

# DMV3D: Denoising Multi-view Diffusion Using **3D Large Reconstruction Model**

ICLR 2024 submission (avg. score 8)

&

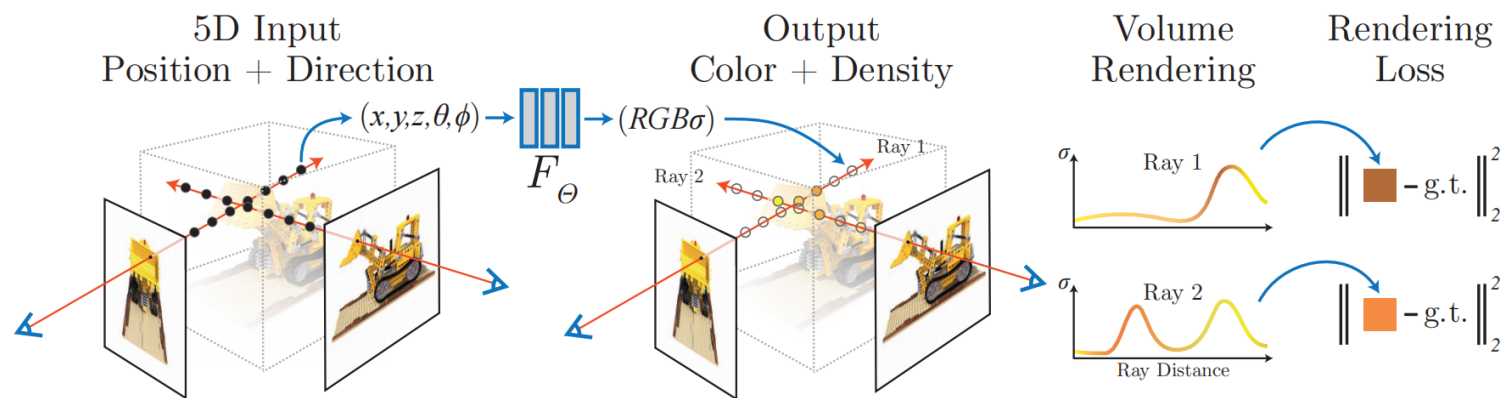
# LRM: **Large Reconstruction Model** for Single Image to 3D

ICLR 2024 submission (avg. score 8.5)

Presenter: Rundong Luo

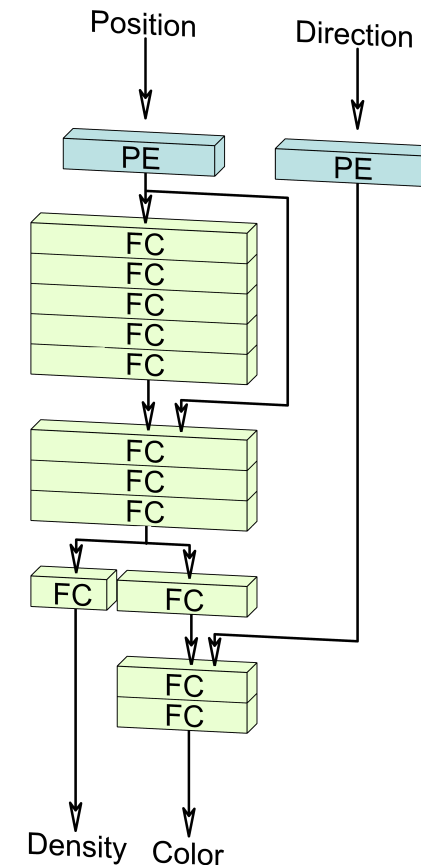
2023.12.10

## Overview



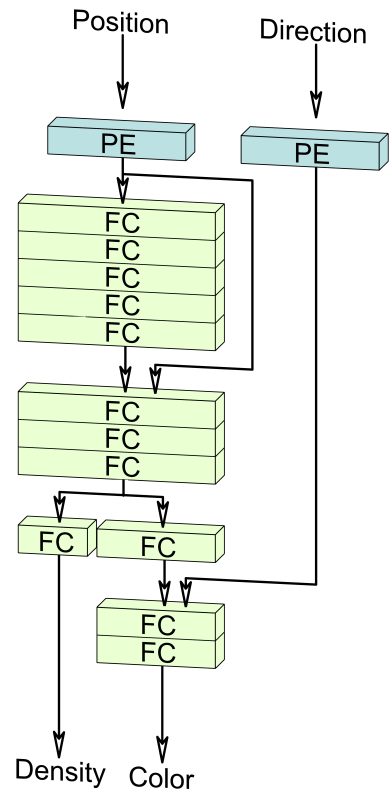
## Volume rendering

$$C(\mathbf{r}) = \int_{t_n}^{t_f} T(t) \sigma(\mathbf{r}(t)) \mathbf{c}(\mathbf{r}(t), \mathbf{d}) dt, \text{ where } T(t) = \exp\left(-\int_{t_n}^t \sigma(\mathbf{r}(s)) ds\right).$$

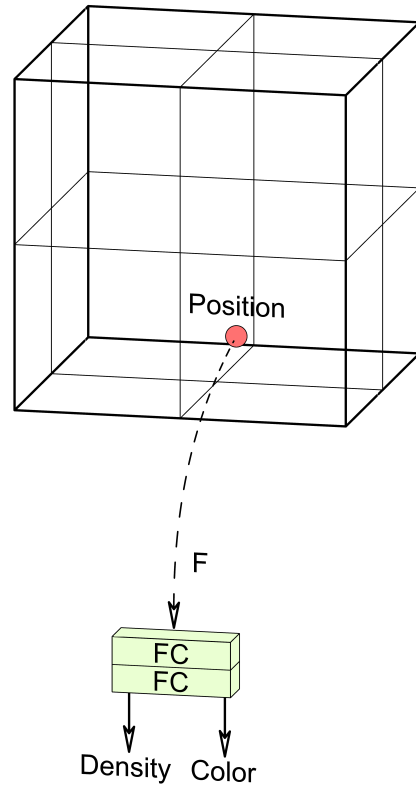


(a) NeRF (Implicit)

