

2D Magic in a 3D World

Songyou Peng

The University of Hong Kong

Feb 22, 2024

Who Am I?

- Senior Researcher 
- Incoming Research Scientist 
- Earned my PhD
 - Marc Pollefeys 
 - Andreas Geiger 
- Internships during PhD
 - 2021: Michael Zollhoefer 
 - 2022: Tom Funkhouser 



pengsongyou.github.io

Research Overview of My PhD

Learn to Reconstruct and Understand 3D World



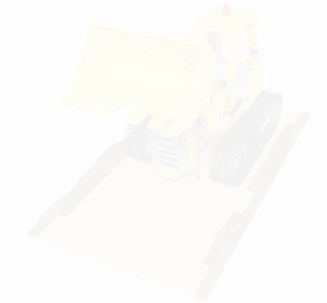
ConvOccNet
ECCV 2020 (Spotlight)



MonoSDF
NeurIPS 2022



Shape As Points
NeurIPS 2021 (Oral)



runs now at 50 fps on a GTX 1080 Ti

KiloNeRF
ICCV 2021



NICE-SLAM
CVPR 2022



NICER-SLAM
3DV 2024 (Oral)



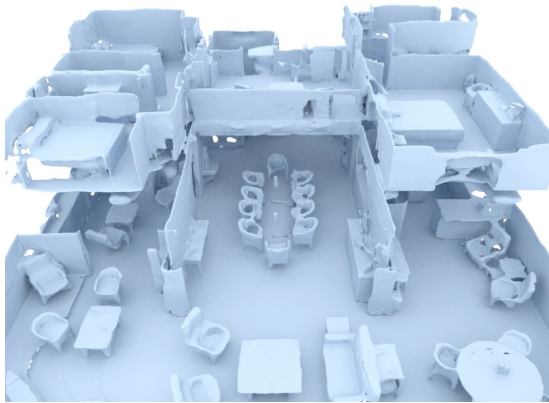
UNISURF
ICCV 2021 (Oral)



OpenScene
CVPR 2023 ₂

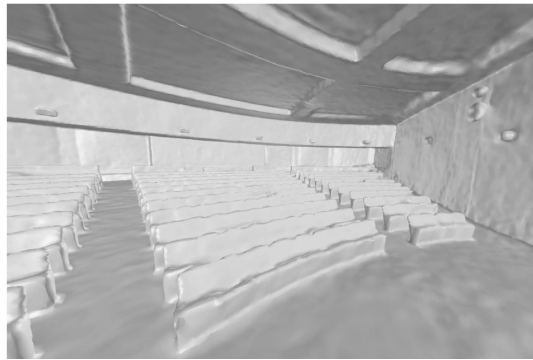
Research Overview of My PhD

Learn to Reconstruct and Understand 3D World



ConvOccNet

ECCV 2020 (Spotlight)

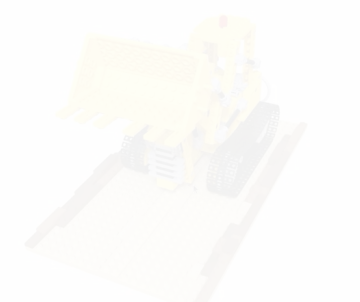


MonoSDF

NeurIPS 2022



Shape As Points
NeurIPS 2021 (Oral)



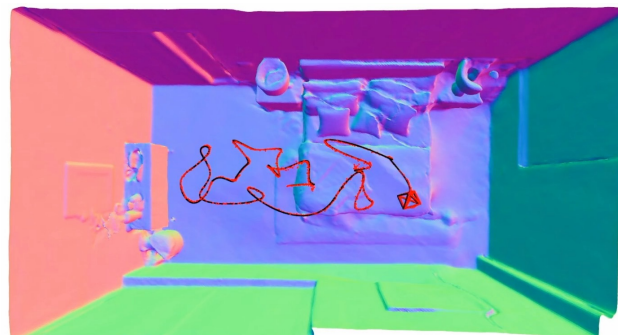
runs now at 50 fps on a GTX 1080 Ti

KiloNeRF
ICCV 2021



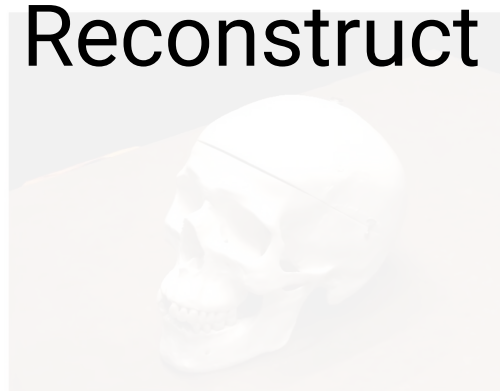
NICE-SLAM

CVPR 2022

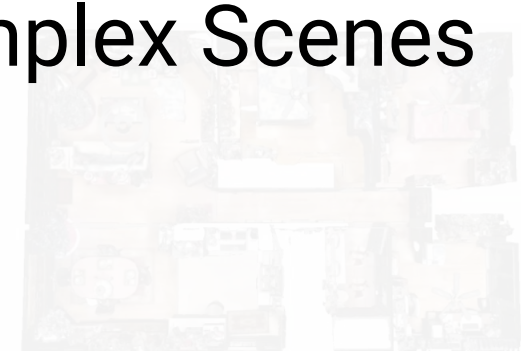


NICER-SLAM

3DV 2024 (Oral)



UNISURF
ICCV 2021 (Oral)



OpenScene
CVPR 2023

Topic #1:
Reconstruct Complex Scenes

Research Overview of My PhD

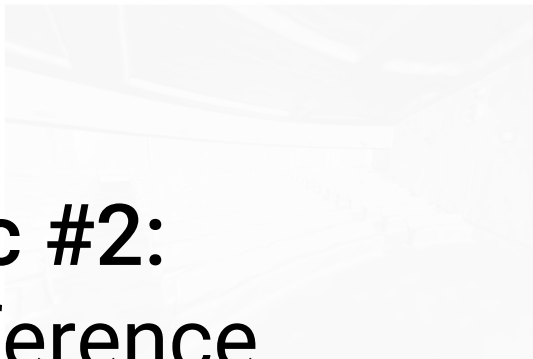
Learn to Reconstruct and Understand 3D World

Topic #2: Fast Inference



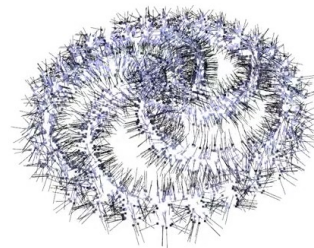
ConvOccNet

ECCV 2020 (Spotlight)



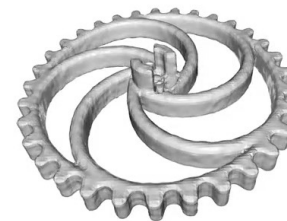
MonoSDF

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

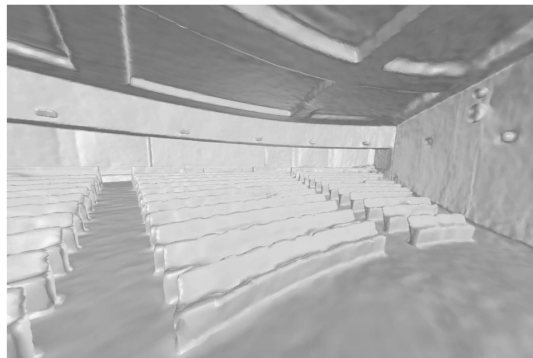
CVPR 2023 ⁴

Research Overview of My PhD

Learn to Reconstruct and Understand 3D World



ConvOccNet
ECCV 2020 (Spotlight)



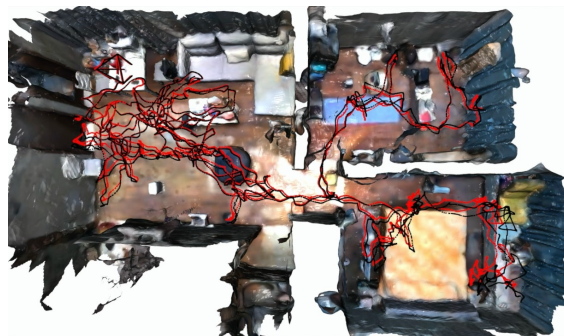
MonoSDF
NeurIPS 2022

Topic #3:
Reconstruct from 2D Observations

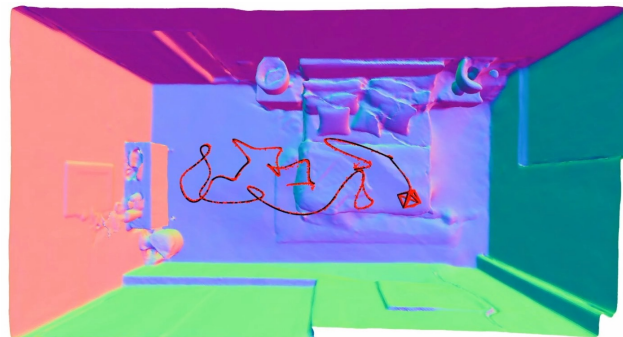
Shape As Points
NeurIPS 2021 (Oral)



KiloNeRF
ICCV 2021



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)



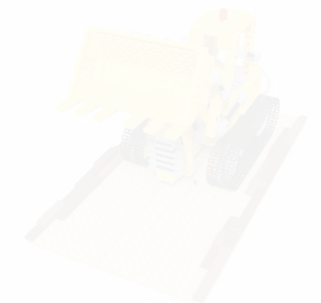
MonoSDF

NeurIPS 2022



Shape As Points

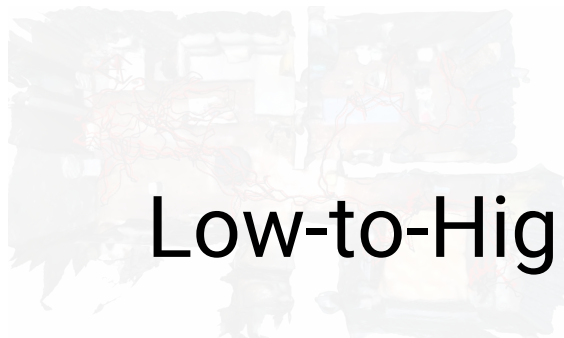
NeurIPS 2021 (Oral)



runs now at 50 fps on a GTX 1080 Ti

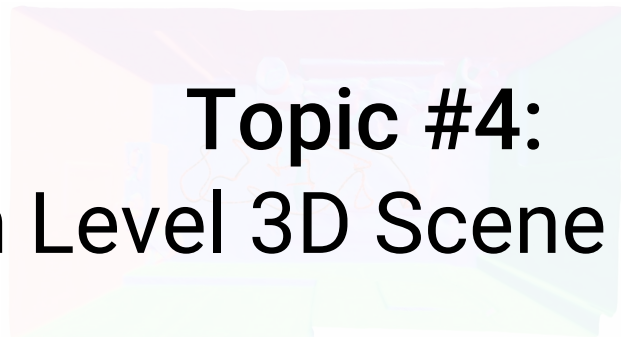
KiloNeRF

ICCV 2021



NICE-SLAM

CVPR 2022

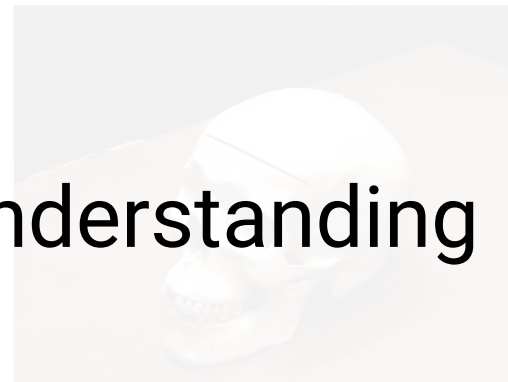


Topic #4:

Low-to-High Level 3D Scene Understanding

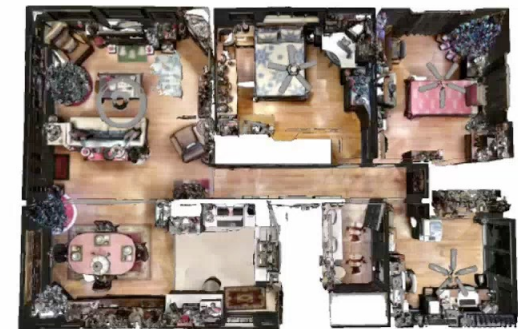
NICER-SLAM

3DV 2024 (Oral)



UNISURF

ICCV 2021 (Oral)

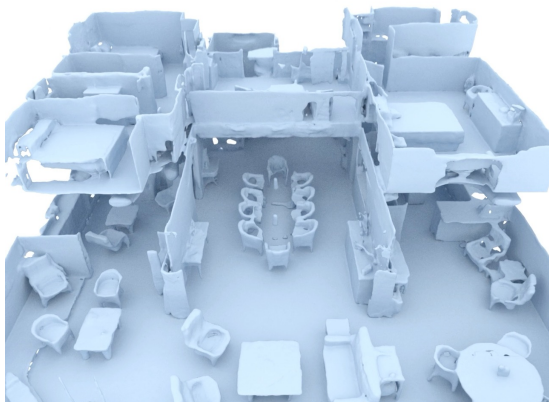


OpenScene

CVPR 2023 6

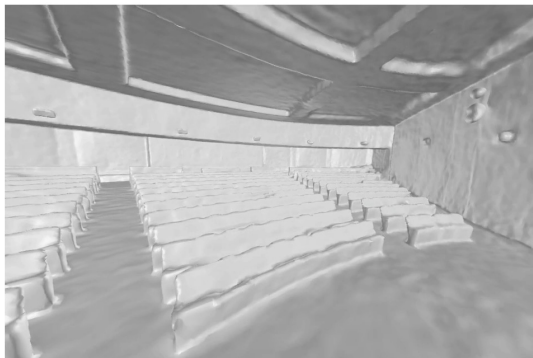
Research Overview of My PhD

Learn to Reconstruct and Understand 3D World



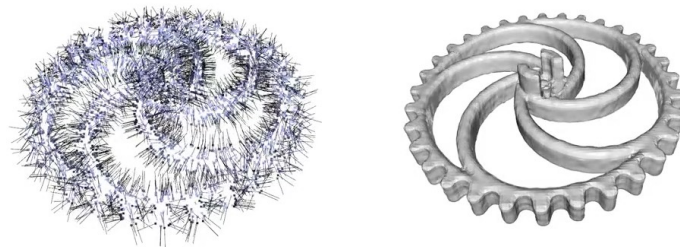
ConvOccNet

ECCV 2020 (Spotlight)



MonoSDF

NeurIPS 2022



Shape As Points

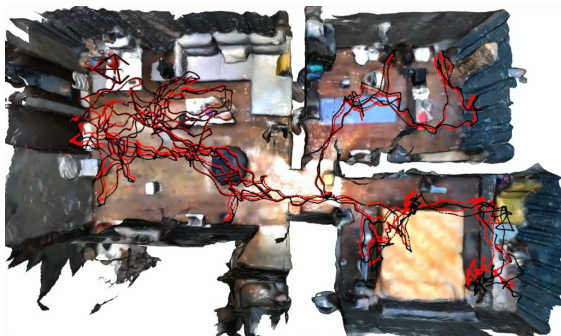
NeurIPS 2021 (Oral)



runs now at 50 fps on a GTX 1080 Ti

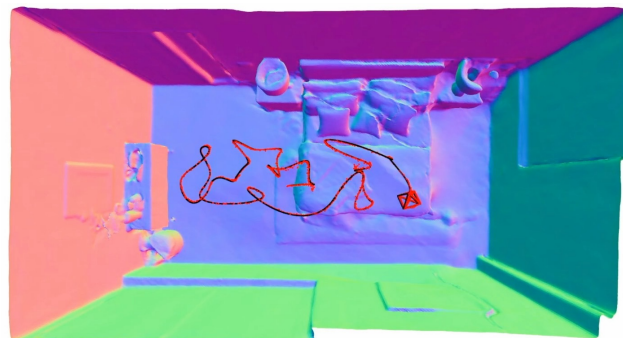
KiloNeRF

ICCV 2021



NICE-SLAM

CVPR 2022



NICER-SLAM

3DV 2024 (Oral)



UNISURF

ICCV 2021 (Oral)

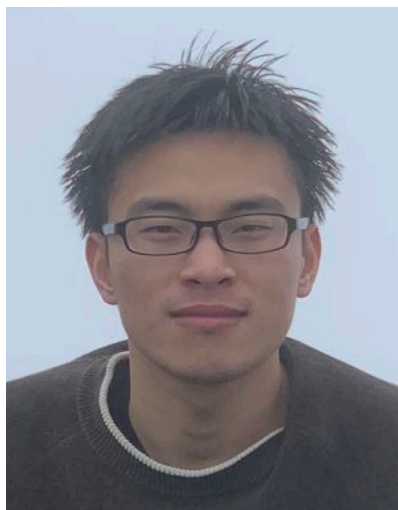


OpenScene

CVPR 2023 7



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

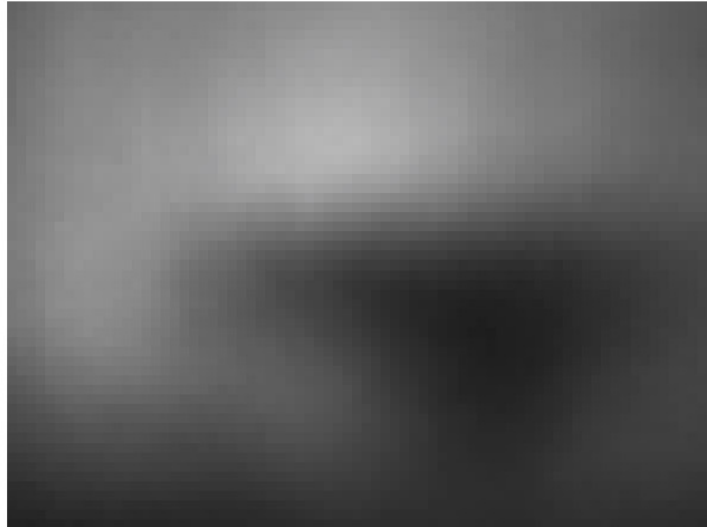
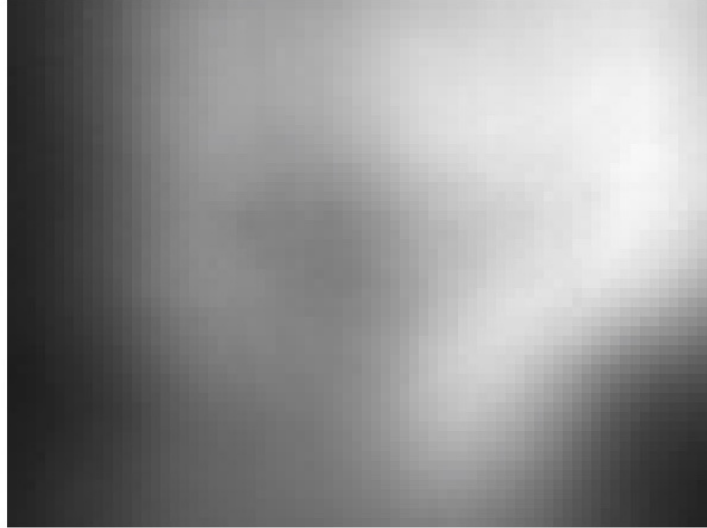
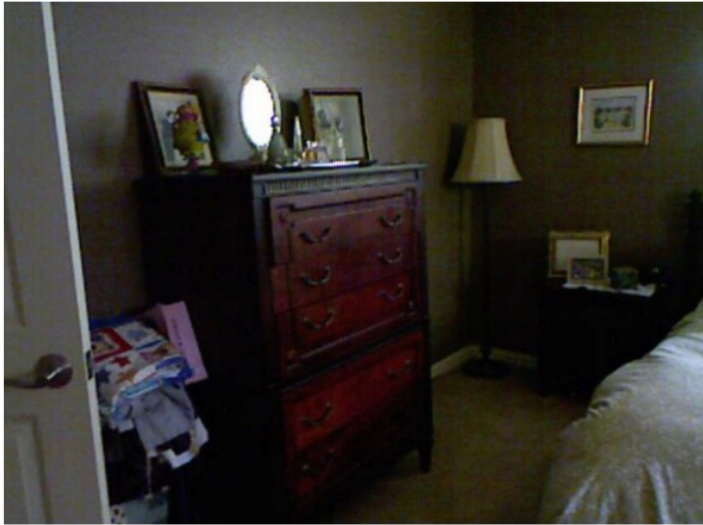


