

2D Magic in a 3D World

Songyou Peng

Imperial College London
Mar 22, 2024

Who Am I?

- Senior Researcher
- Incoming Research Scientist

ETHzürich

Google Research

- Earned my PhD
 - Marc Pollefeys
 - Andreas Geiger

ETHzürich

MAX PLANCK INSTITUTE
FOR INTELLIGENT SYSTEMS



- Internships during PhD
 - 2021: Michael Zollhoefer
 - 2022: Tom Funkhouser

Meta

Google Research



pengsongyou.github.io

Research Overview of My PhD

Learn to Reconstruct and Understand 3D World


ConvOccNet
ECCV 2020 (Spotlight)


MonoSDF
NeurIPS 2022


Shape As Points
NeurIPS 2021 (Oral)

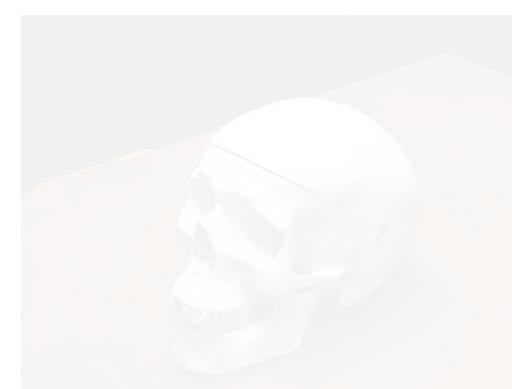

KiloNeRF
ICCV 2021
runs now at 50 fps on a GTX 1080 Ti



NICE-SLAM
CVPR 2022



NICER-SLAM
3DV 2024 (Oral)



UNISURF
ICCV 2021 (Oral)



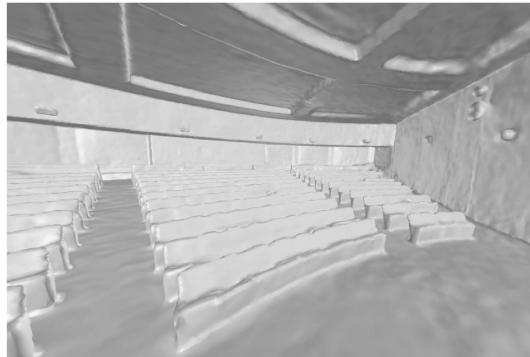
OpenScene
CVPR 2023

Research Overview of My PhD

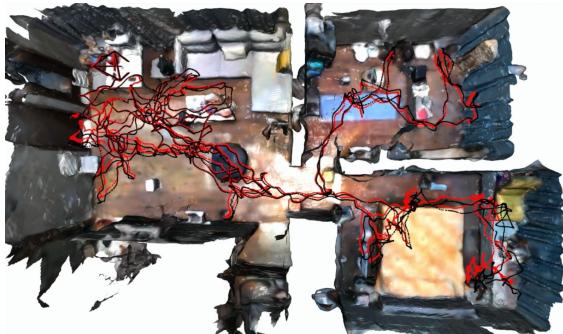
Learn to Reconstruct and Understand 3D World



ConvOccNet
ECCV 2020 (Spotlight)



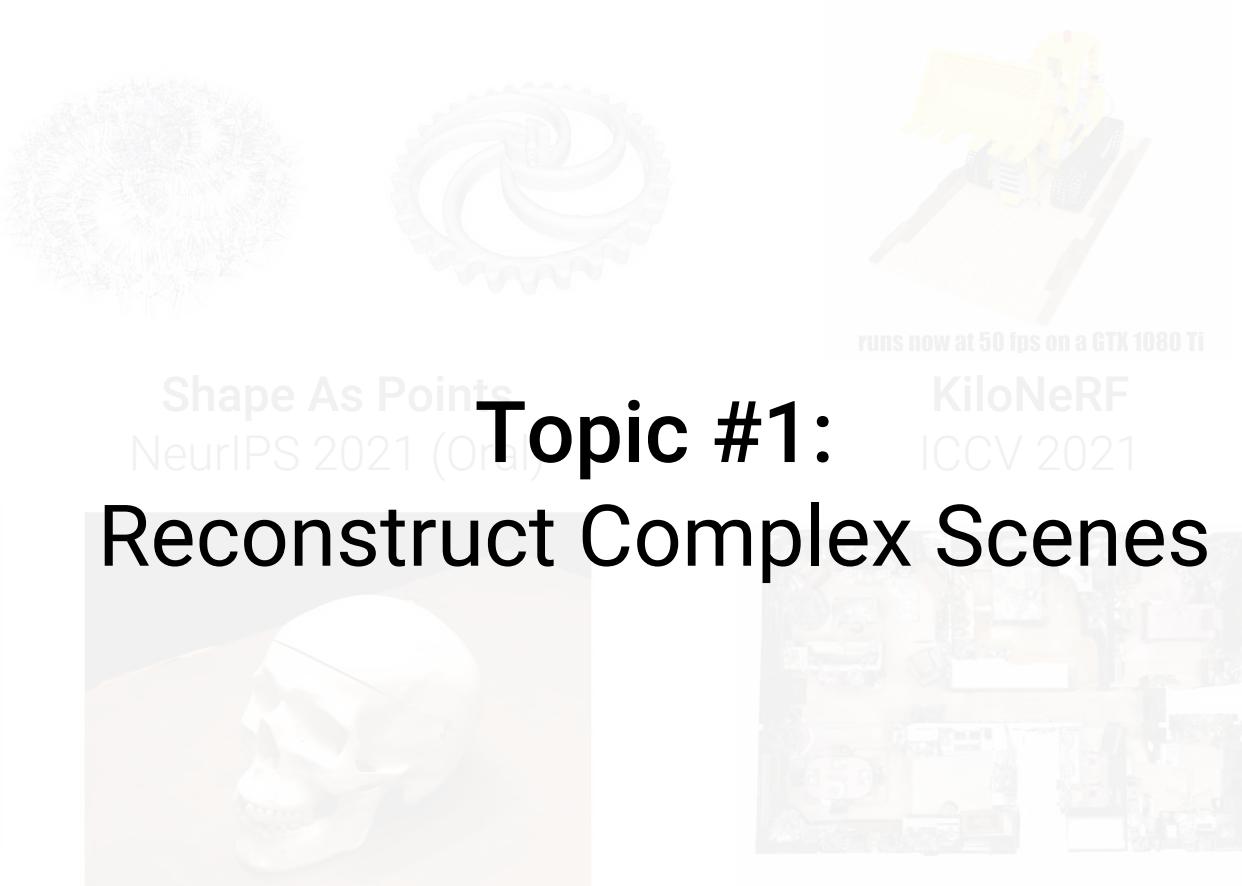
MonoSDF
NeurIPS 2022



NICE-SLAM
CVPR 2022



NICER-SLAM
3DV 2024 (Oral)



Topic #1:
Reconstruct Complex Scenes

UNISURF
ICCV 2021 (Oral)

OpenScene
CVPR 2023

Research Overview of My PhD

Learn to Reconstruct and Understand 3D World

Topic #2: Fast Inference

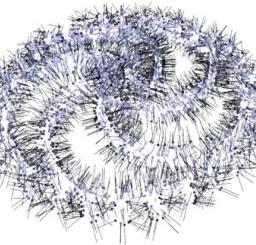
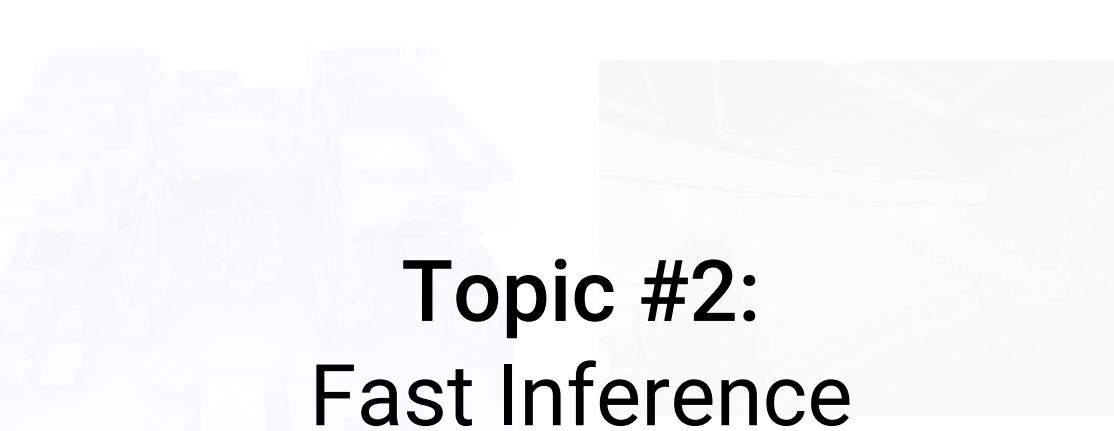
ConvOccNet
ECCV 2020 (Spotlight)



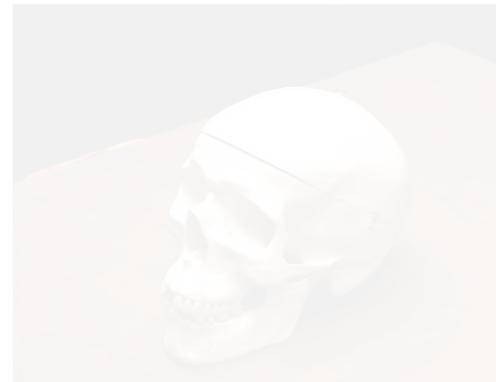
NICE-SLAM
CVPR 2022



NICER-SLAM
3DV 2024 (Oral)



Shape As Points
NeurIPS 2021 (Oral)



UNISURF
ICCV 2021 (Oral)



runs now at 50 fps on a GTX 1080 Ti

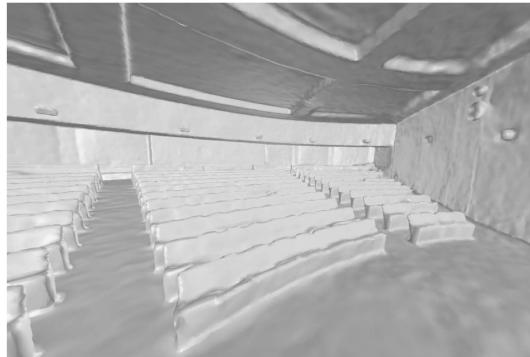
KiloNeRF
ICCV 2021



OpenScene
CVPR 2023

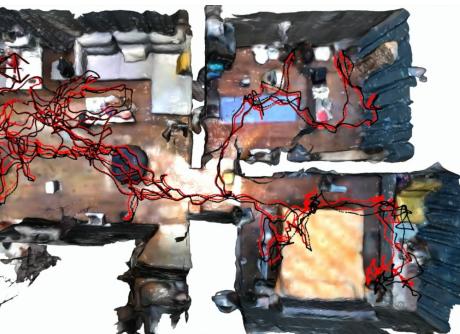
Research Overview of My PhD

Learn to Reconstruct and Understand 3D World

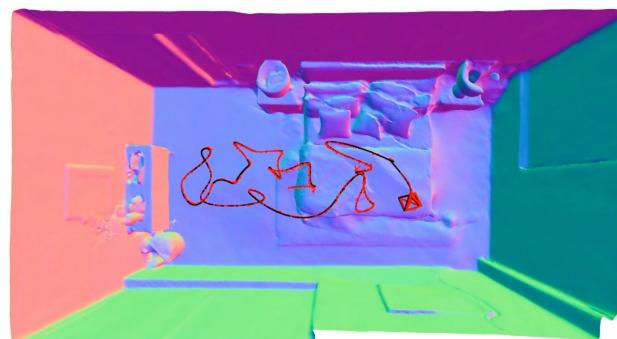


ConvOccNet
ECCV 2020 (Spotlight)

MonoSDF
NeurIPS 2022



NICE-SLAM
CVPR 2022



NICER-SLAM
3DV 2024 (Oral)



UNISURF
ICCV 2021 (Oral)

runs now at 50 ips on a GTX 1080 Ti

KiloNeRF
ICCV 2021



OpenScene
CVPR 2023 5

Research Overview of My PhD

Learn to Reconstruct and Understand 3D World



ConvOccNet
ECCV 2020 (Spotlight)



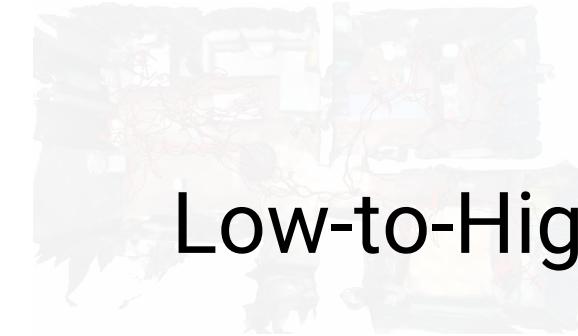
MonoSDF
NeurIPS 2022



Shape As Points
NeurIPS 2021 (Oral)



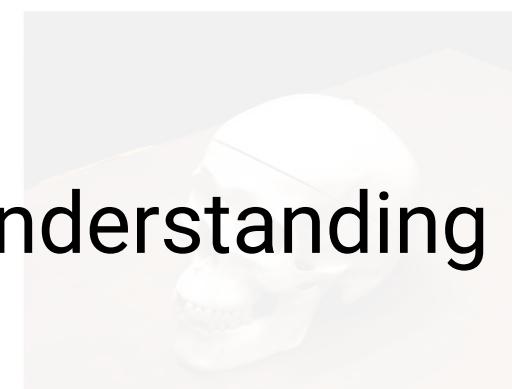
KiloNeRF
ICCV 2021
runs now at 50 fps on a GTX 1080 Ti



Topic #4:
Low-to-High Level 3D Scene Understanding



NICE-SLAM
CVPR 2022



UNISURF
ICCV 2021 (Oral)



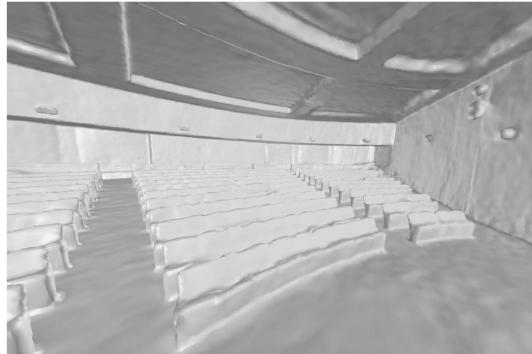
OpenScene
CVPR 2023 6

Research Overview of My PhD

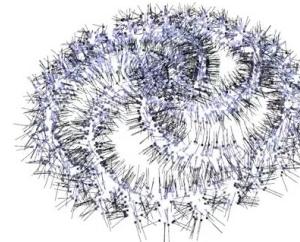
Learn to Reconstruct and Understand 3D World



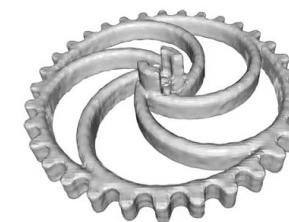
ConvOccNet
ECCV 2020 (Spotlight)



MonoSDF
NeurIPS 2022

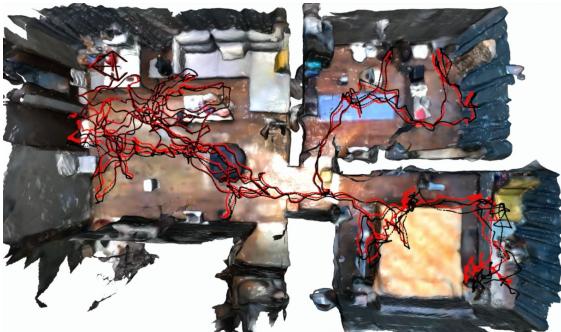


Shape As Points
NeurIPS 2021 (Oral)

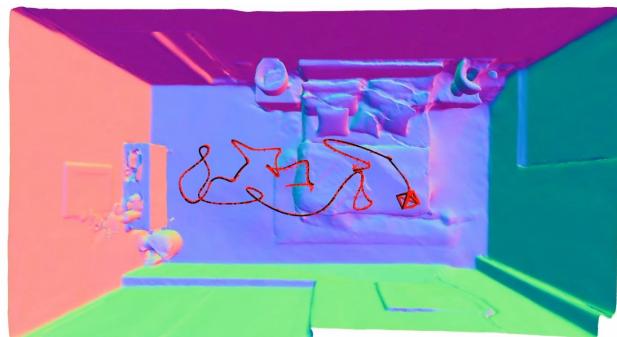


runs now at 50 fps on a GTX 1080 Ti

KiloNeRF
ICCV 2021



NICE-SLAM
CVPR 2022



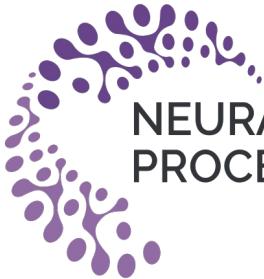
NICER-SLAM
3DV 2024 (Oral)



UNISURF
ICCV 2021 (Oral)



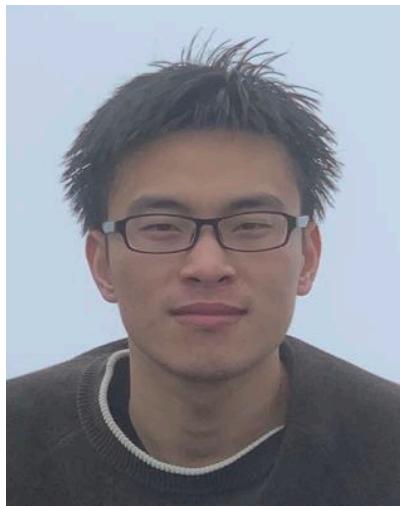
OpenScene
CVPR 2023



NEURAL INFORMATION
PROCESSING SYSTEMS



MonoSDF: Exploring Monocular Geometric Cues for Neural Implicit Surface Reconstruction



Zehao Yu



Songyou Peng



Michael Niemeyer



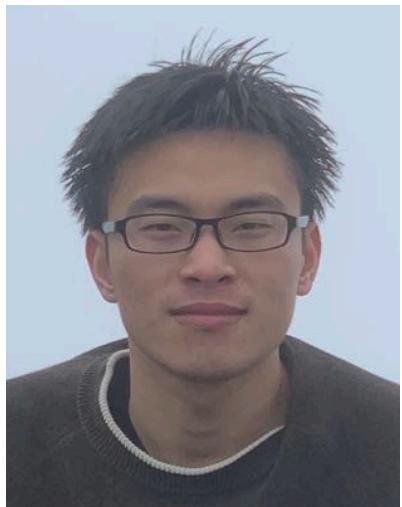
Torsten Sattler



Andreas Geiger



MonoSDF: Exploring **Monocular Geometric Cues** for Neural Implicit Surface Reconstruction



Zehao Yu



Songyou Peng



Michael Niemeyer

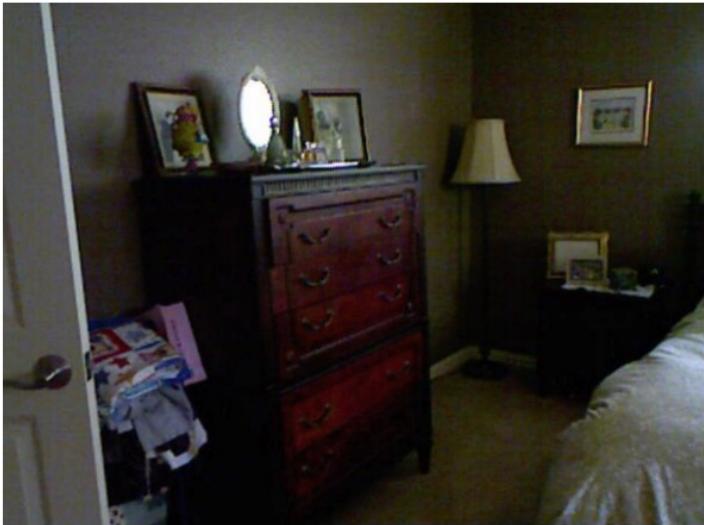


Torsten Sattler



Andreas Geiger

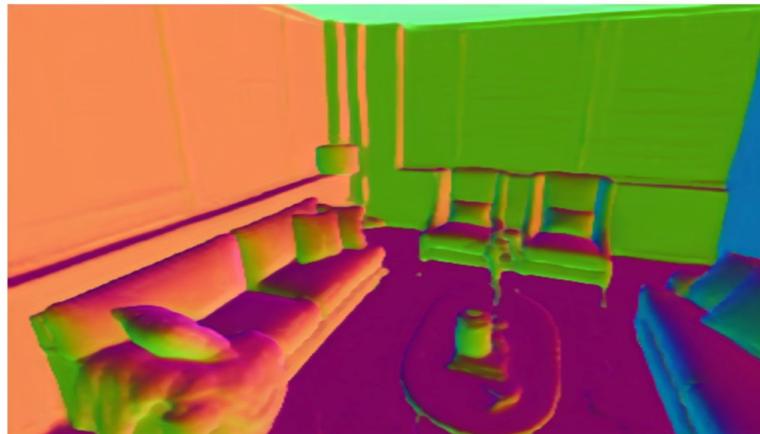
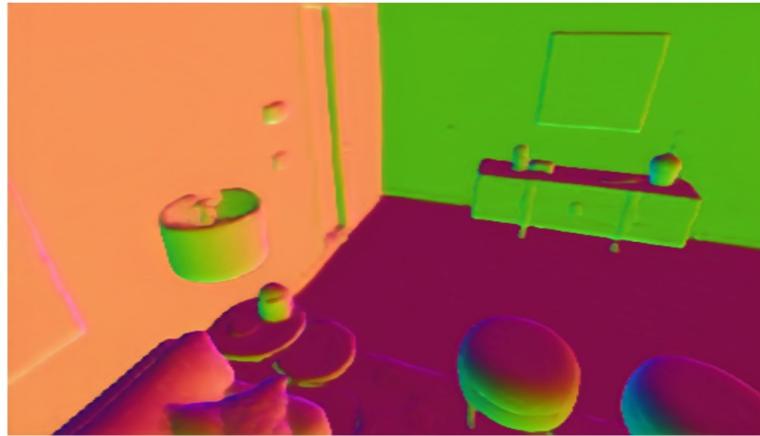
Depth Prediction from a Single Image



