

Digital Humans – Winter 24/25

Lecture 13_3 – Diffusion Models in 3D Reconstruction

Prof. Dr. Gerard Pons-Moll

University of Tübingen / MPI-Informatics

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

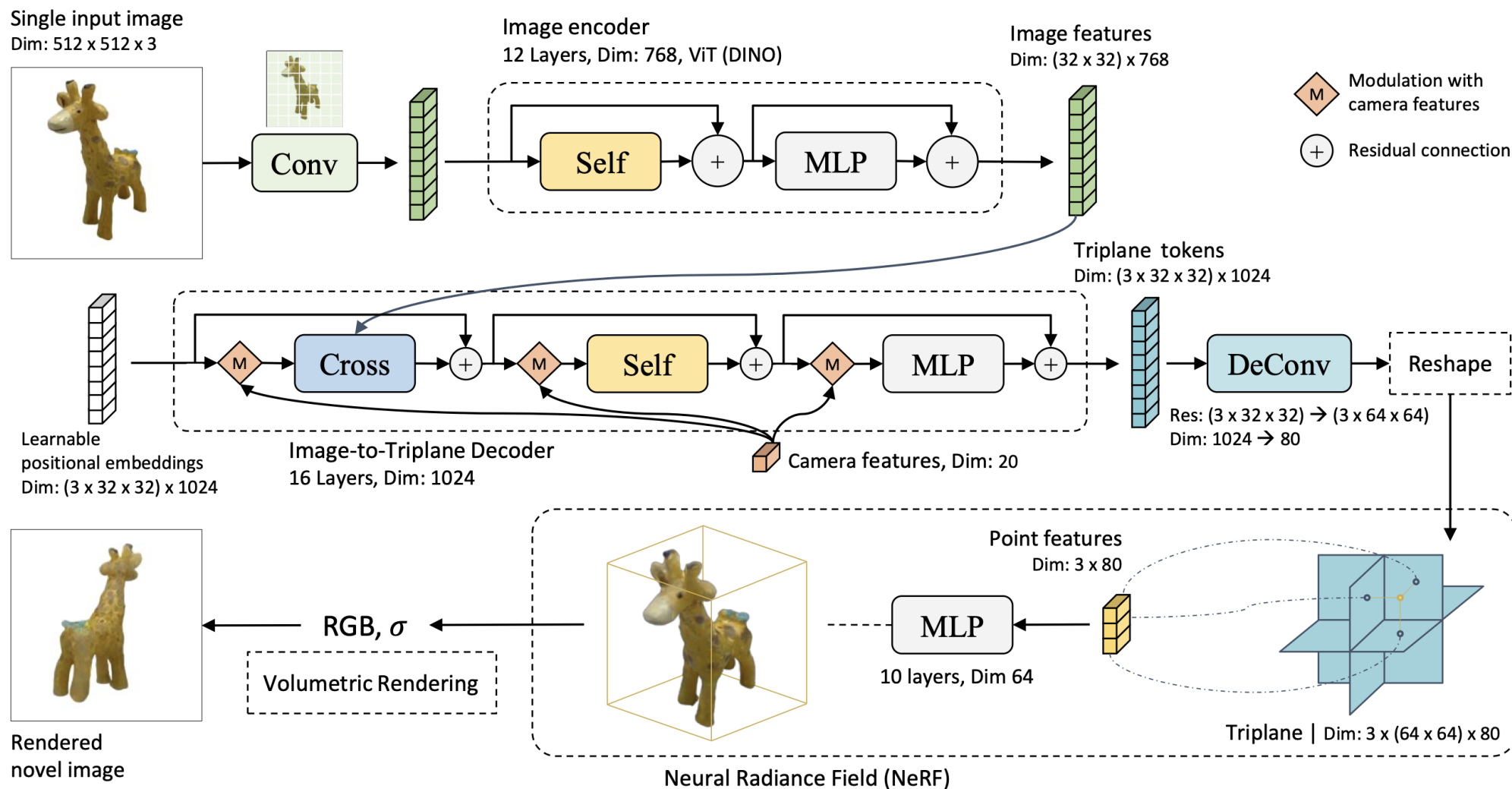


Main contents

- **3D Reconstruction from single Image**
- 2D Diffusion Model for Novel-view Synthesis
 - Novel-view Diffusion Models
 - Multi-view Image Diffusion Models
- Sync 2D Diffusion & 3D Reconstruction

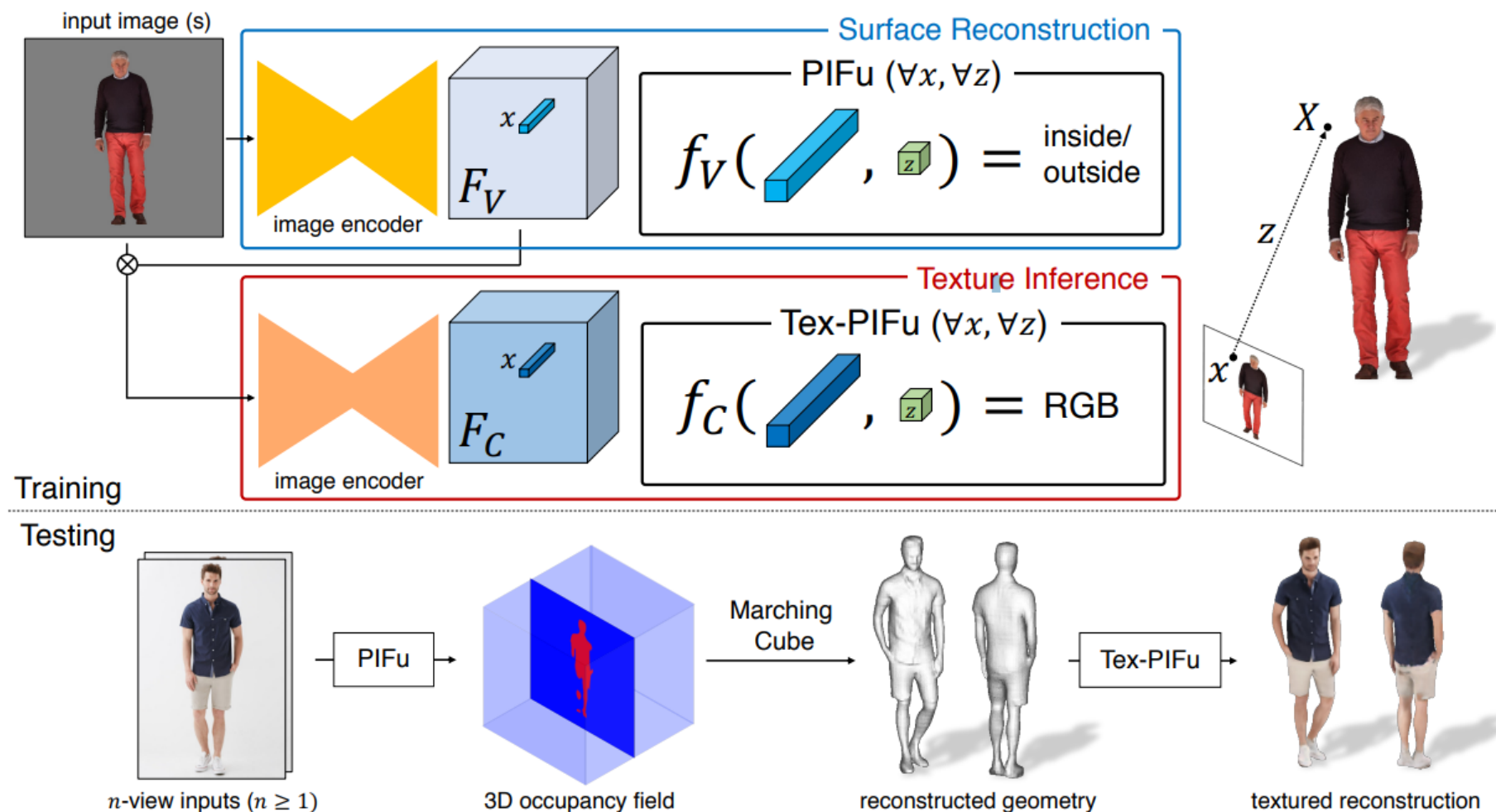
Reconstruct 3D from Single Image

- LRM: Regress NeRF tri-plane features from a RGB image



Reconstruct 3D from Single Image

- PiFu: regress occlusion field from a RGB image.



Reconstruct 3D from Single Image

- Limitation: no generative power, blurry occluded region.



Input Image



OpenLRM, He et al.

Reconstruct 3D from Single Image

- Limitation: no generative power, blurry occluded region.



