

# Large-Scale 3D Scene Reconstruction with NeRF

Songyou Peng

ETH Zurich and Max Planck Institute for Intelligent Systems



MAX PLANCK INSTITUTE  
FOR INTELLIGENT SYSTEMS



Stanford Computational Imaging Lab

Oct 26, 2022

# Who Am I?

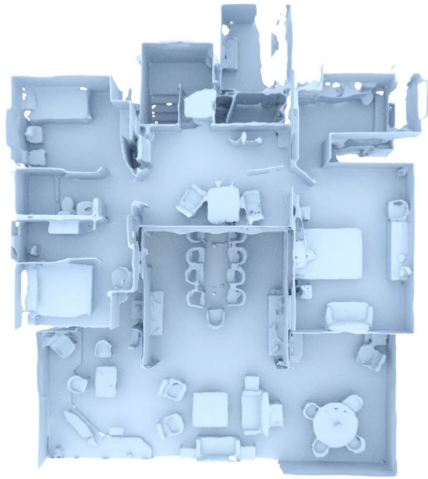
- PhD Student since 2019.09
  - Marc Pollefeys
  - Andreas Geiger
- Internships during PhD
  - 2021: Michael Zollhoefer
  - Ongoing: Tom Funkhouser
- Open to 1:1 chat!

**ETH** zürich

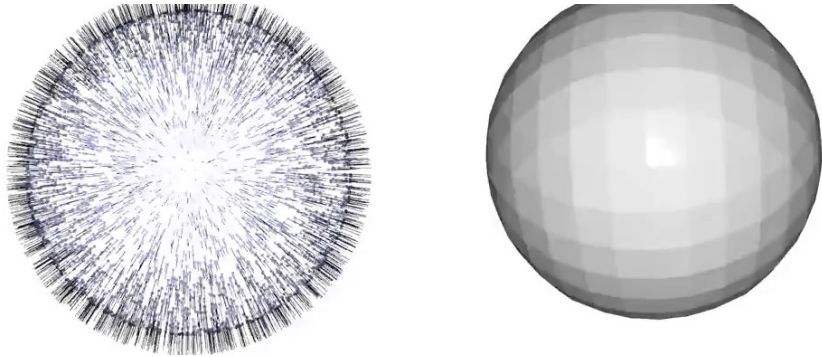


[pengsongyou.github.io](https://pengsongyou.github.io)

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



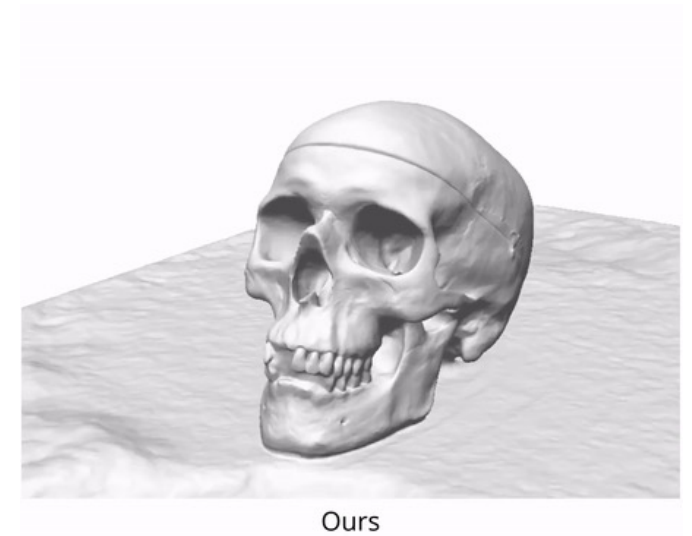
**Convolutional Occupancy Networks**  
ECCV 2020 (Spotlight)



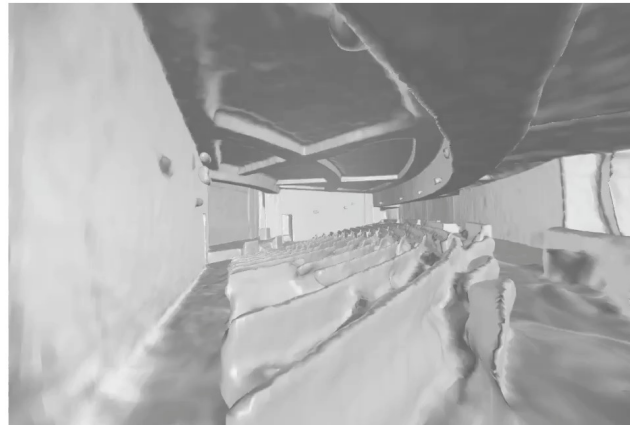
**Shape As Points**  
NeurIPS 2021 (Oral)



**KiloNeRF**  
ICCV 2021



**UNISURF**  
ICCV 2021 (Oral)



Ours  
**MonoSDF**  
NeurIPS 2022

**NICE-SLAM**  
CVPR 2022



# 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

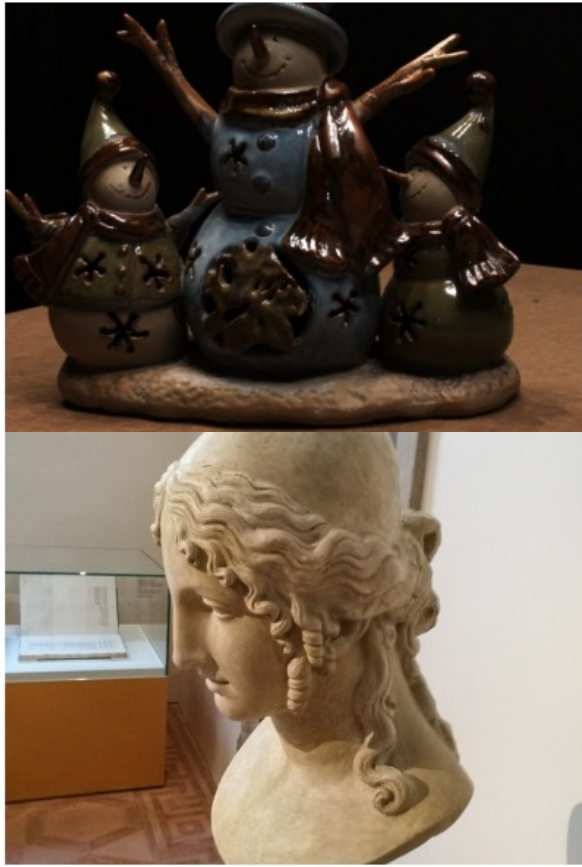


Torsten Sattler



Andreas Geiger

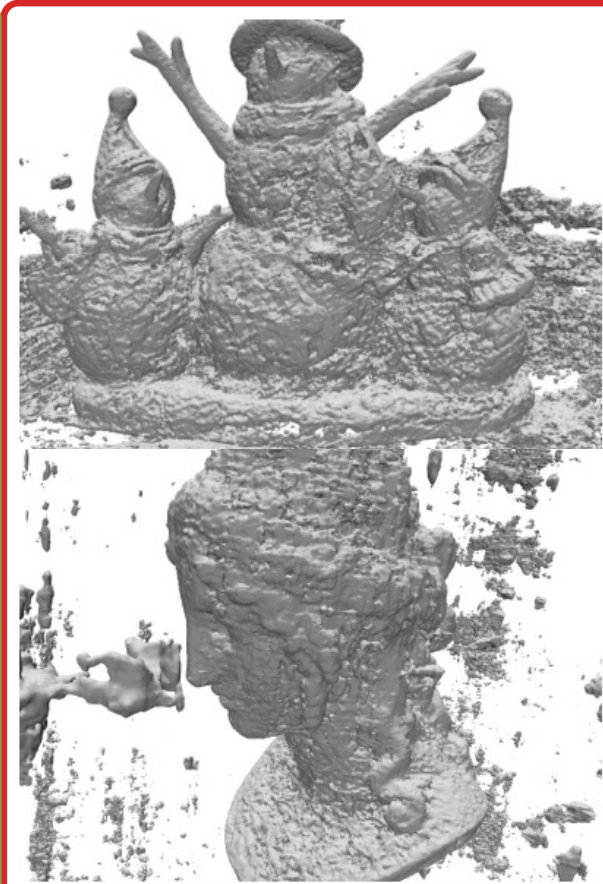
# Neural Implicit Surfaces with Volume Rendering



RGB Images



NeuS/VoISDF/UNISURF



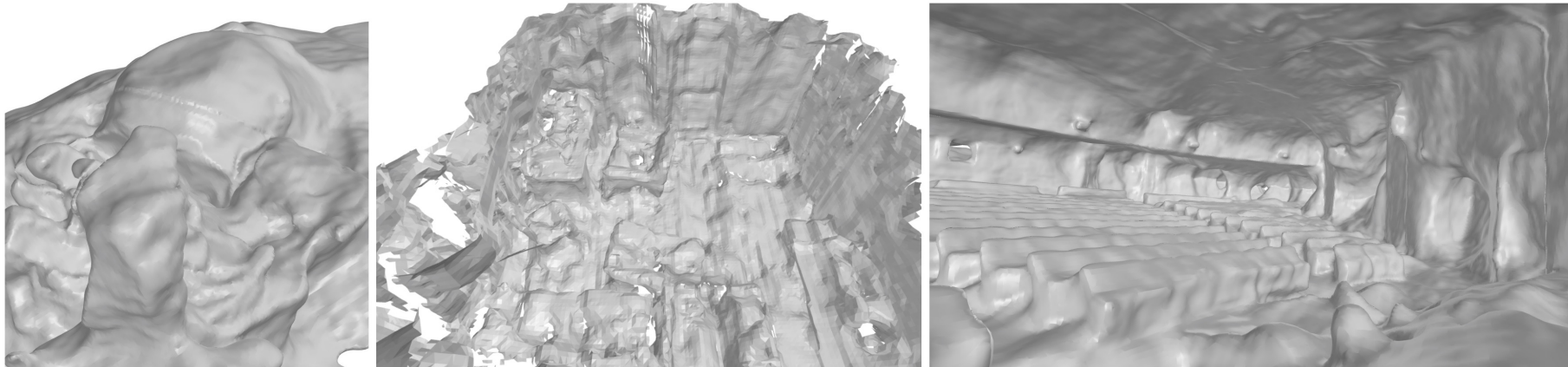
NeRF

- [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

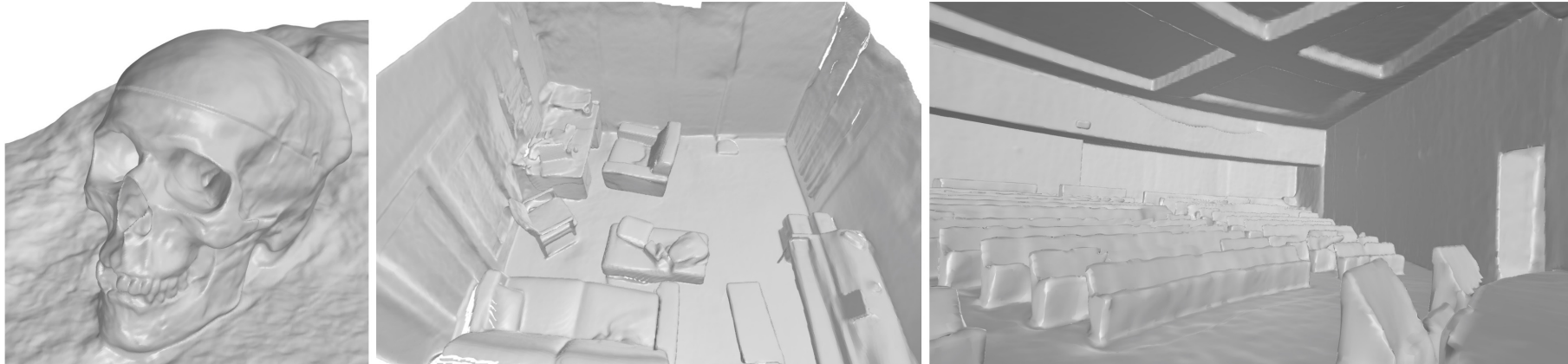


# Neural Implicit Surfaces with Volume Rendering

VoISDF



MonoSDF



DTU (3 views)

