





Intelligent systems interact with **3D environments**

## **3D Reconstruction**

Create digital twins from real scenes

## **3D Scene Understanding**

Analyze the scene digitally

# Key Challenges

**Reconstruct** and **Understand** 3D Environments

- Reconstruct 3D scenes **at scale**
- Reconstruct 3D scenes **at speed**
- Reconstruct purely **from 2D observations**

# Key Challenges

Reconstruct and **Understand** 3D Environments

- Reconstruct 3D scenes **at scale**
  - Reconstruct 3D scenes **at speed**
  - Reconstruct purely **from 2D observations**
- 
- Understand **arbitrary concepts** in a 3D scene
  - Learn to understand **without labeled 3D data**



# Research Overview of My PhD

Learn to **Reconstruct** and **Understand** 3D Environments



ConvOccNet

ECCV 2020 (Spotlight)



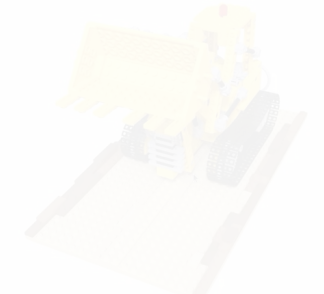
MonoSDF

NeurIPS 2022



Shape As Points

NeurIPS 2021 (Oral)



runs now at 50 fps on a GTX 1080 Ti

KiloNeRF

ICCV 2021



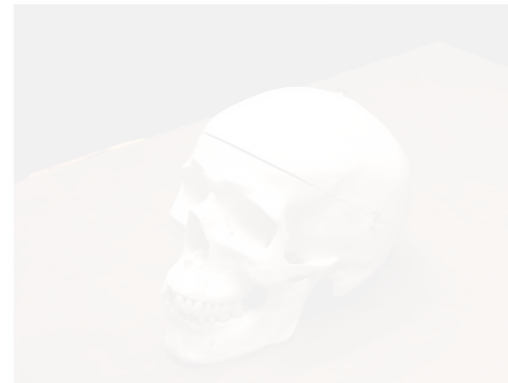
NICE-SLAM

CVPR 2022



NICER-SLAM

3DV 2024 (Oral)



UNISURF

ICCV 2021 (Oral)

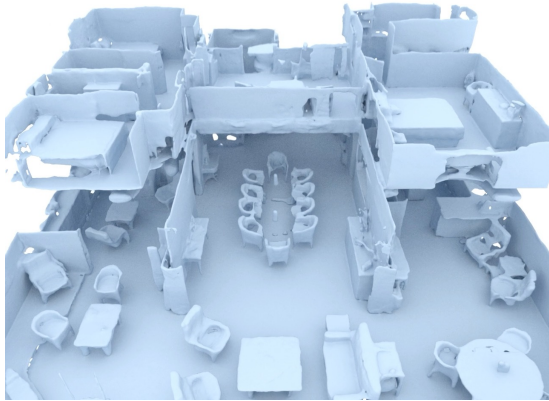


OpenScene

CVPR 2023 <sup>4</sup>

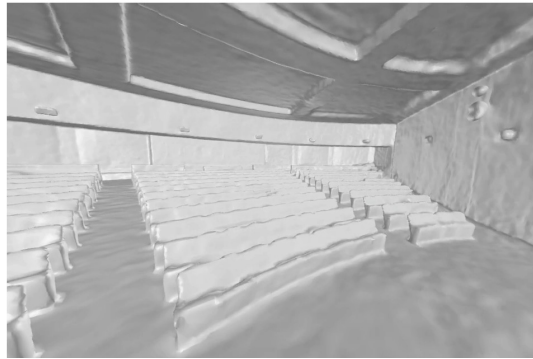
# Research Overview of My PhD

Learn to Reconstruct and Understand 3D Environments



**ConvOccNet**

ECCV 2020 (Spotlight)