Torsten Sattler

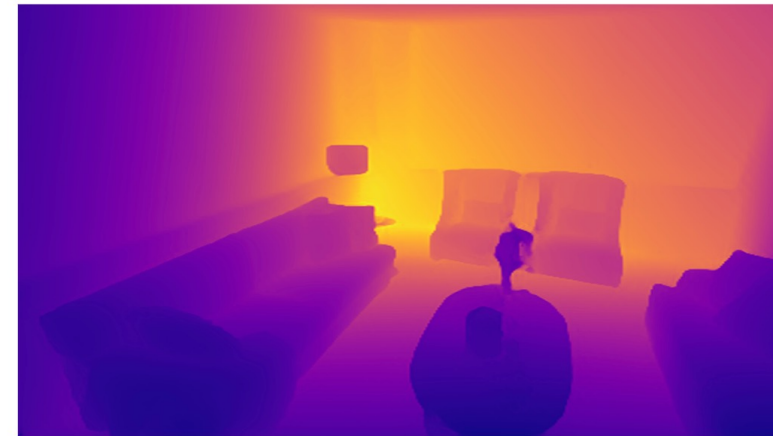
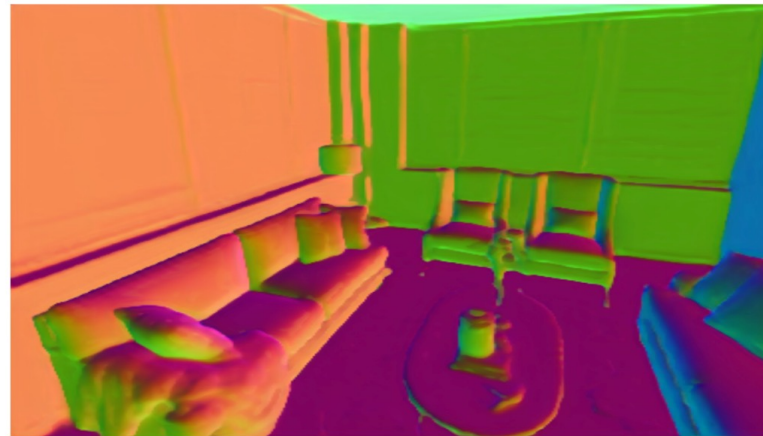
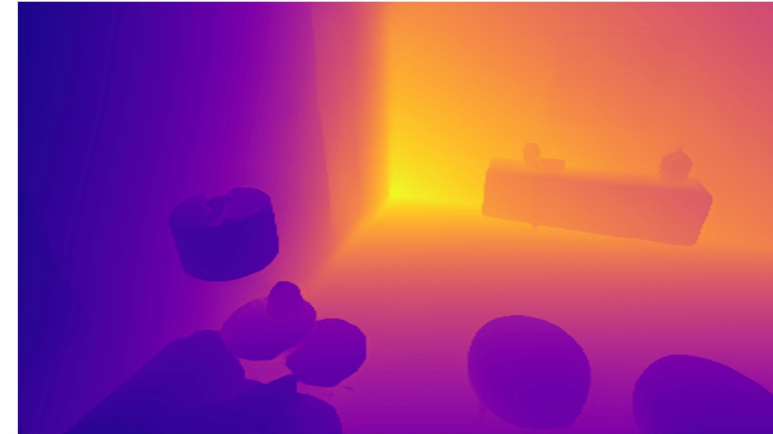
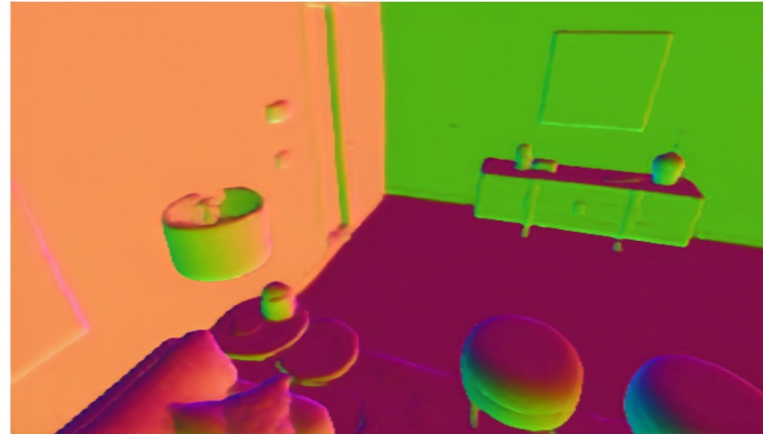