Input Image



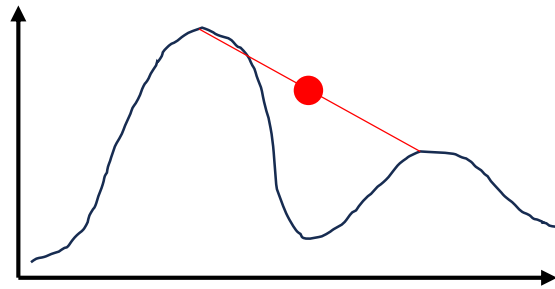
SiTH (CVPR 2024)



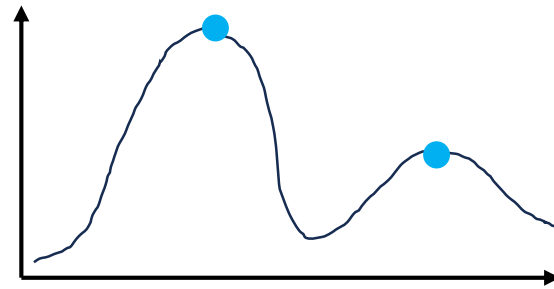
SiFU (CVPR 2024)

Motivation for generative model

- Goal: single view reconstruction is ill-posed.
- Deterministic model might collapse to average value.



Deterministic: learn an average



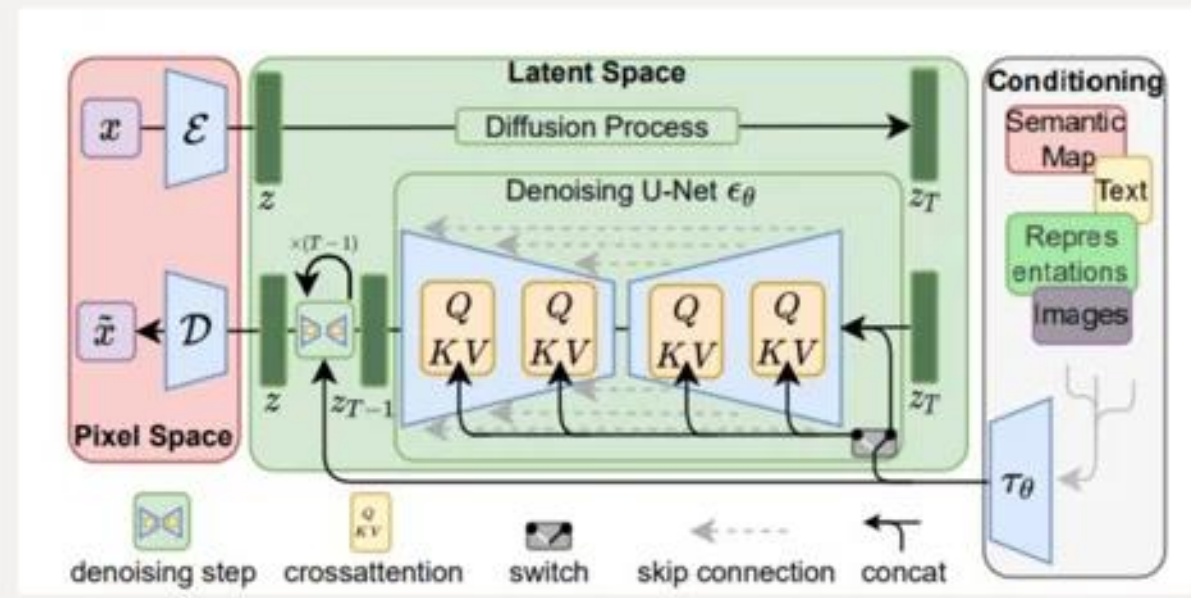
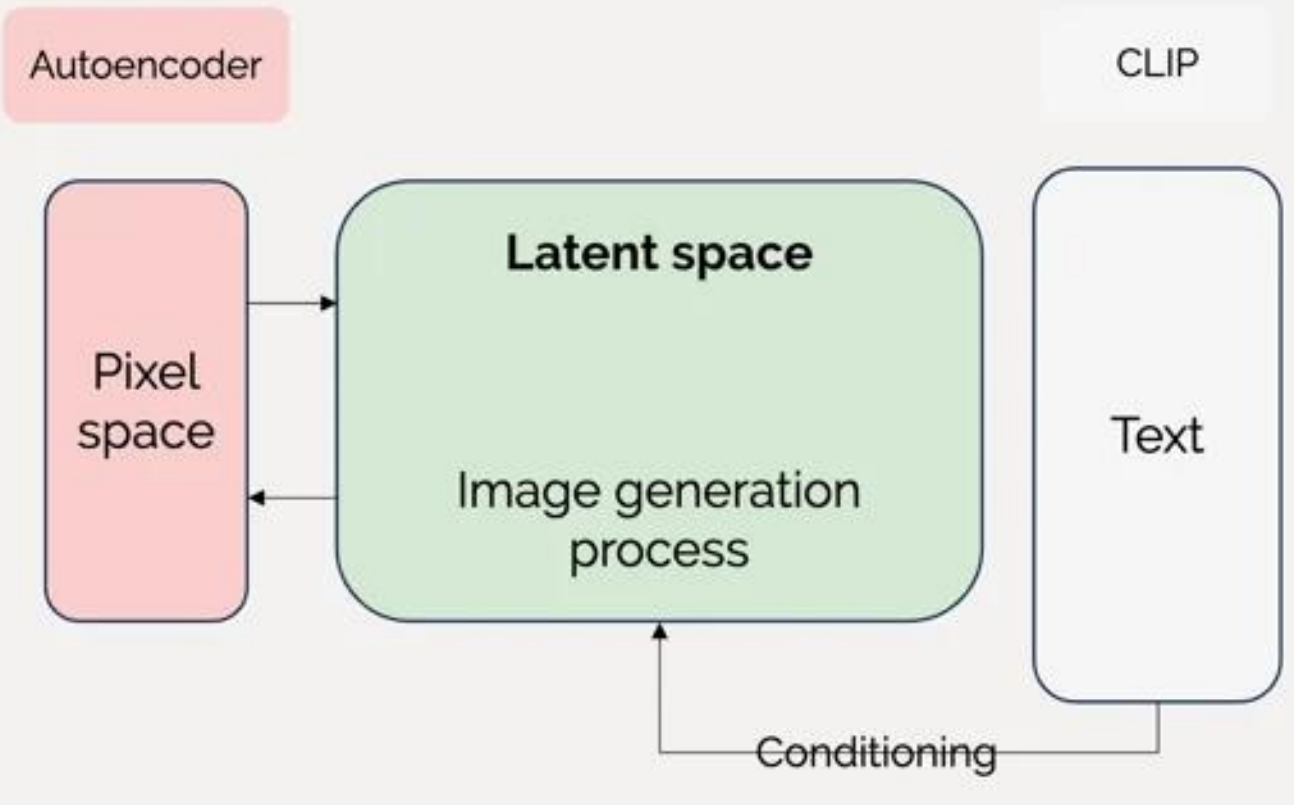
Generative model: learn a distribution

- We should learn a distribution of all possible configurations instead of simply regression.
- Diffusion model for conditional generation!

Main contents

- 3D Reconstruction from single Image
- **2D Diffusion Model for Novel-view Synthesis**
 - Novel-view Diffusion Models
 - Multi-view Image Diffusion Models
- Sync 2D Diffusion & 3D Reconstruction

Image Diffusion Model



High-Resolution Image Synthesis with Latent Diffusion Models, Rombach et al, 2022

Large-scale Training

cat



Conception animale d'illustration de chat de probl...



10+ Times 'Stupid Cat Drawings' Made Everyone Laug...



laik si t pasa
Gato, El Gato, and Laik: otra vez se me ha bugeado...



Dibujos realistas gatito - en la alfombra



Mouser painting by Lynda Nolte



"Oil on Canvas - "Dignan" - Dignan, the loner, w...



Drawings On Black Paper Ideas ; Drawings On Black ...



You bite that box Michael
Michael, Stuff, and Cat: This is Michael. Michael ...



What is this cat doing? - more at megacutie.co.uk ...



42+ Times 'Stupid Cat Drawings' Made



maine-coon-black-cat-portrait



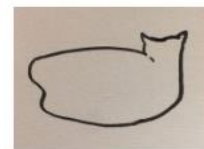
明日はなに描こうかな〜#ねこ#猫#猫のいる暮らし#ねこ部#愛猫#art#アート#作品...



#cats #DailyDoodle



Когда все думают, что ты рисуешь какую-то хрень, а...



