# Learning to Reconstruct and Understand the 3D World

Songyou Peng

ETH Zurich and Max Planck Institute for Intelligent Systems

Microsoft Mixed Reality & AI Lab – Zurich

May 31, 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

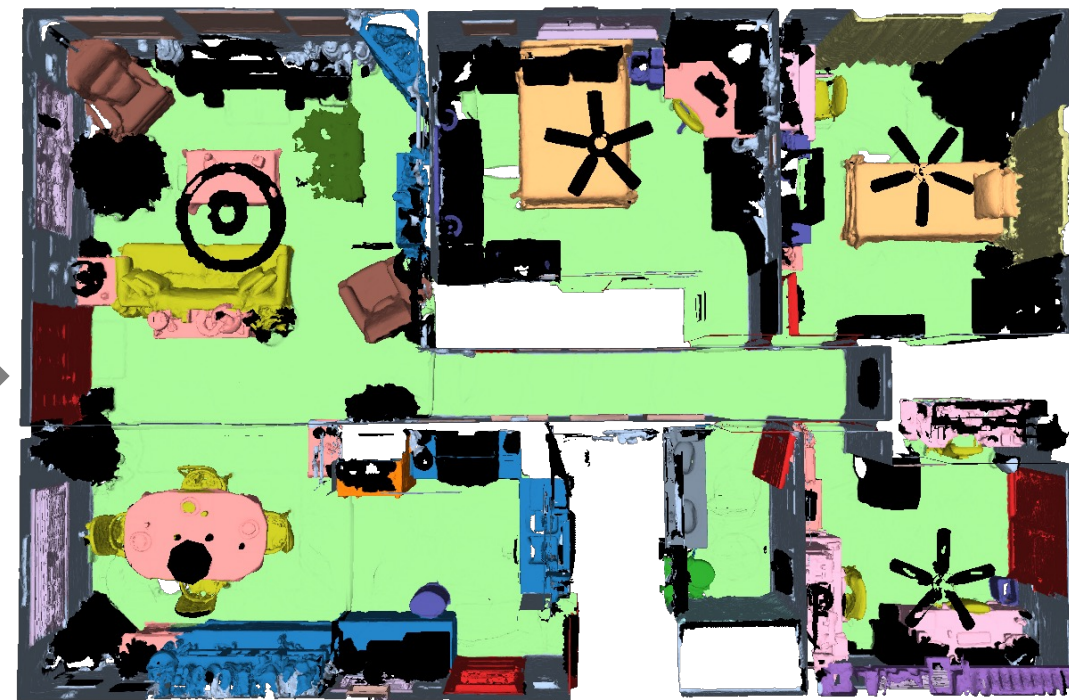# Motivation



Input Images

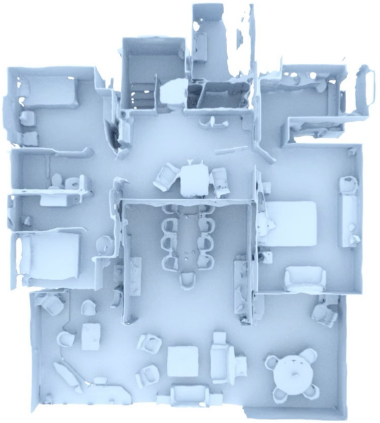**3D Reconstruction**

# Motivation

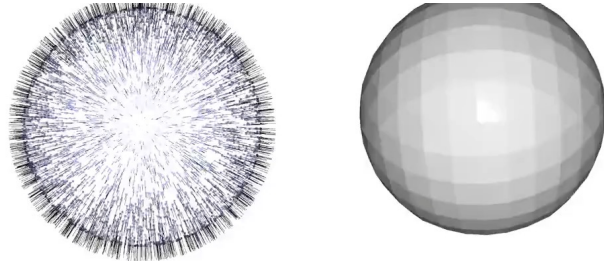

**3D Reconstruction**

**3D Scene Understanding**

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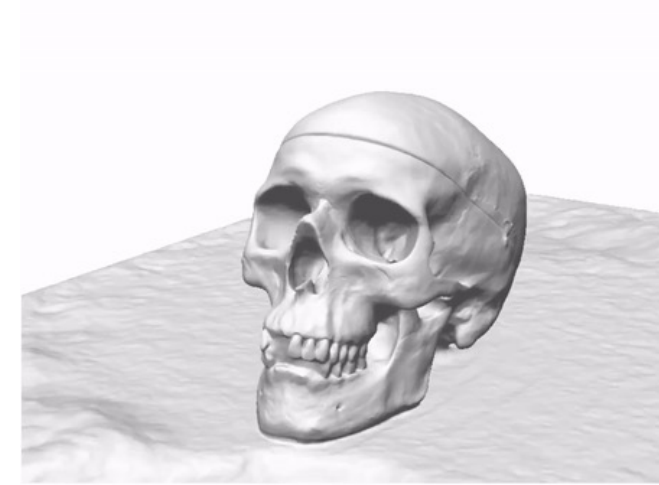## for <u>3D reconstruction</u> and <u>3D scene understanding</u>
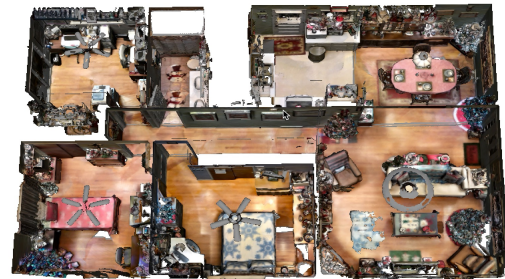


**Convolutional Occupancy Nets**
ECCV 2020 (Spotlight)

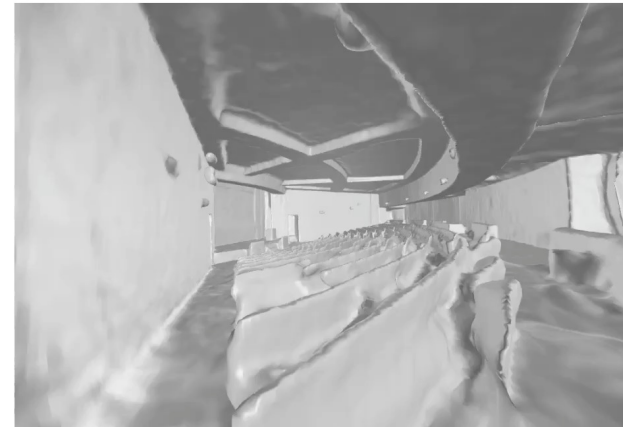**Shape As Points**
NeurIPS 2021 (Oral)
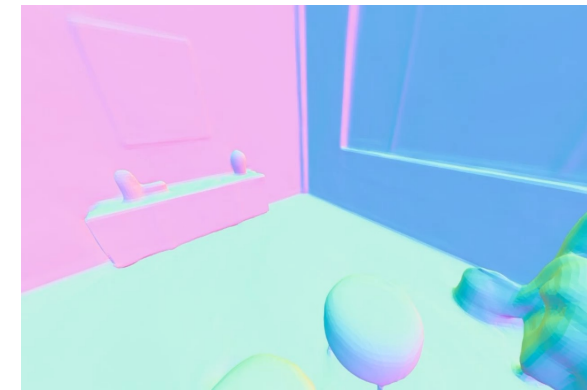
**KiloNeRF**
ICCV 2021