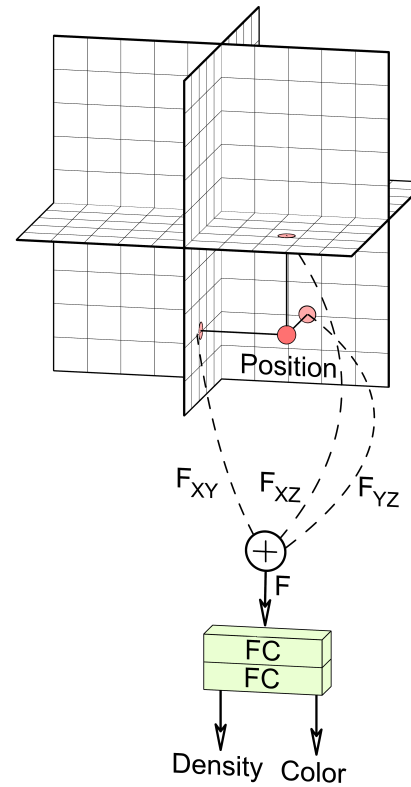
# Background: Triplane Representation for NeRF



(a) NeRF (Implicit)



(b) Voxels (Explicit or Hybrid)



(c) Ours (Hybrid)

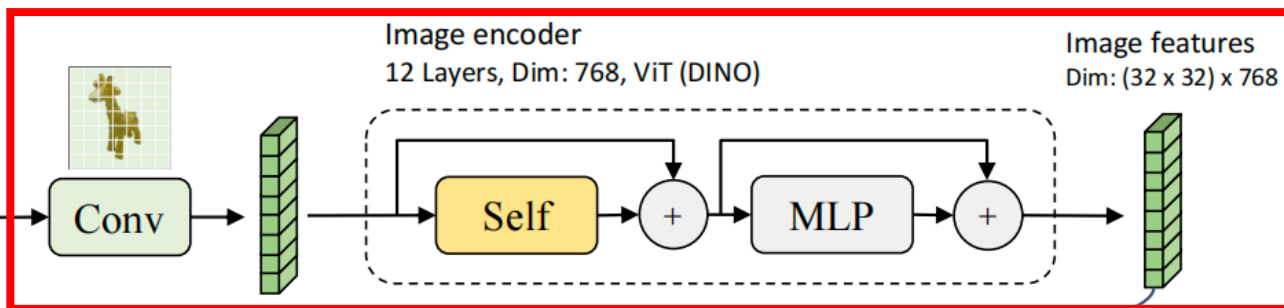
“The primary advantage of this hybrid representation is efficiency.”

- Task: Single-image to 3D (Triplane representation → Mesh)
- Overview: The first large-scale (500M params) 3D reconstruction model
  - Trained on one million 3D shapes and video data across diverse categories
  - Category-agnostic
  - Training objective: simple L2 reconstruction loss
- Performance: Can reconstruct high-fidelity 3D shapes from a wide range of images captured in the real world in five seconds

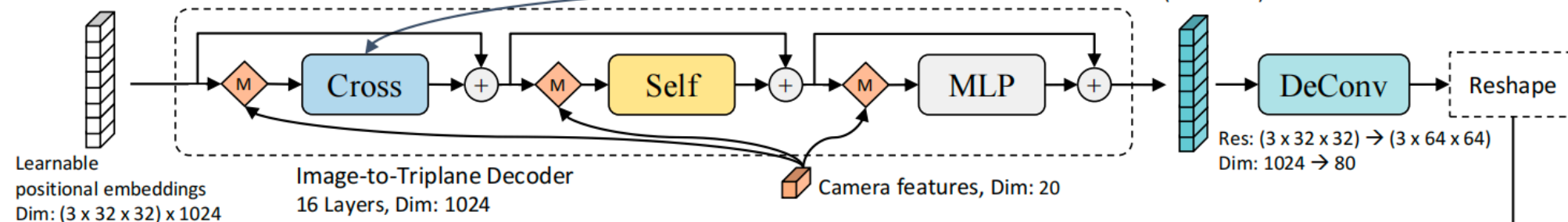
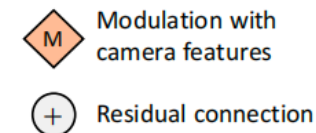


# Large Reconstruction Model: Pipeline

Single input image  
Dim: 512 x 512 x 3

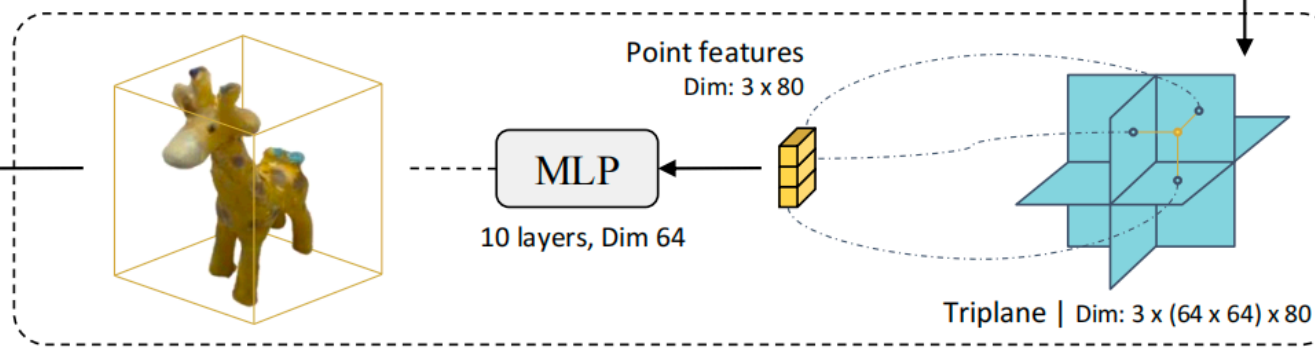


## DINO Encoder



Rendered  
novel image

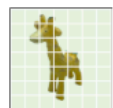
RGB,  $\sigma$   
Volumetric Rendering



Neural Radiance Field (NeRF)

# Large Reconstruction Model: Pipeline

Single input image  
Dim:  $512 \times 512 \times 3$



Conv



Image encoder  
12 Layers, Dim: 768, ViT (DINO)