Poorly Drawn Cat

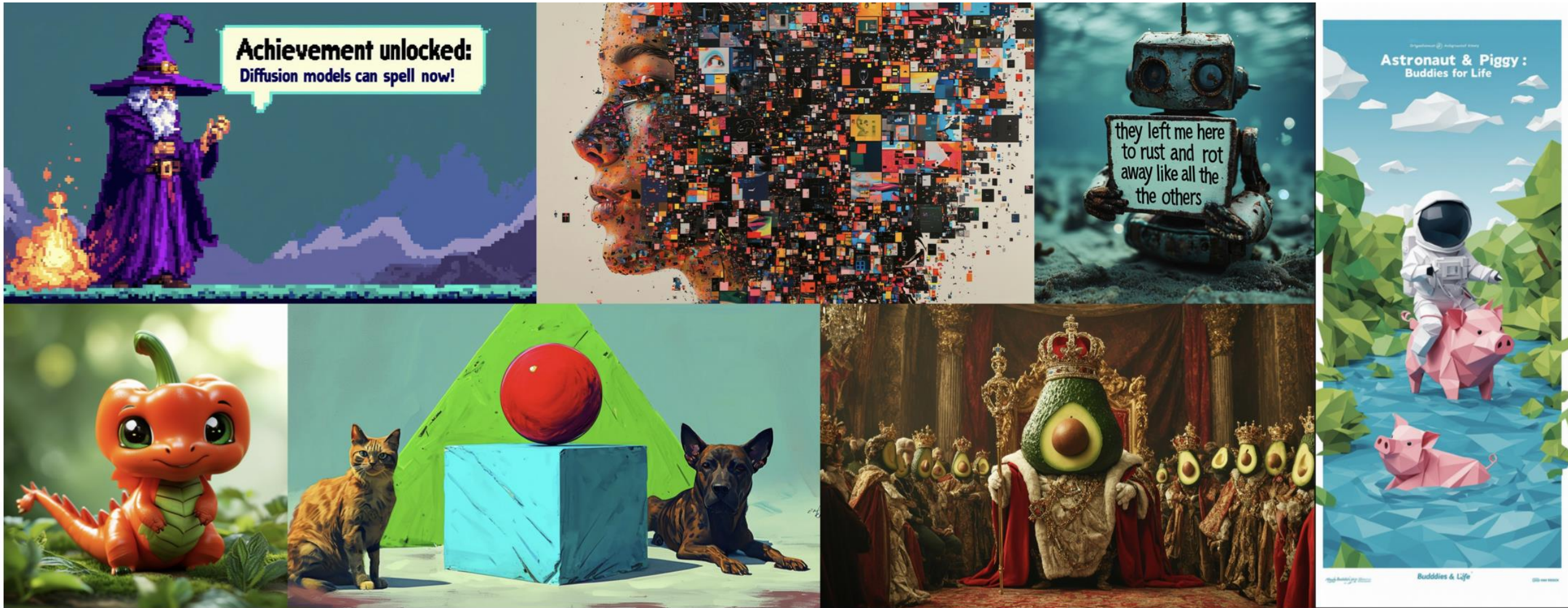


Cat, Neko Poster Fonts, Japanese Poster, Japanese ...

Laion5B:

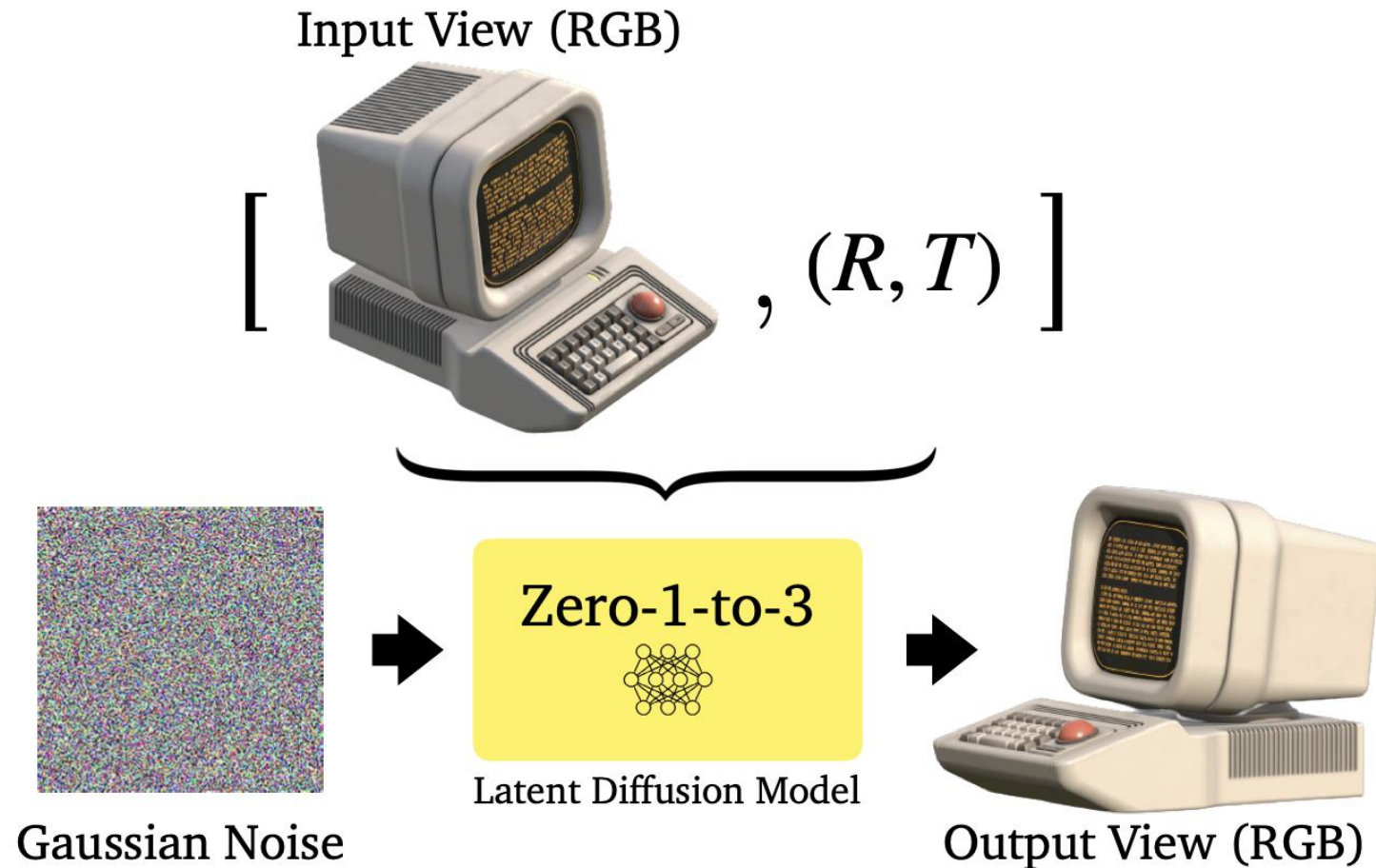
- 5 Billion images
- With Text annotations