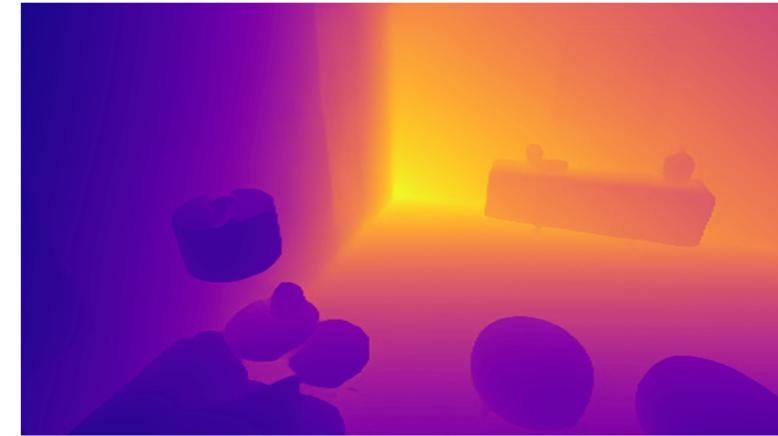
Omnidata



RGB Image



Omnidata Normal



Omnidata Depth

2D Magic in a 3D World

Songyou Peng

Imperial College London
Mar 22, 2024

2D Magic in a 3D World

2D Magic in a 3D World

2D Monocular Cues Benefit 3D Reconstruction

3D Reconstruction Pipeline



3D Reconstruction Pipeline



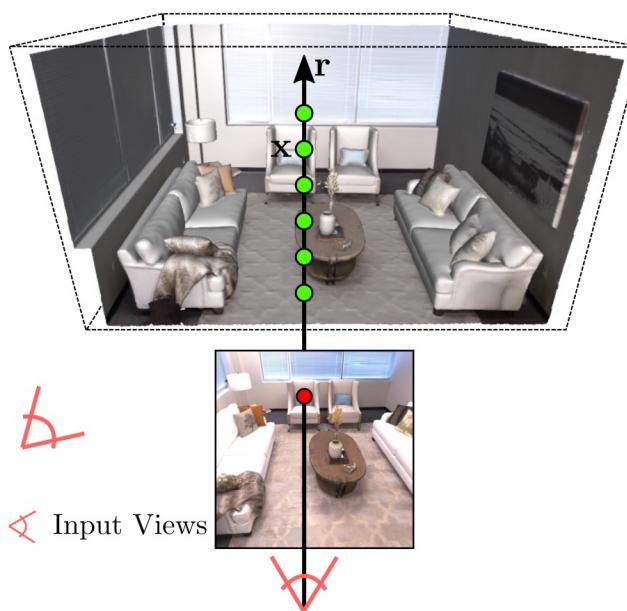
↖



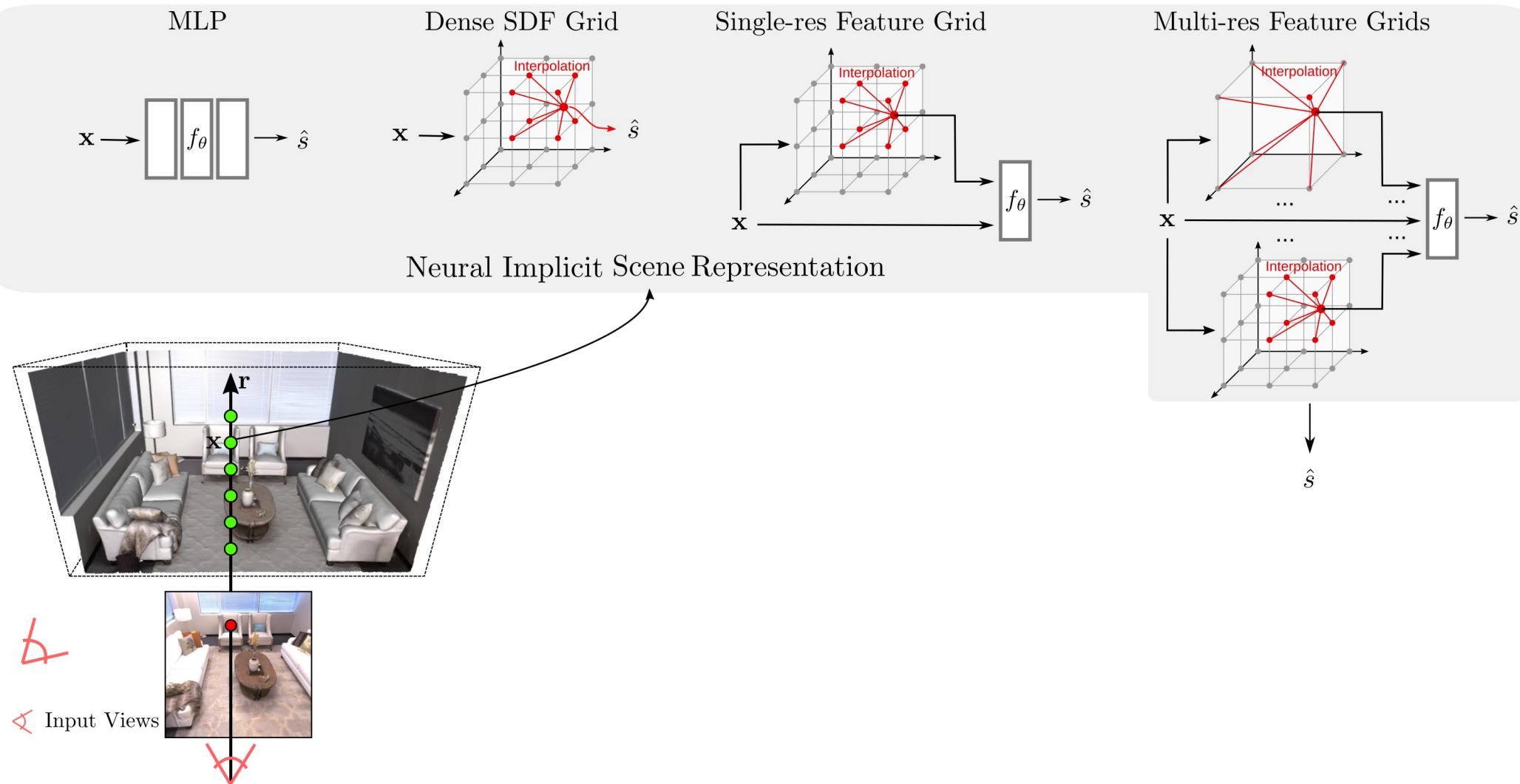
↖ Input Views

↖

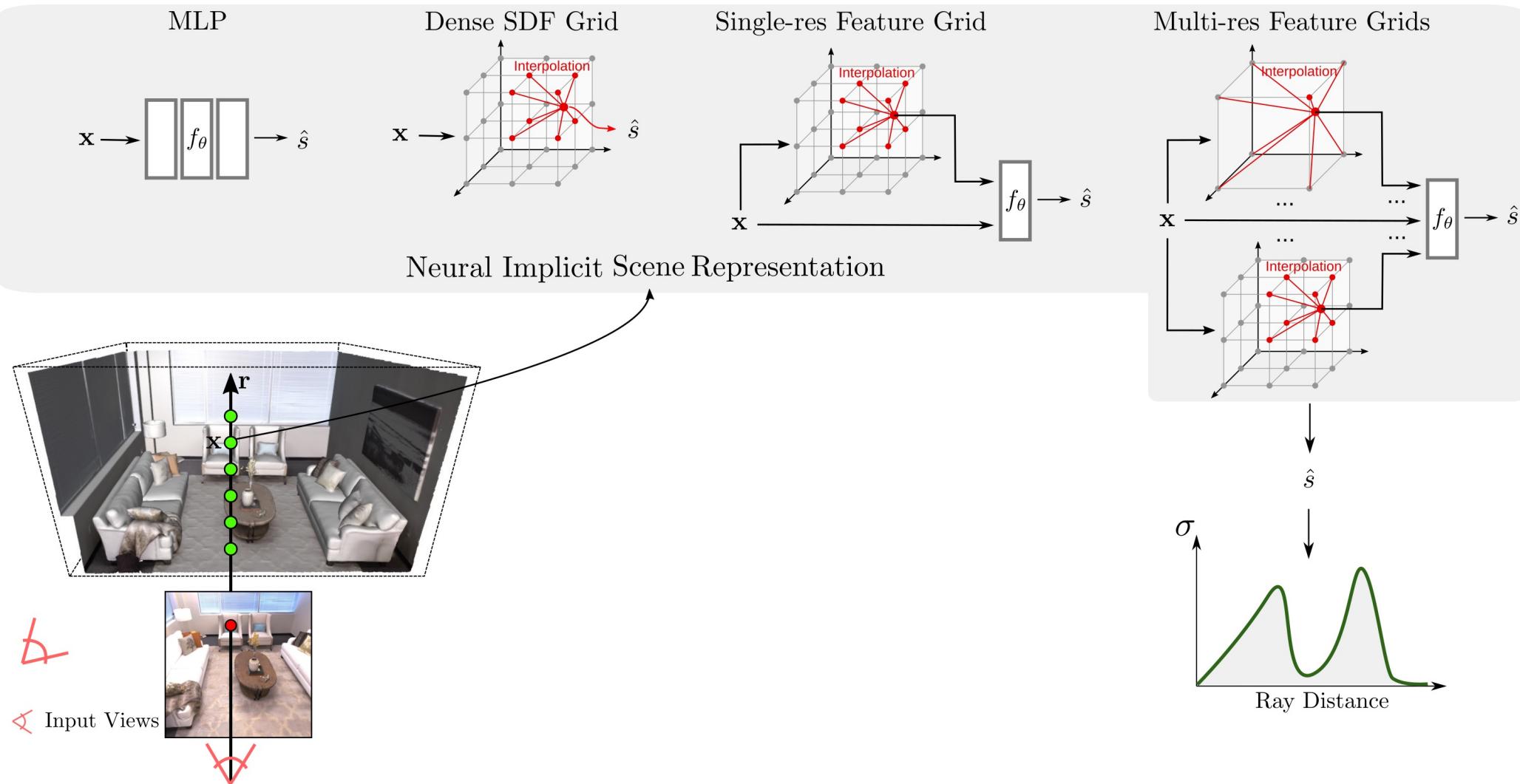
3D Reconstruction Pipeline



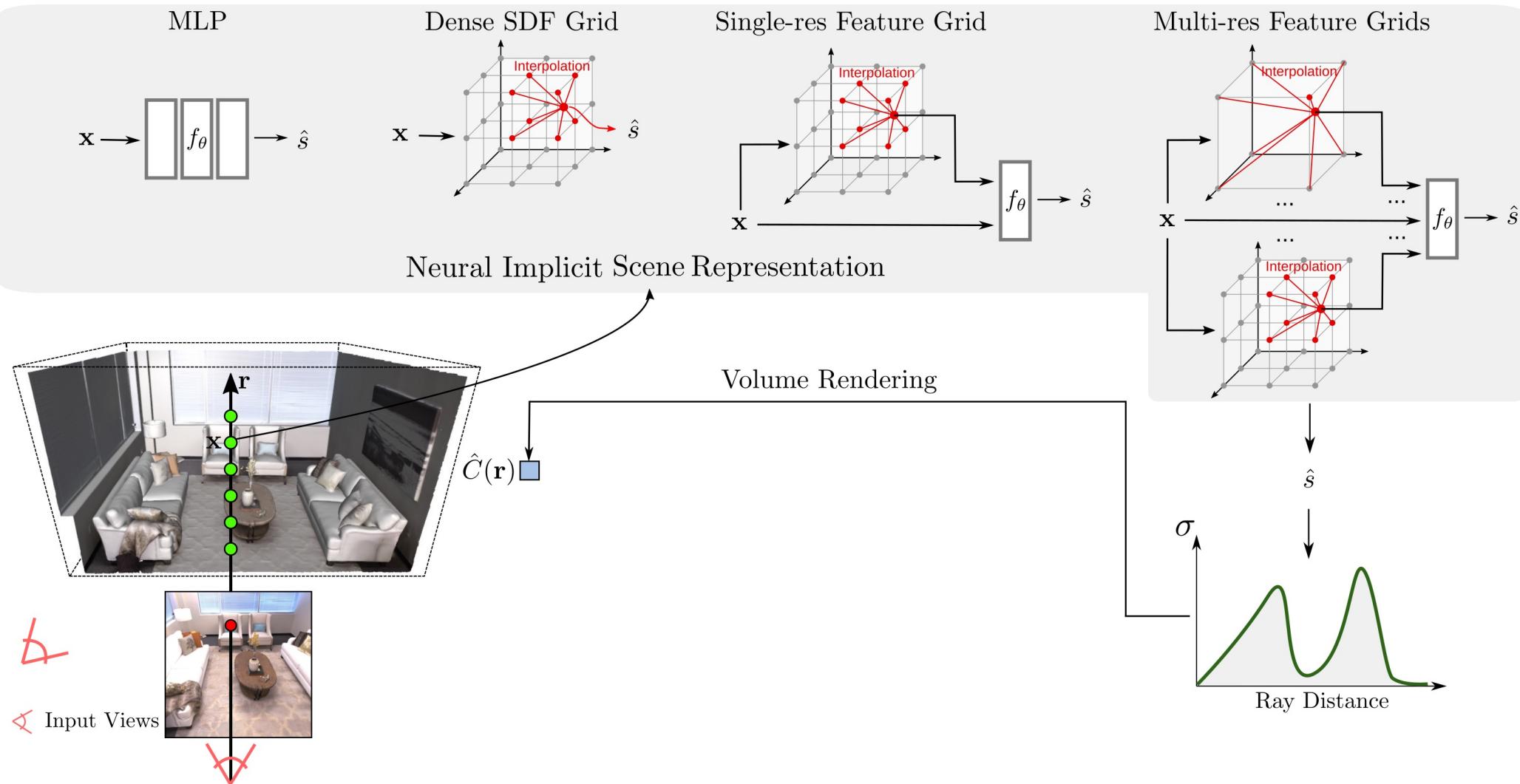
3D Reconstruction Pipeline



3D Reconstruction Pipeline

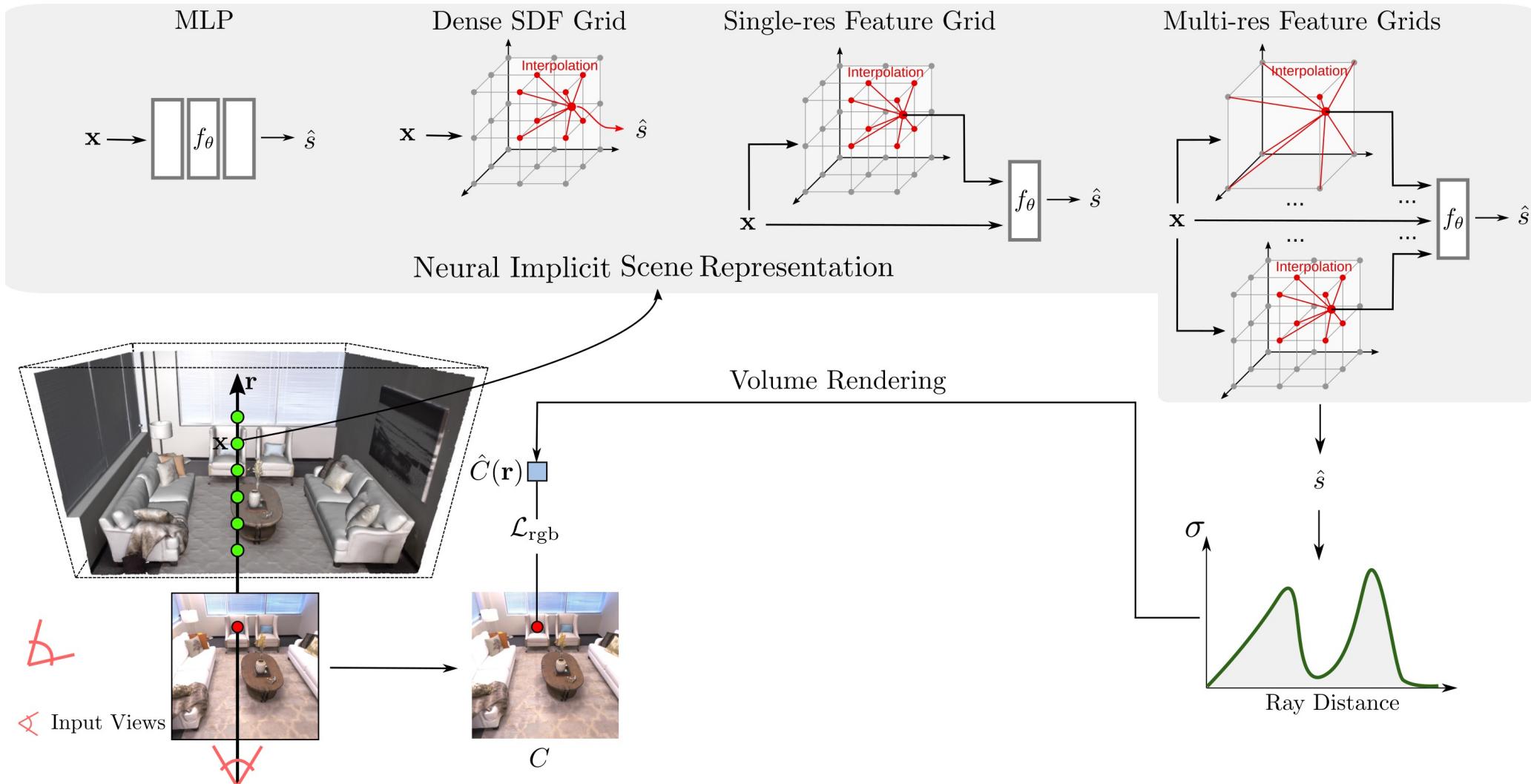


3D Reconstruction Pipeline



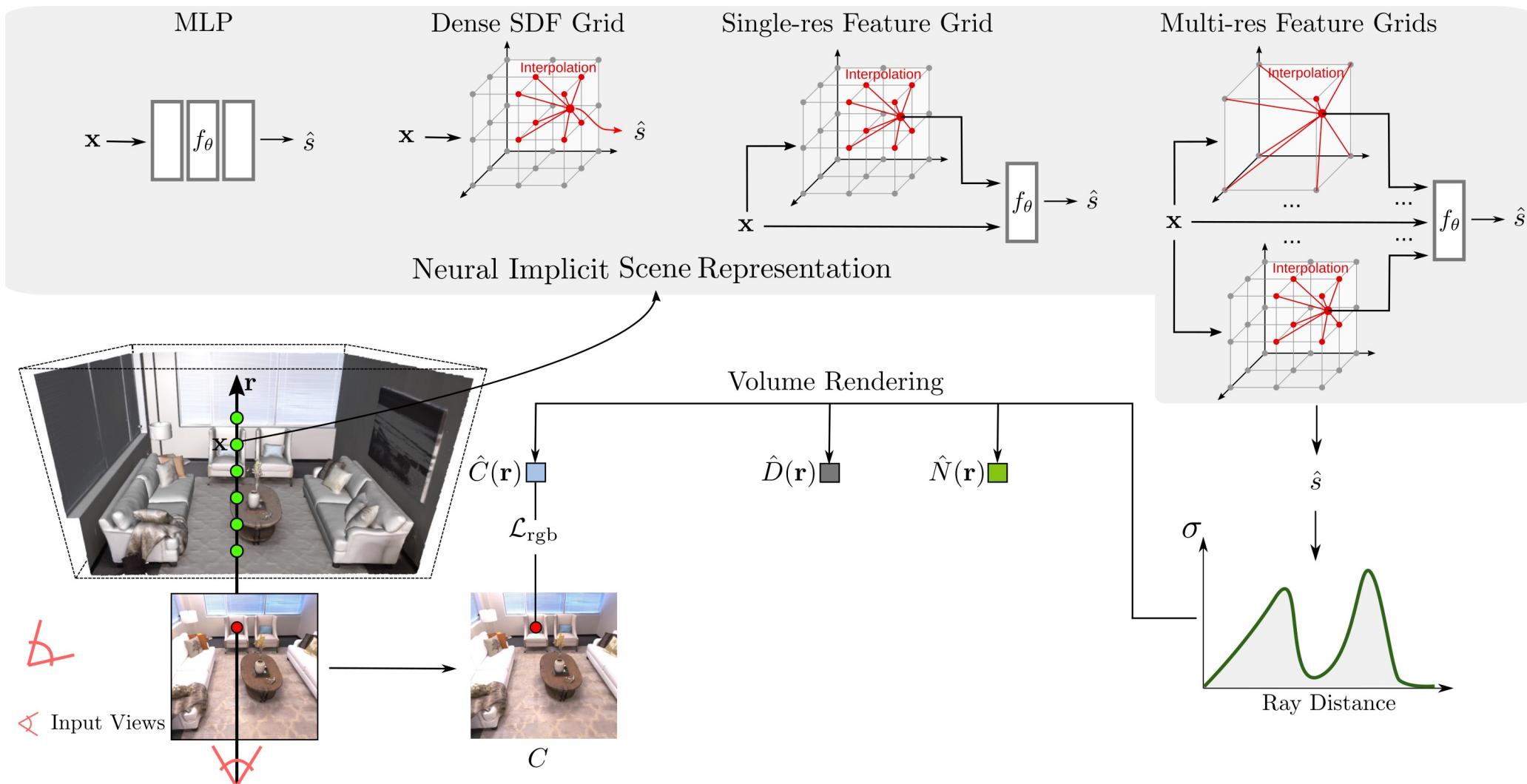
