

iNVS: Repurposing Diffusion Inpainters for Novel View Synthesis

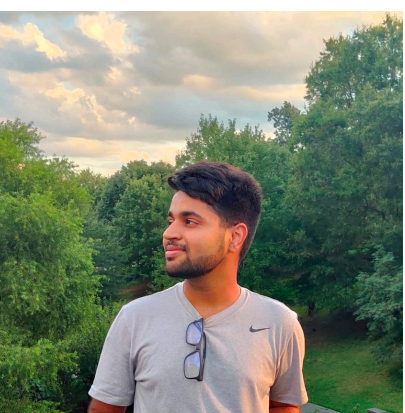
Yash Kant¹, Aliaksandr Siarohin², Michael Vasilkovsky², Riza Alp Guler², Jian Ren², Sergey Tulyakov², Igor Gilitschenski¹



University of Toronto¹



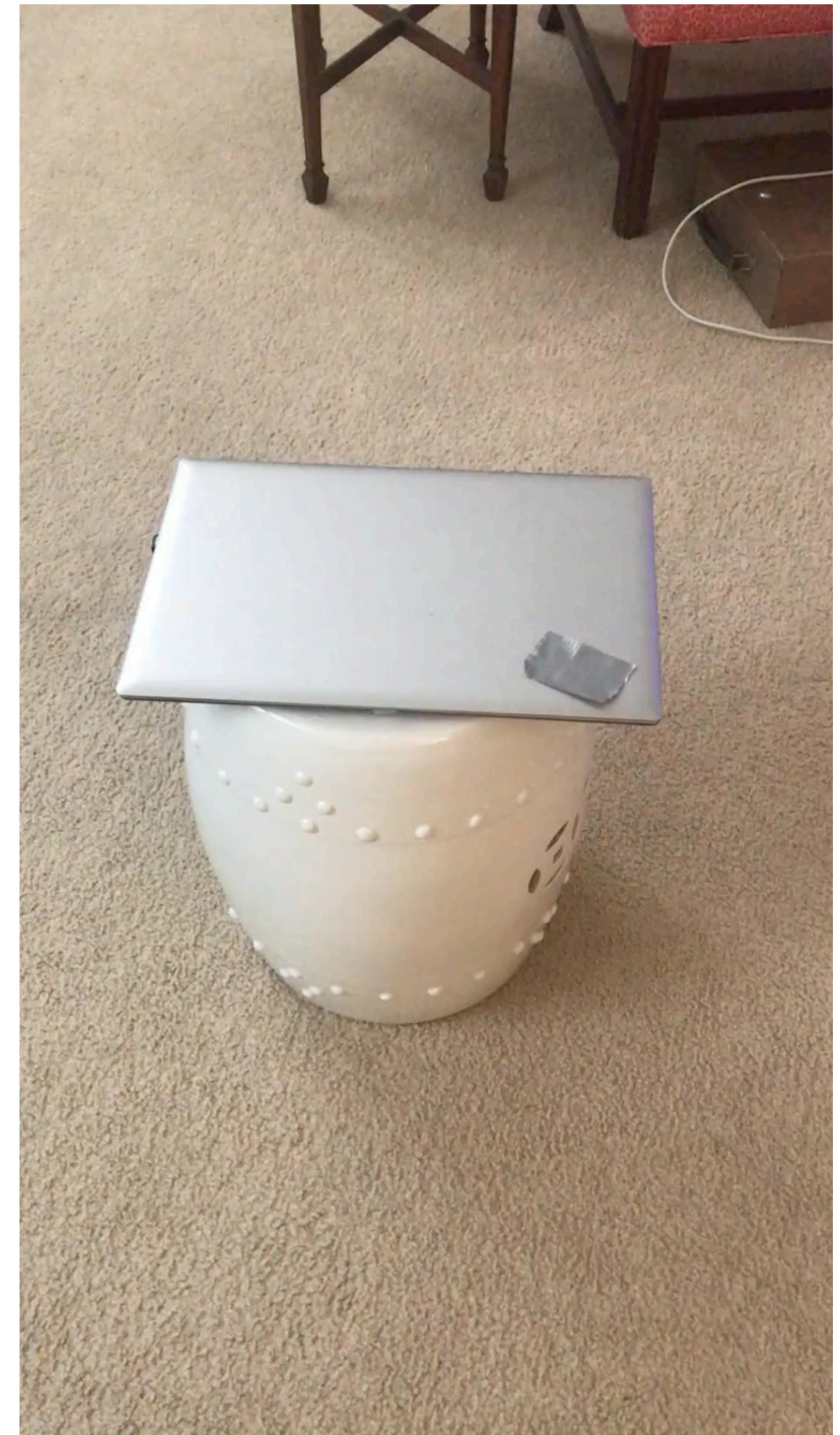
Snap Research²



Task: Given single image of an object, we want to generate it from novel viewpoints.

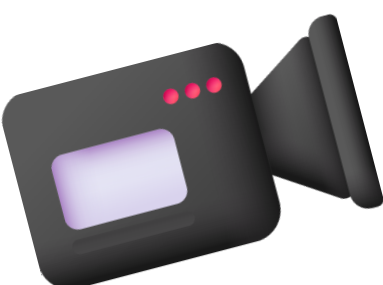


Given Image

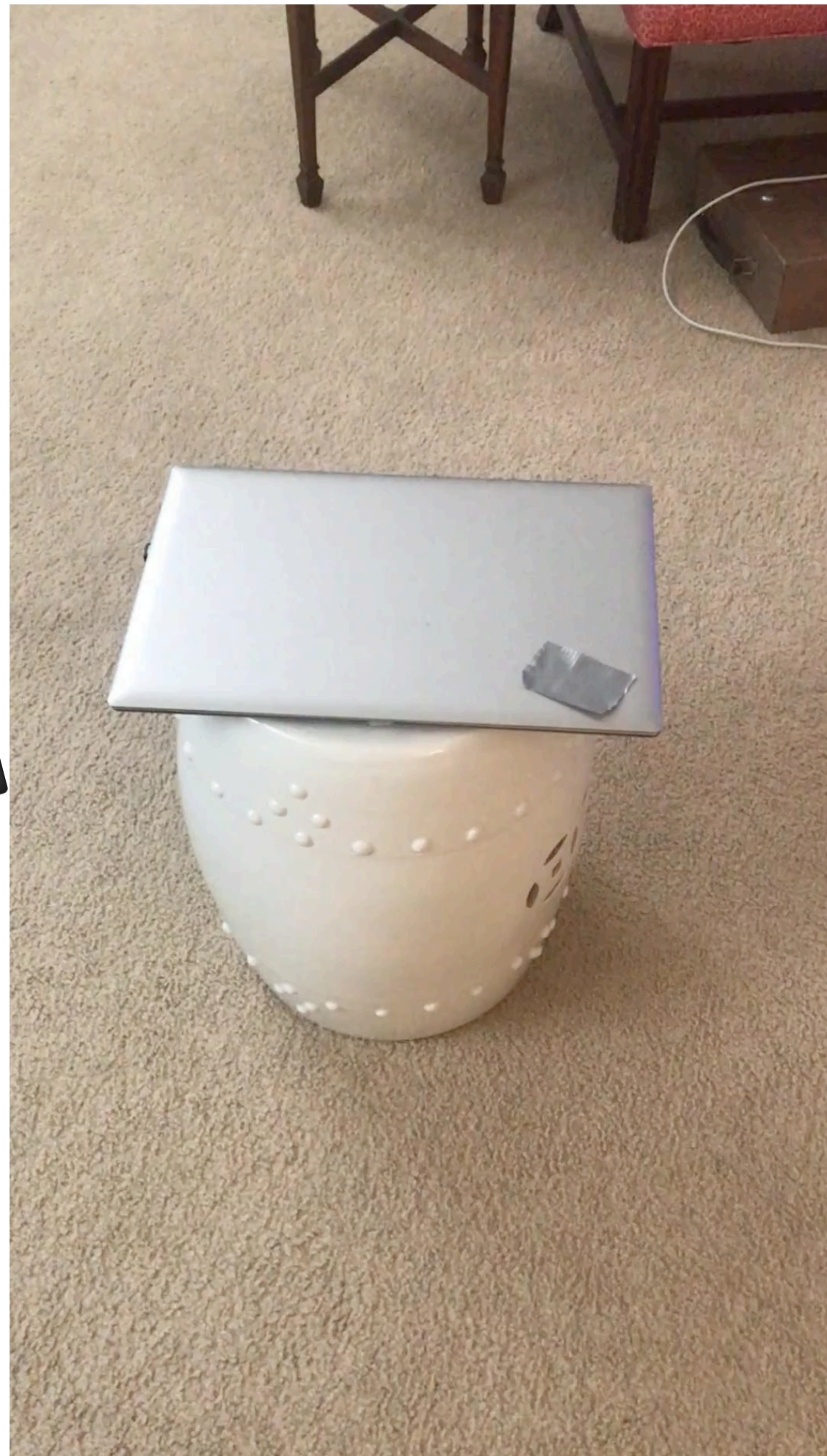


Want Novel Views!

