# Learning Neural Scene Representations for 3D Reconstruction and Understanding

Songyou Peng

ETH Zurich and Max Planck Institute for Intelligent Systems

Shanghai AI Lab

June 15, 2023

# Who Am I?

- **Final-year PhD Student**
  - Marc Pollefeys
  - Andreas Geiger

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

- Before PhD, worked in Singapore, and interned at INRIA and TUM

**ETH** *zürich*

**MAX PLANCK INSTITUTE**
FOR INTELLIGENT SYSTEMS

∞ Meta

Google Research

pengsongyou.github.io

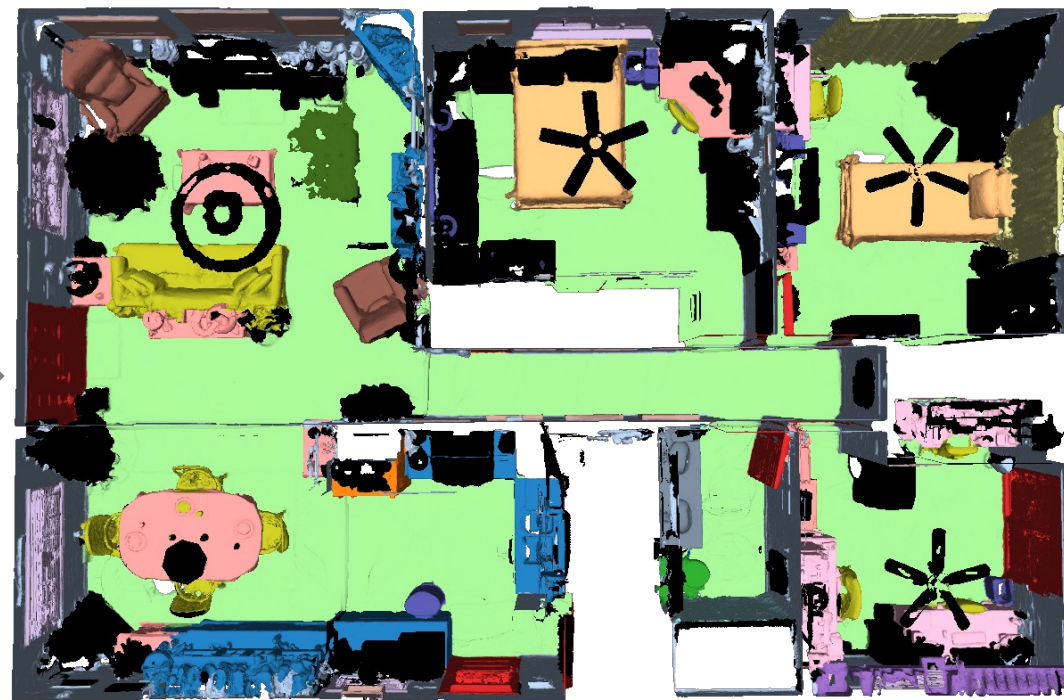# Motivation



Input Images

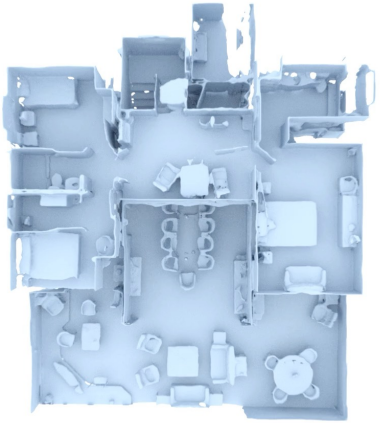**3D Reconstruction**

# Motivation
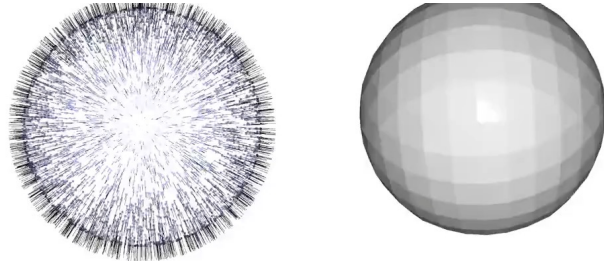


**3D Reconstruction**

**3D Scene Understanding**

# My PhD Topics: Neural Scene Representations
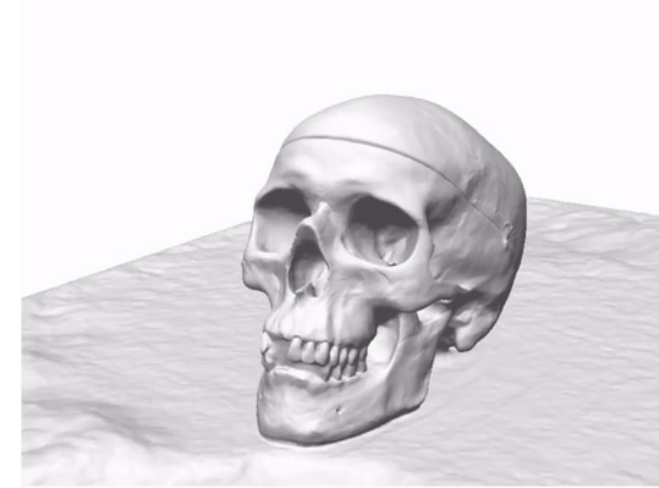## for <u>3D reconstruction</u> and <u>3D scene understanding</u>
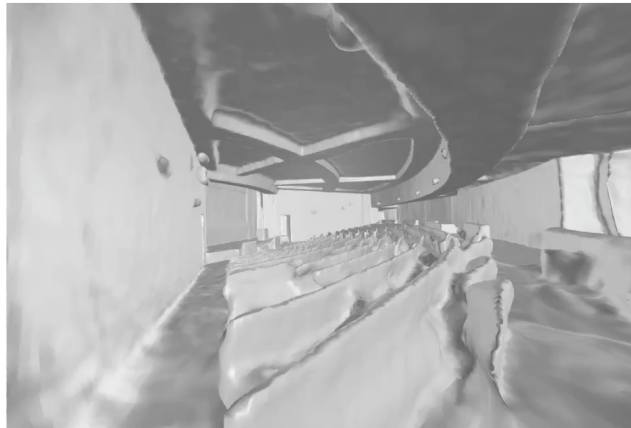


**Convolutional Occupancy Nets**
ECCV 2020 (Spotlight)

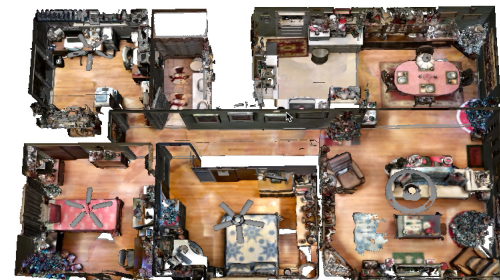**Shape As Points**
NeurIPS 2021 (Oral)
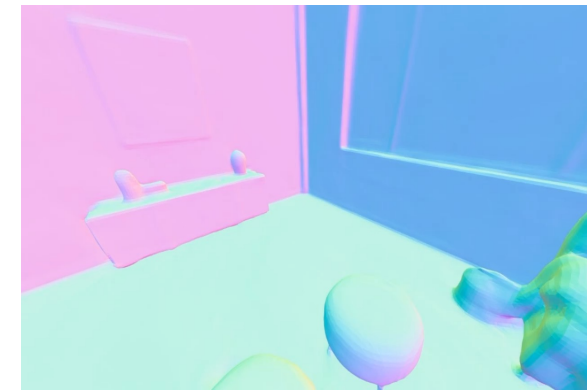
**KiloNeRF**
ICCV 2021