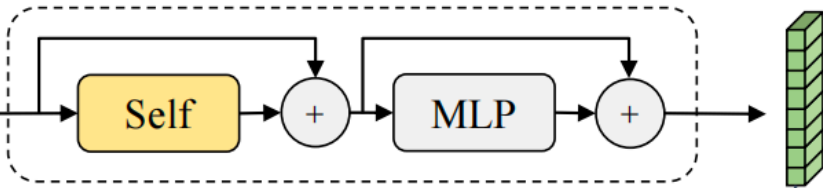
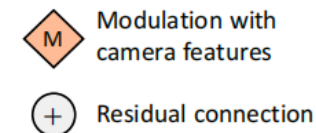
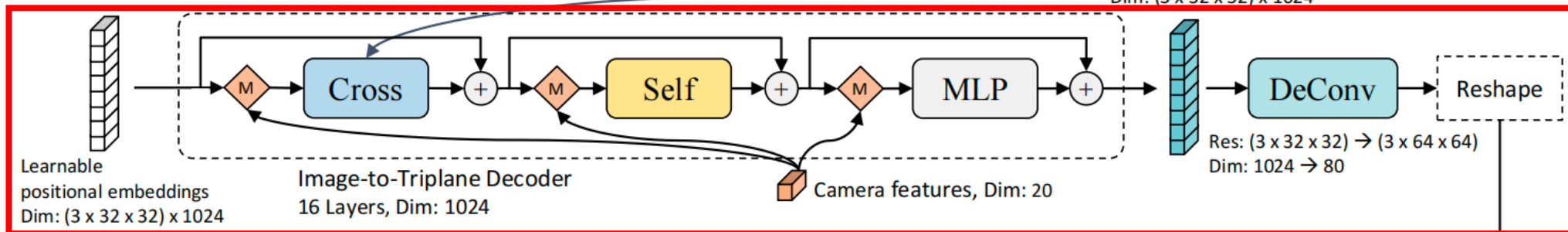


Image features  
Dim:  $(32 \times 32) \times 768$



Triplane tokens  
Dim:  $(3 \times 32 \times 32) \times 1024$

## Triplane Decoder



Rendered  
novel image

RGB,  $\sigma$   
Volumetric Rendering



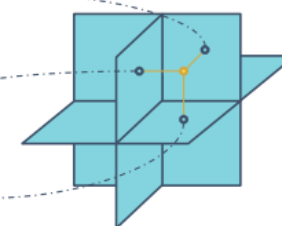
Neural Radiance Field (NeRF)