Superior image generation quality



Novel-view Diffusion Model

- Leverage Image diffusion prior, generate desired novel view image

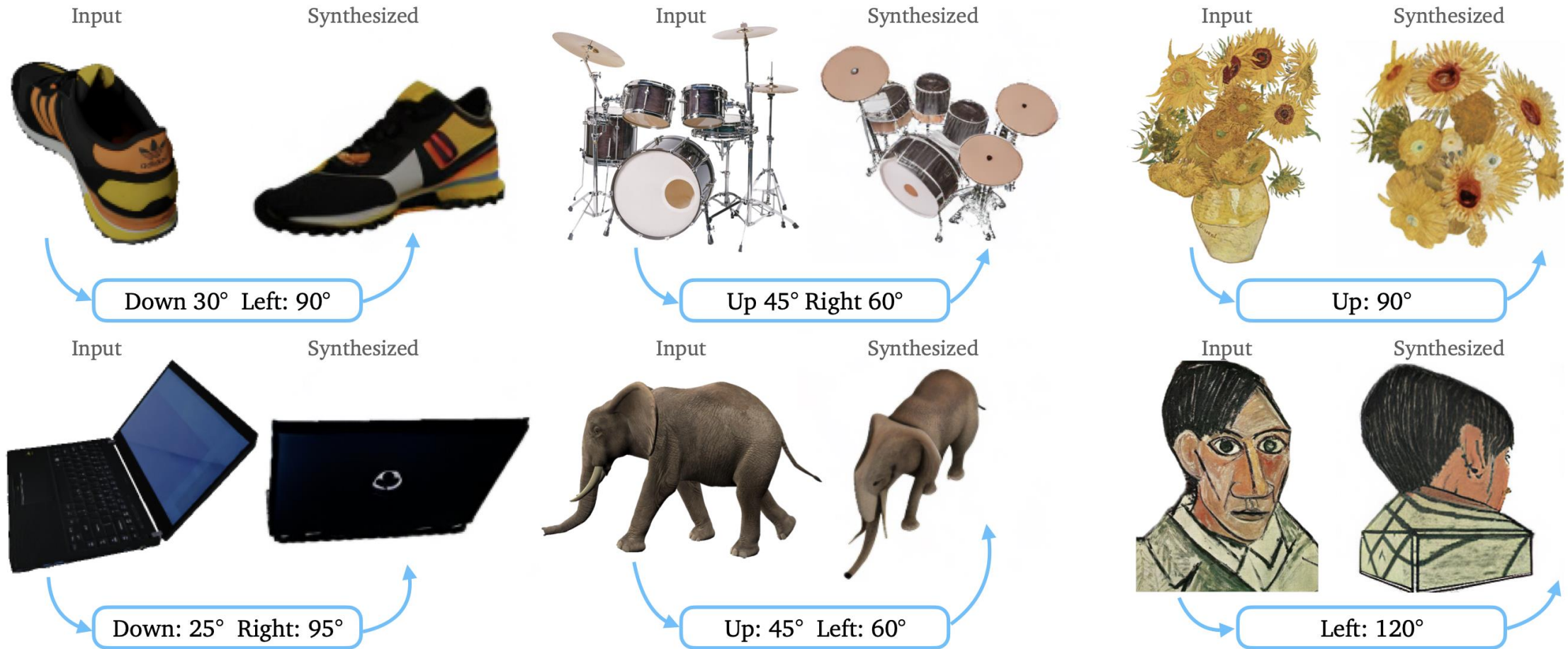


Large-scale Training



Objaverse:
- 800K 3D objects

Novel-view Diffusion Model



Novel-view Diffusion Model

- Limitation: each view generated individually, inconsistent across generation



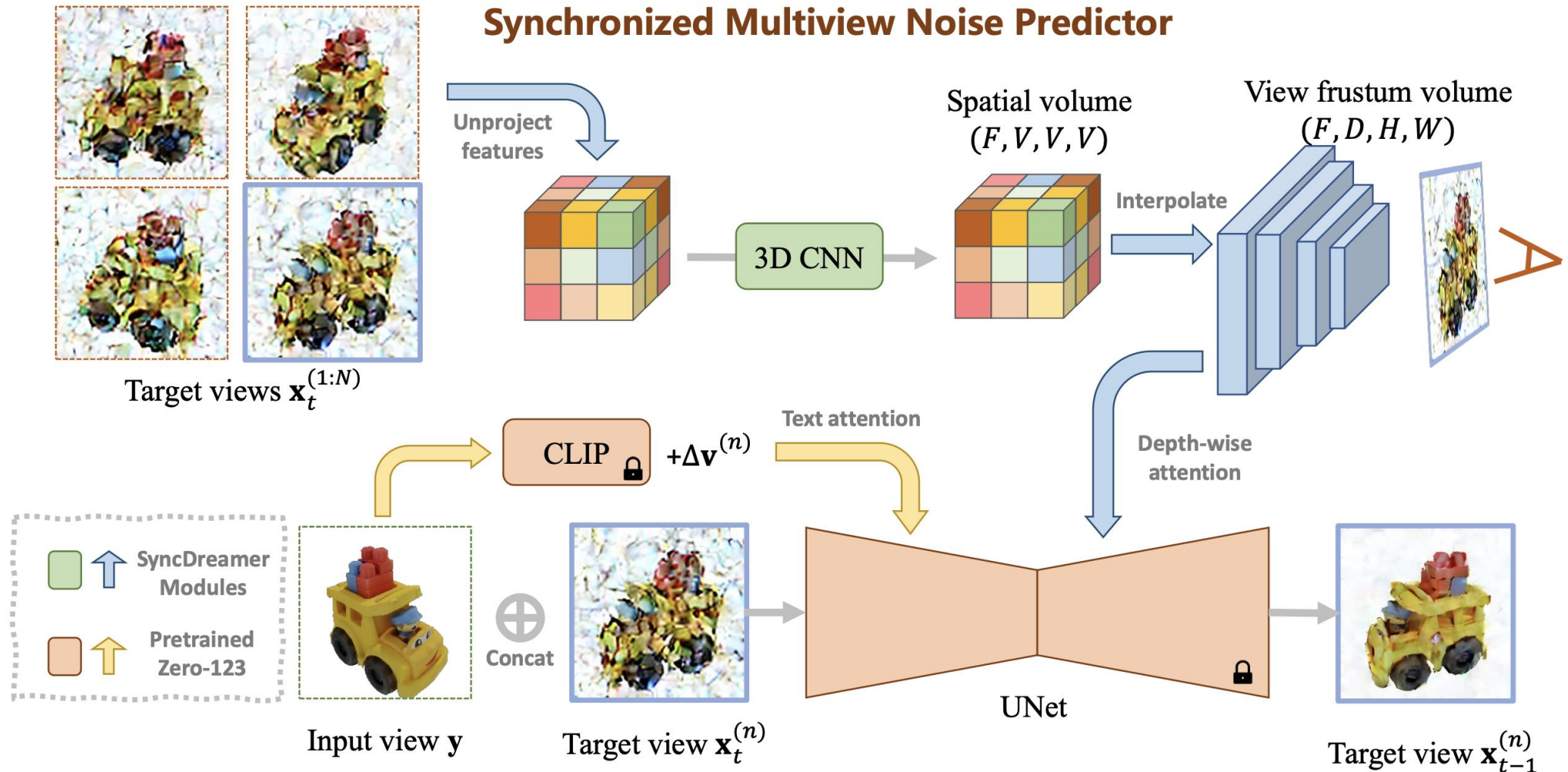
Input Image



Zero123-XL, Liu et al.

Multi-view Diffusion Model

- Generate Multiple Views simultaneously, have more consistency



Multi-view Diffusion Model

- Compared to Single-view Diffusion, the multi-view diffusion models are more consistent across generated novel views