Task: Given single image of an object, we want to generate it from novel viewpoints.



Source
Camera



Source View



Target
Camera-1



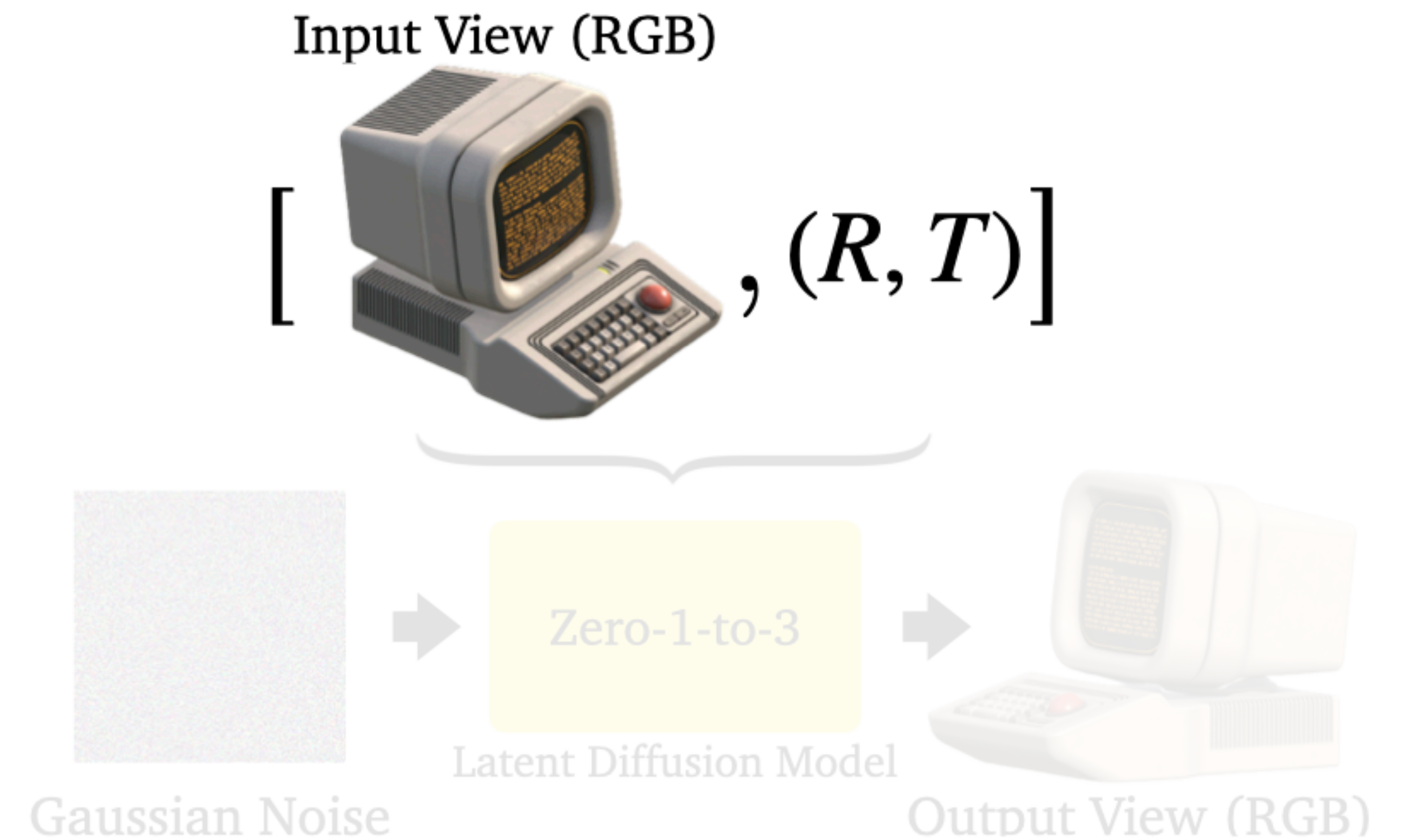
Target
Camera-2



Target Views

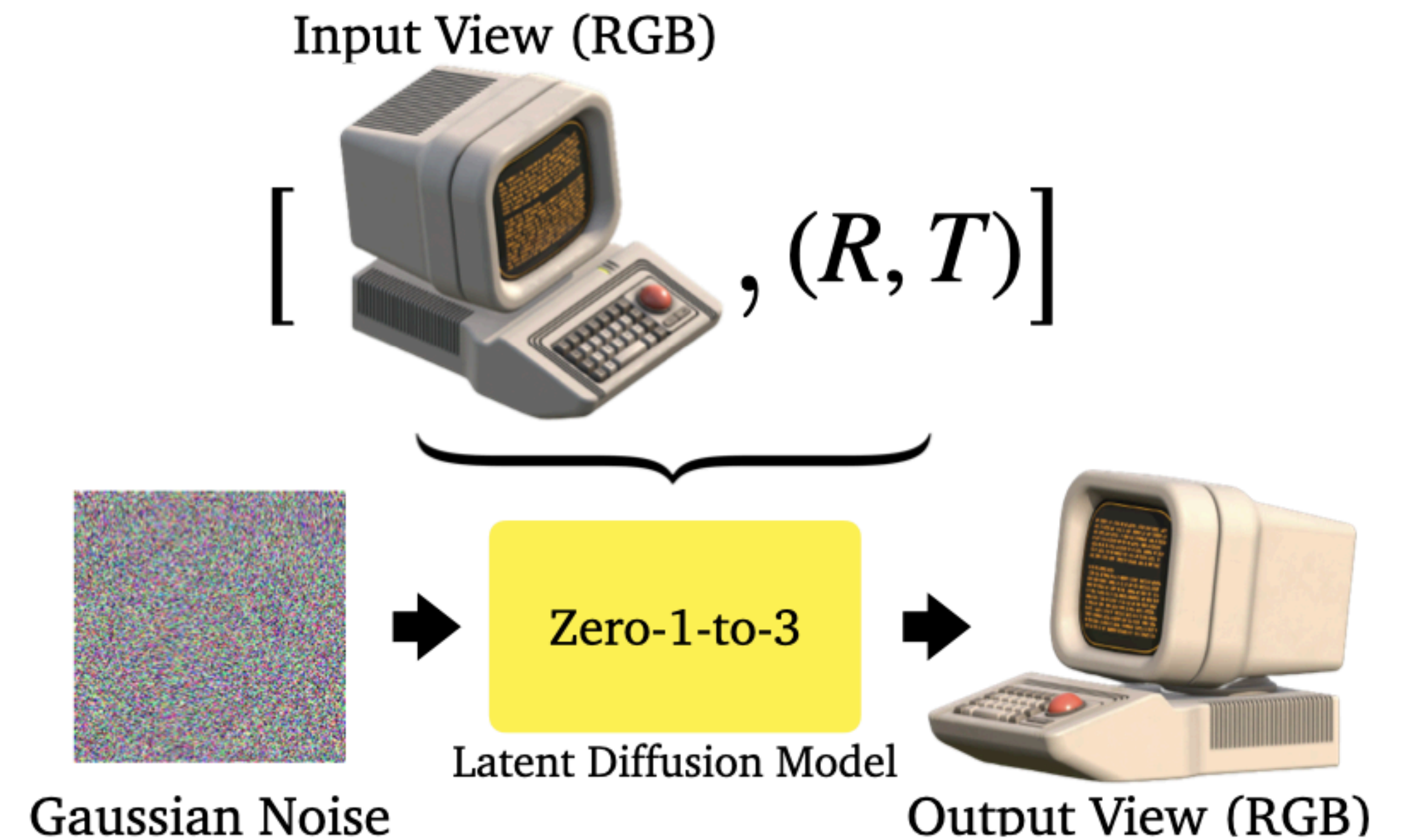
Prior Works: Zero-1-to-3 (ICCV, 2023)

- Encodes the source view using CLIP.
- Relative camera pose is encoded using dense MLP.