Point features  
Dim:  $3 \times 80$

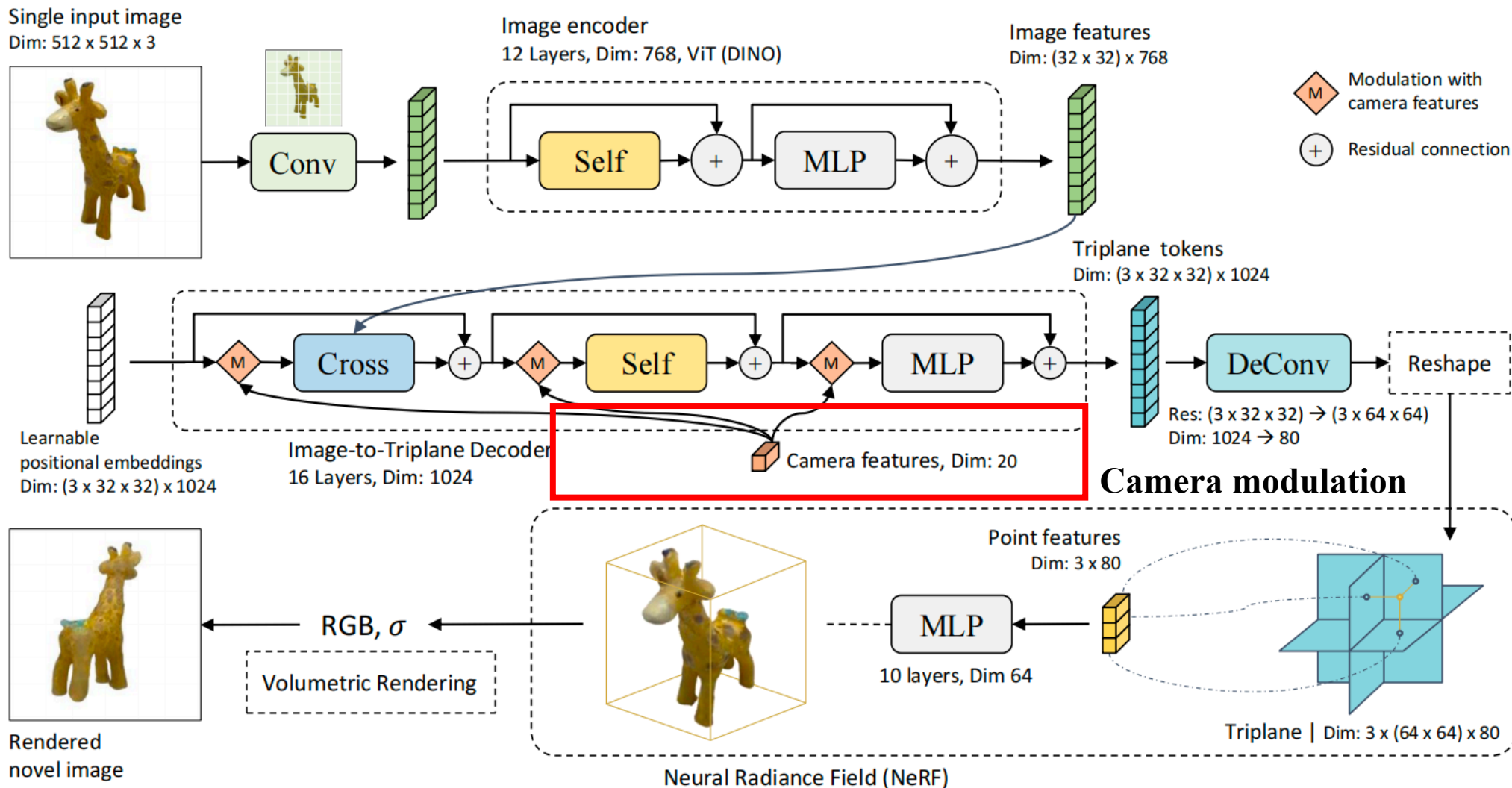
MLP  
10 layers, Dim 64



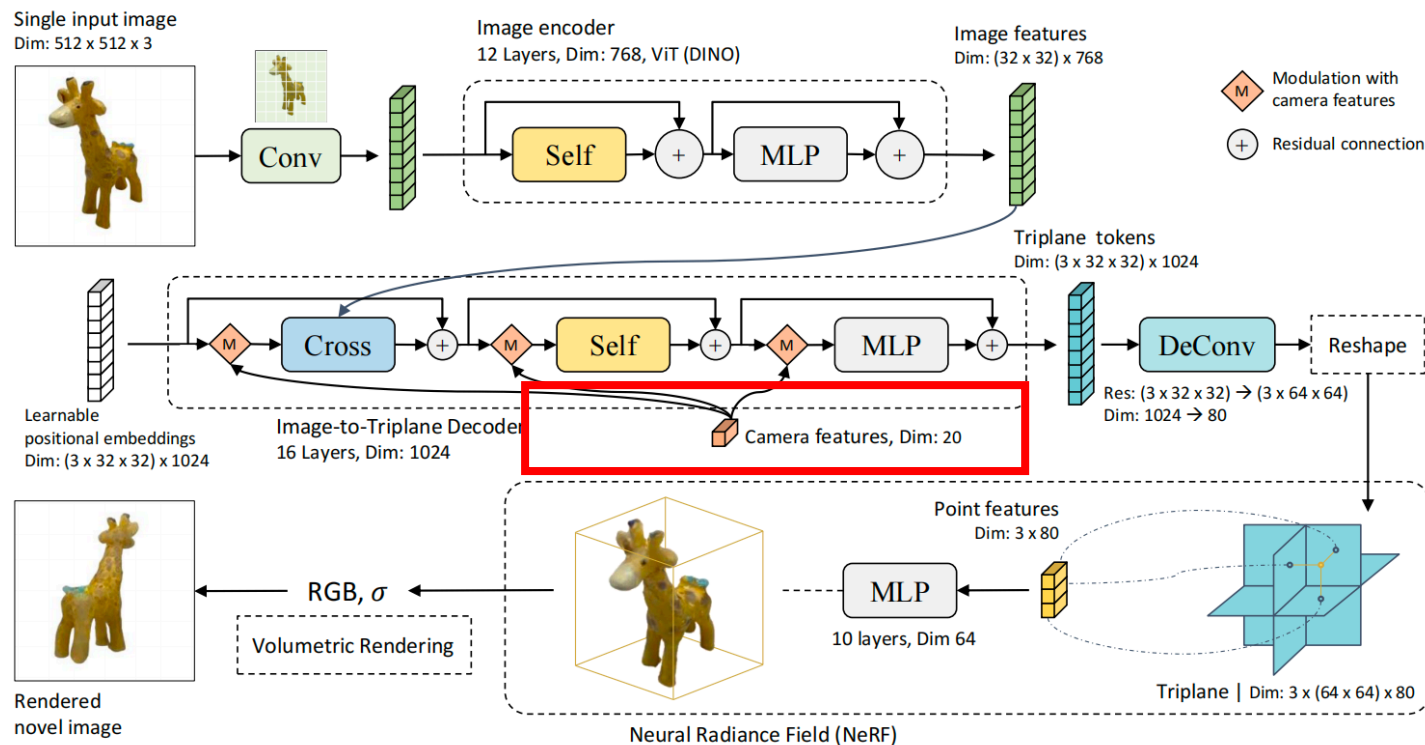
Triplane | Dim:  $3 \times (64 \times 64) \times 80$



# Large Reconstruction Model: Pipeline



# Large Reconstruction Model: Pipeline



## Camera modulation (AdaLN)

$$\mathbf{c} = [\mathbf{E}_{1 \times 16}, f_{oc_x}, f_{oc_y}, pp_x, pp_y]$$

$$\gamma, \beta = \text{MLP}^{\text{mod}}(\tilde{\mathbf{c}})$$

$$\text{ModLN}_c(\mathbf{f}_j) = \text{LN}(\mathbf{f}_j) \cdot (1 + \gamma) + \beta$$

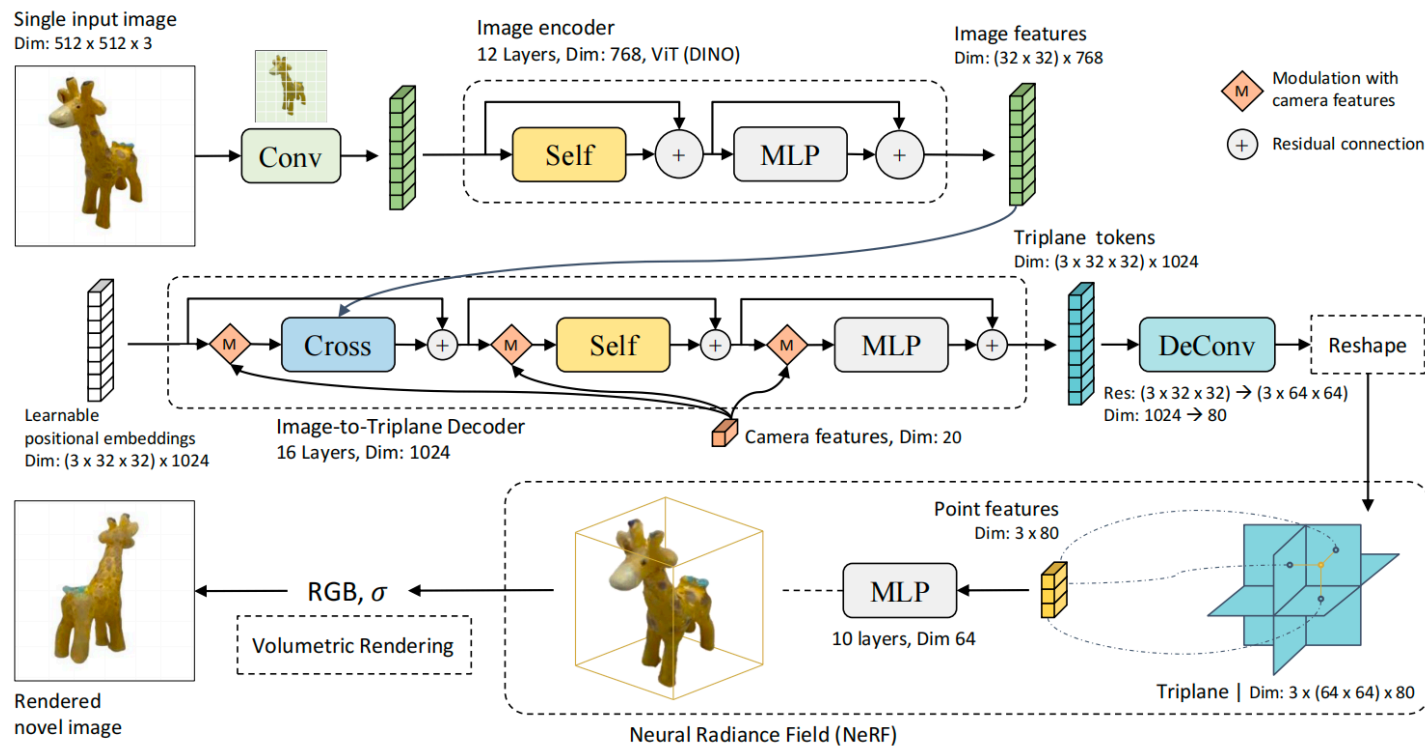
## Decoder (overall)