Andreas Geiger

Depth Prediction from a Single Image



Omnidata



RGB Image

Omnidata Normal

Omnidata Depth

2D Magic in a 3D World

Songyou Peng

The University of Hong Kong

Feb 22, 2024

2D Magic in a 3D World

2D Magic in a 3D World

2D Monocular Cues Benefit 3D Reconstruction

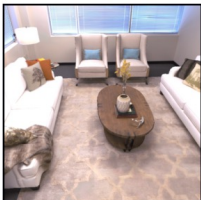
3D Reconstruction Pipeline



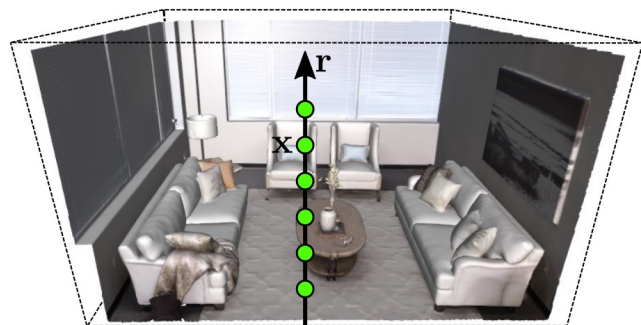
3D Reconstruction Pipeline



Input Views



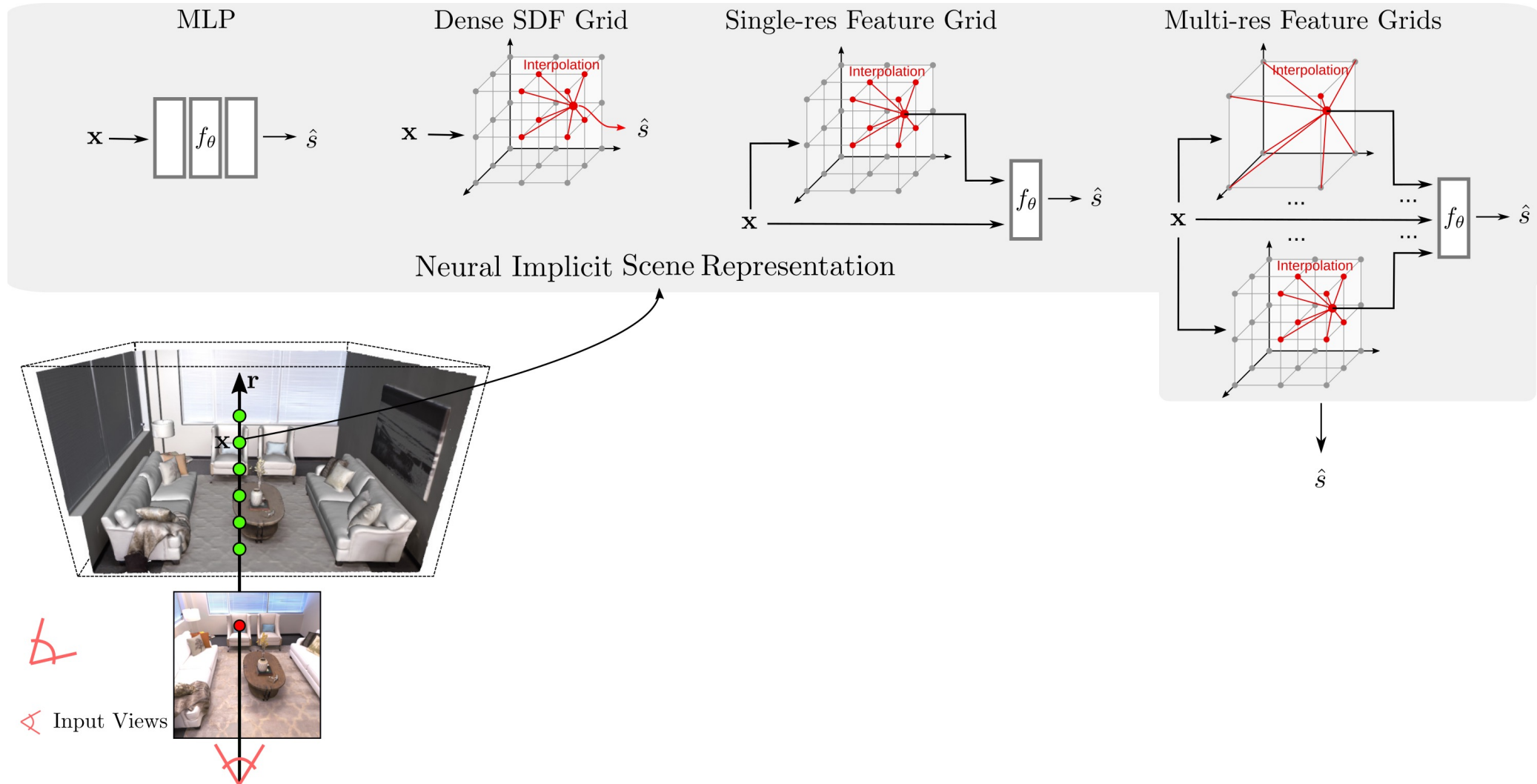
3D Reconstruction Pipeline



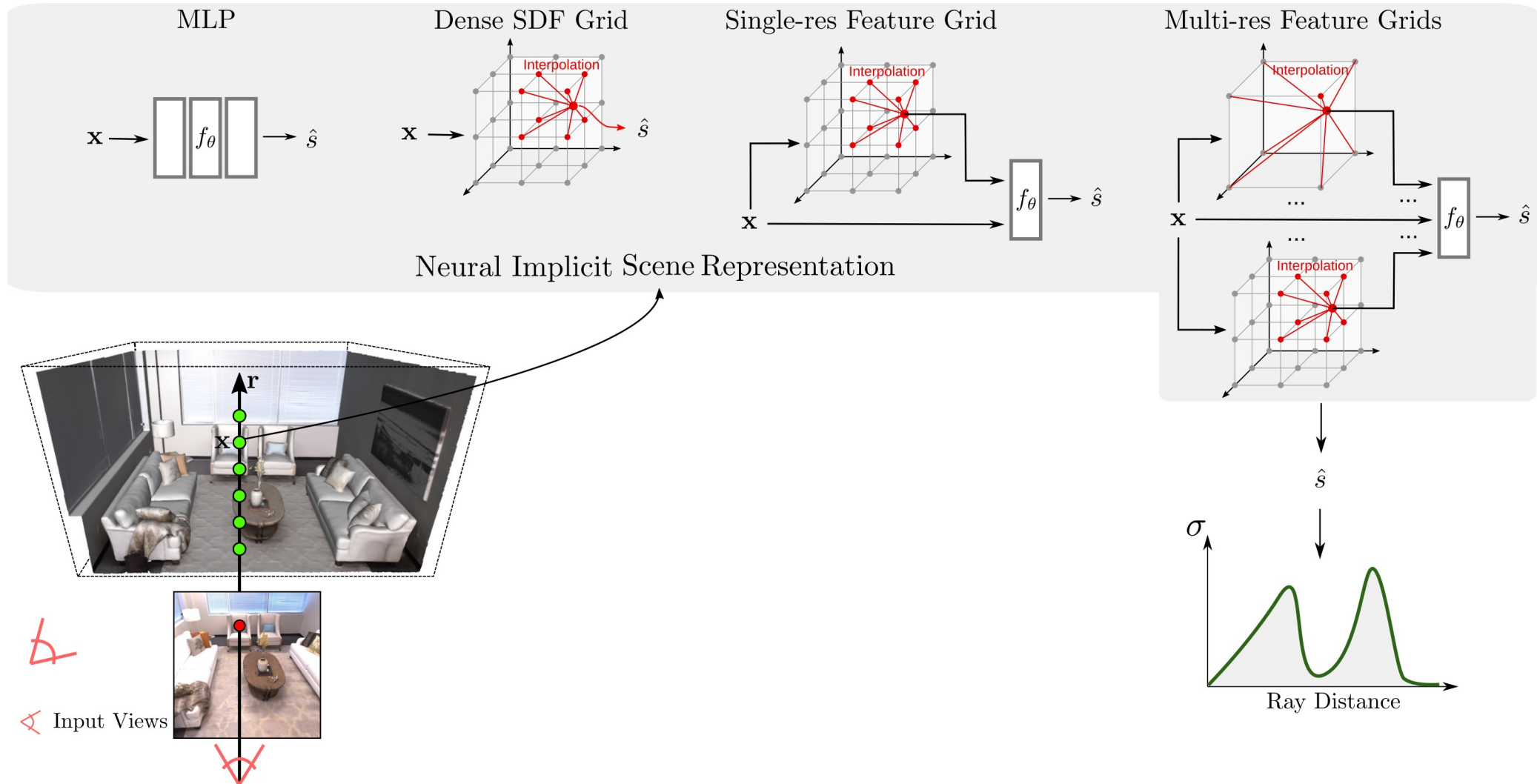
Input Views



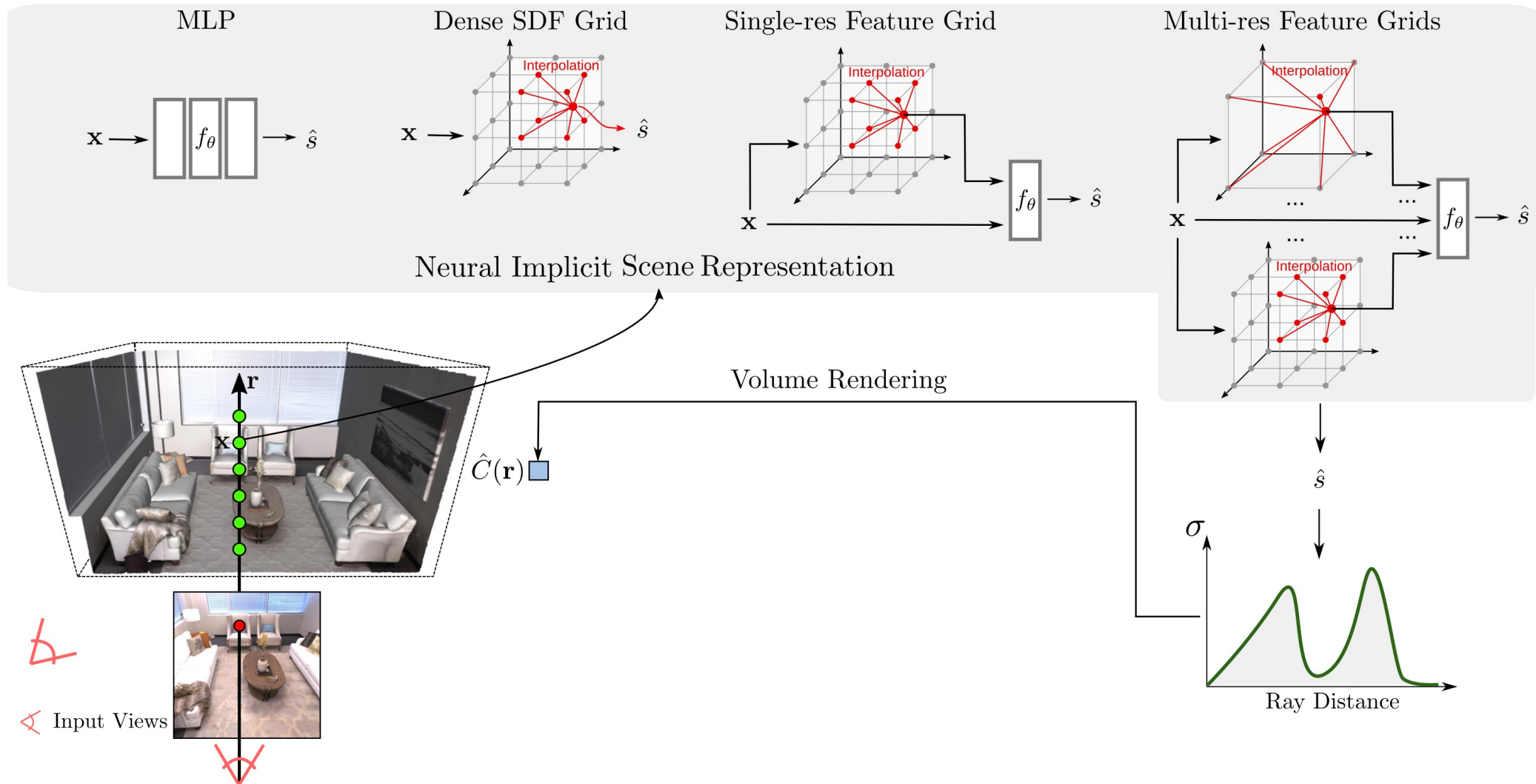
3D Reconstruction Pipeline



3D Reconstruction Pipeline

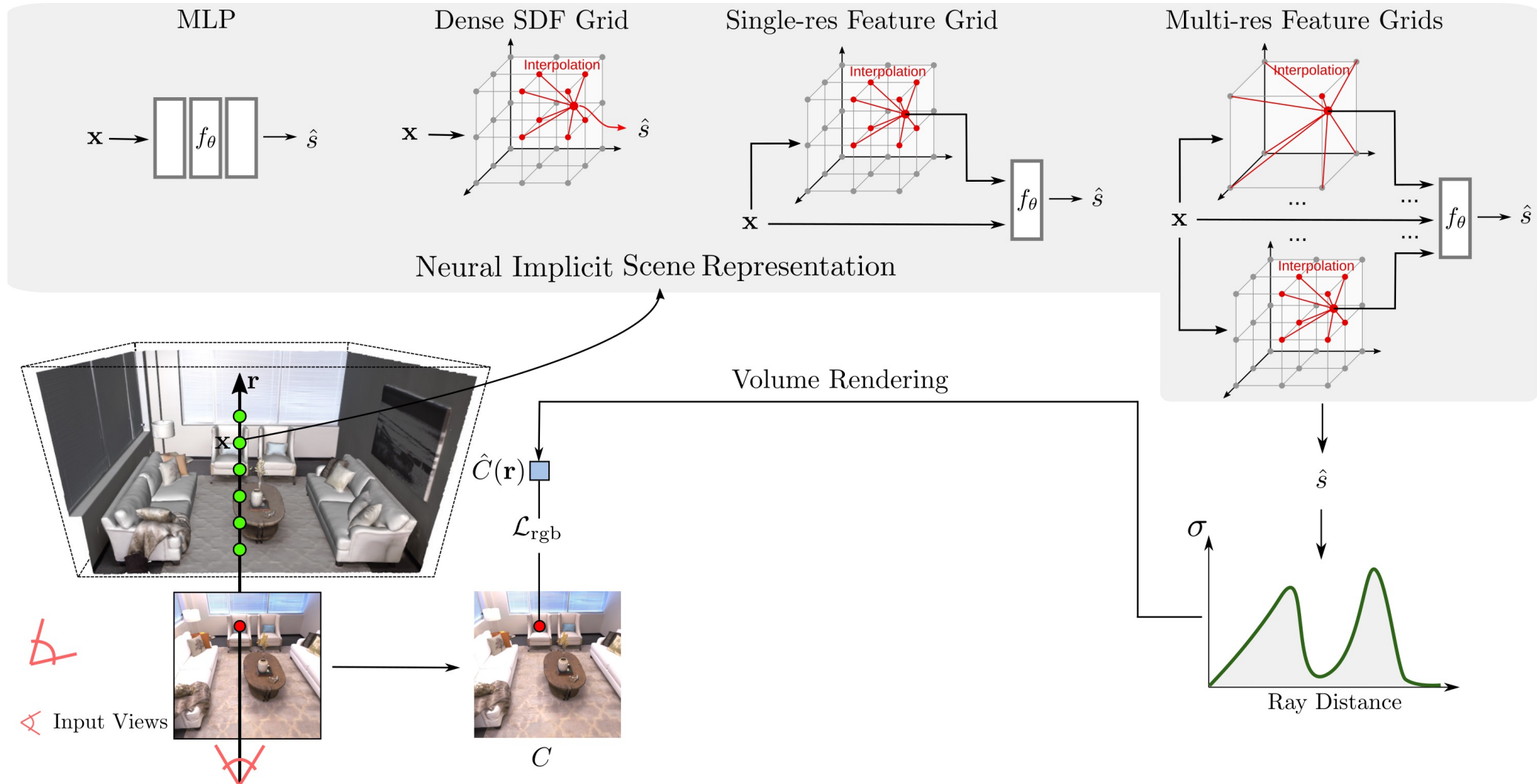


3D Reconstruction Pipeline



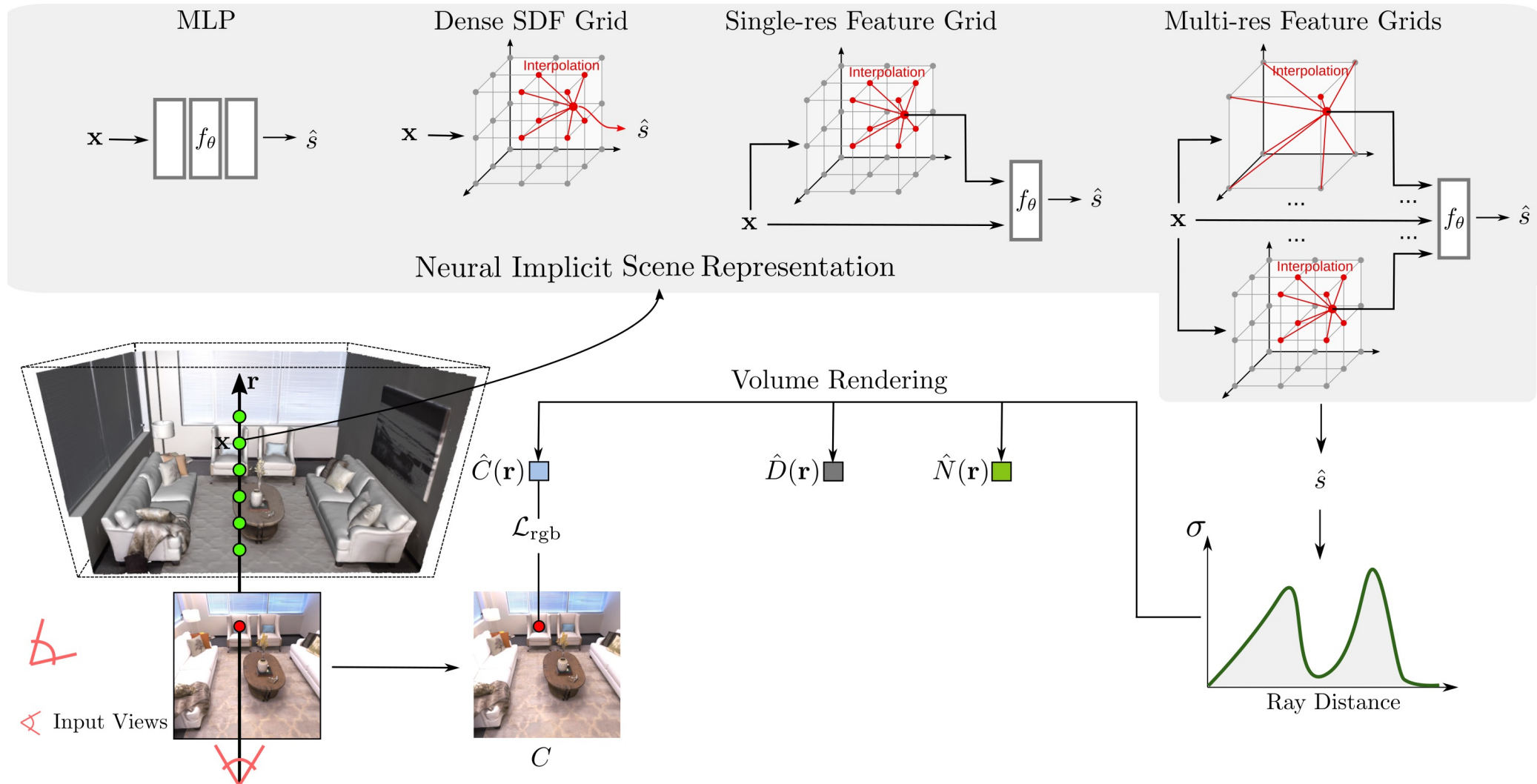
3D Reconstruction Pipeline

VolSDF/NeuS/UNISURF/Neuralangelo



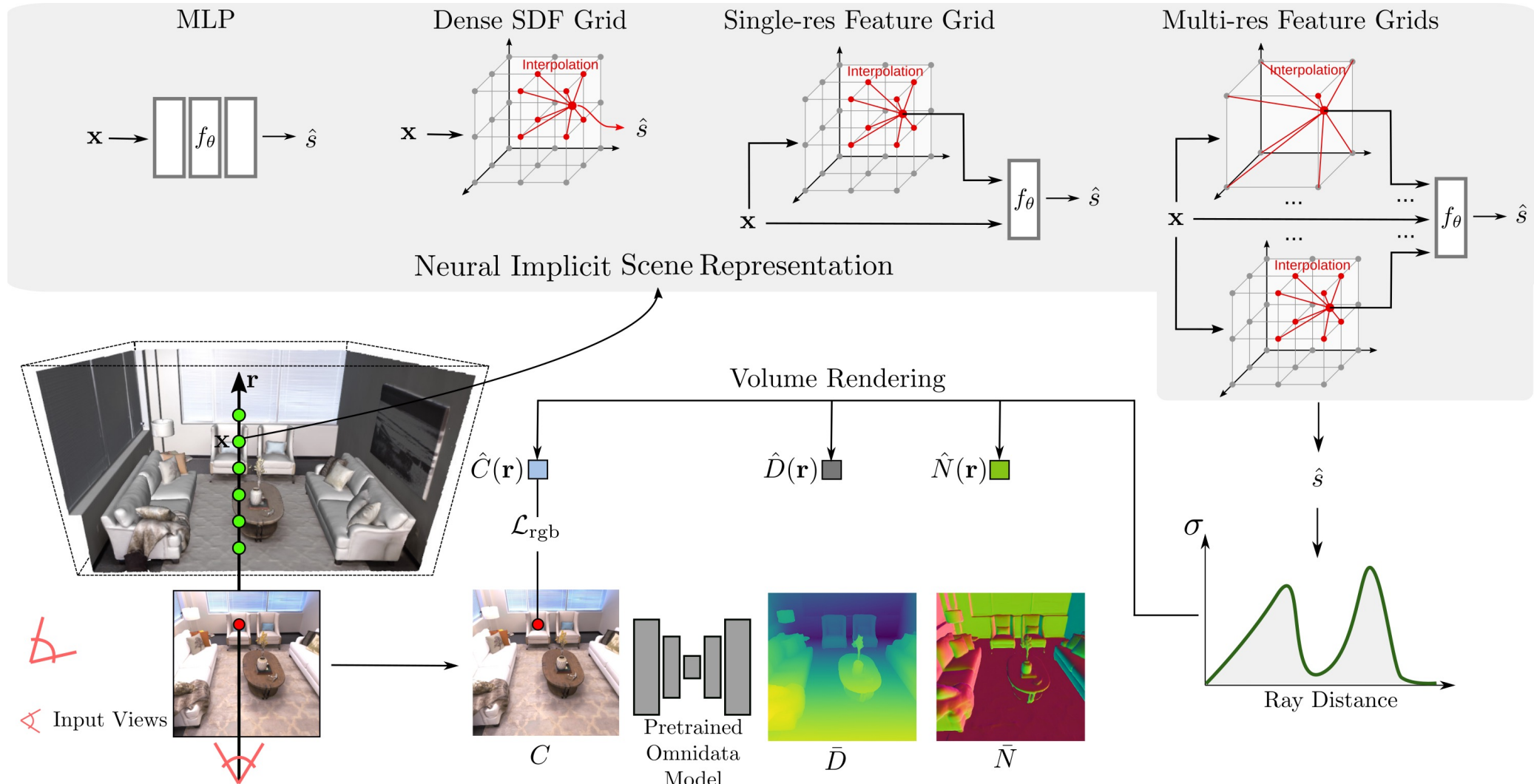
3D Reconstruction Pipeline

MonoSDF



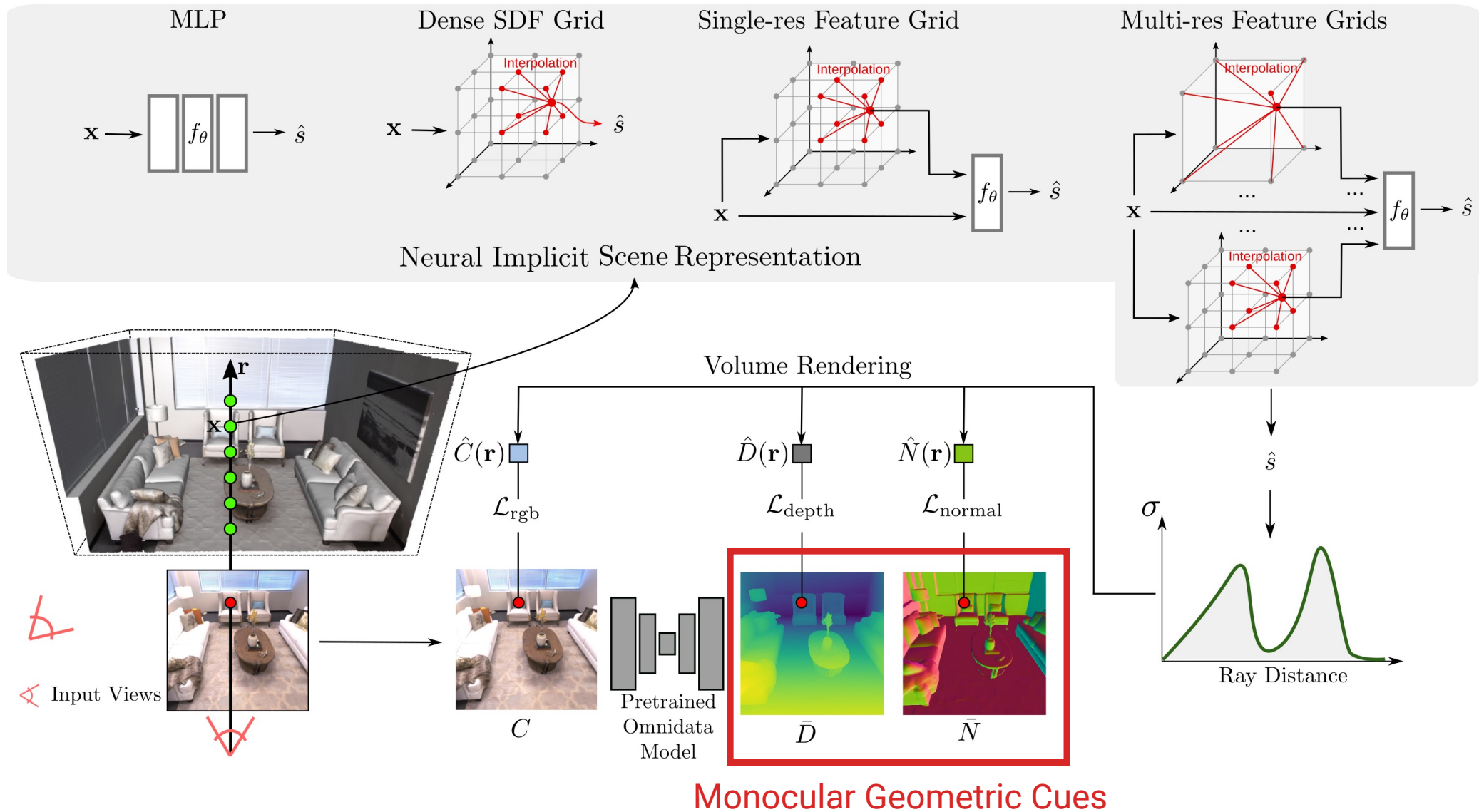
3D Reconstruction Pipeline

MonoSDF



3D Reconstruction Pipeline

MonoSDF



2D Magic in a 3D World

2D Monocular Cues Benefit 3D Reconstruction

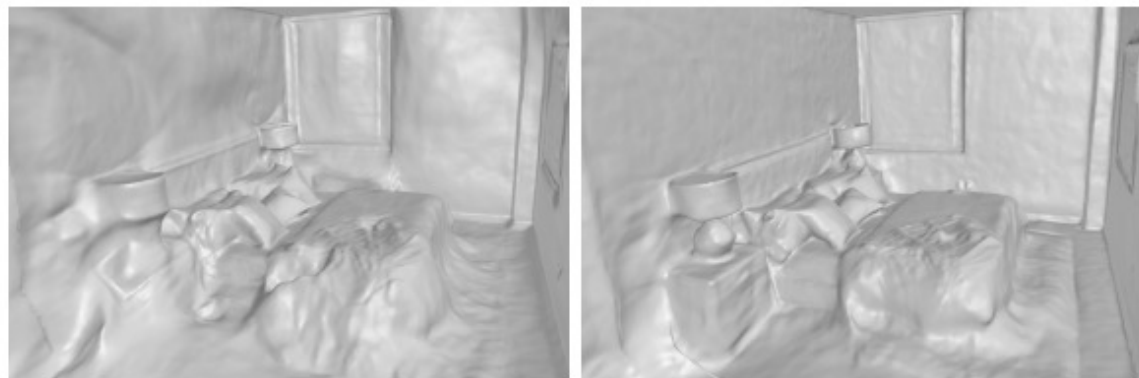


Results

Baseline Comparison

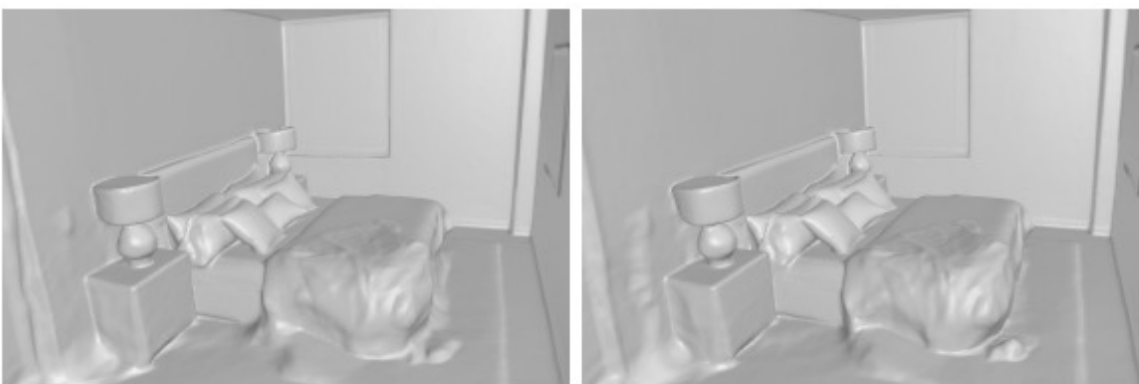


Ours



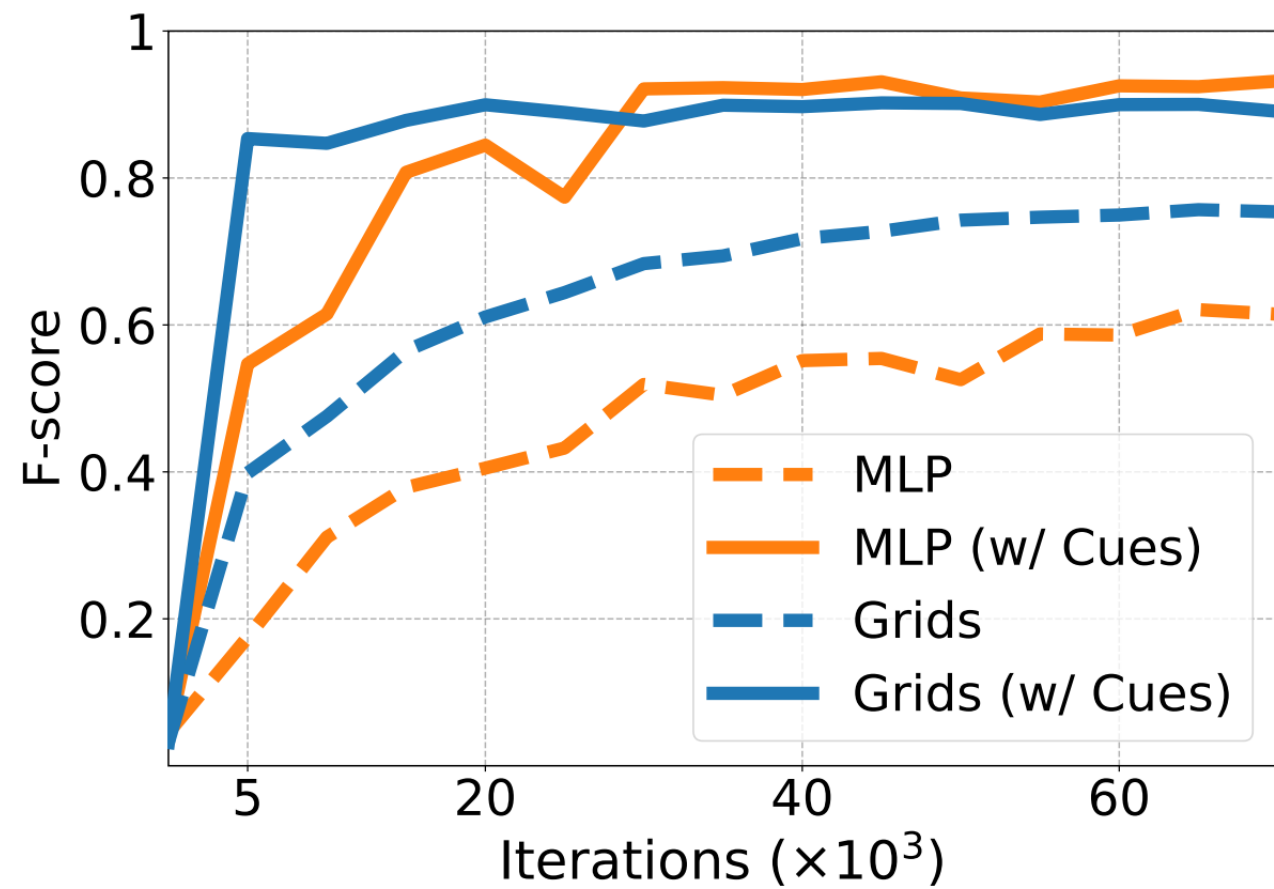
No Cue

+ Depth



+ Normal

+ Both



Large-Scale 3D Scene Reconstruction

Tanks & Temples Dataset



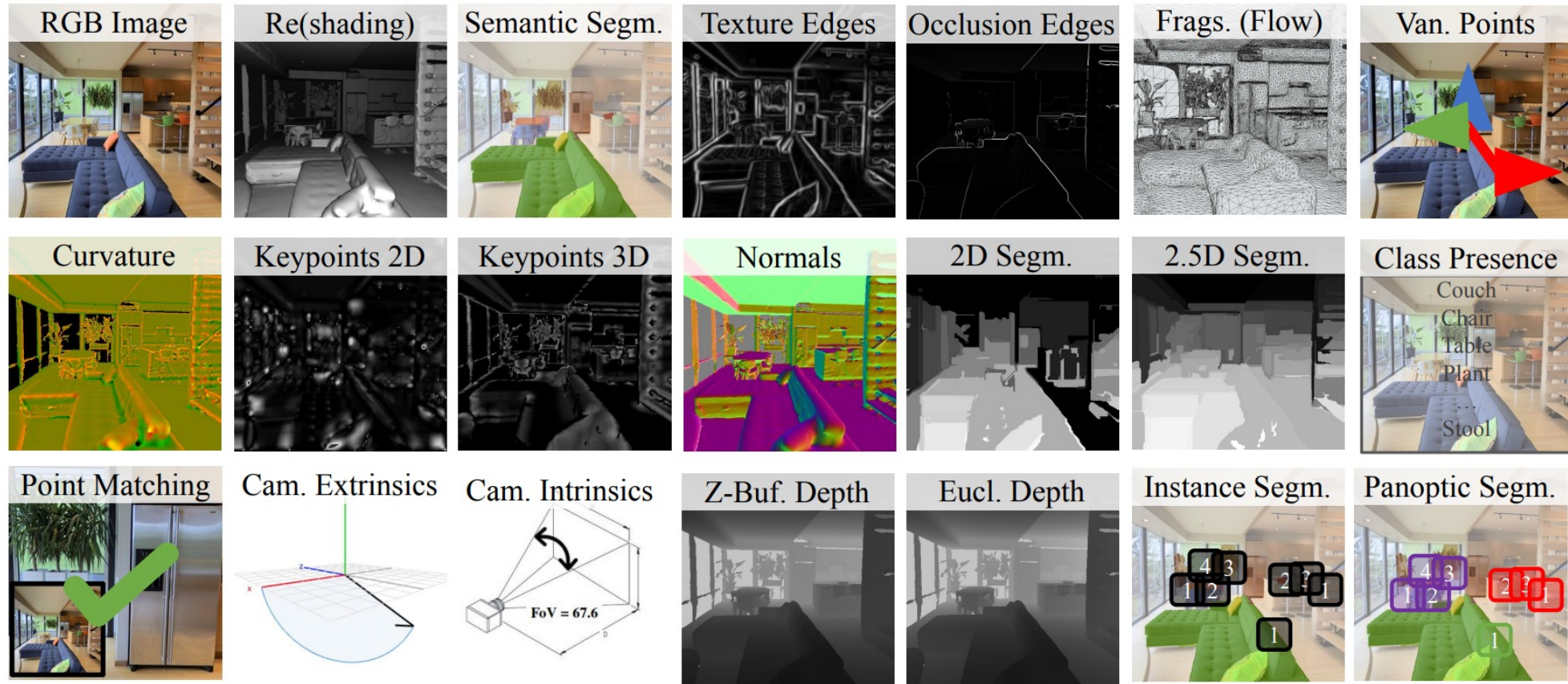
2D Magic in a 3D World

2D Monocular Cues Benefit 3D Reconstruction



Omnidata

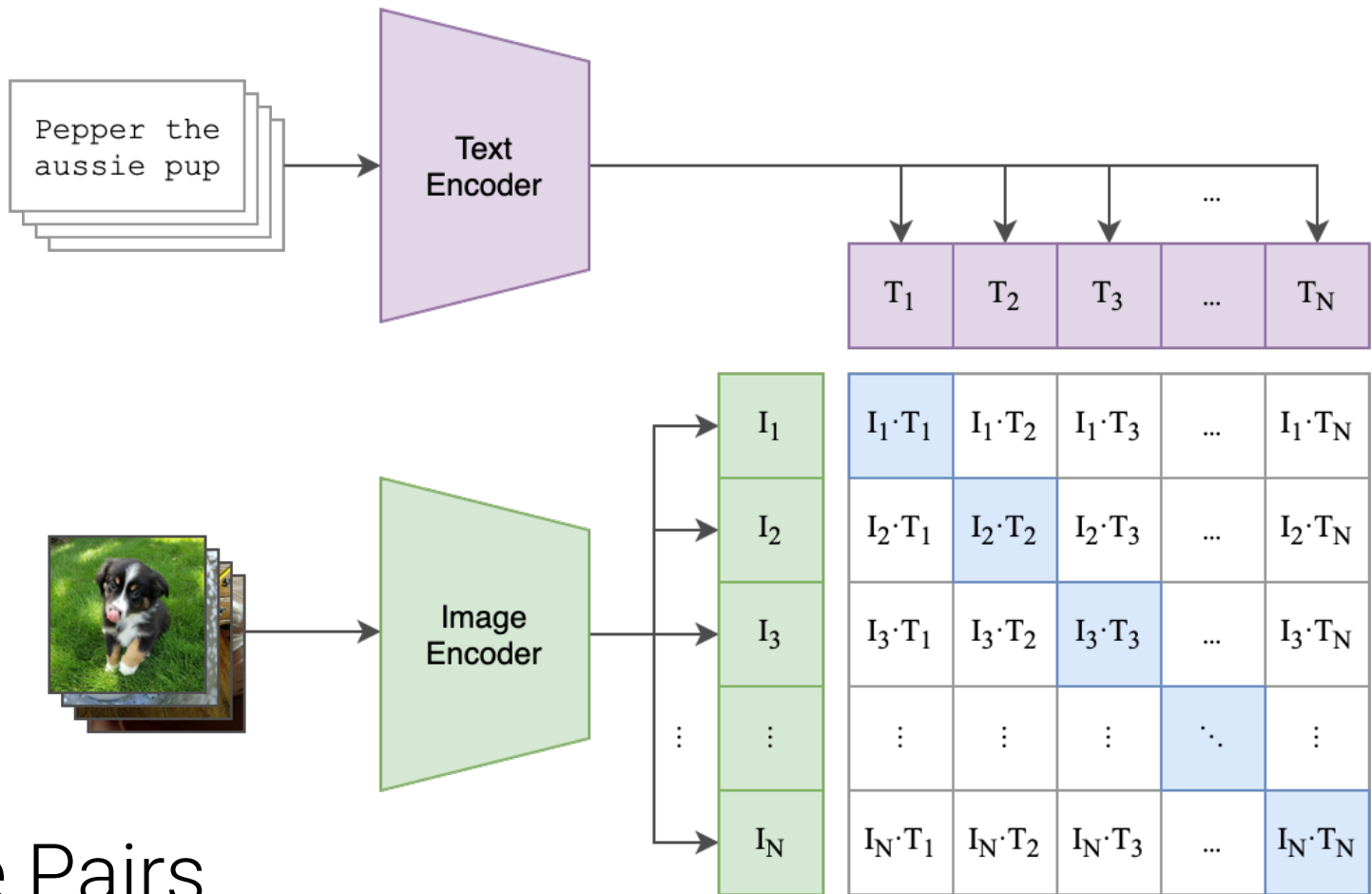
~14M Training Images



What happened since 2021?

