

How do NeRF and CLIP advance 3D Scene Reconstruction and Understanding?

Songyou Peng

ETH Zurich and Max Planck Institute for Intelligent Systems



MAX PLANCK INSTITUTE
FOR INTELLIGENT SYSTEMS



Chinese University of Hong Kong (CUHK), Shenzhen

March 07, 2023

Who Am I?

- 4th Year PhD Student
 - Marc Pollefeys
 - Andreas Geiger
- Internships during PhD
 - 2021: Michael Zollhoefer
 - 2022: Tom Funkhouser
- Open to chat!

ETH zürich

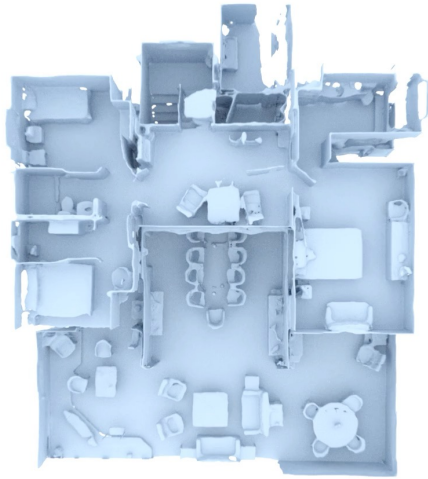


 **Meta**

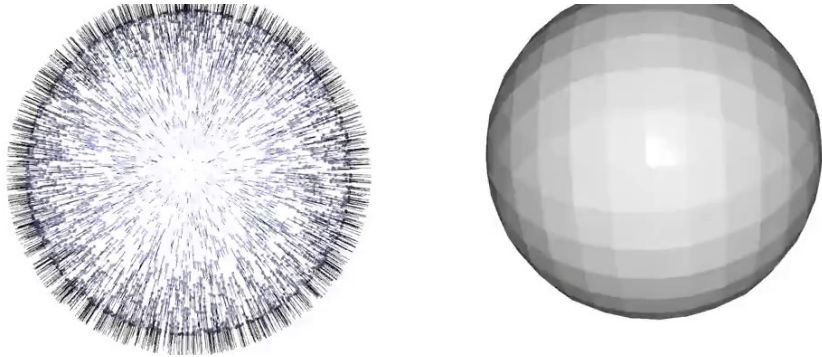



pengsongyou.github.io

My PhD Topics: Neural Scene Representations for 3D reconstruction, novel view synthesis, and SLAM



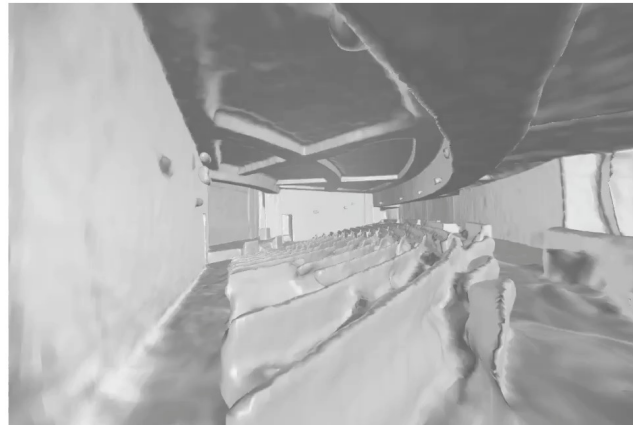
Convolutional Occupancy Networks
ECCV 2020 (Spotlight)



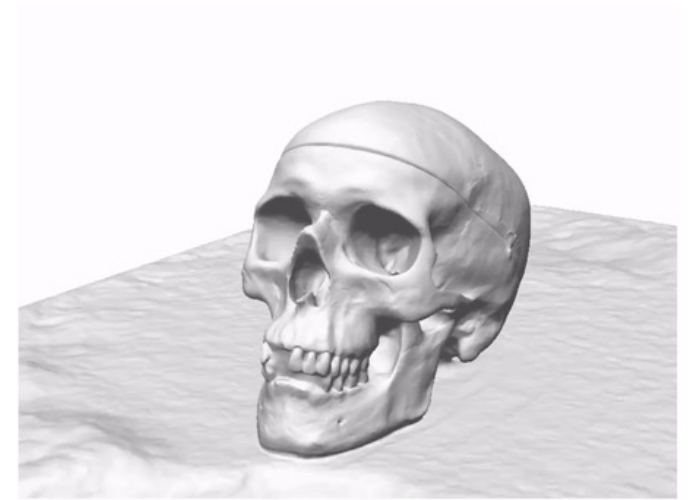
Shape As Points
NeurIPS 2021 (Oral)



KiloNeRF
ICCV 2021



Ours
MonoSDF
NeurIPS 2022



UNISURF
ICCV 2021 (Oral)

NICE-SLAM
CVPR 2022

How do NeRF and CLIP advance
3D Scene Reconstruction and Understanding?

How does NeRF advance 3D Scene Reconstruction?

How does CLIP advance 3D Scene Understanding?

How does NeRF advance 3D Scene Reconstruction?

How does CLIP advance 3D Scene Understanding?

NeRF is awesome!



Some problems still exist...

😓 Poor underlying geometry

😓 Camera poses needed

😊 MonoSDF

😊 NICE-SLAM

MonoSDF: Exploring Monocular Geometric Cues for Neural Implicit Surface Reconstruction



Zehao Yu



Songyou Peng



Michael Niemeyer

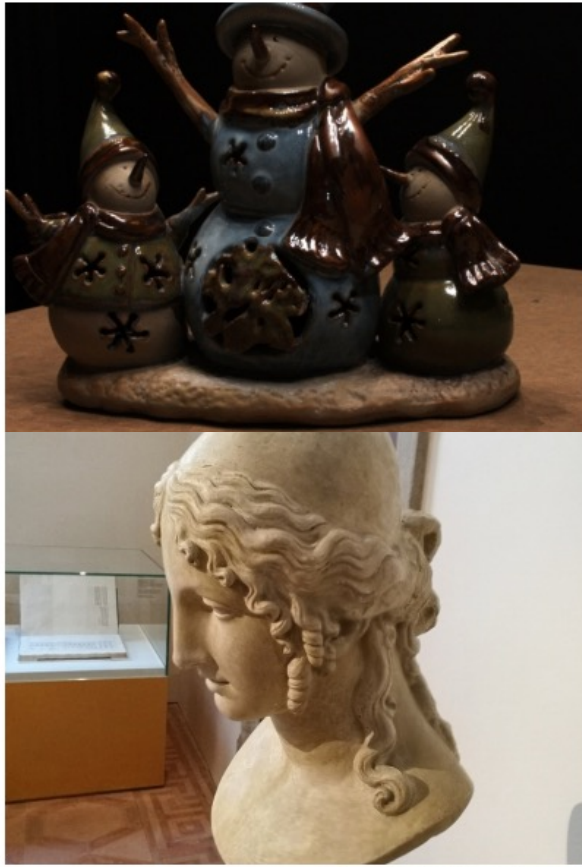


Torsten Sattler



Andreas Geiger

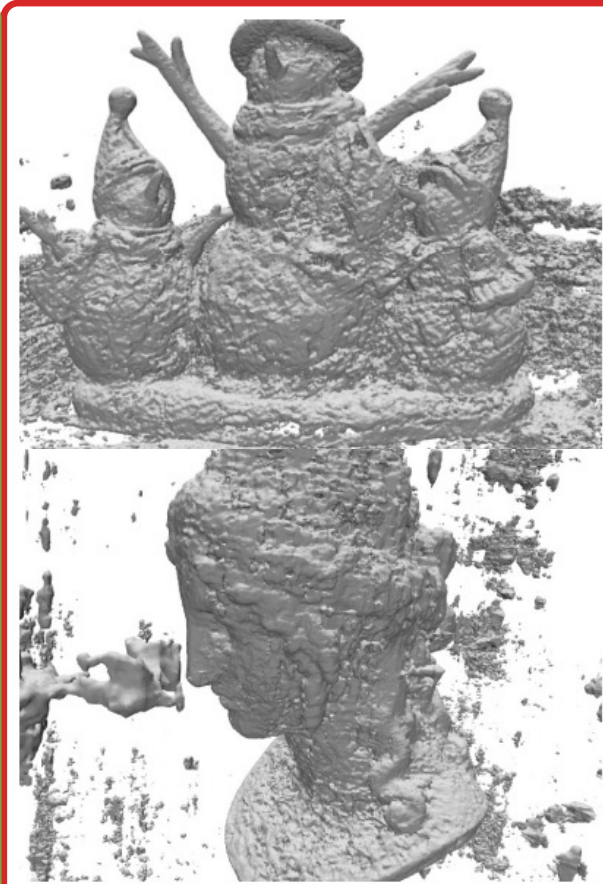
Neural Implicit Surfaces with Volume Rendering



RGB Images



VoISDF/NeuS/UNISURF



NeRF

- [1] Oechsle, Peng, Geiger: UNISURF: Unifying Neural Implicit Surfaces and Radiance Fields for Multi-View Reconstruction. ICCV, 2021
- [2] Wang, Liu, Liu, Theobalt, Komura, Wang: NeuS: Learning Neural Implicit Surfaces by Volume Rendering for Multi-view Reconstruction. NeurIPS, 2021
- [3] Yariv, Gu, Kasten, Lipman: Volume rendering of neural implicit surfaces. NeurIPS, 2021

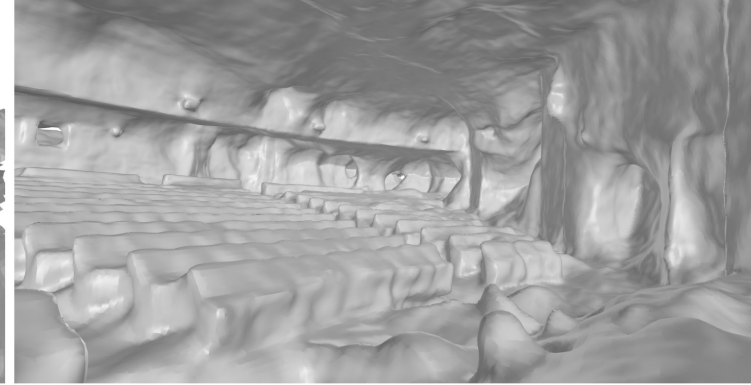
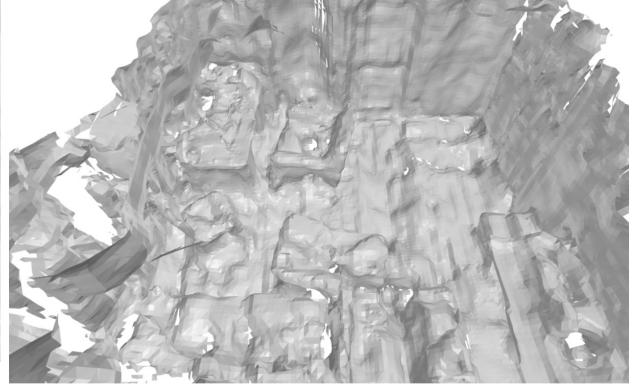
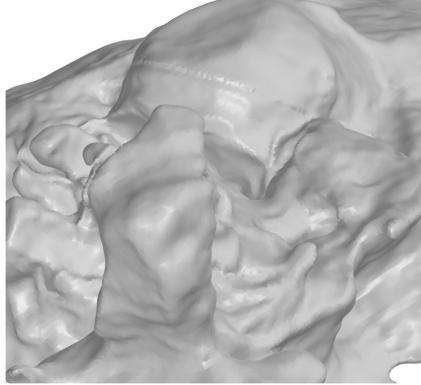
Neural Implicit Surfaces with Volume Rendering

DTU (3 views)

ScanNet (464 views)

Tanks & Temples (298 views)

VoISDF



- Fails with sparse input views
- Poor results in large-scale indoor scenes

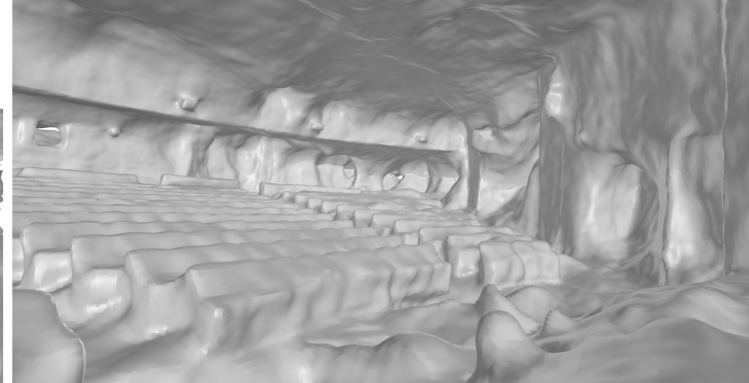
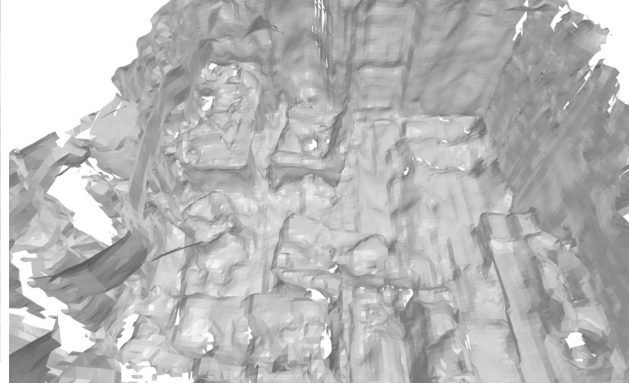
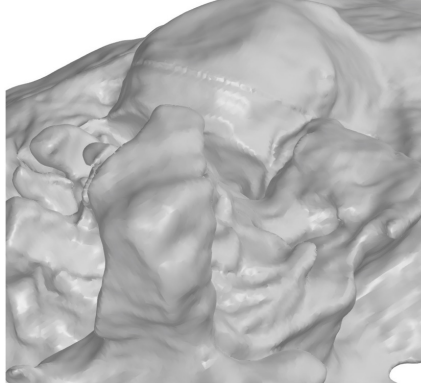
Neural Implicit Surfaces with Volume Rendering

DTU (3 views)