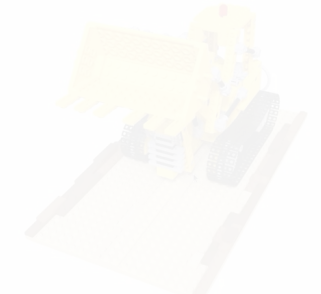


**MonoSDF**

NeurIPS 2022



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

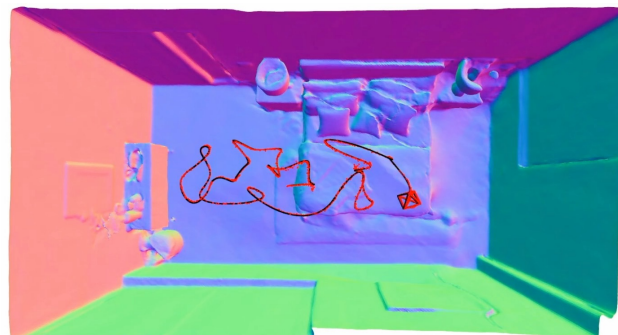


runs now at 50 fps on a GTX 1080 Ti

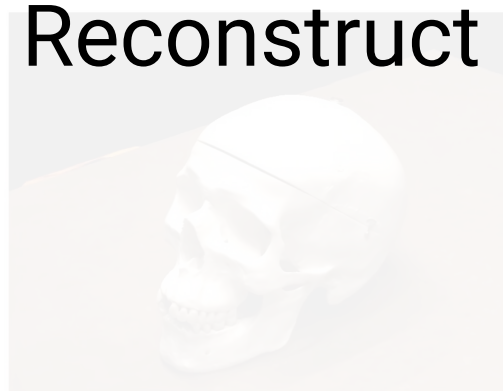
**KiloNeRF**  
ICCV 2021



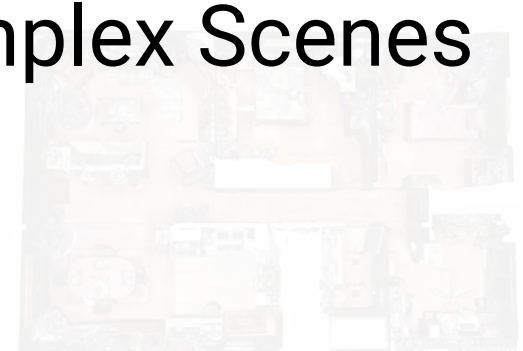
**NICE-SLAM**  
CVPR 2022



**NICER-SLAM**  
3DV 2024 (Oral)



**UNISURF**  
ICCV 2021 (Oral)



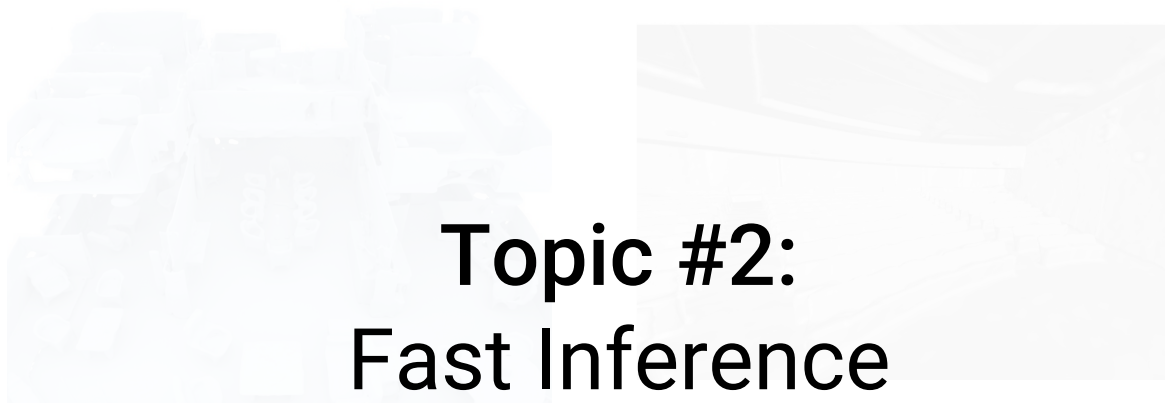
**OpenScene**  
CVPR 2023

**Topic #1:**  
**Reconstruct Complex Scenes**

# Research Overview of My PhD

Learn to Reconstruct and Understand 3D Environments

## Topic #2: Fast Inference

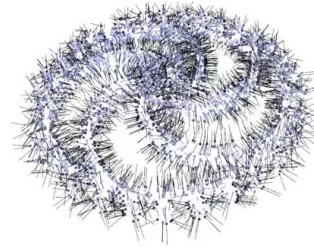


ConvOccNet

ECCV 2020 (Spotlight)