Text-Image Pretraining

CLIP



~12.8B Training Text-Image Pairs

Text-to-Image Generation

DALL·E 3



Imagen



Stable Diffusion

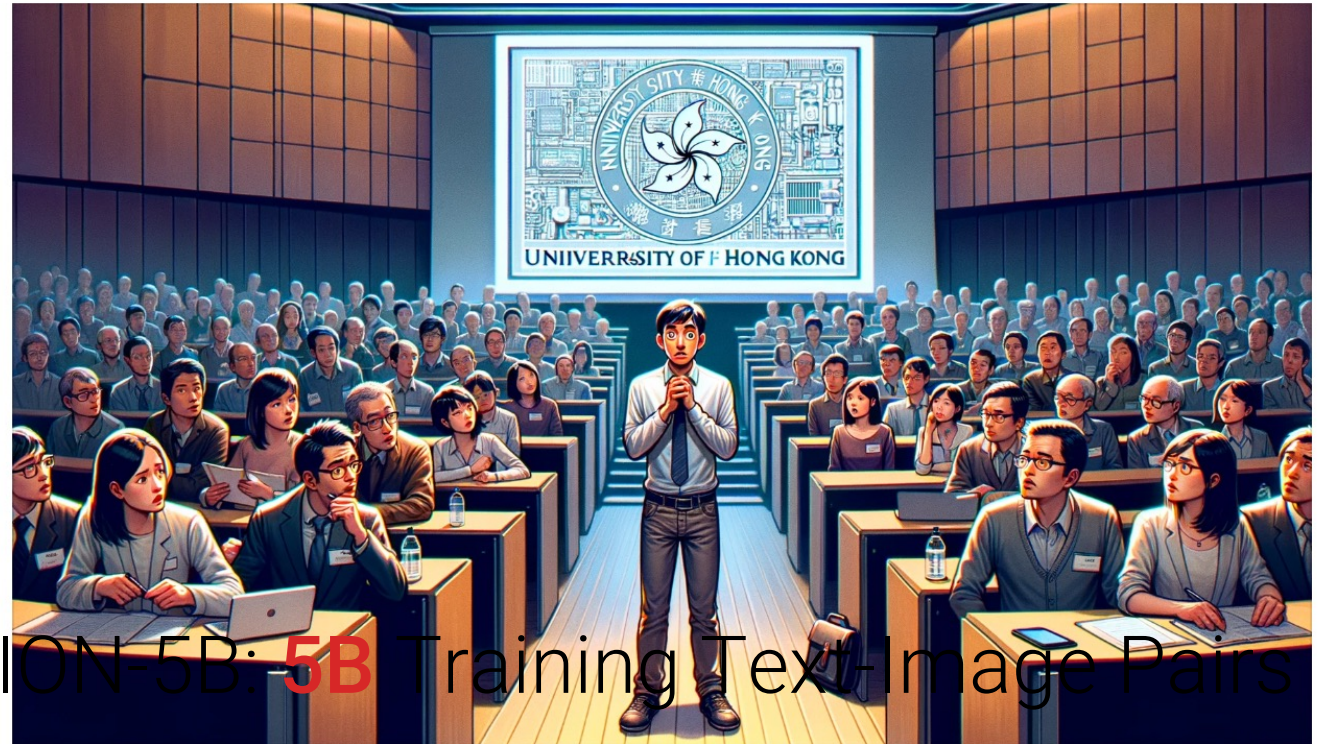


You

Generate an image of "a nervous person presenting in front of many smart researchers at a lecture hall in the University of Hong Kong, with the university logo"



ChatGPT



LAIION-5B: **5B** Training Text-Image Pairs

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

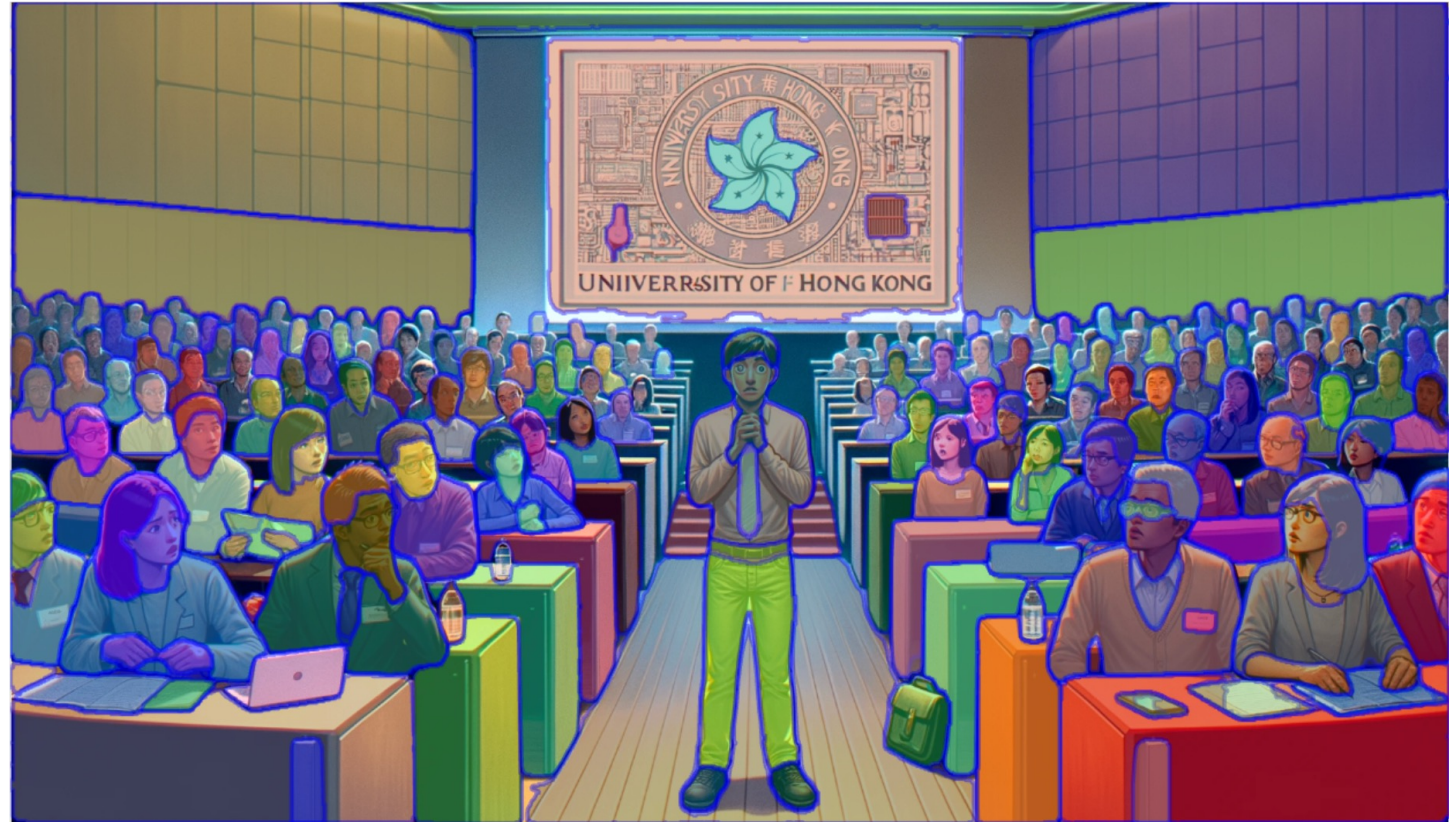
Saharia et al.: [Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding](#). NeurIPS 2022

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

2D Image Segmentation

SAM

 Meta



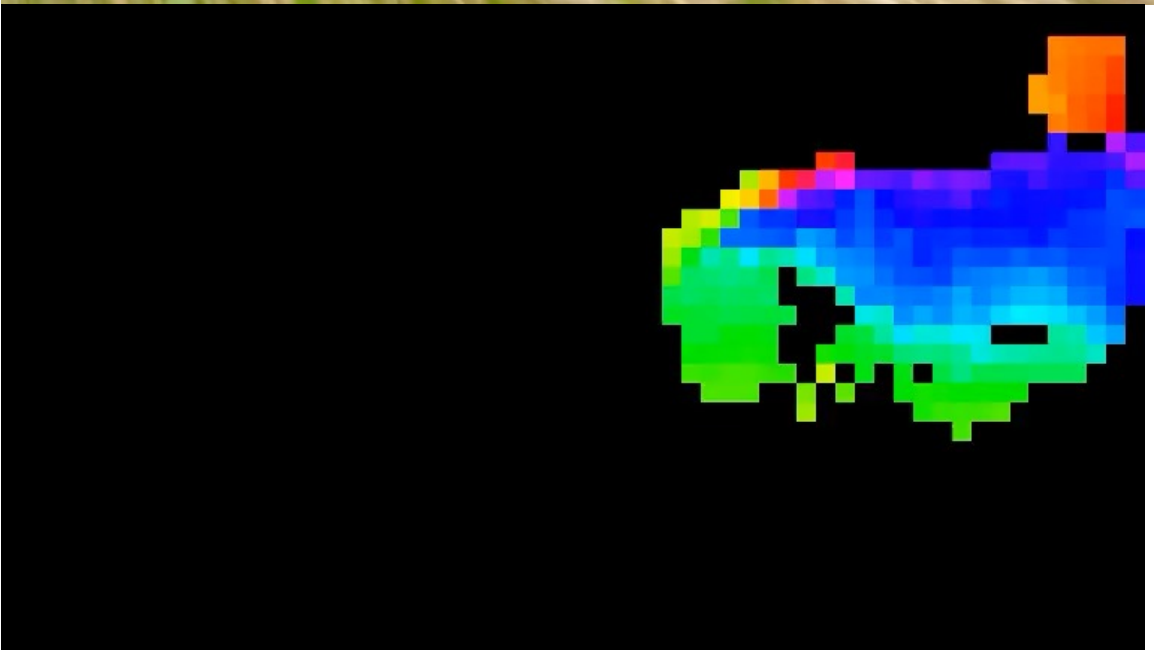
~**1B** Mask + Human-in-the-loop

2D Visual Features

DINO v2



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

The University of Hong Kong

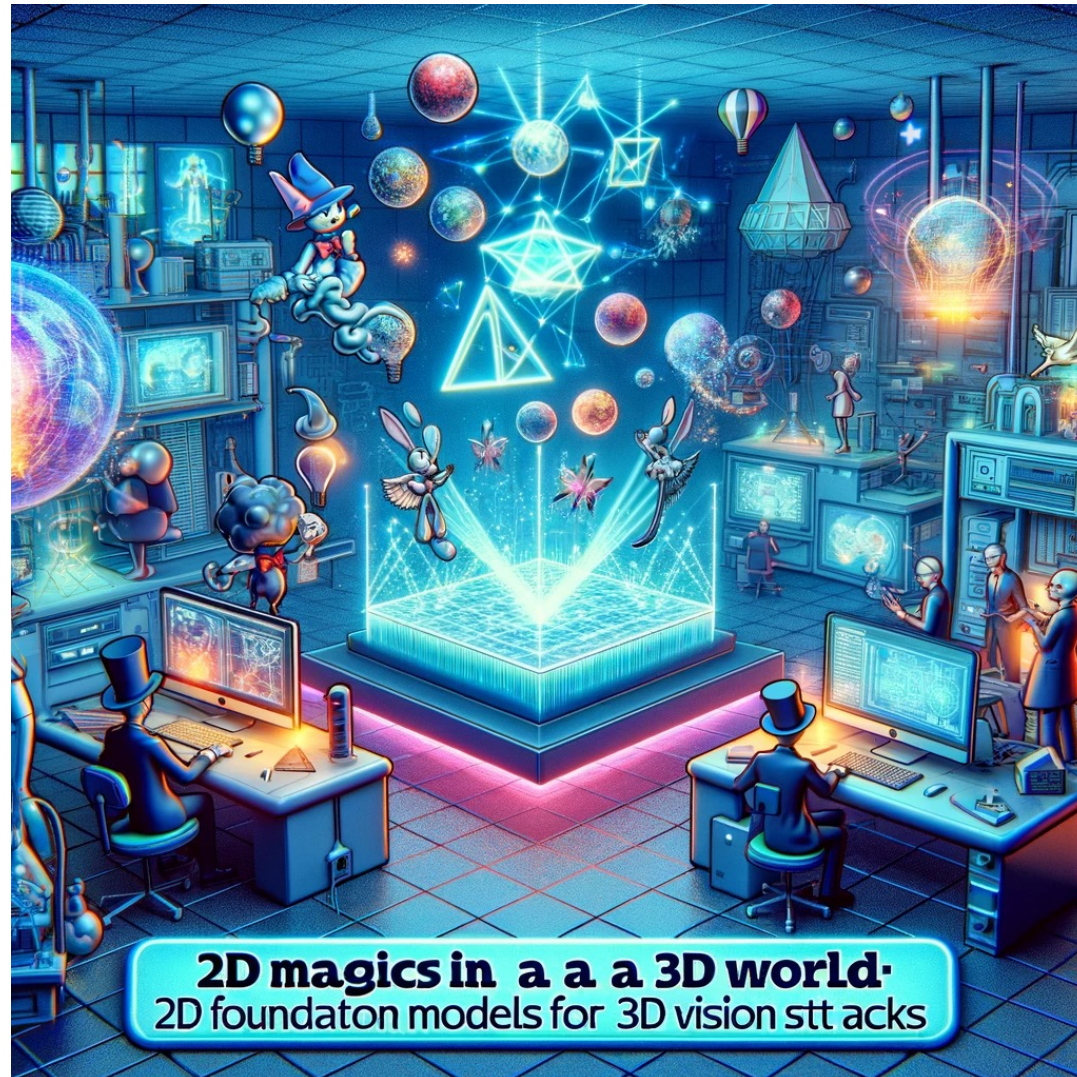
Feb 22, 2024

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



[Generated by DALL·E 3]

2D **Magic** in a 3D World

2D **Foundation Models** for 3D Vision Tasks



2D Magic in a 3D World

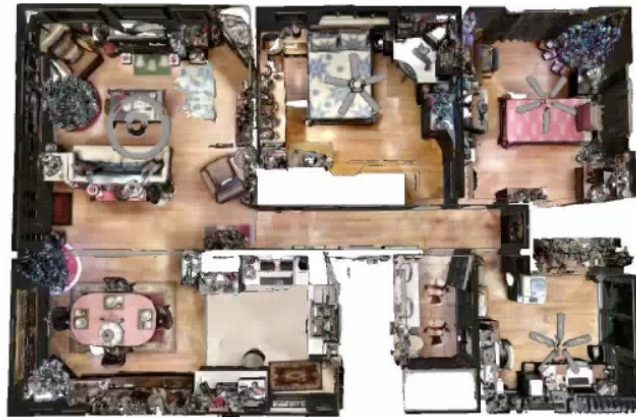
2D Foundation Models for 3D Vision Tasks

3D Reconstruction



NeRF *On-the-go*
(under review)

3D Scene Understanding



OpenScene
CVPR 2023



Segment3D
(under review)

2D Magic in a 3D World

2D Foundation Models for 3D Vision Tasks

3D Reconstruction



NeRF *On-the-go*
(under review)

3D Scene Understanding



OpenScene
CVPR 2023



Segment3D
(under review)

NeRF Is Awesome





THE
PARISIAN
DEPARTMENT
STORE



TOUR EIFFEL - CHAMP-DE-MARS - MUSÉE DU LOUVRE - NOTRE-DAME - MUSÉE D'ORSAY - OPERA GARNIER - CHAMPS-ÉLYSÉES - GRAND PALAIS - TROCADÉROS
BIGBUS PARIS - LES CARROUGES

**BIG
BUS**

HOP-ON HOP-OFF

Motivation

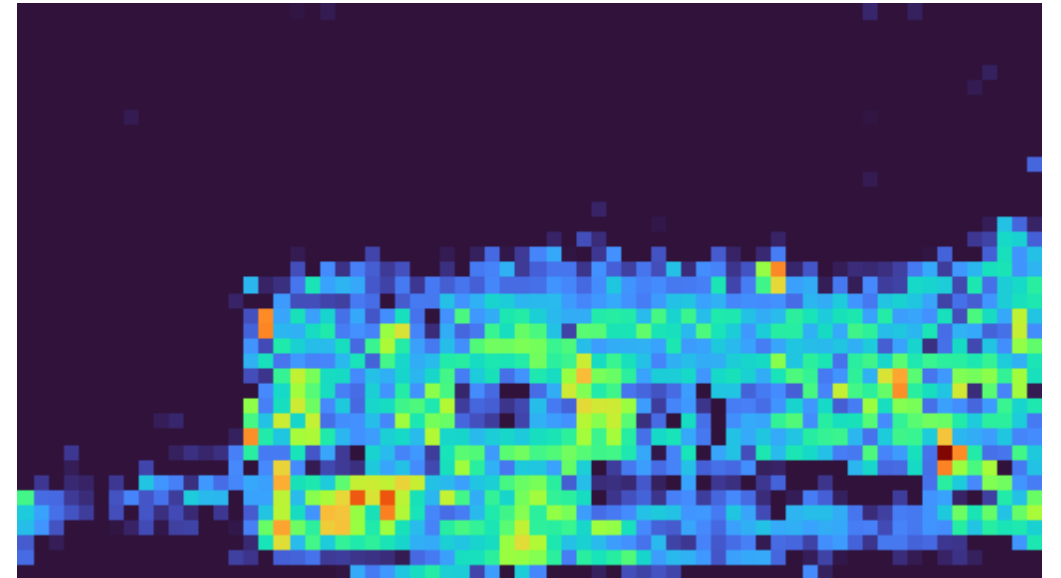


How to obtain **distractor-free NeRFs** from **casually captured sequences**?

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

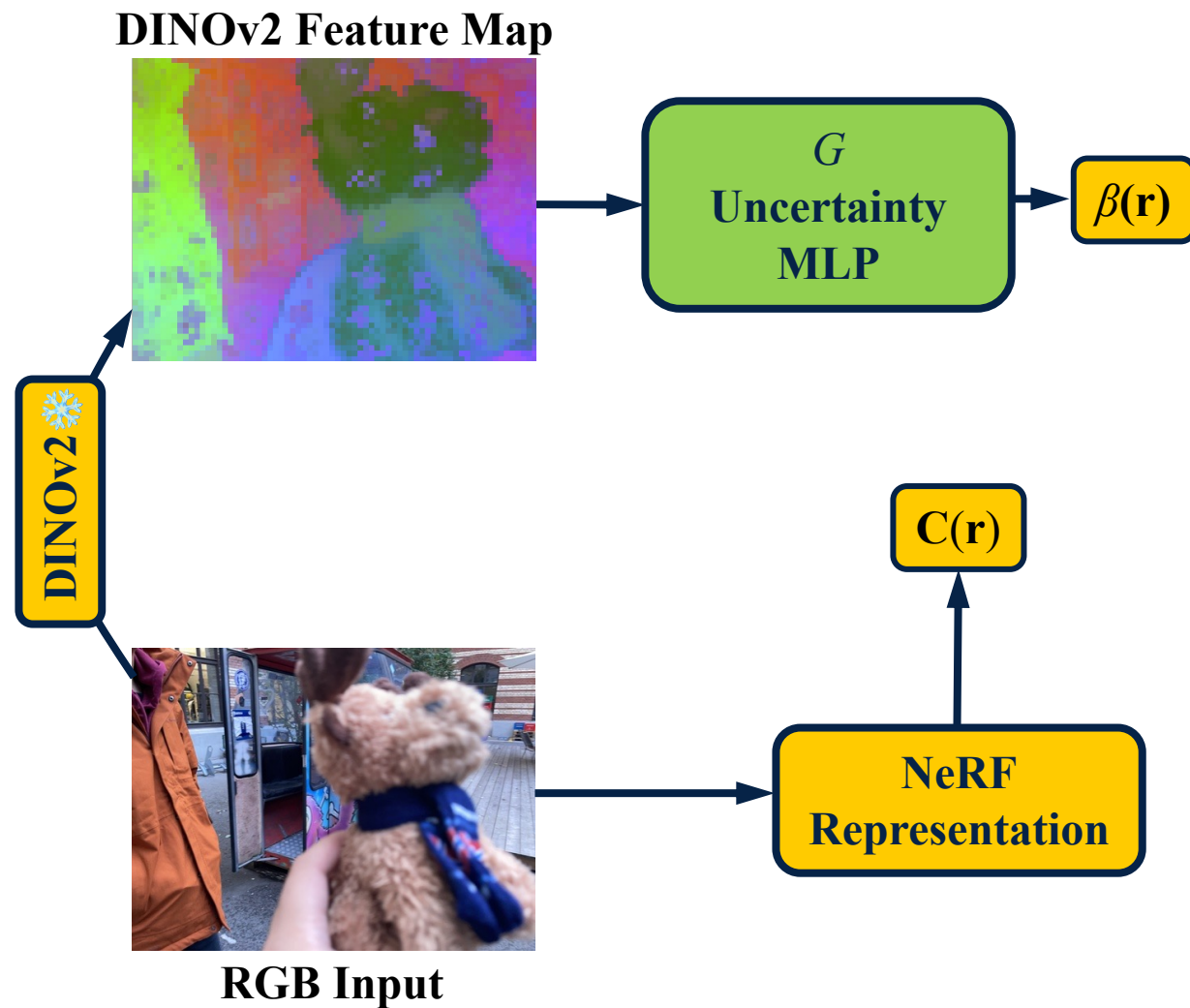


Marc
Pollefeys



Songyou
Peng

Pipeline



To Learn the Uncertainty MLP...