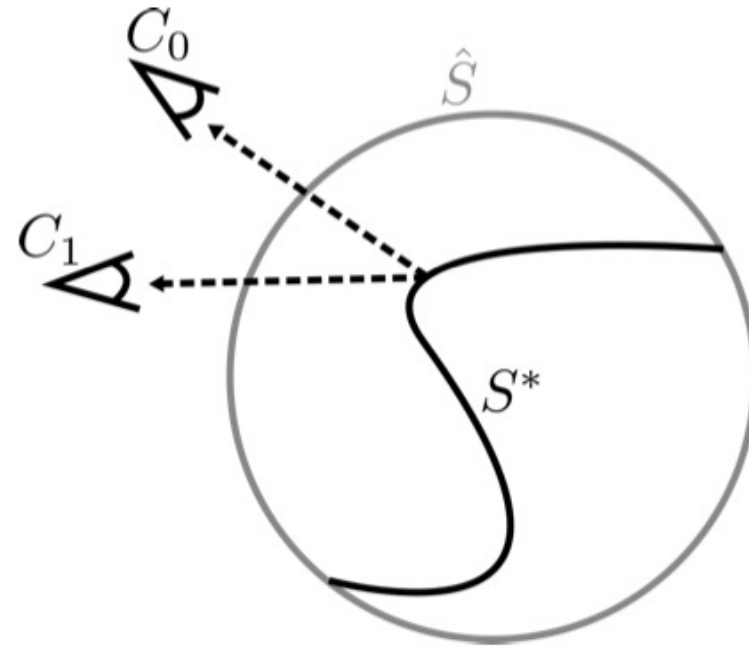
ScanNet (464 views)

Tanks & Temples (298 views)

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

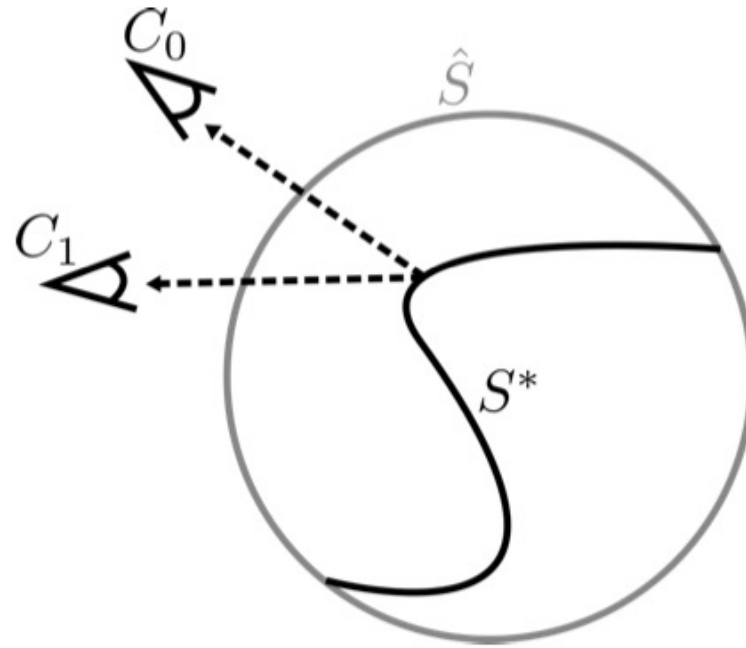


# Shape-Appearance Ambiguity



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

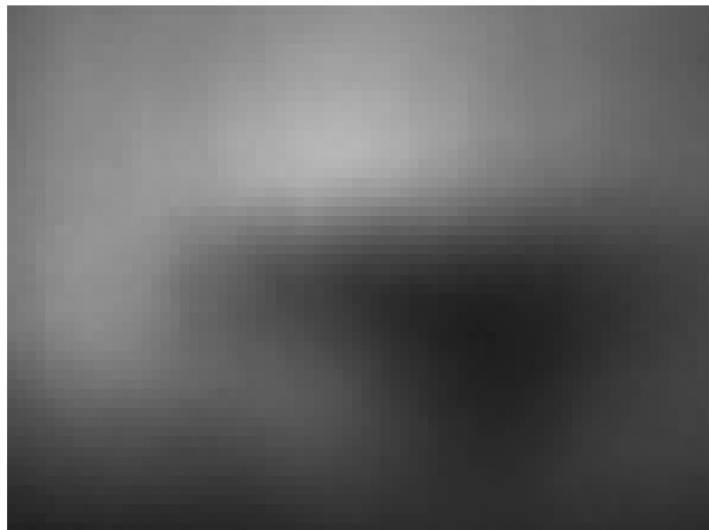
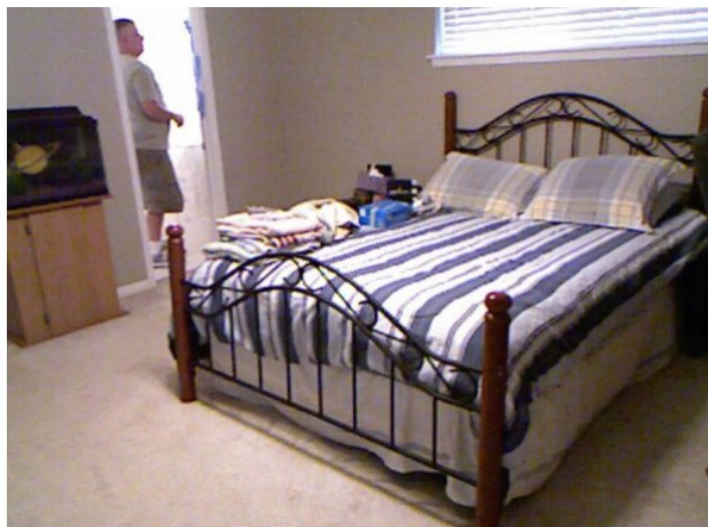
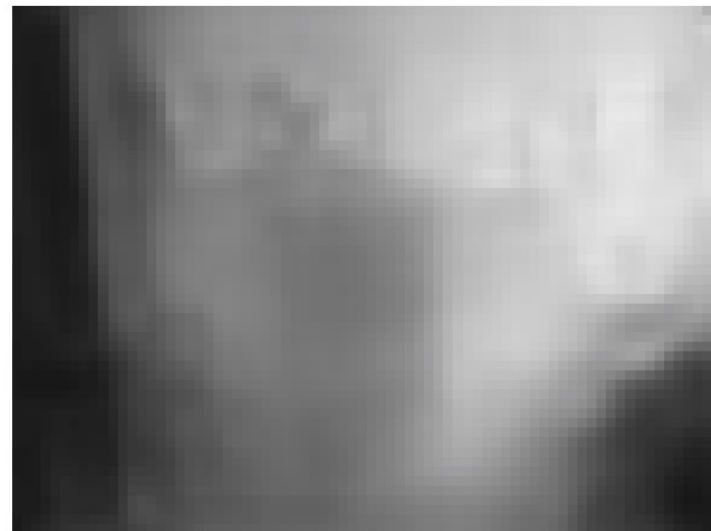
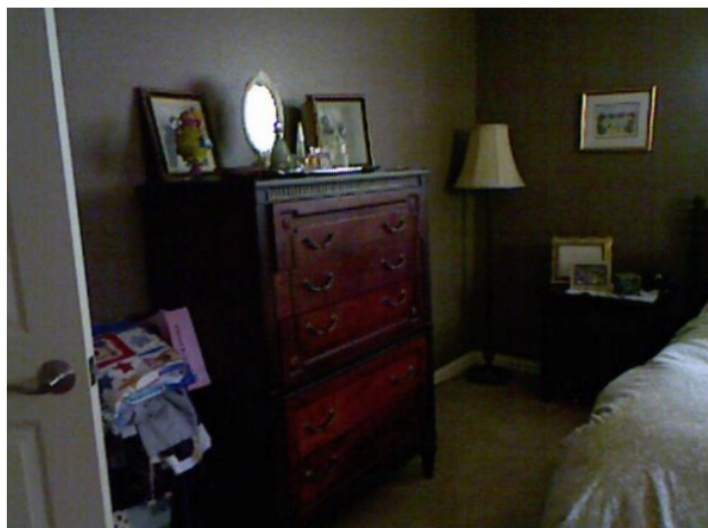
# Shape-Appearance Ambiguity



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

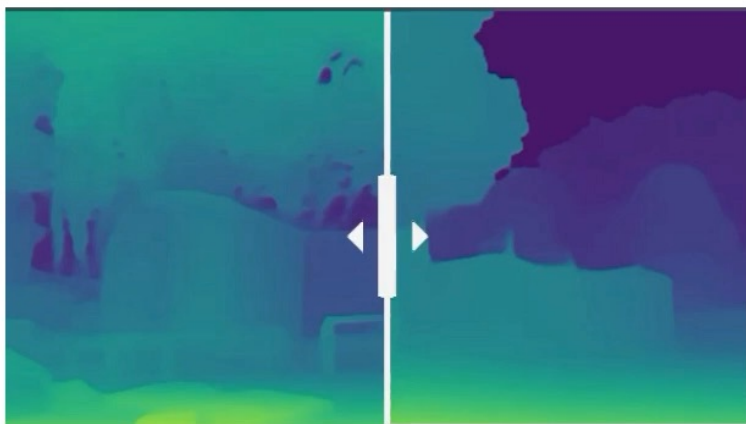
⇒ **Exploit monocular geometric priors**