**UNISURF**
ICCV 2021 (Oral)
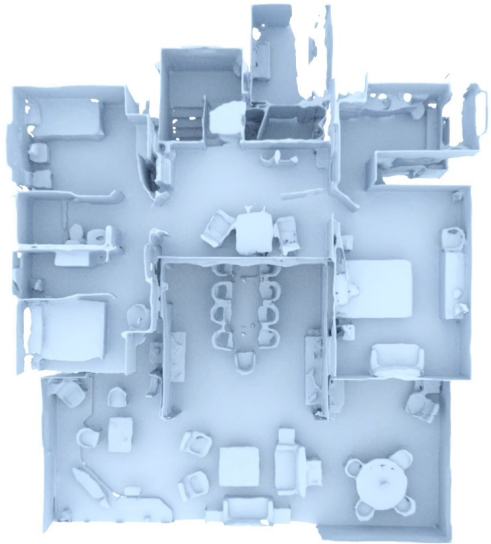
**NICE-SLAM**
CVPR 2022

**OpenScene**
CVPR 2023

**MonoSDF**
NeurIPS 2022

**NICER-SLAM**
arXiv 2023

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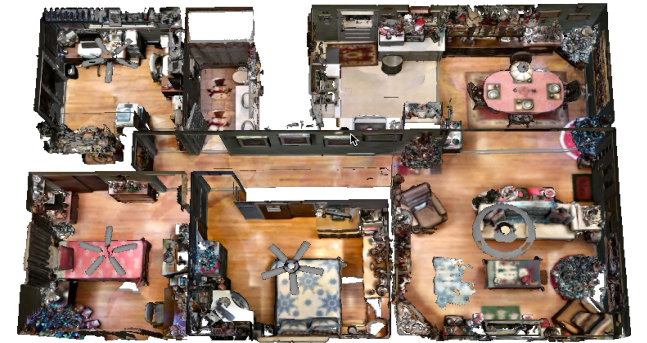## for 3D reconstruction and 3D scene understanding
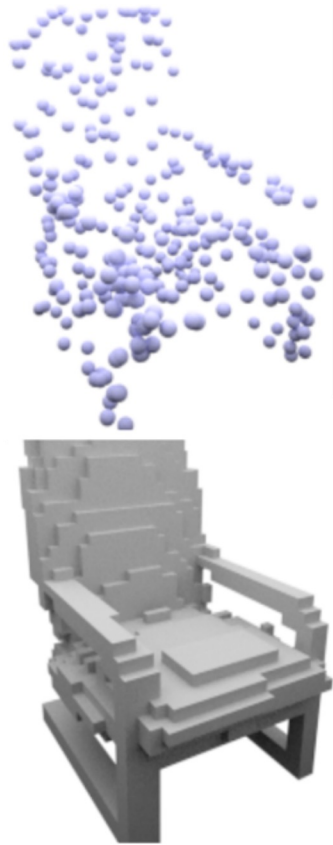


floor

**Convolutional Occupancy Networks**
ECCV 2020 (Spotlight)
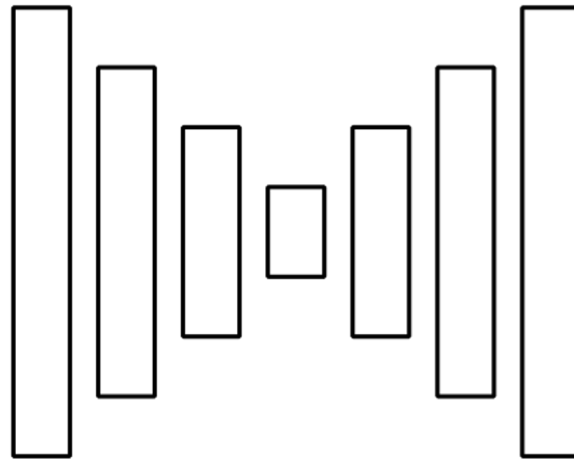
**NICE-SLAM**
CVPR 2022

**OpenScene**
CVPR 2023
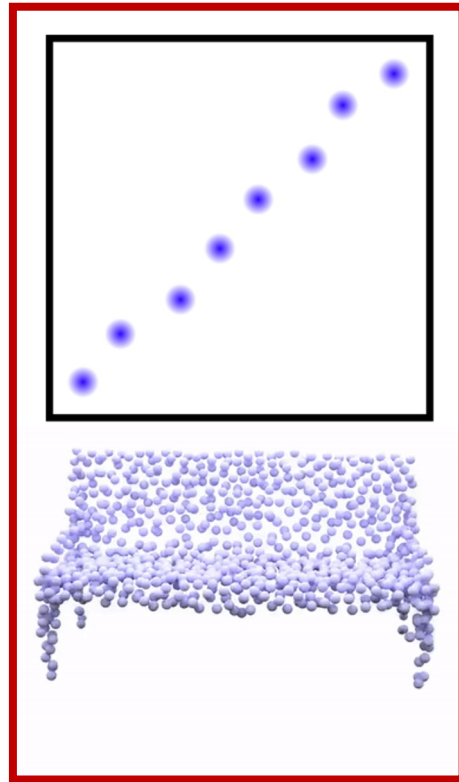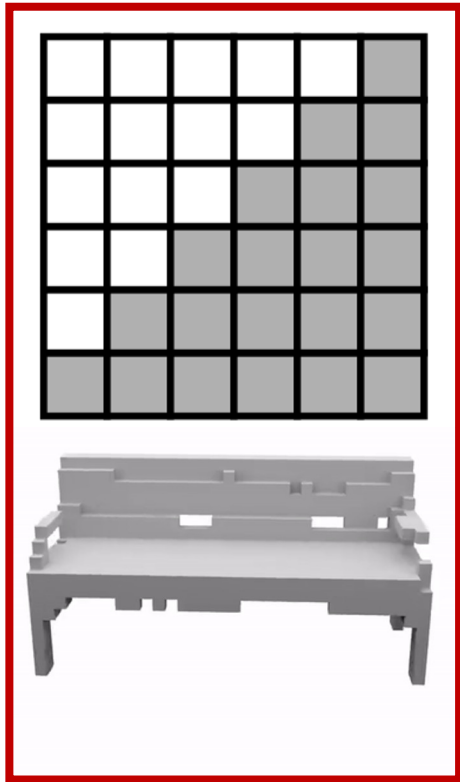
# Learning-based 3D Surface Reconstruction



Input

Neural Network

3D Reconstruction

# What is a good **3D representation?**

# 3D Representations



- Traditional Explicit Representations ⇒ **Discrete**

Mescheder, Oechsle, Niemeyer, Nowozin and Geiger: Occupancy Networks: Learning 3D Reconstruction in Function Space. CVPR, 2019

# 3D Representations



$$f_\theta(p) = \tau$$

- Traditional Explicit Representations ⇒ **Discrete**
- Implicit Neural Representation ⇒ **Continuous**

Mescheder, Oechsle, Niemeyer, Nowozin and Geiger: Occupancy Networks: Learning 3D Reconstruction in Function Space. CVPR, 2019

# Limitations

**Structure of neural implicit representations:**

Mescheder, Oechsle, Niemeyer, Nowozin and Geiger: Occupancy Networks: Learning 3D Reconstruction in Function Space. CVPR, 2019

# Limitations

**Structure of neural implicit representations:**



- Global latent code ⇒ **overly smooth geometry**

Mescheder, Oechsle, Niemeyer, Nowozin and Geiger: Occupancy Networks: Learning 3D Reconstruction in Function Space. CVPR, 2019