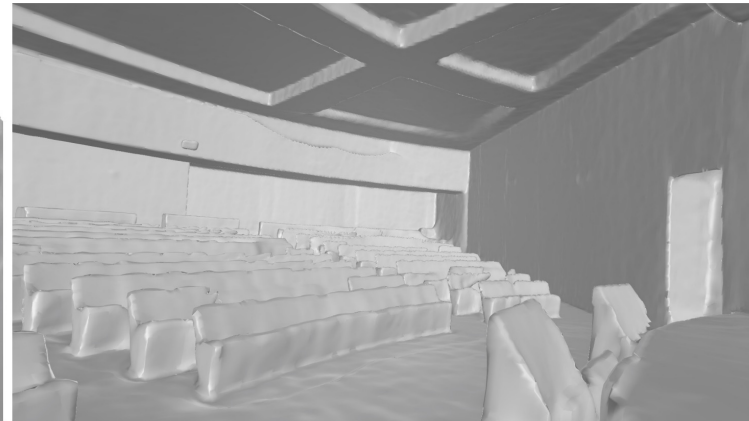
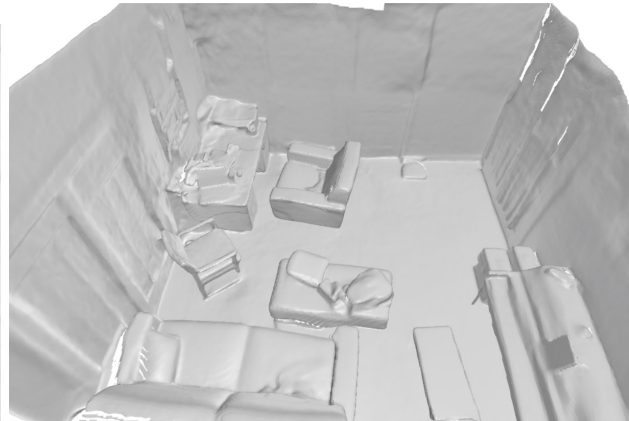
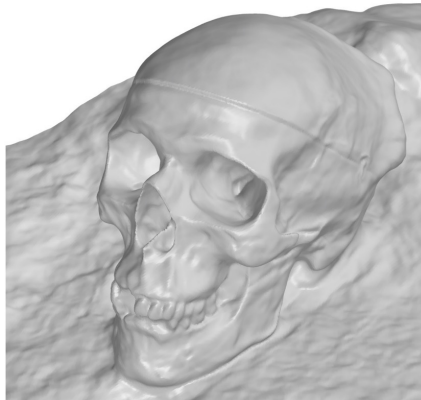
ScanNet (464 views)

Tanks & Temples (298 views)

VoISDF

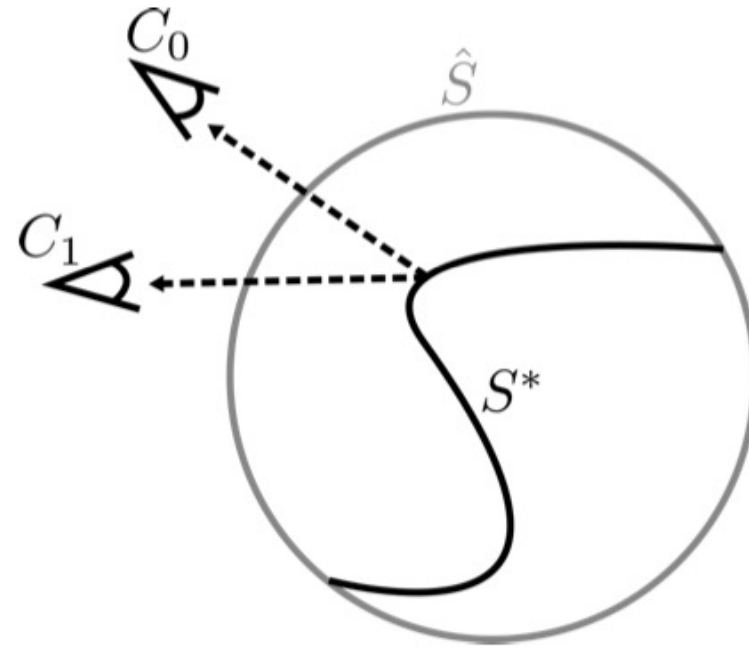


**MonoSDF
(Ours)**



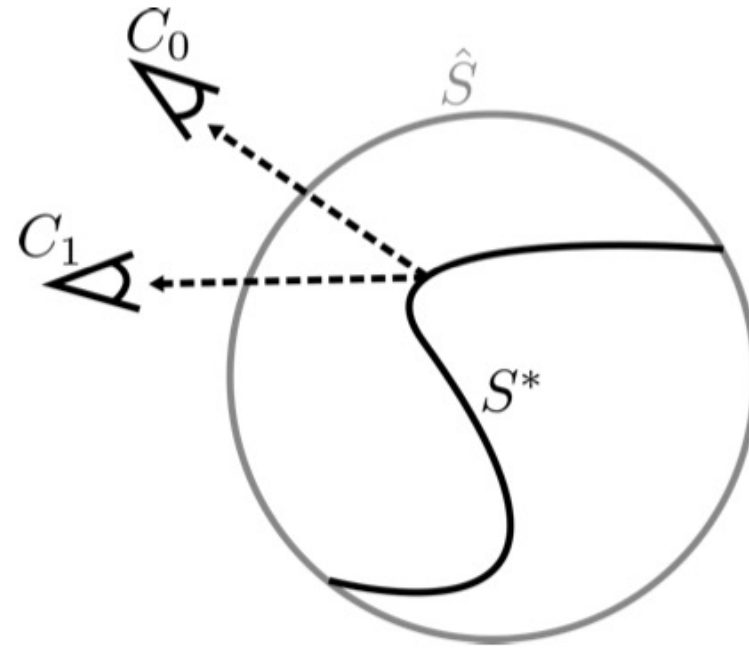
- + Manage to reconstruct with sparse views
- + Nice 3D reconstruction in large-scale indoor scenes

Shape-Appearance Ambiguity



There exists an infinite number of photo-consistent explanations for input images!

Shape-Appearance Ambiguity

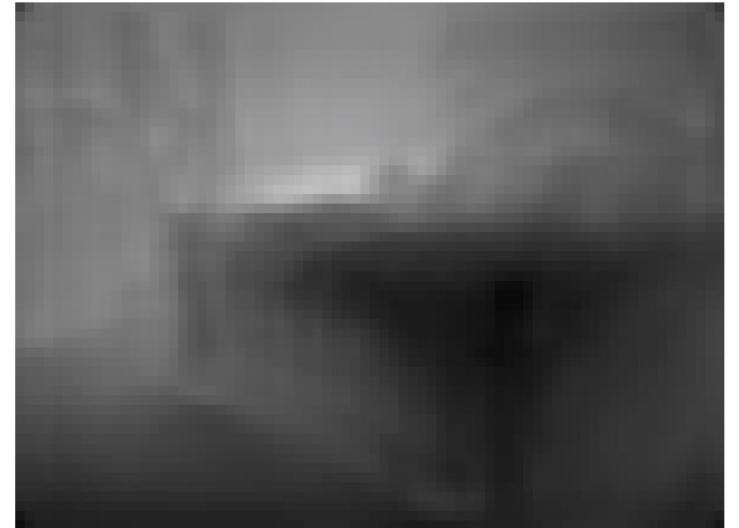
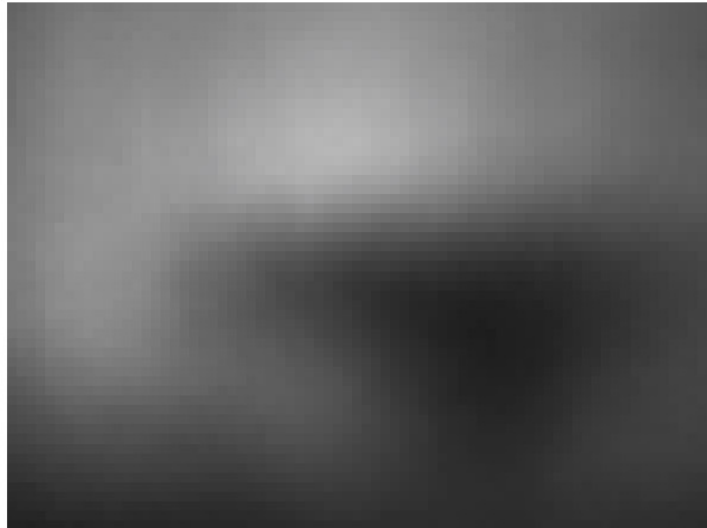
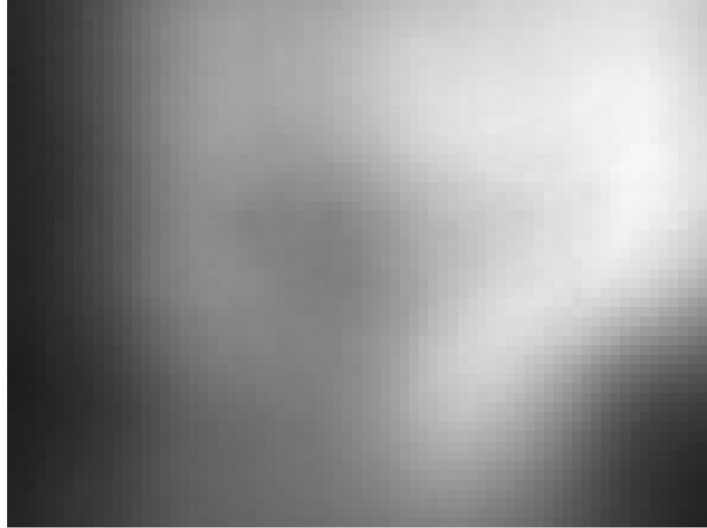
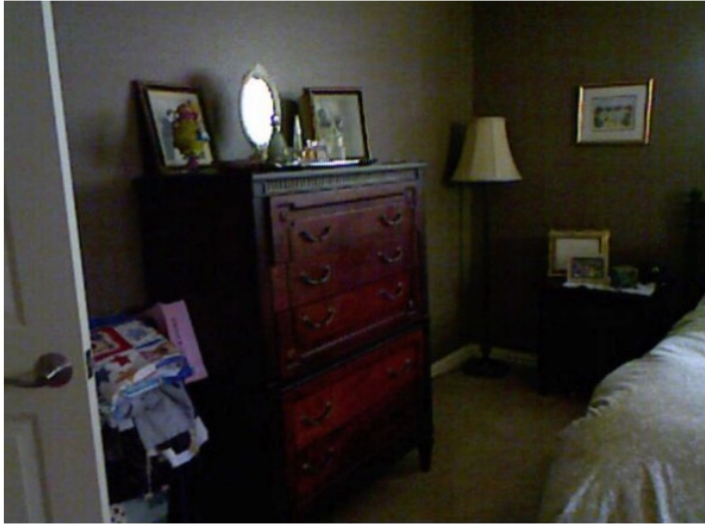


There exists an infinite number of photo-consistent explanations for input images!