# Depth Map Prediction from a Single Image



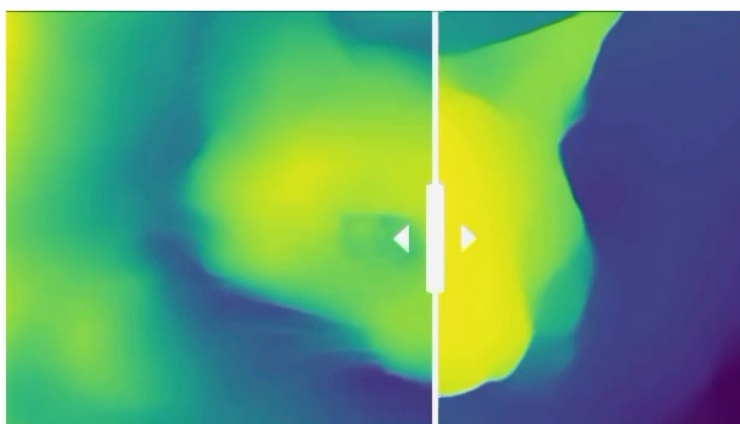


# Omnidata



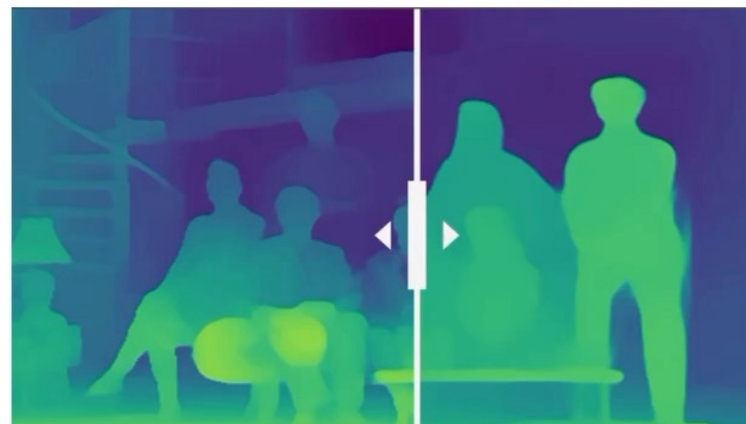
**Ours**

**MiDaS  
DPT-Hybrid**



**Ours**

**MiDaS  
DPT-Hybrid**

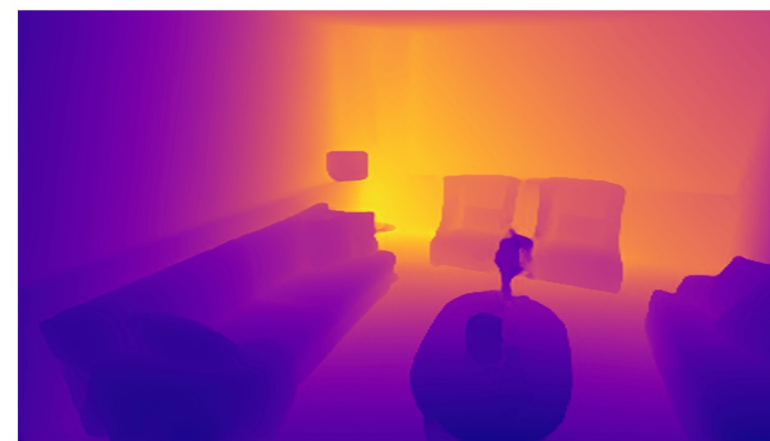
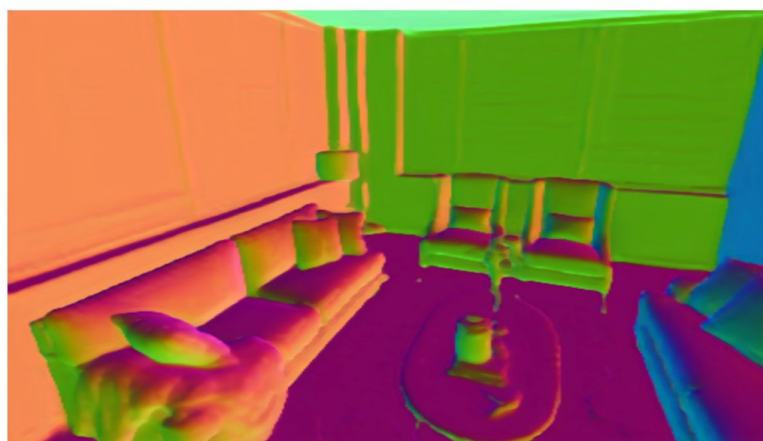
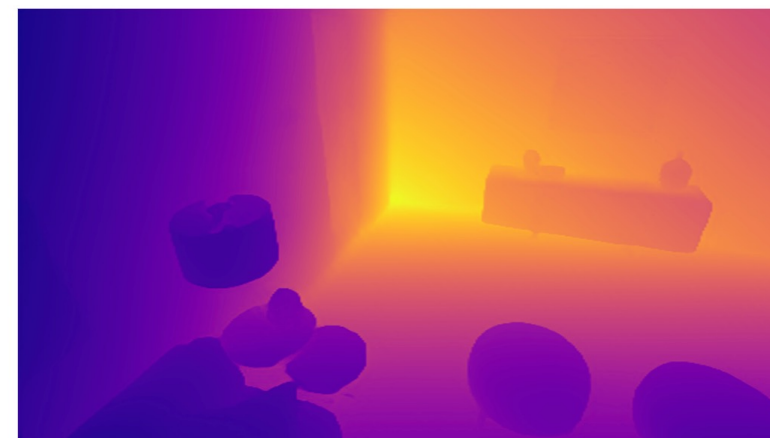
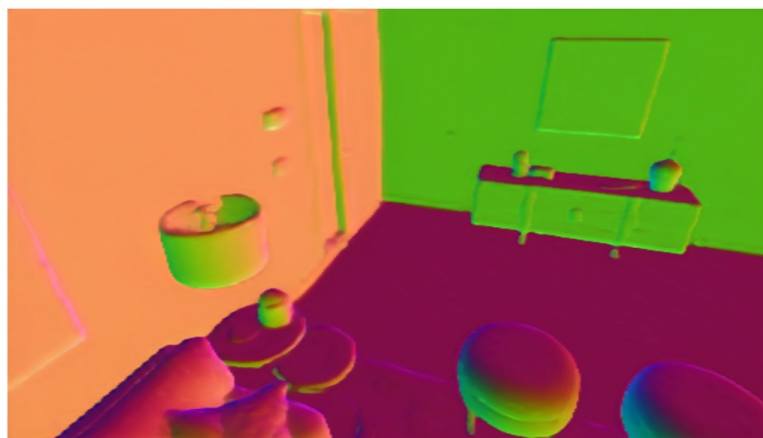


**Ours**

**MiDaS  
DPT-Hybrid**

[Ranftl et al. 2021]