Input Image



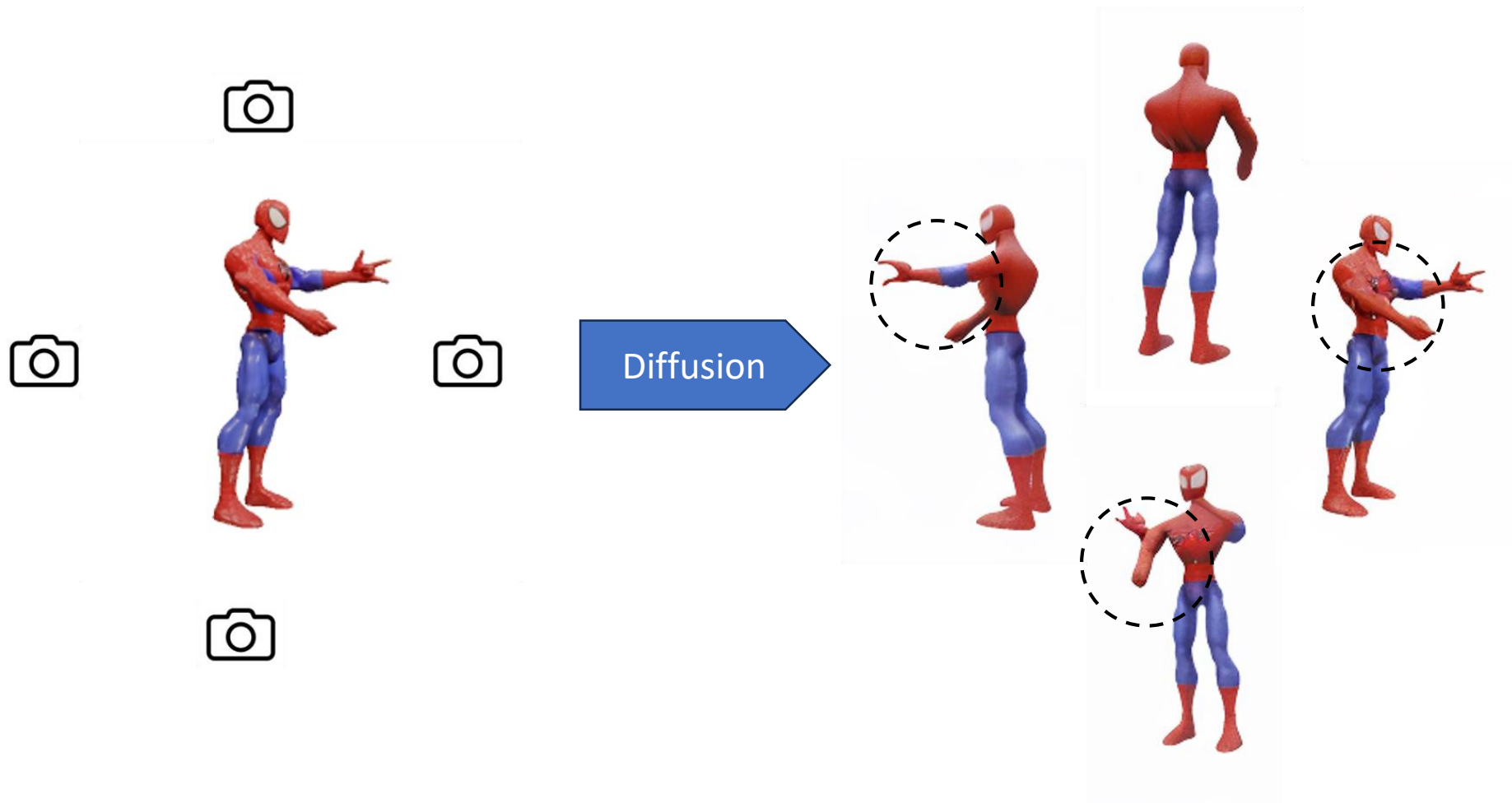
Single-view Diffusion



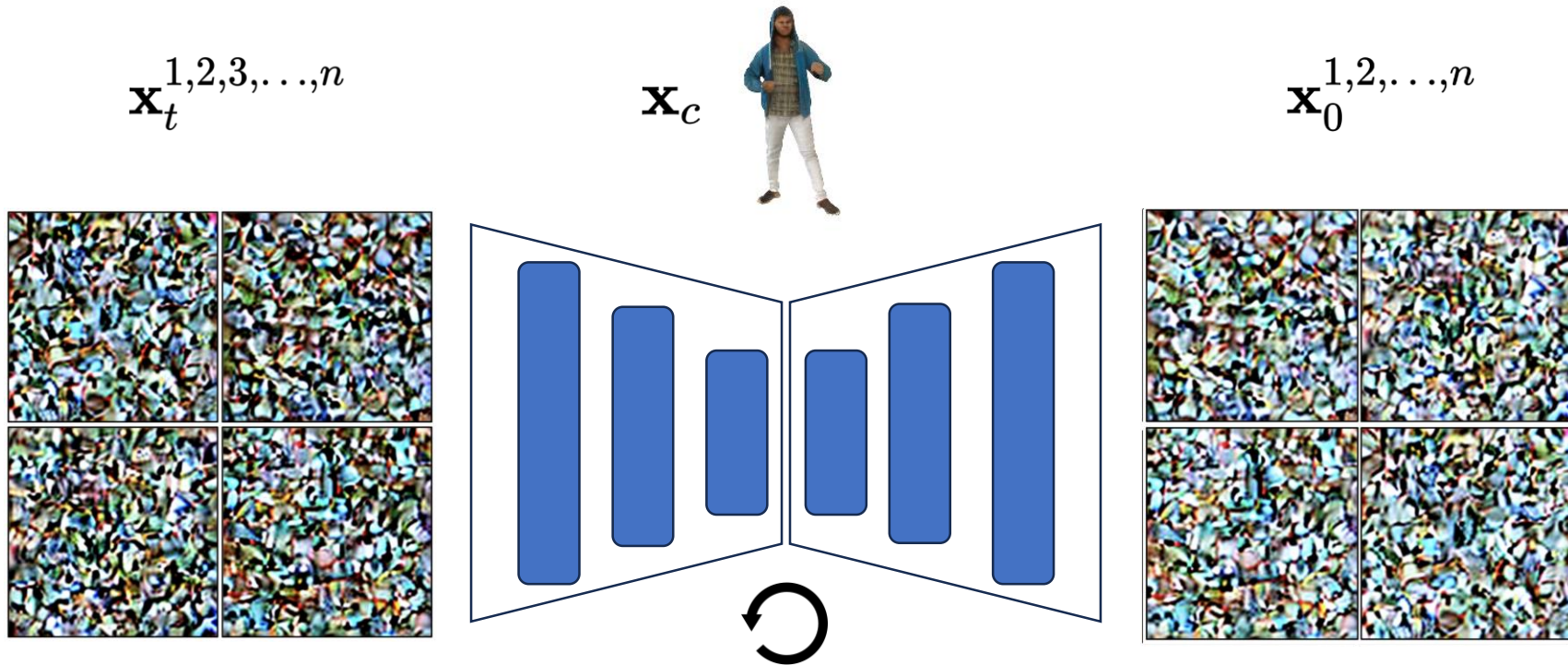
Multi-view Diffusion

3D Consistency in 2D Multi-view Diffusion Model

- 2D Multi-view Diffusion has no explicit 3D representation (e.g. NeRF or 3D-GS). Thus, the 3D consistency of the generated images are not constrained.

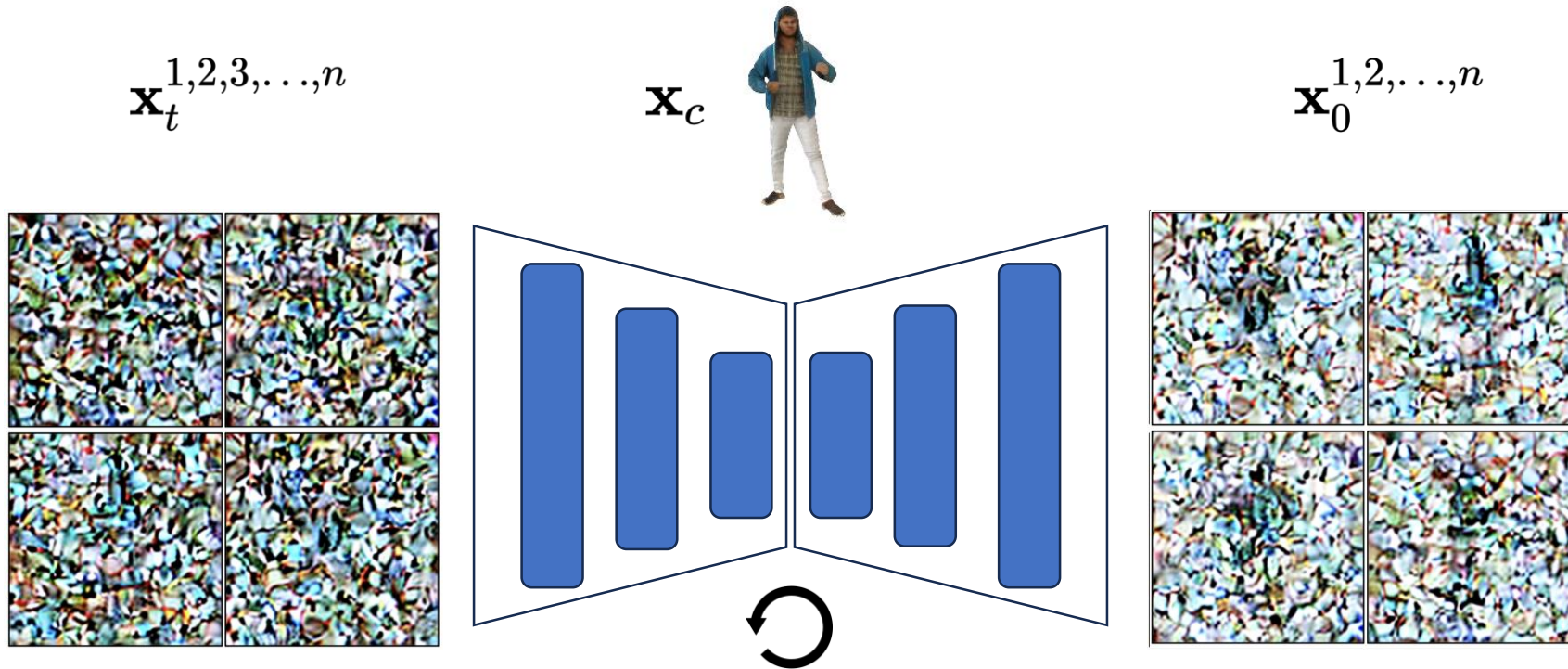


Multi-view Diffusion Model



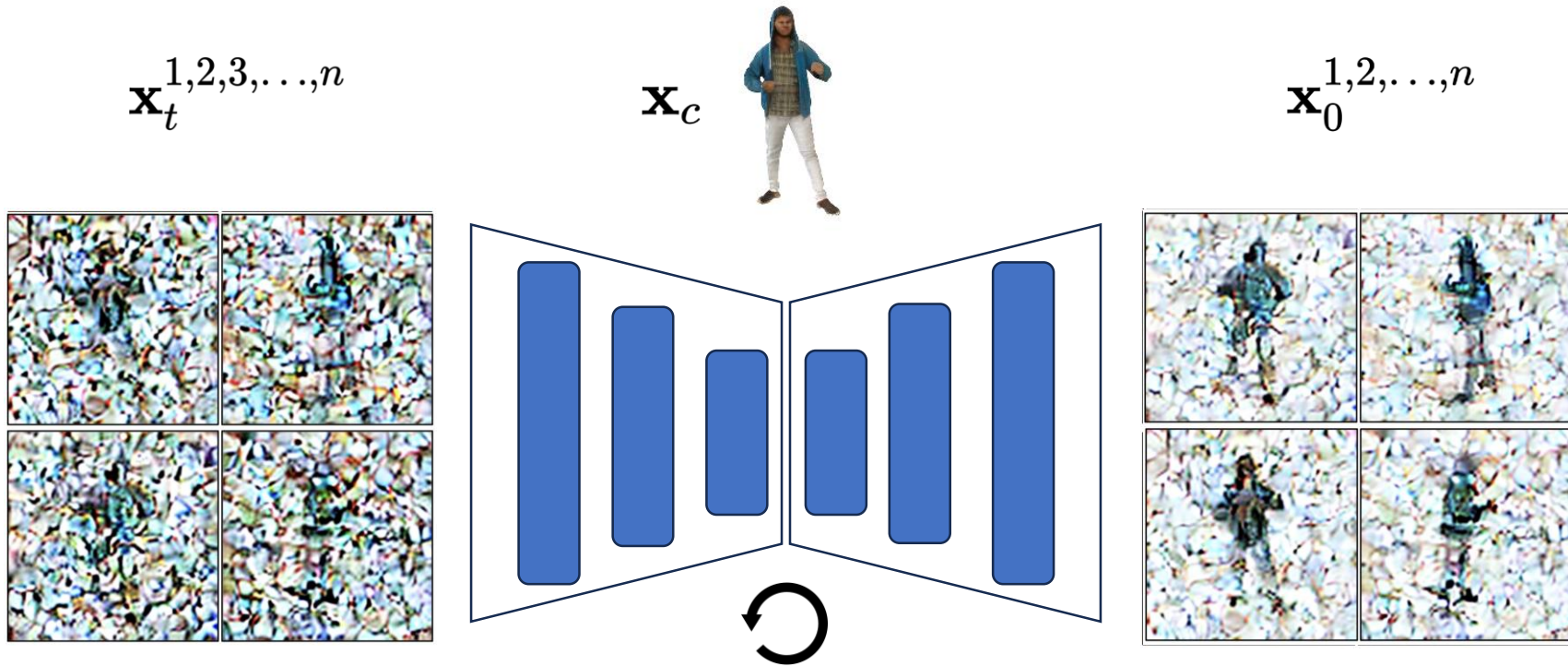
Pre-trained on **5B** 2D images and **800K** 3D objects