# Limitations

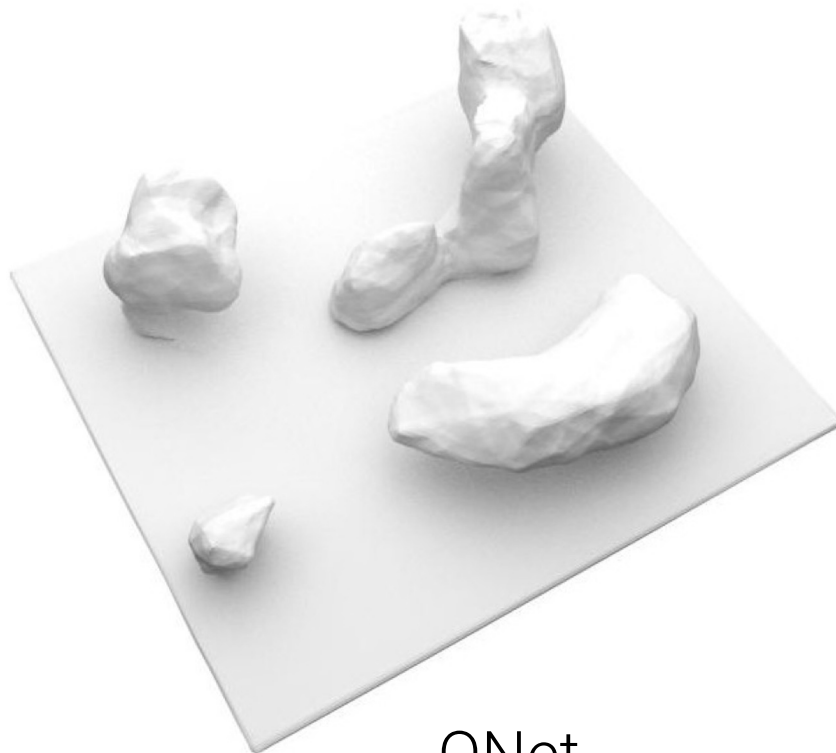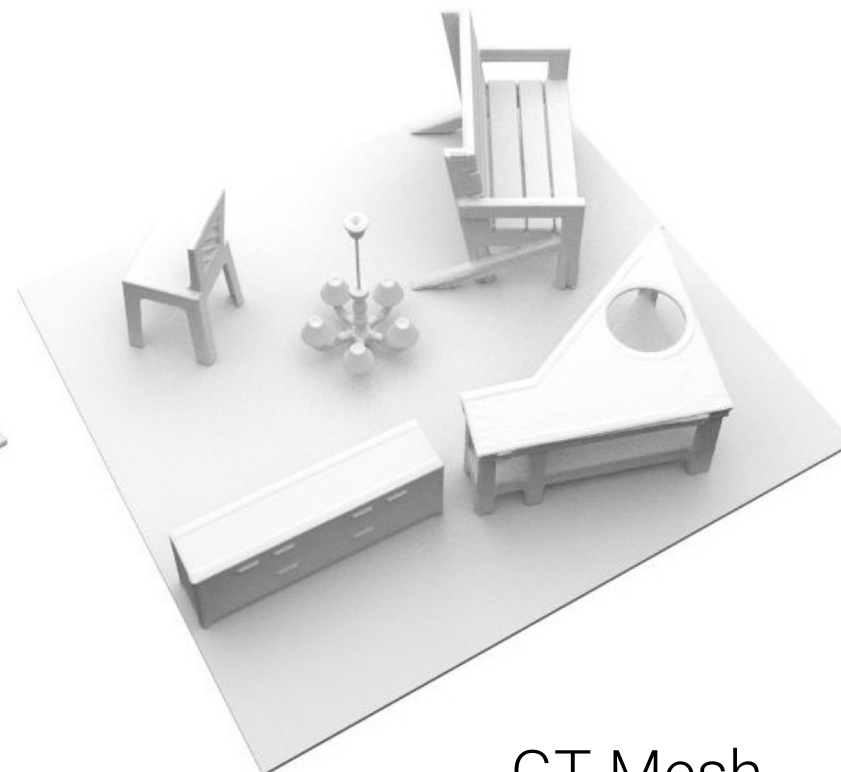## Structure of neural implicit representations:



- Global latent code ⇒ **overly smooth geometry**
- Fully-connected architecture ⇒ **no translation equivariance**

Mescheder, Oechsle, Niemeyer, Nowozin and Geiger: Occupancy Networks: Learning 3D Reconstruction in Function Space. CVPR, 2019

# Limitations

Implicit models work well for **simple objects** but poorly on **complex scenes**:



ONet

GT Mesh

Mescheder, Oechsle, Niemeyer, Nowozin and Geiger: Occupancy Networks: Learning 3D Reconstruction in Function Space. CVPR, 2019

How to reconstruct large-scale 3D scenes with **neural implicit representations**?
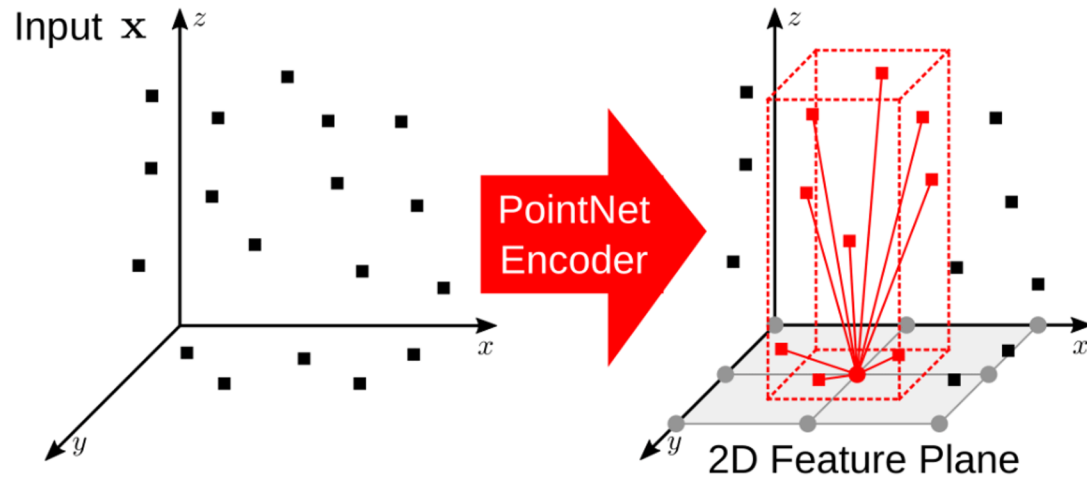
# Main Idea
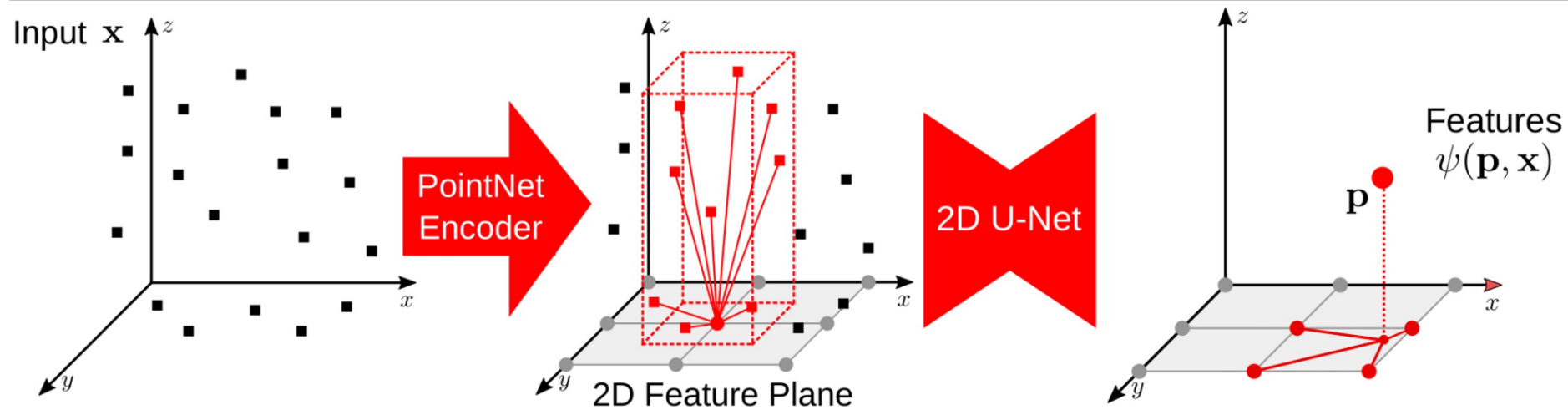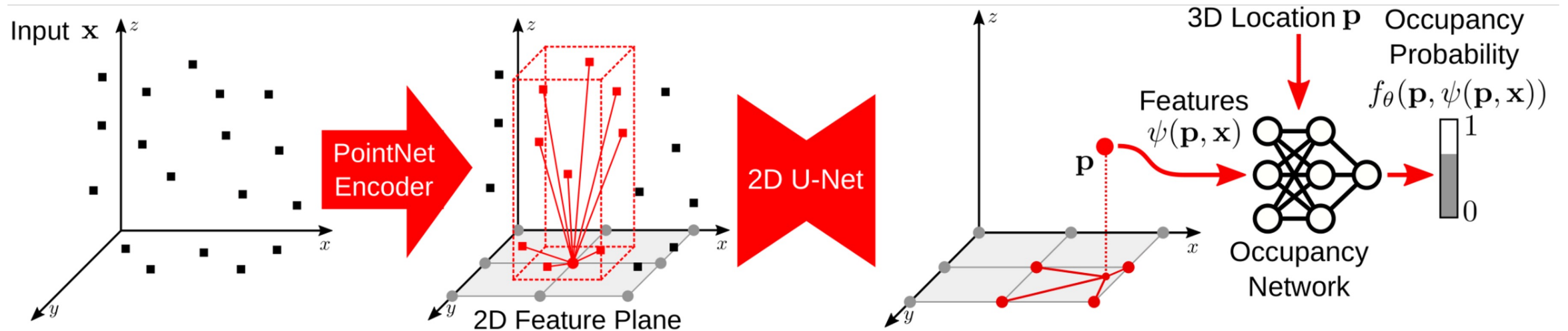