# Omnidata



RGB Image

Omnidata Normal

Omnidata Depth

# MonoSDF

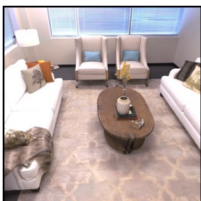




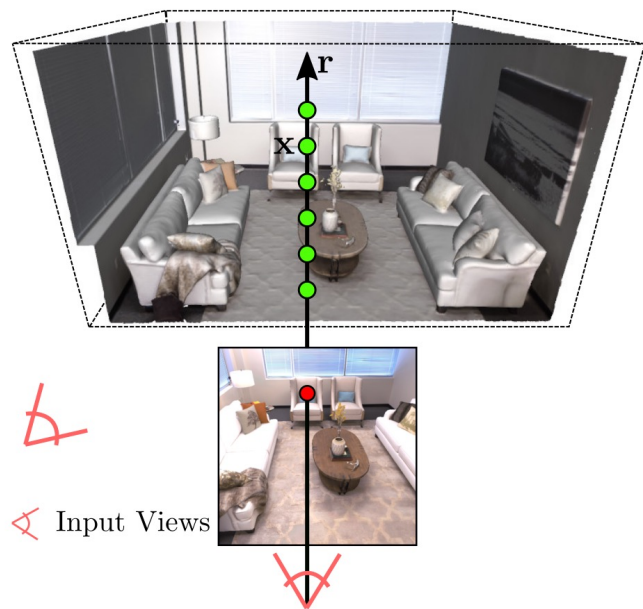
# MonoSDF



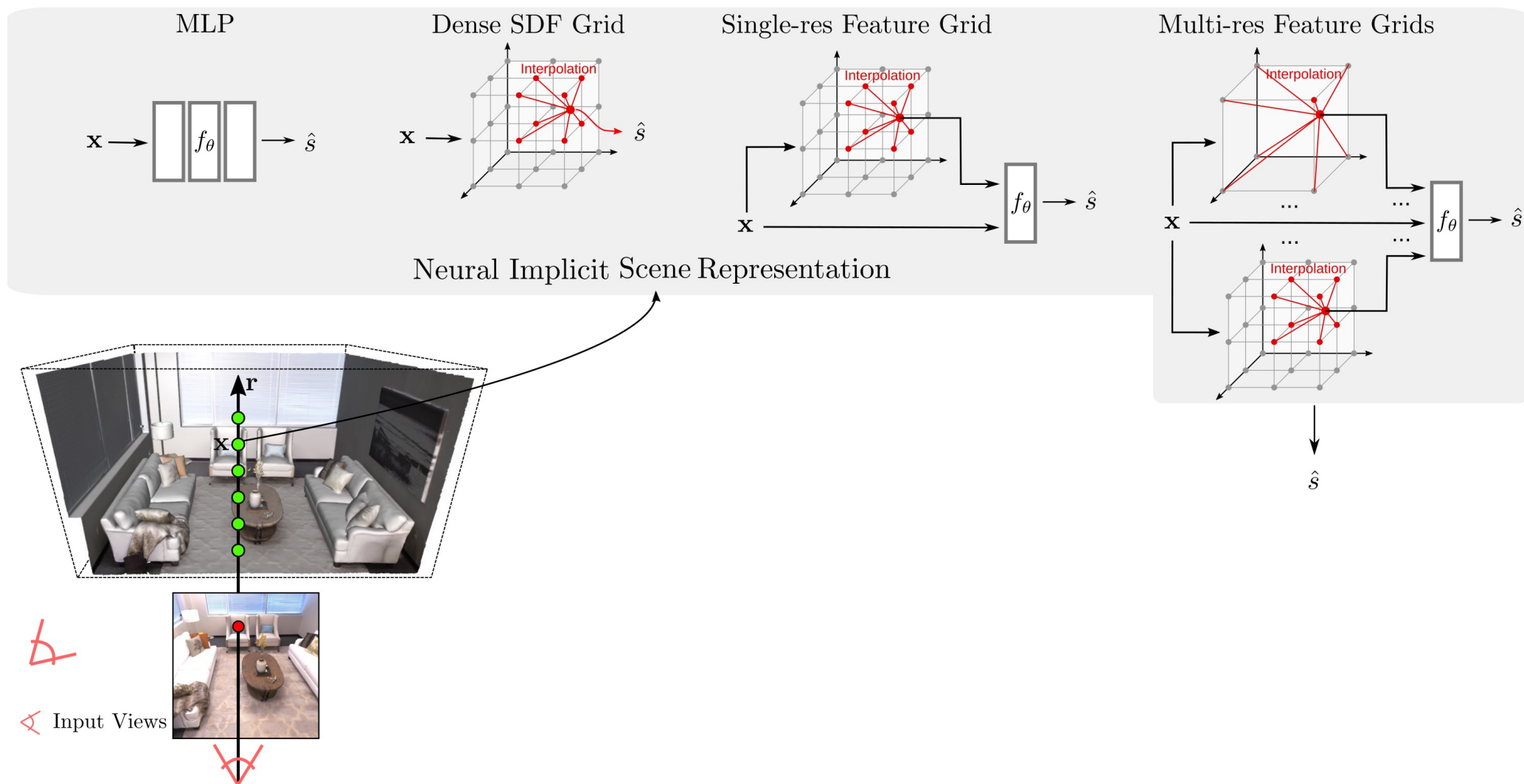
Input Views



# MonoSDF

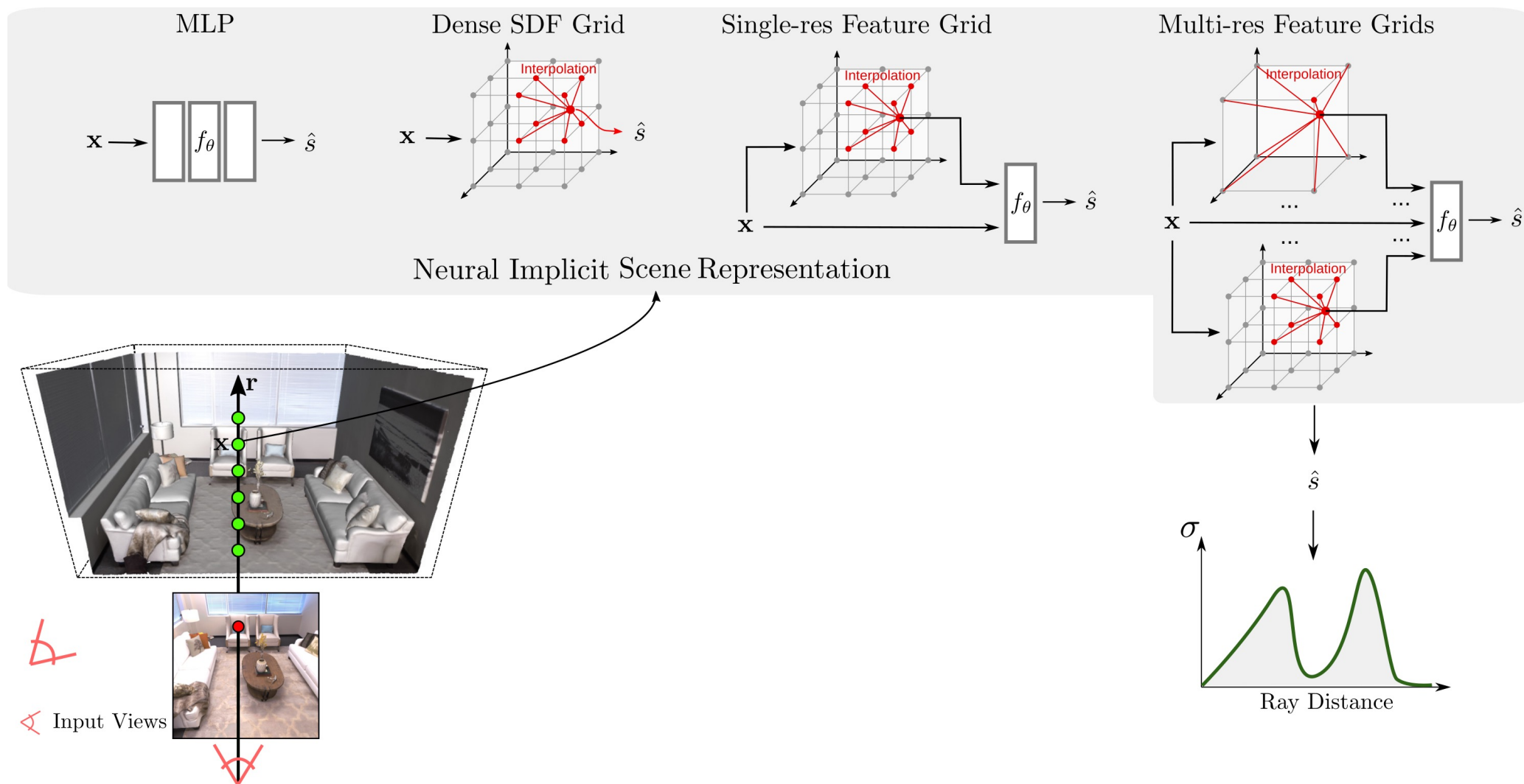


# MonoSDF

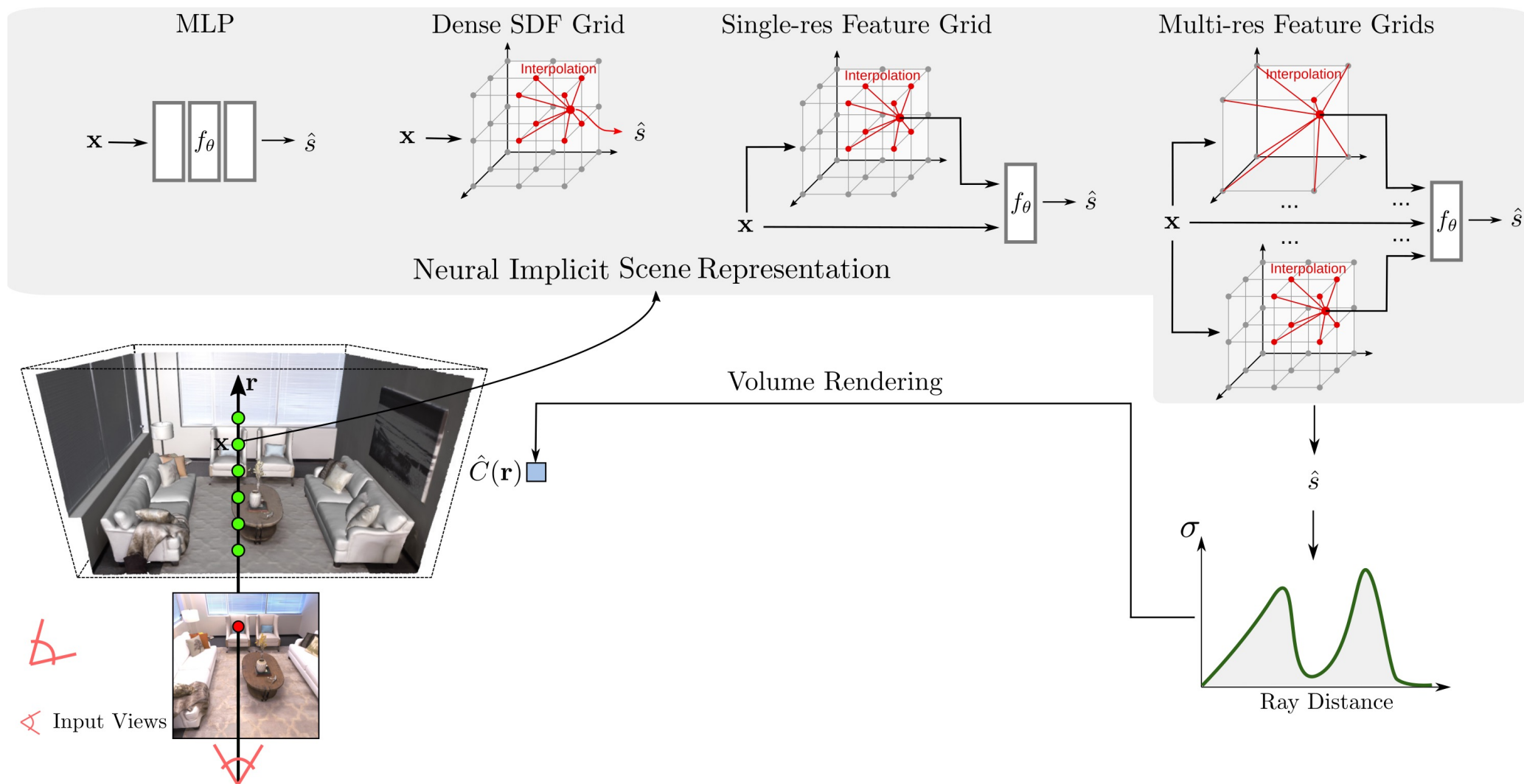




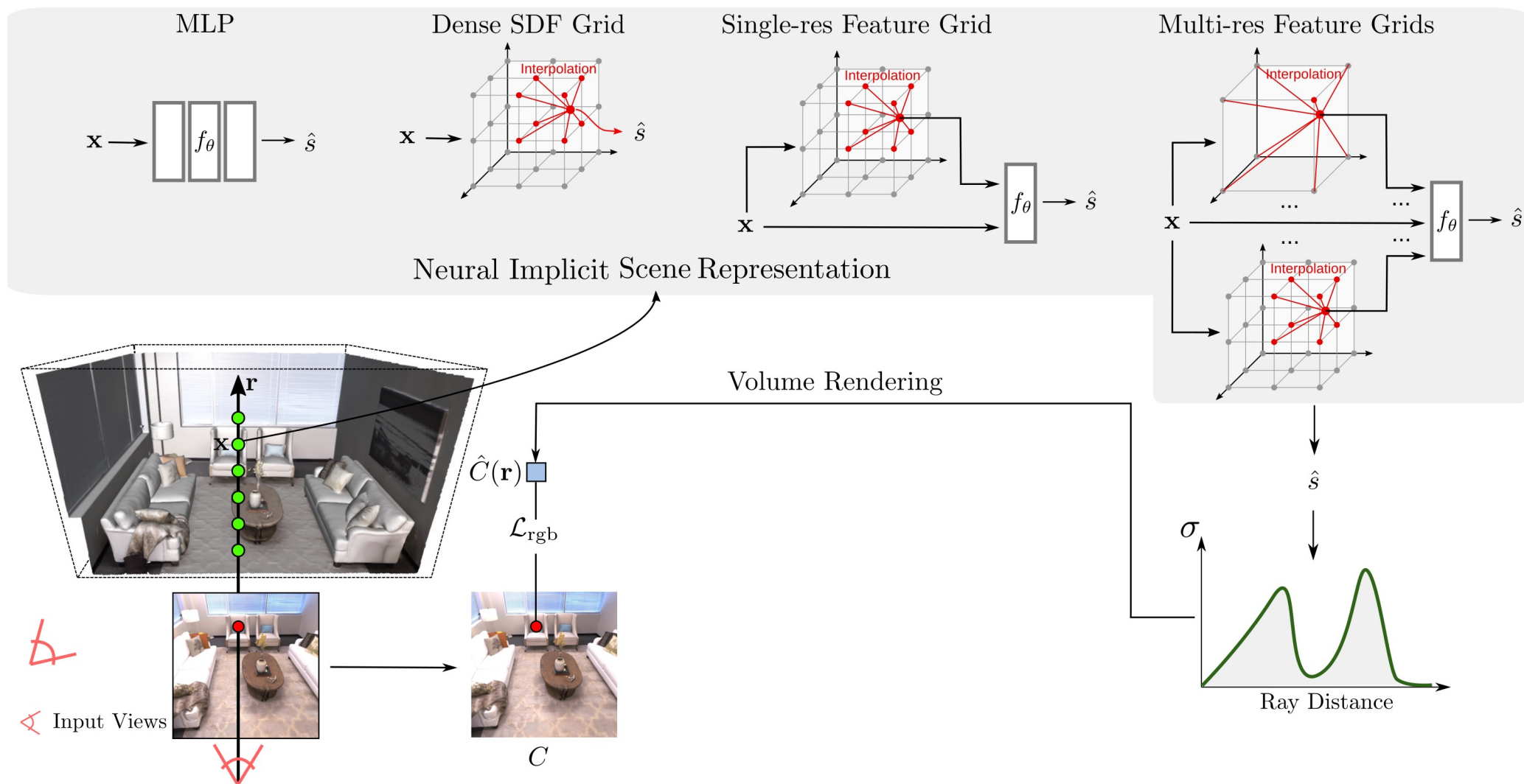
# MonoSDF



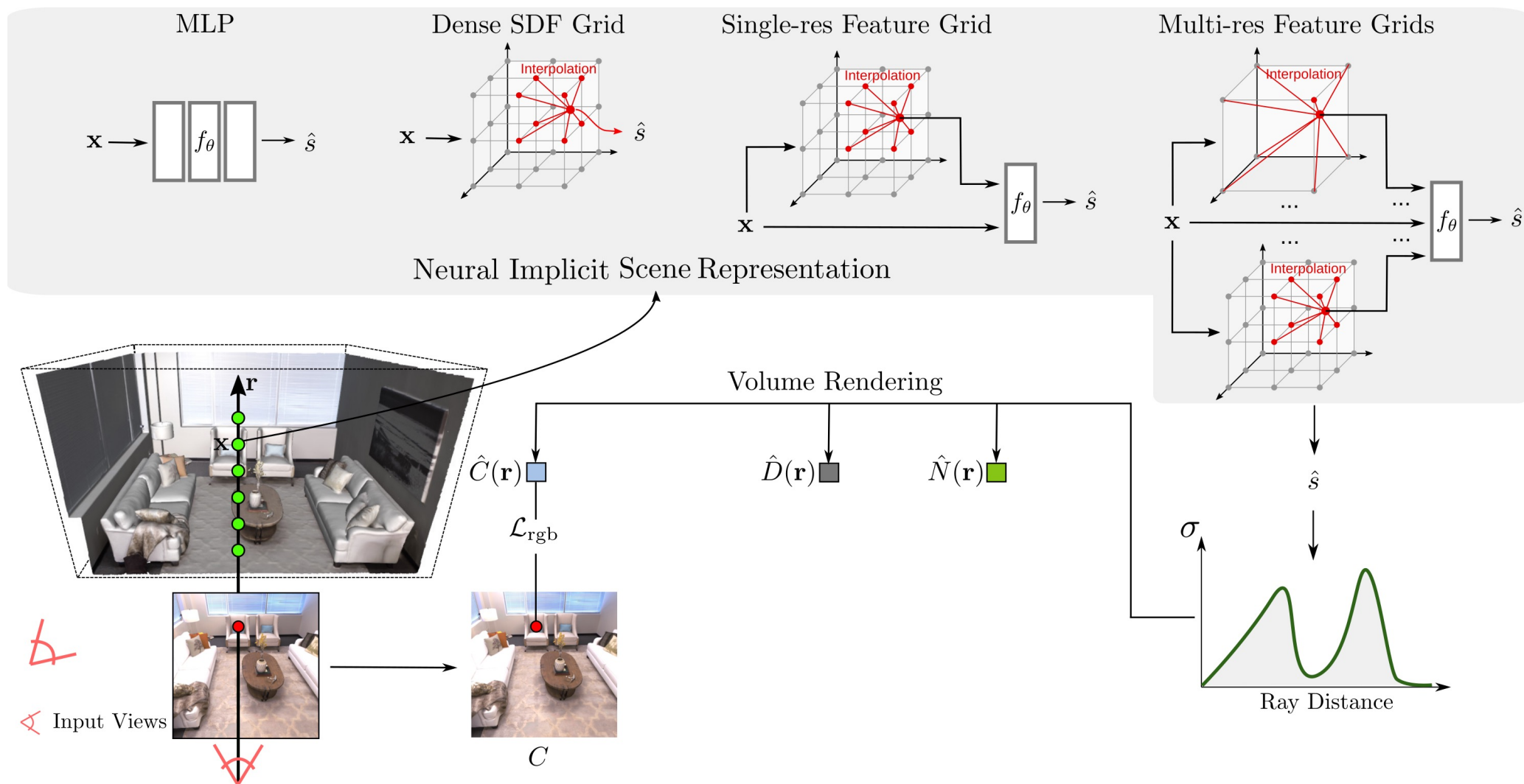