**UNISURF**
ICCV 2021 (Oral)

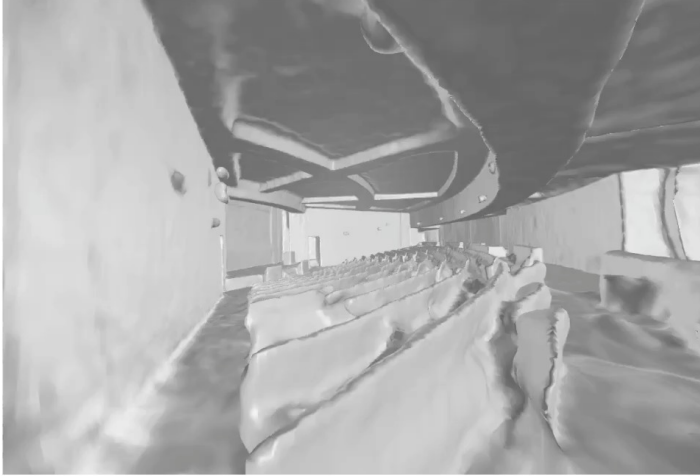**MonoSDF**
NeurIPS 2022

**NICE-SLAM**
CVPR 2022

**OpenScene**
CVPR 2023

**NICER-SLAM**
arXiv 2023

# My PhD Topics: Neural Scene Representations
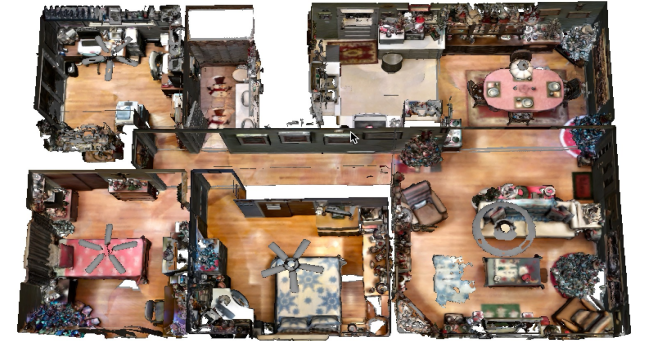## for 3D reconstruction and 3D scene understanding



Ours

**MonoSDF**

NeurIPS 2022

**NICE-SLAM**

CVPR 2022



floor

**OpenScene**

CVPR 2023

# NeRF is awesome!



**Some existing problems…**

😢 Poor underlying geometry
😢 Camera poses needed

😊 MonoSDF
😊 NICE-SLAM

Mildenhall*, Srinivasan*, Tancik* et al: <u>NeRF : Representing Scenes as Neural Radiance Fields for View Synthesis</u>. ECCV 2020
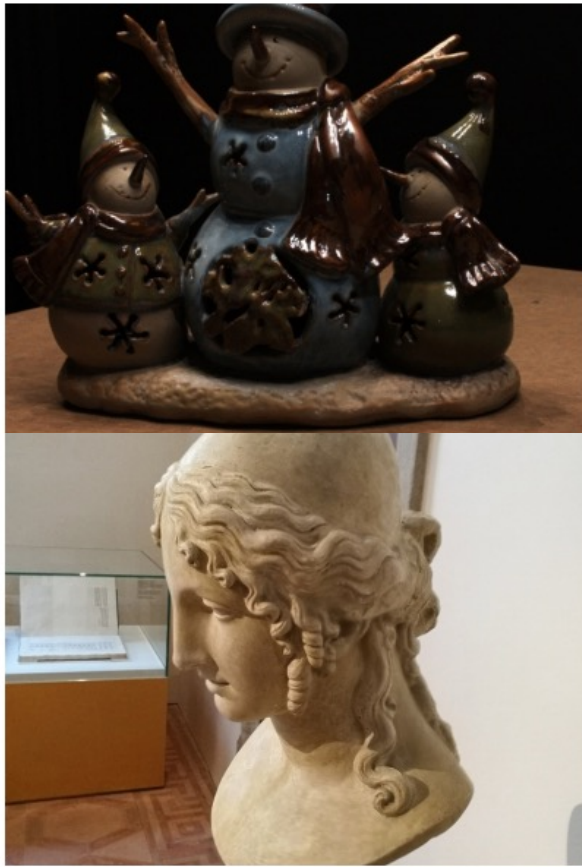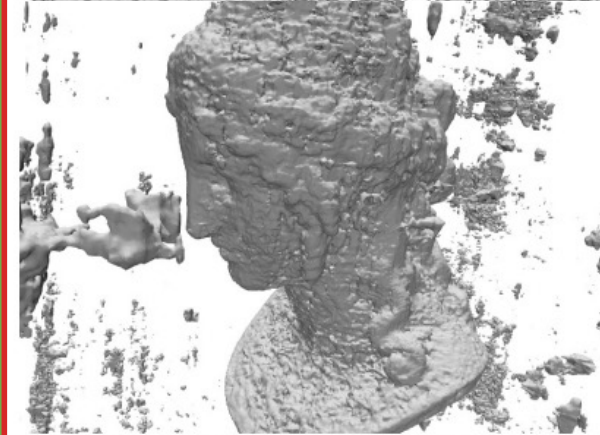
# Neural Implicit Surfaces with Volume Rendering
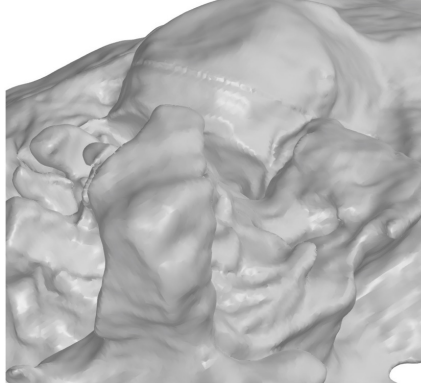


RGB Images      VolSDF/NeuS/UNISURF      NeRF

[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
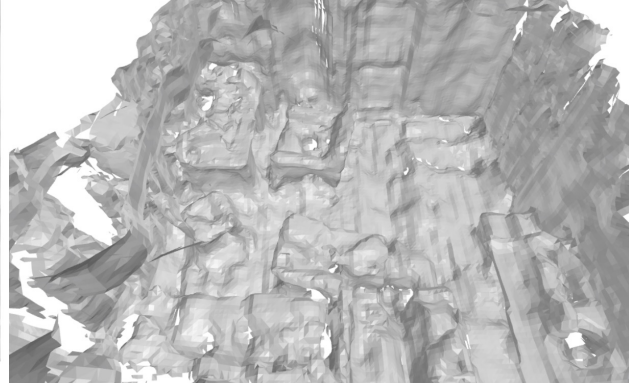[3] Yariv, Gu, Kasten, Lipman: Volume rendering of neural implicit surfaces. NeurIPS, 2021

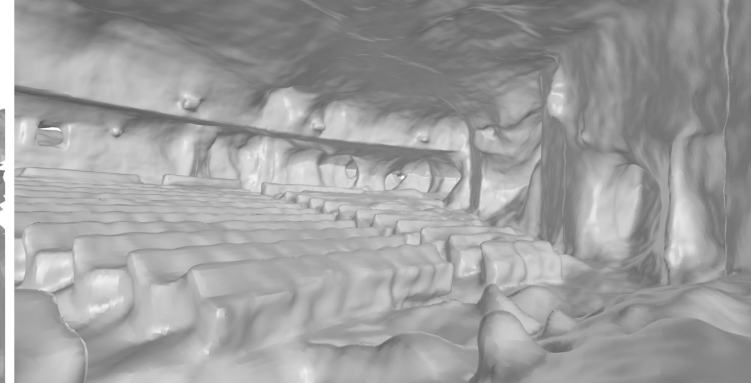# Neural Implicit Surfaces with Volume Rendering

DTU (3 views)  ScanNet (464 views)  Tanks & Temples (298 views)

VolSDF



— Fails with sparse input views

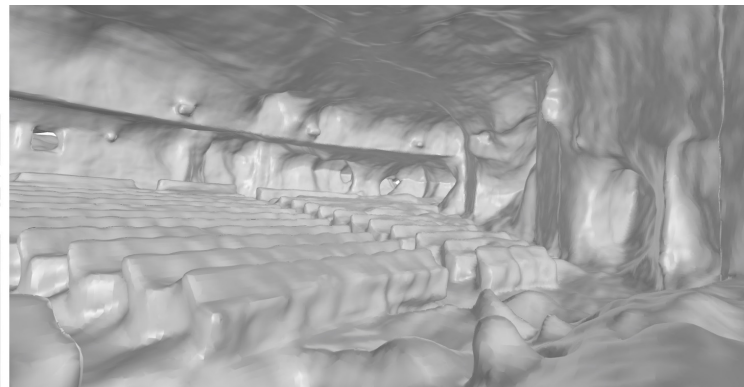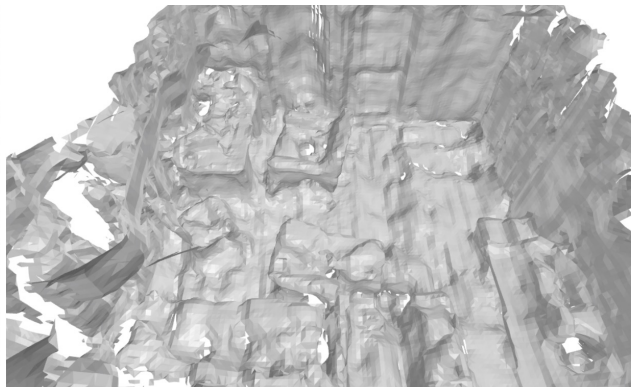— Poor results in large-scale indoor scenes