Prior Works: Zero-1-to-3 (ICCV, 2023)

- Encodes the source view using CLIP.
- Relative camera pose is encoded using dense MLP.
- Trained a conditional Stable Diffusion to denoise novel views on Objaverse.



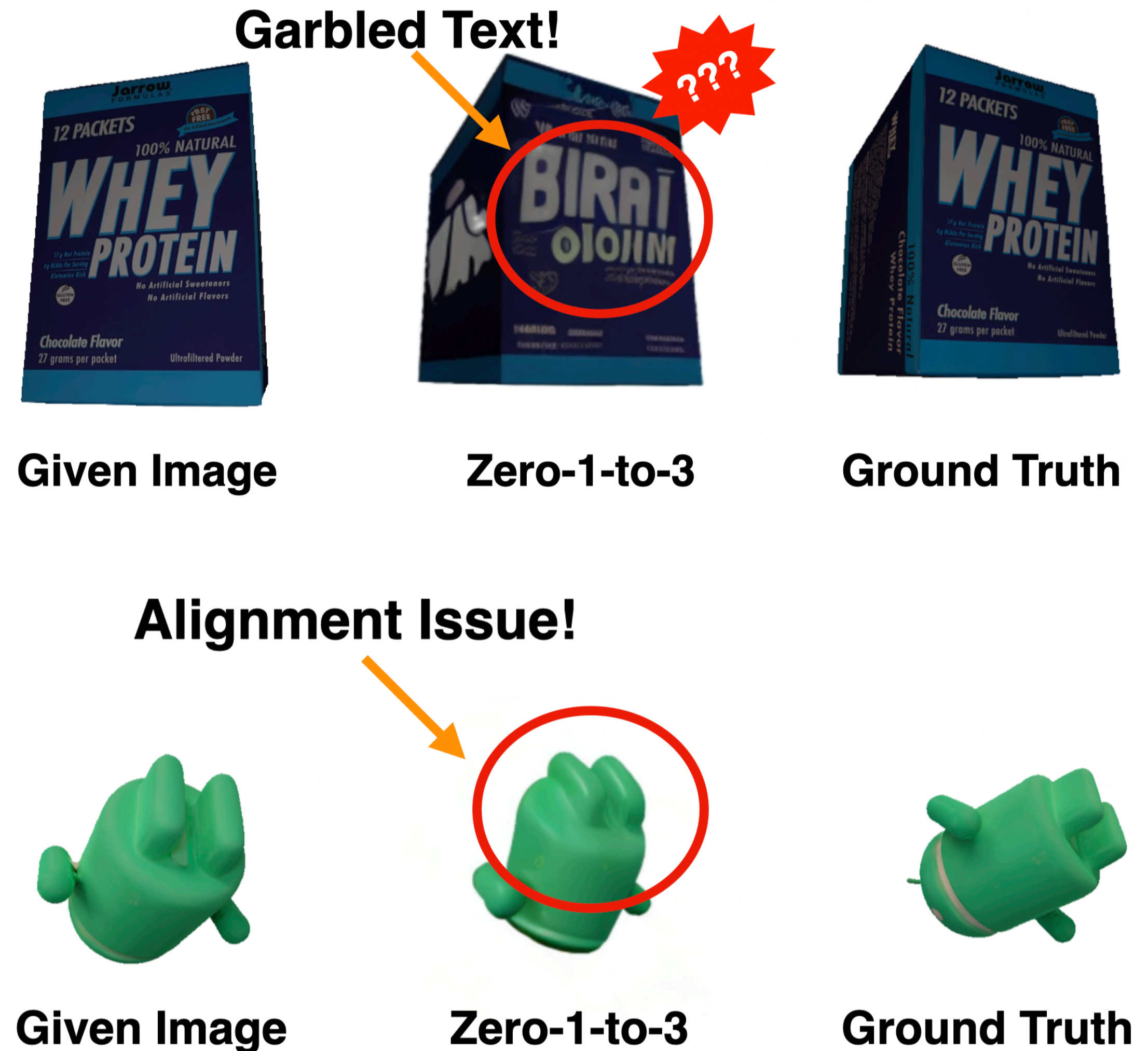
Prior Works: Zero-1-to-3 (ICCV, 2023)

- Inefficient reuse of input pixels — sharp details (text and texture) get garbled.



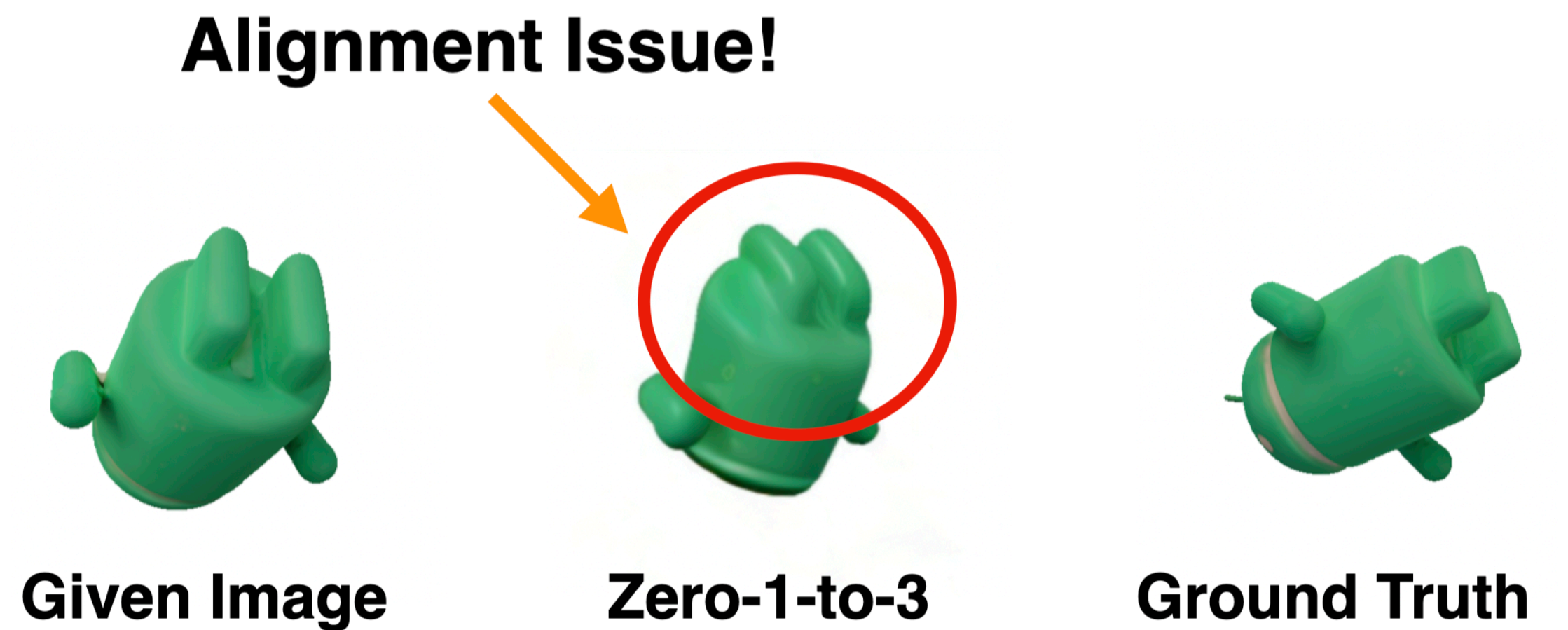
Prior Works: Zero-1-to-3 (ICCV, 2023)

- Inefficient reuse of input pixels — sharp details (text and texture) get garbled.
- Camera is encoded in a dense vector — coarse control.