# MonoSDF



# MonoSDF

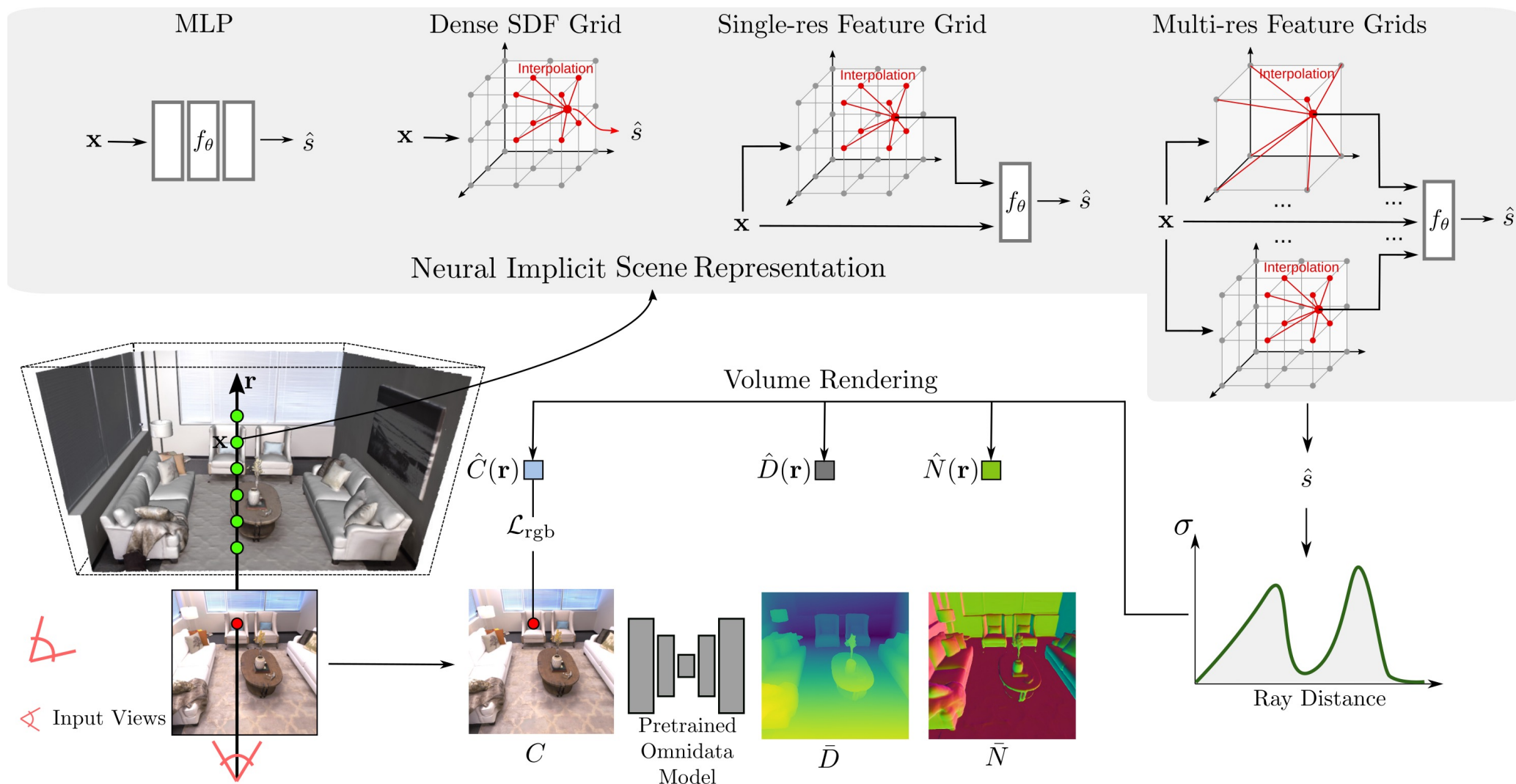


# MonoSDF

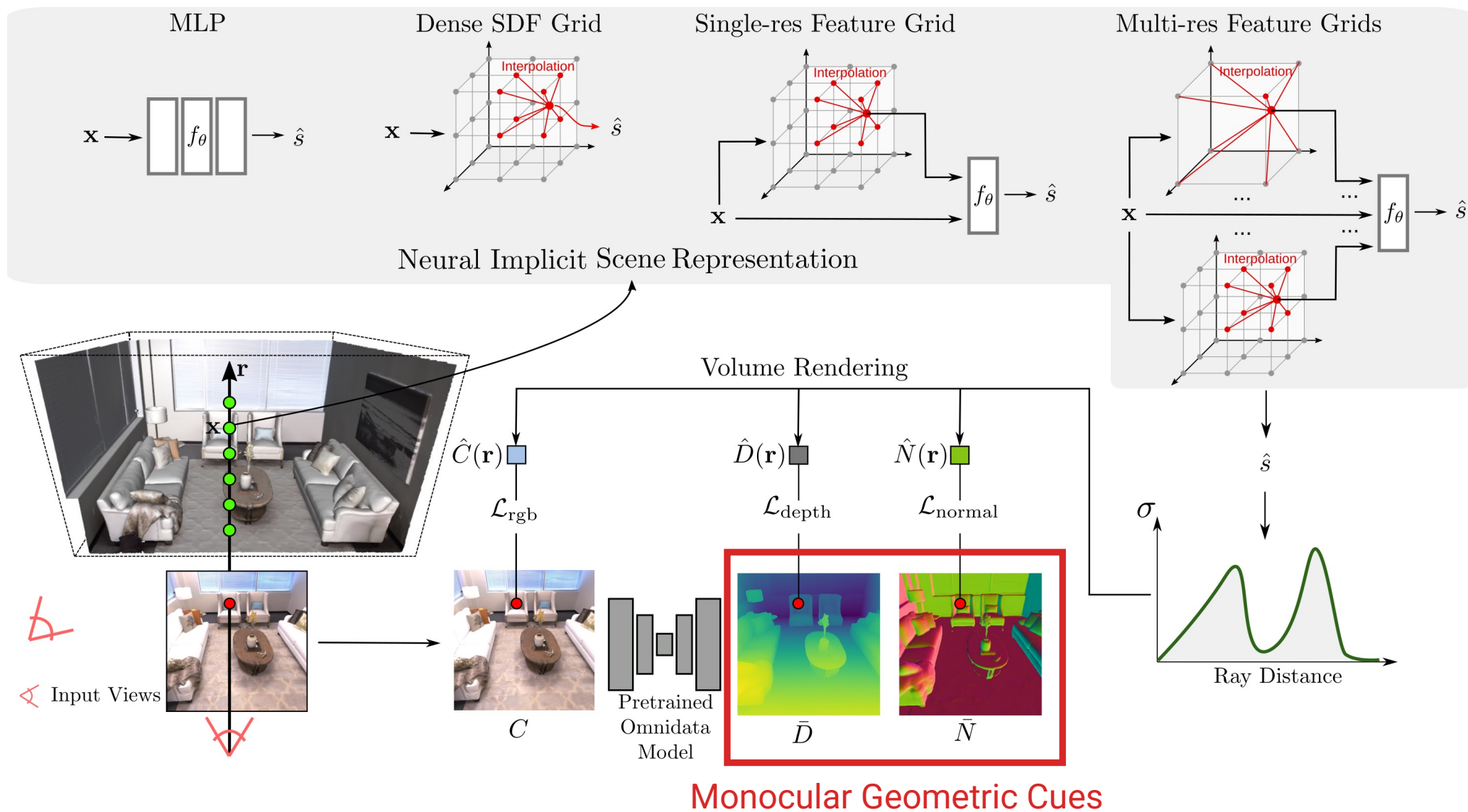




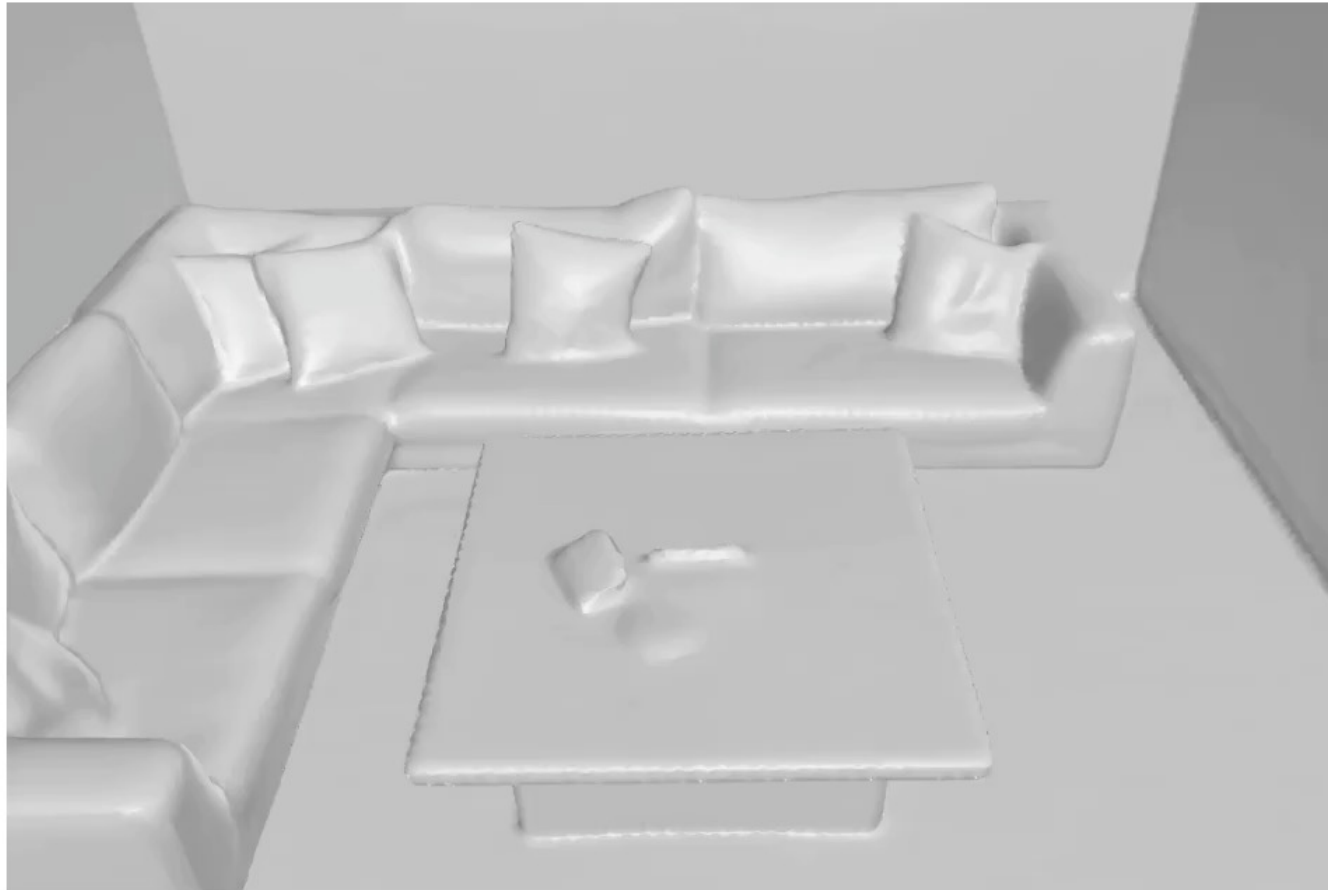
# MonoSDF



# MonoSDF



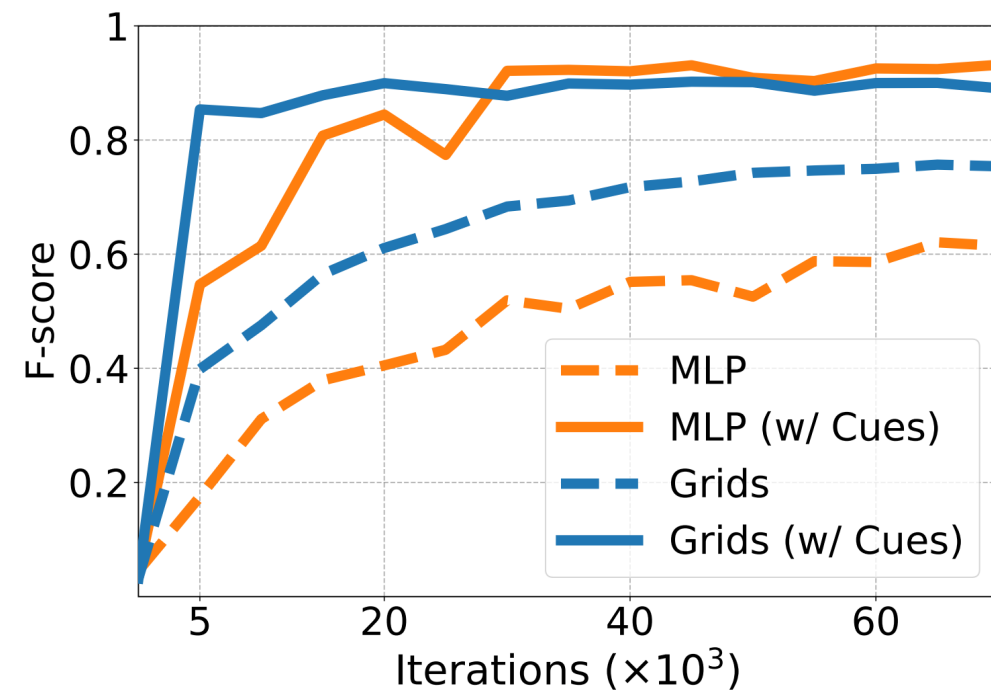
# Ablation Study



Depth & Normal Cues

# Ablation Study

		Normal C. $\uparrow$	Chamfer- $L_1 \downarrow$	F-score $\uparrow$
<b>MLP</b>	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	<b>92.11</b>	<b>2.94</b>	<b>86.18</b>