3D Reconstruction Pipeline

VolSDF/NeuS/UNISURF/Neuralangelo



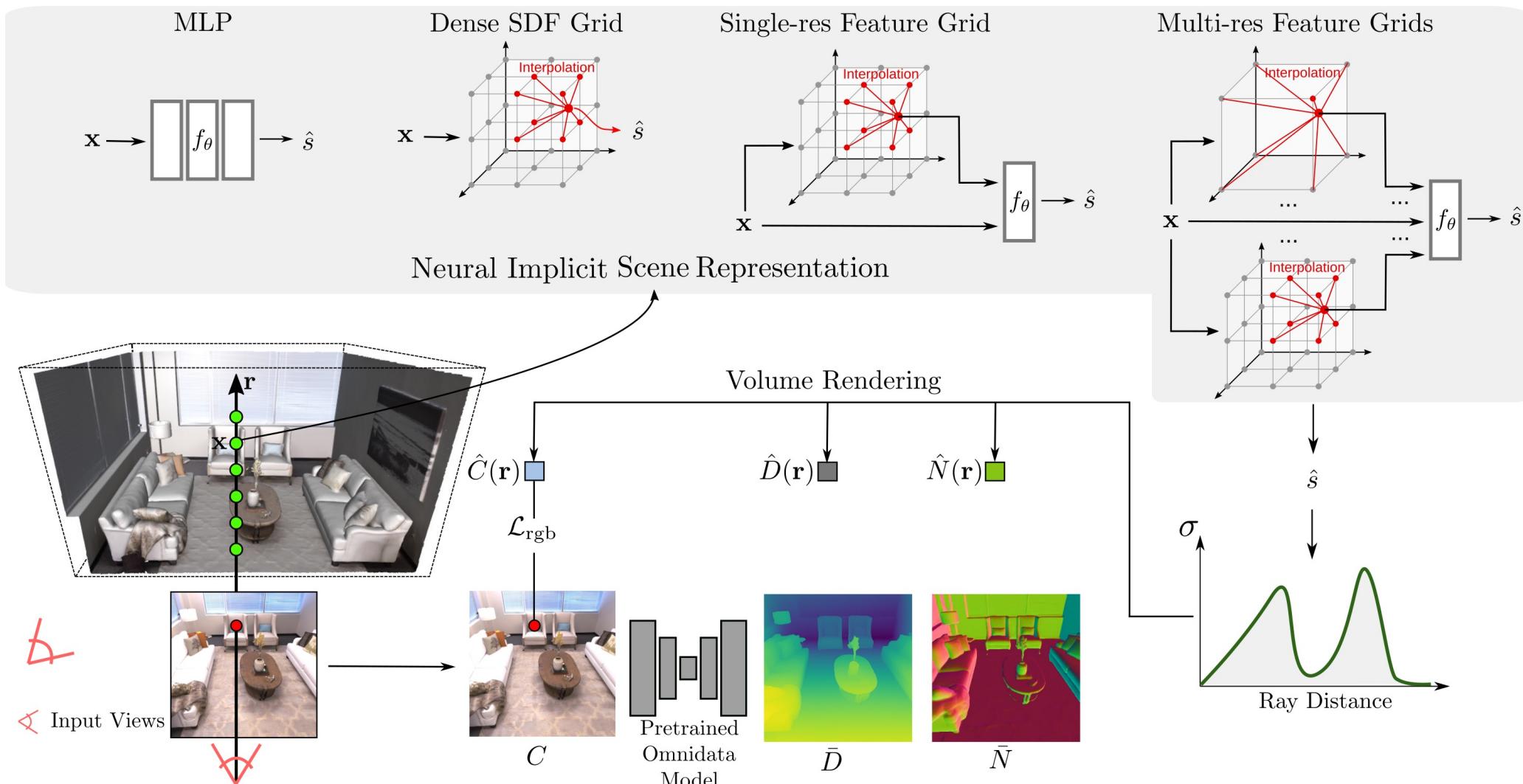
3D Reconstruction Pipeline

MonoSDF



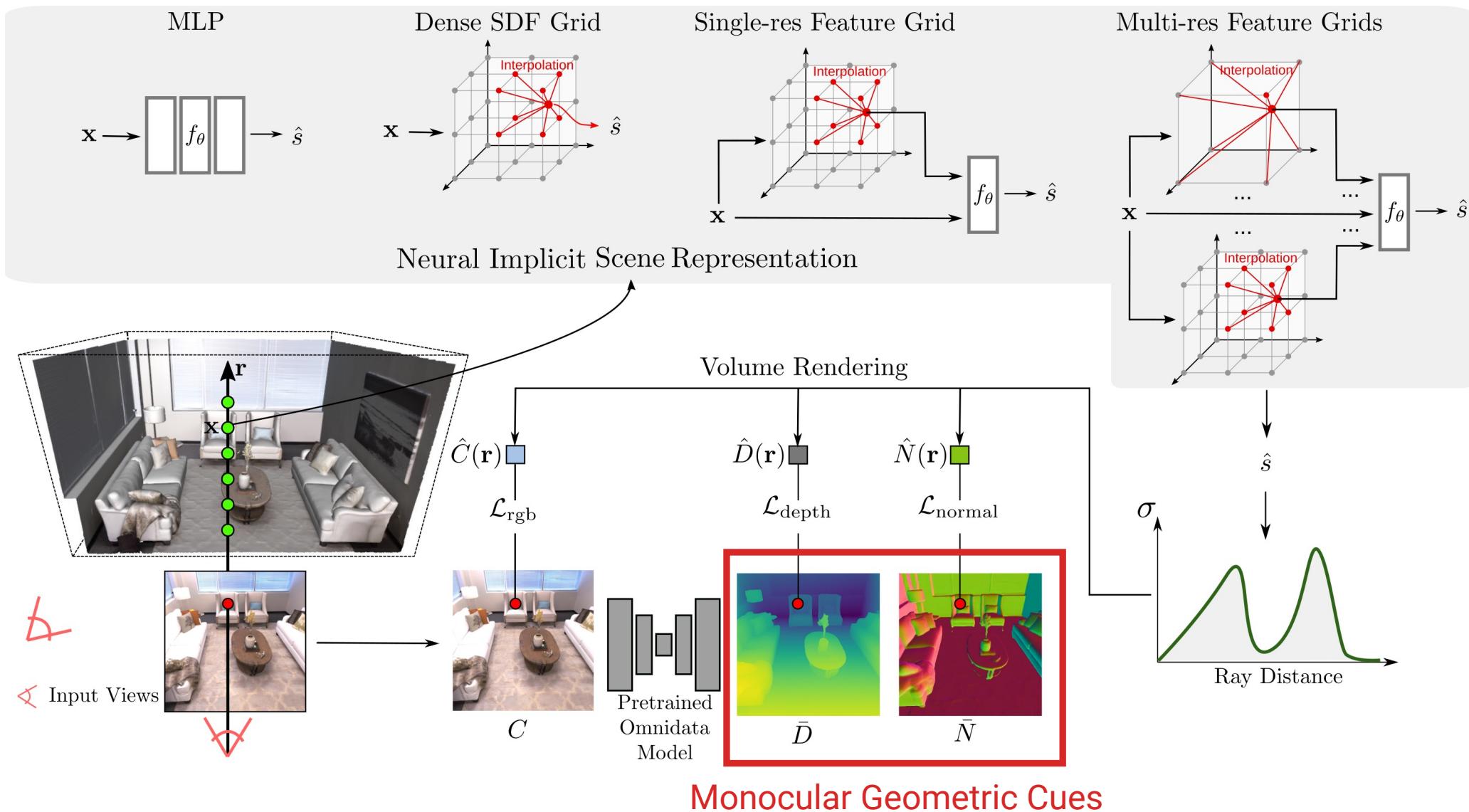
3D Reconstruction Pipeline

MonoSDF



3D Reconstruction Pipeline

MonoSDF



2D Magic in a 3D World

2D Monocular Cues Benefit 3D Reconstruction



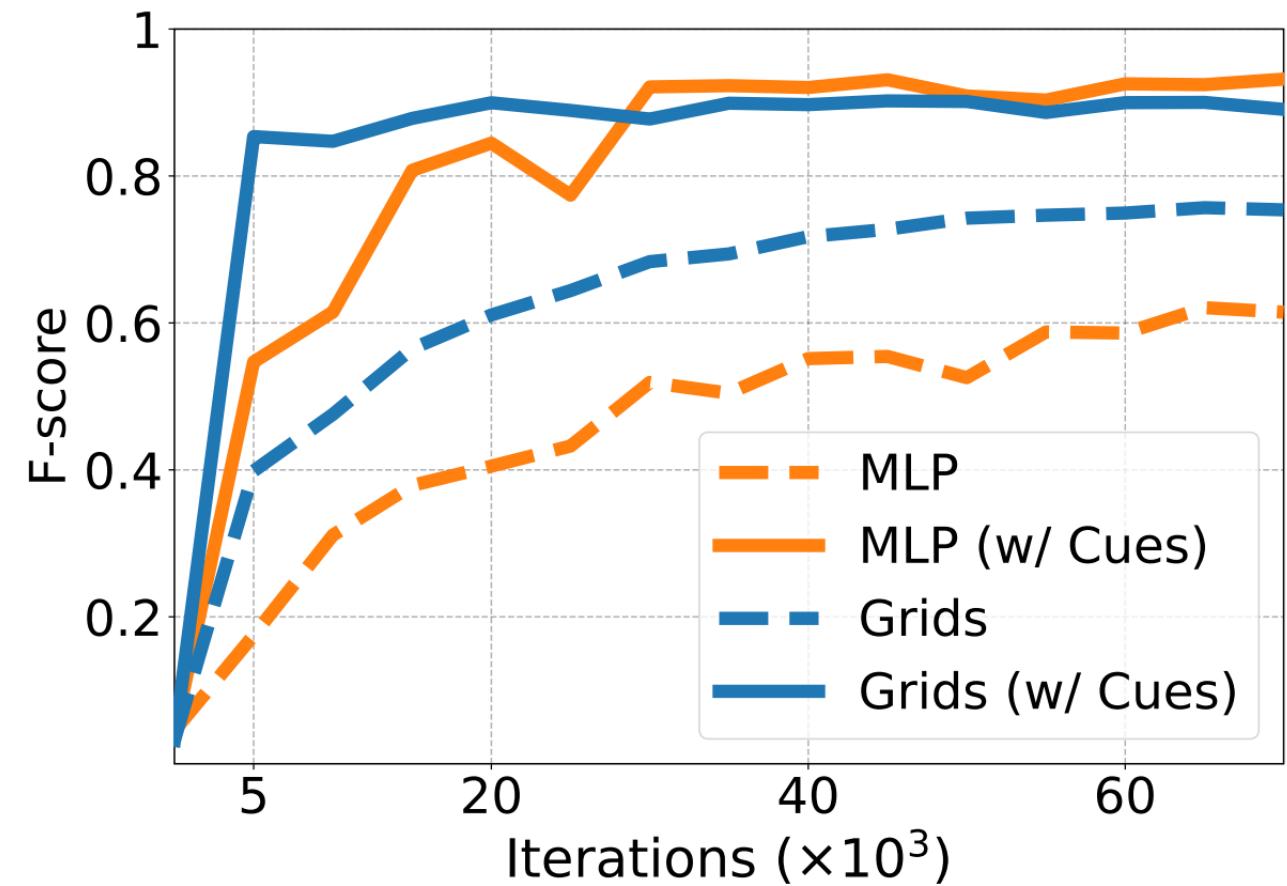
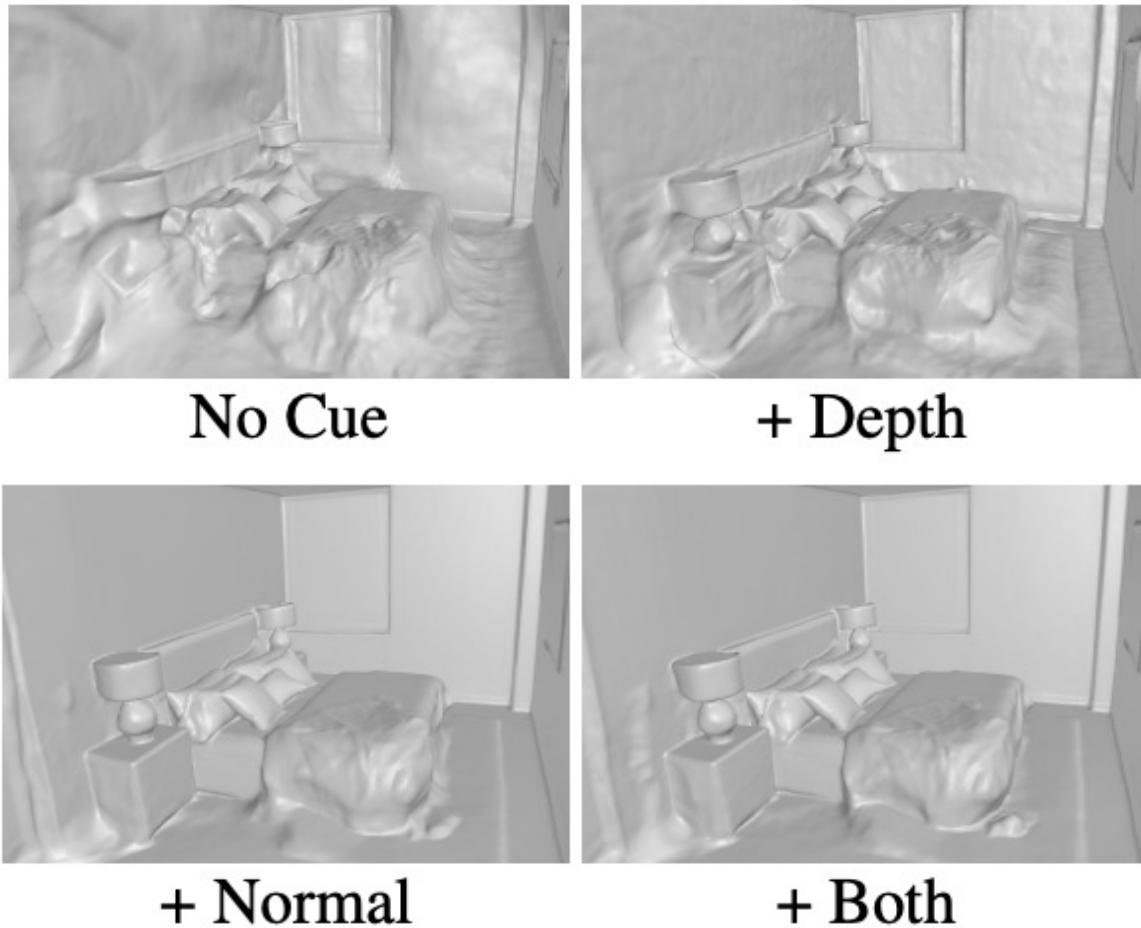
Really?

Results

Baseline Comparison



Ours



Large-Scale 3D Scene Reconstruction

Tanks & Temples Dataset



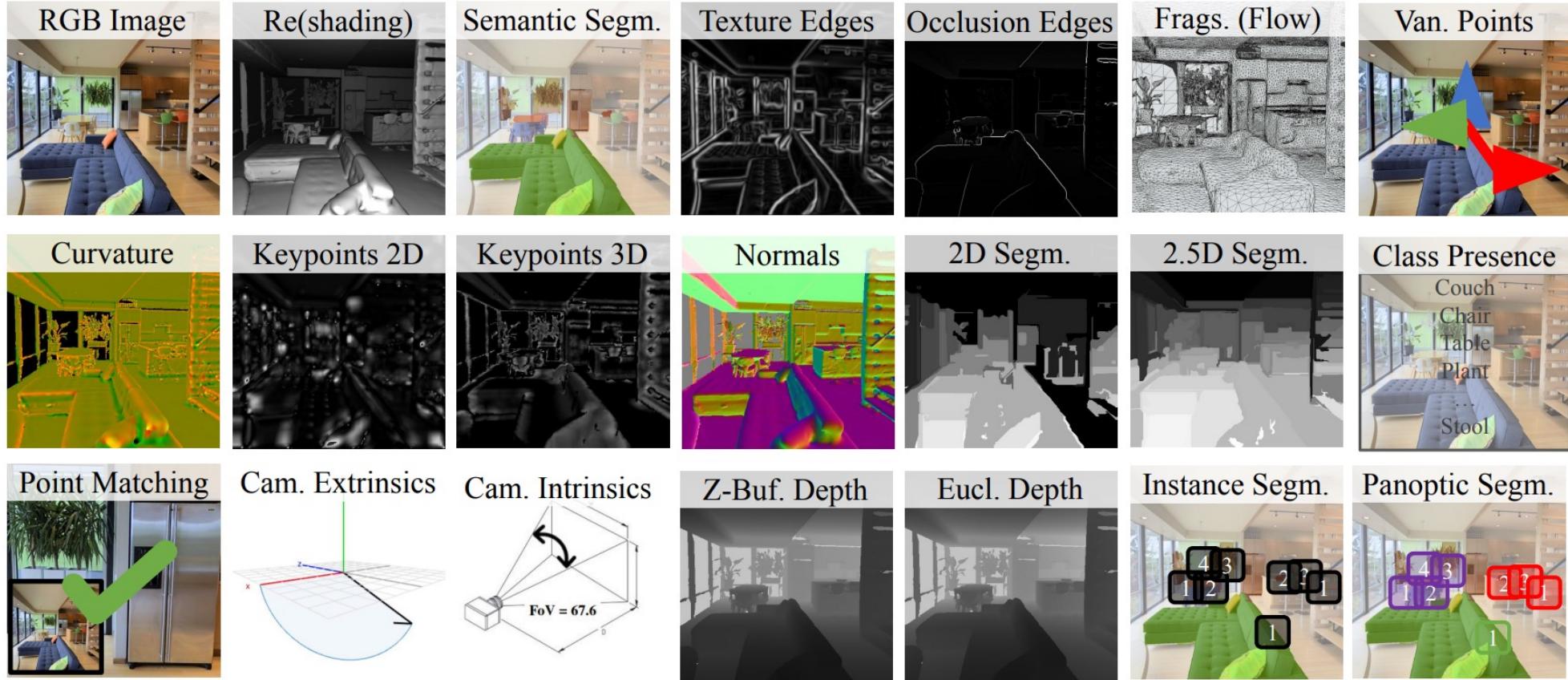
2D Magic in a 3D World

2D Monocular Cues Benefit 3D Reconstruction



Omnidata

~14M Training Images



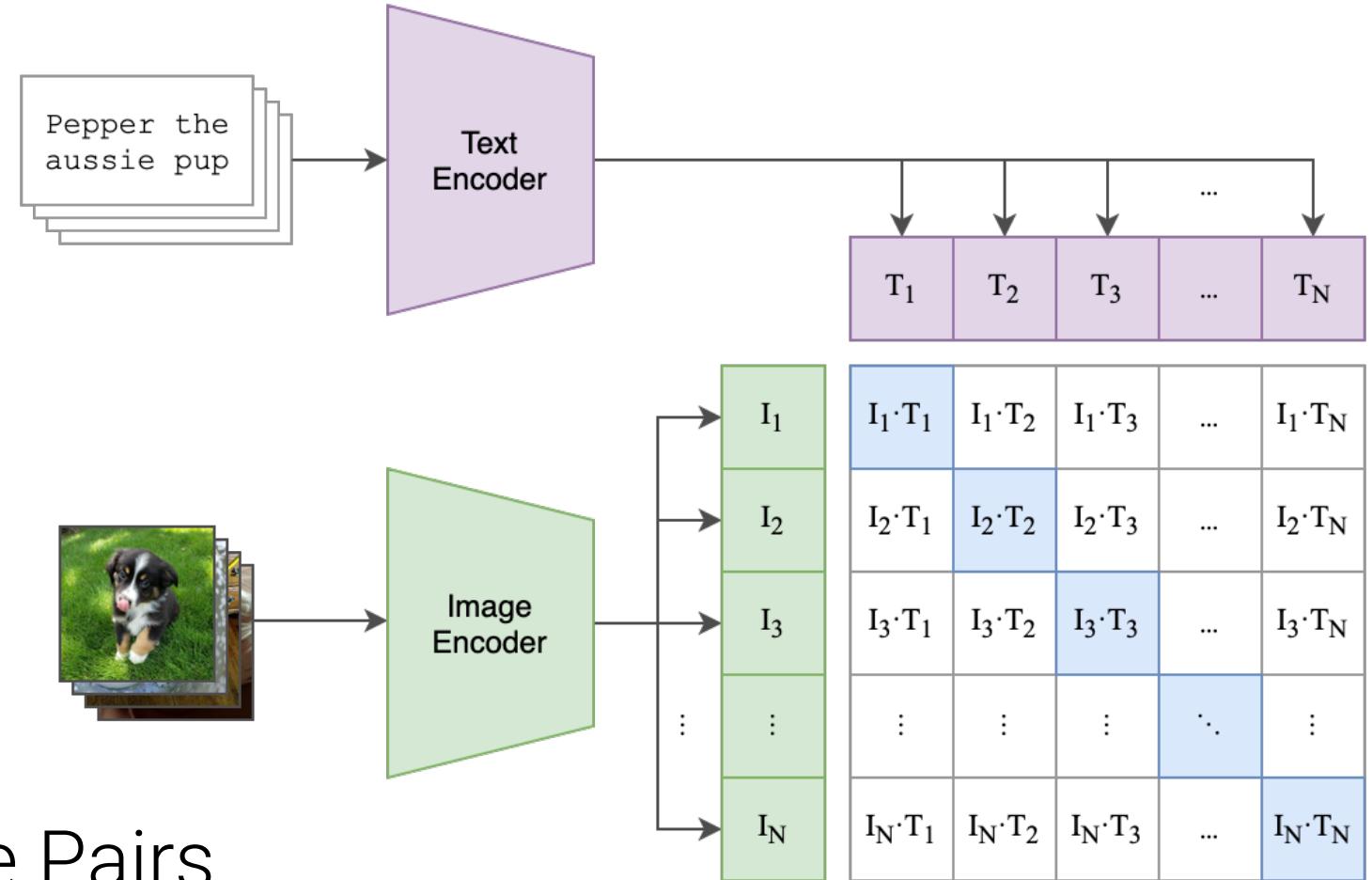
What happened since 2021?