Input **x**

# Main Idea



- **2D Plane Encoder**: Use a local PointNet to process input, project onto canonical plane

# Main Idea



- **2D Plane Encoder**: Use a local PointNet to process input, project onto canonical plane
- **2D Plane Decoder**: Processed by U-Net, query features via bilinear interpolation
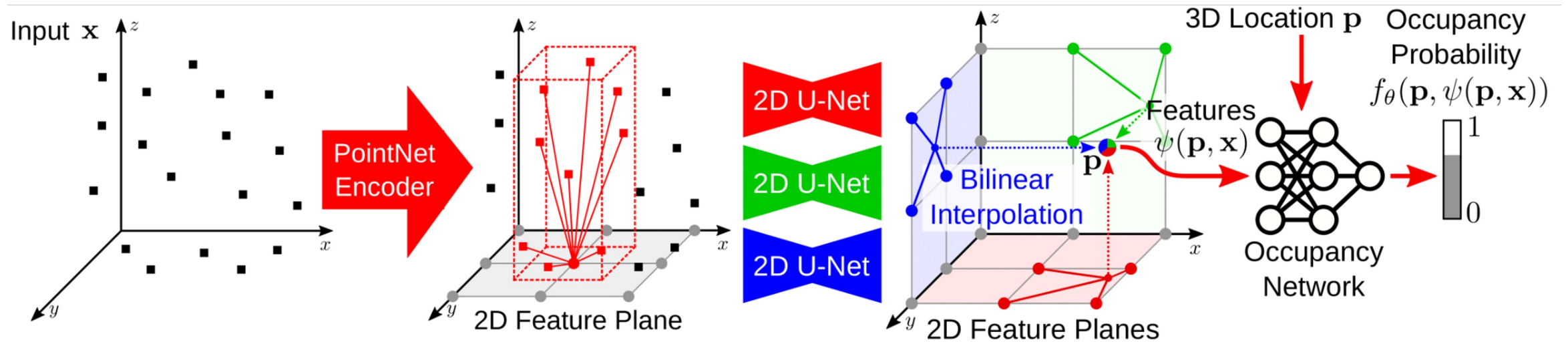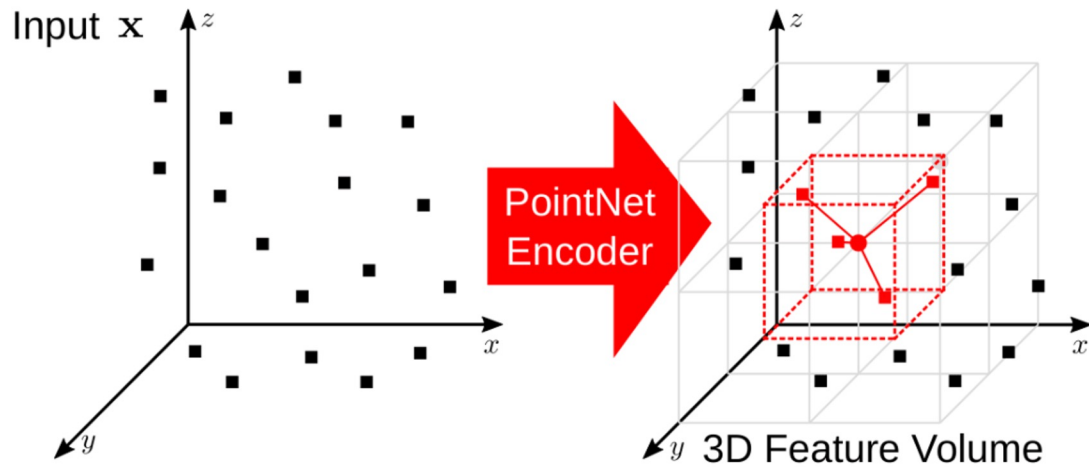
# Main Idea



- **2D Plane Encoder**: Use a local PointNet to process input, project onto canonical plane
- **2D Plane Decoder**: Processed by U-Net, query features via bilinear interpolation
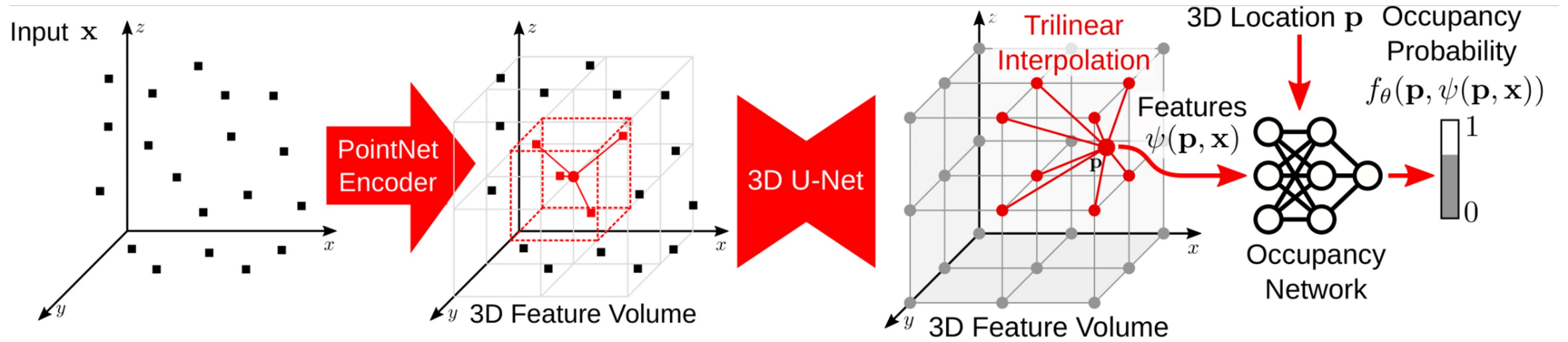- **Occupancy Readout**: Shallow occupancy network $f_\theta(\cdot)$