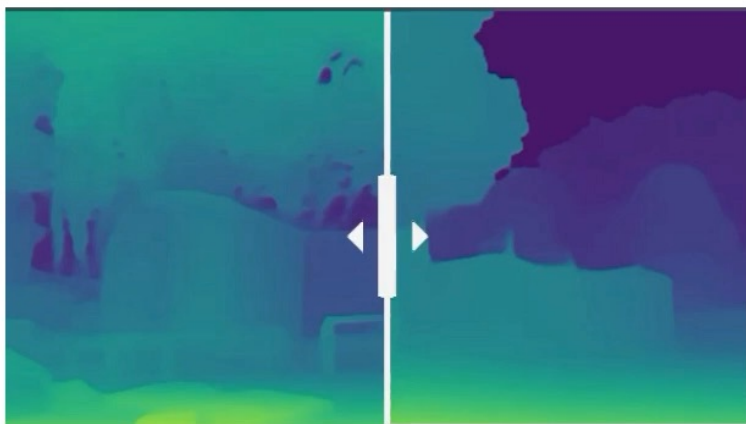


Exploit monocular geometric priors

Depth Map Prediction from a Single Image

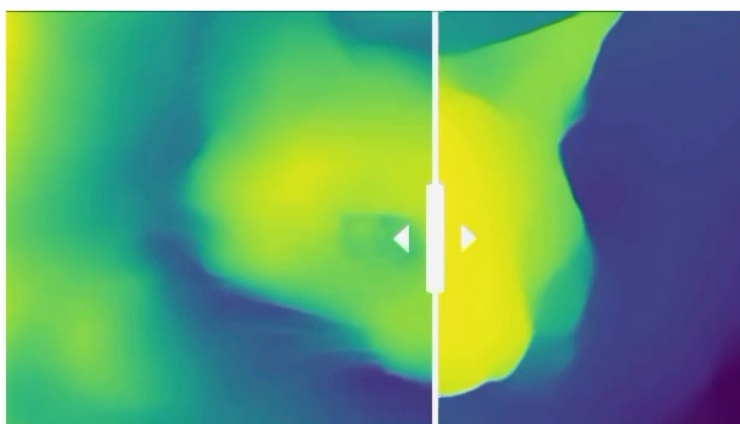


Omnidata



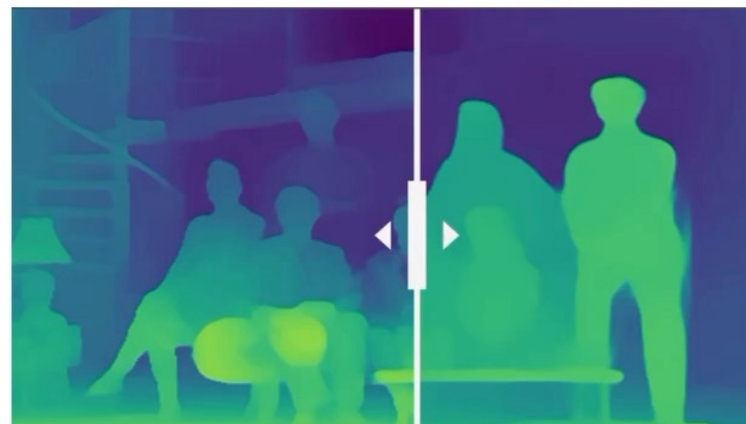
Ours

**MiDaS
DPT-Hybrid**



Ours

**MiDaS
DPT-Hybrid**

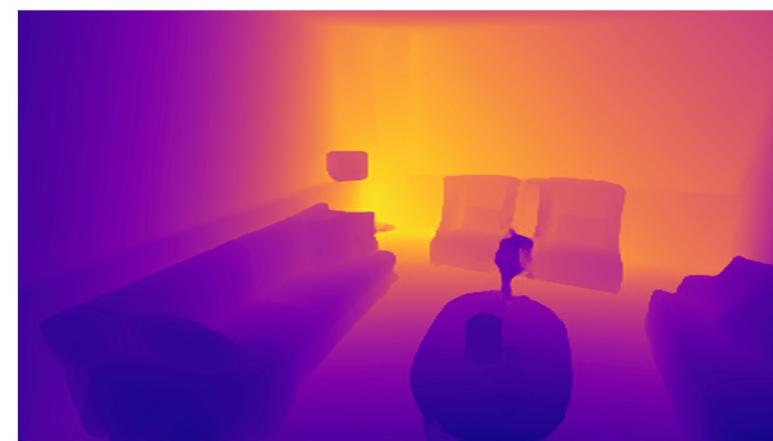
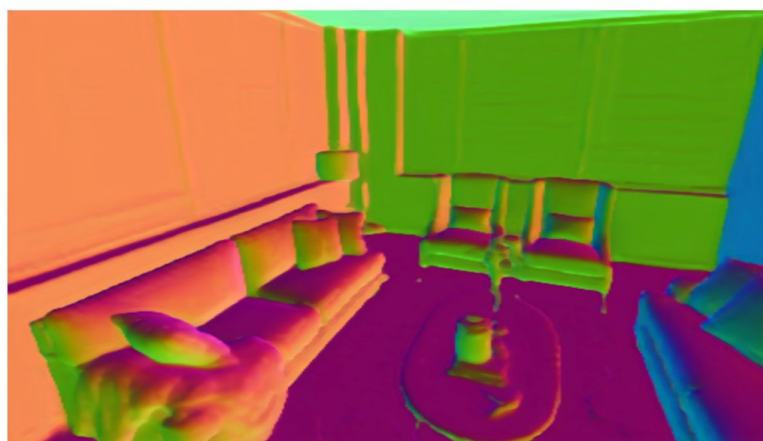
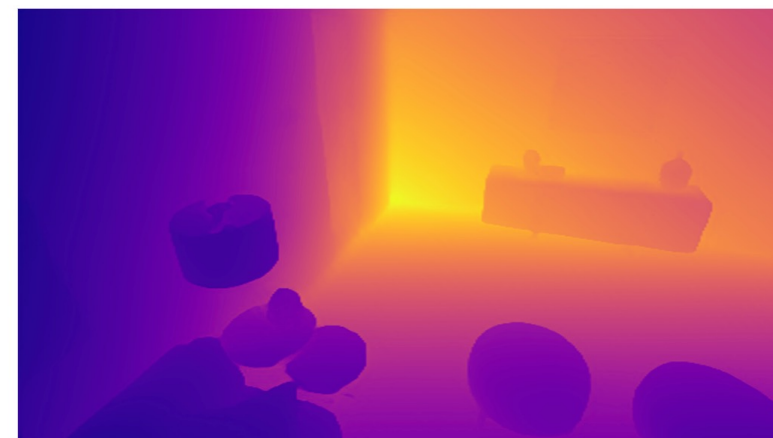
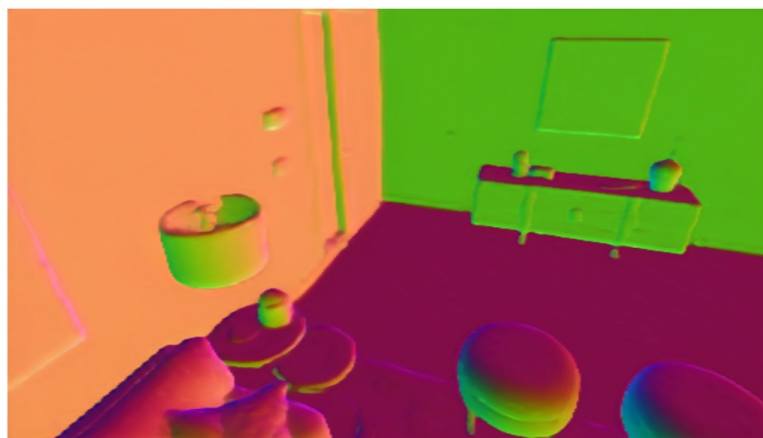


Ours

**MiDaS
DPT-Hybrid**

[Ranftl et al. 2021]