Multi-view Diffusion Model



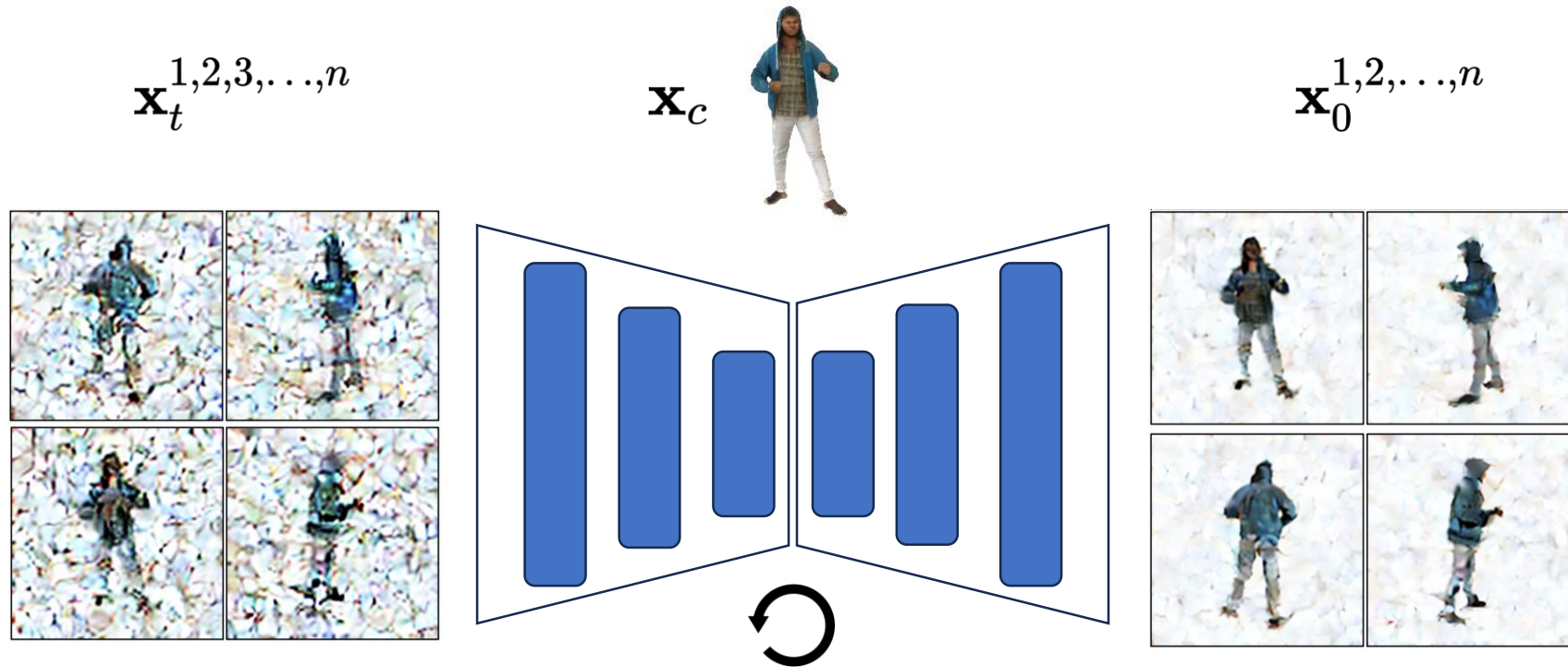
Pre-trained on **5B** 2D images and **800K** 3D objects

Multi-view Diffusion Model



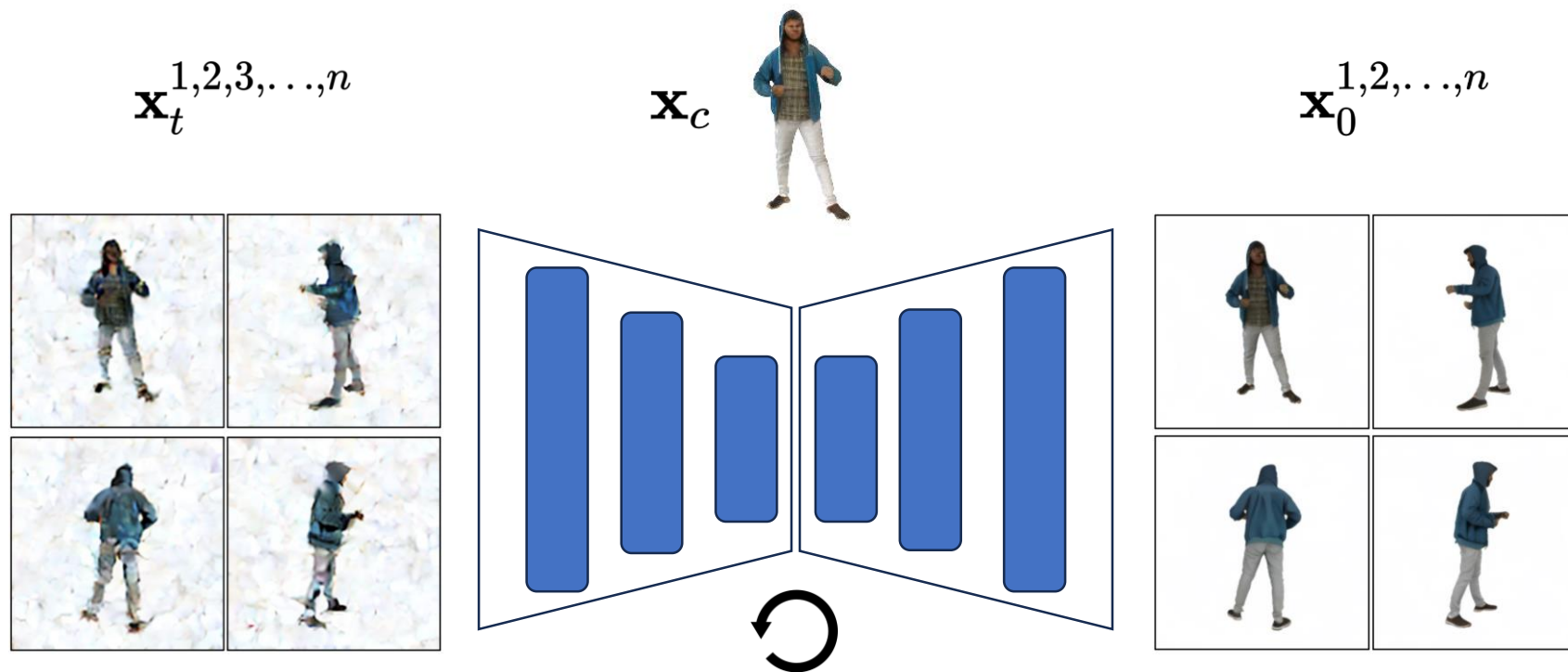
Pre-trained on **5B** 2D images and **800K** 3D objects

Multi-view Diffusion Model



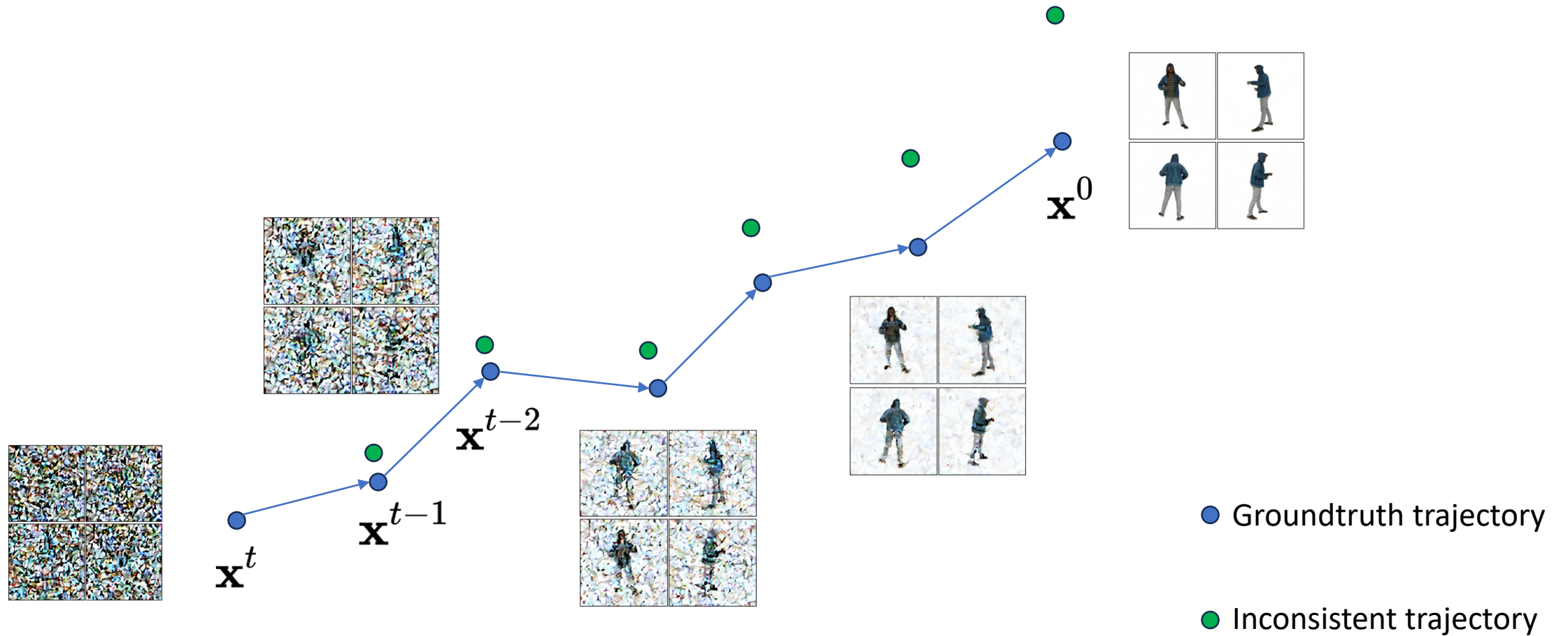
Pre-trained on **5B** 2D images and **800K** 3D objects