Yariv, Gu, Kasten, Lipman: <u>Volume rendering of neural implicit surfaces</u>. NeurIPS, 2021
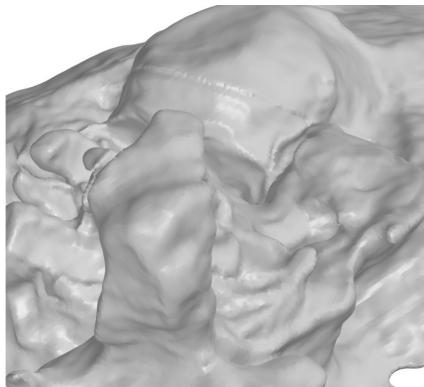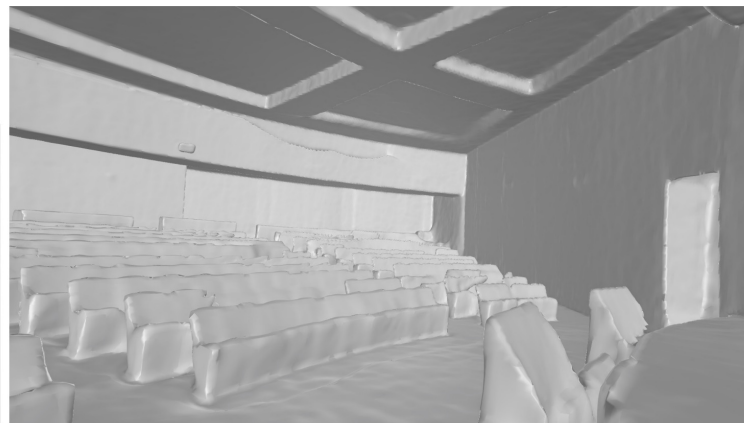
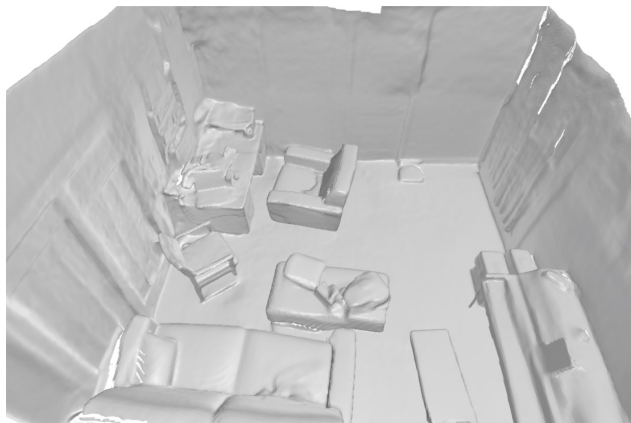# Neural Implicit Surfaces with Volume Rendering



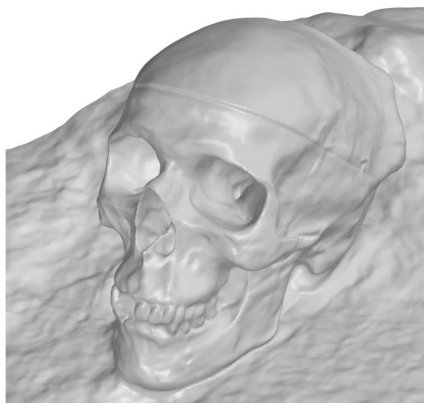DTU (3 views)  ScanNet (464 views)  Tanks & Temples (298 views)
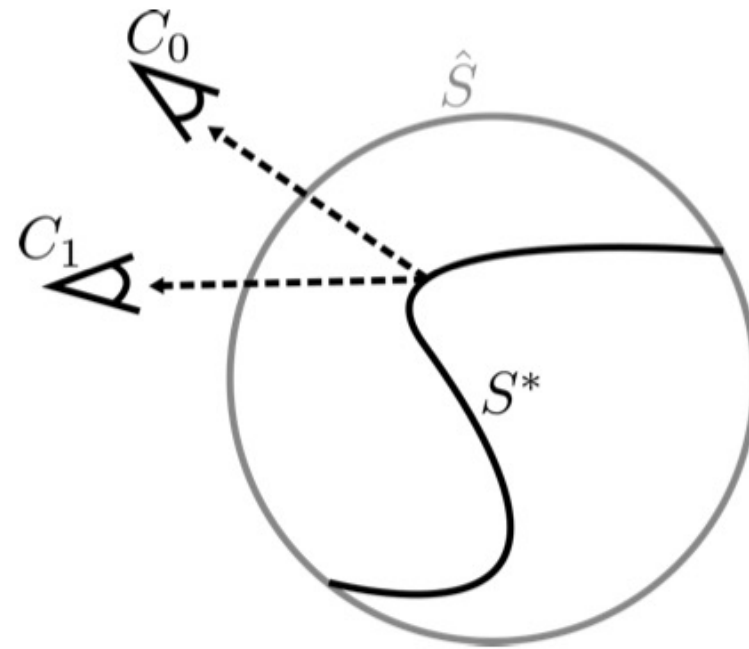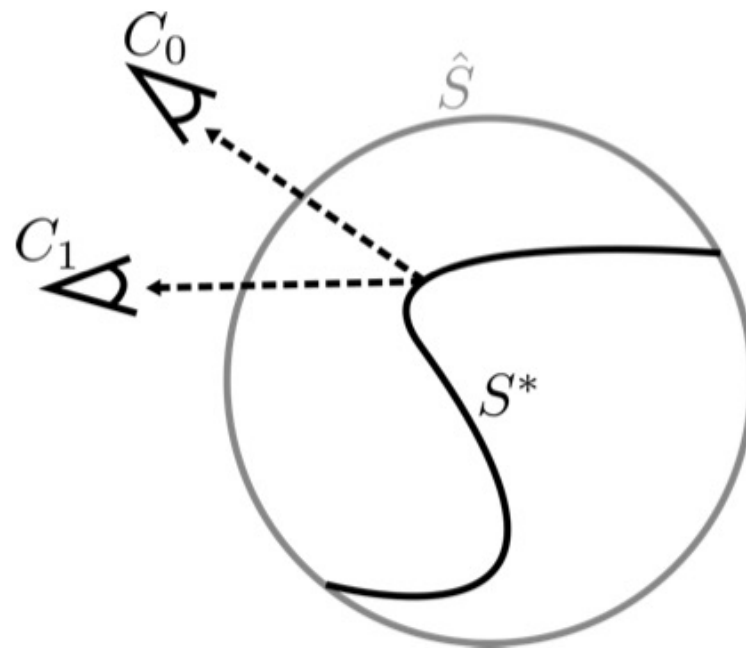
VolSDF

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!