# Main Idea



- **2D Plane Encoder**: Use a local PointNet to process input, project onto **3-canonical planes**
- **2D Plane Decoder**: Processed by U-Net, query features via bilinear interpolation
- **Occupancy Readout**: Shallow occupancy network $f_\theta(\cdot)$

# Main Idea − 3D



- **3D Volume Encoder**: Use a local PointNet to process input, volumetric feature encoding
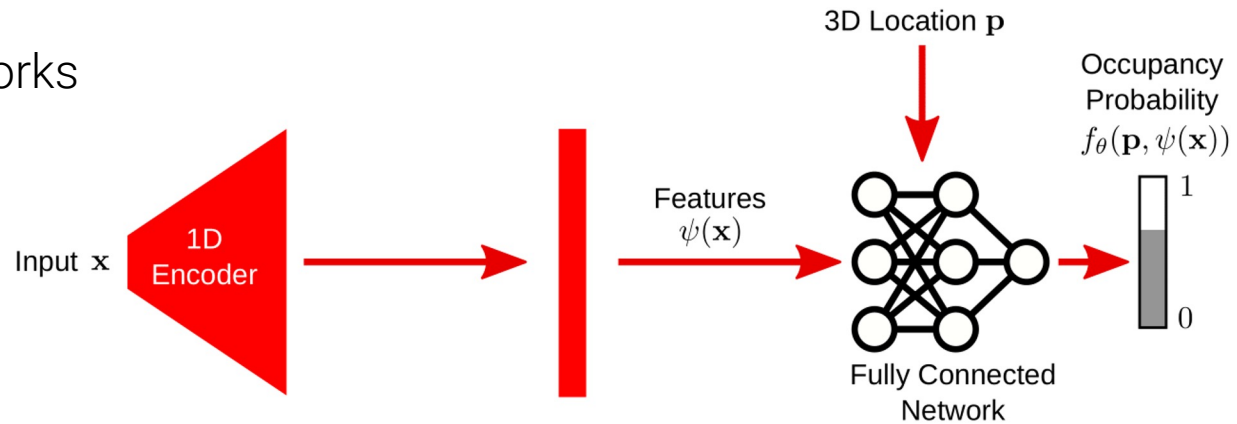
# Main Idea − 3D



- **3D Volume Encoder**: Use a local PointNet to process input, volumetric feature encoding
- **3D Volume Decoder**: Processed by 3D U-Net, query features via trilinear interpolation
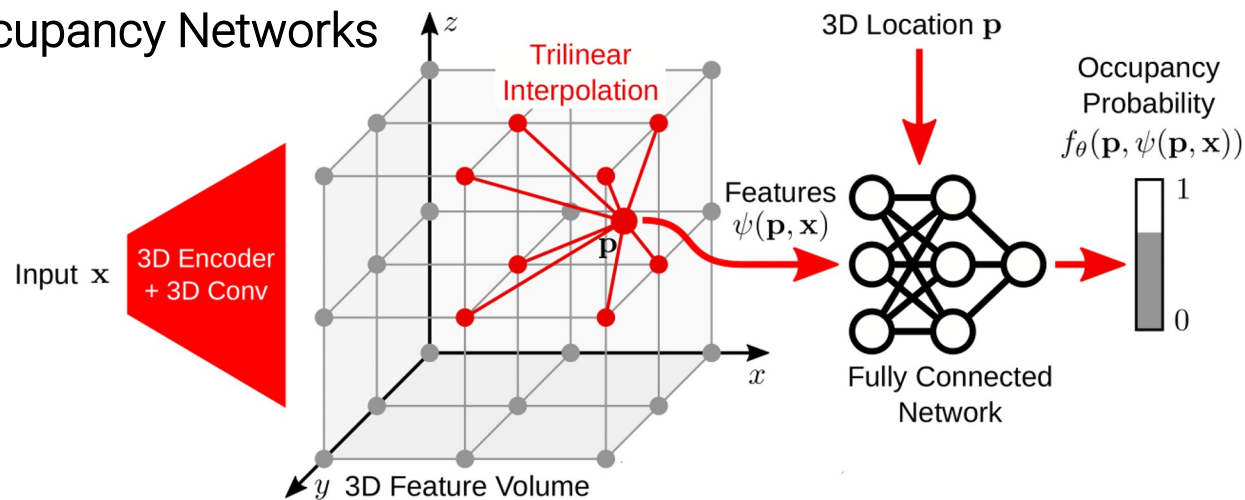- **Occupancy Readout**: Shallow occupancy network $f_\theta(\cdot)$

# Comparison

Occupancy Networks



- global feature
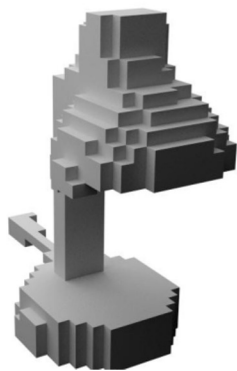- heavy FC network
- no translation equivariance

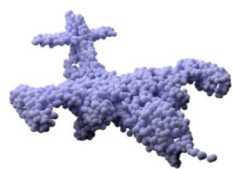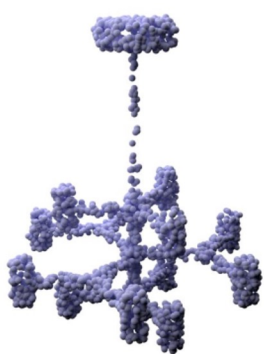Convolutional Occupancy Networks



- local feature
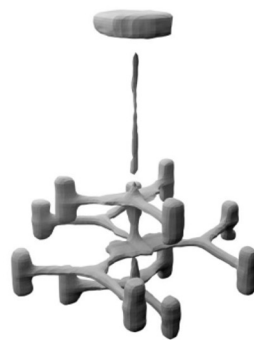- shallow FC network
- translation equivariance

# Results

# Object-Level Reconstruction



Input      ONet      **Ours - 2D**      **Ours - 3D**      GT Mesh

# Training Speed

# Training Speed

# Scene-Level Reconstruction: Synthetic

- Trained and evaluated on synthetic rooms



Input

GT Mesh

# Scene-Level Reconstruction: Synthetic

- ONet fails on room-level reconstruction



Input

ONet

# Scene-Level Reconstruction: Synthetic

- SPSR requires surface normals, output is <span style="color:red">noisy</span>



Input



SPSR

(Screened Poisson Surface Reconstruction)

# Scene-Level Reconstruction: Synthetic

- Our method preserves better details



Input

**Ours**

# Large-Scale Reconstruction

**Scene size**: 15.7m x 12.3m x 4.5m

## Results on Matterport3D

- Fully convolutional model

- Trained on synthetic crops

- Sliding-window evaluation

- Scale to any scene size



**Our reconstruction output**

# Large-Scale Reconstruction

**Scene size**: 15.7m x 12.3m x 4.5m

## Results on Matterport3D

- Fully convolutional model
- Trained on synthetic crops
- Sliding-window evaluation
- Scale to any scene size



**Our reconstruction output**

# Take-home Messages

- Introduce 3 different expressive hybrid representations for neural fields

- CNN's translation equivariance enables to reconstruct large scenes

- The "**tri-plane**" representation became VERY popular

  - Especially in the **NeRF era**, see e.g. EG3D [CVPR'21], TensoRF [ECCV'22]

**Limitations**

- Not rotational equivariance

# NeRF is awesome!



**Some existing problems…**
😢 Poor underlying geometry
😢 Camera poses needed

Mildenhall*, Srinivasan*, Tancik* et al: NeRF : Representing Scenes as Neural Radiance Fields for View Synthesis. ECCV 2020

