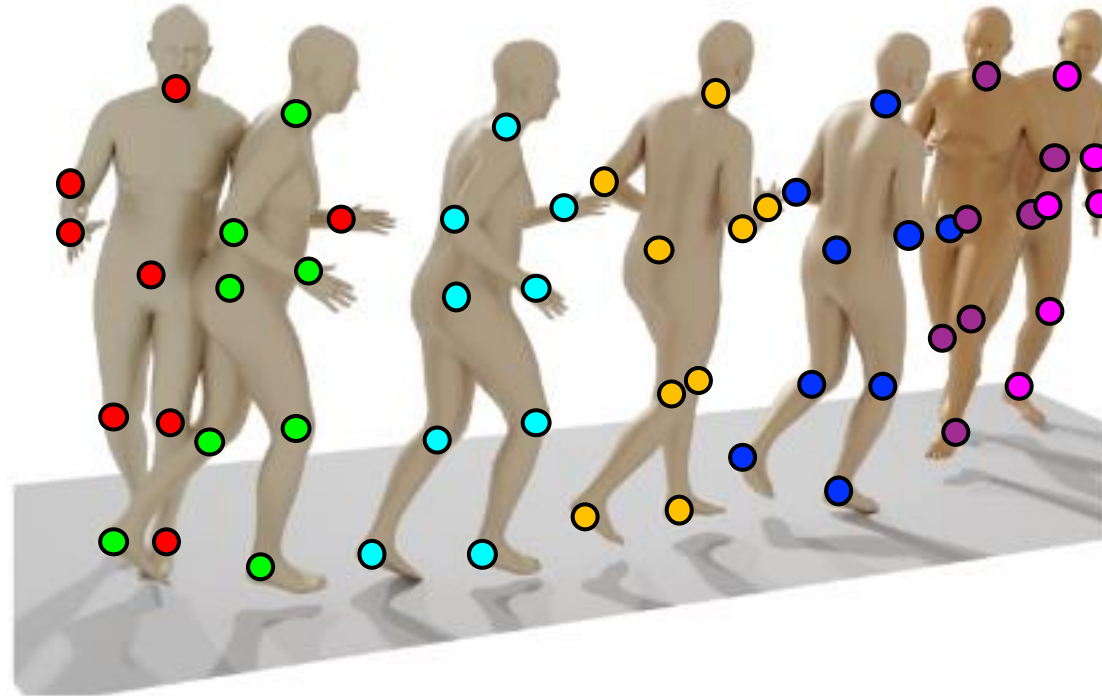
(C) Network predicts SE(3) transform from object and cropped scene