Zhang, Riegler, Snavely, Koltun: NeRF++: Analyzing and Improving Neural Radiance Fields. ArXiv, 2020
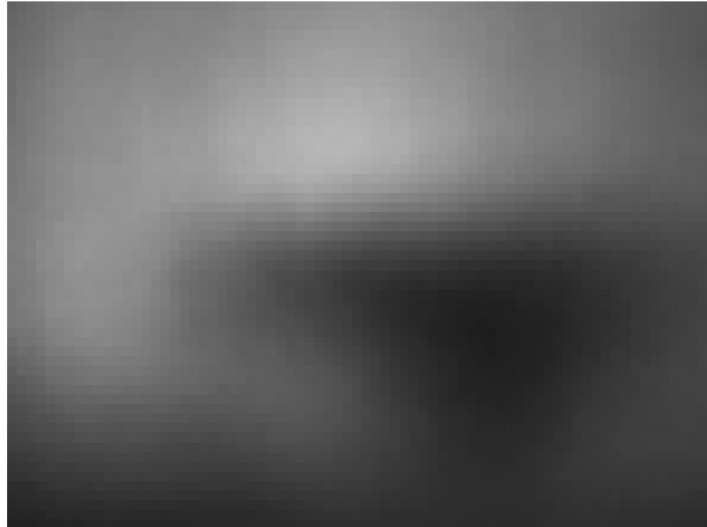
# Shape-Appearance Ambiguity



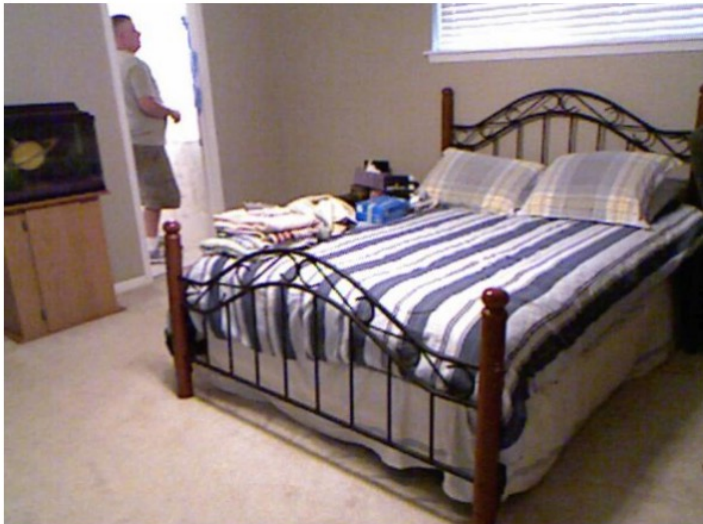There exists an infinite number of photo-consistent explanations for input images!

$\Longrightarrow$ Exploit monocular geometric priors

Zhang, Riegler, Snavely, Koltun: NeRF++: Analyzing and Improving Neural Radiance Fields. ArXiv, 2020

# Depth Map Prediction from a Single Image



Eigen, Puhrsch and Fergus: Depth Map Prediction from a Single Image using a Multi-Scale Deep Network. NIPS, 2014

# Omnidata



[Ranftl et al. 2021]

Eftekhar, Sax, Malik and Zamir: Omnidata: A Scalable Pipeline for Making Multi-Task Mid-Level Vision Datasets from 3D Scans. ICCV, 2021.

# Omnidata



RGB Image

Omnidata Normal

Omnidata Depth

Eftekhar, Sax, Malik and Zamir: Omnidata: A Scalable Pipeline for Making Multi-Task Mid-Level Vision Datasets from 3D Scans. ICCV, 2021.

# MonoSDF

# MonoSDF



Input Views

# MonoSDF



Input Views

# MonoSDF



MLP

Dense SDF Grid

Single-res Feature Grid

Multi-res Feature Grids

Neural Implicit Scene Representation

Input Views

# MonoSDF

# MonoSDF



MLP

Dense SDF Grid

Single-res Feature Grid

Multi-res Feature Grids

Neural Implicit Scene Representation

Volume Rendering

Input Views

Ray Distance

# MonoSDF



MLP

Dense SDF Grid

Single-res Feature Grid

Multi-res Feature Grids

$\mathbf{x} \rightarrow$ $f_\theta$ $\rightarrow \hat{s}$

Interpolation

$\mathbf{x} \rightarrow$ $\rightarrow \hat{s}$

Interpolation

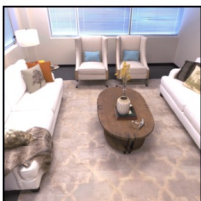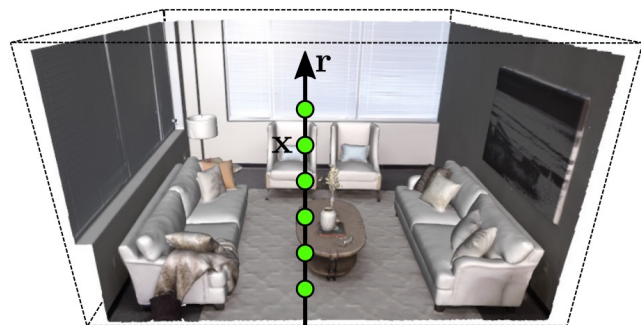$\mathbf{x}$ $f_\theta \rightarrow \hat{s}$

Interpolation

$\mathbf{x}$ ... ... $f_\theta \rightarrow \hat{s}$

Interpolation

Neural Implicit Scene Representation

$\mathbf{r}$

$\mathbf{x}$

Volume Rendering

$\hat{C}(\mathbf{r})$

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

$\hat{s}$

$\sigma$

Input Views

$C$

Ray Distance

# MonoSDF

# MonoSDF

# MonoSDF



Neural Implicit Scene Representation

MLP

Dense SDF Grid

Single-res Feature Grid

Multi-res Feature Grids

Volume Rendering

Input Views

Pretrained
Omnidata
Model

Monocular Geometric Cues

Ray Distance

# Ablation Study

| | | Normal C.↑ | Chamfer-$L_1$ ↓ | F-score ↑ |
|---|---|---|---|---|
| **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

DTU (3 views)

ScanNet

Tanks and Temples

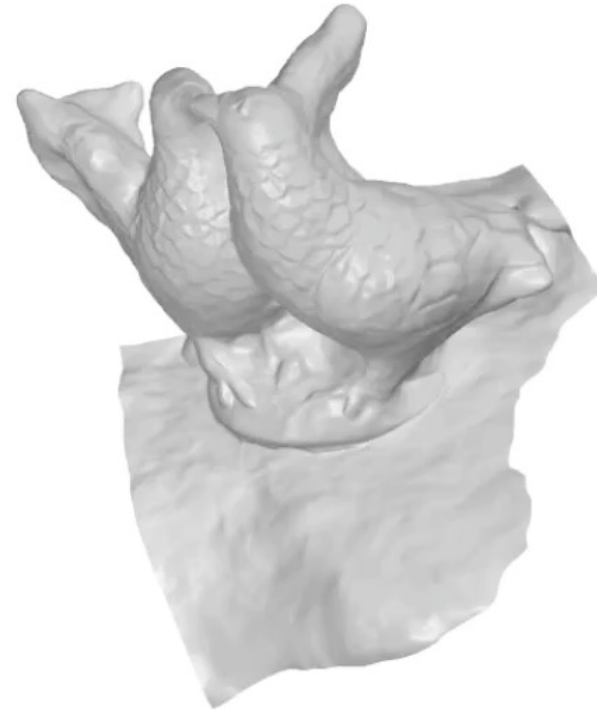! **Monocular cues** improve reconstruction results and speed up optimization

! Inspire applications in other fields [GOOD, ICLR 2023]

! <u>Limitation</u>: Still require camera poses given :(

# RGB-D Sequences





**4**0x Speed

# 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

# iMAP
[Sucar et al., ICCV'21]



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

# iMAP
[Sucar et al., ICCV'21]



A single MLP

▬ Fail when scaling up to larger scenes

▬ Global update → Catastrophic forgetting

▬ Slow convergence

Predicted Poses
GT Poses

# NICE-SLAM



Feature grids + tiny MLPs

➕ Applicable to **large-scale scenes**

➕ Local update → **No forgetting problem**

➕ **Fast** convergence

— Predicted Poses
— GT Poses

38

# Pipeline

# Results

# iMAP*
(our re-implementation of iMAP)

# NICE-SLAM

4x Speed

| | |
|---|---|
| <span style="color:red">———</span> | Predicted Poses |
| ——— | GT Poses |

# iMAP*
### (our re-implementation of iMAP)

# NICE-SLAM

**10x Speed**

Note: Runtime evaluation setting from iMAP paper, not the best-performing setting

# Take-home Message

- A NICE NeRF-based SLAM system for indoor scenes

- Hierarchical feature grids + a tiny MLP **seems to be a trend**!