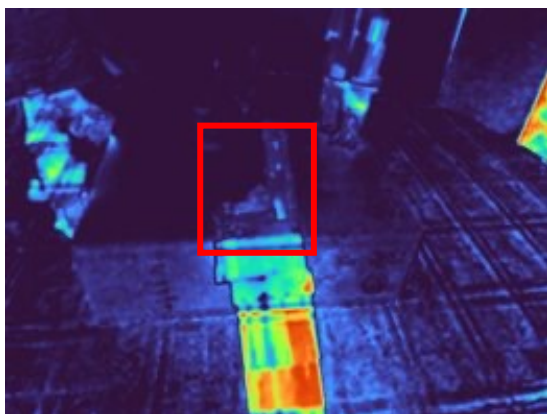
Rendered RGB



Train RGB

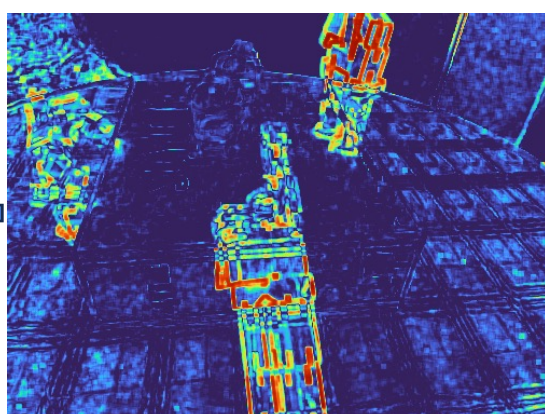
Why SSIM?

Leverage structure information when RGB is similar!



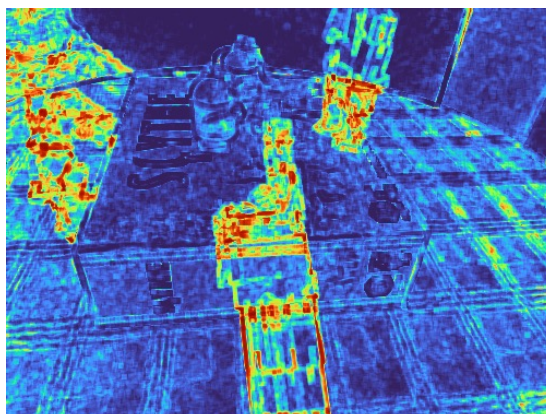
Luminance

+



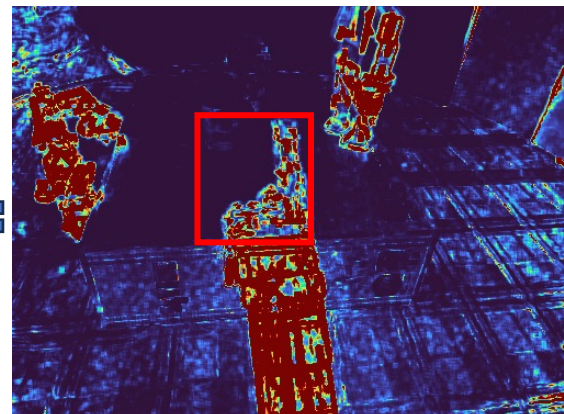
Contrast

+



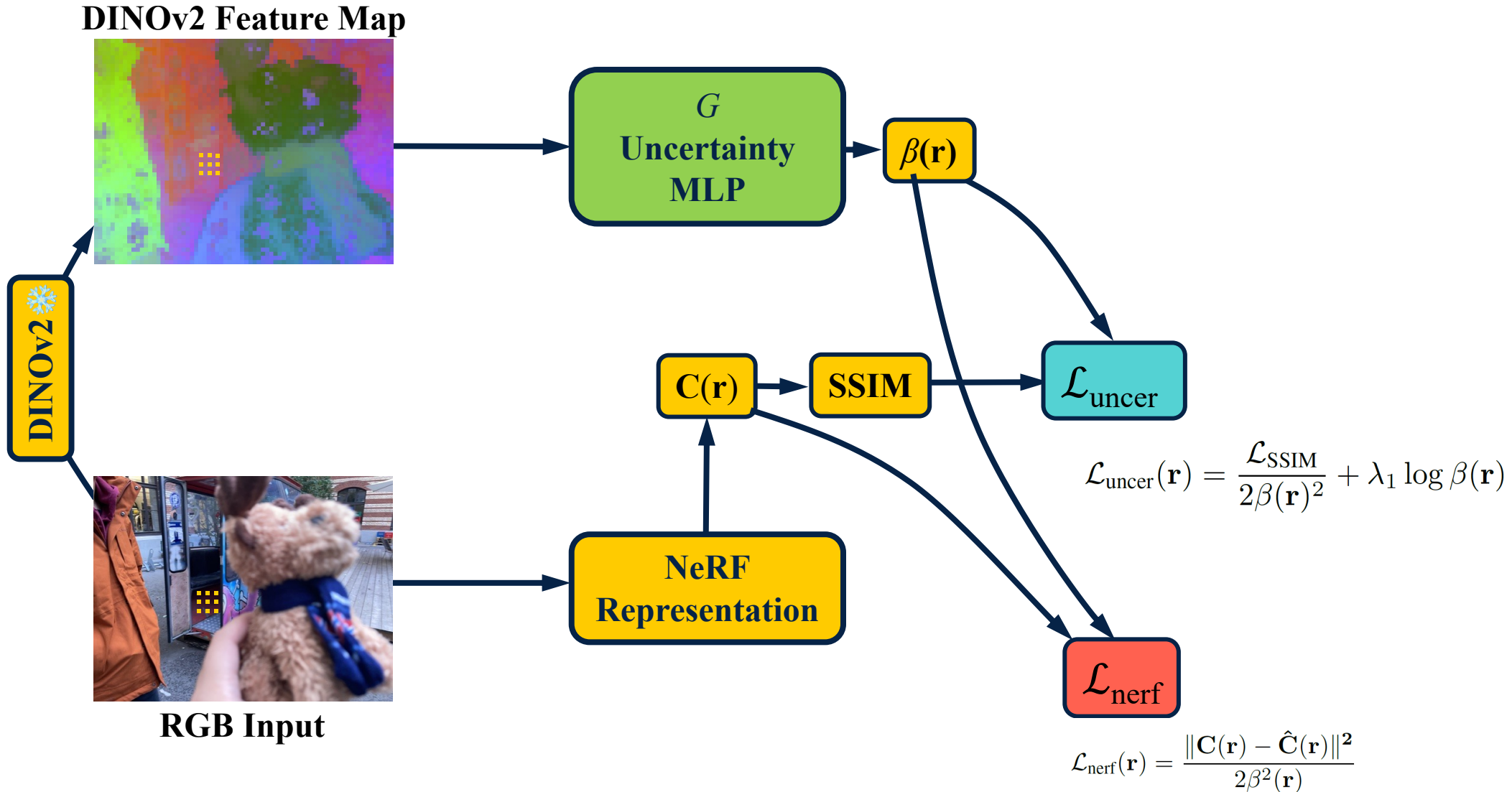
Structure

=

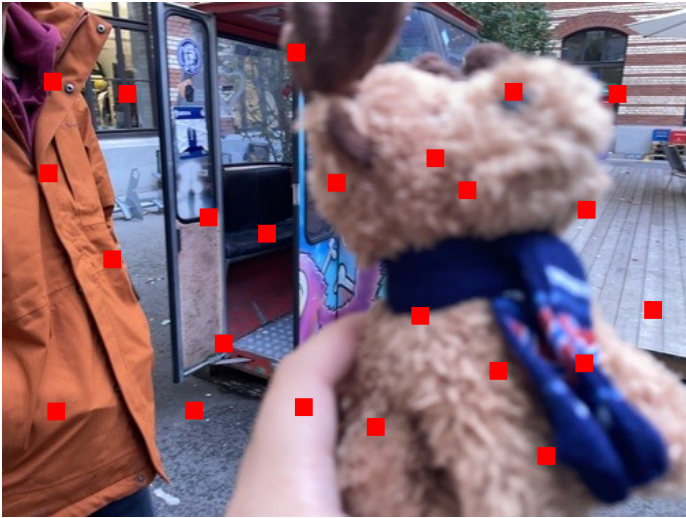


SSIM Error

Pipeline



Sampling Strategy



Random
(NeRF)



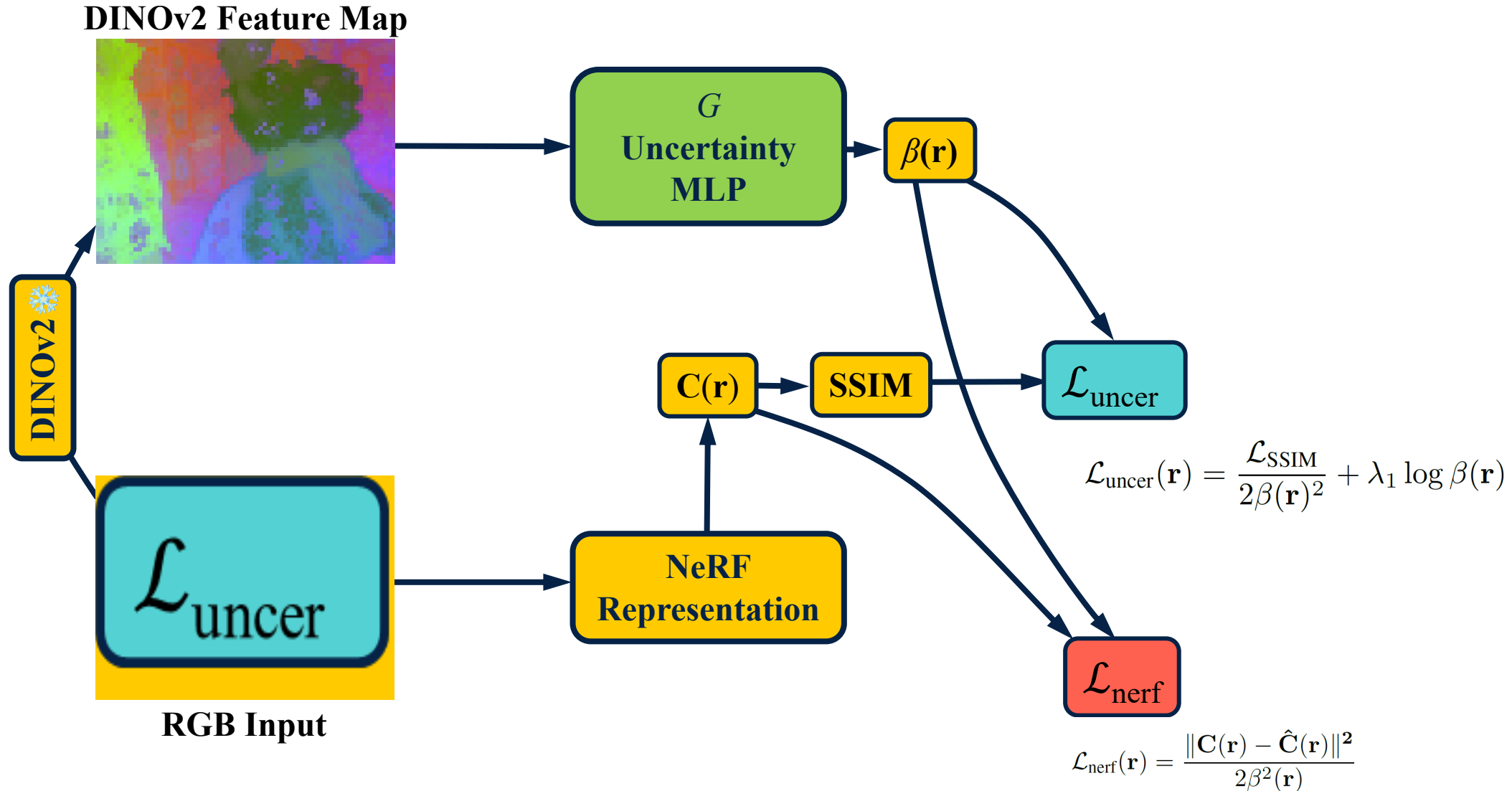
Patch
(RobustNeRF)



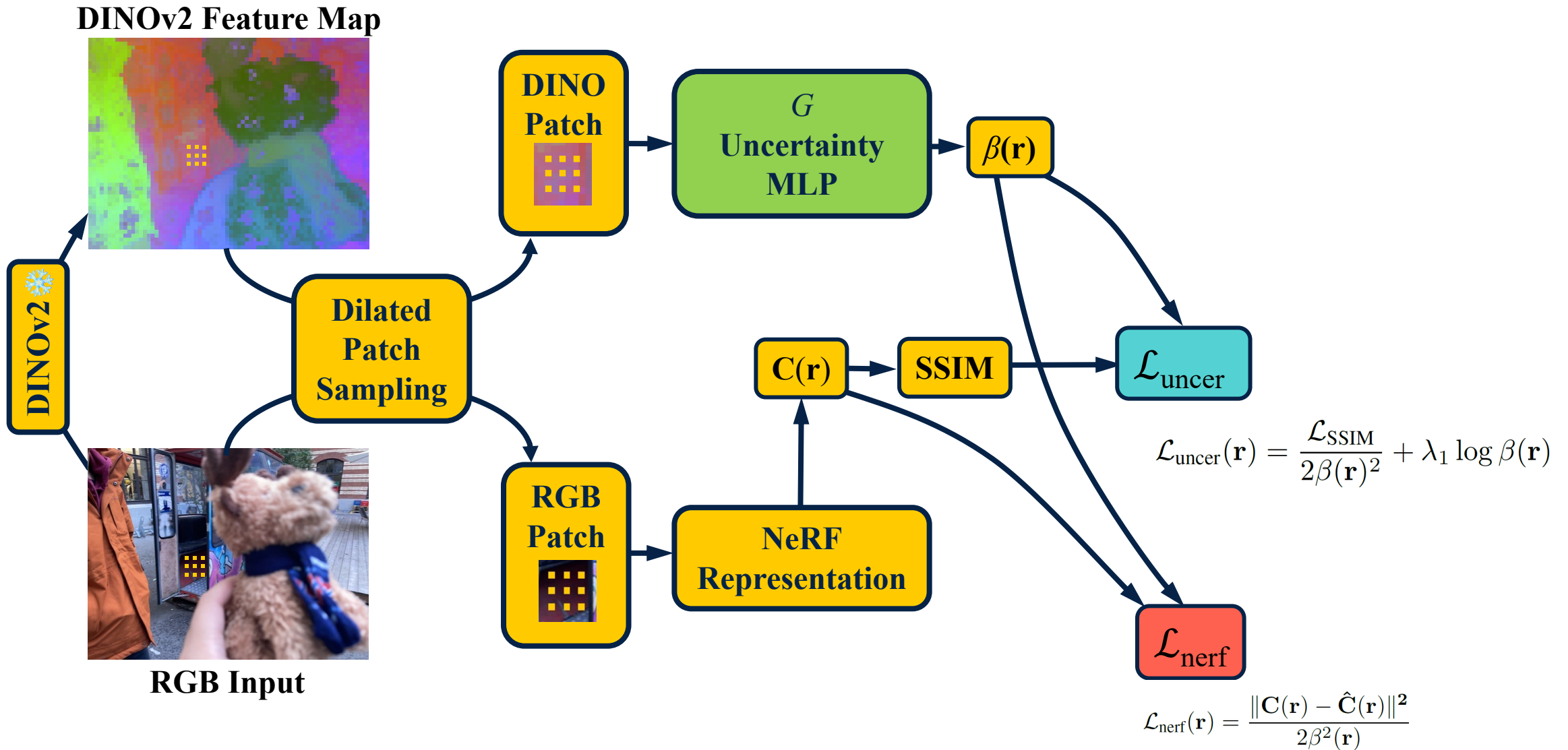
Dilated Patch
(Ours)

- ✓ **Larger Perceptive Field:** Improve efficiency, reconstruction quality
- ✓ **More Local Information:** Better distractor removal

Pipeline



Pipeline

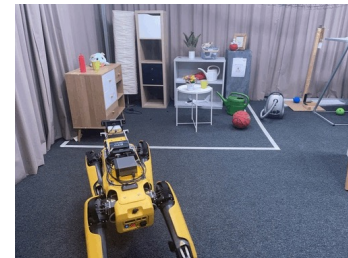


Results

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



Train Station - Rendering Comparisons

Occlusion
Ratio: **High**



Patio-High - Input



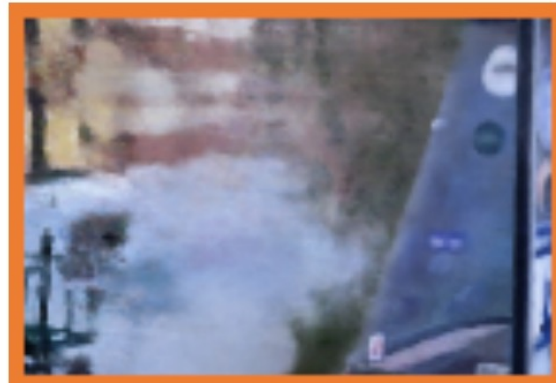
NeRF On-the-go
(Ours)

Analysis

Analysis - Efficiency



25K



50K



100K



250K

RobustNeRF

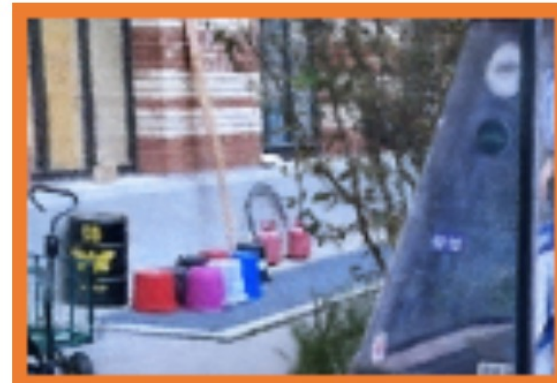
Analysis - Efficiency



25K



50K



100K



250K

NeRF On-the-go
(Ours)

Analysis – Static Scene



RobustNeRF

Ours

MipNeRF 360

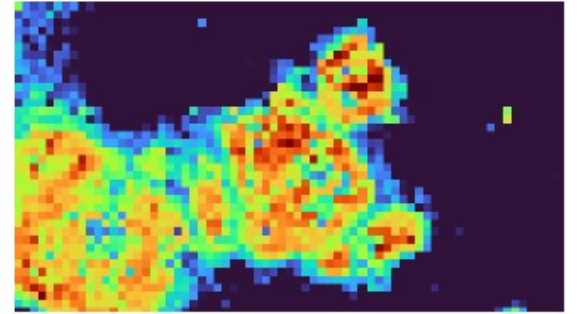
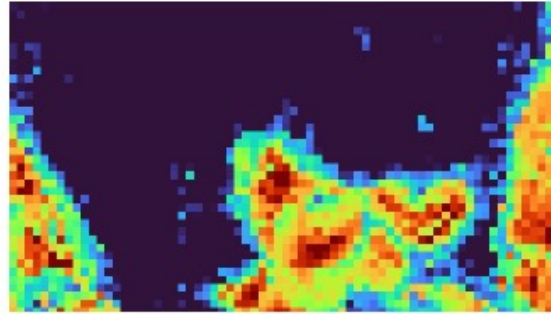
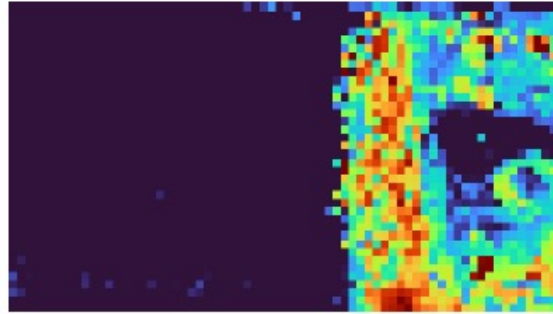
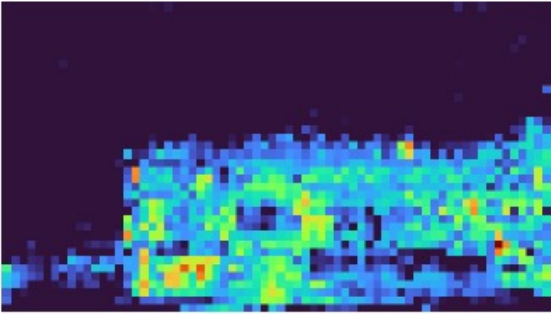
GT