Text-Image Pretraining

CLIP



~**12.8B** Training Text-Image Pairs



Text-to-Image Generation

DALL·E 3



Imagen

Google Research

Stable Diffusion stability.ai

SO

You

Generate an image of “a nervous and nerdy guy presenting in front of many smart researchers at a lecture hall at Dyson Robotics Lab in Imperial College London”



ChatGPT



Betker et al.: [Improving Image Generation with Better Captions](#). 2023

Saharia et al.: [Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding](#). NeurIPS 2022

Rombach et al.: [High-Resolution Image Synthesis with Latent Diffusion Models](#). CVPR 2022

2D Image Segmentation

SAM

∞ Meta

~1B Mask + Human-in-the-loop

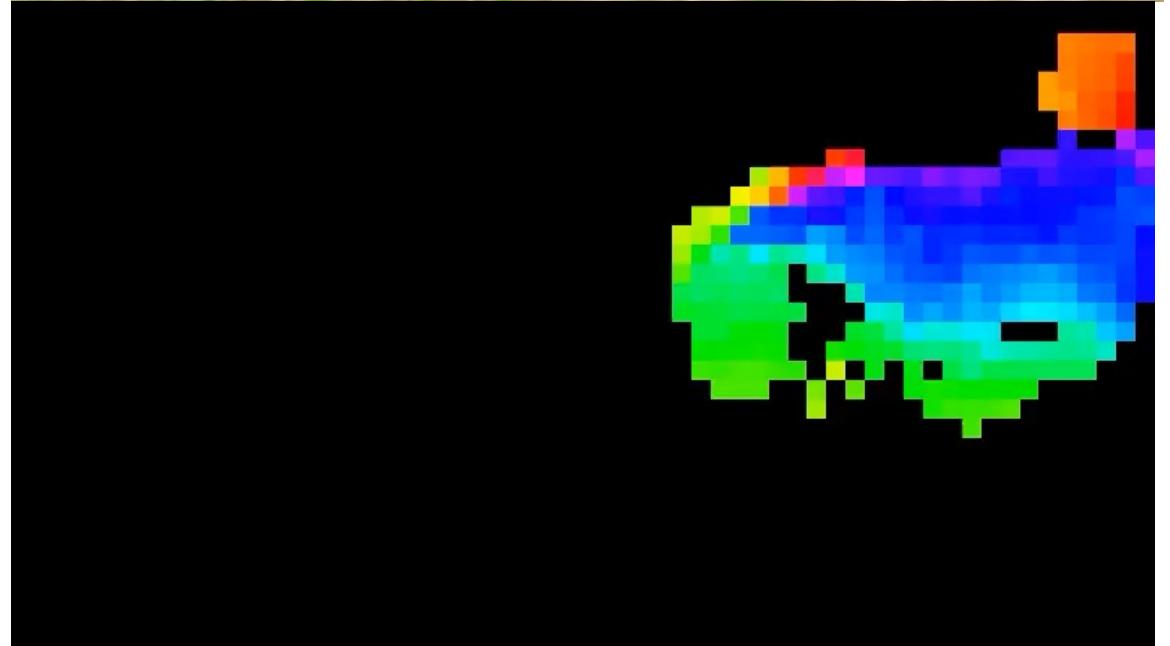


2D Visual Features

DINO v2

 Meta

1.2B Training Images



Text-to-Video Generation

Sora



OpenAI

??? B Training Videos?



Prompt: A movie trailer featuring the adventures of the 30 year old space man wearing a red wool knitted motorcycle helmet, blue sky, salt desert, cinematic style, shot on 35mm film, vivid colors.

2D Magic in a 3D World

Songyou Peng

Imperial College London
Mar 22, 2024

2D Magic in a 3D World

2D Foundation Models for 3D Vision Tasks

2D Magic in a 3D World

2D Foundation Models for 3D Vision Tasks



2D magics in a a a 3D world-
2D foundation models for 3D vision stt acks

[Generated by DALL·E 3]

2D Magic in a 3D World

2D Foundation Models for 3D Vision Tasks

3D Reconstruction



NeRF On-the-go
CVPR 2024

3D Scene Understanding



OpenScene
CVPR 2023



Segment3D
arXiv 2024

2D Magic in a 3D World

2D Foundation Models for 3D Vision Tasks

3D Reconstruction



NeRF On-the-go
CVPR 2024

3D Scene Understanding



OpenScene
CVPR 2023



Segment3D
arXiv 2024

NeRF Is Awesome





TOUR EIFFEL · CHAMP-DE-MARS · MUSÉE DU LOUVRE · NOTRE-DAME · MUSÉE D'ORSAY · OPÉRA GARNIER · CHAMPS-ÉLYSÉES · GRAND PALAIS · TROCADÉRO

BIGBUS PARIS·LESCARS ROUGES

**BIG
BUS**

HOP-ON HOP-O

Motivation

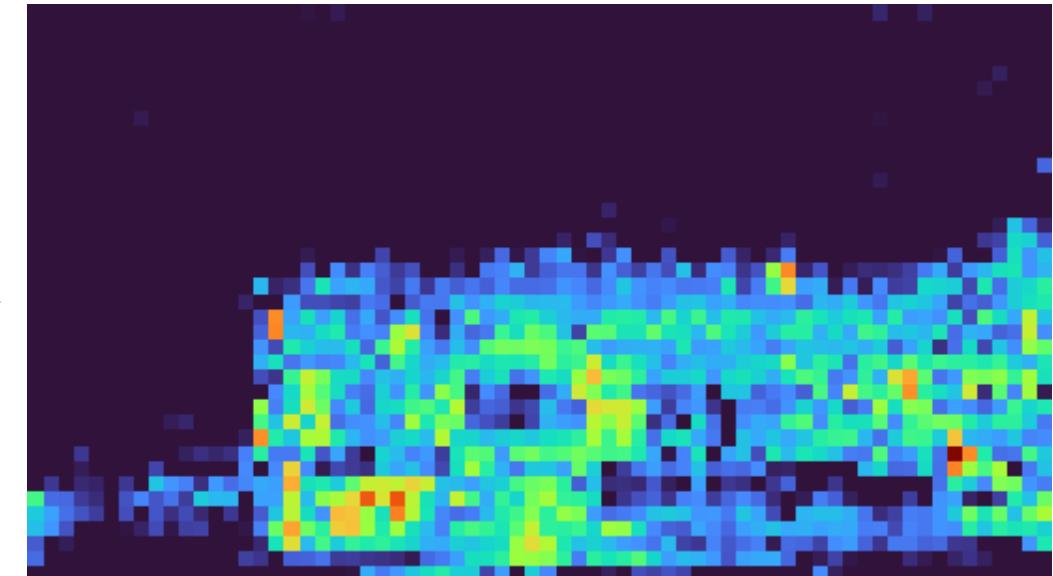


How to obtain **distractor-free NeRFs** from **casually captured sequences**?

Uncertainty



Input RGB



Uncertainty Map

How to learn a good uncertainty map?

DINO v2



- A 2D foundation model producing **universal features**
- Preserve temporal-spatial consistency

How to Leverage the **2D Foundation Model** for **Distractor-free NeRF?**



NeRF *On-the-go*

Exploiting Uncertainty for Distractor-free NeRFs in the Wild



Weining
Ren*



Zihan
Zhu*



Boyang
Sun



Julia
Chen



Marc
Pollefeys



Songyou
Peng

Pipeline

DINOv2 Feature Map



DINOv2



RGB Input

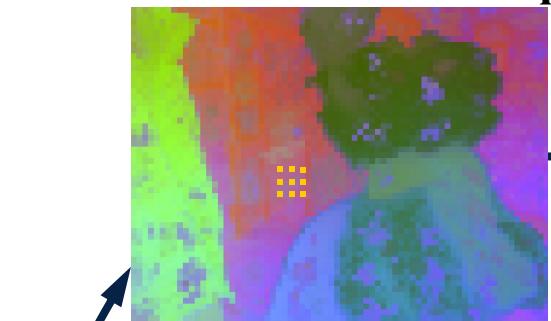
G
Uncertainty
MLP

$$\beta(r)$$

$C(r)$
NeRF
Representation

Pipeline

DINOv2 Feature Map



DINOv2



RGB Input

G
Uncertainty
MLP

$\beta(r)$

$C(r)$

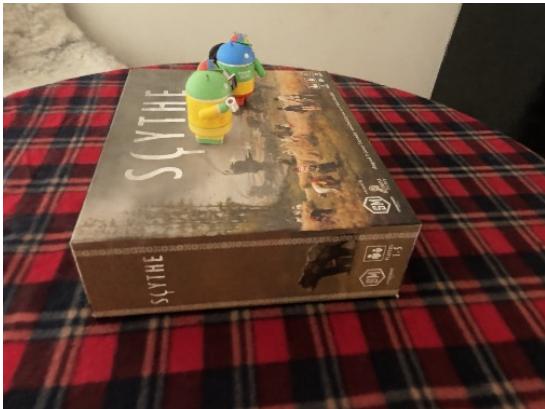
SSIM

$\mathcal{L}_{\text{uncer}}$

$$\mathcal{L}_{\text{uncer}}(r) = \frac{\mathcal{L}_{\text{SSIM}}}{2\beta(r)^2} + \lambda_1 \log \beta(r)$$

NeRF
Representation

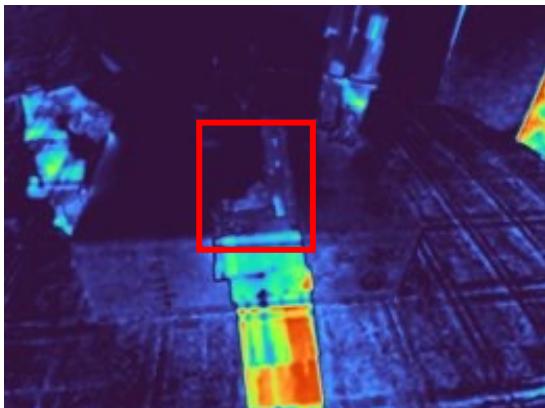
To Learn the Uncertainty MLP...



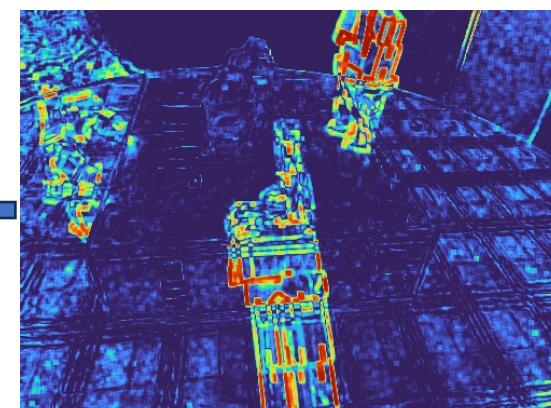
Rendered RGB



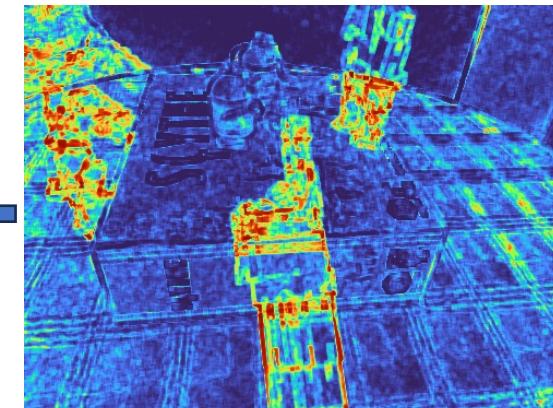
Train RGB



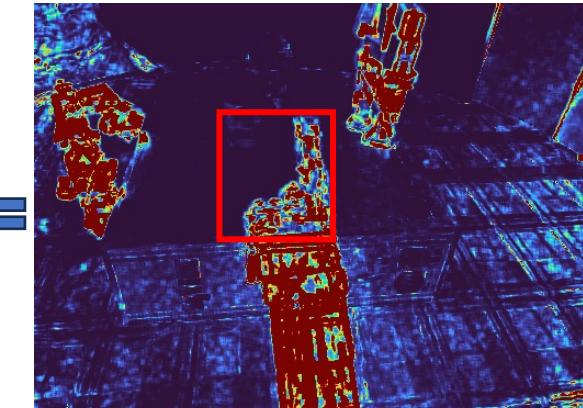
Luminance



Contrast



Structure



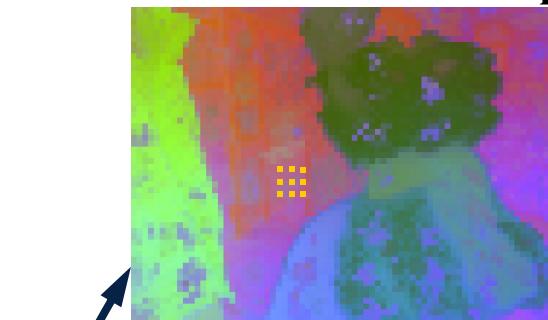
SSIM Error

Why SSIM?

Leverage structure information when RGB is similar!

Pipeline

DINOv2 Feature Map



DINOv2



RGB Input

G
Uncertainty
MLP

$C(r)$

NeRF
Representation

SSIM

$\beta(r)$

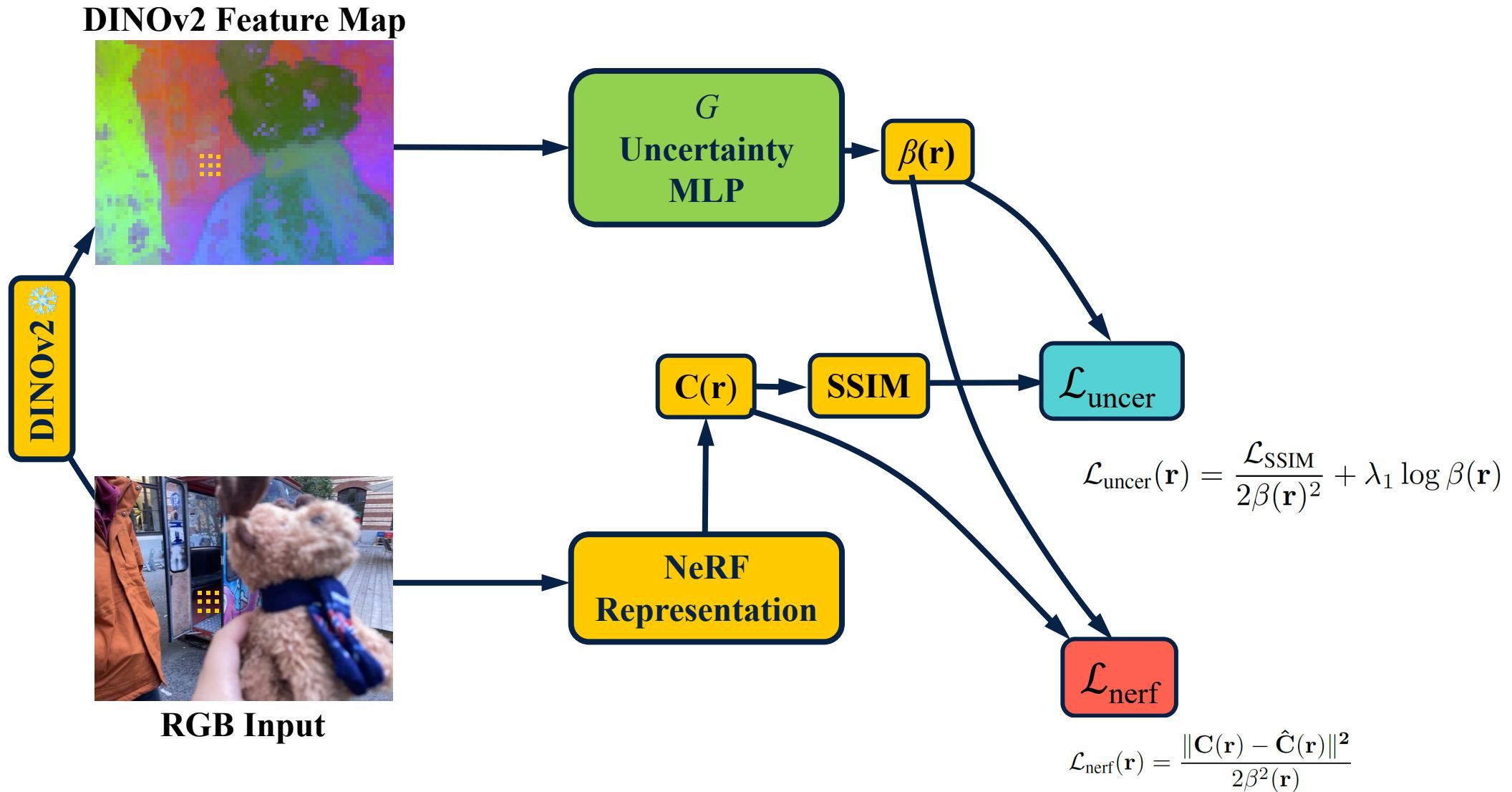
$\mathcal{L}_{\text{uncer}}$

$\mathcal{L}_{\text{nerf}}$

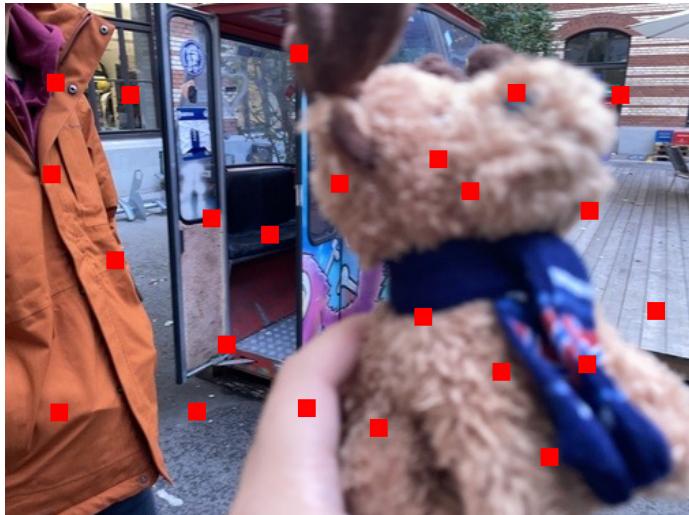
$$\mathcal{L}_{\text{uncer}}(r) = \frac{\mathcal{L}_{\text{SSIM}}}{2\beta(r)^2} + \lambda_1 \log \beta(r)$$

$$\mathcal{L}_{\text{nerf}}(r) = \frac{\|\mathbf{C}(r) - \hat{\mathbf{C}}(r)\|^2}{2\beta^2(r)}$$

Pipeline



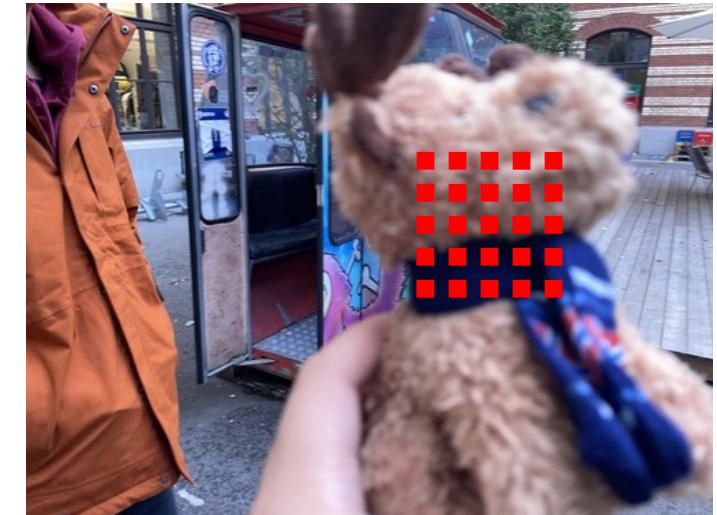
Sampling Strategy



Random
(NeRF)



Patch
(RobustNeRF)

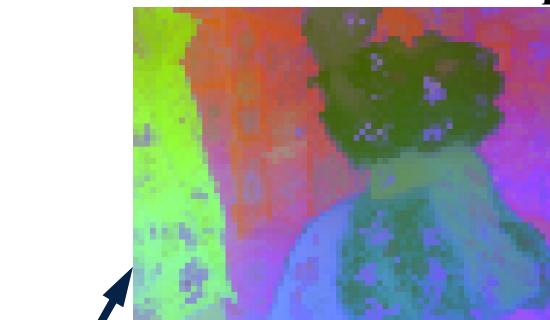


Dilated Patch
(Ours)

✓ **Larger Perceptive Field:** Improve optimization efficiency & reconstruction quality