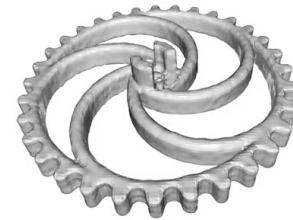
MonoSDF

NeurIPS 2022



Shape As Points

NeurIPS 2021 (Oral)



runs now at 50 fps on a GTX 1080 Ti

KiloNeRF

ICCV 2021



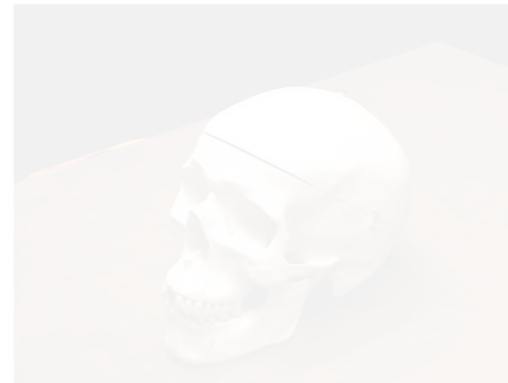
NICE-SLAM

CVPR 2022



NICER-SLAM

3DV 2024 (Oral)



UNISURF

ICCV 2021 (Oral)



OpenScene

CVPR 2023 <sub>6</sub>



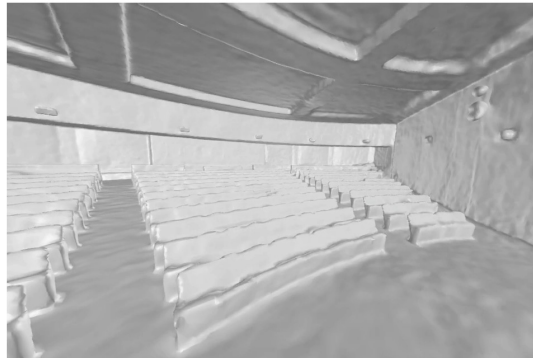
# Research Overview of My PhD

Learn to Reconstruct and Understand 3D Environments



**ConvOccNet**

ECCV 2020 (Spotlight)



**MonoSDF**

NeurIPS 2022

**Topic #3:**  
**Reconstruct from 2D Observations**

runs now at 50 fps on a GTX 1080 Ti

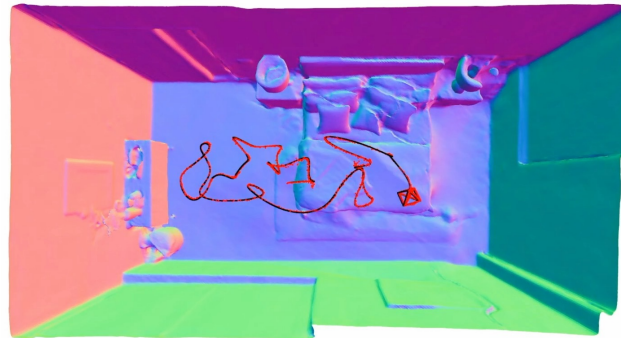
**KiloNeRF**

ICCV 2021



**NICE-SLAM**

CVPR 2022



**NICER-SLAM**

3DV 2024 (Oral)



**UNISURF**

ICCV 2021 (Oral)



**OpenScene**

CVPR 2023



# Research Overview of My PhD

Learn to Reconstruct and Understand 3D Environments



ConvOccNet

ECCV 2020 (Spotlight)