Analysis – Handle Large Occlusions

Input



Uncertainty



Rendering

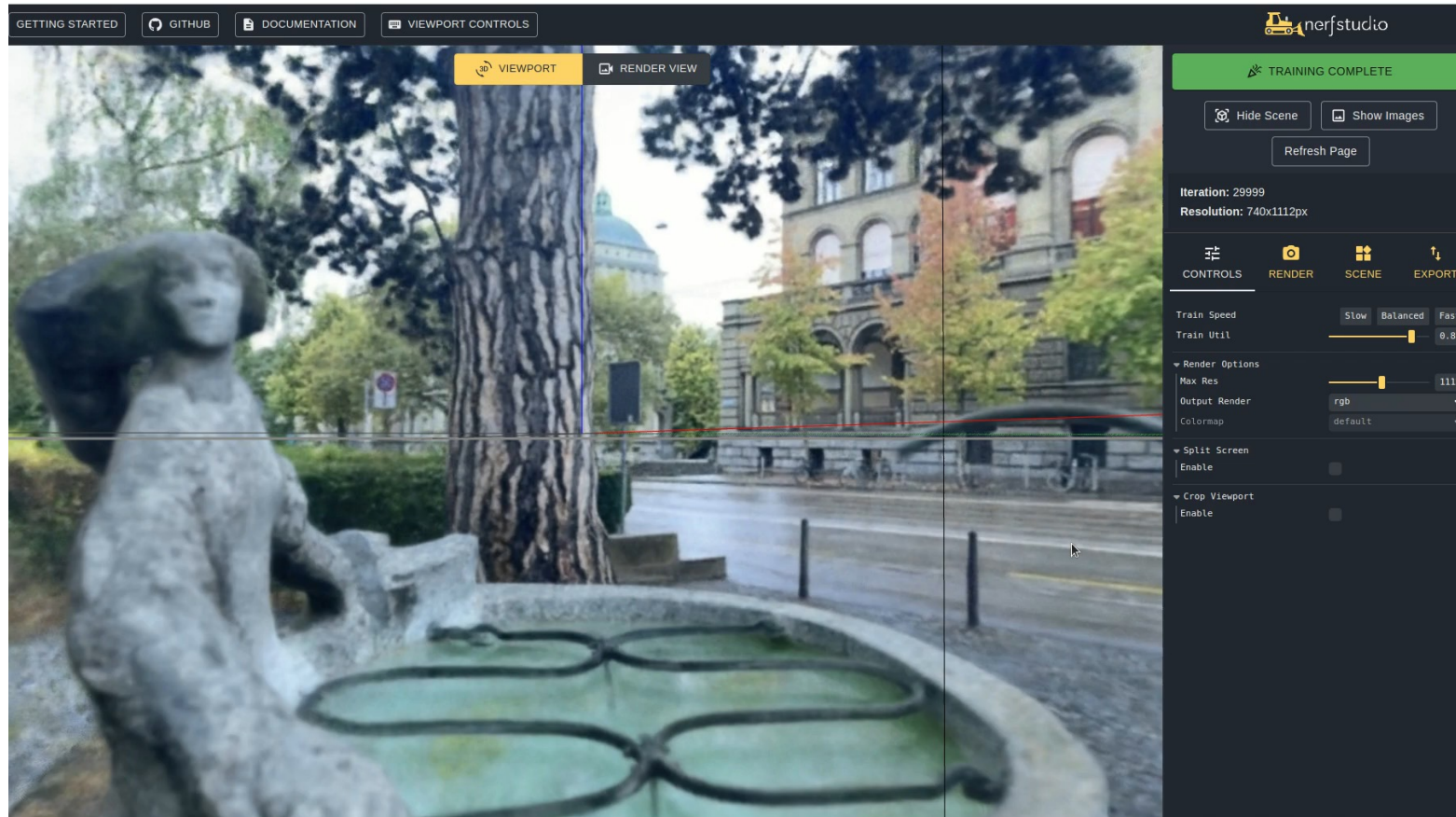


Arc de Triomphe

Patio-High

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?

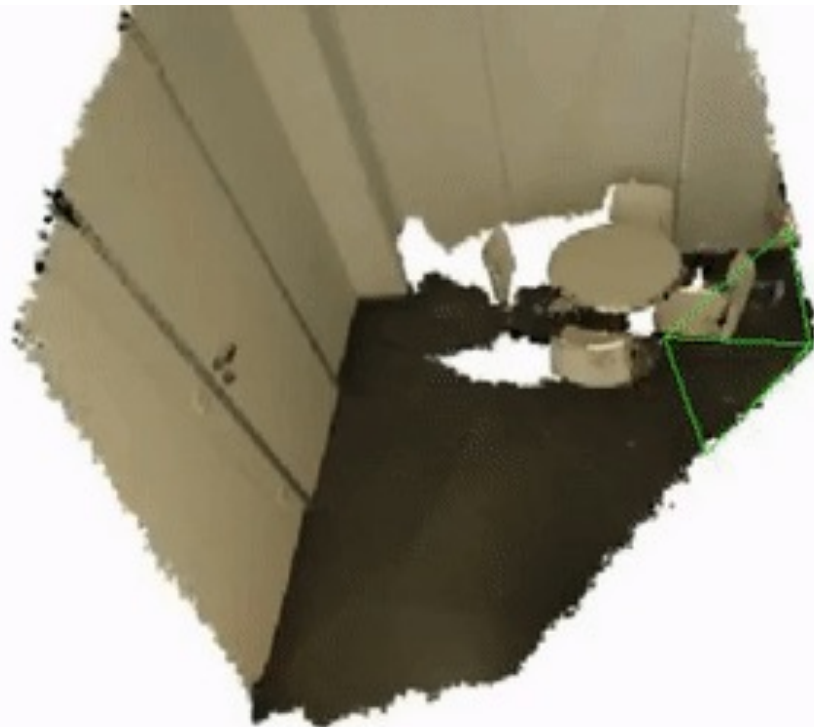
NeRF *On-the-go*

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*
(under review)

3D Scene Understanding



OpenScene
CVPR 2023



Segment3D
(under review)

2D Magic in a 3D World

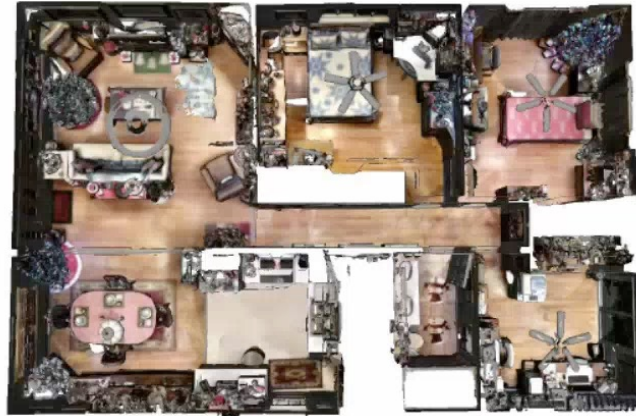
2D Foundation Models for 3D Vision Tasks

3D Reconstruction



NeRF *On-the-go*
(under review)

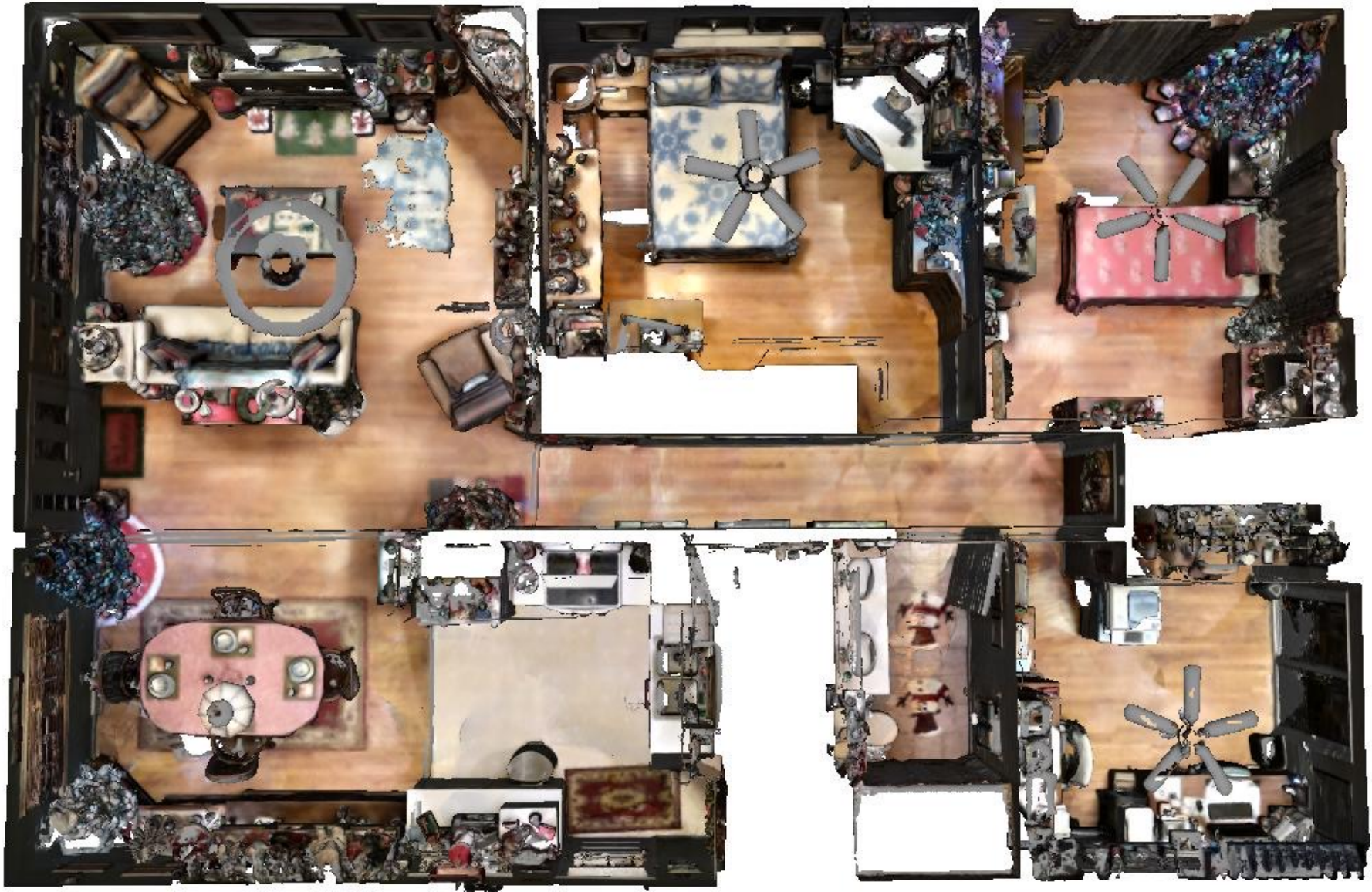
3D Scene Understanding



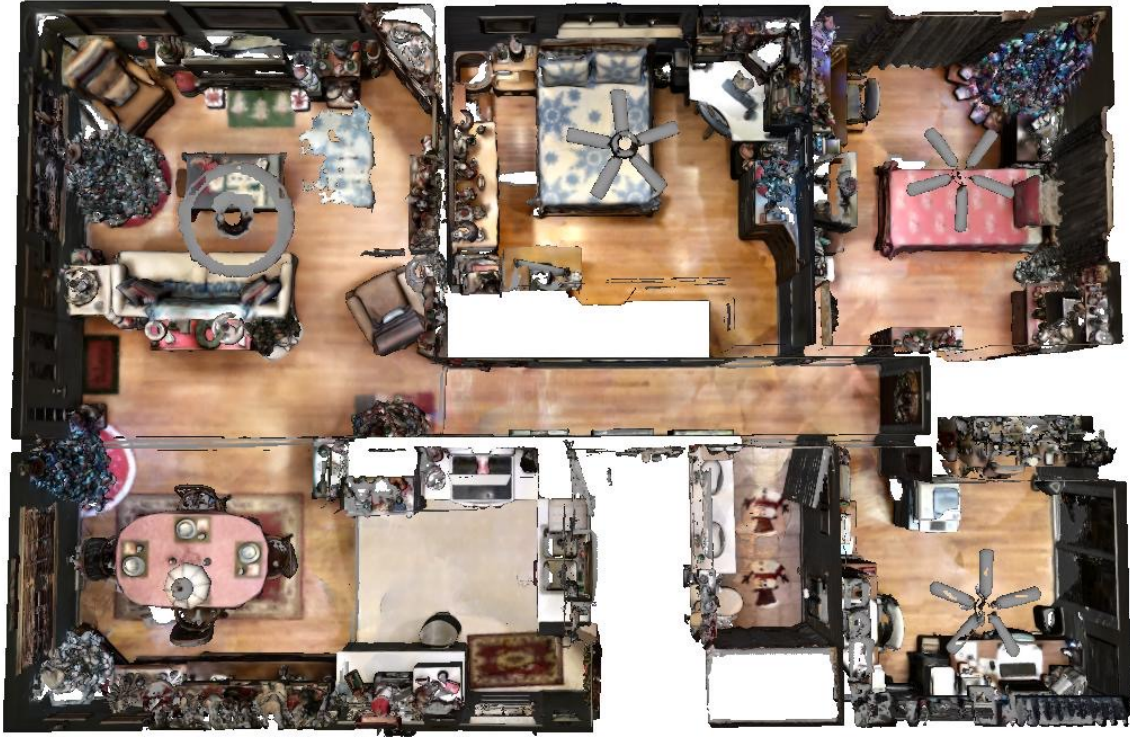
OpenScene
CVPR 2023



Segment3D
(under review)

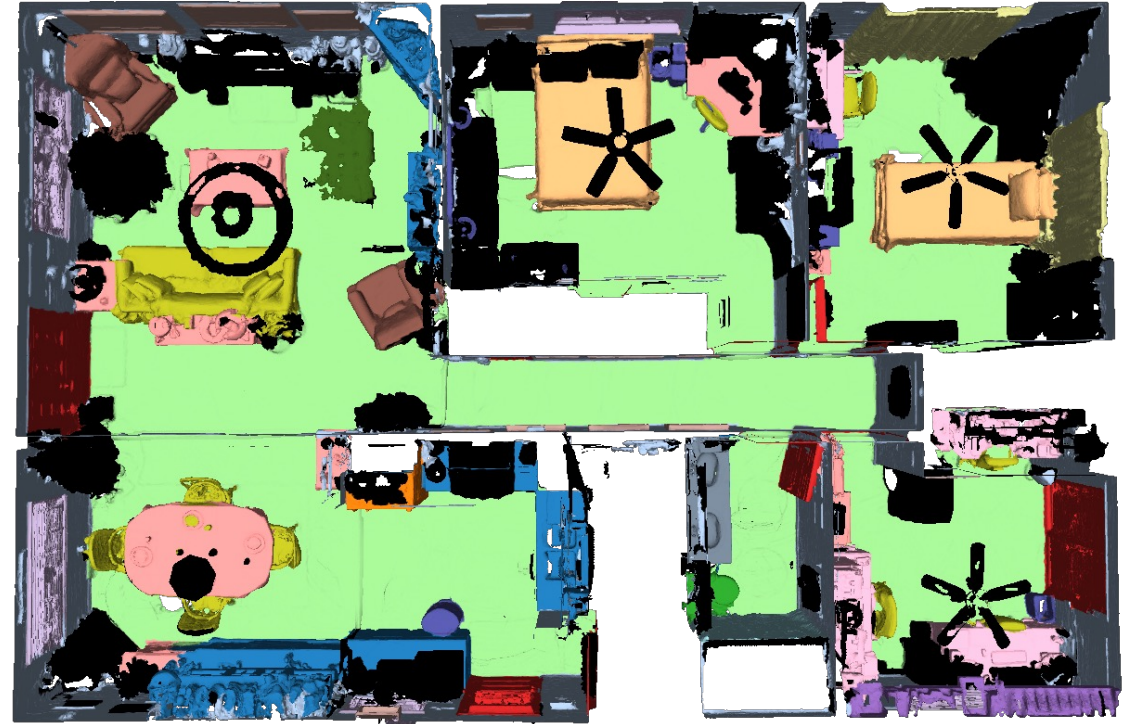


Input 3D Geometry



Input 3D Geometry

■ 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

- Affordance prediction



Input 3D Geometry

3D Scene Understanding Tasks **w/o** Labels

- Affordance prediction



Example: “where can I sit?”

3D Scene Understanding Tasks **w/o** Labels



Input 3D Geometry

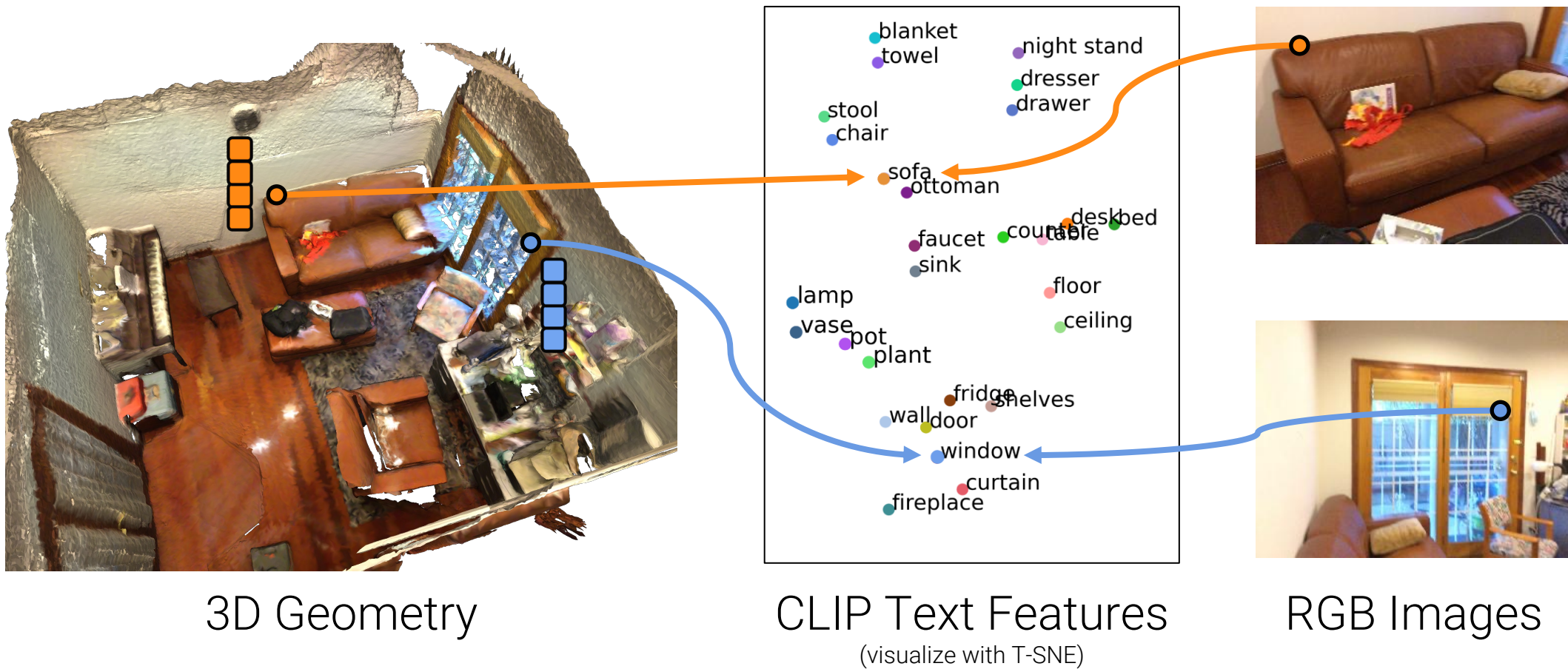
- Affordance prediction
- Material identification
- Physical property estimation
- Rare object retrieval
- Activity site prediction
- Fine-grained semantic segmentation
- Many more...