  - Instant-NGP [SIGGRAPH'22 Best Paper]

**Limitations**

- <u>Requires depths as input</u>

- Only bounded scenes

- Still not real-time

# NICER-SLAM: Neural Implicit Scene Encoding for RGB SLAM

Zihan Zhu[1*]   Songyou Peng[1,2*]   Viktor Larsson[3]   Zhaopeng Cui[4]
Martin R. Oswald[1,5]   Andreas Geiger[6]   Marc Pollefeys[1,7]

[1]ETH Zürich   [2]MPI for Intelligent Systems, Tübingen   [3]Lund University
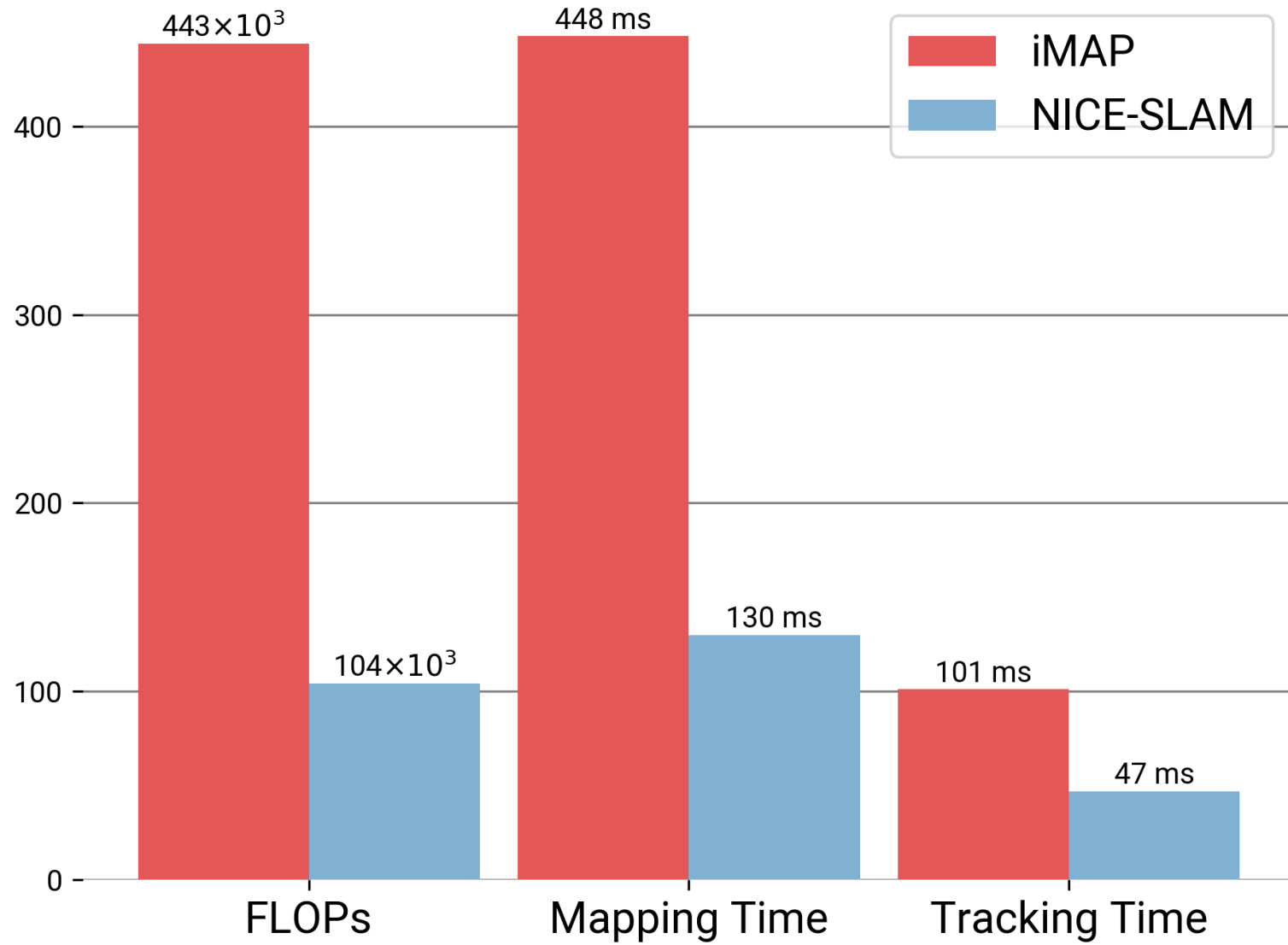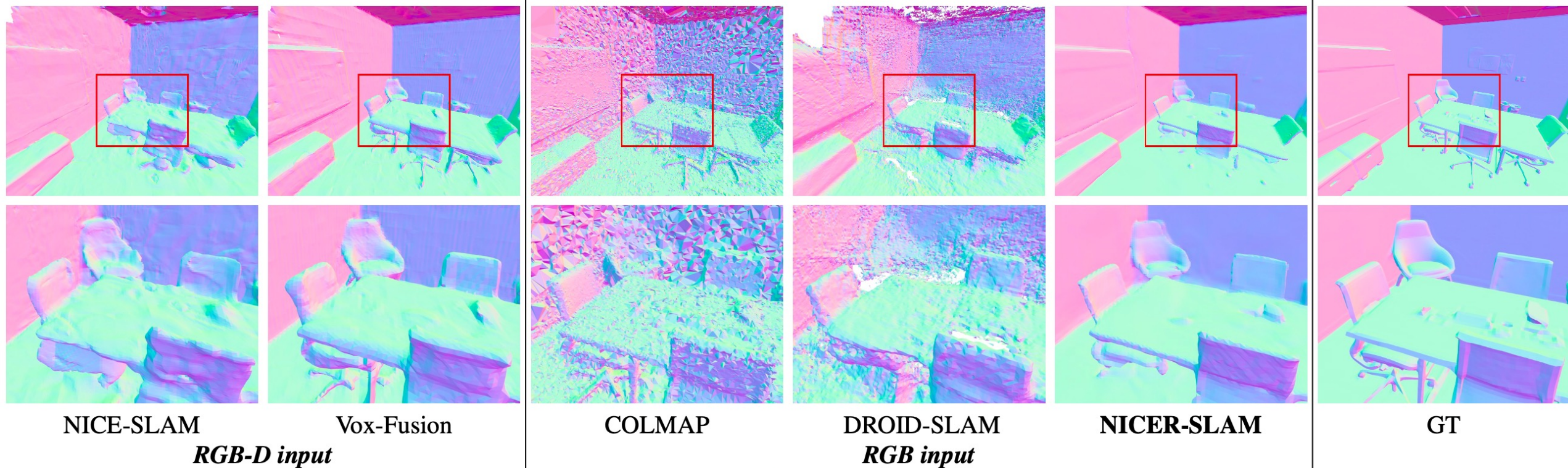[4]State Key Lab of CAD&CG, Zhejiang University   [5]University of Amsterdam
[6]University of Tübingen, Tübingen AI Center   [7]Microsoft

| NICE-SLAM | Vox-Fusion | COLMAP | DROID-SLAM | **NICER-SLAM** | GT |

*RGB-D input*   *RGB input*

https://arxiv.org/abs/2302.03594

Input 3D Geometry

Input 3D Geometry

Traditional Semantic Segmentation

Only train and test on a few common classes

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

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



**CLIP**: Contrastive Language-Image Pre-Training

Radford et al.: Learning Transferable Visual Models From Natural Language Supervision. ICML 2021

# 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



3D Geometry

CLIP Text Features
(visualize with T-SNE)

RGB Images

Note: bold word embeddings are approximate

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

# Step 1: Multi-view Feature Fusion



3D Geometry

$\mathbf{f}^{2D}$

Per-pixel Features
(visualize with PCA)

$\mathcal{E}^{2D}$

OpenSeg [1]
LSeg [2]

RGB Images

[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



3D Geometry

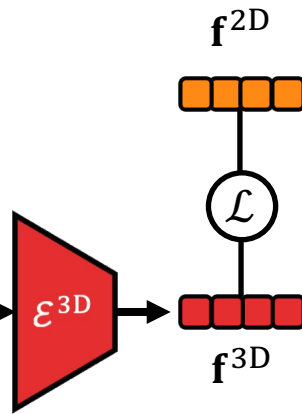$$\mathbf{s}_n^{2D} = \cos(\mathbf{f}^{2D}, \mathbf{t}_n)$$
$$\mathbf{s}_n^{3D} = \cos(\mathbf{f}^{3D}, \mathbf{t}_n)$$