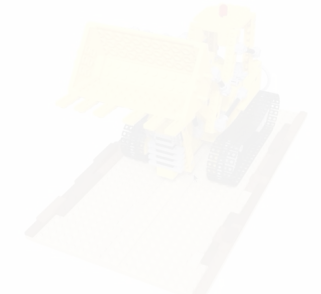
MonoSDF

NeurIPS 2022



Shape As Points

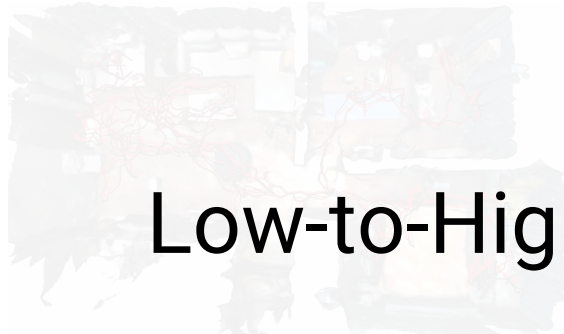
NeurIPS 2021 (Oral)



runs now at 50 fps on a GTX 1080 Ti

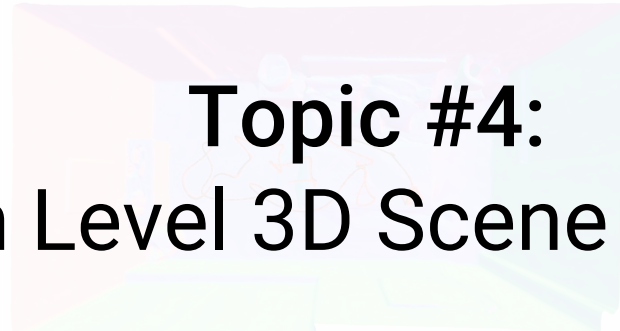
KiloNeRF

ICCV 2021



NICE-SLAM

CVPR 2022

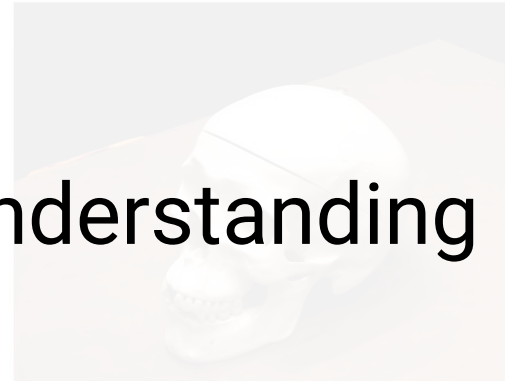


**Topic #4:**

**Low-to-High Level 3D Scene Understanding**

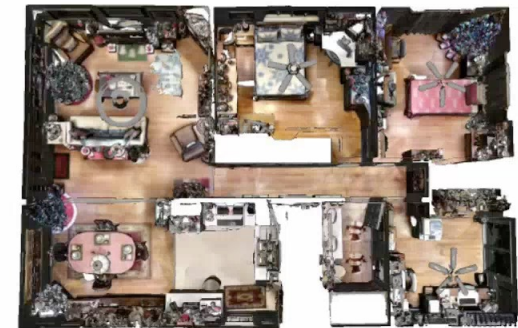
NICER-SLAM

3DV 2024 (Oral)



UNISURF

ICCV 2021 (Oral)