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



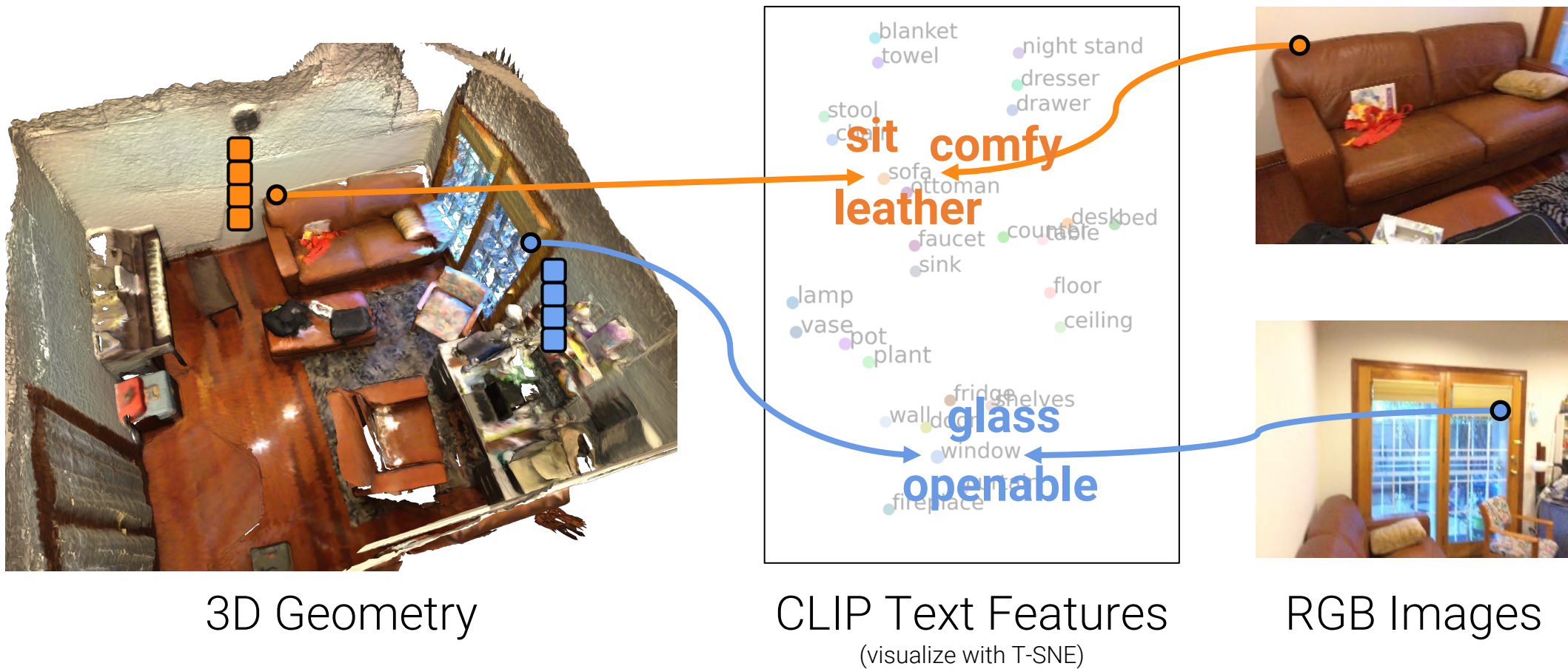
3D Geometry

CLIP Text Features
(visualize with T-SNE)

RGB Images

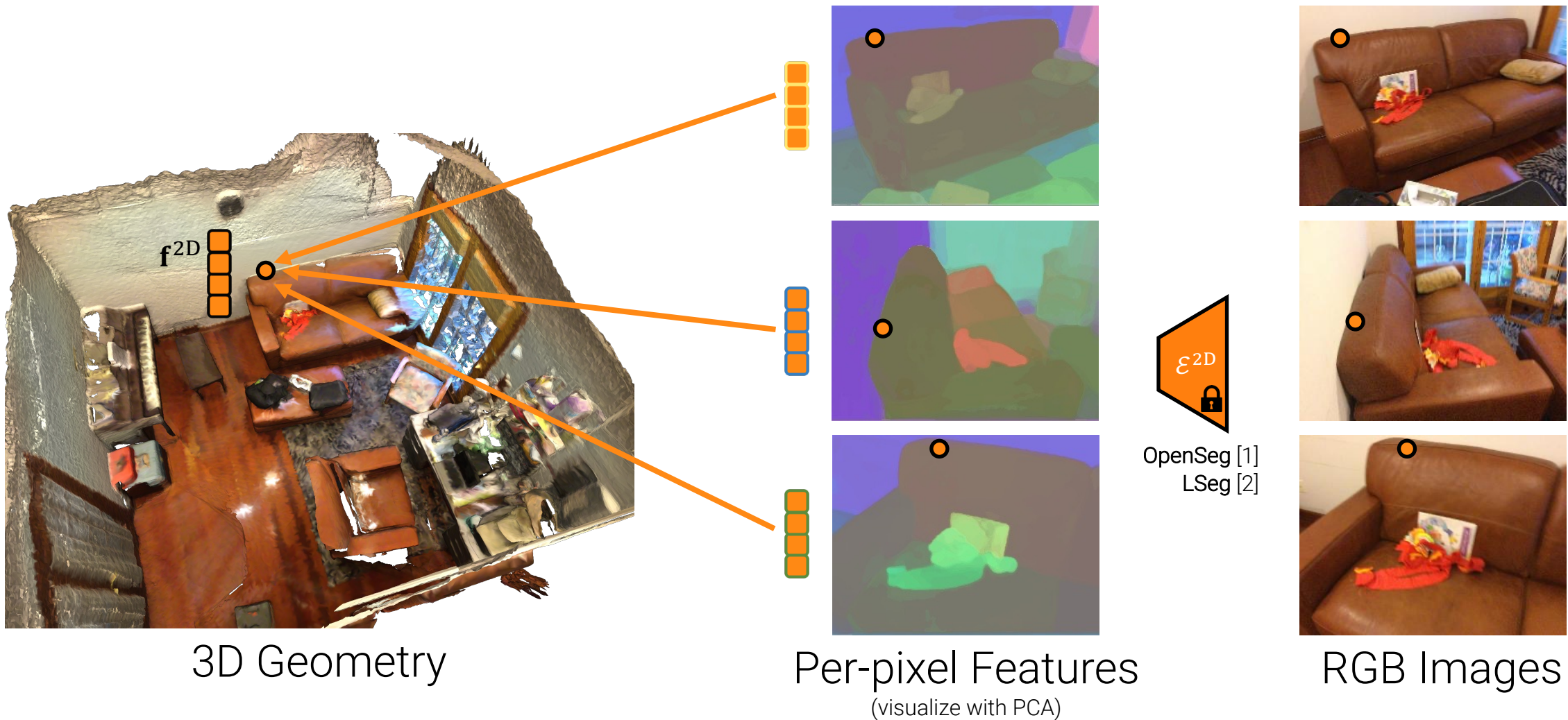
Key Idea

Co-embed 3D Features with CLIP Features



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

Step 1: Multi-view Feature Fusion



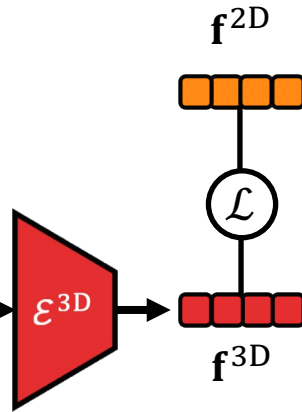
[1] Ghiasi, Gu, Cui, Lin: [Scaling Open-Vocabulary Image Segmentation with Image-Level Labels](#). ECCV 2022

[2] Li, Weinberger, Belongie, Koltun, Ranftl: [Language-driven Semantic Segmentation](#). ICLR 2022

Step 2: 3D Feature Distillation

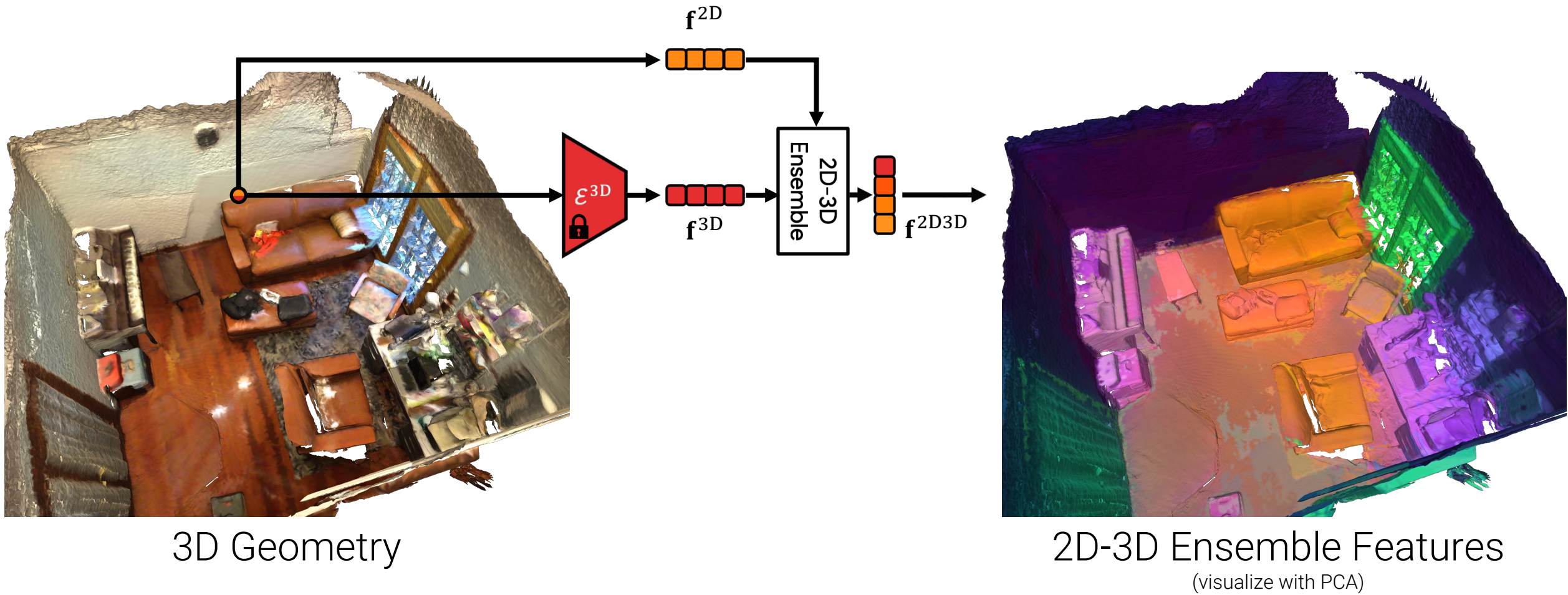


3D Geometry



$$\mathcal{L} = 1 - \cos(\mathbf{f}^{2D}, \mathbf{f}^{3D})$$

Inference: 2D-3D Ensemble



3D Geometry

2D-3D Ensemble Features
(visualize with PCA)

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	■ bathtub	■ dresser	■ stand	■ clock	■ tissue box	■ furniture	■ soap	■ cup	■ hanger	■ urn	■ paper towel dispenser	■ toy
■ door	■ curtain	■ night stand	■ desk	■ book	■ rug	■ drawer	■ stove	■ tv stand	■ air conditioner	■ thermostat	■ ladder	■ candlestick	■ plate	■ lamp shade	■ foot rest
■ ceiling	■ table	■ toilet	■ box	■ air vent	■ ottoman	■ container	■ washing machine	■ shoe	■ fire extinguisher	■ radiator	■ garage door	■ light	■ car	■ soap dish	■ cleaner
■ floor	■ plant	■ column	■ coffee table	■ faucet	■ bottle	■ light switch	■ shower curtain	■ heater	■ kitchen island	■ paper towel	■ board	■ scale	■ jacket	■ drum	■ computer
■ picture	■ mirror	■ banister	■ counter	■ photo	■ refridgerator	■ purse	■ curtain rod	■ headboard	■ printer	■ sheet	■ rope	■ display case	■ bottle of soap	■ whiteboard	■ knob
■ window	■ towel	■ stairs	■ bench	■ toilet paper	■ bookshelf	■ door way	■ bin	■ telephone	■ bucket	■ glass	■ ball	■ toilet paper holder	■ water cooler	■ whiteboard	■ computer
■ chair	■ sink	■ stool	■ garbage bin	■ fan	■ wardrobe	■ basket	■ chest	■ blanket	■ microwave	■ candle	■ exercise equipment	■ tea pot	■ range hood	■ paper	■ projector
■ pillow	■ shelves	■ vase	■ fireplace	■ railing	■ pipe	■ chandelier	■ microwave	■ flower pot	■ blinds	■ handle	■ tray	■ stuffed animal	■ candelabra		

Image-based 3D Scene Query



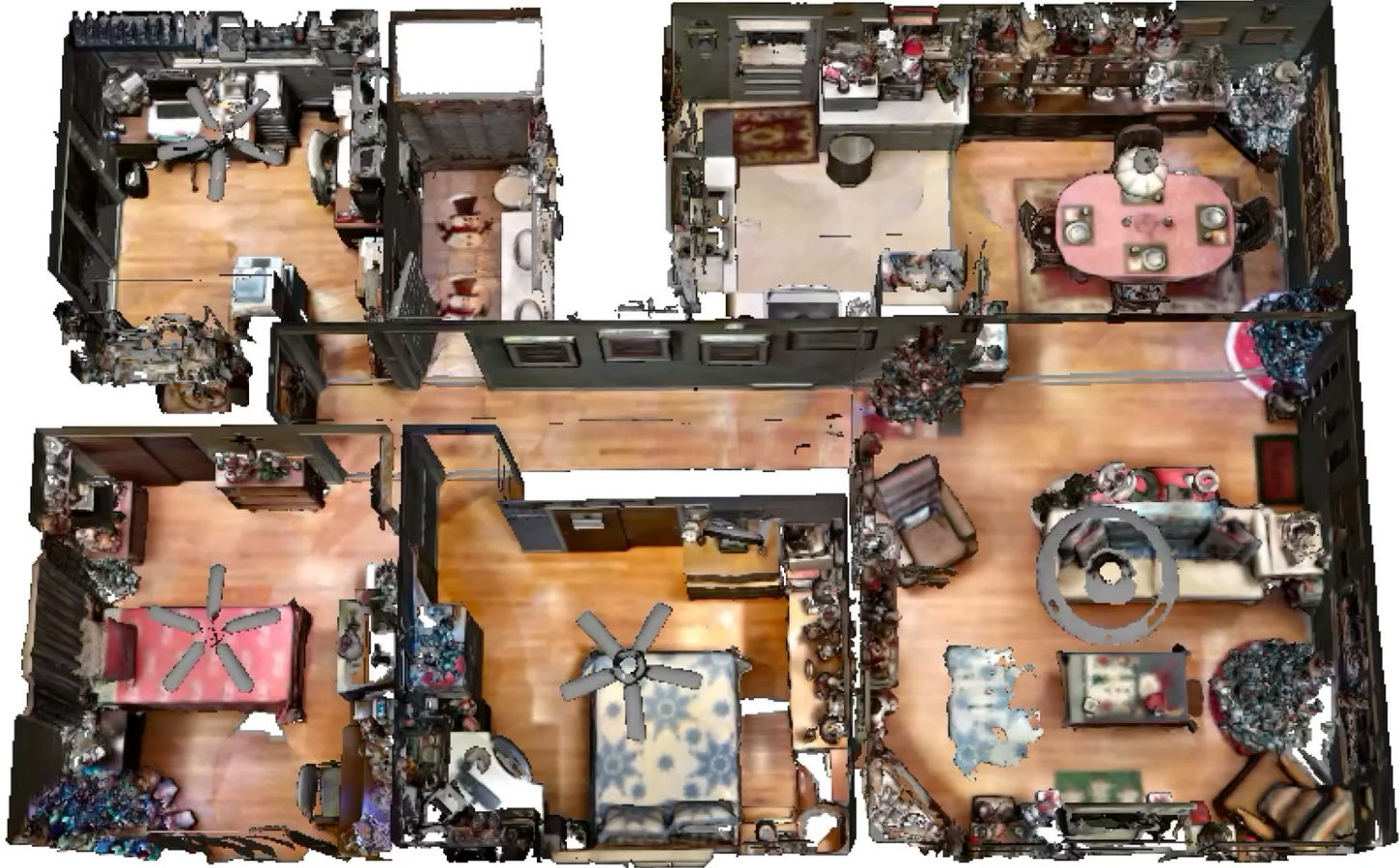
Image Queries

Given 3D Geometry

Interactive Demo

Open-vocabulary 3D Scene Exploration

Text queries:

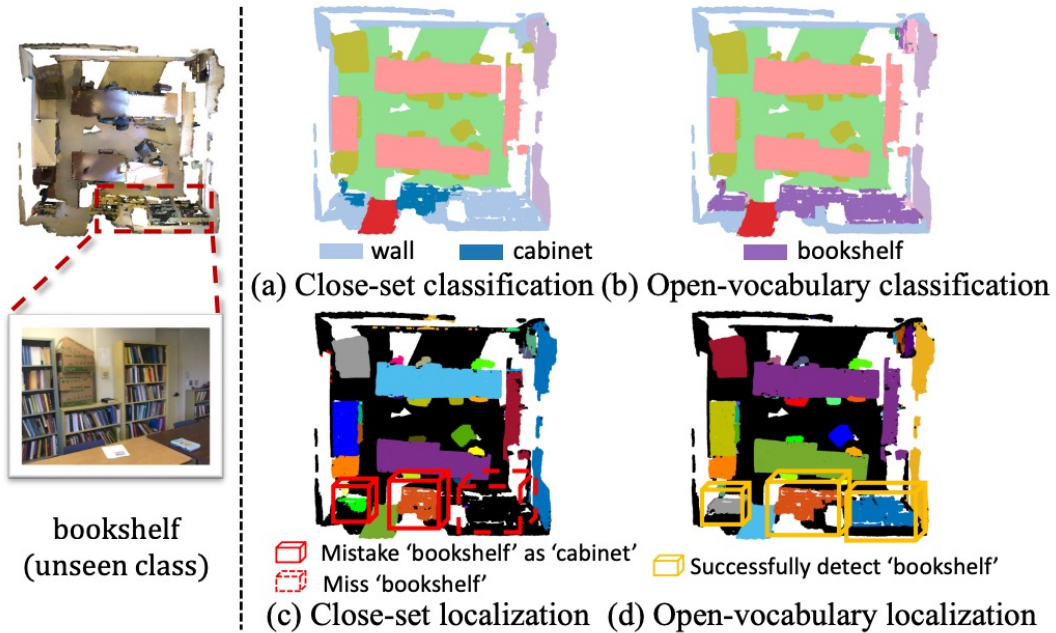


Take-home Messages

- + 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

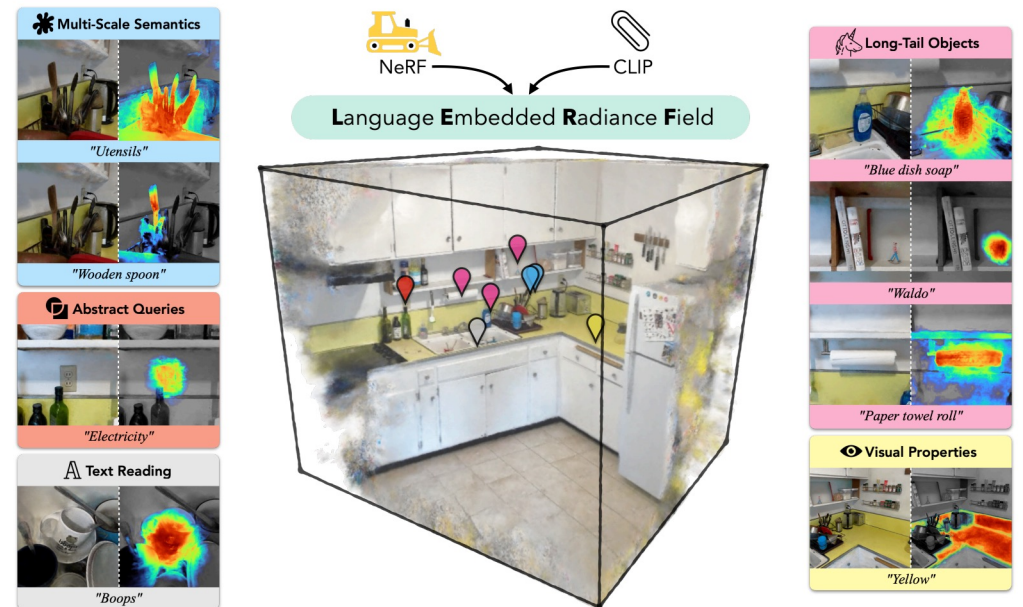
PLA

Ding*, Yang*, ..., Qi. CVPR 2023

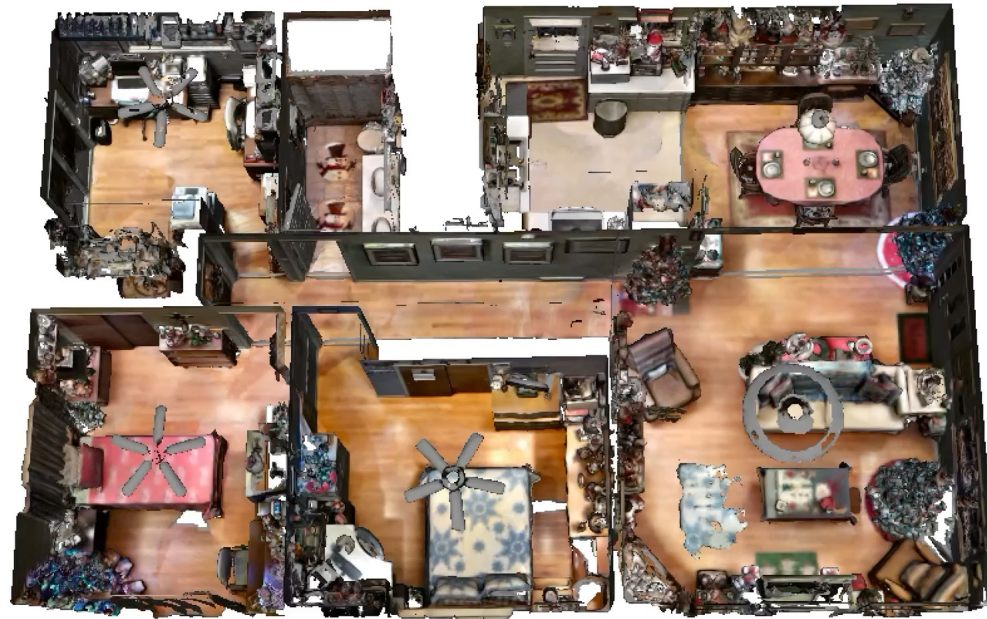


LERF

Kerr*, Kim*, et al. ICCV 2023



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**



清華大學

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



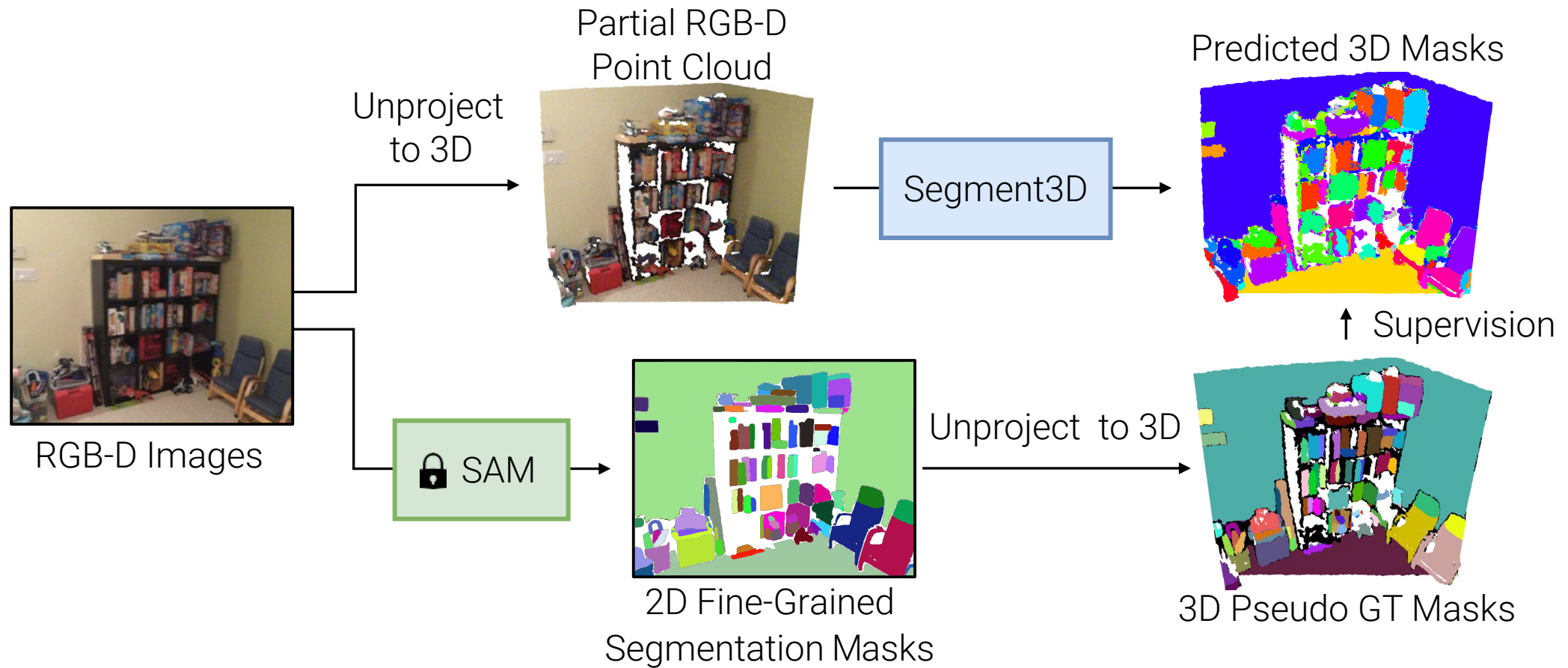
Gao Huang



Francis Engelmann

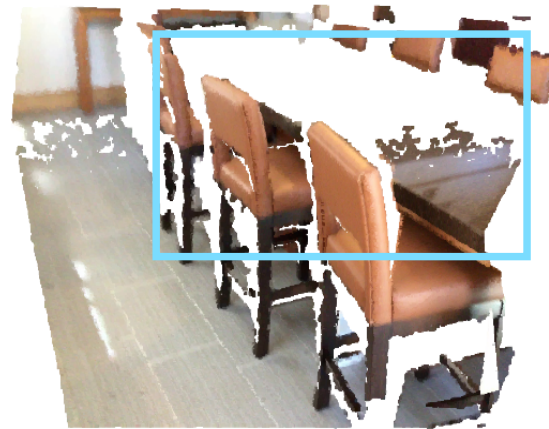
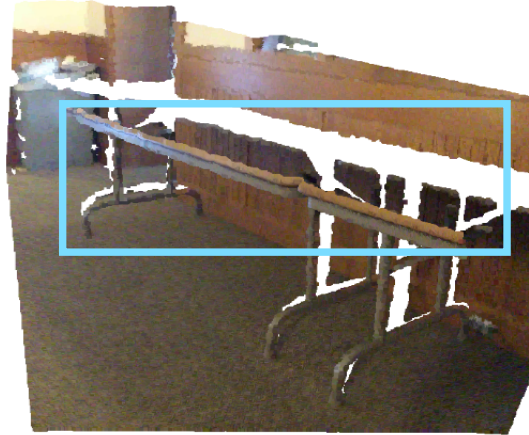
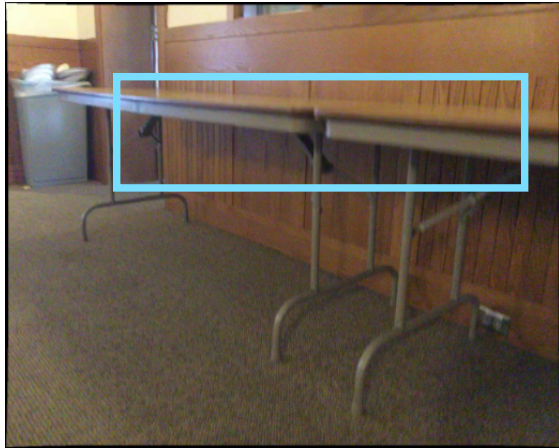
Segment3D

Stage 1: Pre-training on Partial Point Clouds



Segment3D

Domain Gap Between Partial and Full Point Clouds



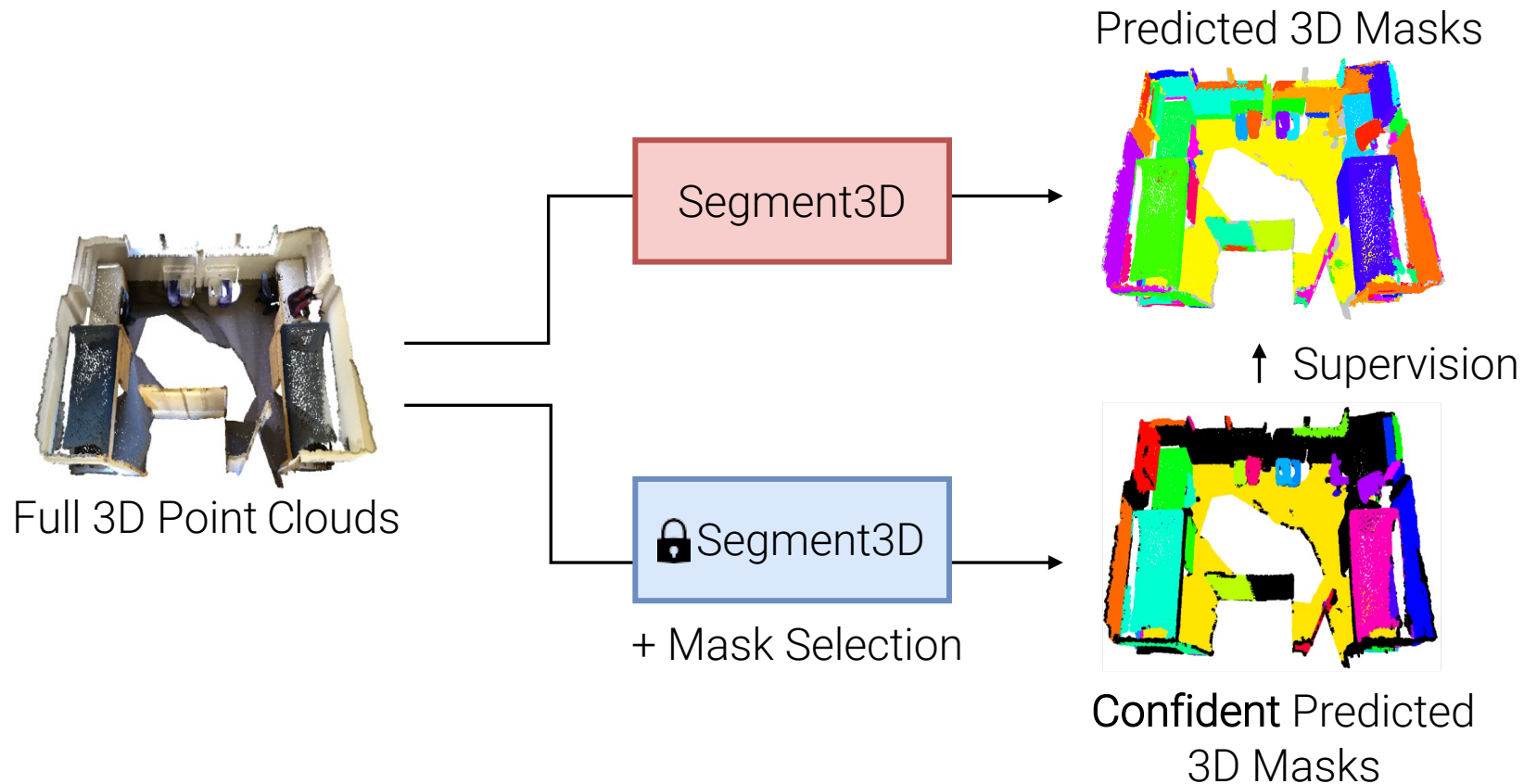
RGB Images

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

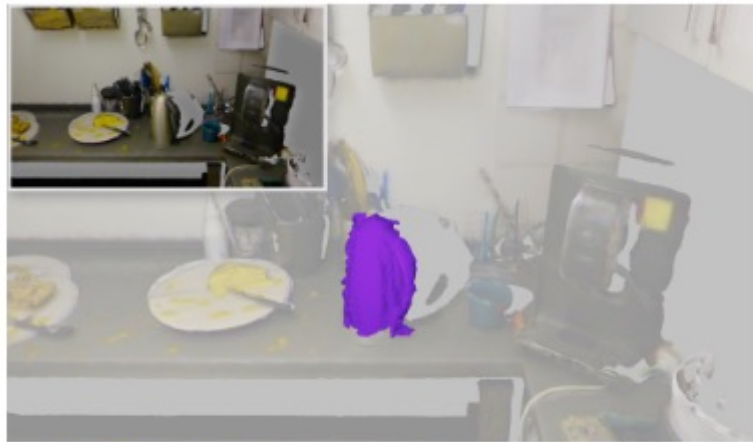
Model	Ground Truth Labels	<i>without post-processing</i>			<i>with post-processing</i>		
		AP	AP ₅₀	AP ₂₅	AP	AP ₅₀	AP ₂₅
SAM3D [49]	✗	3.9	9.3	22.1	8.4	16.1	30.0
Felzenszwalb <i>et al.</i> [17]	✗	5.8	11.6	27.2	–	–	–
Mask3D [40]	ScanNet200 [39]	8.7	15.5	27.2	14.3	21.3	29.9
Mask3D [40]	ScanNet [12]	9.4	16.8	28.7	15.4	22.7	31.6
Segment3D (Ours)	✗	12.0	22.7	37.8	19.0	29.7	41.6
		(+27.7%)	(+35.1%)	(+31.7%)	(+23.4%)	(+30.8%)	(+31.6%)

Effect of Two-Stage Training

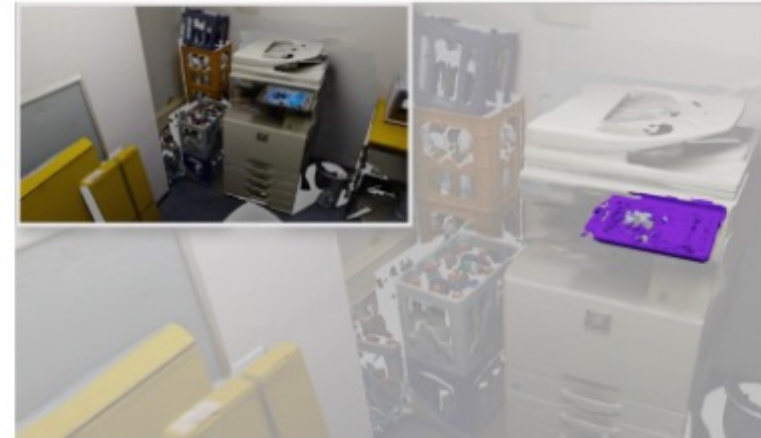
Training Stages	AP	AP₅₀	AP₂₅
Pre-Training (Stage 1)	7.4	15.2	31.2
+ Fine-Tuning (Stage 2)	12.0 (+62%)	22.7 (+49%)	37.8 (+21%)

Open-Vocabulary Segmentation

Mask3D



Segment3D



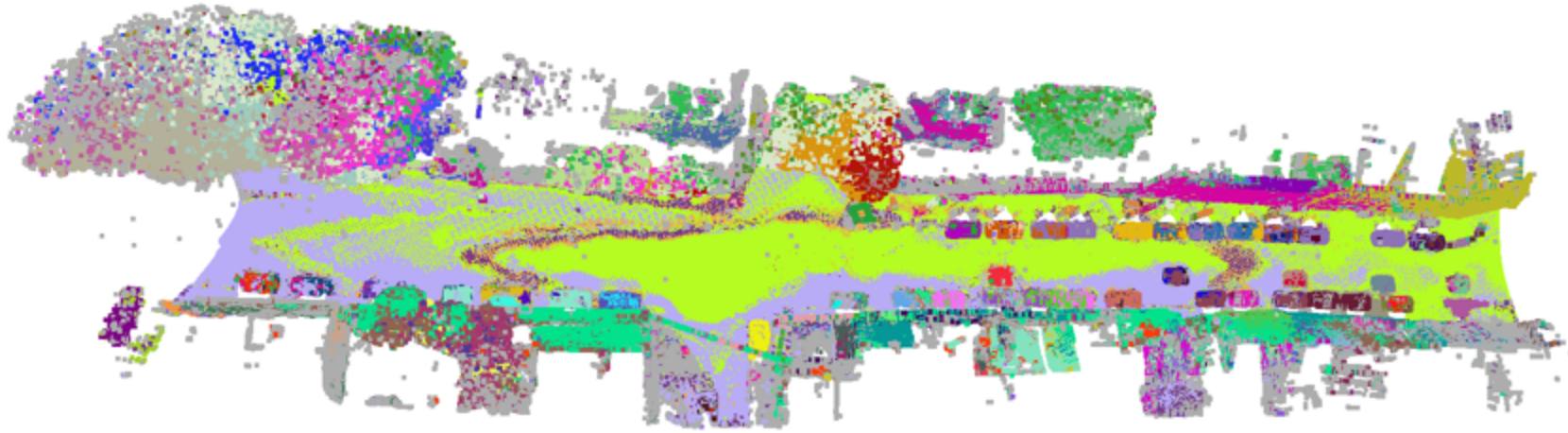
“a black eraser”

“kettle handle”

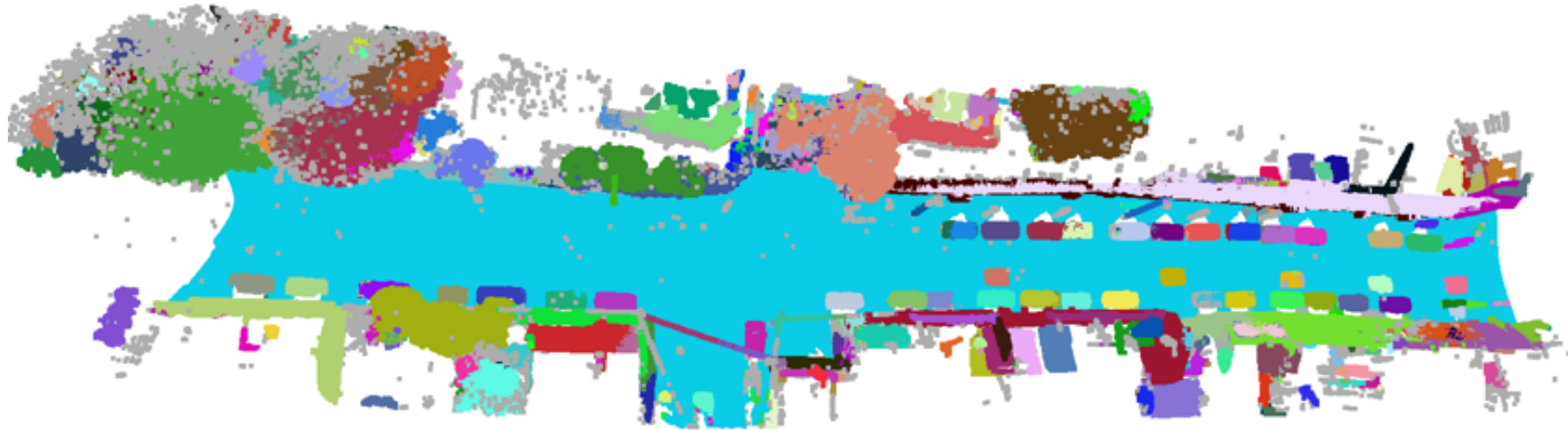
“copier control screen”

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

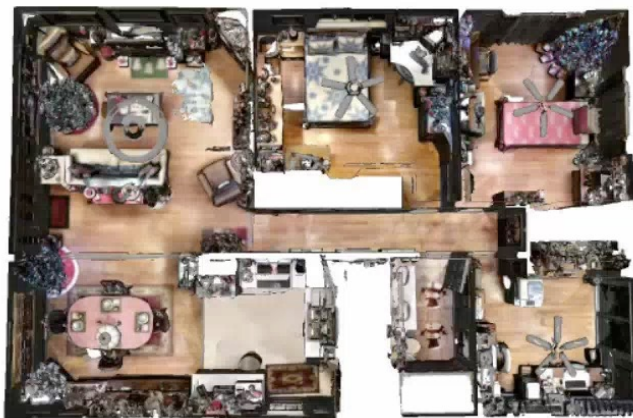
2D Foundation Models for 3D Vision Tasks

3D Reconstruction



NeRF *On-the-go*
(under review)

3D Scene Understanding



OpenScene
CVPR 2023



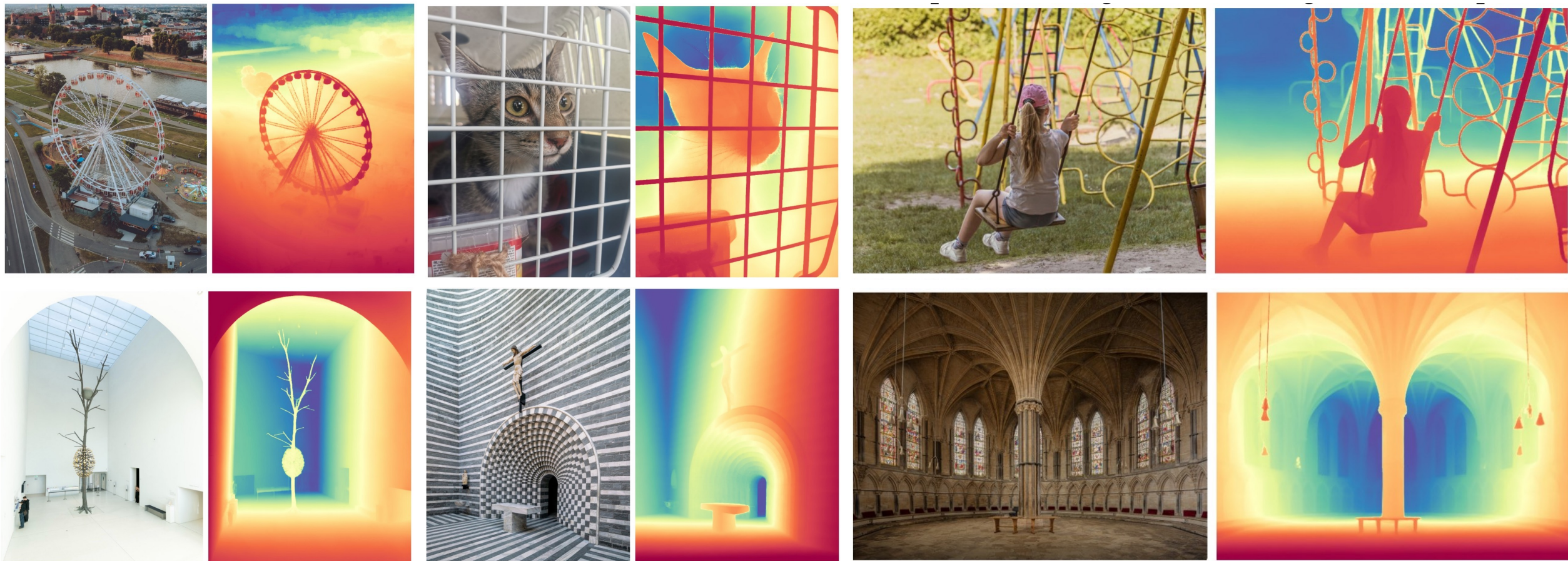
Segment3D
(under review)

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

So, what is next?

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**
- ...

**So far, we only talked about
2D foundation models...**

Current Interests

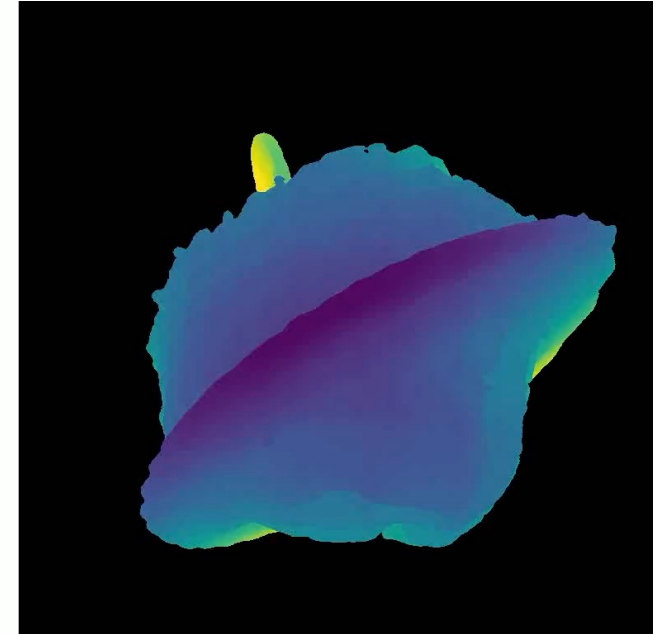
3D Foundation Models



DreamFusion
[ICLR'23]
1.5 Hours



Large Reconstruction Models
[ICLR'24]
5 Seconds!



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

Current Interests

3D Foundation Models for Large-Scale Scenes

- **What data?** RealEstate (80K videos), Ego4D (3000-hour video)...
- **What should be the representations?** Efficient & compact
- **What tasks?** Generation, reconstruction, understanding?
- **Many open questions:** How to learn from pure 2D inputs? How to jointly train with limited 3D data?



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,.....

Do we really need text to 3D?

How to inject 3D to help Sora?

What Sora cannot do?

“It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness”

2D Magic in a 3D World

Q&A