# RGB-D Sequences





40x 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

# Pipeline

# 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

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
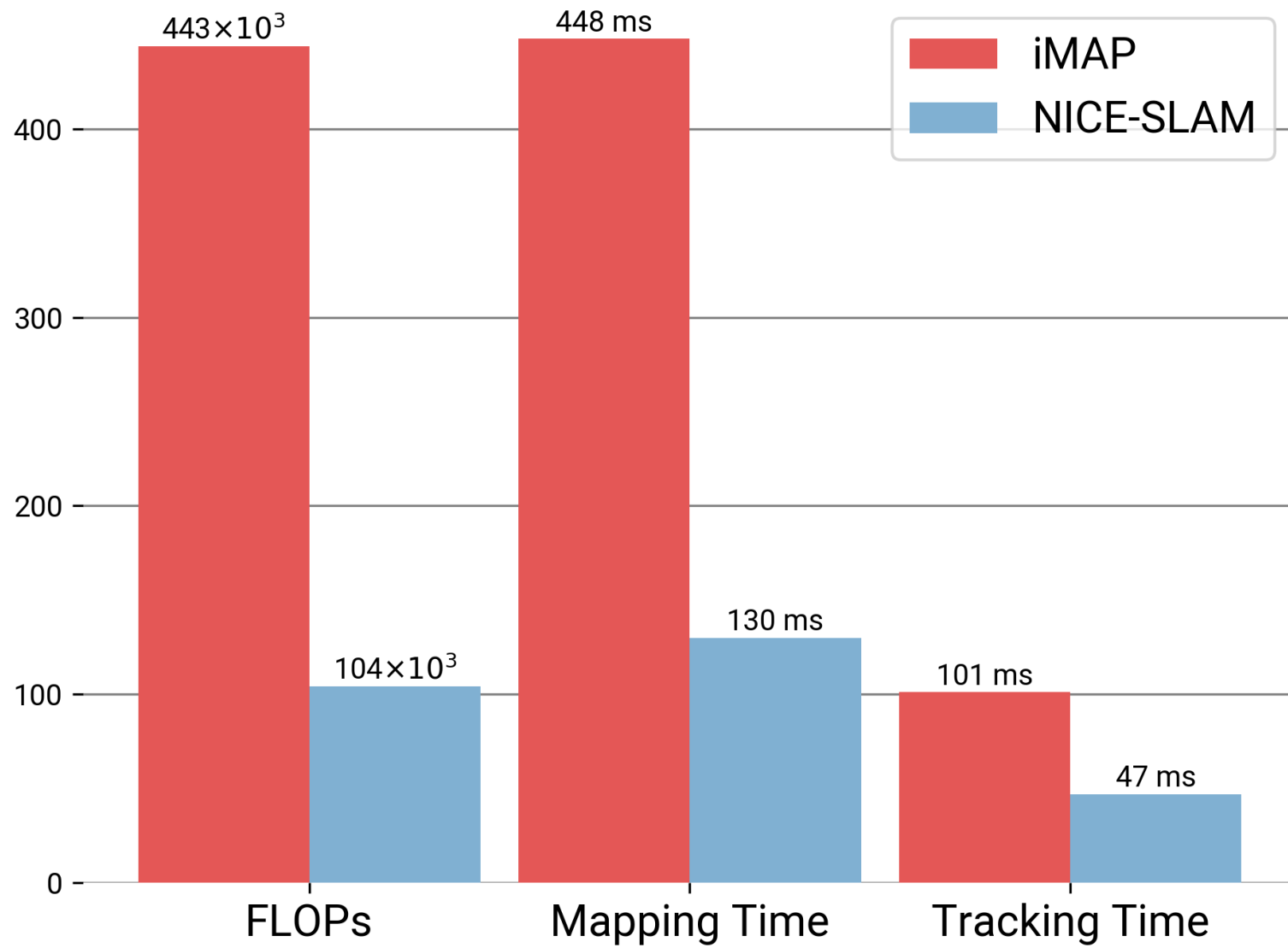
[6]University of Tübingen, Tübingen AI Center       [7]Microsoft

| 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



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

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

OpenSeg [1]
LSeg [2]

**3D Geometry**

**Per-pixel Features**
(visualize with PCA)

**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

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

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)