<b>Multi-Res. Grids</b>	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	<b>90.93</b>	<b>3.23</b>	<b>85.91</b>



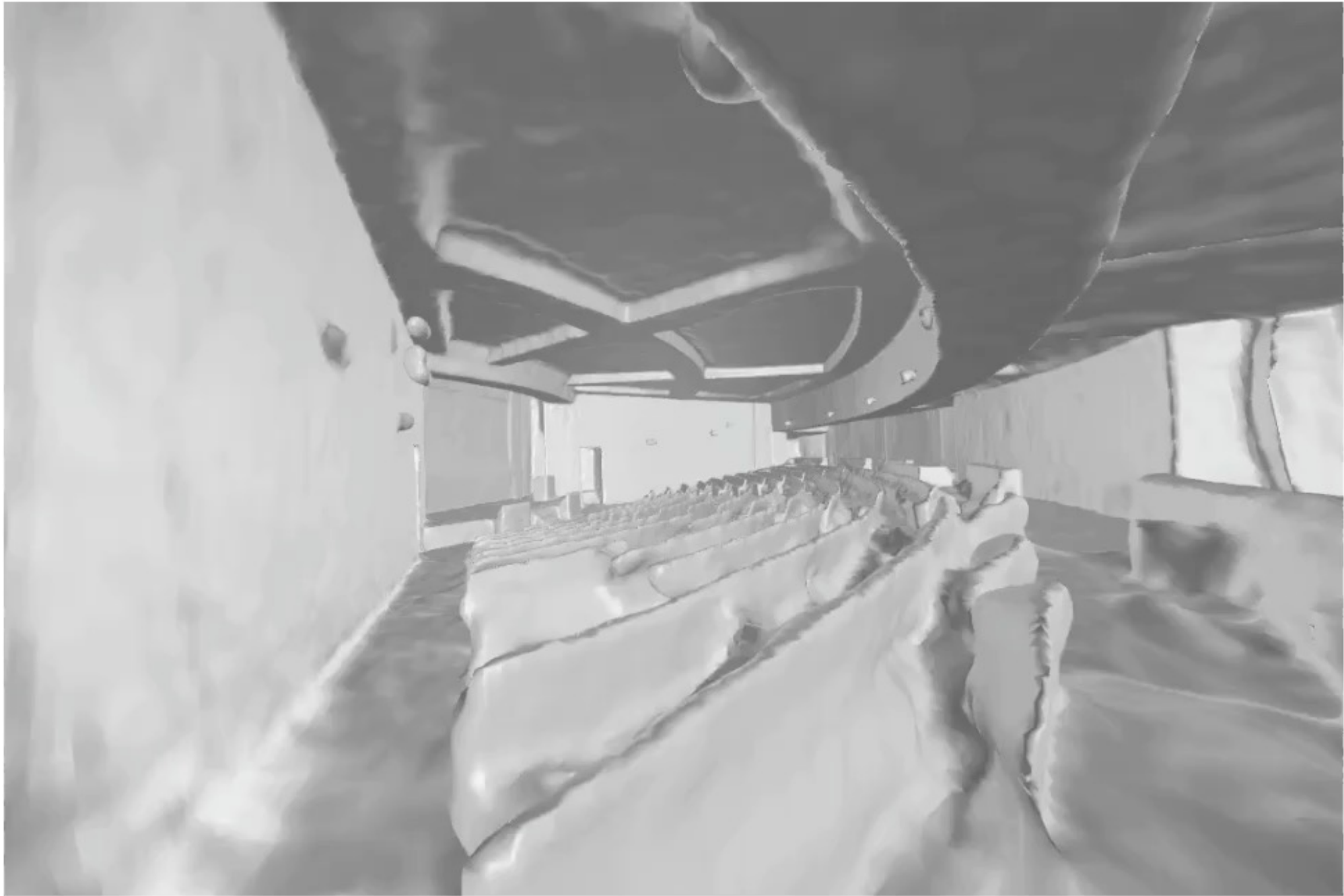
- ! Monocular cues improve reconstruction results significantly
- ! Combining **depth & normal** leads to best performance
- ! Monocular cues can improve **convergence speed**



# Baseline Comparisons on ScanNet

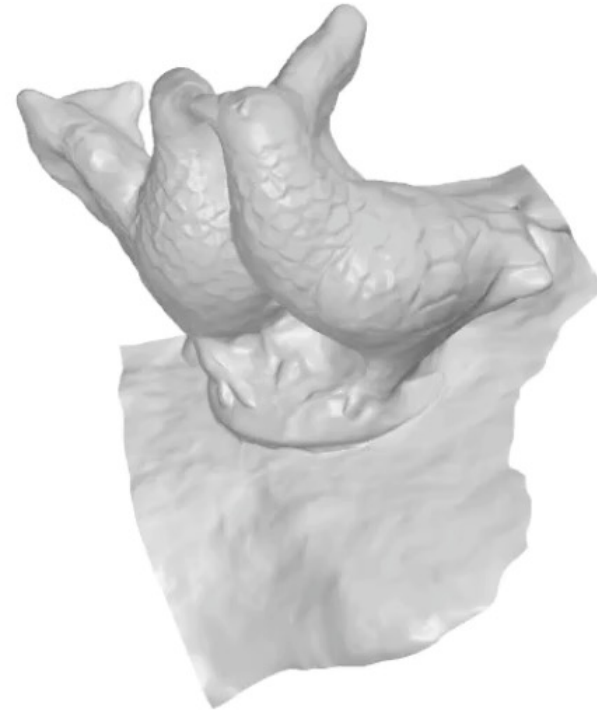
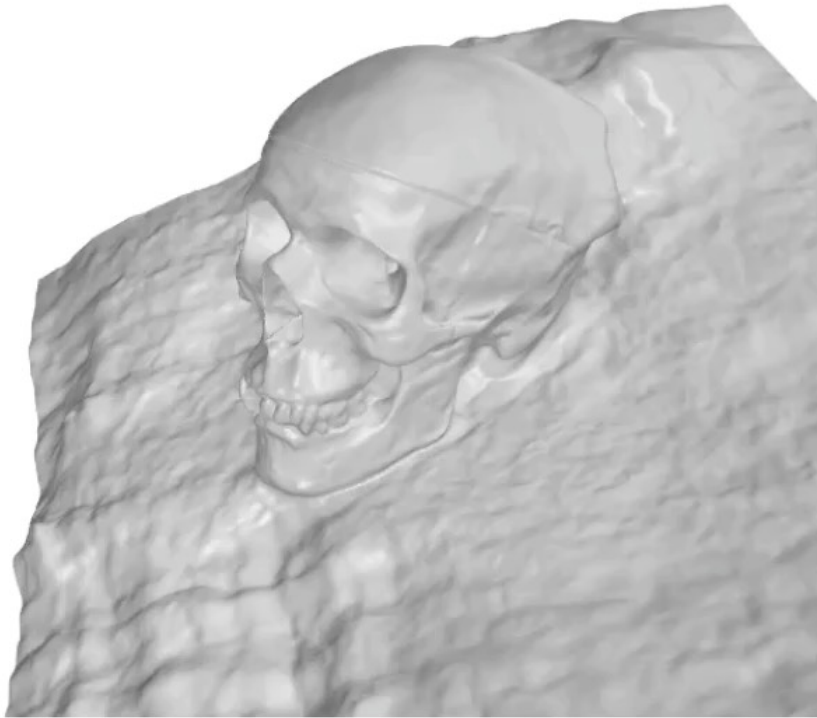