Choose the feature with
the highest max score among all prompts

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

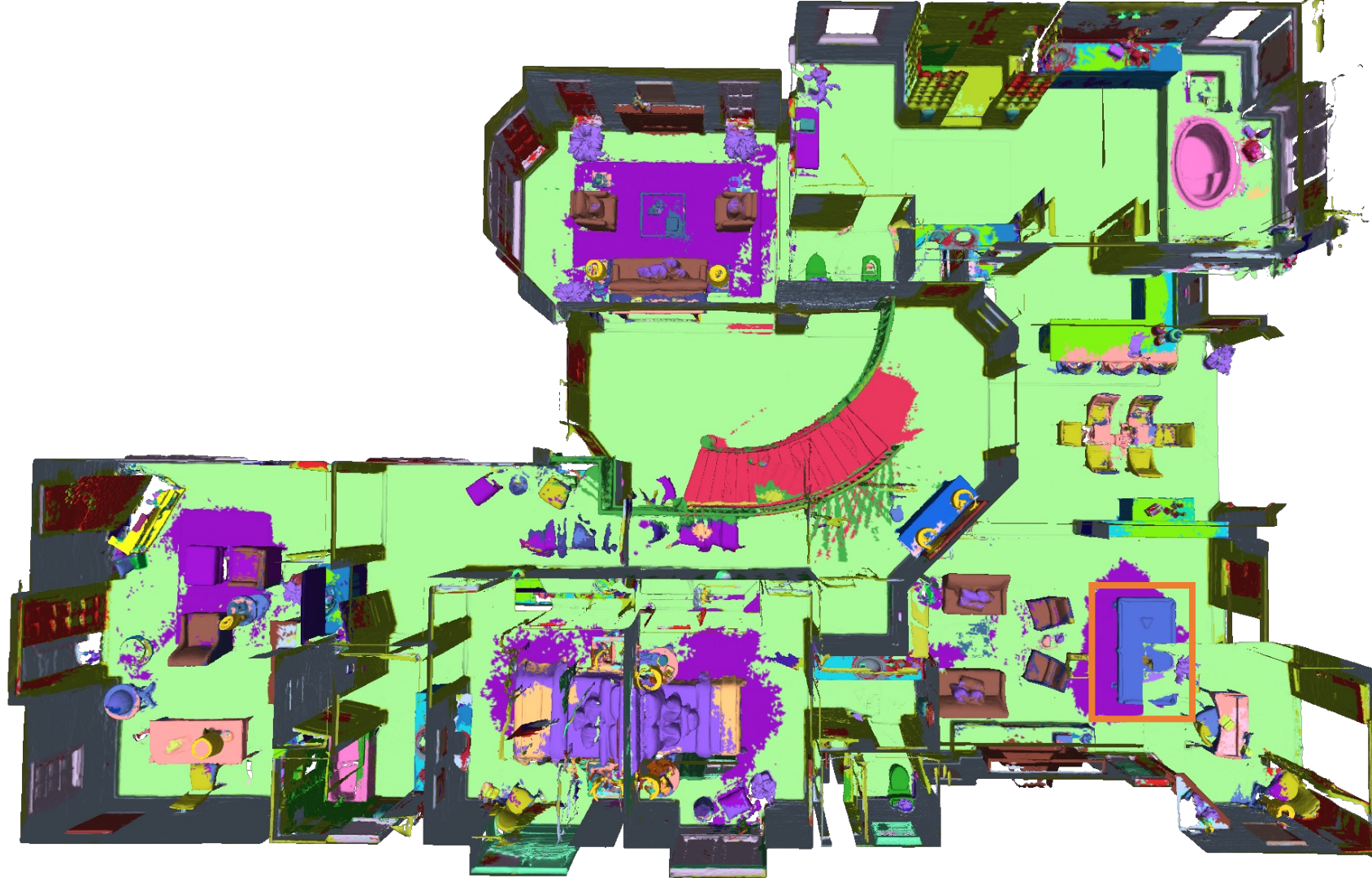# Open-Vocabulary, Zero-shot
## 3D Semantic Segmentation

Input 3D Geometry
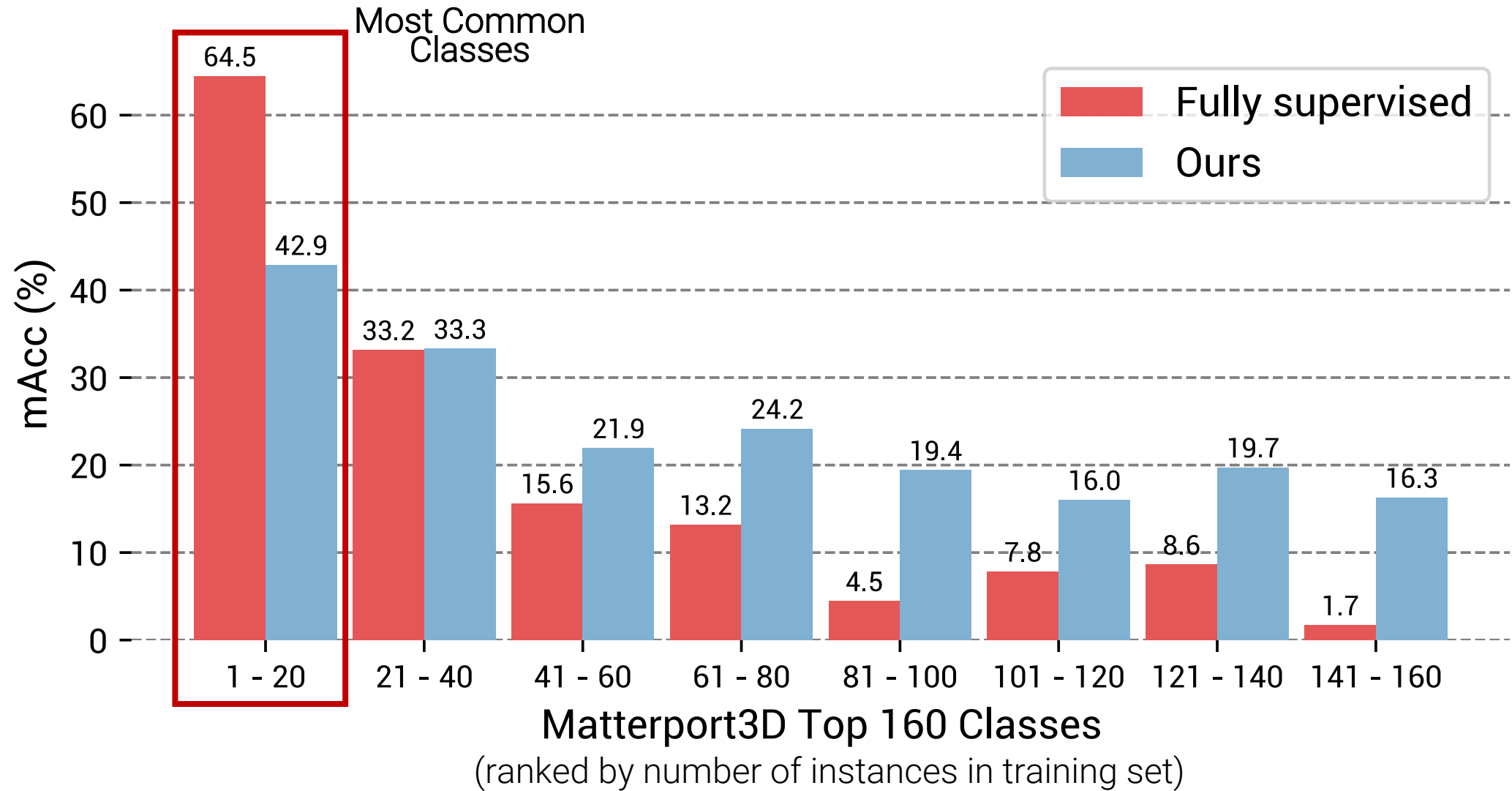
## Our Zero-shot 3D Segmentation
### (20 classes)