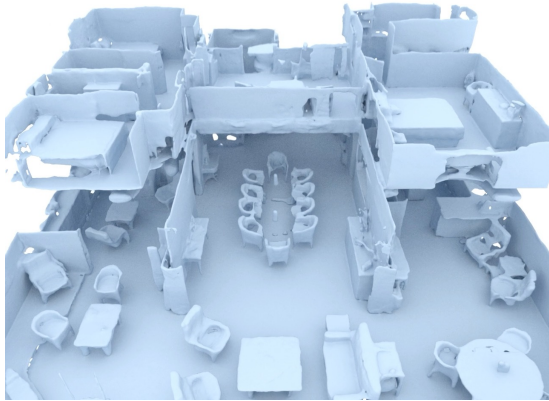


OpenScene

CVPR 2023 8

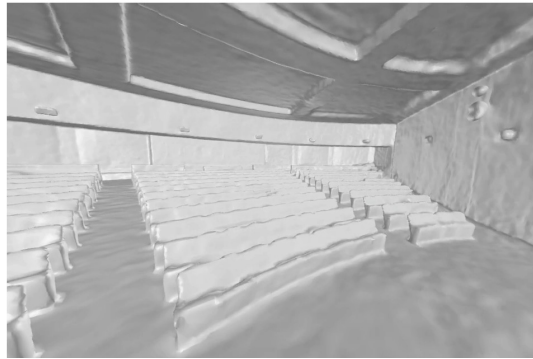
# Research Overview of My PhD

Learn to Reconstruct and Understand 3D Environments



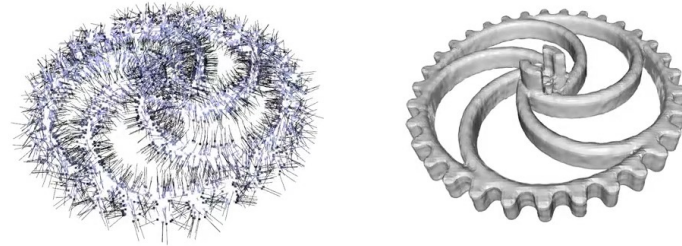
**ConvOccNet**

ECCV 2020 (Spotlight)



**MonoSDF**

NeurIPS 2022



**Shape As Points**

NeurIPS 2021 (Oral)



runs now at 50 fps on a GTX 1080 Ti

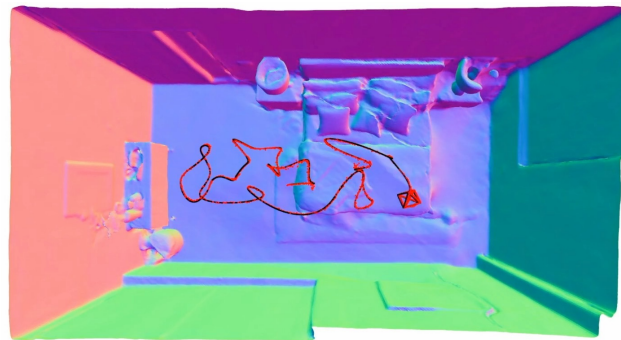
**KiloNeRF**

ICCV 2021



**NICE-SLAM**

CVPR 2022



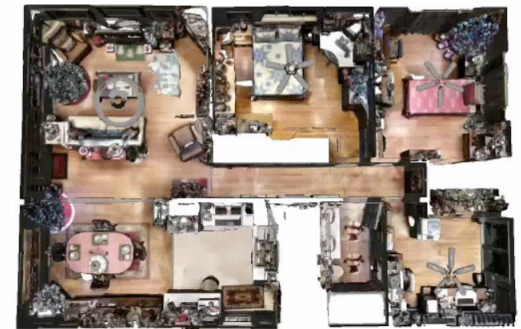
**NICER-SLAM**

3DV 2024 (Oral)



**UNISURF**

ICCV 2021 (Oral)



**OpenScene**

CVPR 2023 9

# This Thesis

## Develop 3D Neural Scene Representations

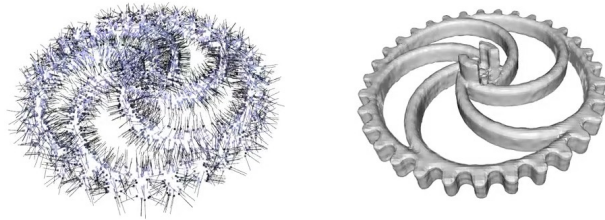
for **3D Reconstruction** and **3D Scene Understanding**

### 1. Complex Scenes



**ConvOccNet**  
ECCV 2020 (Spotlight)

### 2. Fast Inference



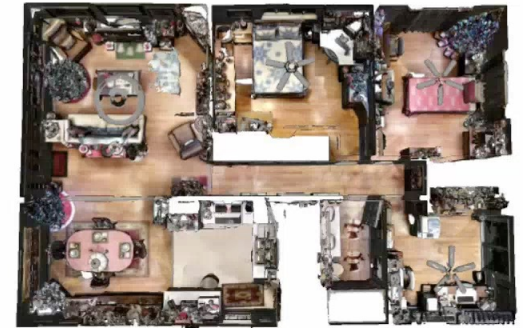
**Shape As Points**  
NeurIPS 2021 (Oral)