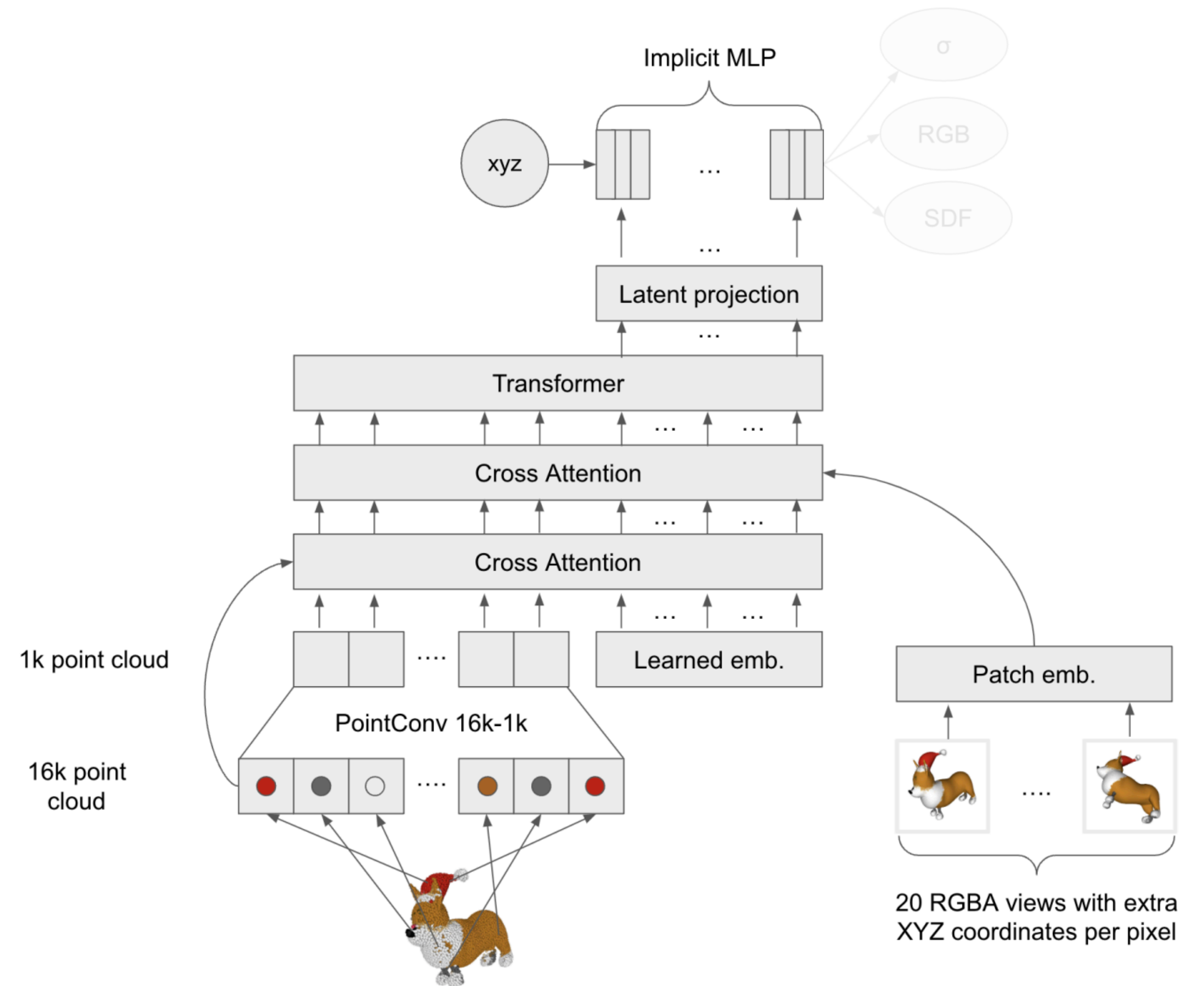
Prior Works: Zero-1-to-3 (ICCV, 2023)

- Inefficient reuse of input pixels — sharp details (text and texture) get garbled.
- Camera is encoded in a dense vector — coarse control.
- 3D consistency is not guaranteed.



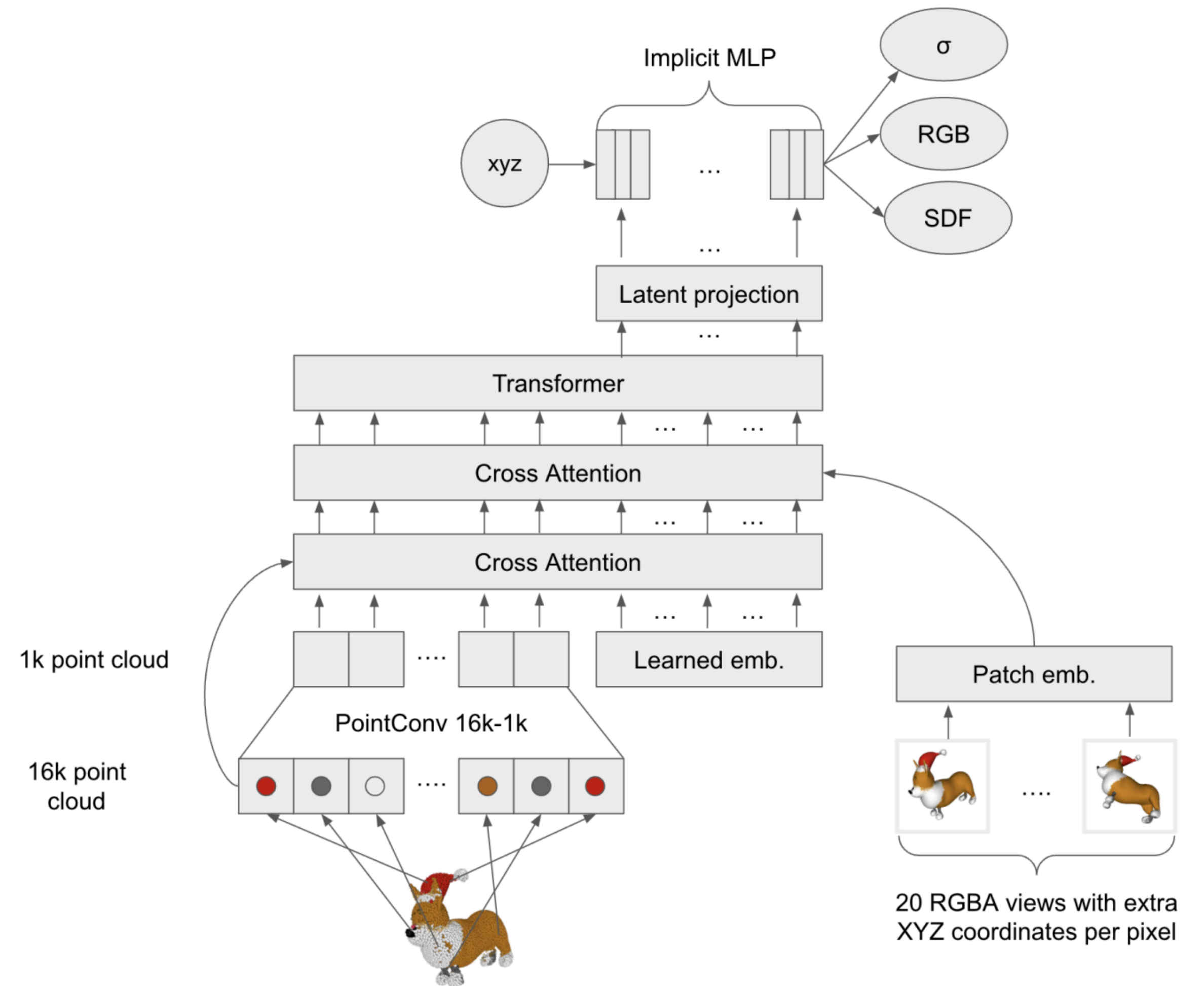
Prior Works: Shap-E (OpenAI, 2023)

- Hypernetwork-based approach which generates weights of an MLP conditioned on input image.