Ours



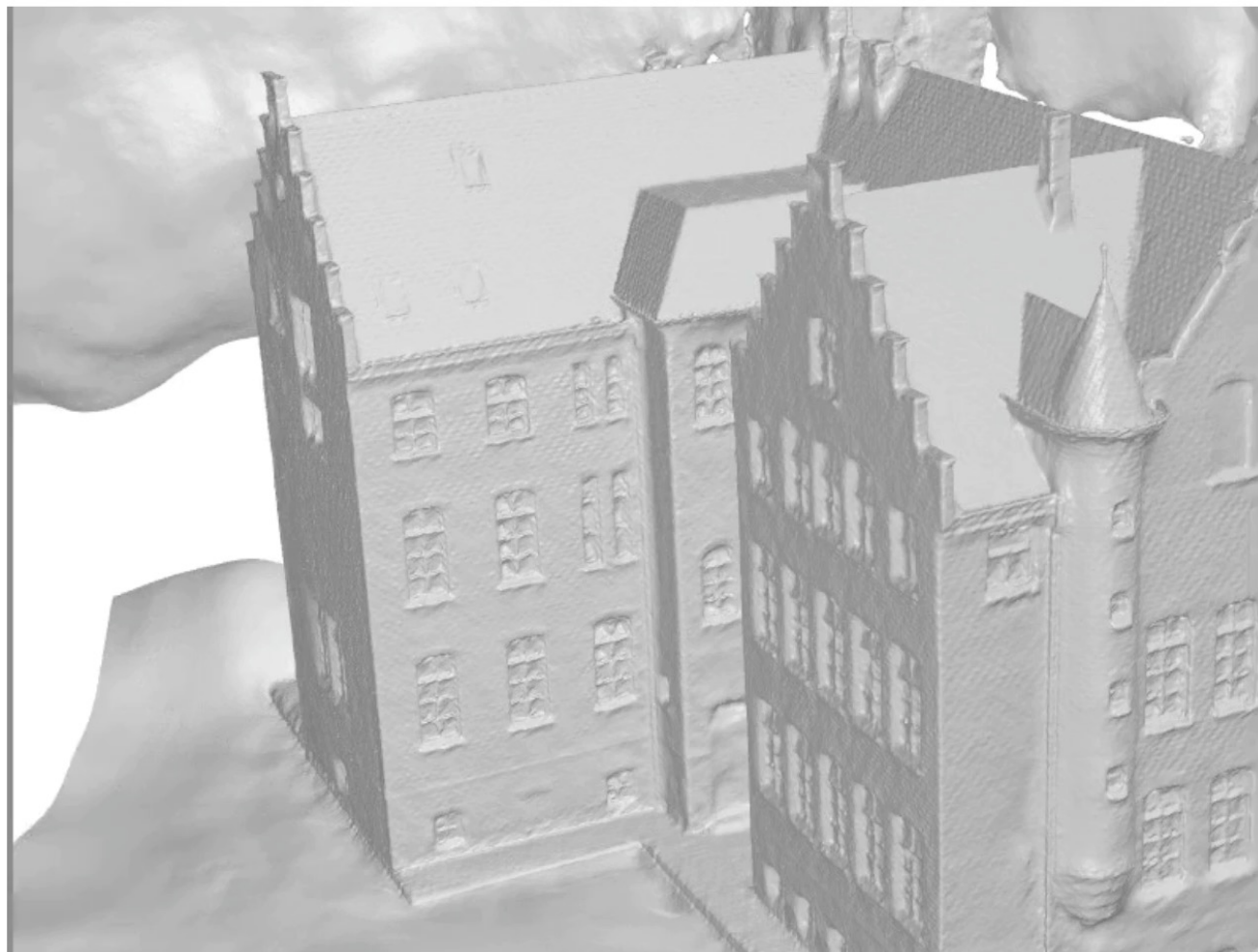
Ours

# Baseline Comparisons on DTU (3-views)



Ours

# Baseline Comparisons on DTU (all views)



Ours (Grids)

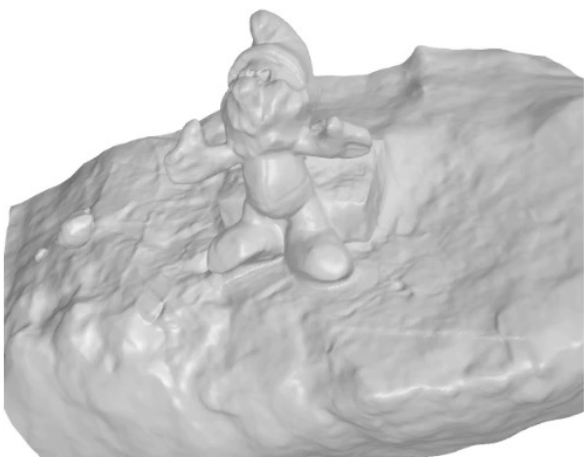


# Multi-Res. Feature Grids with High-Res. Cues

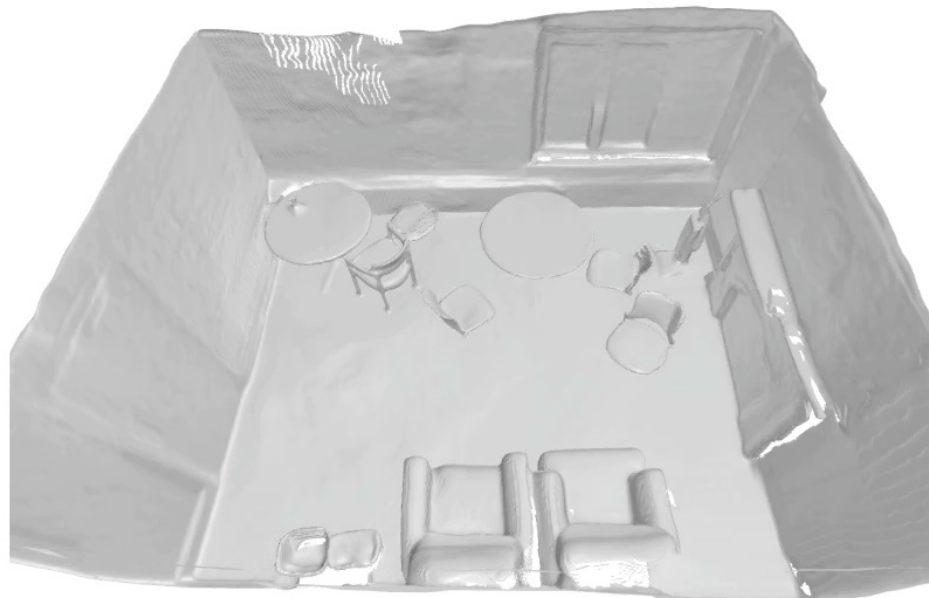


# Take-home Message

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



DTU (3 views)



ScanNet



Tanks and Temples

- ! Monocular cues improve reconstruction results and speed up optimization
- ! Analysis and investigate multiple scene representations
- ! **Limitation:** Still require camera poses given :(



# NICE-SLAM

## Neural Implicit Scalable Encoding for SLAM

CVPR 2022

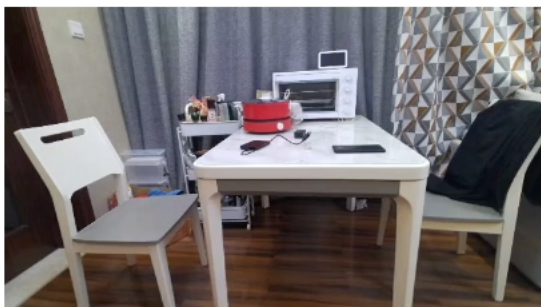
Zihan Zhu\*   Songyou Peng\*   Viktor Larsson   Weiwei Xu   Hujun Bao  
Zhaopeng Cui   Martin R. Oswald   Marc Pollefeys

\* Equal Contributions

**ETH** zürich



## RGB-D Sequences



40x Speed



# iMAP

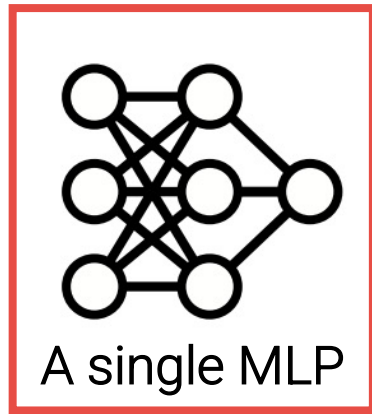
[Sucar et al., ICCV'21]



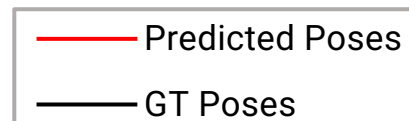
First neural implicit-based **online** SLAM system

# iMAP

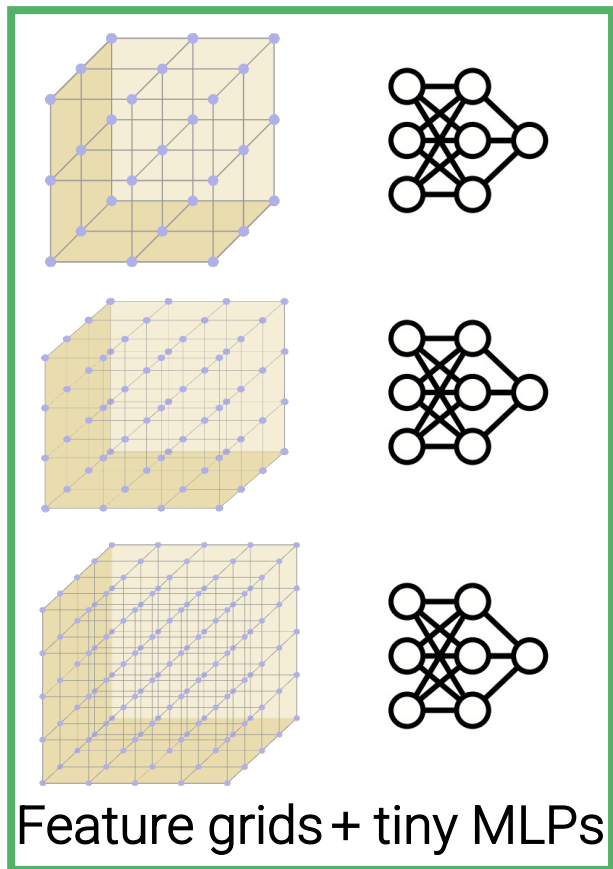
[Sucar et al., ICCV'21]