Omnidata



RGB Image

Omnidata Normal

Omnidata Depth

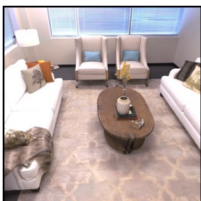
MonoSDF



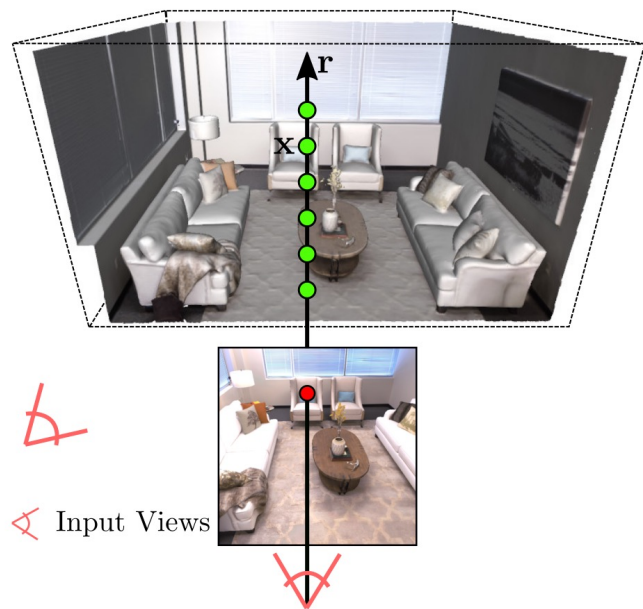
MonoSDF



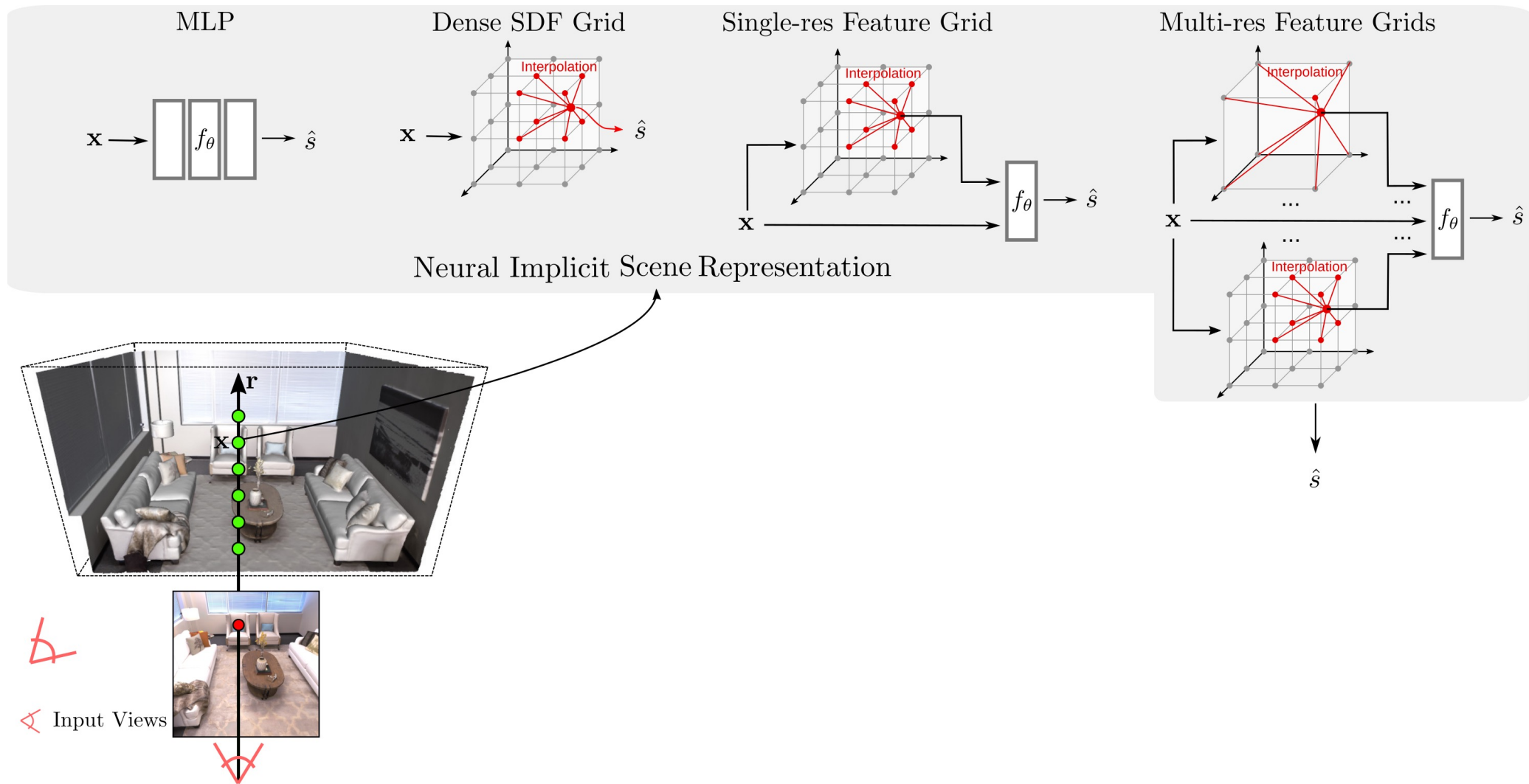
Input Views



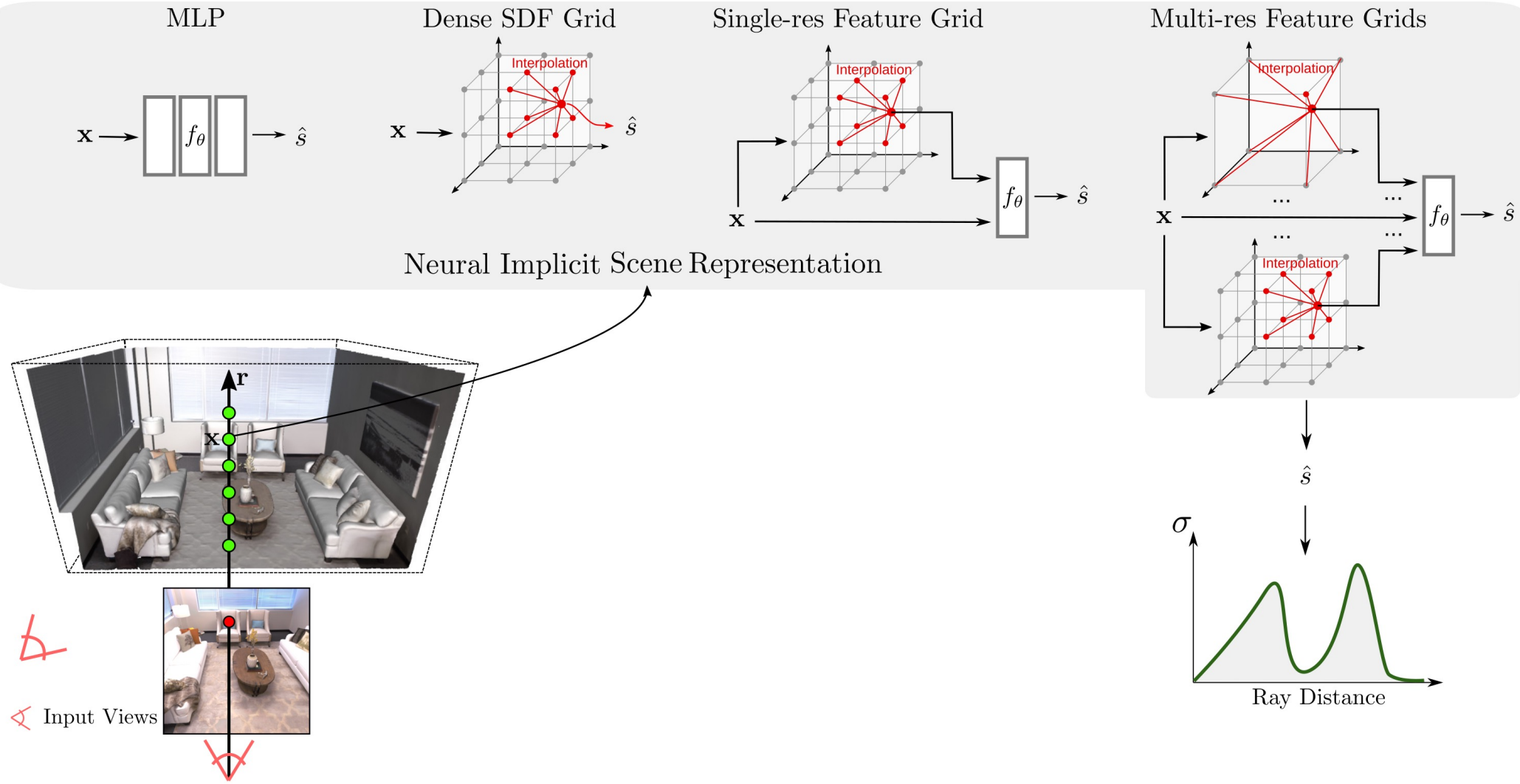
MonoSDF



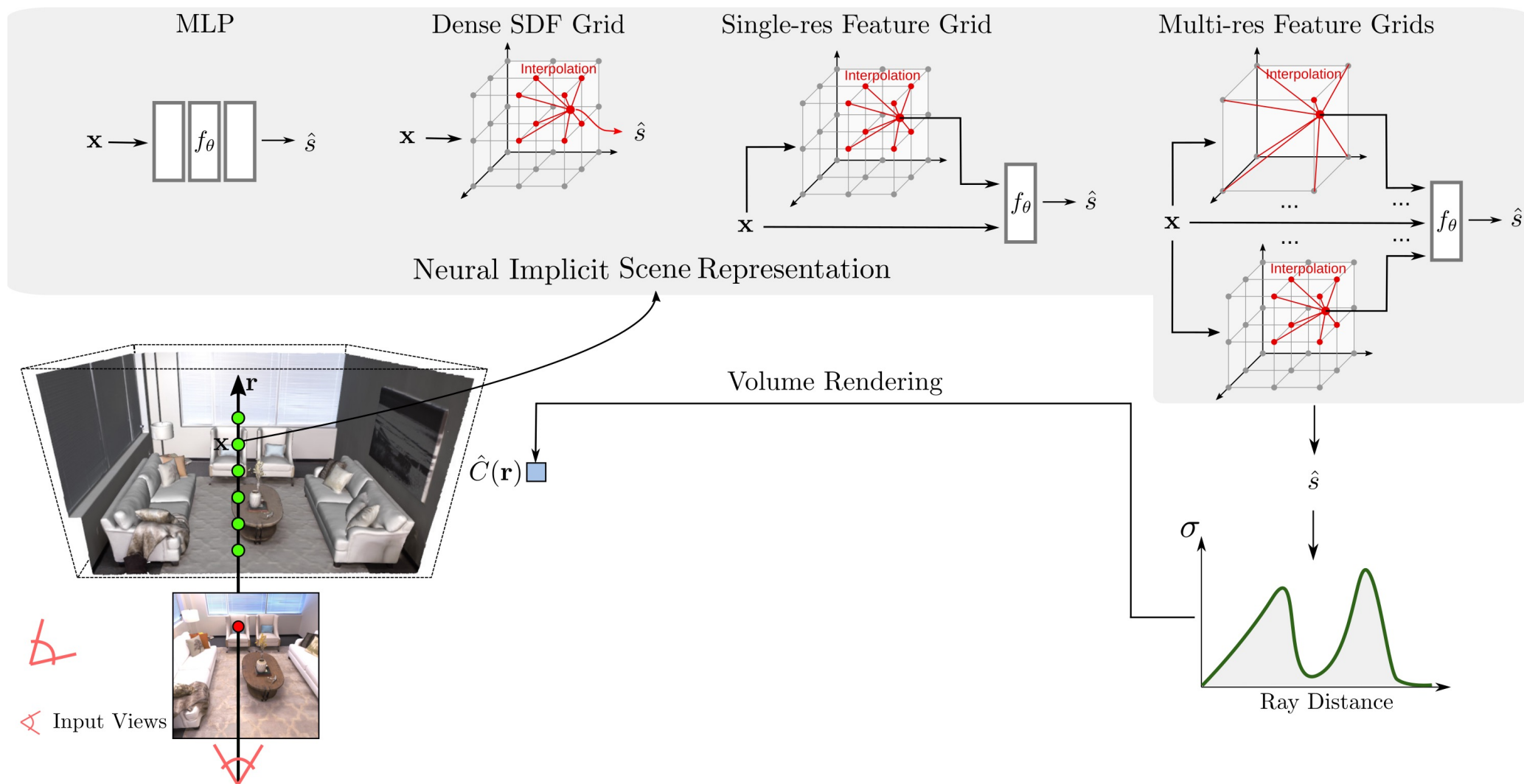
MonoSDF



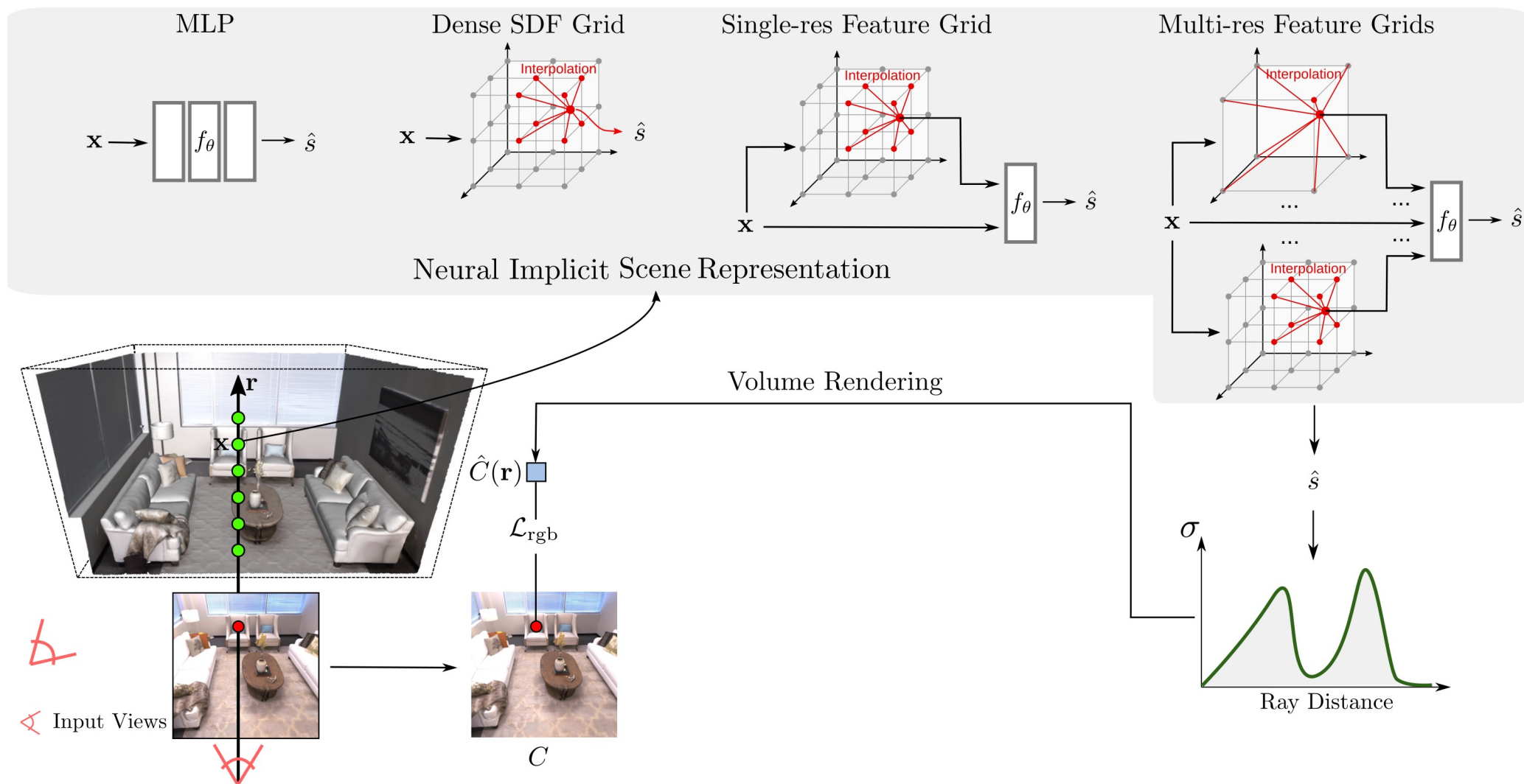
MonoSDF



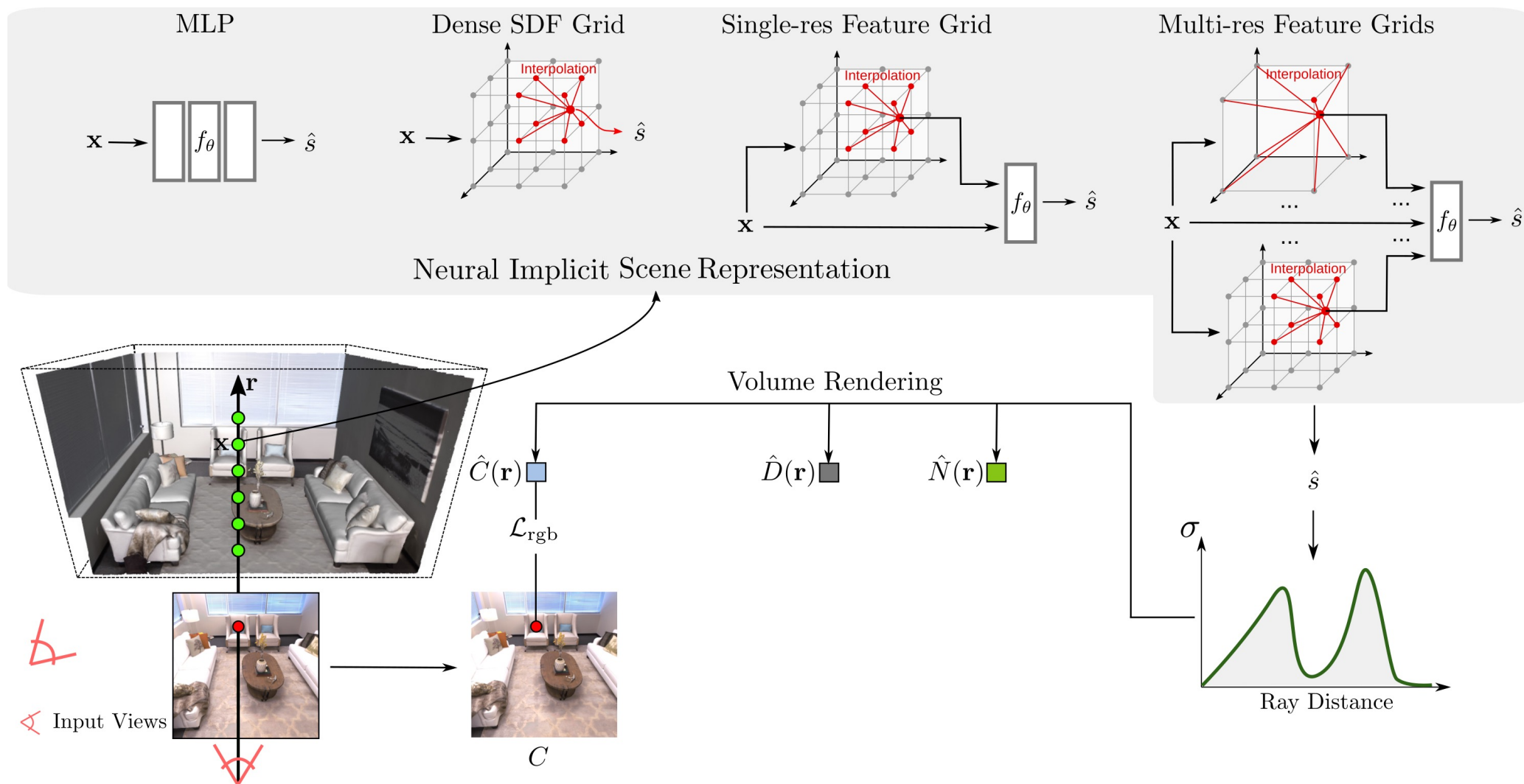
MonoSDF



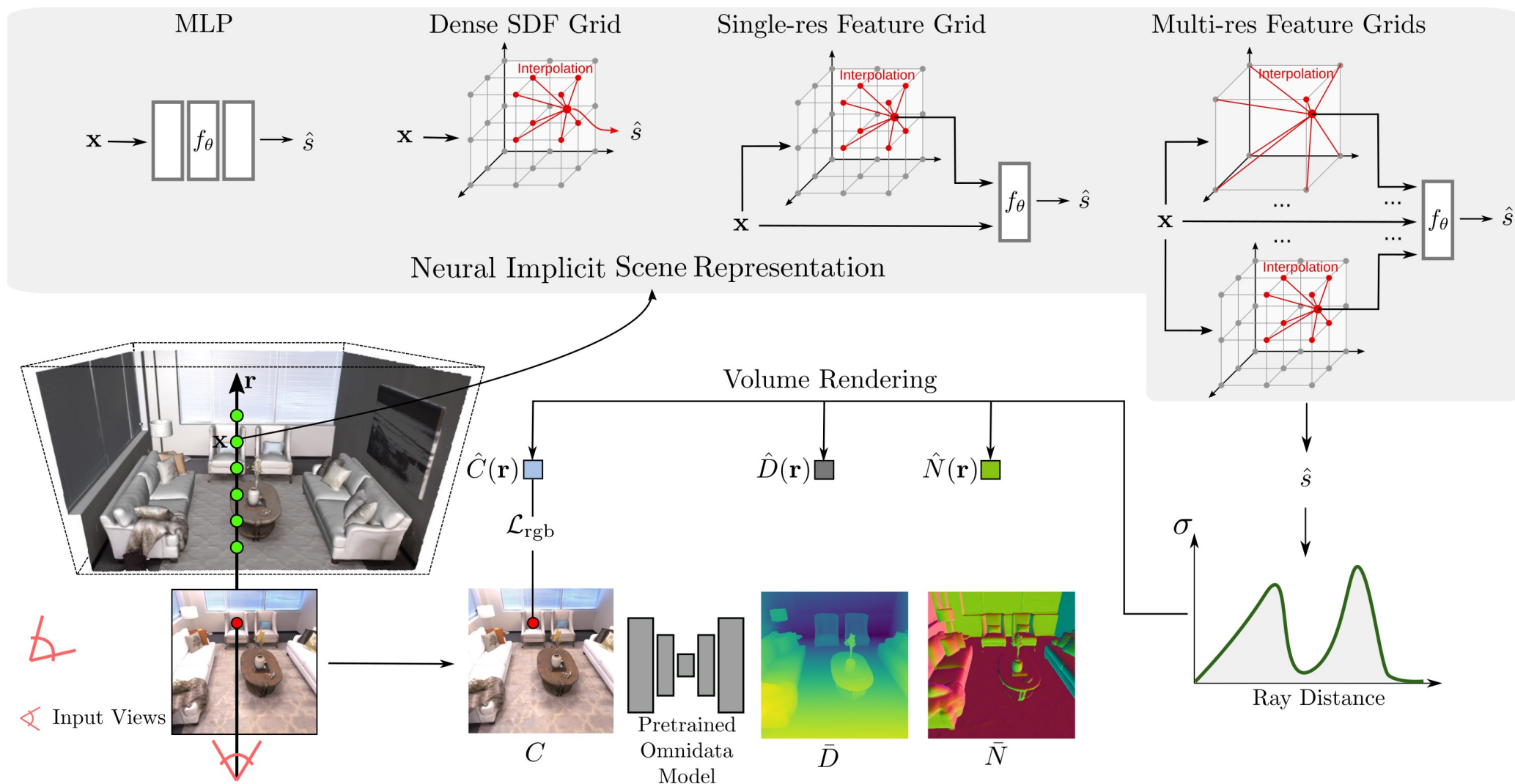
MonoSDF



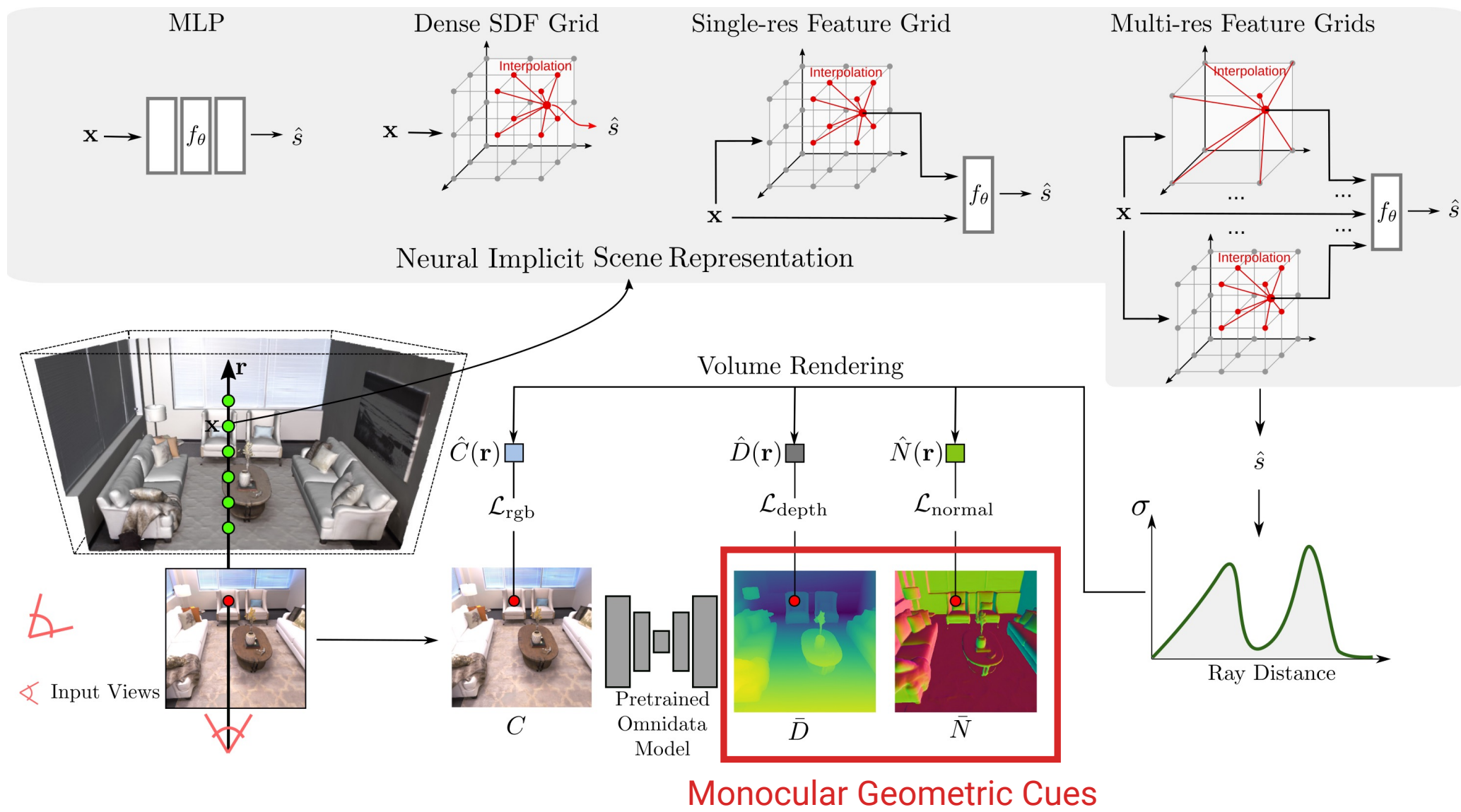
MonoSDF



MonoSDF

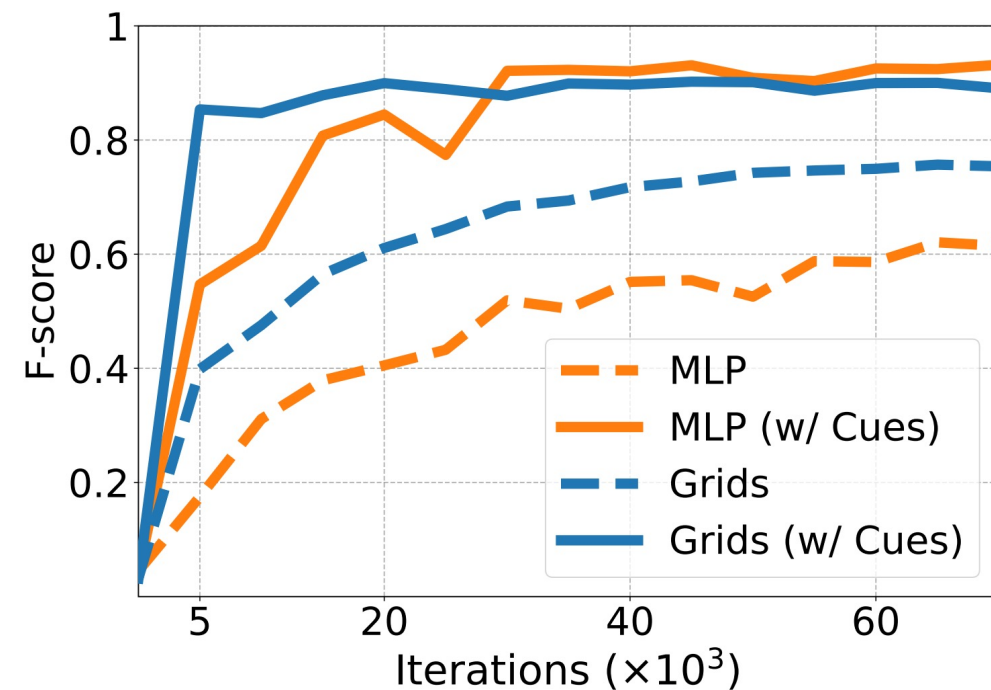


MonoSDF



Ablation Study

		Normal C. \uparrow	Chamfer- $L_1 \downarrow$	F-score \uparrow
MLP	No Cues	86.48	6.75	66.88
	Only Depth	90.56	4.26	76.42
	Only Normal	91.35	3.19	85.84
	Both Cues	92.11	2.94	86.18
Multi-Res. Grids	No Cues	87.95	5.03	78.38
	Only Depth	90.87	3.75	80.32
	Only Normal	89.90	3.61	81.28
	Both Cues	90.93	3.23	85.91



- ! Monocular cues improve reconstruction results significantly
- ! Combining **depth & normal** leads to best performance
- ! Monocular cues can improve **convergence speed**

Baseline Comparisons on ScanNet

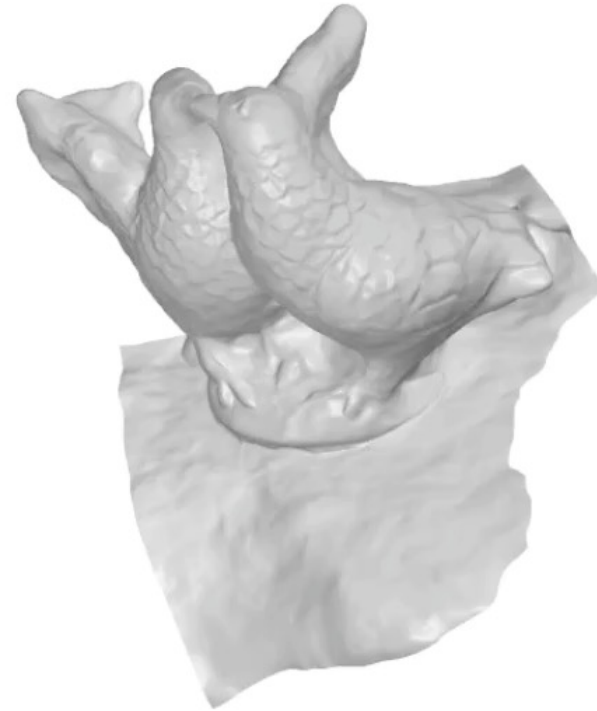
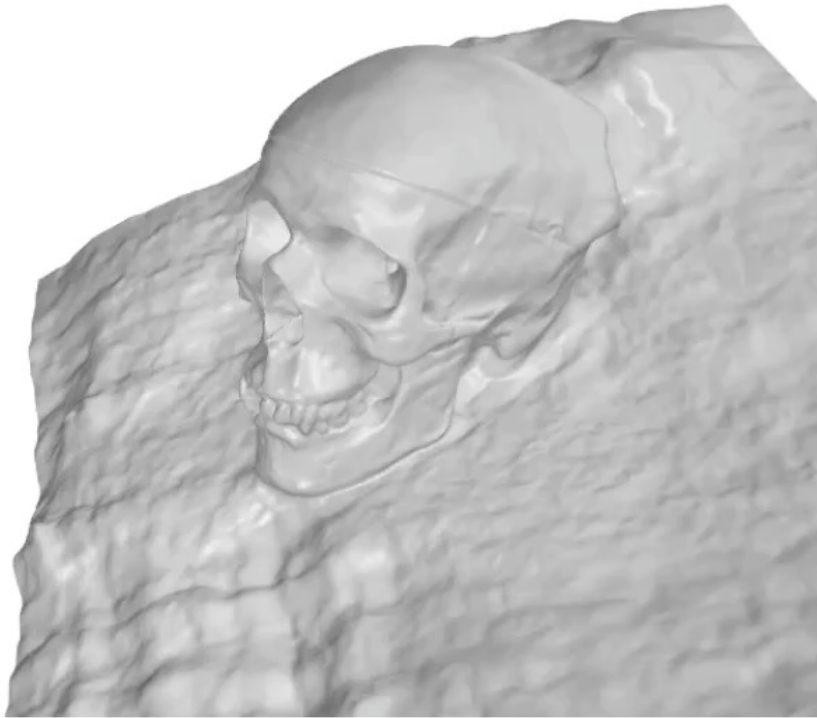


Ours

Multi-Res. Feature Grids with High-Res. Cues