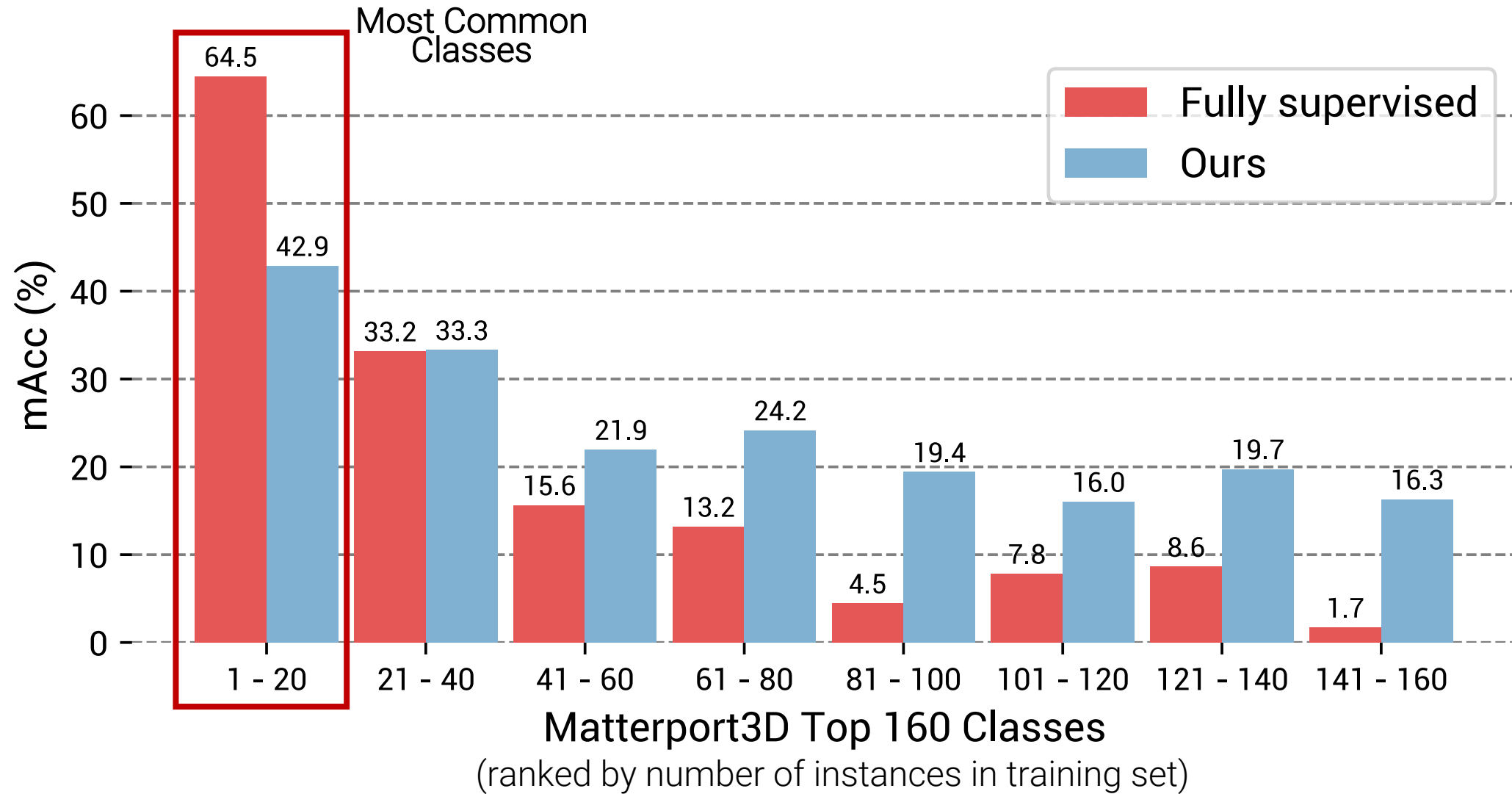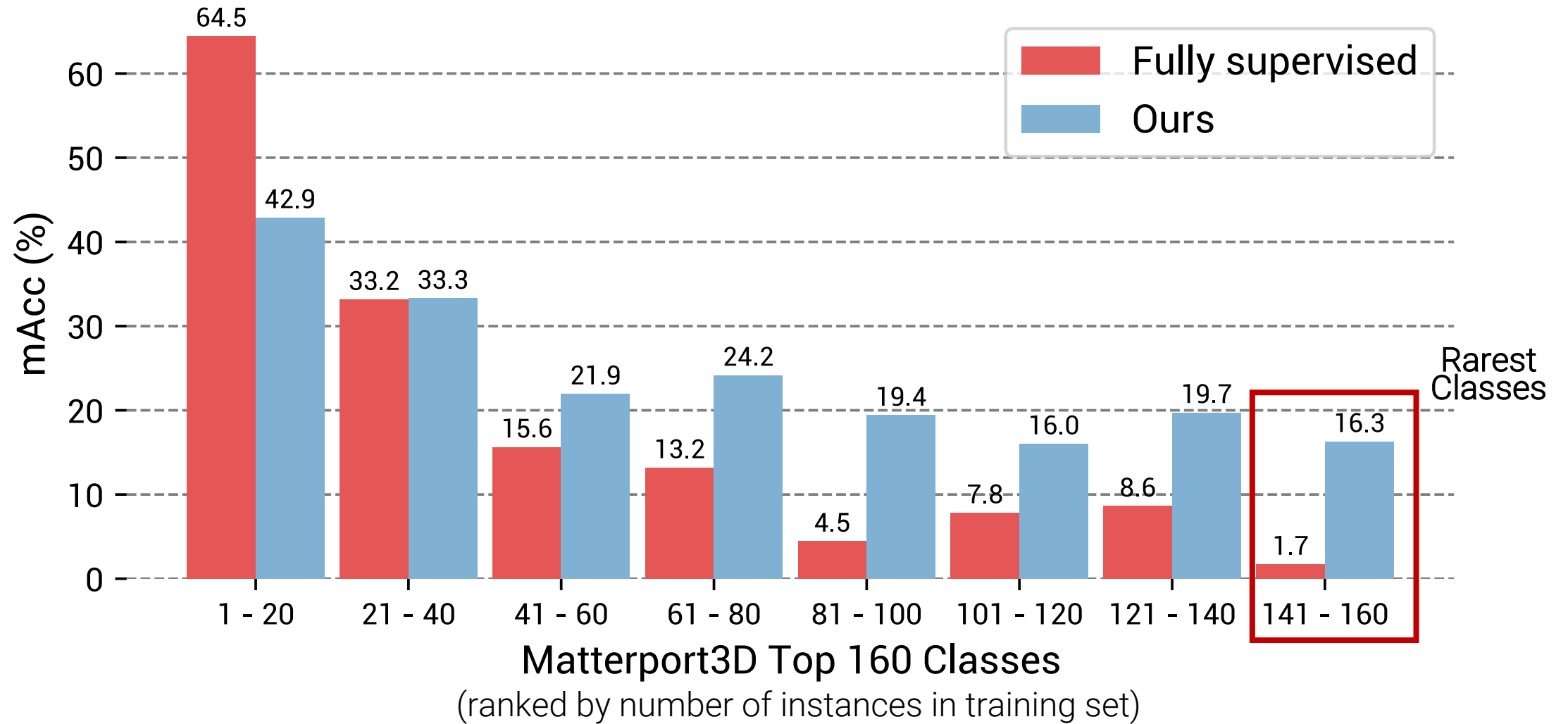| | wall | | floor | | cabinet | | bed | | chair | | sofa | | table | | door | | window | | bookshelf | | picture | | counter | | desk | | curtain | | refrigerator | | shower curtain | | toilet | | sink | | bathtub | | other |

## 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 | rug | drawer | stove | tv stand | air conditioner | thermostat | ladder | candlestick | decorative plate | lamp shade | foot rest |
| ceiling | table | toilet | box | book | air vent | washing machine | shoe | fire extinguisher | radiator | garage door | light | pool table | car | soap dish |
| floor | plant | coffee table | ottoman | bottle | light switch | shower curtain | heater | curtain rod | kitchen island | piano | scale | bag | bottle of soap | toilet brush | cleaner |
| picture | mirror | counter | photo | refridgerator | purse | bin | headboard | printer | paper towel | board | bag | display case | water cooler | drum | computer |
| window | towel | stairs | bench | toilet paper | bookshelf | door way | 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 |

63

# Comparison

# Comparison



Matterport3D Top 160 Classes
(ranked by number of instances in training set)
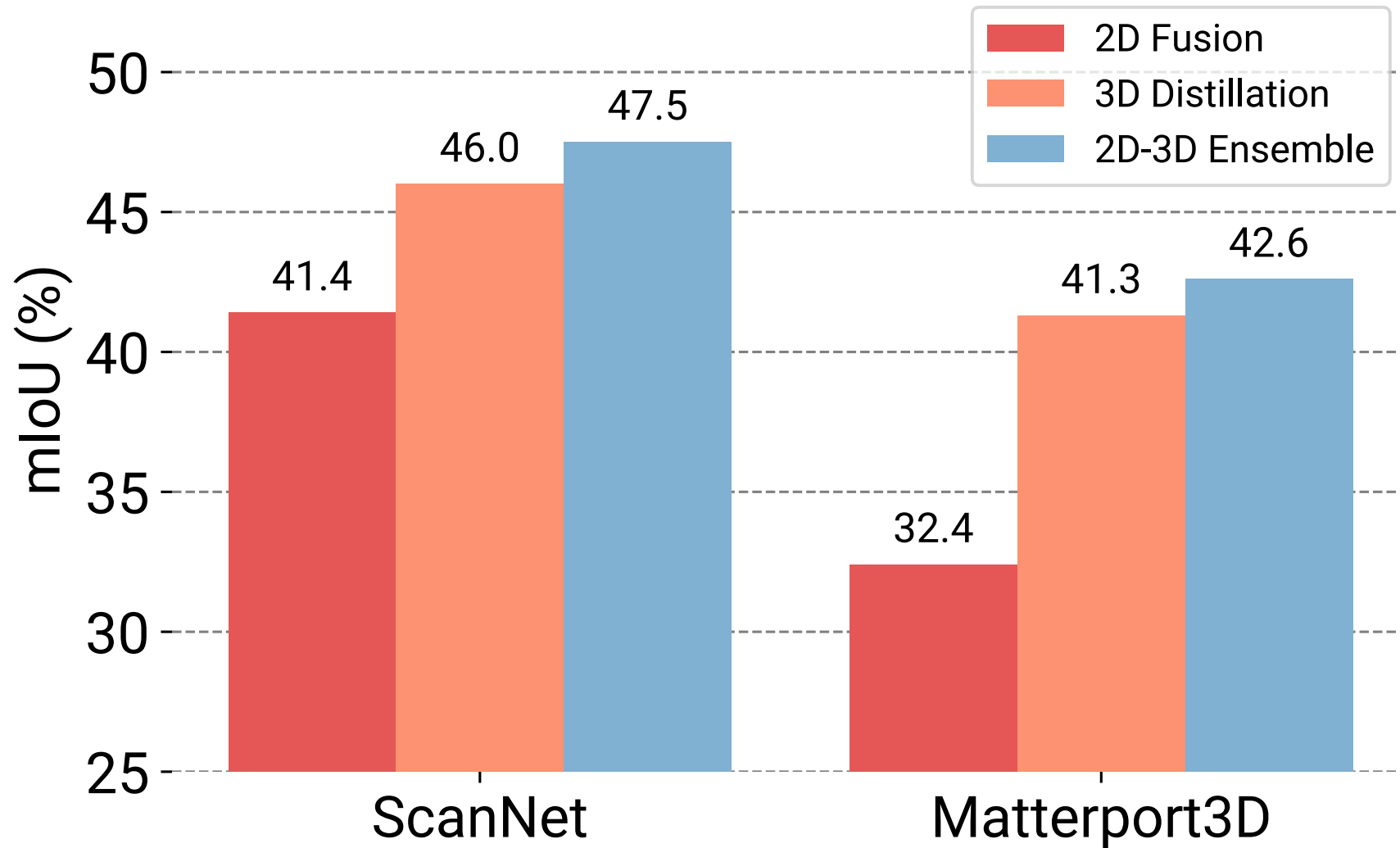
# Ablation