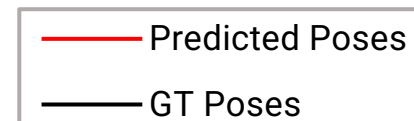
- Fail when scaling up to larger scenes
- Global update → Catastrophic forgetting
- Slow convergence



# NICE-SLAM

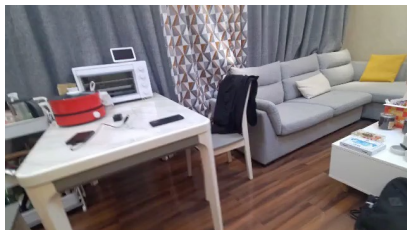
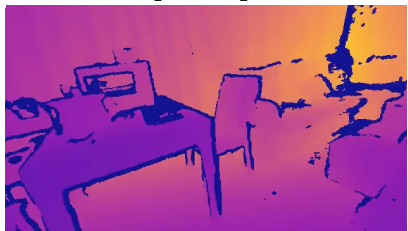


- + Applicable to large-scale scenes
- + Local update → No forgetting problem
- + Fast convergence

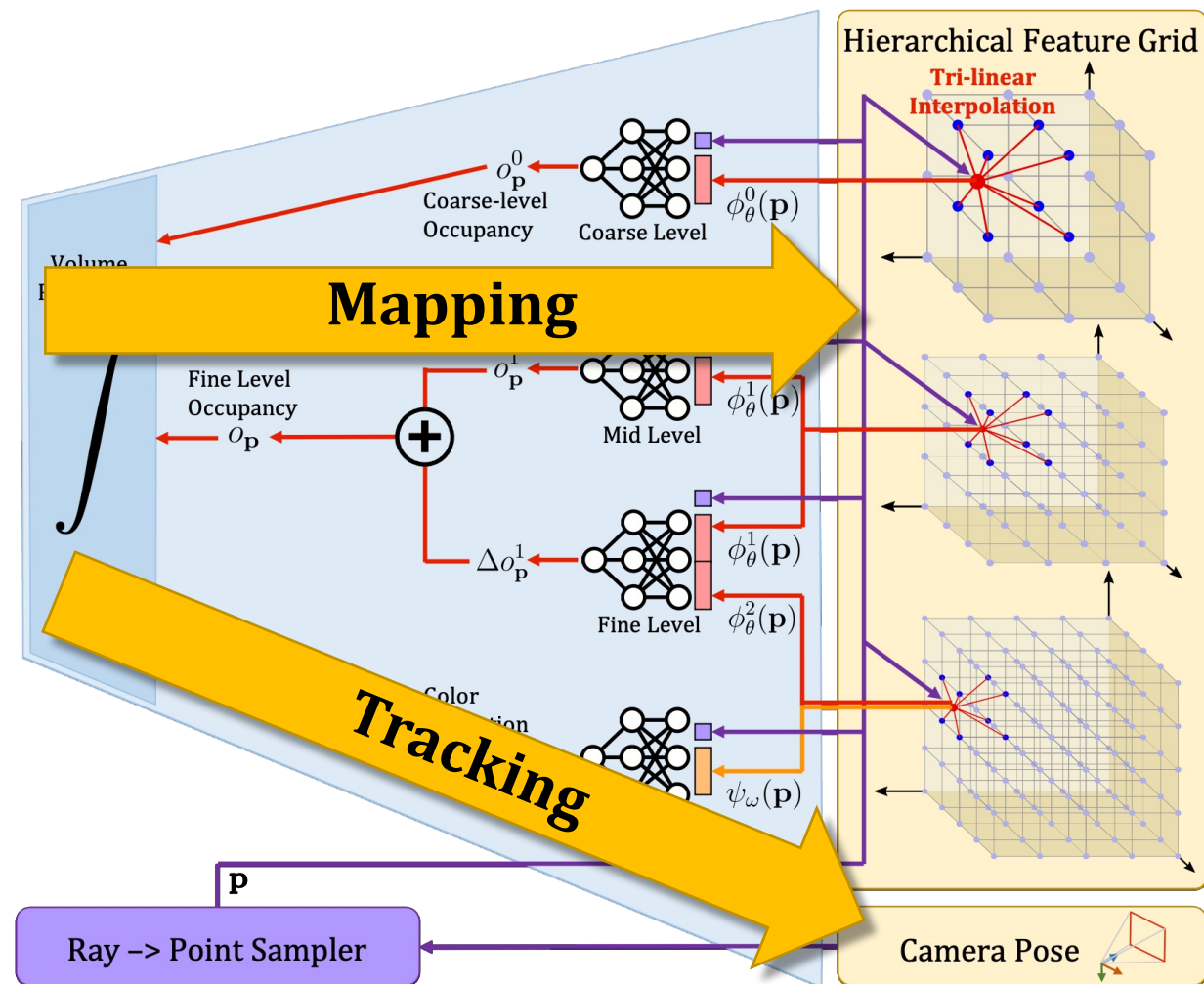


# Pipeline

Input Depth



Input RGB



# Results



# iMAP\*

(our re-implementation of iMAP)

# NICE-SLAM

4x Speed

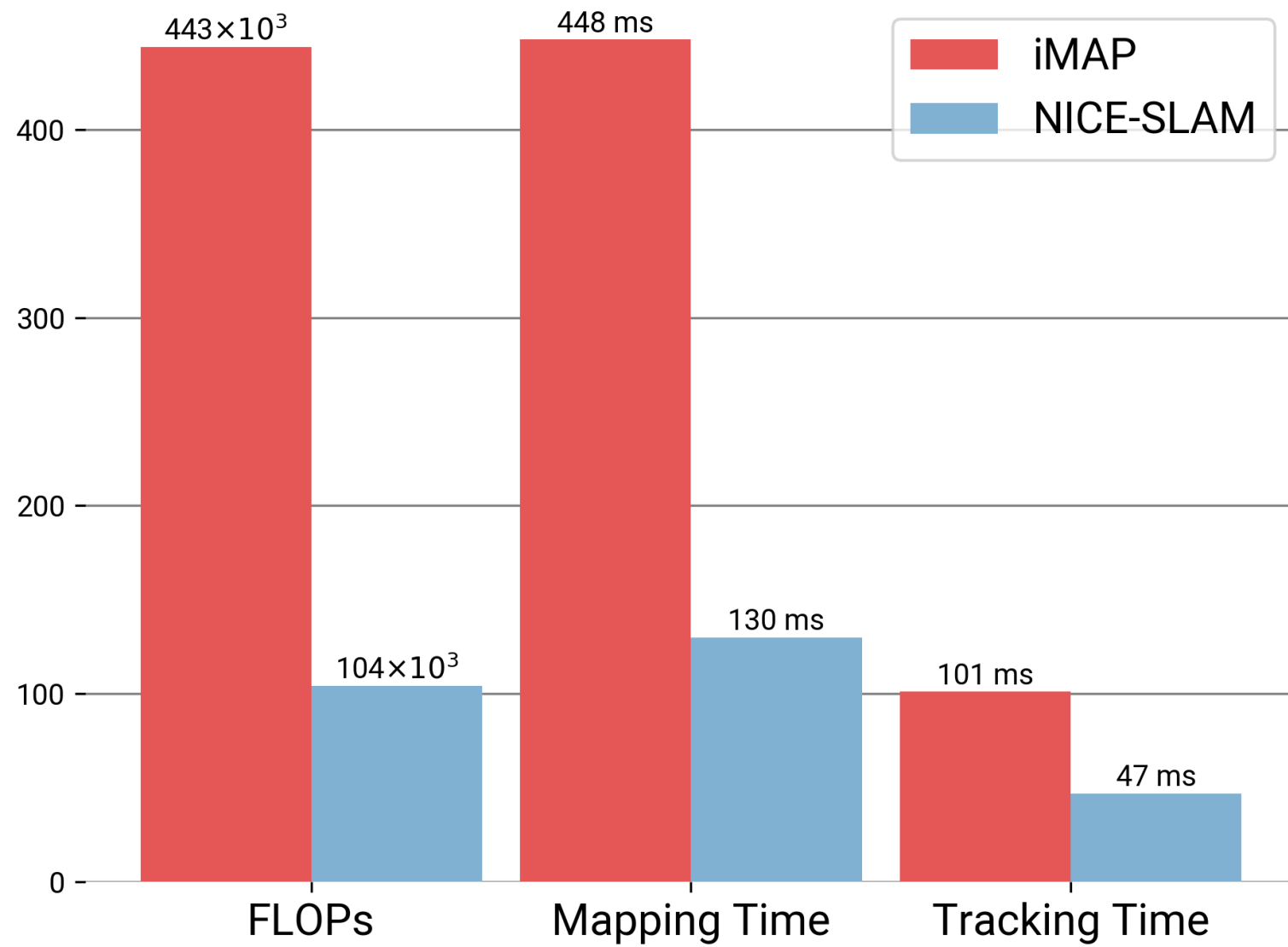
— Predicted Poses  
— GT Poses

# 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 [TOG'22]

## Limitations

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

# Final Remarks

- NeRF-based multi-view surface reconstruction still has rooms to improve
- A completely COLMAP-free NeRF pipeline?
- What is THE representation?

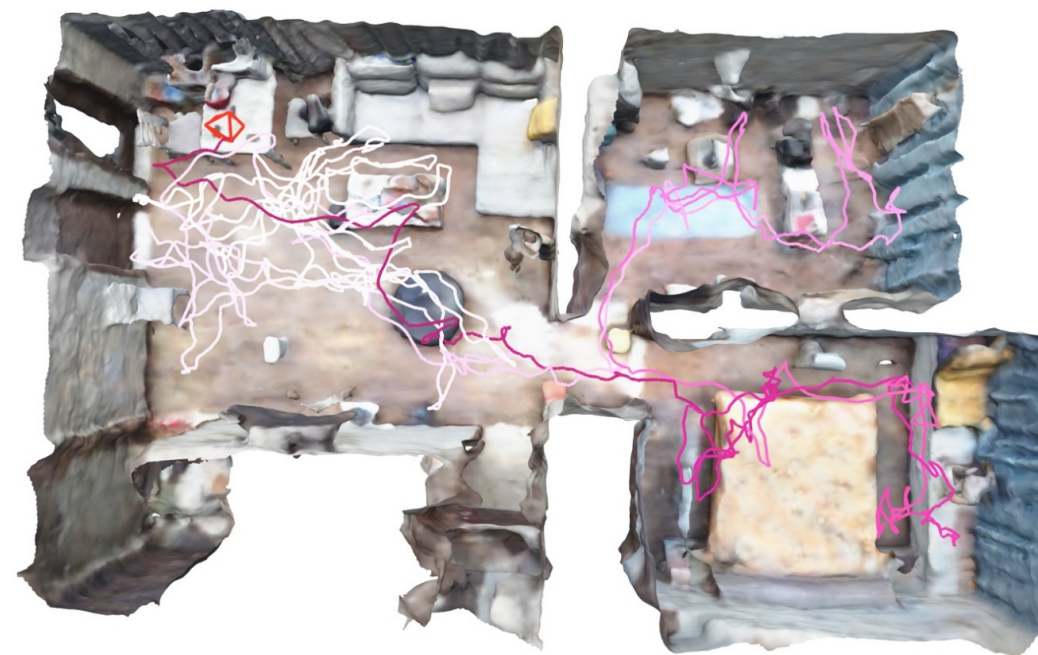


# Large-scale Scene Reconstruction with NeRF



MonoSDF

[github.com/autonomousvision/monosdf](https://github.com/autonomousvision/monosdf)



NICE-SLAM

[github.com/cvg/nice-slam](https://github.com/cvg/nice-slam)

Thank you!