# Results

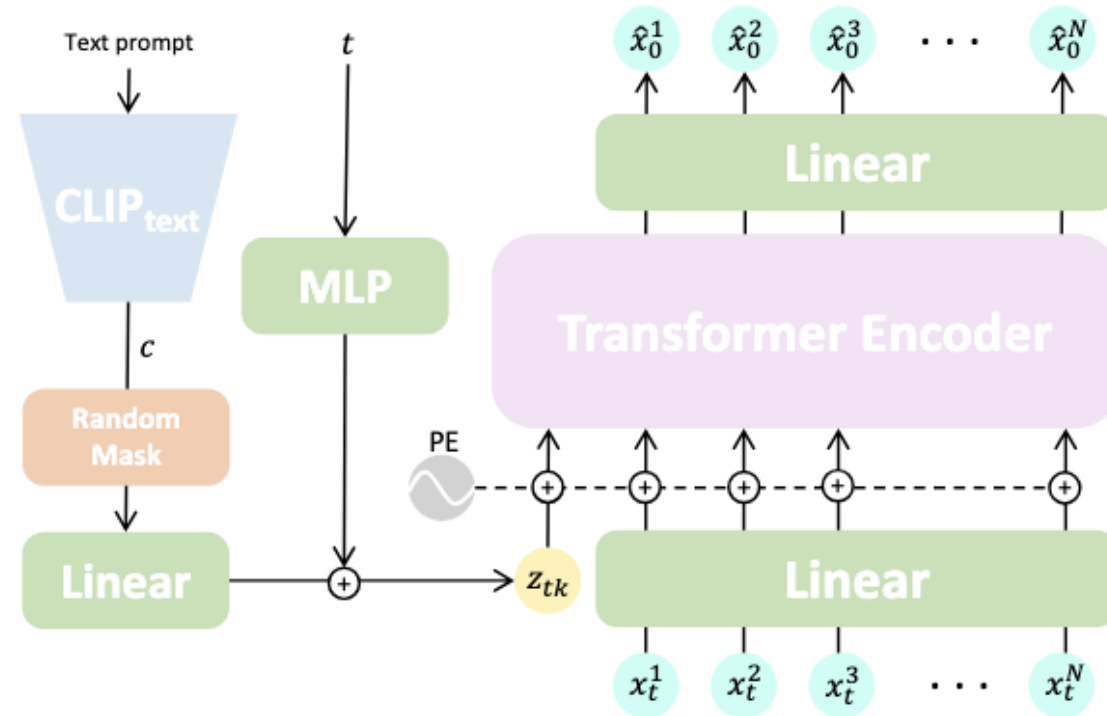


# Diffusion Models that Generate Particle Motion



- Idea: Generate human poses  $x^{1:N} = (x^1, x^2, x^3, x^4, \dots, x^N)$  represented by either joint rotations or positions

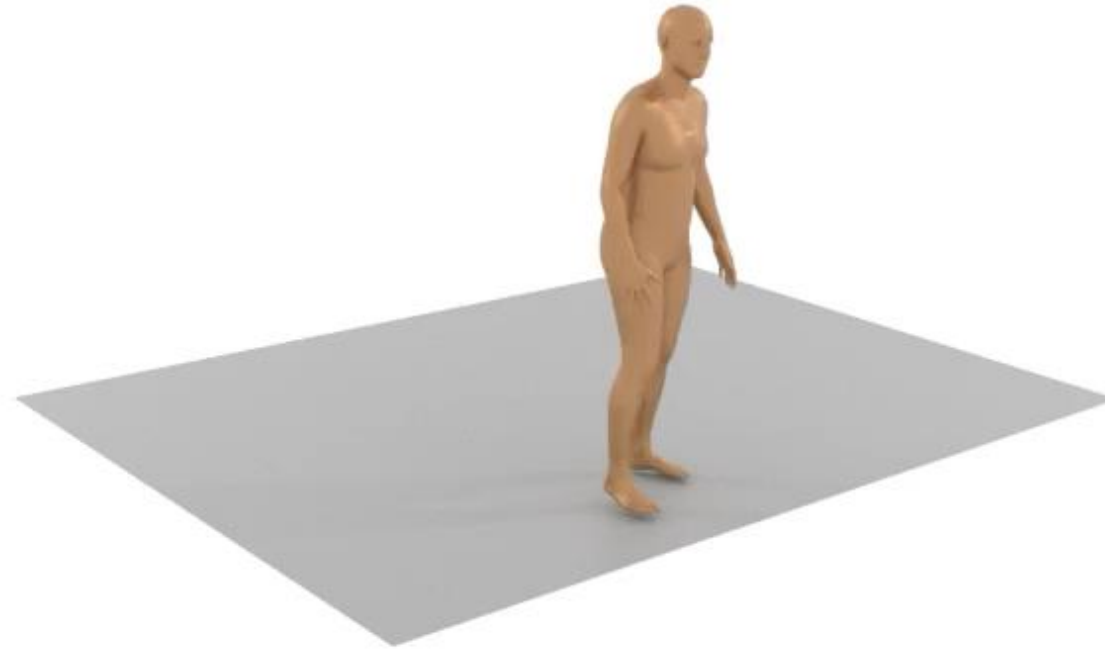
# Diffusion Models that Generate Particle Motion



- Idea: Generate human poses  $x^{1:N} = (x^1, x^2, x^3, x^4, \dots, x^N)$  represented by either joint rotations or positions

# Results

“a person turns to his right and paces back and forth.”



# Are Our Generated 4D Worlds or Motions Physically Correct?



# What We Will Cover the Next Week

- Physical modeling
  - Rigid-body motion