Prior Works: Shap-E (OpenAI, 2023)

- Hypernetwork-based approach which generates weights of an MLP conditioned on input image.
- This MLP is then queried using points in 3D space to generate occupancy and colour (NeRF).



Prior Works: Shap-E (OpenAI, 2023)

- Instability during training; since the space of possible solutions (weights) is very high dimensional.

Prior Works: Shap-E (OpenAI, 2023)

- Instability during training; since the space of possible solutions (weights) is very high dimensional.
- Inefficient reuse of input pixels — sharp details (text and texture) get garbled.

Given Image



Generated 3D



Prior Works: Shap-E (OpenAI, 2023)

- Instability during training; since the space of possible solutions (weights) is very high dimensional.
- Inefficient reuse of input pixels — sharp details (text and texture) get garbled.
- Suffers from janus artefacts.

Given Image



Generated 3D



***iNVS*: Reframing NVS as completion task**

Inspired by video generation approaches such as InfiniteNature, we ask —

“can we reframe novel view synthesis as a image completion task rather than generation from scratch?”



InfiniteNature [Andrew Liu, et al. ICCV, 2021]

***iNVS*: Reframing NVS as completion task**

Inspired by video generation approaches such as Infinite Nature —

“can we reframe novel view synthesis as a image completion task rather than generation from scratch?”



InfiniteNature [Andrew Liu, et al. ICCV, 2021]

***iNVS*: Creating Partial Views**

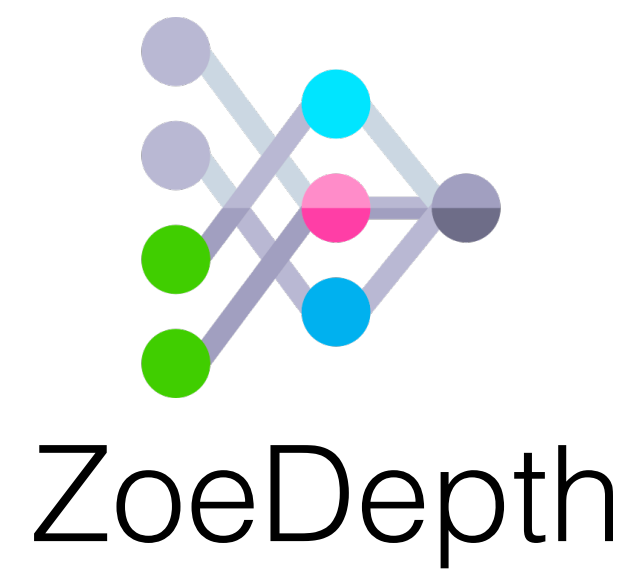
We reuse source image pixels to create a partial image.



Source View

iNVS: Creating Partial Views

We use monocular depth [ZoeDepth] to unproject source pixels in 3D.



Source View

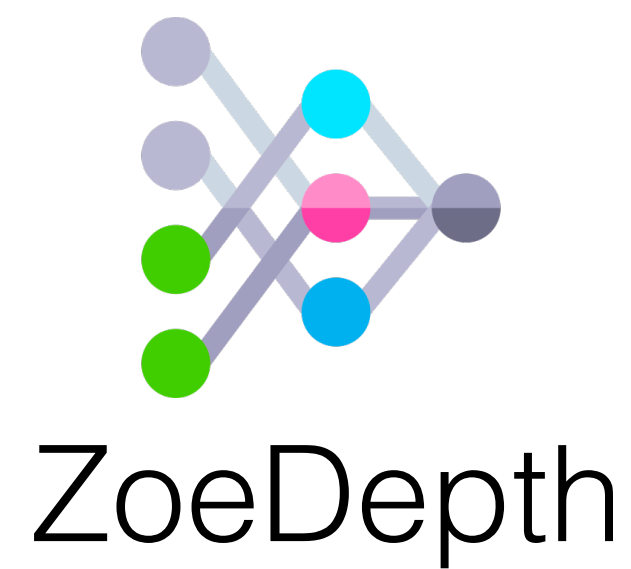
Depth Map

iNVS: Creating Partial Views

We re-project these 3D points back to target view using softmax-splatting and create partial target view.



Source View



Depth Map

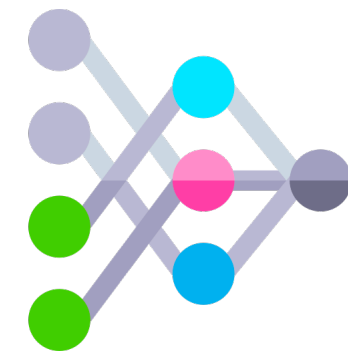


Partial Target View

***iNVS*: Inpainting Partial Views for NVS**

iNVS: Inpainting Partial Views for NVS

- Stable Diffusion
Inpainter fills in newly
discovered regions.



SD
Inpainter

iNVS: Inpainting Partial Views for NVS

- Stable Diffusion Inpainter fills in newly discovered regions.
- We train the Inpainter on Objaverse dataset, to learn 3D completion priors.



Partial Target View



Inpainted View

iNVS: Epipolar and Pose-aware Inpainting Mask

We further constrain the inpainting to the areas occluded in source view using epipolar lines.

Source View

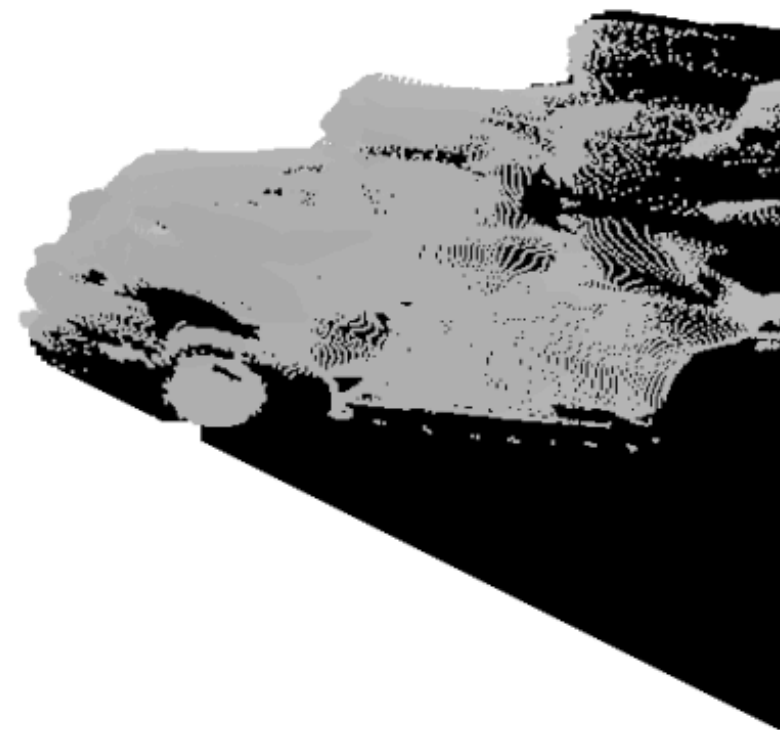


Partial View



Occluded region

Epipolar Mask



iNVS



iNVS: Epipolar and Pose-aware Inpainting Mask

We further constrain the inpainting to the areas occluded in source view using epipolar lines.