Pipeline

DINOv2 Feature Map



DINOv2



RGB Input

G
Uncertainty
MLP

$\beta(r)$

$C(r)$

SSIM

$\mathcal{L}_{\text{uncer}}$

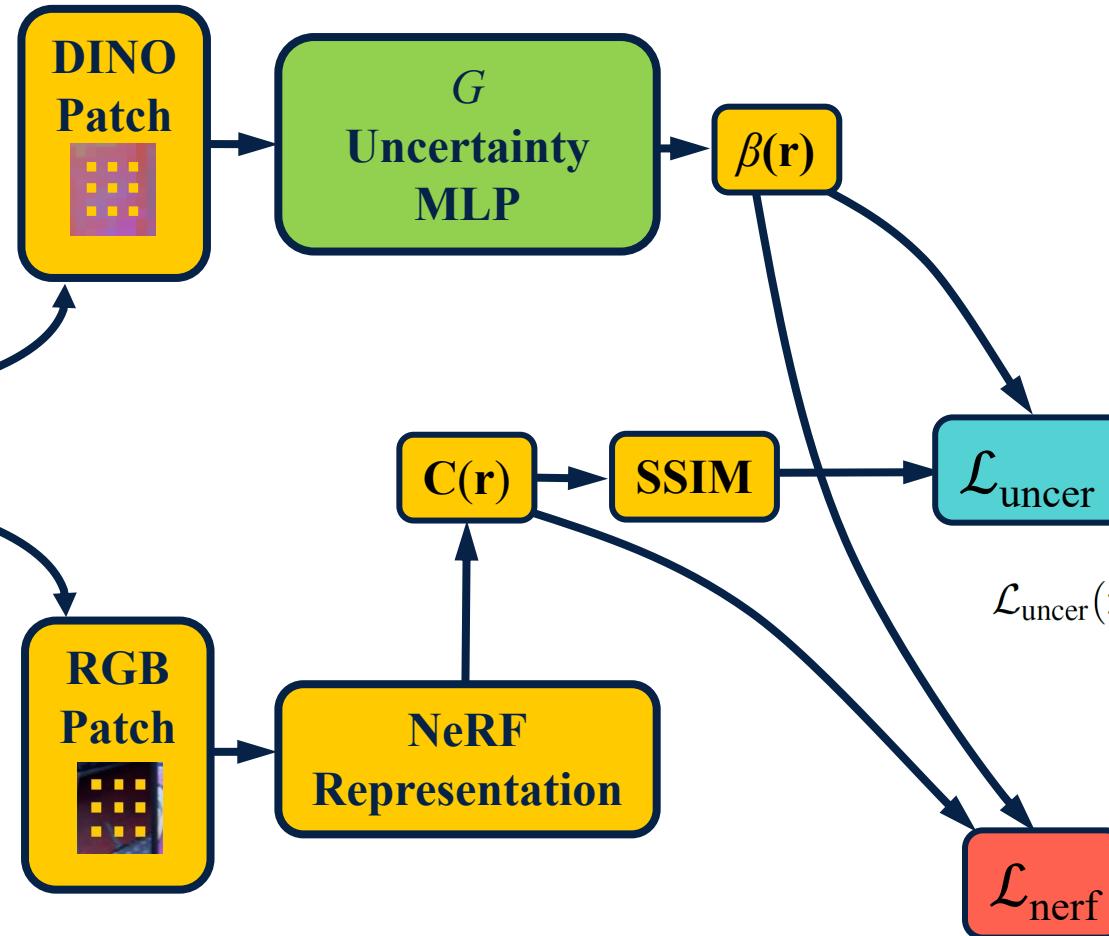
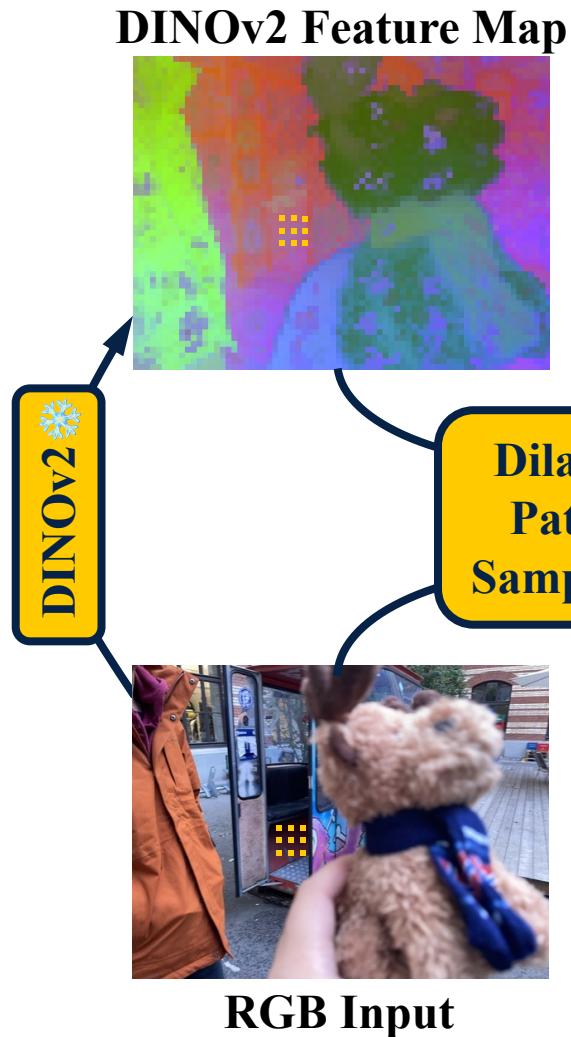
NeRF
Representation

$\mathcal{L}_{\text{nerf}}$

$$\mathcal{L}_{\text{uncer}}(r) = \frac{\mathcal{L}_{\text{SSIM}}}{2\beta(r)^2} + \lambda_1 \log \beta(r)$$

$$\mathcal{L}_{\text{nerf}}(r) = \frac{\|C(r) - \hat{C}(r)\|^2}{2\beta^2(r)}$$

Pipeline



$$\mathcal{L}_{\text{uncer}}(\mathbf{r}) = \frac{\mathcal{L}_{\text{SSIM}}}{2\beta(\mathbf{r})^2} + \lambda_1 \log \beta(\mathbf{r})$$

$$\mathcal{L}_{\text{nerf}}(\mathbf{r}) = \frac{\|\mathbf{C}(\mathbf{r}) - \hat{\mathbf{C}}(\mathbf{r})\|^2}{2\beta^2(\mathbf{r})}$$

Results

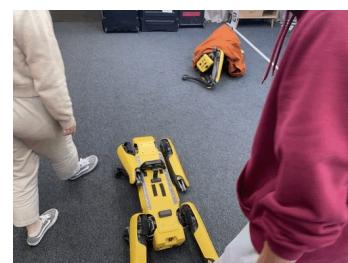
On-the-go Dataset



Low Occlusion (5% ~ 10%)



Medium Occlusion (15% ~ 20%)



High Occlusion (~30%)

Occlusion
Ratio: **Low**



Statue - Input

RobustNeRF



Statue - Rendering Comparisons



Train Station - Input Images



NeRF On-the-go
(Ours)

Train Station - Rendering Comparisons

Occlusion
Ratio: **High**



Patio-High - Input



NeRF On-the-go
(Ours)

Patio-High - Rendering Comparisons

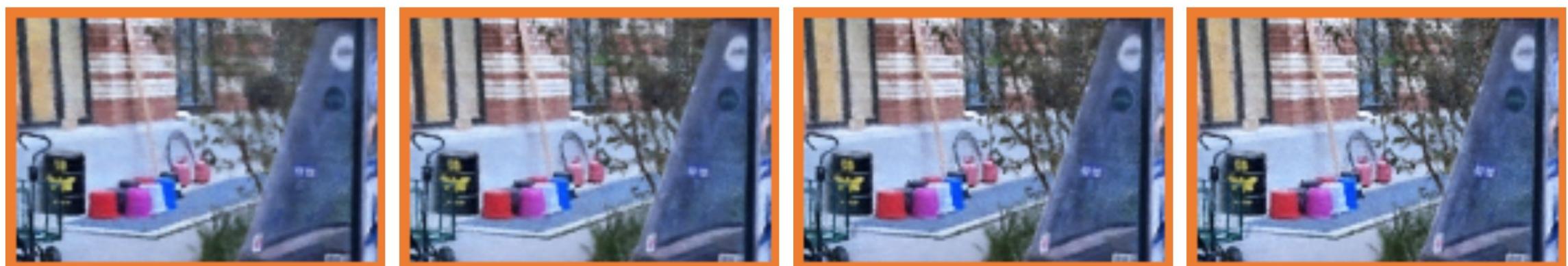
Analysis

Analysis - Efficiency



RobustNeRF

Analysis - Efficiency



25K

50K

100K

250K

**NeRF On-the-go
(Ours)**

Analysis – Static Scene



0.447



0.376



0.374



RobustNeRF



Ours



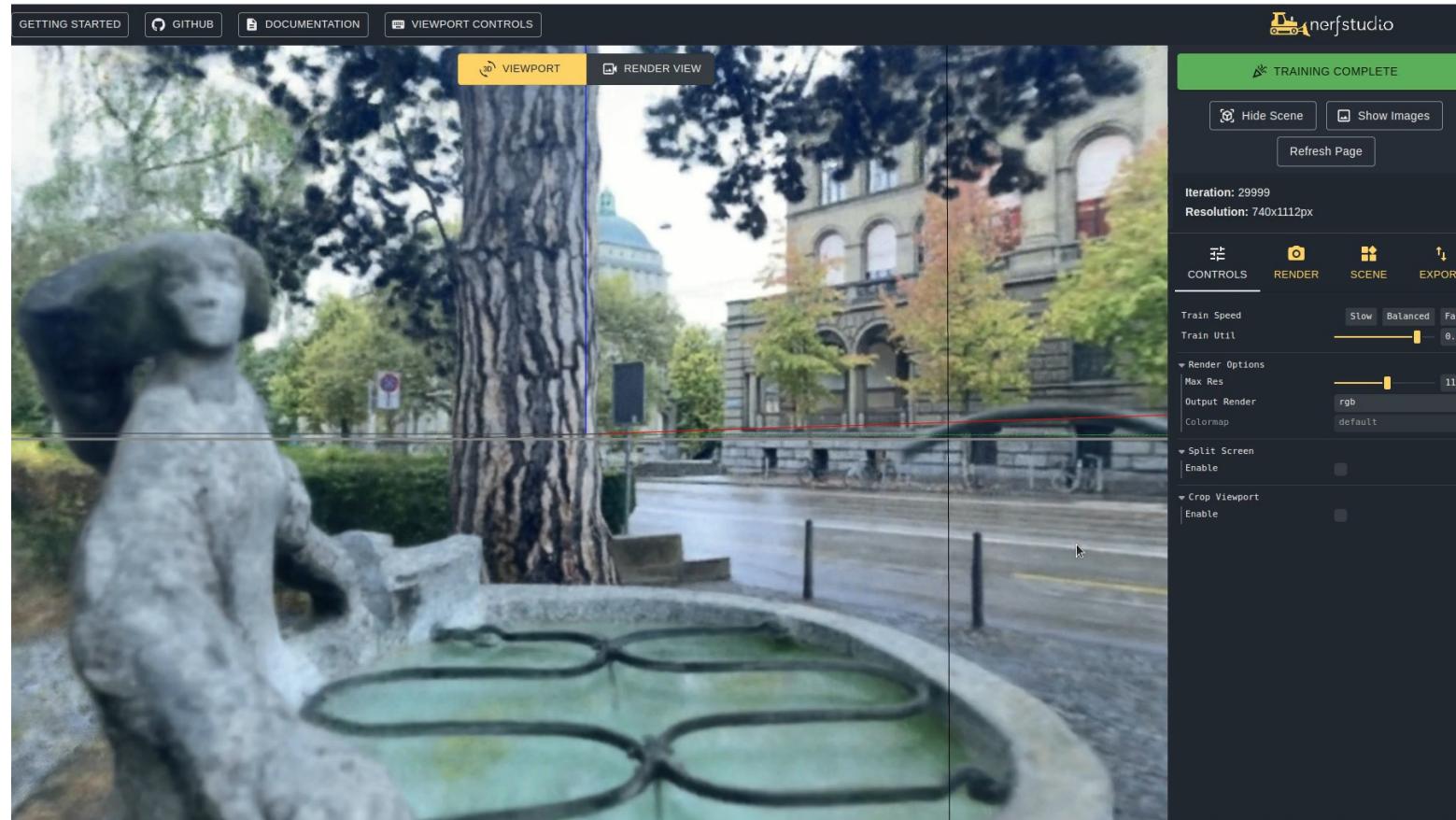
MipNeRF 360



GT

Take-home Messages

- ***On-the-go*** module is plug-and-play for all NeRF methods
 - Integrated into NeRFStudio



Take-home Messages

- ***On-the-go*** module is plug-and-play for all NeRF methods
 - Integrated into NeRFStudio
- **2D foundation model** (DINOv2) rocks!

How to improve upon it?

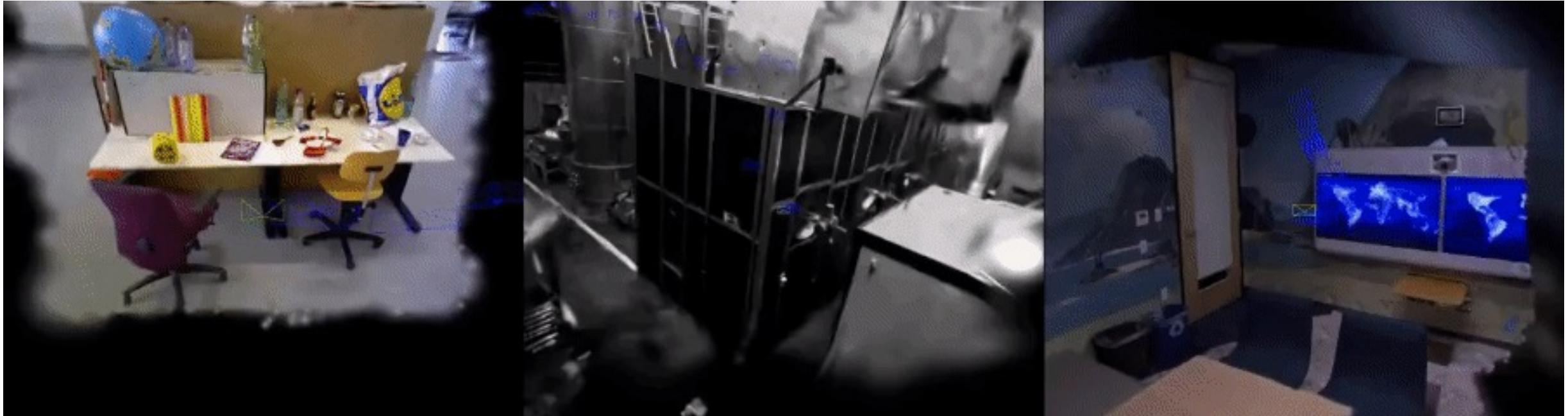
NeRF On-the-go

for VERY Large Urban Scenes



NeRF On-the-go

Without COLMAP



2D Magic in a 3D World

2D Foundation Models for 3D Vision Tasks

3D Reconstruction



NeRF On-the-go
CVPR 2024

3D Scene Understanding



OpenScene
CVPR 2023



Segment3D
arXiv 2024

2D Magic in a 3D World

2D Foundation Models for 3D Vision Tasks

3D Reconstruction



NeRF On-the-go
CVPR 2024

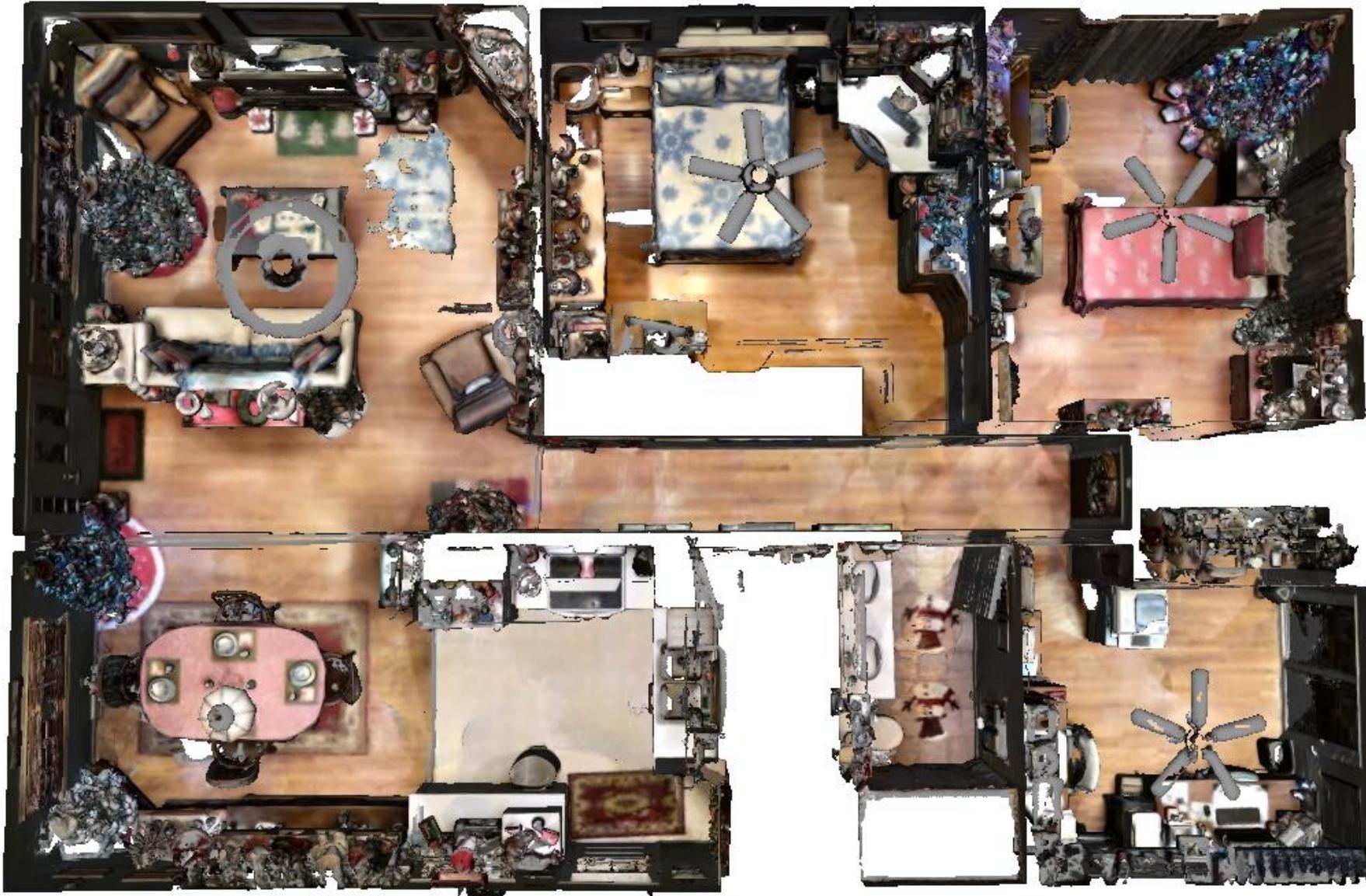
3D Scene Understanding



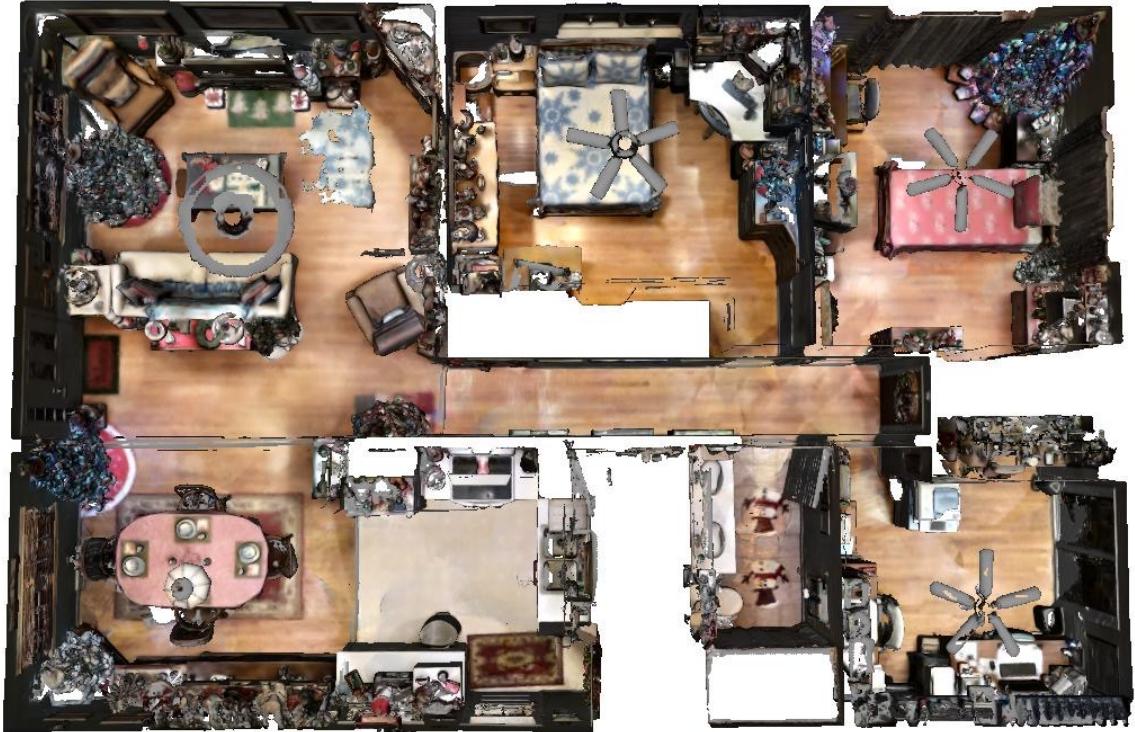
OpenScene
CVPR 2023



Segment3D
arXiv 2024



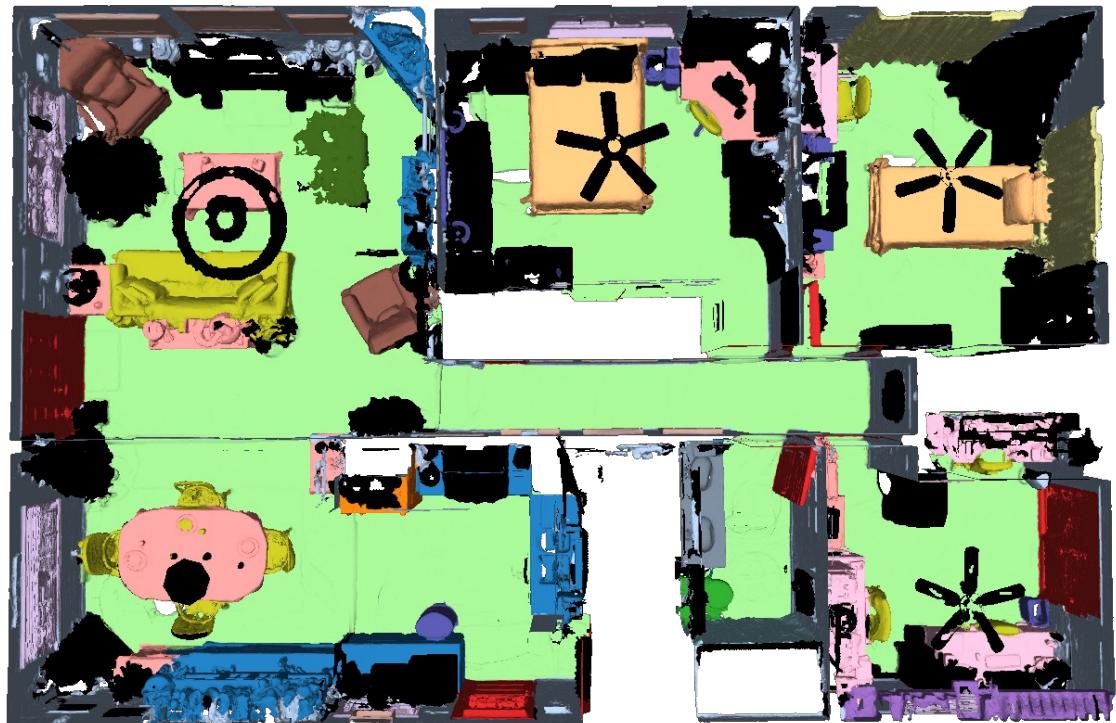
Input 3D Geometry



Input 3D Geometry

Legend:

- wall
- floor
- cabinet
- bed
- chair
- sofa
- table
- door
- window
- counter
- curtain
- toilet
- sink
- bathtub
- other
- unlabeled



Traditional 3D Scene Understanding
(e.g. Semantic Segmentation)
Only train and test on a few common classes

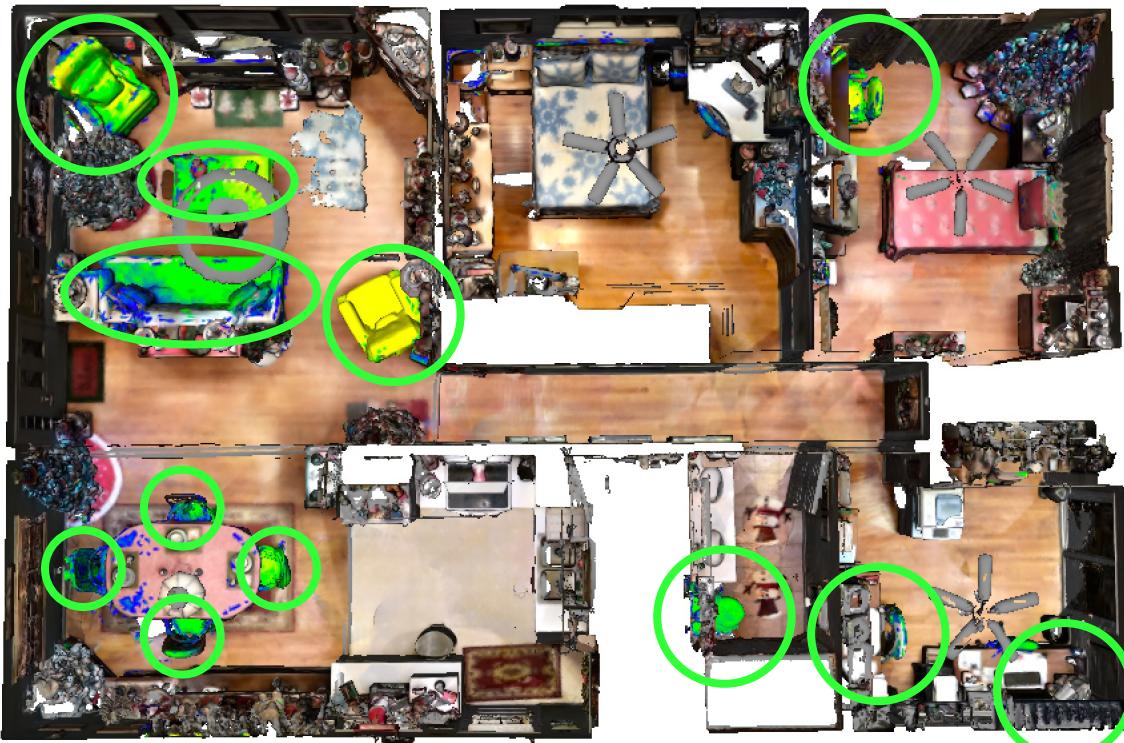
3D Scene Understanding Tasks w/o Labels



Input 3D Geometry

- Affordance prediction

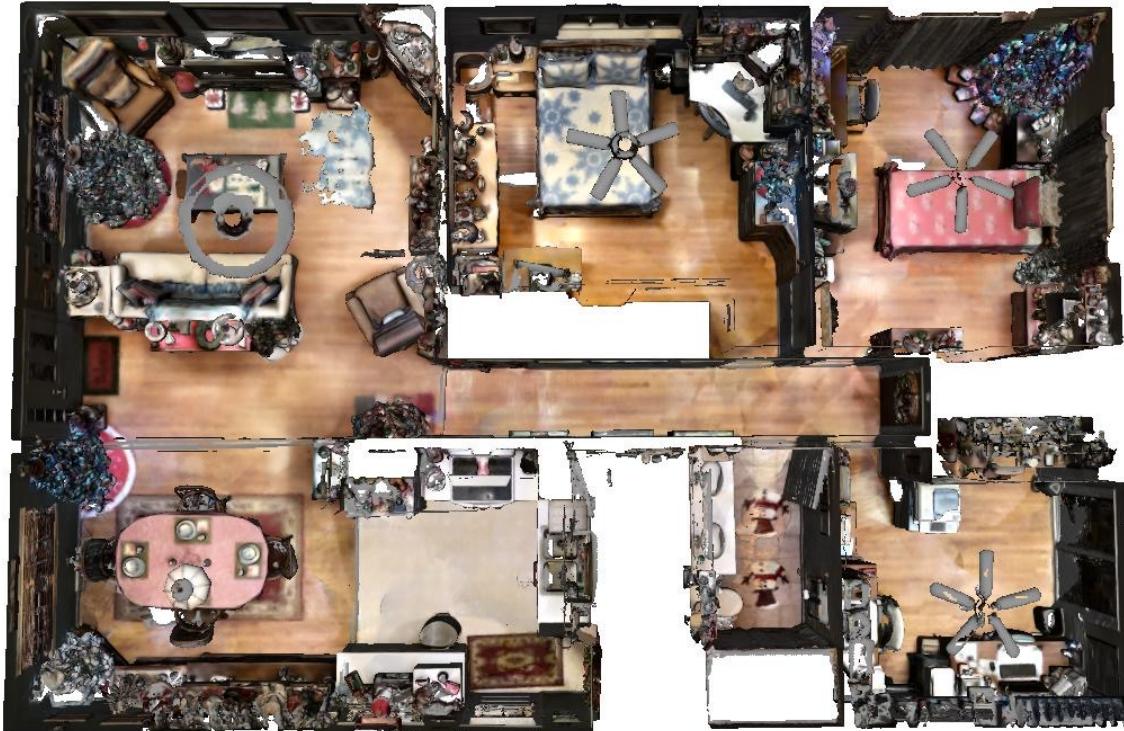
3D Scene Understanding Tasks w/o Labels



- Affordance prediction

Example: “where can I sit?”

3D Scene Understanding Tasks w/o Labels



Input 3D Geometry

- Affordance prediction
- Material identification
- Physical property estimation
- Rare object retrieval
- Activity site prediction
- Fine-grained semantic segmentation
- Many more...

How to have a single model for all these 3D tasks
without any labeled 3D data?

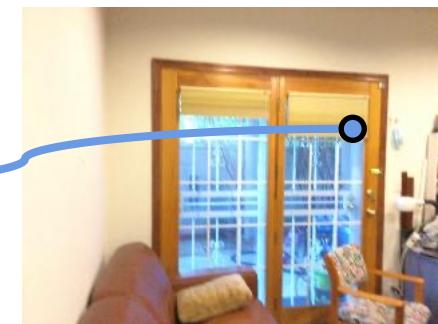
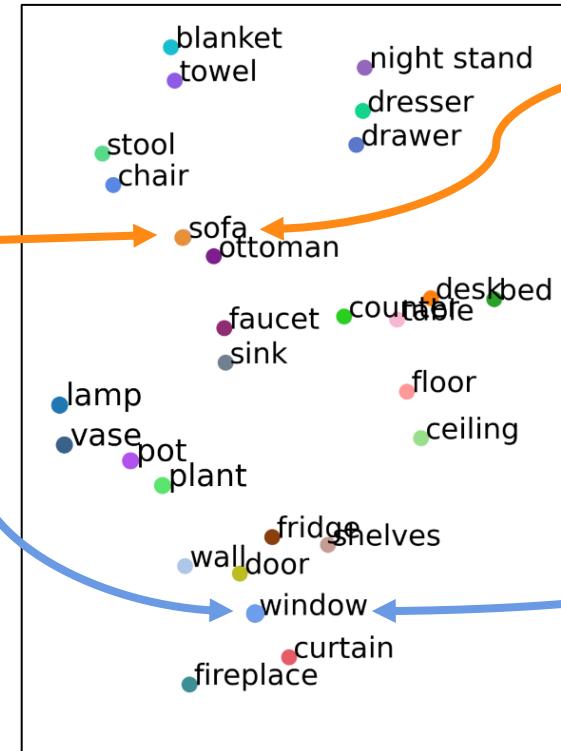
Leverage 2D foundation models

Key Idea

Co-embed 3D Features with CLIP Features



3D Geometry



RGB Images

Key Idea

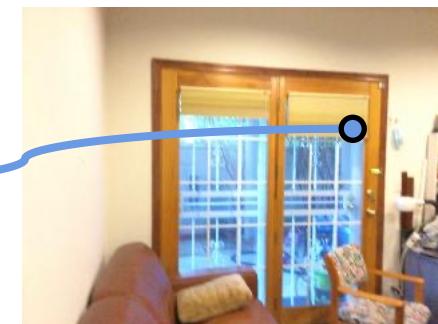
Co-embed 3D Features with CLIP Features



3D Geometry



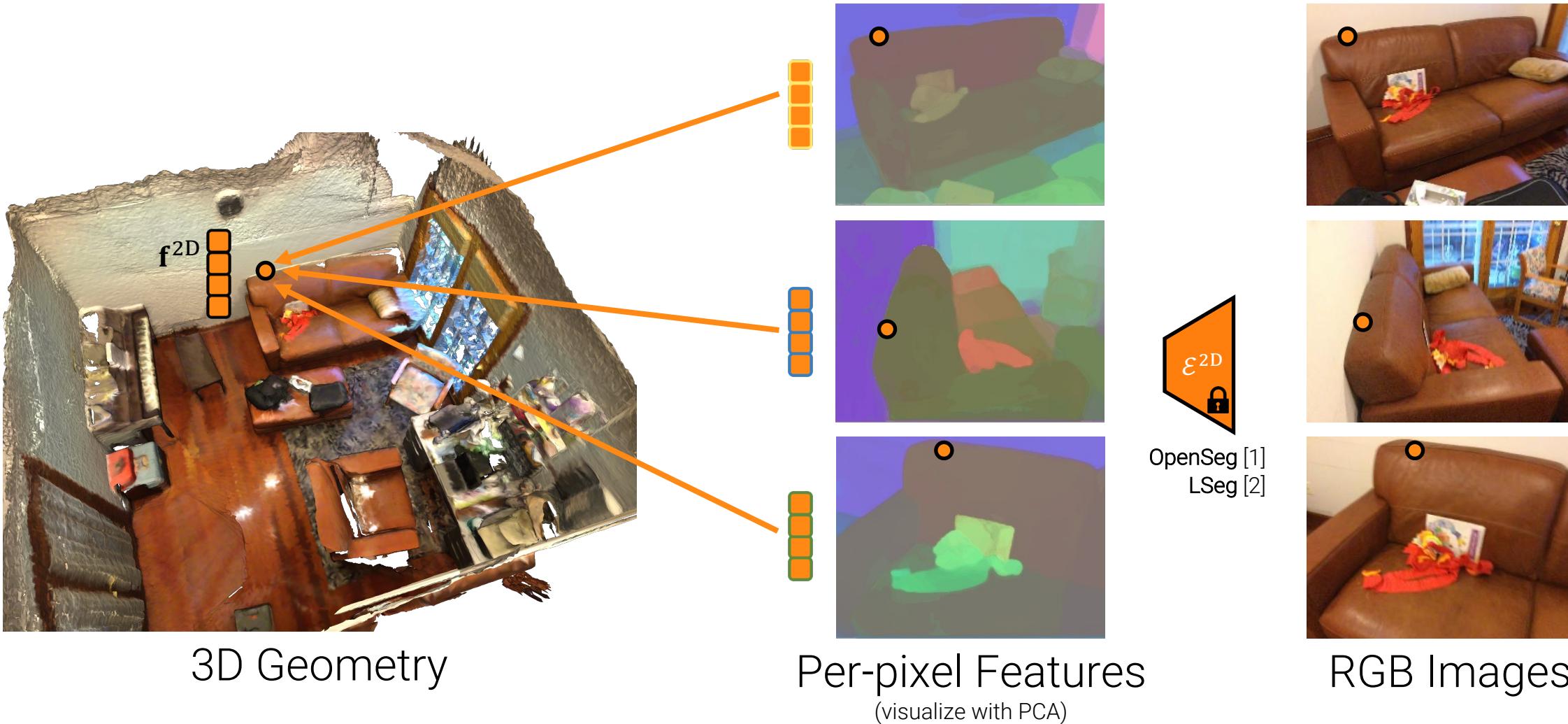
CLIP Text Features
(visualize with T-SNE)



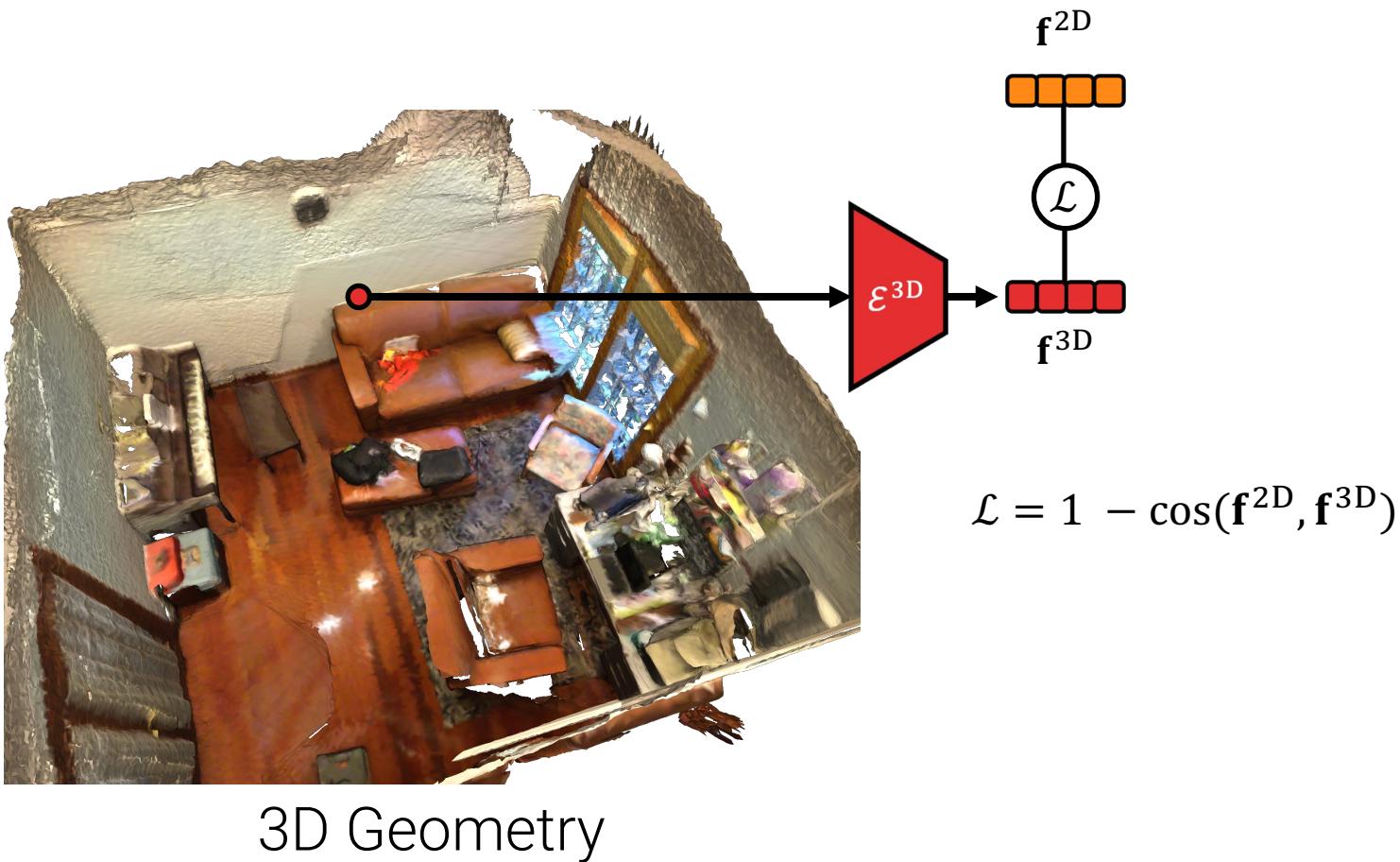
RGB Images

How to Learn Such **Text-Image-3D Co-Embeddings?**

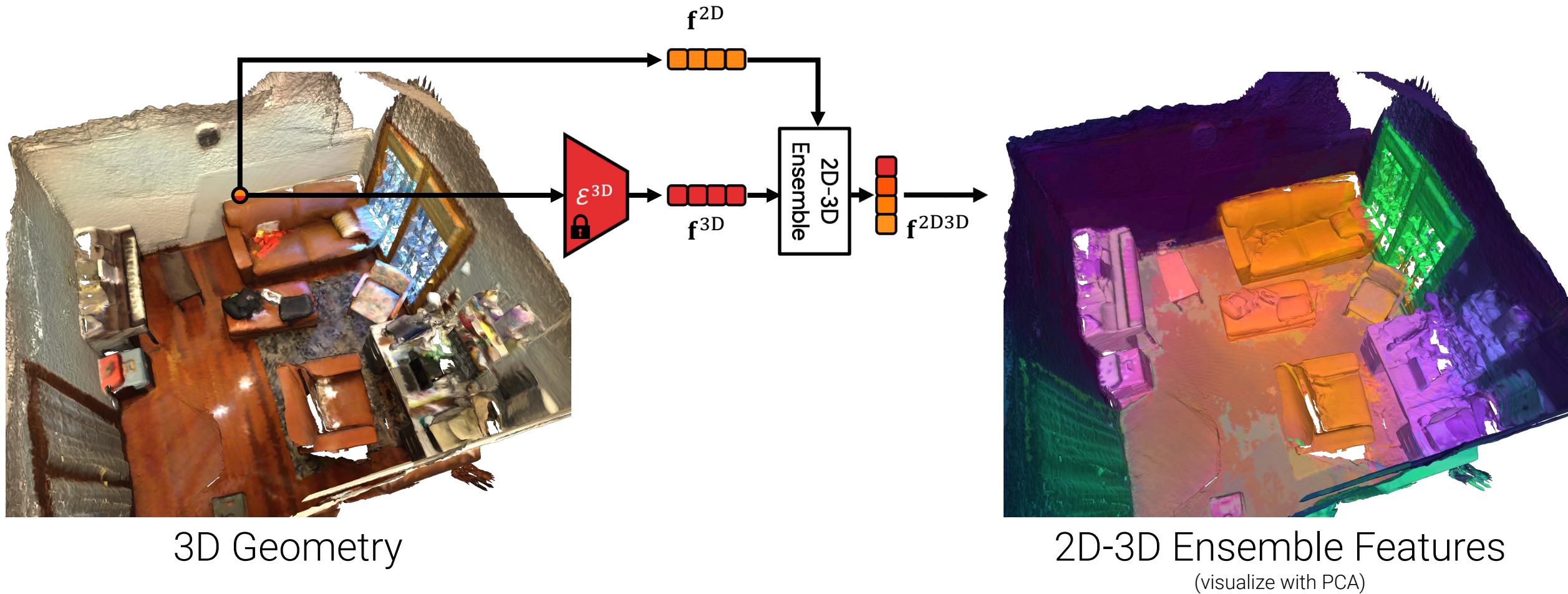
Step 1: Multi-view Feature Fusion



Step 2: 3D Feature Distillation



Inference: 2D-3D Ensemble



Open-Vocabulary, Zero-shot

3D Semantic Segmentation



Input 3D Geometry

■ wall ■ floor ■ cabinet ■ bed ■ chair ■ sofa ■ table ■ door ■ window ■ bookshelf ■ picture ■ counter ■ desk ■ curtain ■ refrigerator ■ shower curtain ■ toilet ■ sink ■ bathtub ■ other