$$\mathbf{f}_j^{\text{cross}} = \text{CrossAttn}(\text{ModLN}_c(\mathbf{f}_j^{\text{in}}); \{\mathbf{h}_i\}_{i=1}^n) + \mathbf{f}_j^{\text{in}}$$

$$\mathbf{f}_j^{\text{self}} = \text{SelfAttn}(\text{ModLN}_c(\mathbf{f}_j^{\text{cross}}); \{\text{ModLN}_c(\mathbf{f}_{j'}^{\text{cross}})\}_{j'}) + \mathbf{f}_j^{\text{cross}}$$

$$\mathbf{f}_j^{\text{out}} = \text{MLP}^{\text{tfm}}(\text{ModLN}_c(\mathbf{f}_j^{\text{self}})) + \mathbf{f}_j^{\text{self}}$$

# Large Reconstruction Model: Pipeline



## Training objective

$$\mathcal{L}_{\text{recon}}(\mathbf{x}) = \frac{1}{V} \sum_{v=1}^V (\mathcal{L}_{\text{MSE}}(\hat{\mathbf{x}}_v, \mathbf{x}_v^{GT}) + \lambda \mathcal{L}_{\text{LPIPS}}(\hat{\mathbf{x}}_v, \mathbf{x}_v^{GT}))$$


# Large Reconstruction Model: Experiments

---

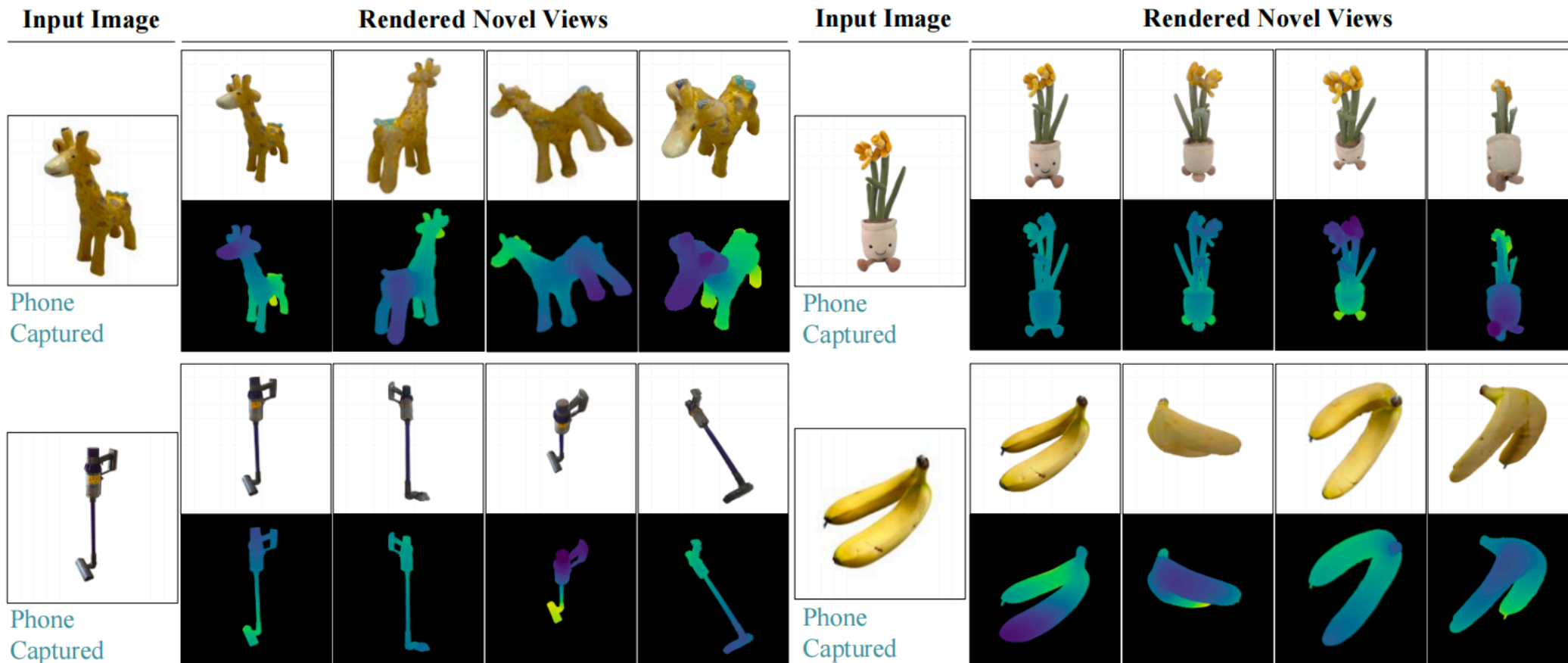
## Datasets

- Objaverse (~730k object meshes)
- MVImgNet (~220k object-centric videos)

Pre-processing: remove background

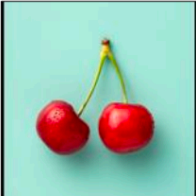
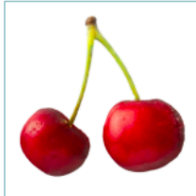
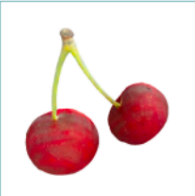


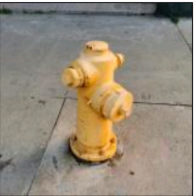



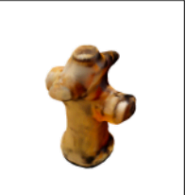






























GPU: 128 A100, 3 days 

# Large Reconstruction Model: Experiments





# Large Reconstruction Model: Experiments

Input Image	Ours				One-2-3-45				Input Image	Ours				One-2-3-45			
																	
																	
																	
																	



- Task: Single-image/Text to 3D (Triplane representation → Mesh)
- Overview: Single-stage framework that leverages multi-view 2D image diffusion model to achieve 3D generation;
  - Trained on one million 3D shapes and video data across diverse categories
  - Training objective: simple L2 reconstruction loss
  - Probabilistic approach, *i.e.*, multiple reasonable 3D outputs given the same input
- Performance: High-quality text-to-3D generation and single-image reconstruction through direct model inference within 30 seconds on an A100 GPU.