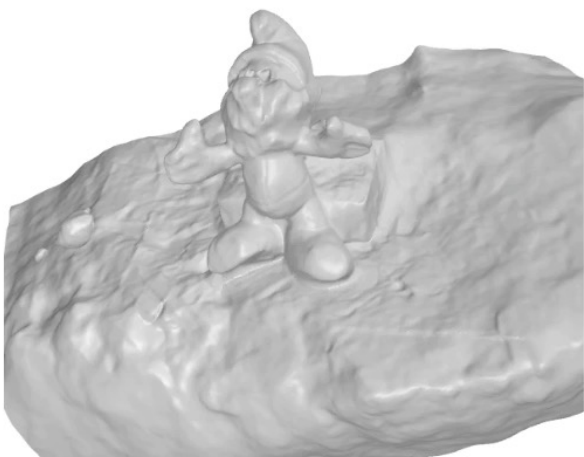
Baseline Comparisons on DTU (3-views)



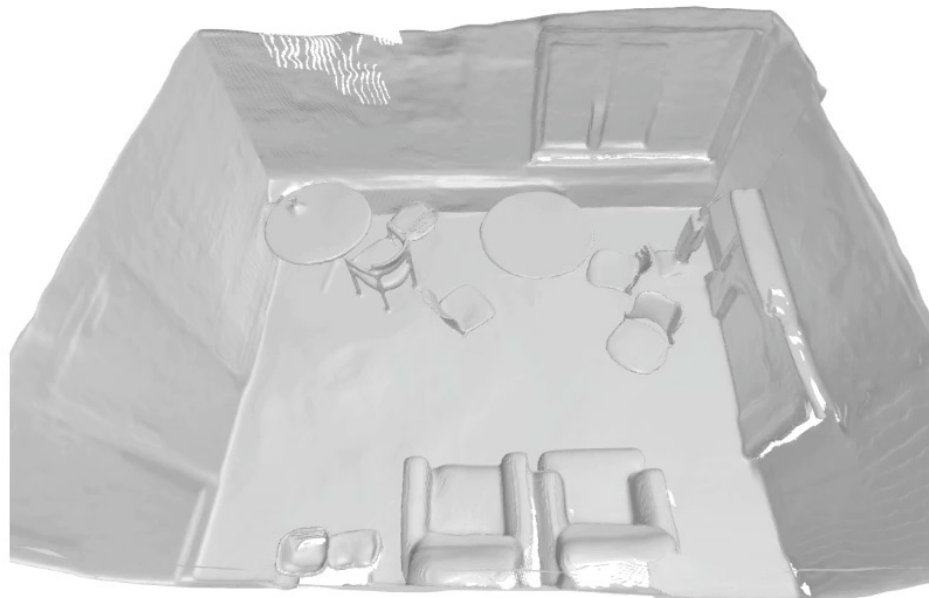
Ours

Take-home Message

<https://niujinshuchong.github.io/monosdf/>



DTU (3 views)



ScanNet



Tanks and Temples

- ! Monocular cues improve reconstruction results and speed up optimization
- ! Inspire Haiwen & Dan's ICLR 2023 paper GOOD ☺
- ! Limitation: Still require camera poses given :(



NICE-SLAM

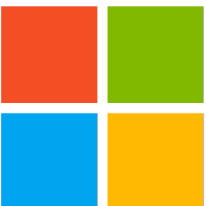
Neural Implicit Scalable Encoding for SLAM

CVPR 2022

Zihan Zhu* Songyou Peng* Viktor Larsson Weiwei Xu Hujun Bao
Zhaopeng Cui Martin R. Oswald Marc Pollefeys

* Equal Contributions

ETH zürich



RGB-D Sequences



40x Speed

iMAP

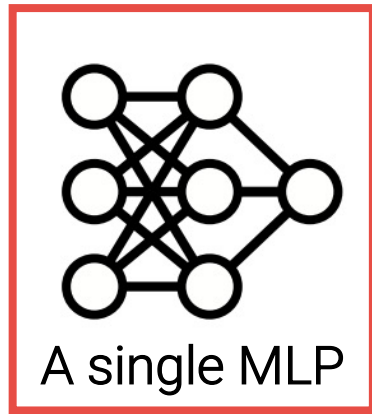
[Sucar et al., ICCV'21]



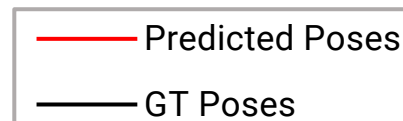
First neural implicit-based **online** SLAM system

iMAP

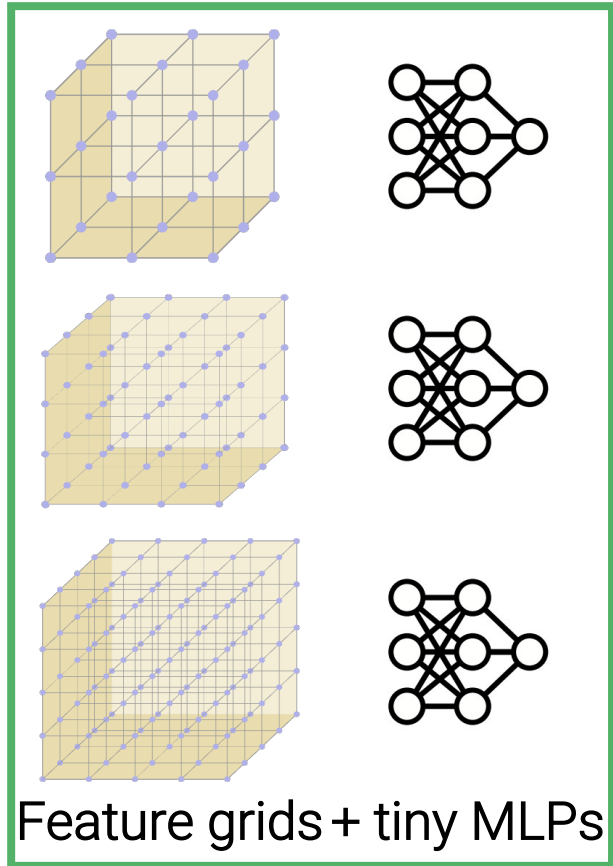
[Sucar et al., ICCV'21]



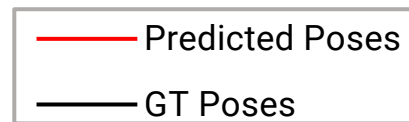
- Fail when scaling up to larger scenes
- Global update → Catastrophic forgetting
- Slow convergence



NICE-SLAM

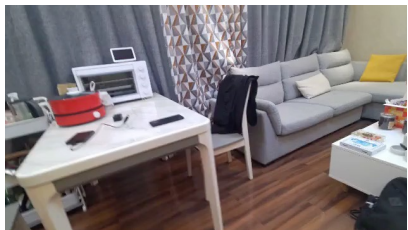
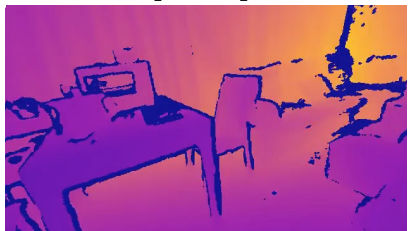


- + Applicable to large-scale scenes
- + Local update → No forgetting problem
- + Fast convergence