Our Zero-shot 3D Segmentation
(20 classes)

■ wall ■ floor ■ cabinet ■ bed ■ chair ■ sofa ■ table ■ door ■ window ■ bookshelf ■ picture ■ counter ■ desk ■ curtain ■ refrigerator ■ shower curtain ■ toilet ■ sink ■ bathtub ■ other



Our Zero-shot 3D Segmentation
(160 classes)

wall	cabinet	bed	pot	bathub	dresser	stand	clock	tissue box	furniture	soap	cup	hanger	urn	paper towel dispenser	toy
door	curtain	night stand	desk	book	rug	drawer	stove	air vent	air conditioner	thermostat	ladder	candlestick	decorative plate	foot rest	
ceiling	floor	table	toilet	air vent	ottoman	container	washing machine	faucet	fire extinguisher	fire extinguisher	garage door	light	car	soap dish	
picture	plant	column	column	coffee table	photo	bottle	light switch	shower curtain	radiator	piano	scale	drum	computer	cleaner	
mirror	mirror	banister	counter	counter	refridgerator	refridgerator	purse	bin	curtain rod	paper towel	board	jacket	whiteboard	computer	
window	towel	stairs	bench	bench	bookshelf	bookshelf	fan	telephone	printer	sheet	rope	bottle of soap	water cooler	knob	
chair	sink	stool	garbage bin	garbage bin	wardrobe	wardrobe	fan	bucket	headboard	paper towel	ball	bag	drum	range hood	
pillow	shelves	vase	fireplace	railing	pipe	chandelier	railing	microwave	handle	sheet	excercise equipment	toilet paper holder	paper	candelabra	

Image-based 3D Scene Query



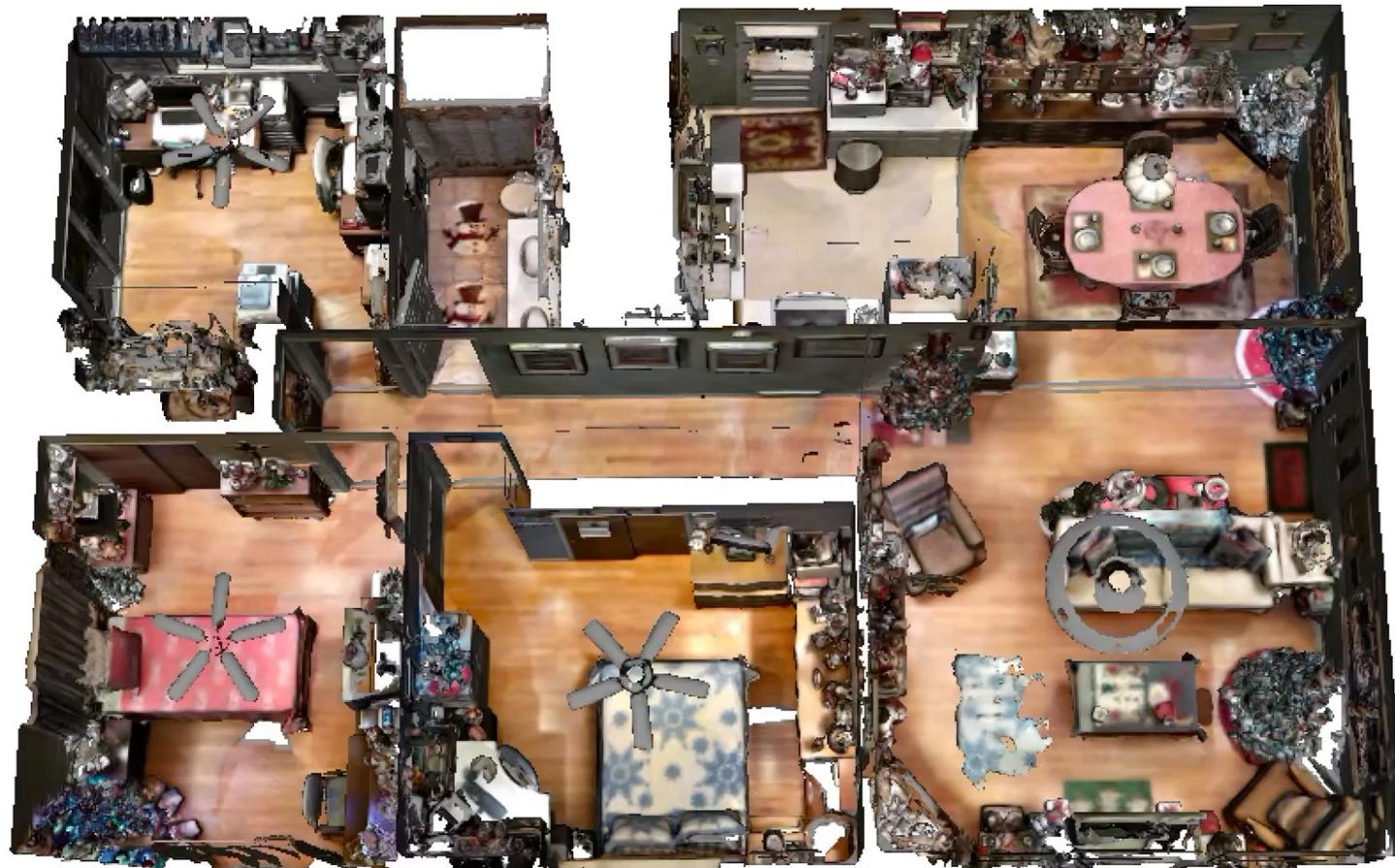
Image Queries

Given 3D Geometry

Interactive Demo

Open-vocabulary 3D Scene Exploration

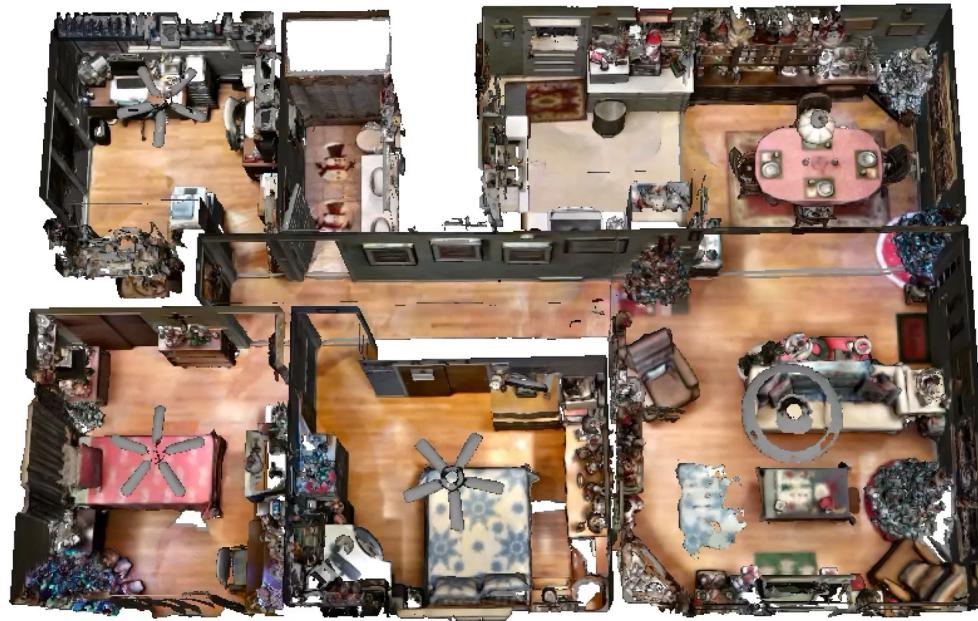
Text queries:



Take-home Messages

- + Open up a **wide range of applications** by leveraging large 2D vision-language models
- + Inspire future works to shift to open-vocabulary tasks
- Segmentation quality is quite limited

OpenScene



Accurately segment and understand 3D scenes is essential!

Motivation

Instance Segmentation Methods Requires 3D Manual Labels



Input 3D Scene



3D Semantic Instances

- ❗ **Expensive** and **challenging** to annotate 3D masks
- ❗ Perform poorly in scenes **out of training distribution**

Motivation

2D Foundation Model Rocks!



- ☀️ SAM exhibits **extraordinary ability to generalize**
- 🤔 Only applicable to **2D data**

How to obtain accurate 3D segmentation
without any manual 3D labels?

Leverage **2D foundation models**



清华大学
Tsinghua University

ETH zürich

Google



Microsoft

Segment3D

Learning Fine-Grained Class-Agnostic 3D Segmentation without Manual Labels



Rui
Huang



Songyou
Peng



Ayça
Takmaz



Federico
Tombari



Marc
Pollefeys



Shiji
Song



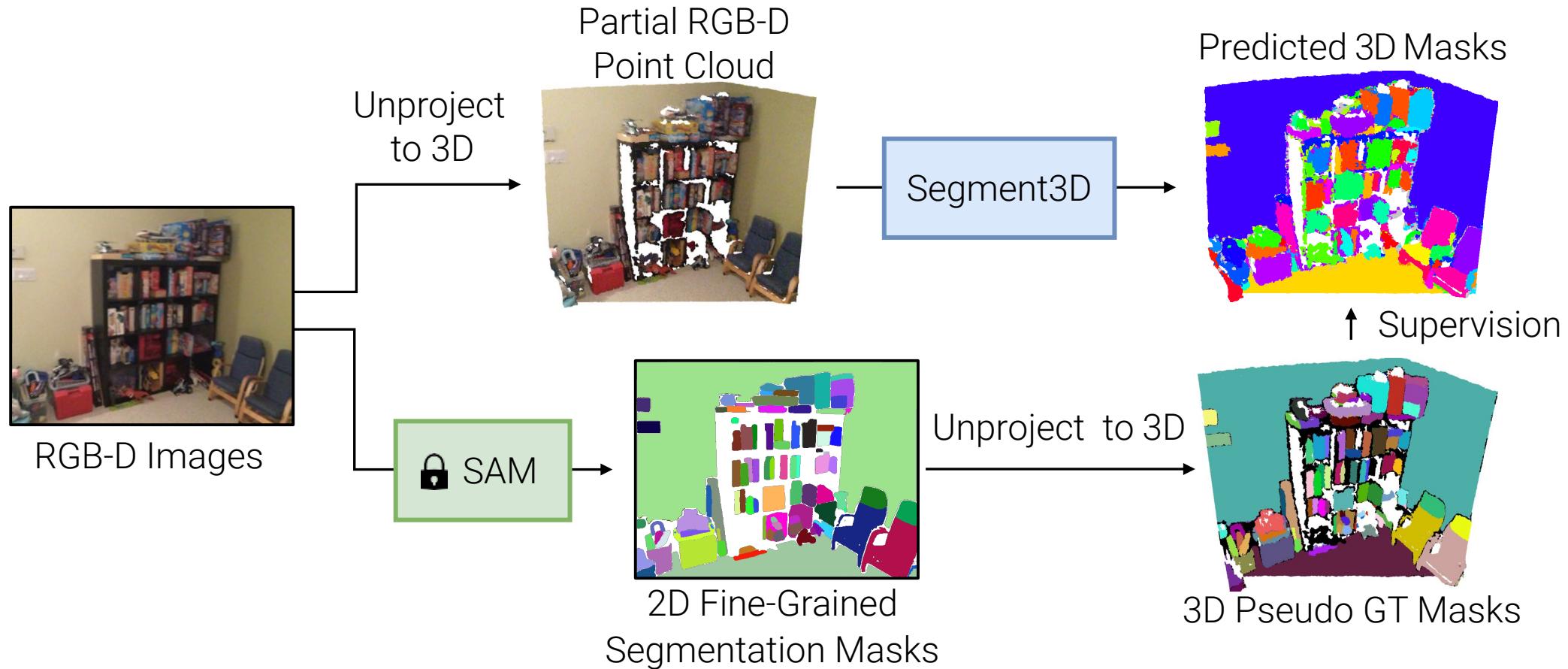
Gao
Huang



Francis
Engelmann

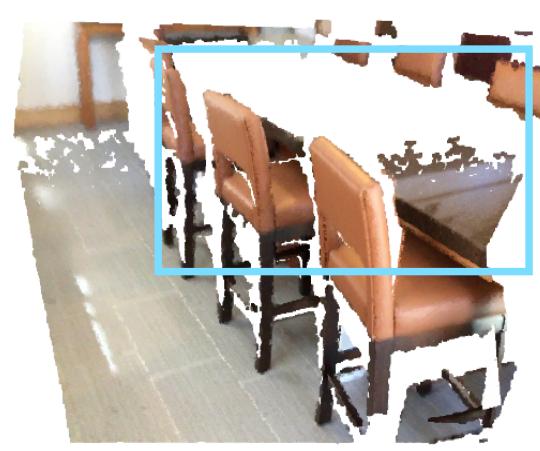
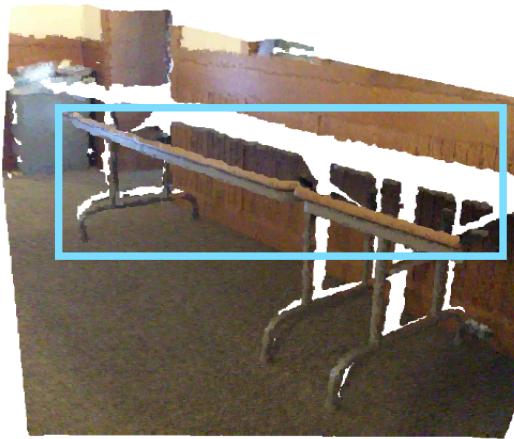
Segment3D

Stage 1: Pre-training on Partial Point Clouds



Segment3D

Domain Gap Between Partial and Full Point Clouds



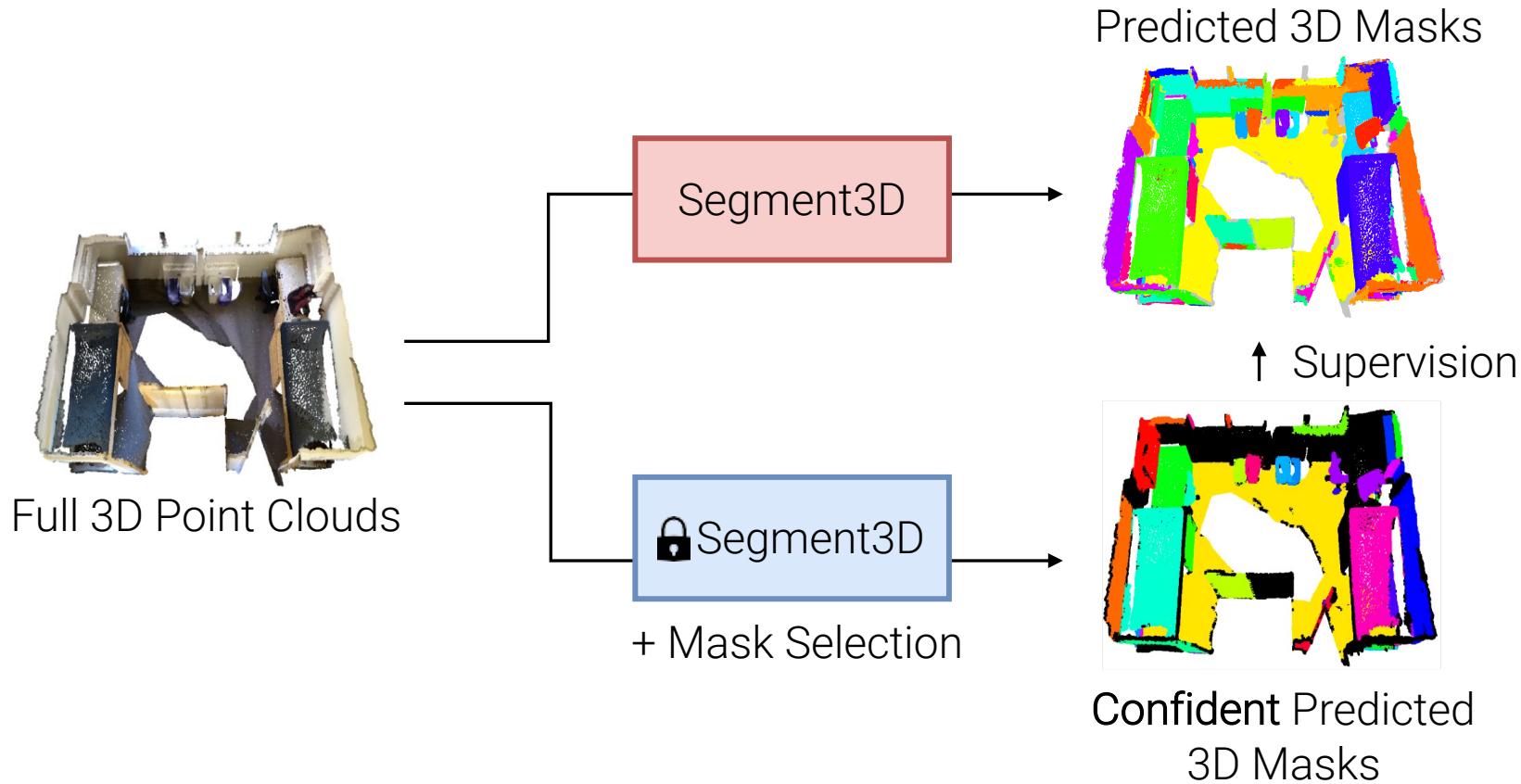
RGB Images

Partial Point Clouds

Full Point Clouds

Segment3D

Stage 2: Fine-tune on Full Point Clouds

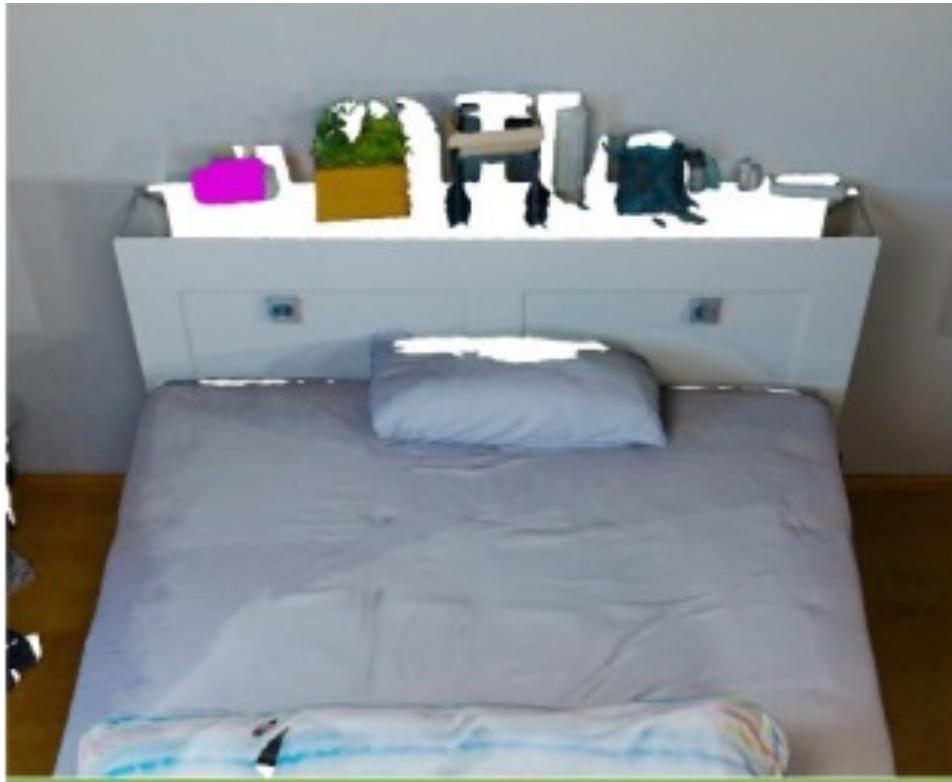


No manual labels are needed at all!