# Denoising Multi-view Diffusion: Task Description

Input (posed multi-view images):  $\mathcal{I} = \{\mathbf{I}_1, \dots, \mathbf{I}_N\}$   $\mathcal{C} = \{\mathbf{c}_1, \dots, \mathbf{c}_N\}$

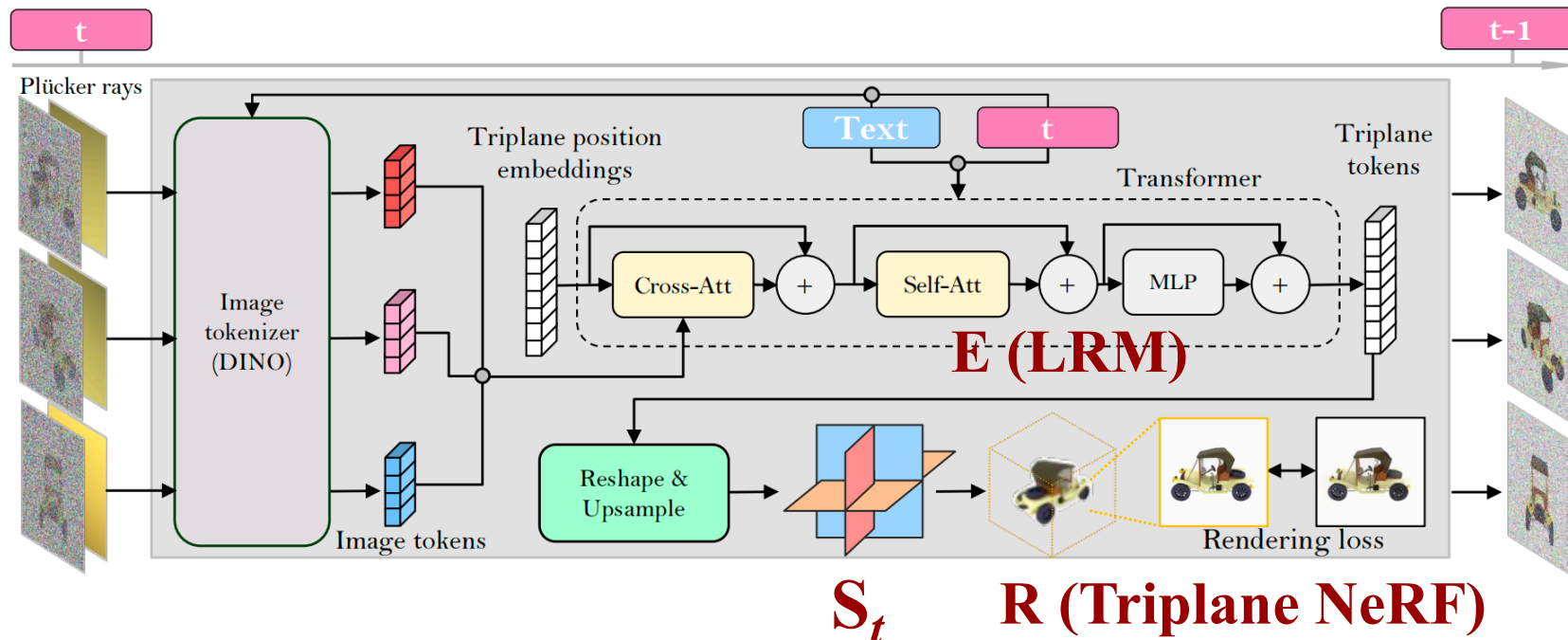
Diffusion process (forward):  $\mathcal{I}_t = \{\sqrt{\bar{\alpha}_t}\mathbf{I} + \sqrt{1 - \bar{\alpha}_t}\epsilon_{\mathbf{I}} | \mathbf{I} \in \mathcal{I}\}$

Denoising process:  $\mathbf{I}_{r,t} = \mathbf{R}(\mathbf{S}_t, \mathbf{c}), \quad \mathbf{S}_t = \mathbf{E}(\mathcal{I}_t, t, \mathcal{C})$