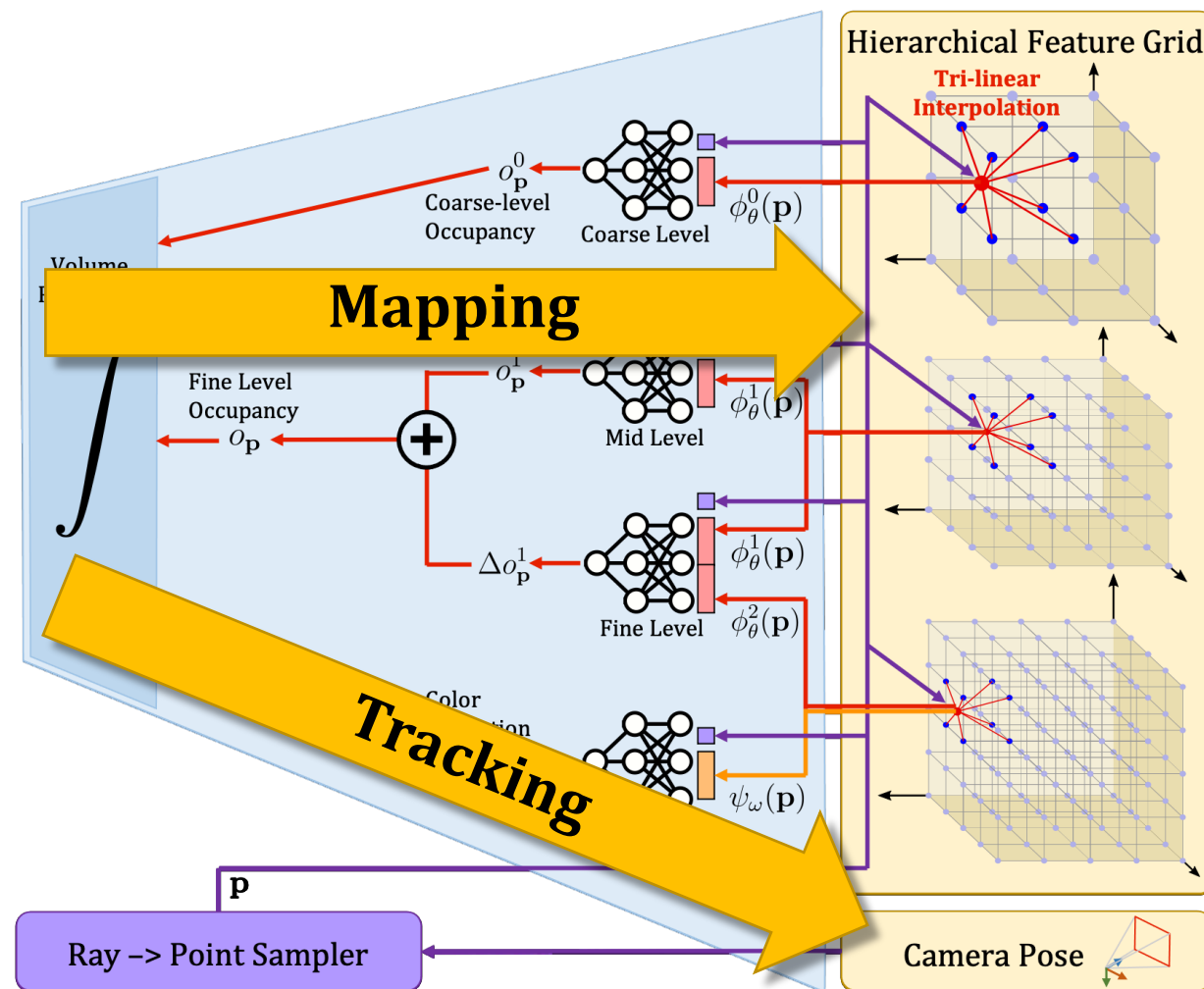


Pipeline

Input Depth



Input RGB



Results

iMAP*

(our re-implementation of iMAP)

NICE-SLAM

4x Speed

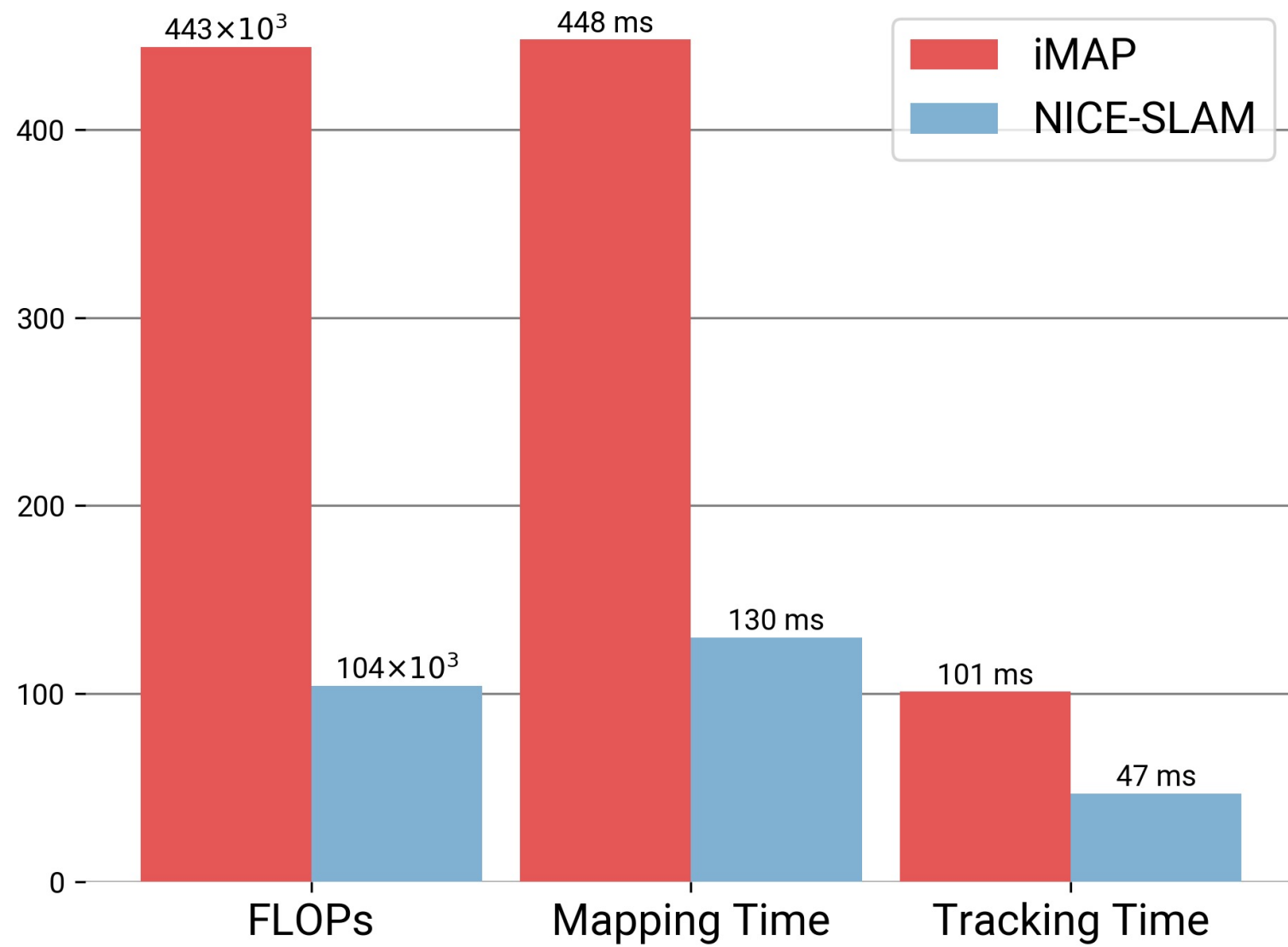
— Predicted Poses
— GT Poses

iMAP*

(our re-implementation of iMAP)

NICE-SLAM

10x Speed



Take-home Message

- A NICE online implicit SLAM system for indoor scenes
- Hierarchical feature grids + a tiny MLP seems to be a trend!
 - Instant-NGP [SIGGRAPH'22 Best Paper]

Limitations

- Requires depths as input
- Only bounded scenes
- Still not real-time

NICER-SLAM: Neural Implicit Scene Encoding for RGB SLAM

Zihan Zhu^{1*}

Songyou Peng^{1,2*}

Viktor Larsson³

Zhaopeng Cui⁴

Martin R. Oswald^{1,5}

Andreas Geiger⁶

Marc Pollefeys^{1,7}

¹ETH Zürich

²MPI for Intelligent Systems, Tübingen

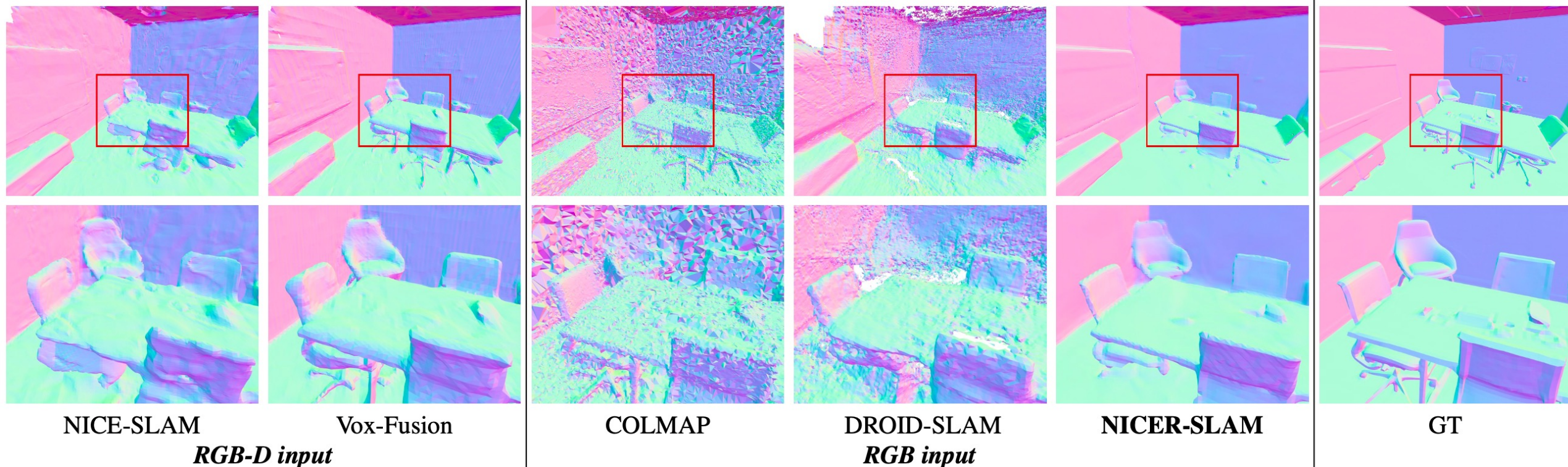
³Lund University

⁴State Key Lab of CAD&CG, Zhejiang University

⁵University of Amsterdam

⁶University of Tübingen, Tübingen AI Center

⁷Microsoft



<https://arxiv.org/abs/2302.03594>

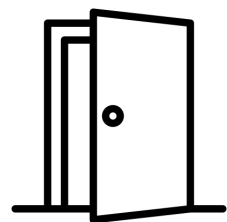
How does NeRF advance 3D Scene Reconstruction?

How does CLIP advance 3D Scene Understanding?

How does NeRF advance 3D Scene Reconstruction?

How does CLIP advance 3D Scene Understanding?