Multi-view Diffusion Model



Pre-trained on **5B** 2D images and **800K** 3D objects

Inconsistency accumulates along trajectory

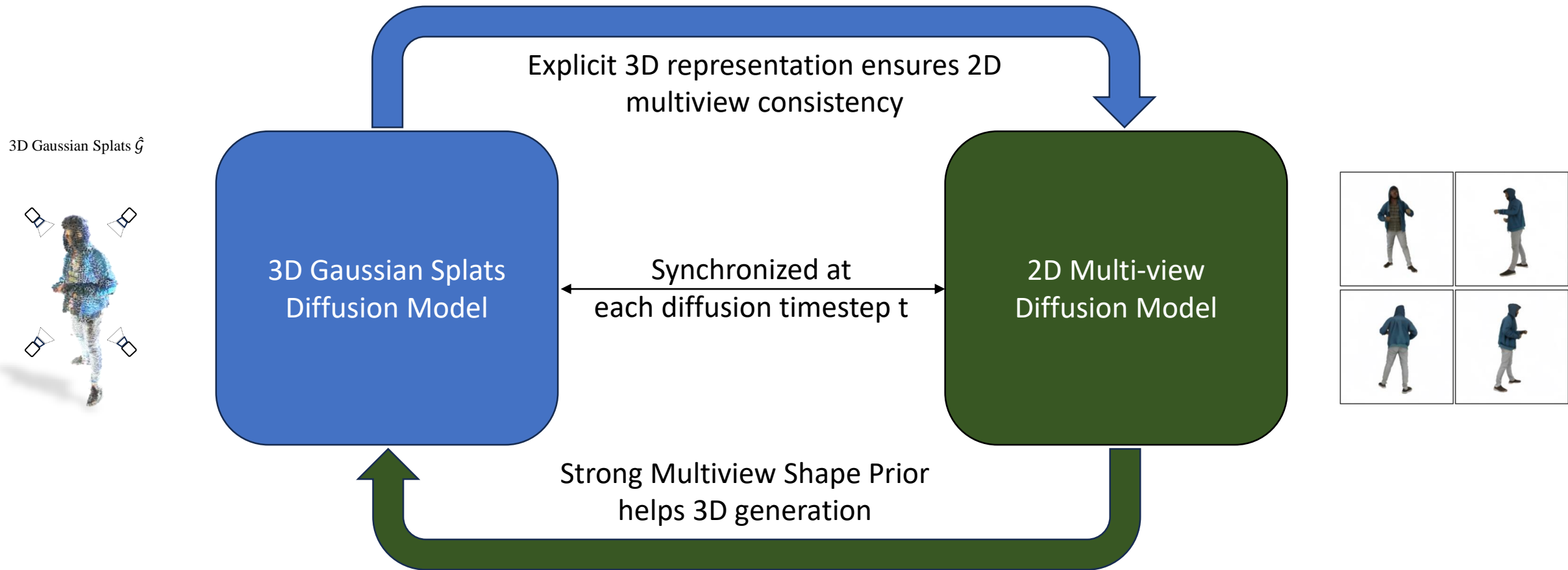


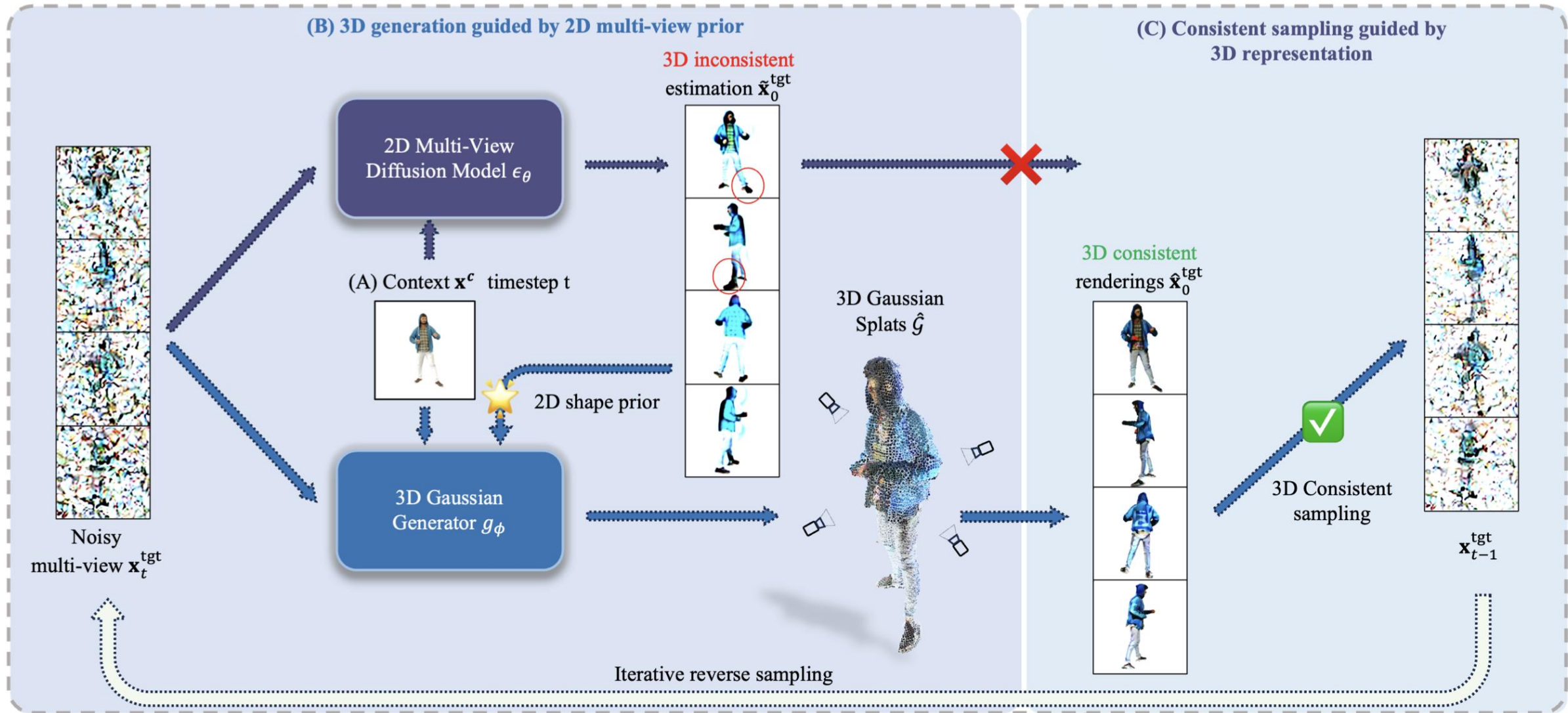
Main contents

- 3D Reconstruction from single Image
- 2D Diffusion Model for Novel-view Synthesis
 - Novel-view Diffusion Models
 - Multi-view Image Diffusion Models
- **Sync 2D Diffusion & 3D Reconstruction**

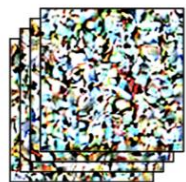
Gen-3Diffusion: Sync 2D Diffusion & 3D Recon

- 2D Diffusion leverages Image Diffusion Prior
- 3D Reconstructor provides 3D representation





Gaussian Noise



3D Consistent Reverse Sampling Step



(D) 3D Gaussian Splats \mathcal{G}



Algorithm

Algorithm 1 Joint 2D & 3D Diffusion Training

Input: Dataset of posed multi-view images $\mathbf{x}_0^{\text{tgt}}, \pi^{\text{tgt}}, \mathbf{x}_0^{\text{novel}}, \pi^{\text{novel}}$, a context image \mathbf{x}^c , text description y

Output: Optimized 2D multi-view diffusion model ϵ_θ and 3D-GS generative model g_ϕ