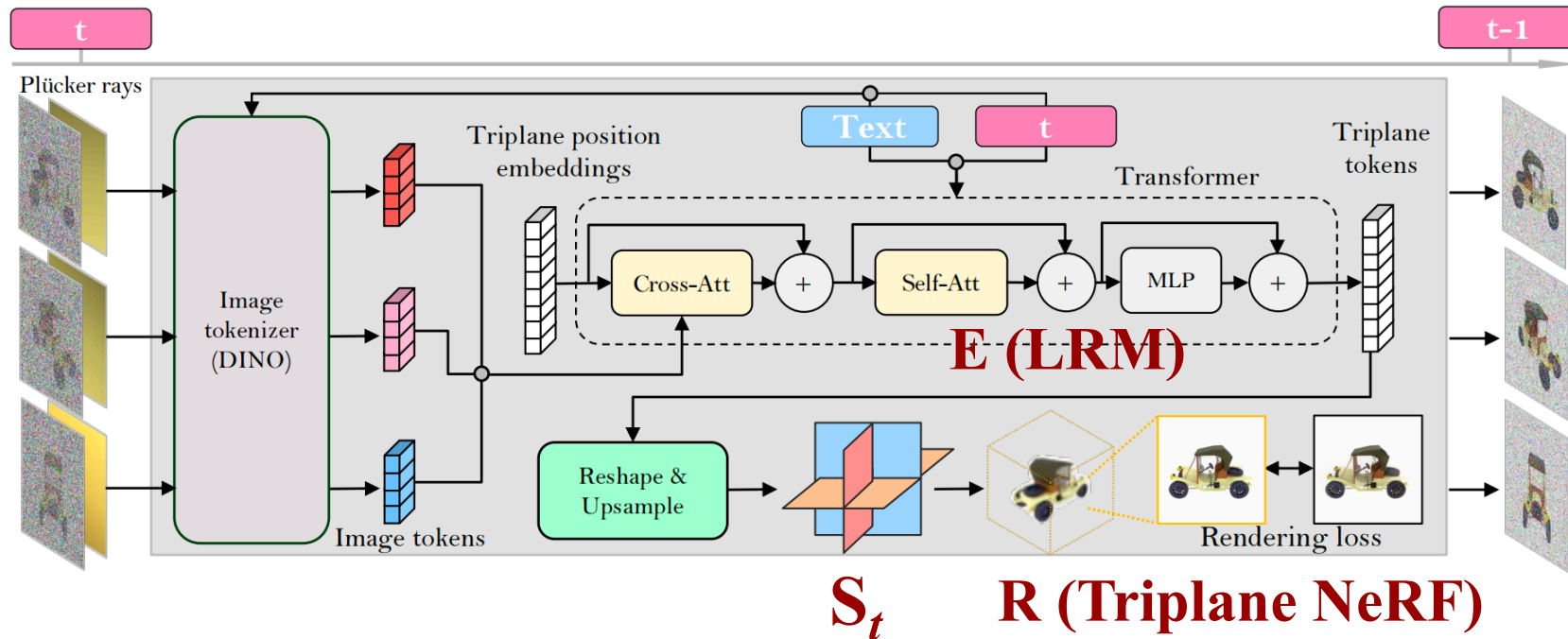
E: 3D reconstruction module

$\mathbf{S}_t$ : 3D representations

R: Rendering module



# Denoising Multi-view Diffusion: Method

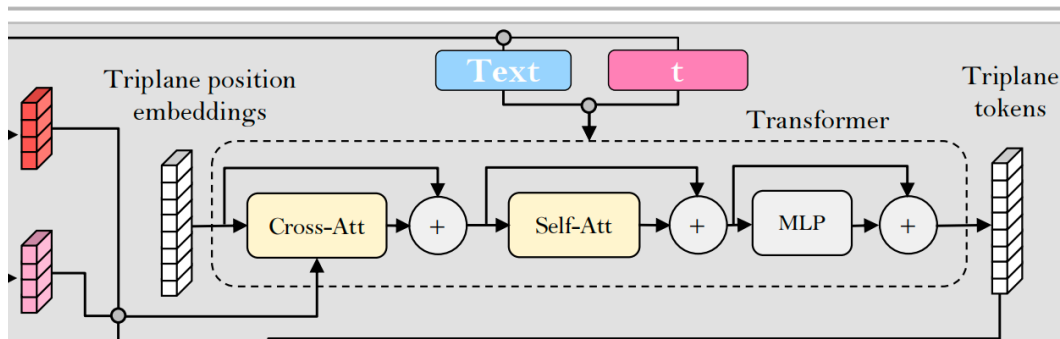


- Time conditioning: AdaLN block
- Camera conditioning: concatenate Plucker Coordinates with input pixels

$$r = (o \times d, d)$$

## Conditional generation

- Single-image condition
  - Keep the condition image noise-free
  - Align the Triplane coordinates with the condition view's coordinates
  - Normalize input camera view as LRM does
- Text condition
  - Use the CLIP text encoder to obtain text embeddings.
  - Inject them into the denoiser using cross-attention.



## Conditional generation

- Single-image condition
  - Keep the condition image noise-free
  - Align the Triplane coordinates with the condition view's coordinates
  - Normalize input camera view as LRM does
- Text condition
  - Use the CLIP text encoder to obtain text embeddings.
  - Inject them into the denoiser using cross-attention.
- Training objective:

$$\mathbf{L} = \mathbb{E}_{t \sim U[1, T], (\mathbf{I}, \mathbf{c}) \sim (\mathcal{I}_{full}, \mathcal{C}_{full})} \ell(\mathbf{I}, \mathbf{R}(\mathbf{E}(\mathcal{I}_t, t, \mathcal{D}, y), \mathbf{c}))$$

## Datasets

- Objaverse (~730k object meshes)
- MVImgNet (~220k object-centric videos)
- Cap3D (~660k image & caption pairs)