Google **ETH** zürich



OpenScene

3D Scene Understanding with Open Vocabularies

CVPR 2023

Songyou Peng



Kyle Genova



Chiyu "Max" Jiang



Andrea Tagliasacchi



Marc Pollefeys



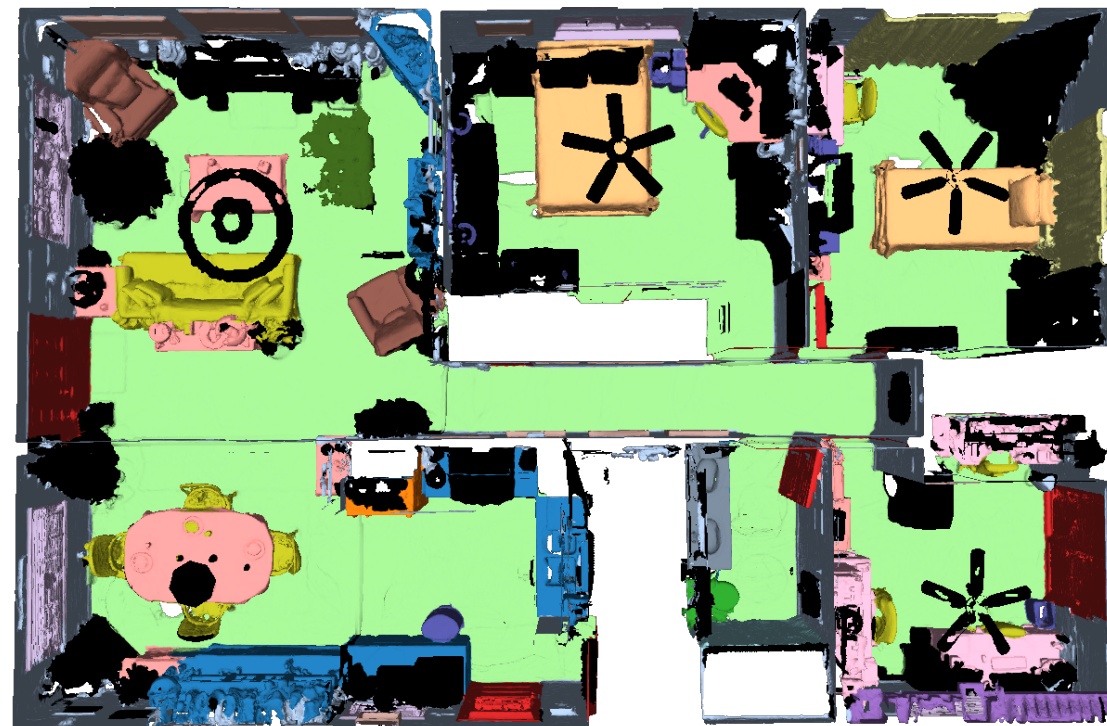
Tom Funkhouser





Input 3D Geometry

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



Traditional Semantic Segmentation

Only train and test on a few common classes

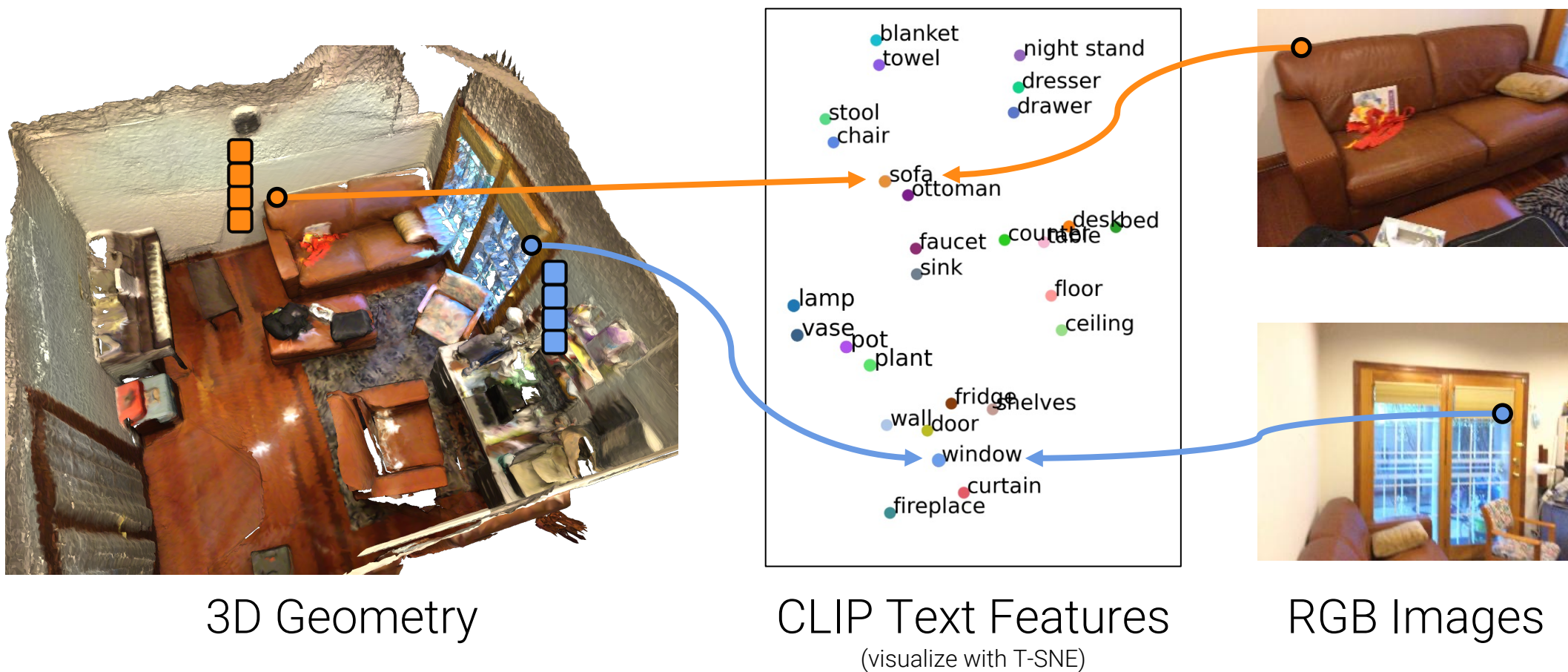


Input 3D Geometry

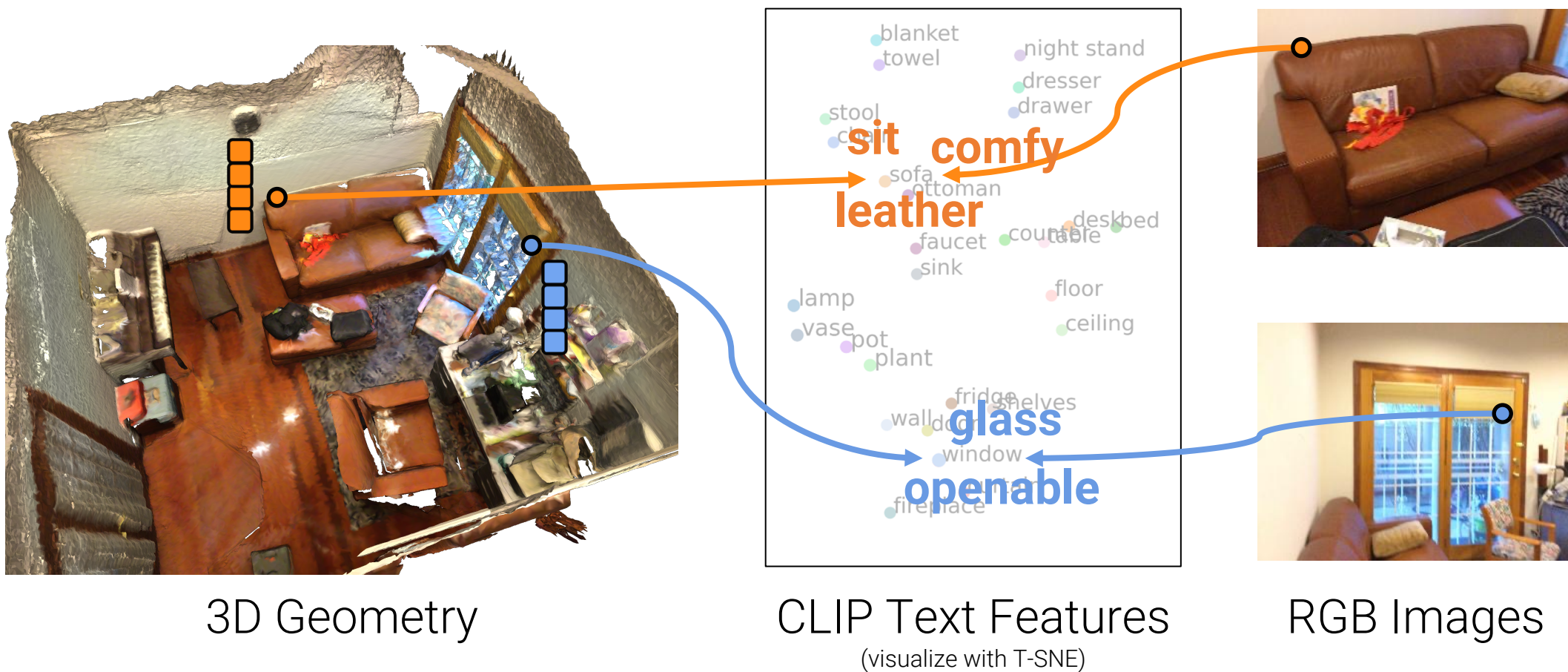
- Affordance prediction
- Material identification
- Physical property estimation
- Rare object retrieval
- Activity site prediction
- Fine-grained semantic segmentation
- Many more...

3D Scene Understanding Tasks w/o Labels

Key Idea: Co-embed 3D features with CLIP features



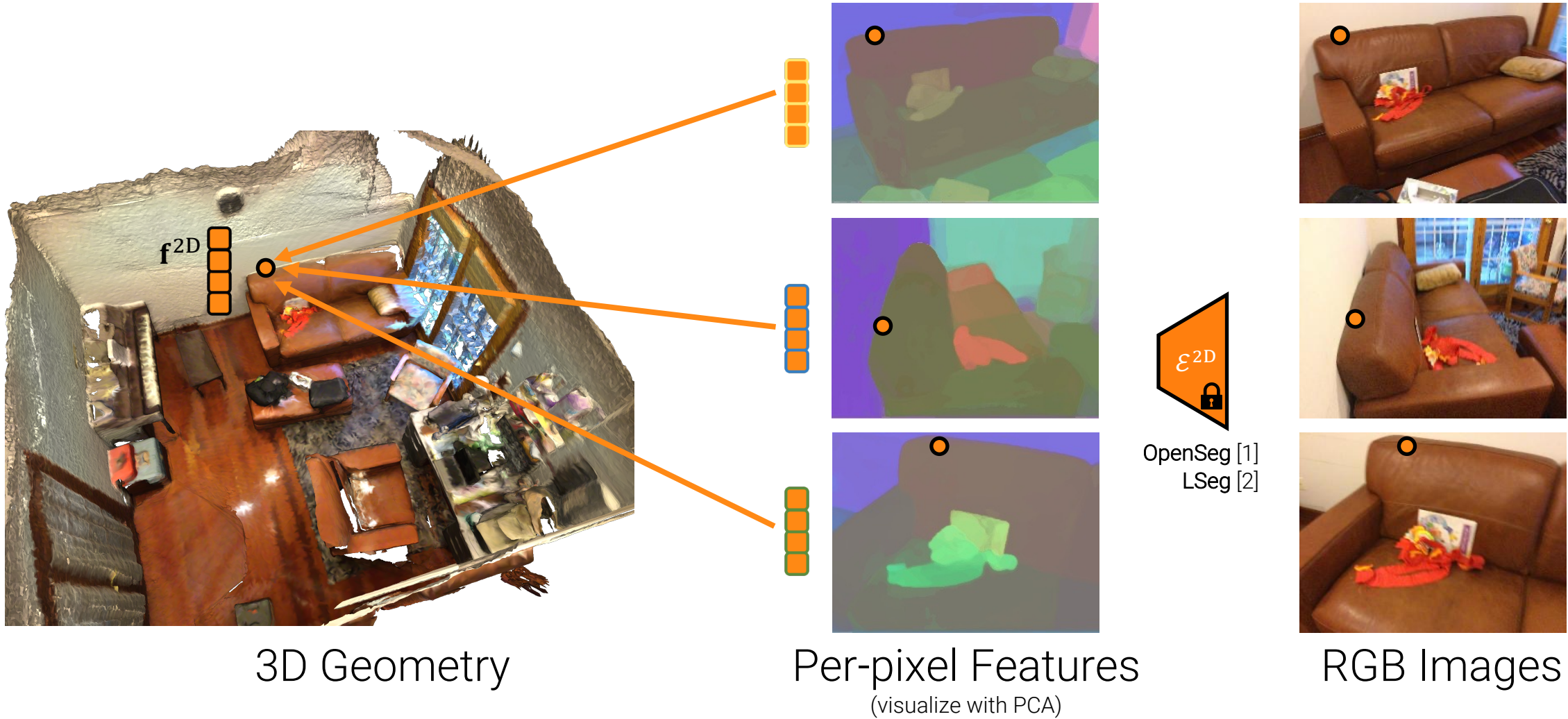
Key Idea: Co-embed 3D features with CLIP features



Note: bold word embeddings are approximate

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

Step 1: Multi-view Feature Fusion



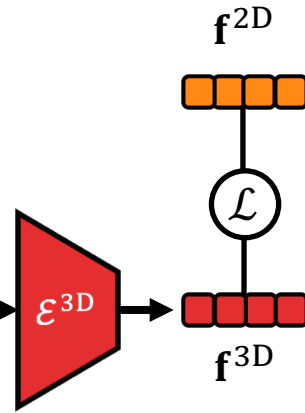
[1] Ghiasi, Gu, Cui, Lin: [Scaling Open-Vocabulary Image Segmentation with Image-Level Labels](#). ECCV 2022

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

Step 2: 3D Distillation

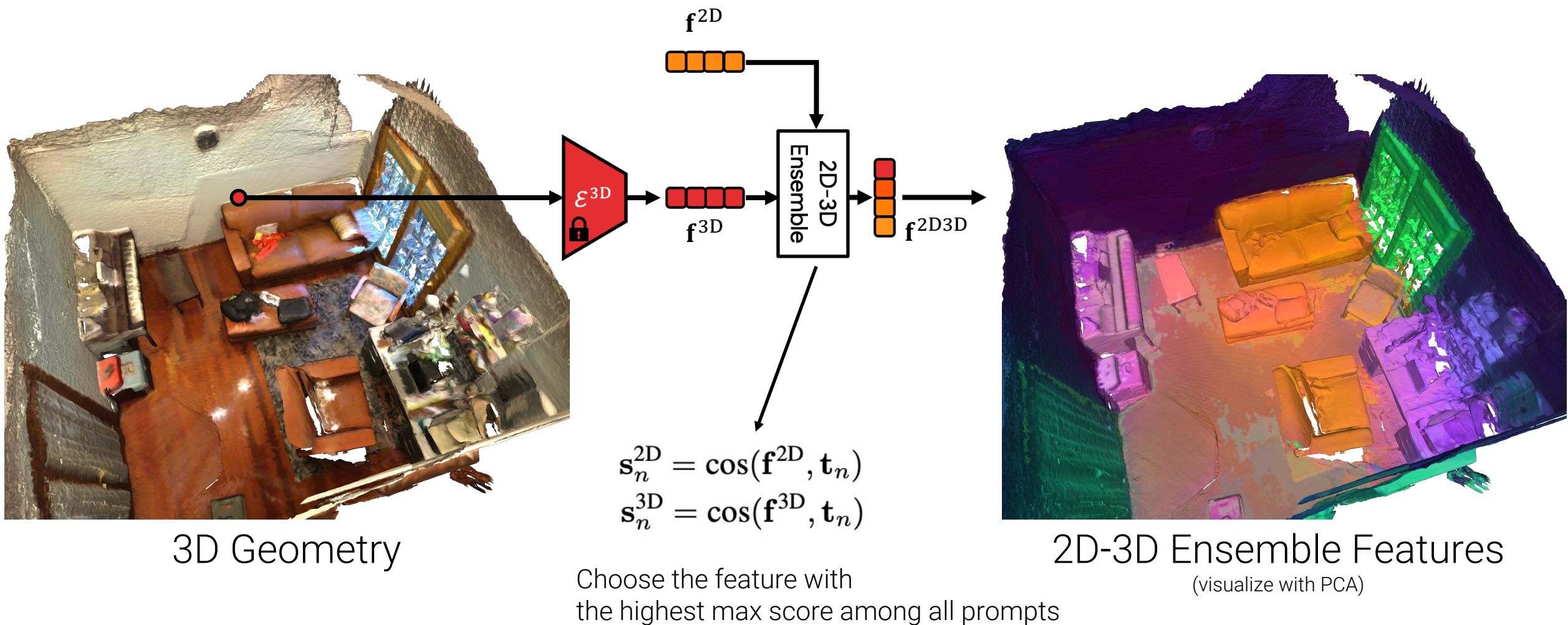


3D Geometry

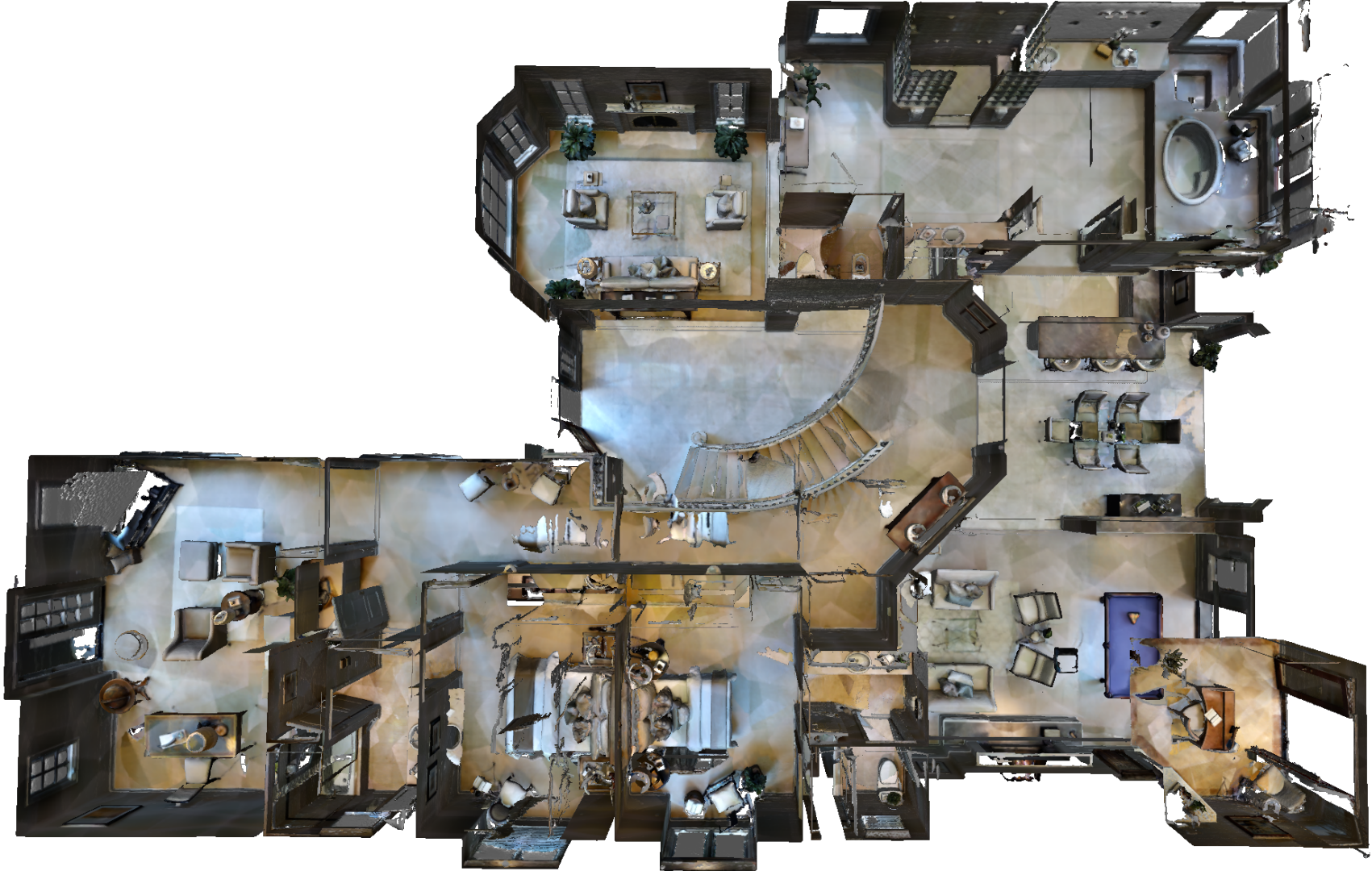


$$\mathcal{L} = 1 - \cos(\mathbf{f}^{2D} - \mathbf{f}^{3D})$$

Step 3: 2D-3D Ensemble



Open-Vocabulary, Zero-shot 3D Semantic Segmentation



Input 3D Geometry



Our Zero-shot 3D Segmentation
(20 classes)