### 3. From 2D Observations



**NICE-SLAM**  
CVPR 2022

### 4. Arbitrary Queries



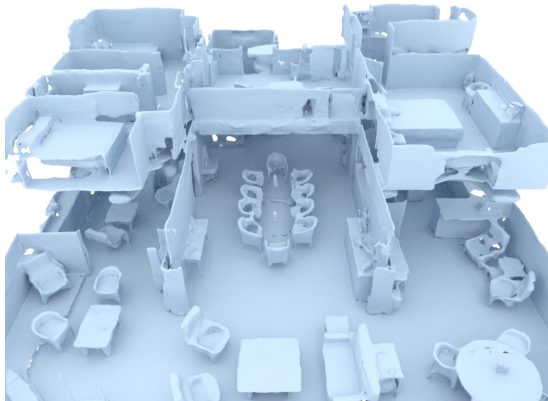
**OpenScene**  
CVPR 2023



# This Thesis

Develop 3D Neural Scene Representations  
for **3D Reconstruction** and **3D Scene Understanding**

## 1. Complex Scenes



**ConvOccNet**  
ECCV 2020 (Spotlight)

## 2. Fast Inference



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

## 3. From 2D Observations



**NICE-SLAM**  
CVPR 2022

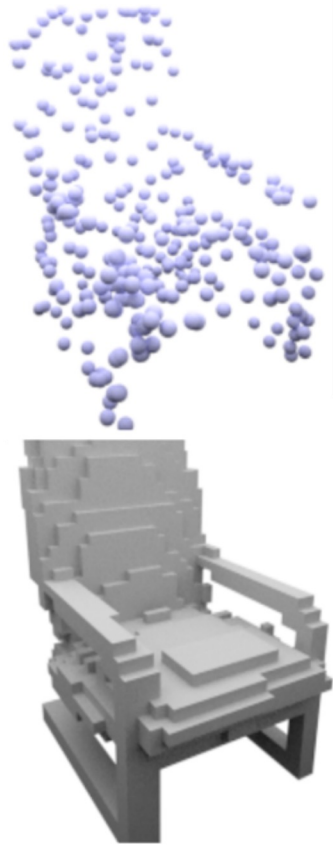
## 4. Arbitrary Queries



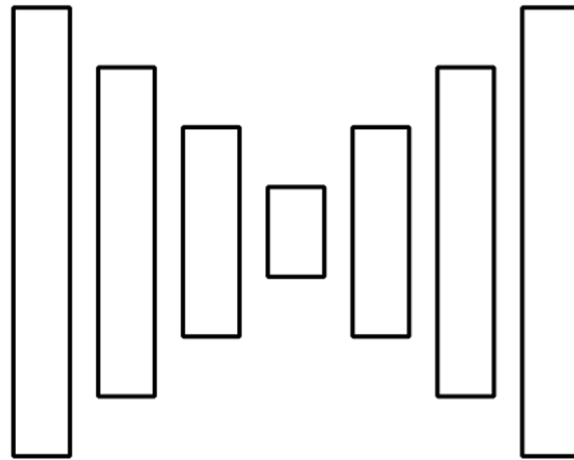
**OpenScene**  
CVPR 2023



# Learning-based 3D Reconstruction



Input