wall ■ floor ■ cabinet ■ bed ■ chair ■ sofa ■ table ■ door ■ window ■ bookshelf ■ picture ■ counter ■ desk ■ curtain ■ refrigerator ■ shower curtain ■ toilet ■ sink ■ bathtub ■ other
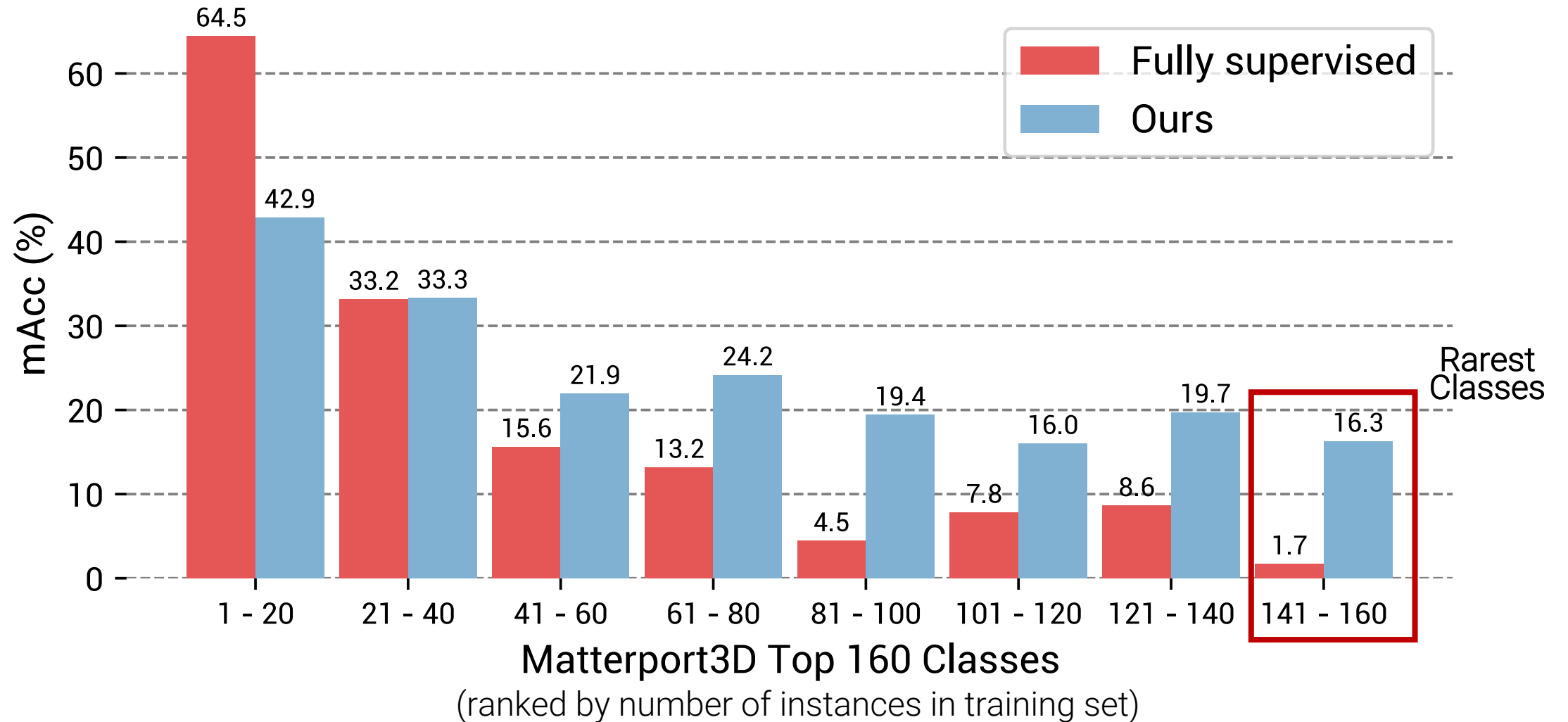
## Our Zero-shot 3D Segmentation
### (160 classes)

| | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 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 | pool table | car | soap dish |
| floor | plant | coffee table | faucet | bottle | refridgerator | light switch | shower curtain | heater | curtain rod | kitchen island | piano | scale | jacket | bottle of soap | cleaner |
| picture | mirror | counter | photo | purse | bookshelf | door way | bin | headboard | printer | paper towel | board | bag | water cooler | toilet brush | computer |
| window | towel | bench | toilet paper | toilet paper | chest | telephone | sheet | rope | display case | whiteboard | knob |
| chair | sink | stool | garbage bin | fan | wardrobe | basket | microwave | candle | blanket | glass | ball | toilet paper holder | tea pot | range hood | paper |
| pillow | shelves | vase | fireplace | railing | pipe | chandelier | blinds | flower pot | handle | dishwasher | excercise equipment | tray | stuffed animal | candelabra | projector |

61

# Comparison



62

# Comparison

# Ablation

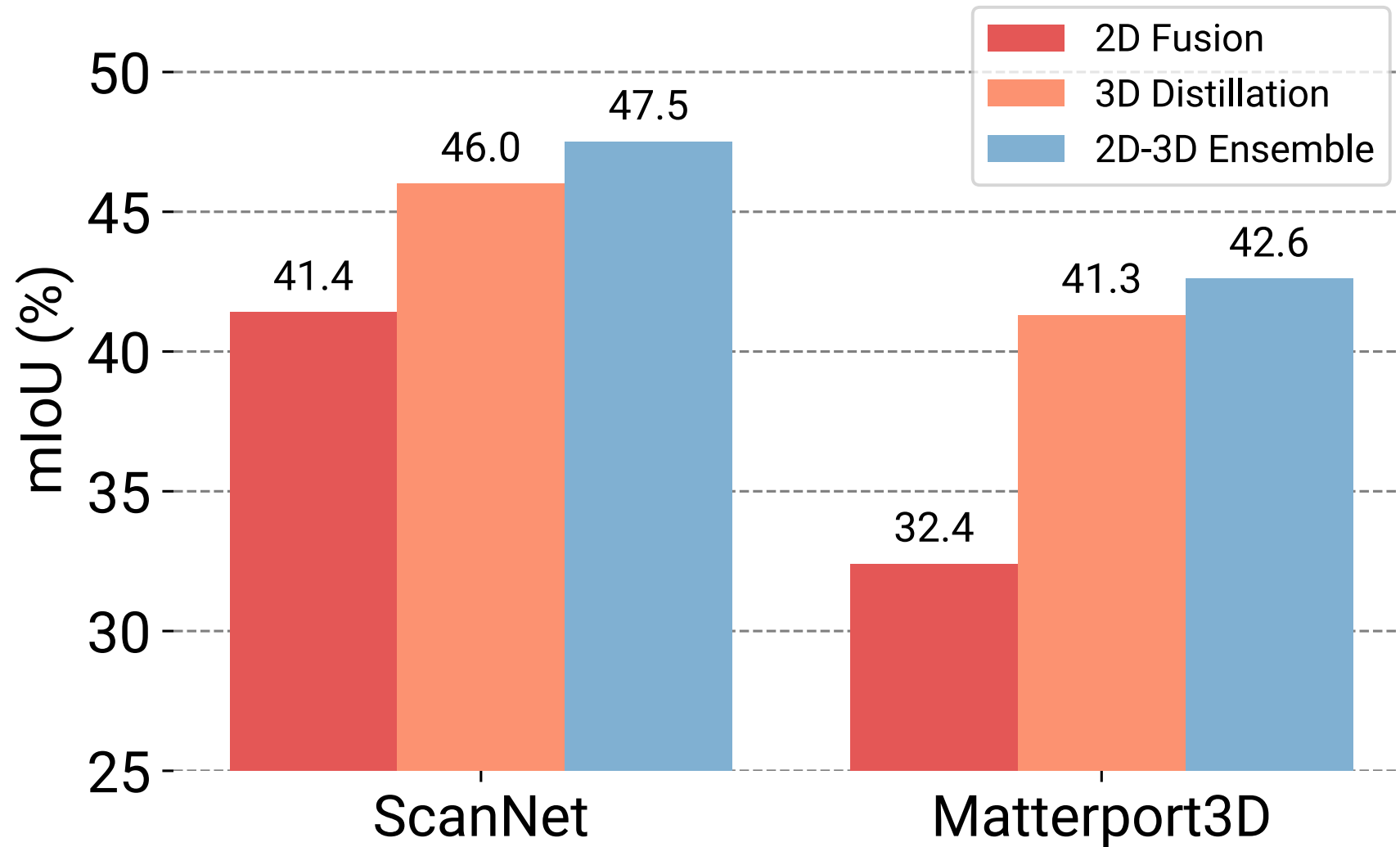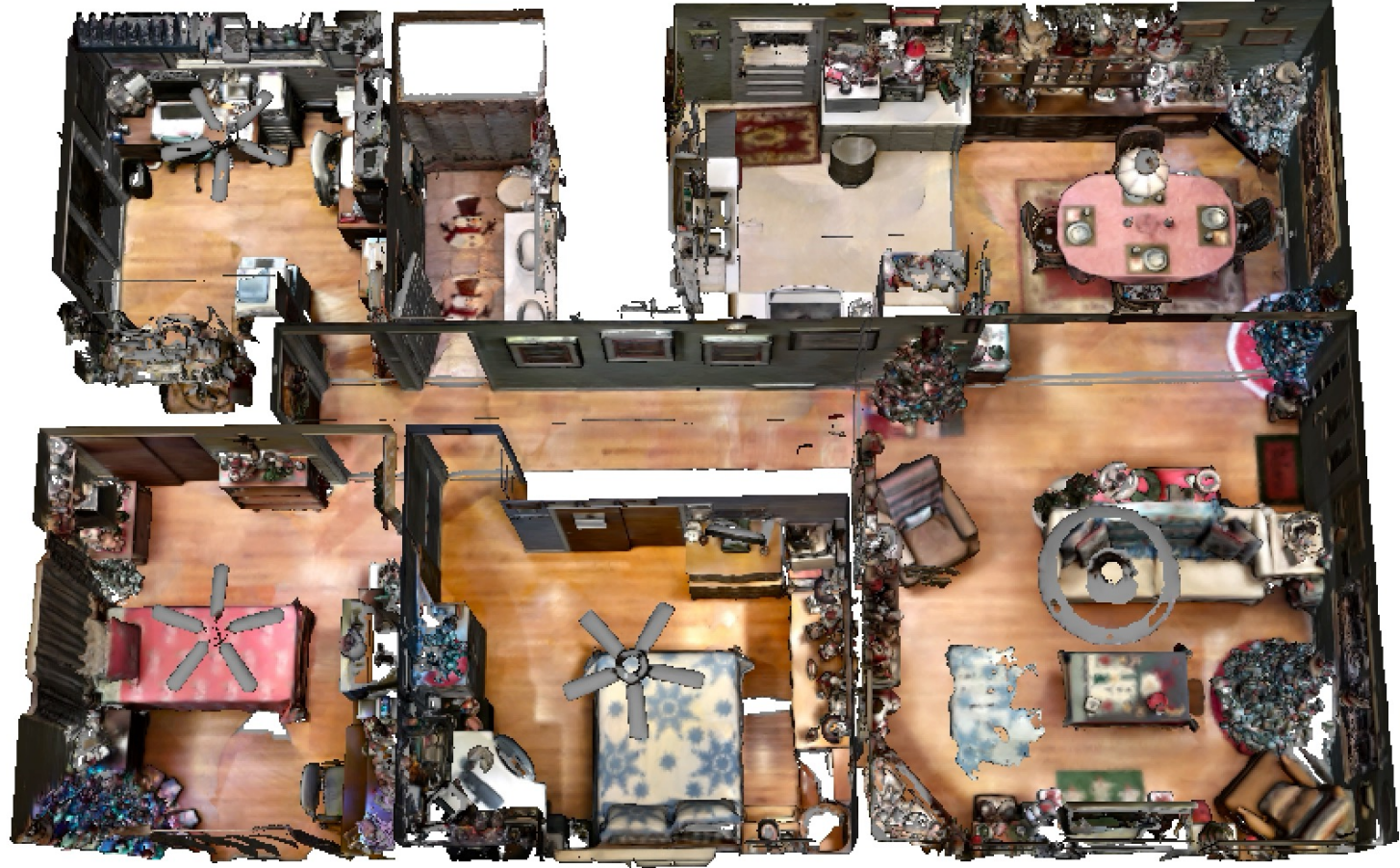# Image-based 3D Scene Query