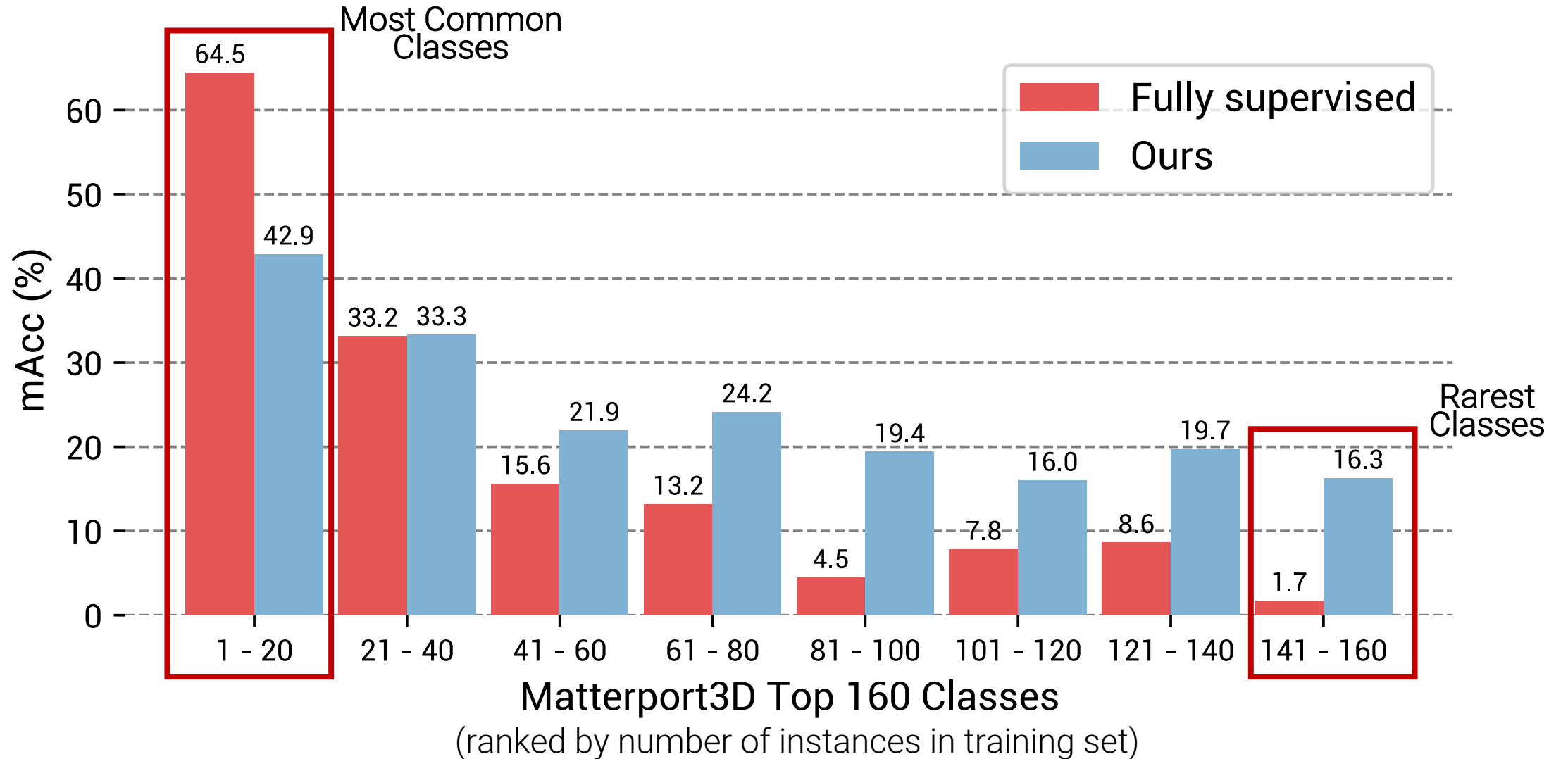
■ wall
 ■ floor
 ■ cabinet
 ■ bed
 ■ chair
 ■ sofa
 ■ table
 ■ door
 ■ window
 ■ bookshelf
 ■ picture
 ■ counter
 ■ desk
 ■ curtain
 ■ refrigerator
 ■ shower curtain
 ■ toilet
 ■ sink
 ■ bathtub
 ■ other



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	decorative plate	lamp shade	foot rest
ceiling	table	toilet	box	air vent	ottoman	container	washing machine	shoe	fire extinguisher	radiator	garage door	light	jacket	car	soap dish
floor	plant	column	coffee table	faucet	bottle	light switch	shower curtain	heater	curtain rod	kitchen island	piano	scale	bottle of soap	toilet brush	cleaner
picture	mirror	banister	counter	photo	refridgerator	purse	bin	headboard	printer	paper towel	board	bag	water cooler	drum	computer
window	towel	stairs	bench	toilet paper	bookshelf	door way	chest	telephone	sheet	bucket	rope	display case	tea pot	whiteboard	knob
chair	sink	stool	garbage bin	fan	wardrobe	basket	microwave	candle	blanket	glass	ball	toilet paper holder	tray	range hood	paper
pillow	shelves	vase	fireplace	railing	pipe	chandelier	blinds	flower pot	handle	dishwasher	exercise equipment	stuffed animal		candelabra	projector

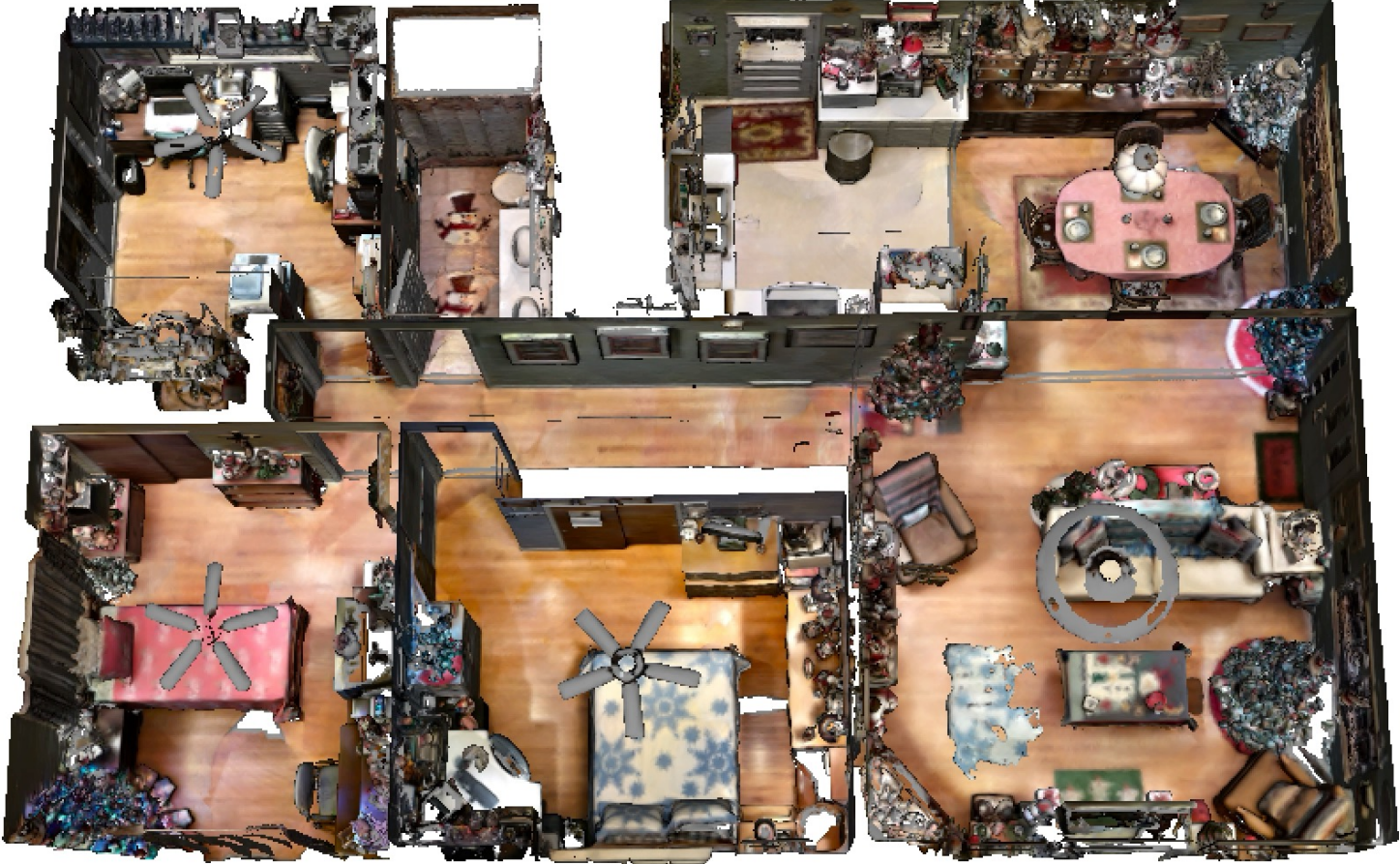
Comparison



Interactive Demo

Open-vocabulary 3D Scene Exploration

Text queries:



Take-home Message

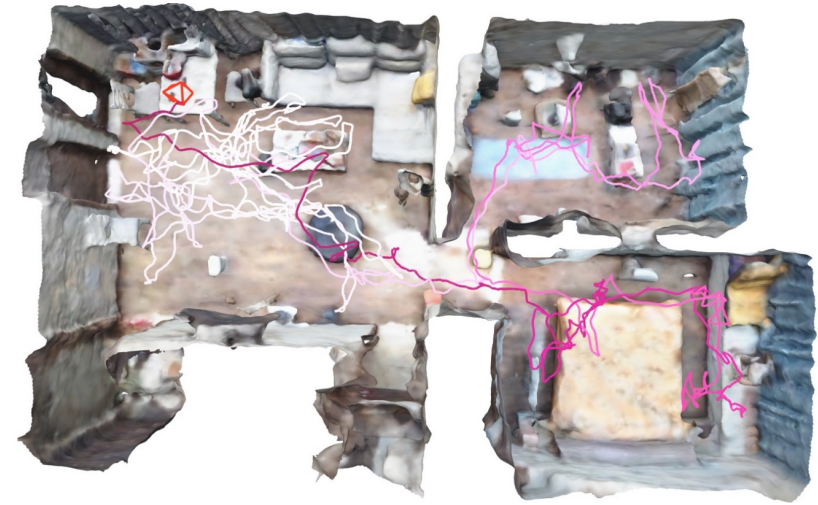
- We enable a wide range of applications by open-vocabulary queries
- This can hopefully influence how people train 3D scene understanding systems in the future
- The project can be improved in many aspects
 - Better feature fusion strategy than simple averaging
 - Combine CLIP features with NeRF/SLAM
 - [concept-fusion.github.io](https://github.com/ConceptFusion/concept-fusion)

How does NeRF advance 3D Scene Reconstruction?



[NeurIPS'22] **MonoSDF**

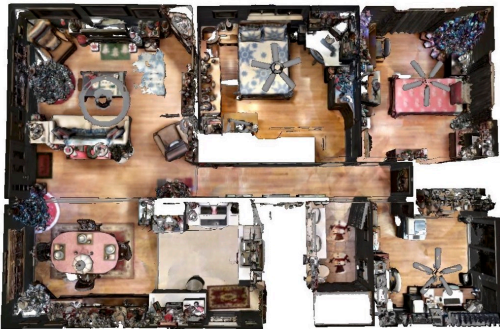
github.com/autonomousvision/monosdf



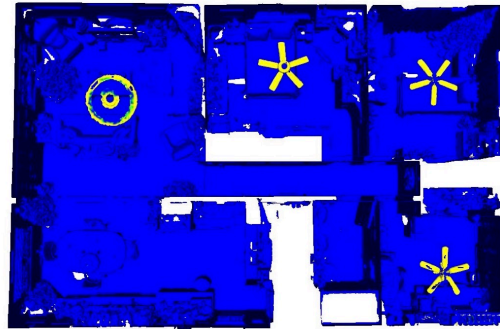
[CVPR'22] **NICE-SLAM**

pengsongyou.github.io/nice-slam

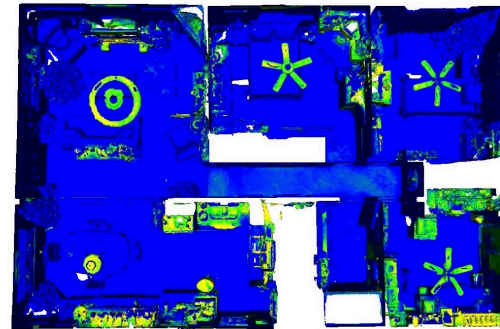
How does CLIP advance 3D Scene Understanding?



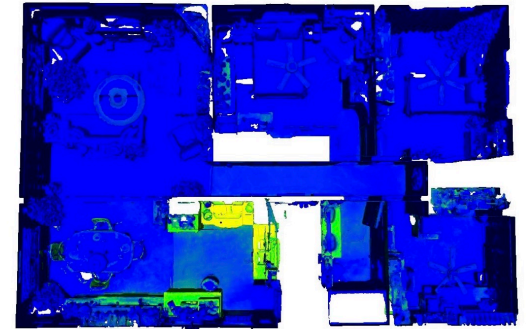
Input 3D Geometry



"fan" - Object



"metal" - Material



"kitchen" - Room Type

[CVPR'23] **OpenScene**

pengsongyou.github.io/openscene