Source View

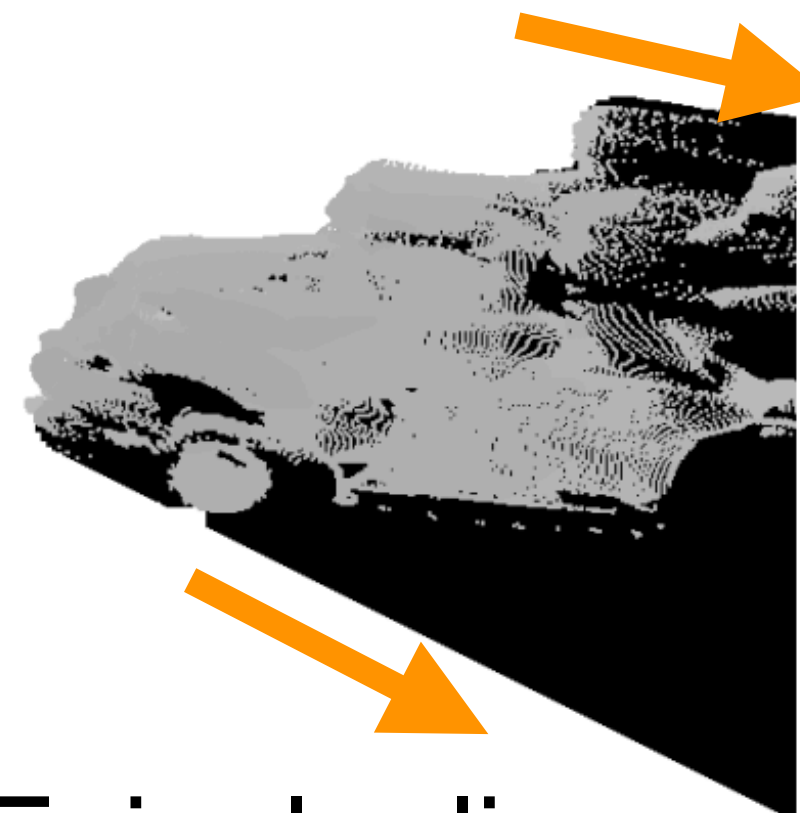


Partial View



Occluded region

Epipolar Mask



Epipolar lines

iNVS



iNVS: Epipolar and Pose-aware Inpainting Mask

Additionally, we also use a soft-valued (0,1) inpainting mask that conveys the relative angle (0,180) between source and target camera ray.