Image Queries

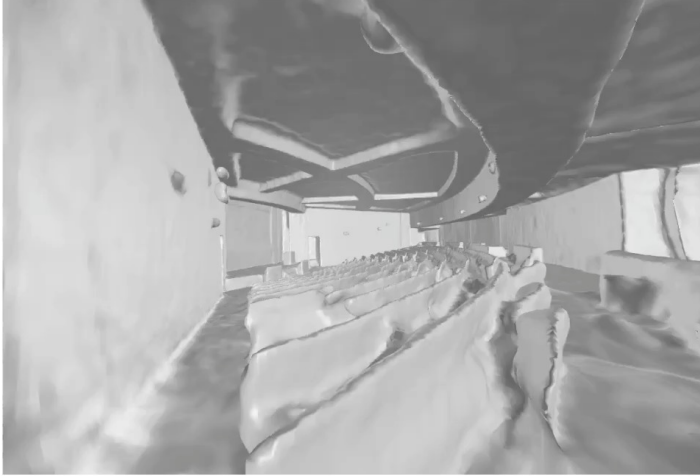Given 3D Geometry

# Interactive Demo

Open-vocabulary 3D Scene Exploration

# 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

- Our real-time demo already shows the **possibility to directly apply to AR/VR**

# My PhD Topics: Neural Scene Representations
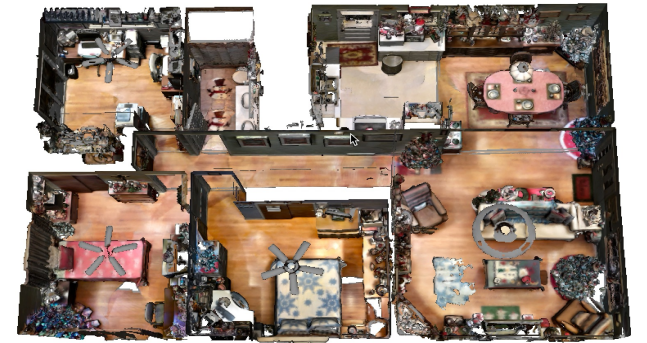## for <u>3D reconstruction</u> and <u>3D scene understanding</u>



Ours

floor

**MonoSDF**
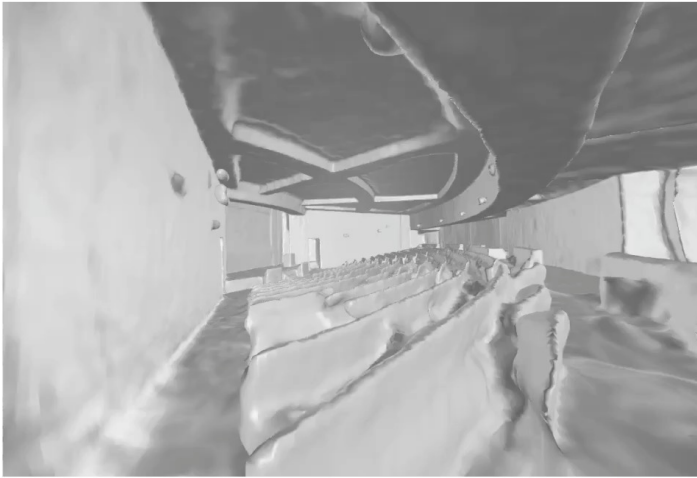NeurIPS 2022

**NICE-SLAM**
CVPR 2022

**OpenScene**
CVPR 2023

# Learning Neural Scene Representations for 3D Reconstruction and Understanding

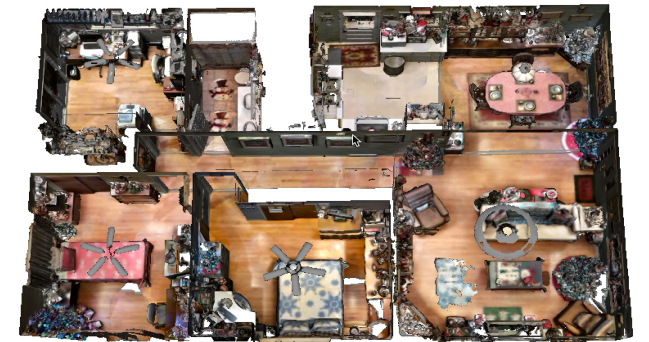Songyou Peng



Ours

**MonoSDF**
NeurIPS 2022
niujinshuchong.github.io/monosdf/

**NICE-SLAM**
CVPR 2022
pengsongyou.github.io/nice-slam

**OpenScene**
CVPR 2023
pengsongyou.github.io/openscene

# Thank you!