- 1: **repeat**
 - 2: $\{\mathbf{x}_0^{\text{tgt}}, \mathbf{x}_0^{\text{novel}}, \mathbf{x}^c, y\} \sim q(\{\mathbf{x}_0^{\text{tgt}}, \mathbf{x}_0^{\text{novel}}, \mathbf{x}^c, y\})$
 - 3: $t \sim \text{Uniform}(\{1, \dots, T\}); \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 4: $\mathbf{x}_t^{\text{tgt}} = \sqrt{\bar{\alpha}_t} \mathbf{x}_0^{\text{tgt}} + \sqrt{1 - \bar{\alpha}_t} \epsilon$
 - 5: $\tilde{\mathbf{x}}_0^{\text{tgt}} = \frac{1}{\sqrt{\bar{\alpha}_t}} (\mathbf{x}_t^{\text{tgt}} - \sqrt{1 - \bar{\alpha}_t} \epsilon_\theta(\mathbf{x}_t^{\text{tgt}}, \mathbf{x}^c, y, t))$
 - 6: $\hat{\mathcal{G}} = g_\phi(\mathbf{x}_t^{\text{tgt}}, t, \mathbf{x}^c, \tilde{\mathbf{x}}_0^{\text{tgt}})$ // Enhance conditional 3D generation with 2D diffusion prior $\tilde{\mathbf{x}}_0^{\text{tgt}}$ from ϵ_θ
 - 7: $\{\hat{\mathbf{x}}_0^{\text{tgt}}, \hat{\mathbf{x}}_0^{\text{novel}}\} = \text{renderer}(\hat{\mathcal{G}}, \{\pi^{\text{tgt}}, \pi^{\text{novel}}\})$
 - 8: Compute loss $\mathcal{L}_{\text{total}}$ (Eq. (9))
 - 9: Gradient step to update ϵ_θ, g_ϕ
 - 10: **until** converged
-

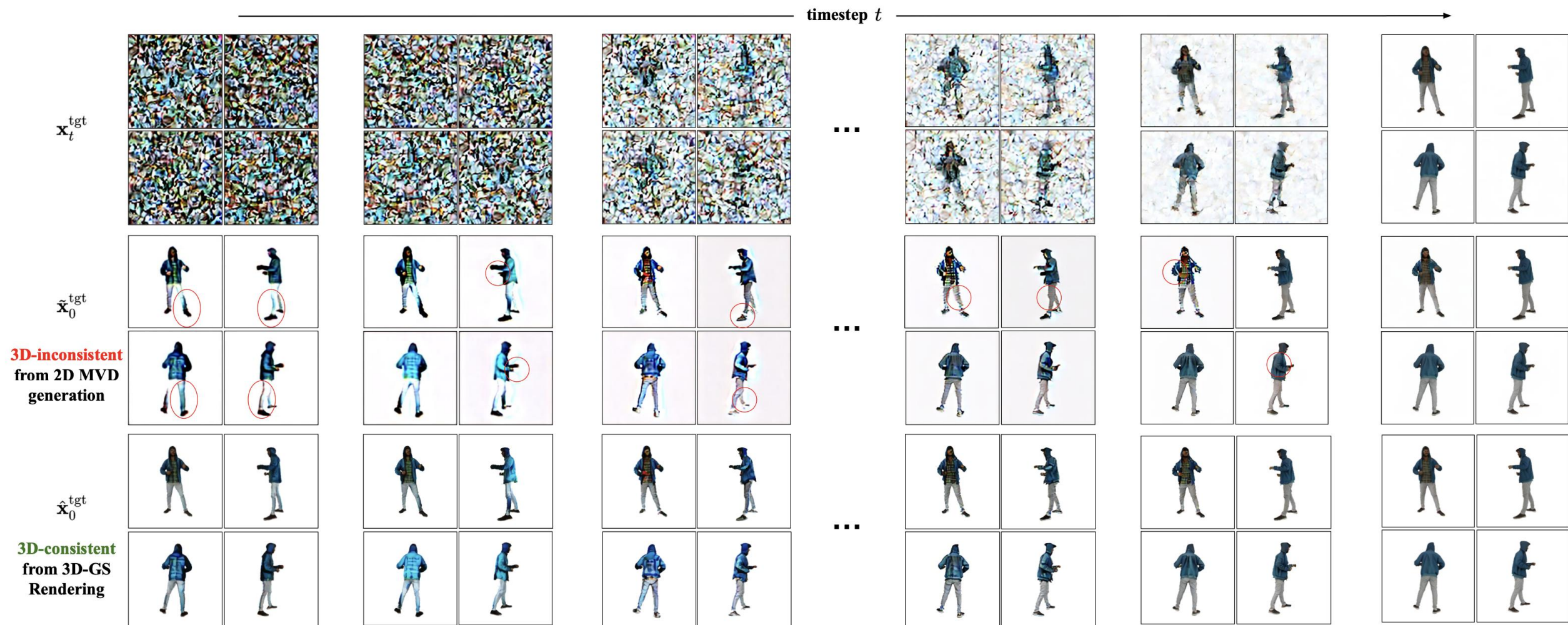
Algorithm 2 3D Consistent Guided Sampling

Input: A context image \mathbf{x}^c and text y ; Converged 2D diffusion model ϵ_θ and 3D generative model g_ϕ

Output: 3D Gaussian Splats \mathcal{G} of the 2D image \mathbf{x}^c

- 1: $\mathbf{x}_T^{\text{tgt}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 2: **for** $t = T, \dots, 1$ **do**
 - 3: $\tilde{\mathbf{x}}_0^{\text{tgt}} = \frac{1}{\sqrt{\bar{\alpha}_t}} (\mathbf{x}_t^{\text{tgt}} - \sqrt{1 - \bar{\alpha}_t} \epsilon_\theta(\mathbf{x}_t^{\text{tgt}}, \mathbf{x}^c, y, t))$
 - 4: $\hat{\mathcal{G}} = g_\phi(\mathbf{x}_t^{\text{tgt}}, t, \mathbf{x}^c, \tilde{\mathbf{x}}_0^{\text{tgt}})$
 - 5: $\hat{\mathbf{x}}_0^{\text{tgt}} = \text{renderer}(\hat{\mathcal{G}}, \pi^{\text{tgt}})$
 - 6: $\mu_{t-1}(\mathbf{x}_t^{\text{tgt}}, \hat{\mathbf{x}}_0^{\text{tgt}}) = \frac{\sqrt{\alpha_t(1-\bar{\alpha}_{t-1})}}{1-\bar{\alpha}_t} \mathbf{x}_t^{\text{tgt}} + \frac{\sqrt{\bar{\alpha}_{t-1}\beta_t}}{1-\bar{\alpha}_t} \hat{\mathbf{x}}_0^{\text{tgt}}$ // Guide 2D sampling with 3D consistent multi-view renderings
 - 7: $\mathbf{x}_{t-1}^{\text{tgt}} \sim \mathcal{N}(\mathbf{x}_{t-1}^{\text{tgt}}; \tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t^{\text{tgt}}, \hat{\mathbf{x}}_0^{\text{tgt}}), \tilde{\boldsymbol{\beta}}_{t-1} \mathbf{I})$
 - 8: **end for**
 - 9: **return** $\mathcal{G} = g_\phi(\mathbf{x}_0^{\text{tgt}}, \tilde{\mathbf{x}}_0^{\text{tgt}}, \mathbf{x}^c, t = 0)$
-

Explicit 3D-GS helps 2D Diffusion



Results: Object Reconstruction

Reconstruction w/ novel view diffusion



Input Image



Gen-3Diffusion



Single-view Diffusion



Multi-view Diffusion

Reconstruction w/ novel view diffusion



Input Image



Gen-3Diffusion



Single-view Diffusion



Multi-view Diffusion

Reconstruction w/ novel view diffusion



Input Image



Gen-3Diffusion



Single-view Diffusion



Multi-view Diffusion

Reconstruction w/ novel view diffusion



Input Image



Gen-3Diffusion



Single-view Diffusion



Multi-view Diffusion

Reconstruction w/ direct 3d reconstruction



Input Image



Gen-3Diffusion



TripoSR



LGM

Reconstruction w/ direct 3d reconstruction



Input Image



Gen-3Diffusion



TripoSR



LGM

Reconstruction w/ direct 3d reconstruction



Input Image



Gen-3Diffusion



TripoSR



LGM

Reconstruction w/ direct 3d reconstruction



Input Image



Gen-3Diffusion



TripoSR



LGM

Results: Avatar Reconstruction

Reconstruction avatar appearance



Input Image



Gen-3Diffusion

SiTH

SiFU

Reconstruction avatar appearance



Input Image



Gen-3Diffusion

SiTH

SiFU

Children



Input Image



Gen-3Diffusion



SiTH



SiFU

Reconstruction avatar geometry



Input Image



Gen-3Diffusion



ICON



ECON

Reconstruction avatar geometry



Input Image



Gen-3Diffusion



ICON



ECON

Strong generalization



Input Image



Gen-3Diffusion

SiTH

SiFU

ECON

Strong generalization

3D-GS Rendering



3D-GS Rendering



3D-GS Rendering



3D-GS Rendering



3D-GS Rendering



3D-GS Rendering



Take away messages

- With training on massive data, image diffusion models achieve superior image generation quality
- The image diffusion prior can be applied to 3D tasks, e.g. generate novel views
- Novel view diffusion models lack 3D consistency because they don't have an explicit 3D representation
- 2D Diffusion Models and 3D Reconstruction Models can be combined to achieve excellent reconstruction capability