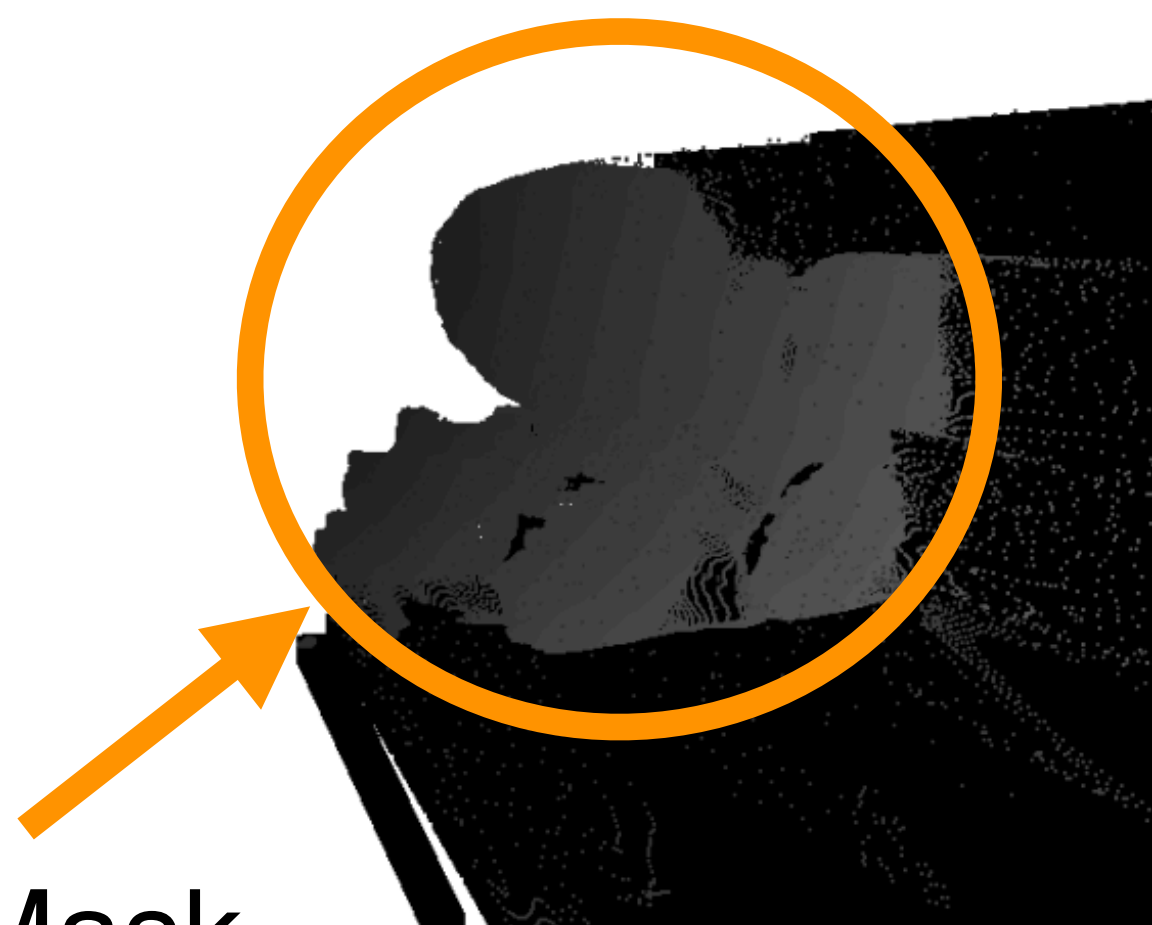
Source View



Partial View



Epipolar Mask



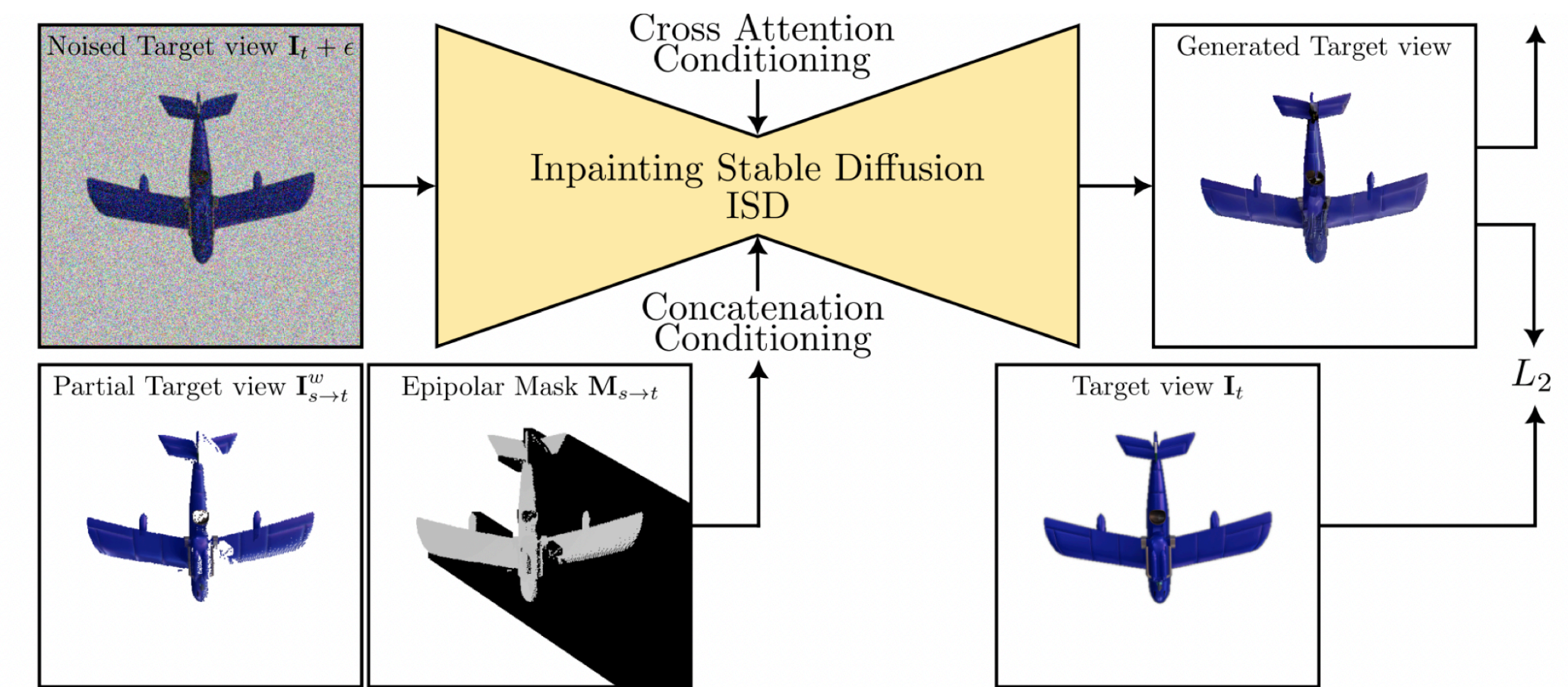
iNVS



Soft Inpainting Mask

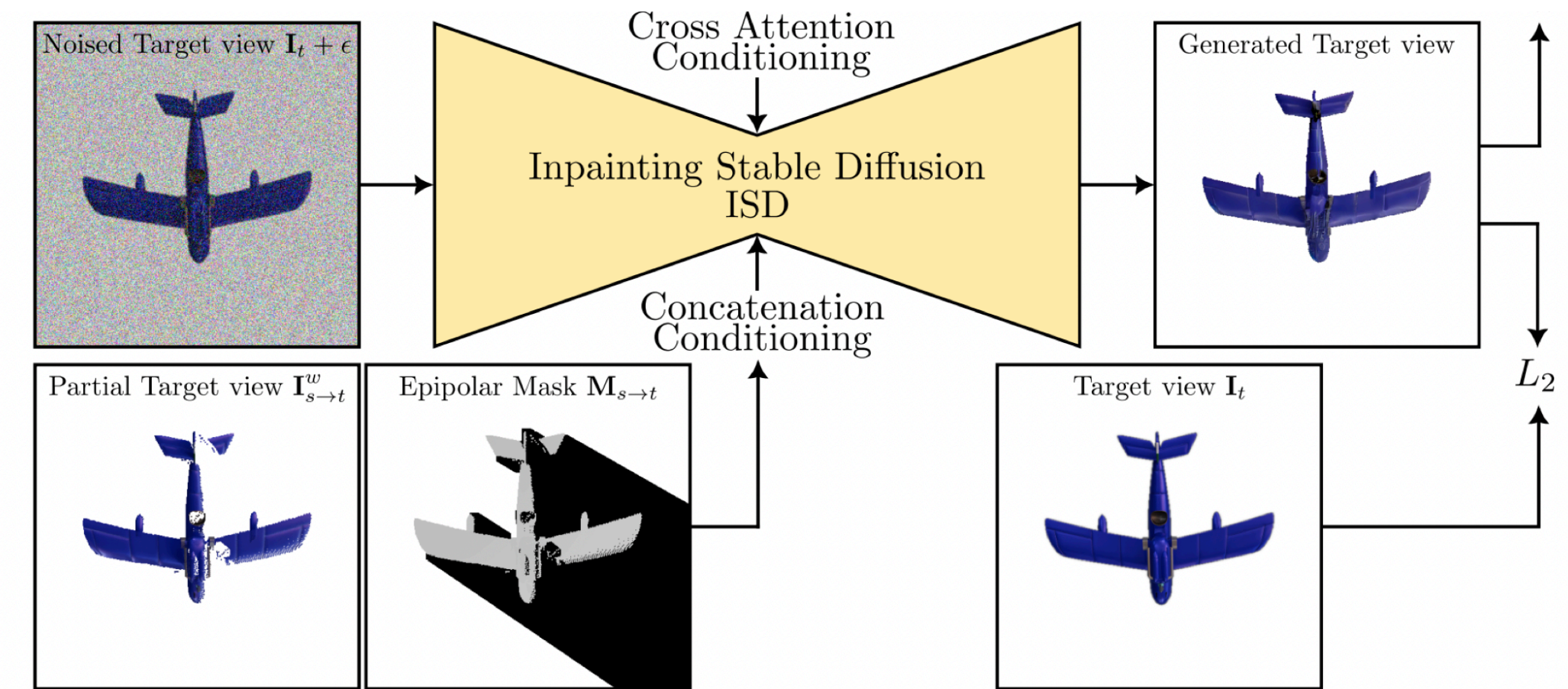
iNVS: Training Details

- We fine-tune SD Inpainter network on 96 A100 GPUs for two weeks, on $\sim 20\text{M}$ rendered images from Objaverse.



iNVS: Training Details

- We fine-tune SD Inpainter network on 96 A100 GPUs for two weeks, on $\sim 20\text{M}$ rendered images from Objaverse.
- Object boundary appears as early as 10% of denoising steps; hence, we sample timesteps with a bias during training.



Early Inference Steps



10%

50%

iNVS: Qualitative Comparison



Given Image

iNVS: Qualitative Comparison



Given Image



**Baseline
(Zero-1-to-3)**

iNVS: Qualitative Comparison

Garbled Text!



Given Image



Baseline
(Zero-1-to-3)

iNVS: Qualitative Comparison

Garbled Text!



Given Image



Baseline
(Zero-1-to-3)



iNVS
(Ours)

iNVS: Qualitative Comparison

Garbled Text!

Text remains
intact!



Given Image



Baseline
(Zero-1-to-3)



iNVS
(Ours)

***iNVS*: Qualitative Comparison**



Given Image

***iNVS*: Qualitative Comparison**



Given Image



**Baseline
(Zero-1-to-3)**

iNVS: Qualitative Comparison

Toaster with Grills



Given Image



**Baseline
(Zero-1-to-3)**

iNVS: Qualitative Comparison

Toaster with Grills



Given Image



**Baseline
(Zero-1-to-3)**



***iNVS*
(Ours)**

iNVS: Qualitative Comparison

Toaster with Grills



Given Image



**Baseline
(Zero-1-to-3)**

**Toaster slots
preserved!**



***iNVS*
(Ours)**

iNVS: Qualitative Comparison

Resolution:
256x256



Given Image



**Baseline
(Zero-1-to-3)**

Resolution:
512x512



***iNVS*
(Ours)**

iNVS: Quantitative Comparison

iNVS outperforms Zero-1-to-3 on 2/3 metrics on GSO (synthetic)

Method	PSNR ↑	SSIM ↑	LPIPS ↓
<i>iNVS</i>	18.95	0.30	0.24
Zero-1-to-3	14.74	0.34	0.25

Google Scanned Objects

iNVS: Quantitative Comparison

iNVS outperforms Zero-1-to-3 on 2/3 metrics on GSO (synthetic) and CO3D (real-world) datasets.

Method	PSNR ↑	SSIM ↑	LPIPS ↓
<i>iNVS</i>	18.95	0.30	0.24
Zero-1-to-3	14.74	0.34	0.25

Google Scanned Objects

Method	PSNR ↑	SSIM ↑	LPIPS ↓
<i>iNVS</i>	17.58	0.33	0.36
Zero-1-to-3	12.32	0.33	0.42

Common Objects in 3D

Failure Mode

Investigating the lower Structural Similarity (SSIM) score, we find some common failure modes.

Failure Mode

Investigating the lower Structural Similarity (SSIM) score, we find some common failure modes.

iNVS struggles most when monocular depth estimator generates inaccurate depth.