GPU: 128 A100, 7 days



# Denoising Multi-view Diffusion : Experiments

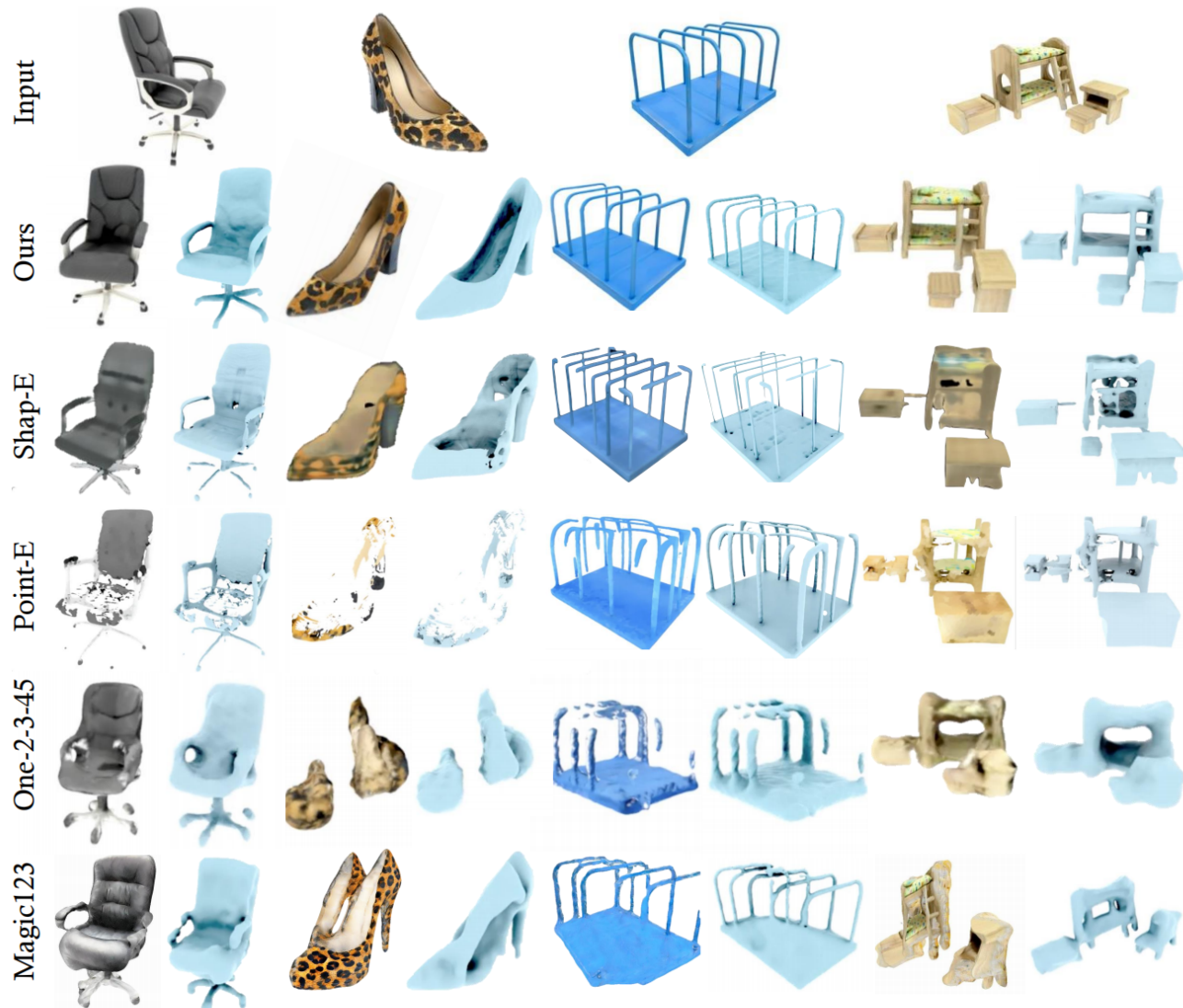


Figure 4: Qualitative comparisons on single-image reconstruction.



# Denoising Multi-view Diffusion : Experiments

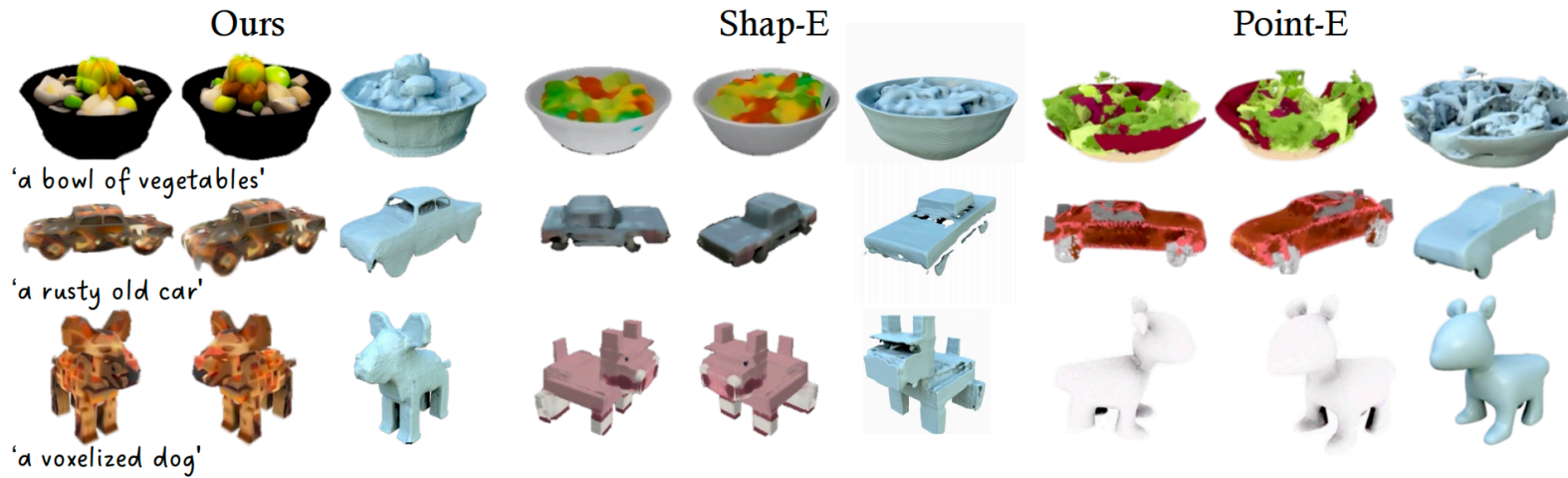


Figure 5: Qualitative comparisons on Text-to-3D .

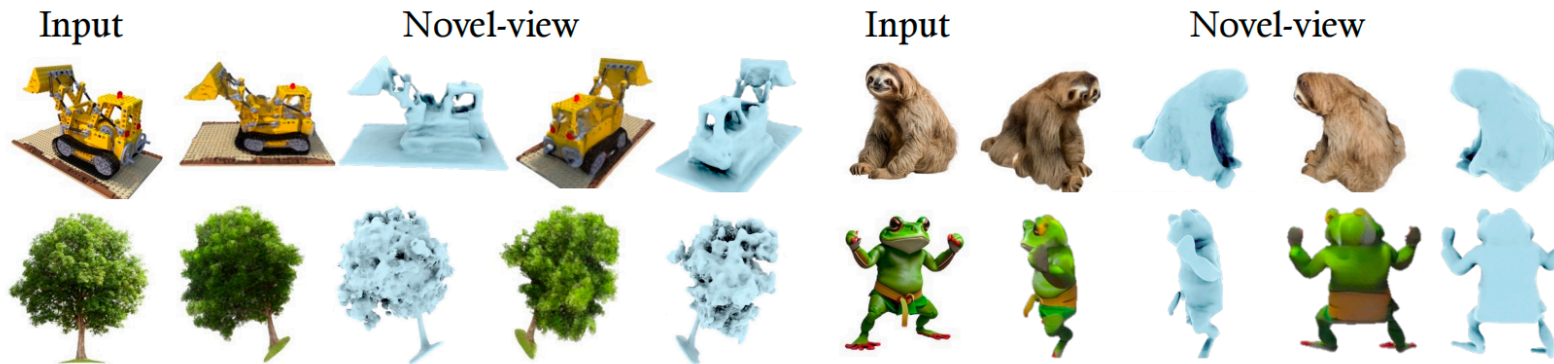


Figure 6: **Robustness to out-of-domain inputs:** synthetic (top left), real (bottom left, top right), and generated images (bottom right).



Thanks for listening!

Presenter: Rundong Luo  
2023.12.10