Results

Class-Agnostic 3D Segmentation

ScanNet++ Validation Set



Input Point Clouds



Mask3D [1]

Class-Agnostic 3D Segmentation

ScanNet++ Validation Set



Input Point Clouds



Segment3D (Ours)

Class-Agnostic 3D Segmentation

ScanNet++ Validation Set



Segment3D (Ours)



GT

Class-Agnostic 3D Segmentation

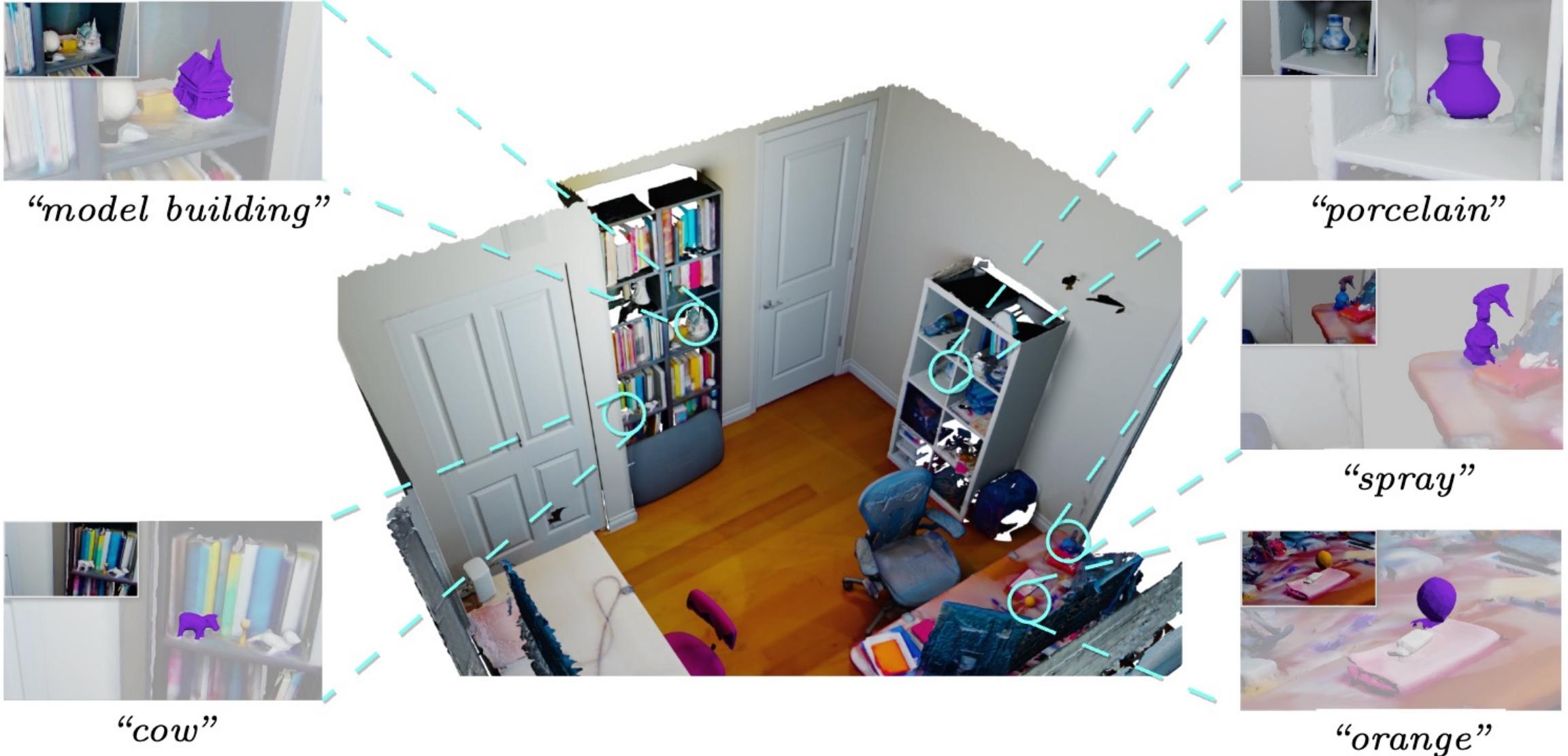
ScanNet++ Validation Set

Model	Avg. Inference Time/s	without post-processing			with post-processing		
		AP	AP ₅₀	AP ₂₅	AP	AP ₅₀	AP ₂₅
Mask3D [41]	0.7	8.7	15.5	27.2	14.3	21.3	29.9
SAM3D [50]	386.7	3.9	9.3	22.1	8.4	16.1	30.0
Felzenszwalb <i>et al.</i> [16]	12.6	5.8	11.6	27.2	—	—	—
Segment3D (Ours)	0.7	13.0	23.8	38.3	20.2	30.9	42.7

Effect of Two-Stage Training

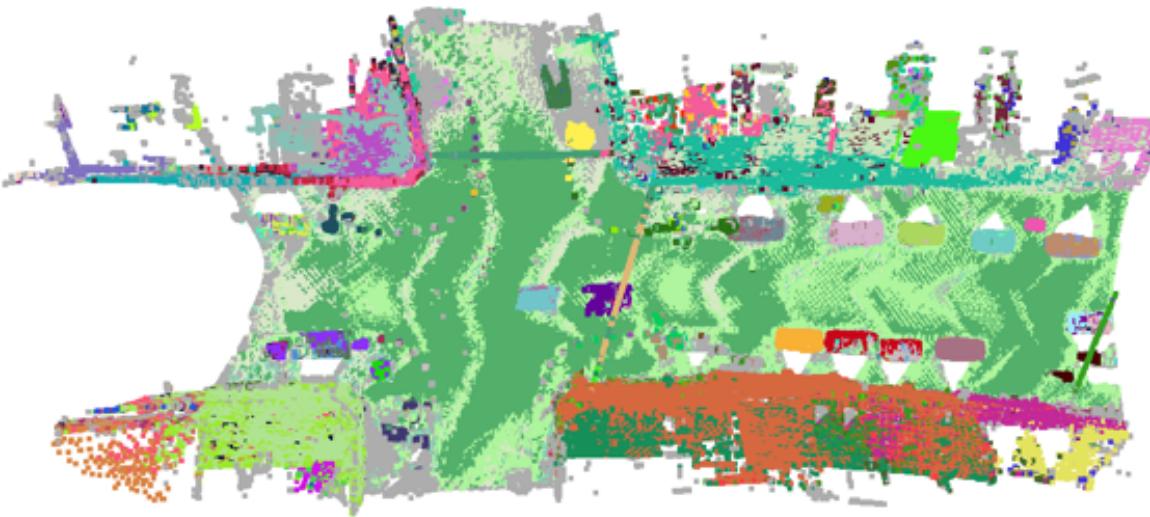
Training Stages	AP	AP ₅₀	AP ₂₅
Pre-Training (Stage 1)	6.8	14.2	30.6
+ Fine-Tuning (Stage 2)	13.0 (+ 6.2)	23.8 (+ 9.6)	38.3 (+ 7.7)

Open-Vocabulary Segmentation



Outdoor Scenes In-the-Wild

Mask3D



Segment3D



Take-home Messages

- No 3D manual labels are used for training at all!
- **2D foundation model** (SAM) rocks!

Future work

- Unify as a single-stage pipeline
- Single pipeline for open-vocabulary 3D segmentation

2D Magic in a 3D World

2D Foundation Models for 3D Vision Tasks

3D Reconstruction



NeRF On-the-go
CVPR 2024

3D Scene Understanding



OpenScene
CVPR 2023

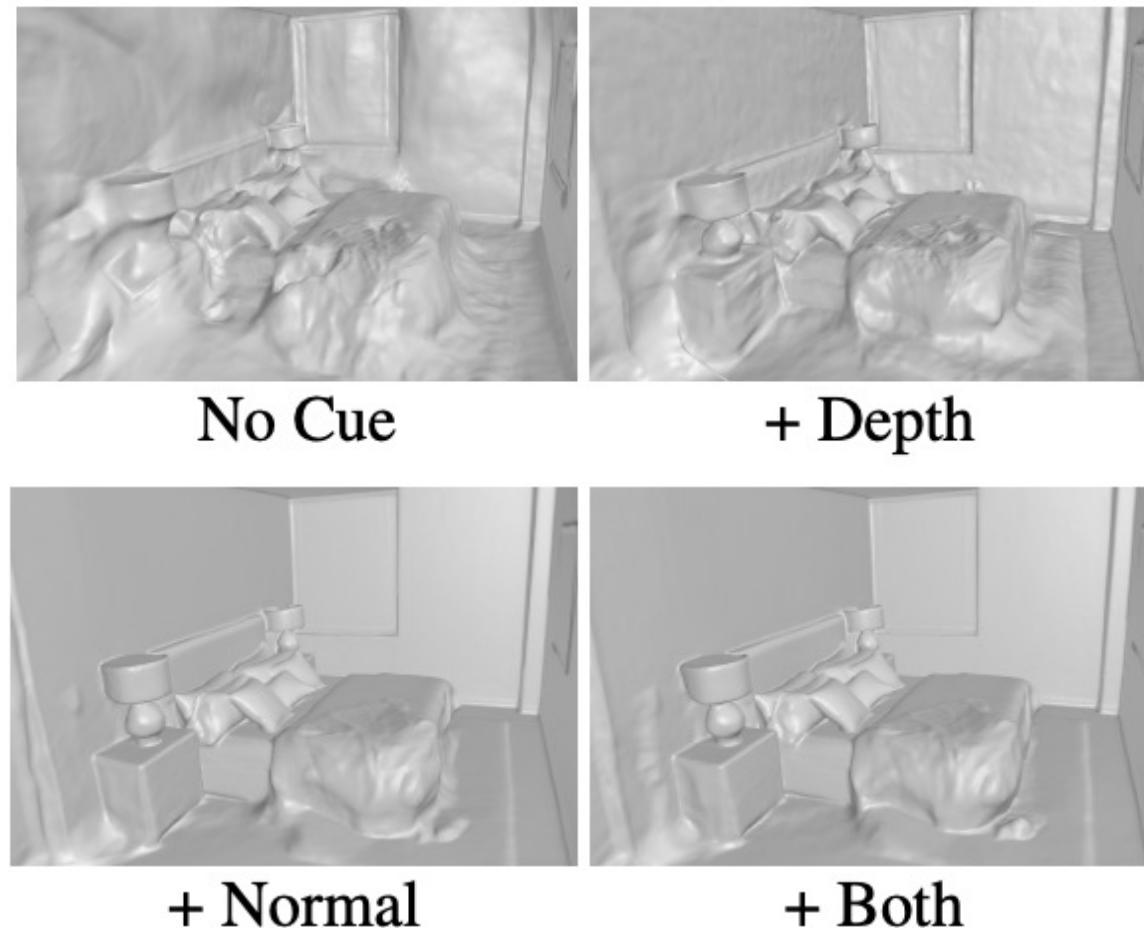
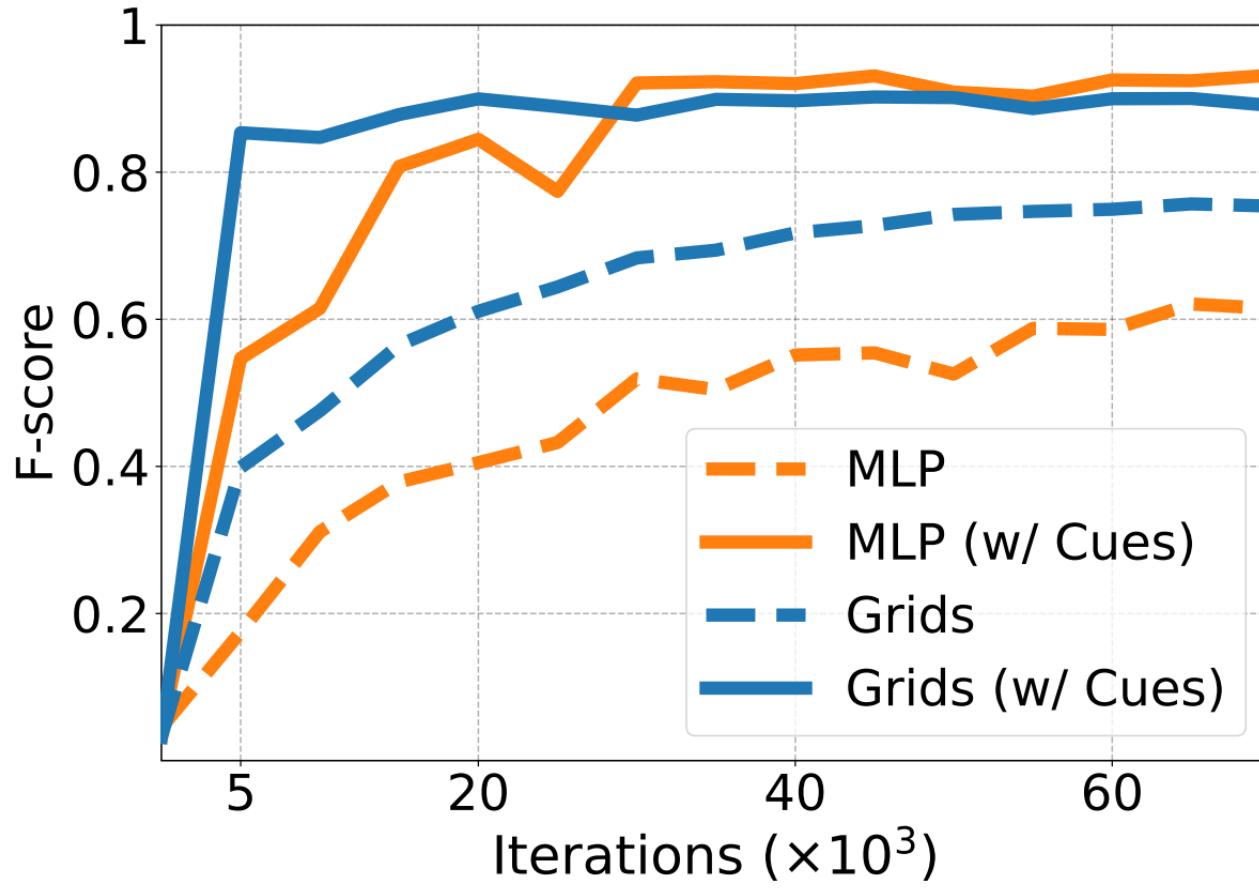


Segment3D
arXiv 2024

This talk focuses on how to **leverage**
2D foundation models for 3D tasks

So, what is next?

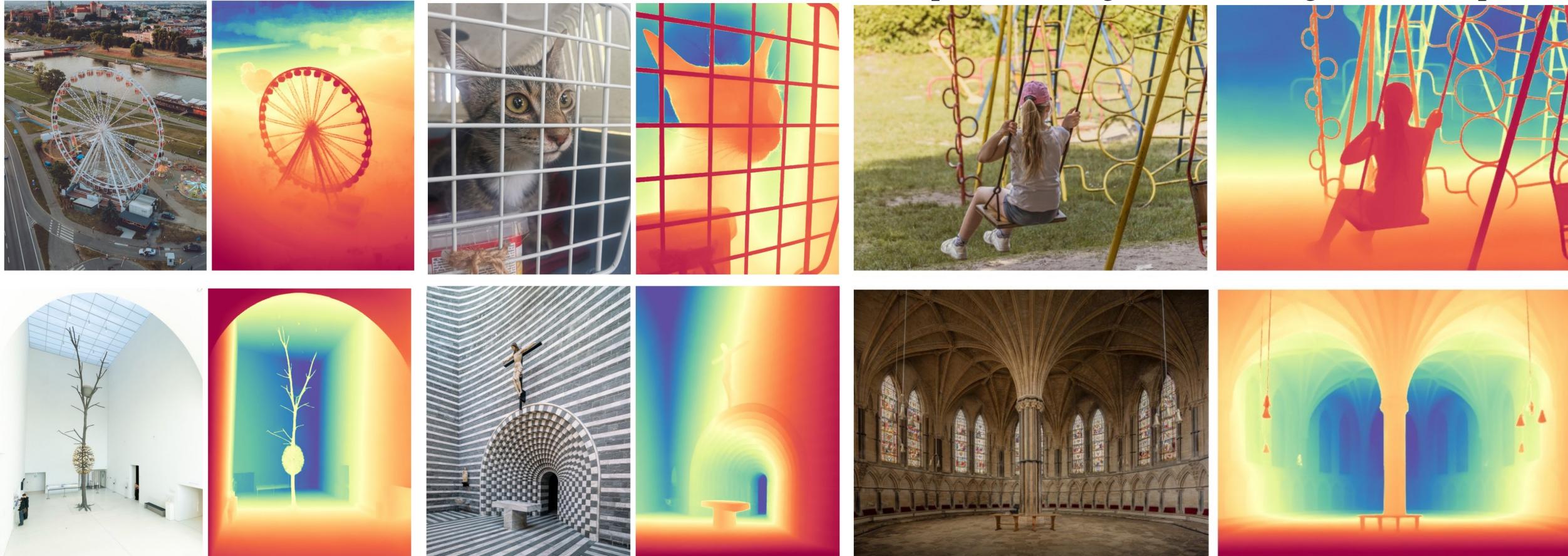
Monocular Cues Help!



Performance is bounded by the monocular cues!

Current Interests

Next-Generation Monocular Predictor



Marigold: Stable Diffusion-Based Monocular Predictor

Current Interests

Next-Generation Monocular Predictor

- Other **modalities** like surface normals, uncertainty, etc...
- Inference **speed**
- **Video depth** predictor (temporal-consistent)
- **Metric depths**
- ...

This talk focuses on how to **leverage**
2D foundation models for 3D tasks

Current Interests

3D Foundation Models



DreamFusion

[ICLR'23]

1.5 Hours

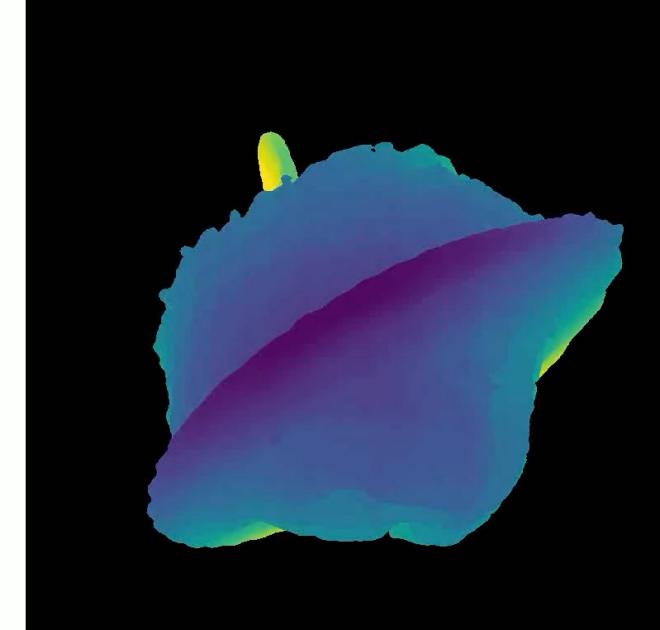
So far, only **object-level** 3D foundation models



Large Reconstruction Models

[ICLR'24]

5 Seconds!



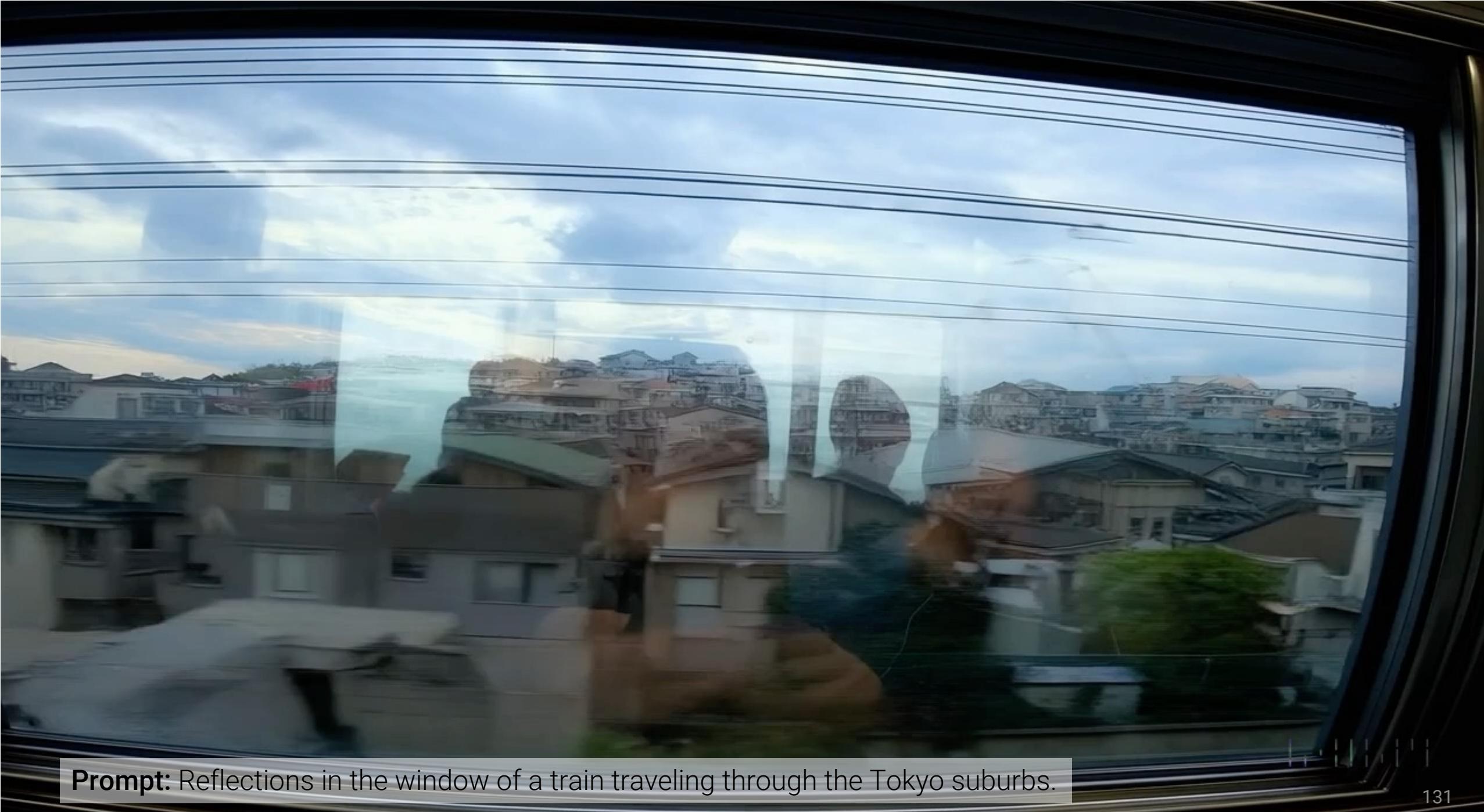
Current Interests

3D Foundation Models

- How to leverage internet-level in-the-wild data?
- How to train in a *smarter* way?
 - Multi-view generator + sparse-view reconstruction?
 - How to better leverage geometry under the data-driven pipeline?

Then... Sora is here





Prompt: Reflections in the window of a train traveling through the Tokyo suburbs.



Prompt: Prompt: The camera follows behind a white vintage SUV with a black roof rack as it speeds up a steep dirt road surrounded by pine trees on a steep mountain slope, dust kicks up from it's tires,.....

Open Questions...

- Do we really need a 3D foundation model?
- How to inject 3D to help video foundation models?
- How to leverage Sora for 3D generation?
- How to leverage Sora to model dynamic+static scenes?
- Controllability/Local editability in generative models
- In this era, how to *survive* as a PhD student in universities?

2D Magic in a 3D World

Songyou Peng

pengsongyou.github.io

Q&A