Failure Mode #1: Deformed Partial View

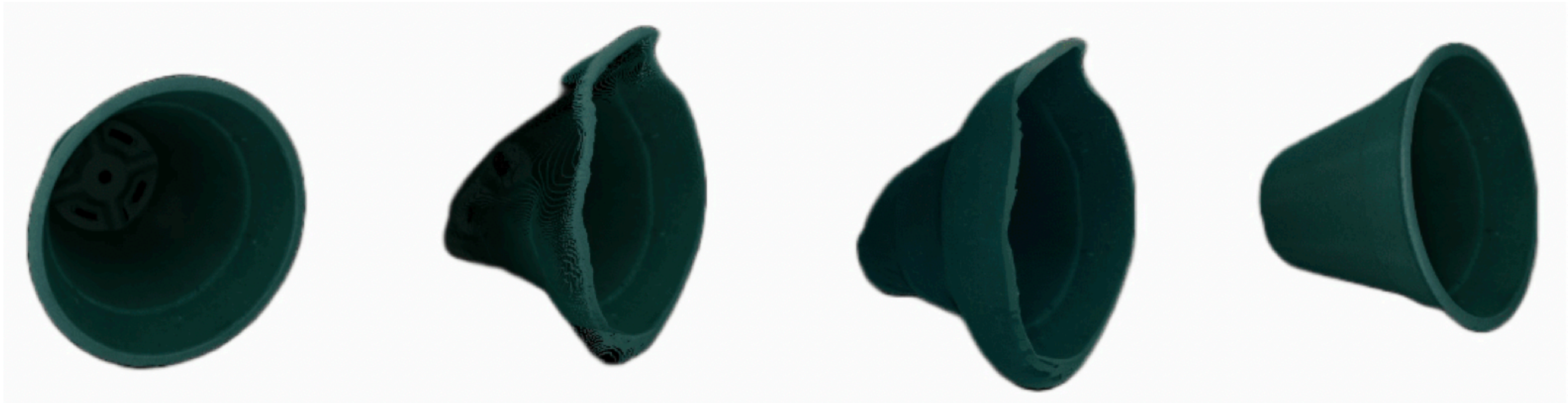
Imprecise depth leads to deformed partial view — difficult to recover from during inpainting.

Source View

**Unprojected
Partial View**

**Inpainted
view by *iNVS***

Ground Truth



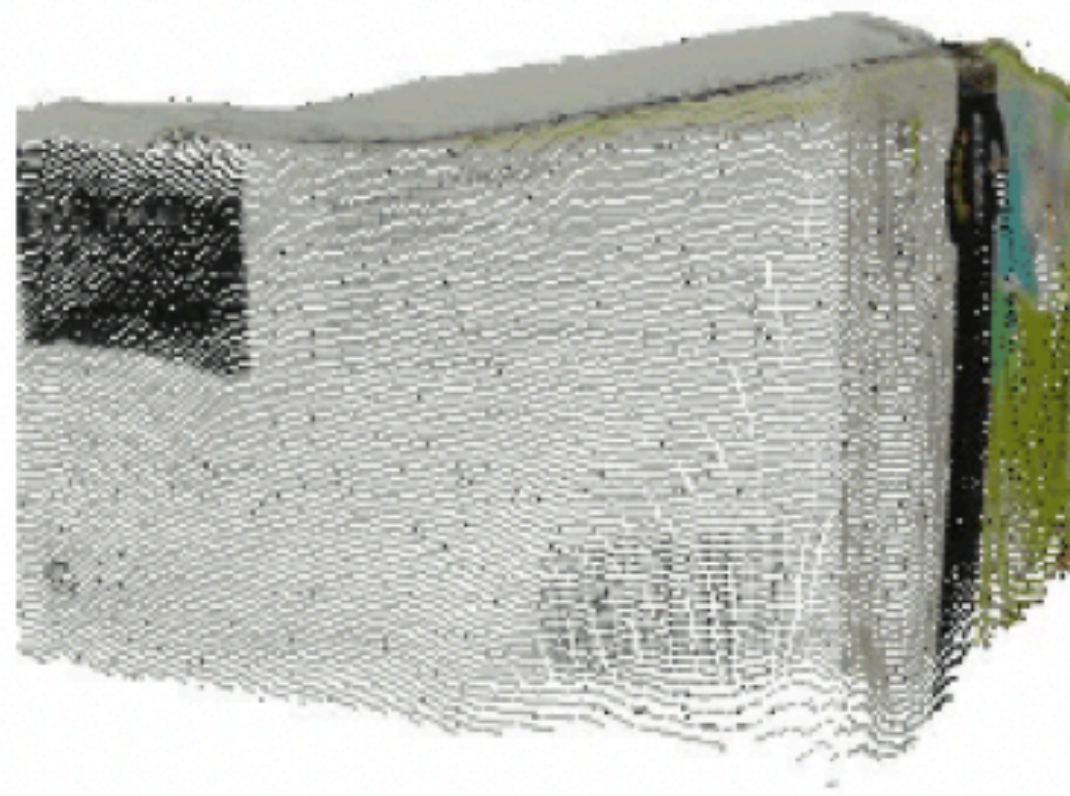
Failure Mode #2: Tiny holes can blend into texture.

Imprecise depth and downsampled inpainting mask occasionally cause reprojection holes to blend in.

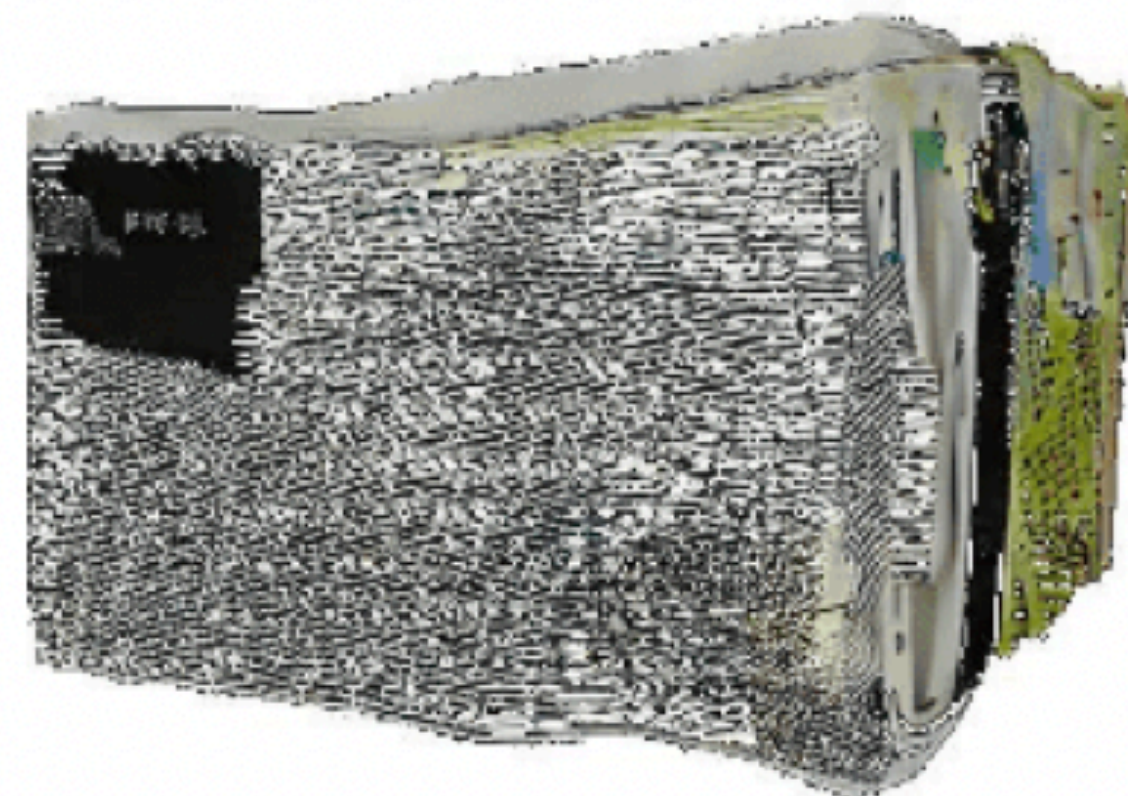
Source View



**Unprojected
Partial View**



**Inpainted
view by *iNVS***



Ground Truth



Failure Mode #3: Flipped pixels throw-off Inpainter.

Under large-viewpoint changes; we rely on inpainting mask to detect large ray angle changes, but it may fail.

Source View

**Unprojected
Partial View**

**Inpainted
view by *iNVS***

Ground Truth



Thanks for listening.

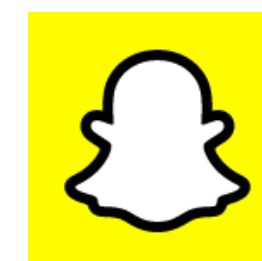
Poster Today @ 6PM!

Webpage: <https://yashkant.github.io/invs>

Yash Kant¹, Aliaksandr Siarohin², Michael Vasilkovsky², Riza Alp Guler², Jian Ren², Sergey Tulyakov², Igor Gilitschenski¹



University of Toronto¹



Snap Research²