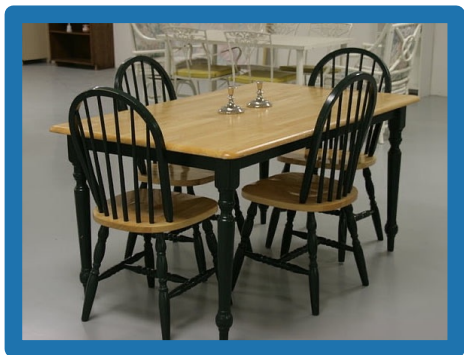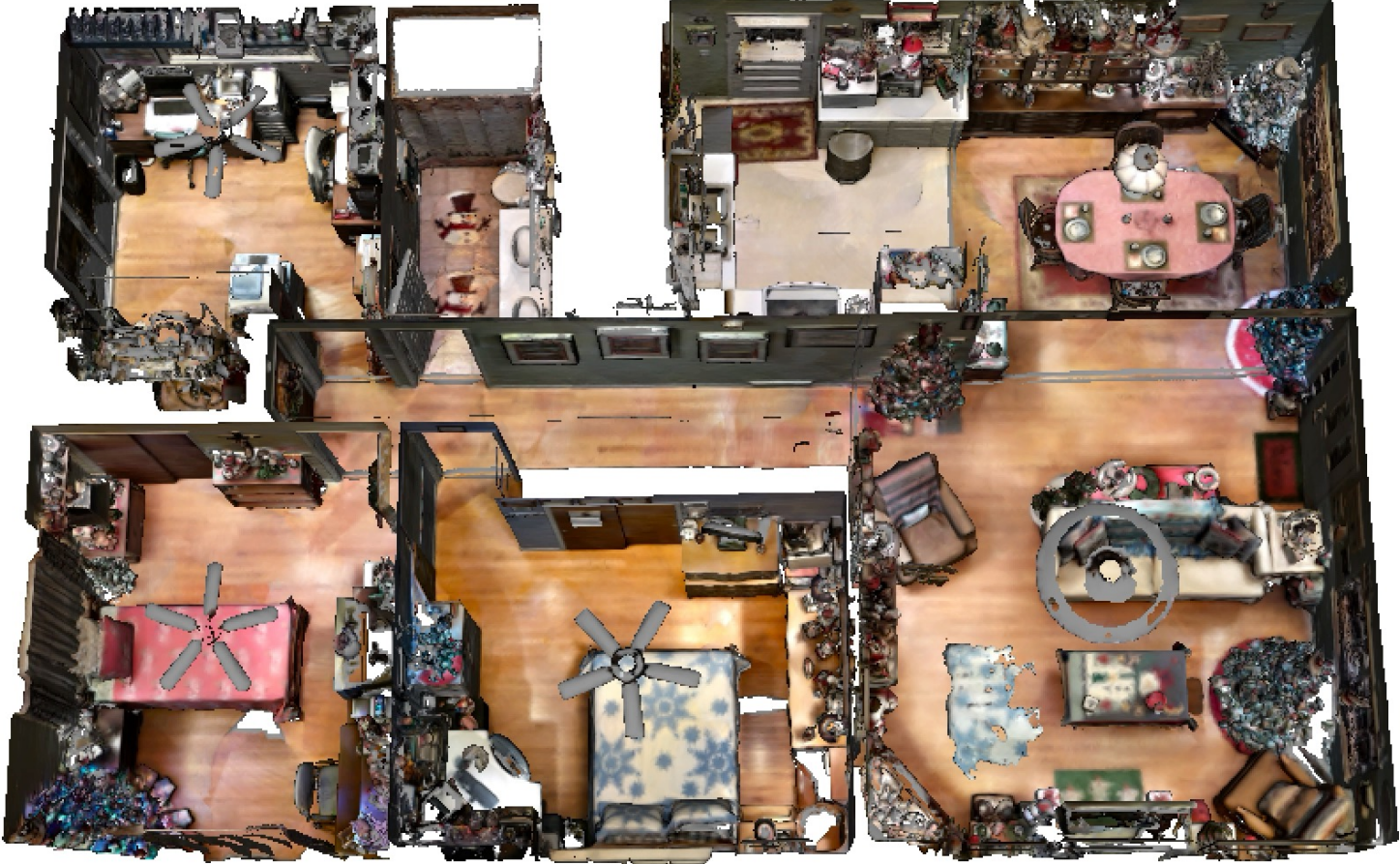# Image-based 3D Scene Query

Image Queries

Given 3D Geometry

# **Interactive Demo**

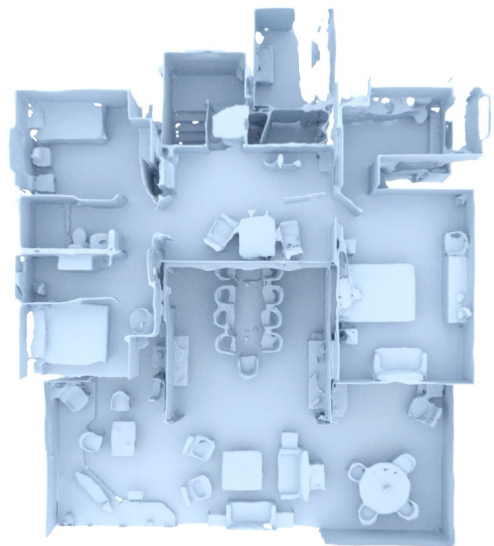Open-vocabulary 3D Scene Exploration

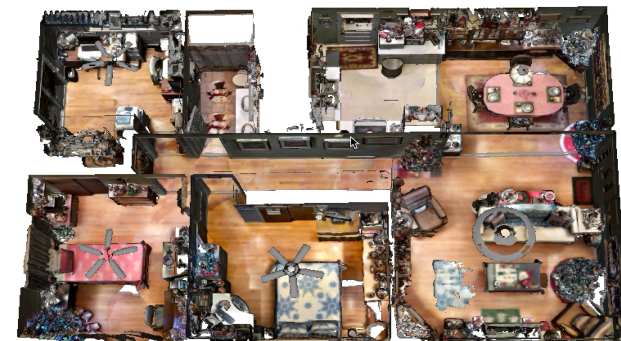Text queries:

# 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

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

# Learn to Reconstruct and Understand the 3D World

Songyou Peng



**Convolutional Occupancy Networks**
ECCV 2020 (Spotlight)
pengsongyou.github.io/conv_onet

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

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

# Thank you!