Neural Network

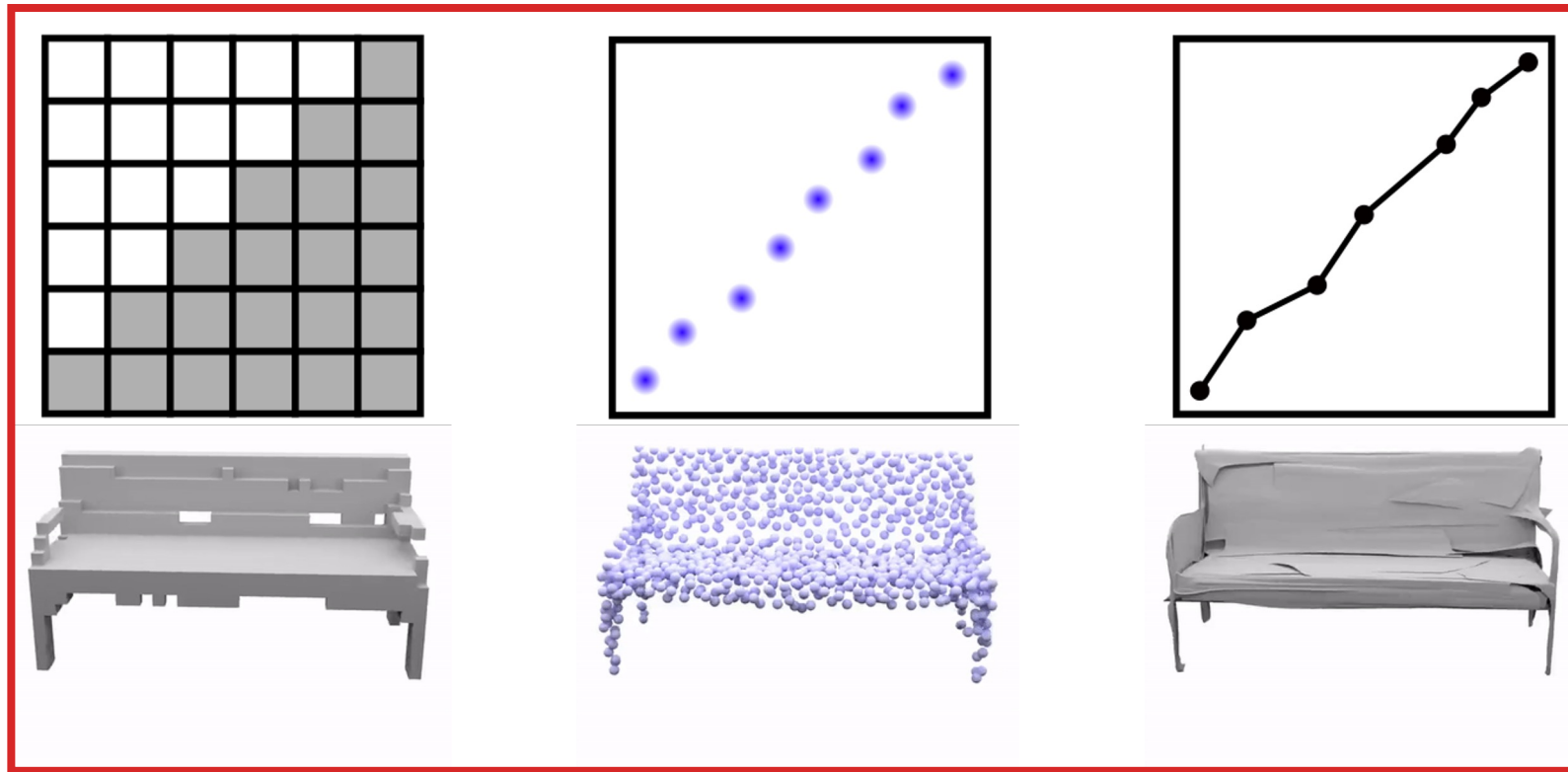


3D Reconstruction

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

# 3D Representations

## Traditional Explicit Representations



## Discretization

# 3 Seminal Papers at the Same CVPR!

## Neural Implicit Representations

### Occupancy Networks: Learning 3D Reconstruction in Function Space

Lars Mescheder<sup>1</sup> Michael Oechsle<sup>1,2</sup> Michael Niemeyer<sup>1</sup> Sebastian Nowozin<sup>3†</sup> Andreas Geiger<sup>1</sup>

<sup>1</sup>Autonomous Vision Group, MPI for Intelligent Systems and University of Tübingen

<sup>2</sup>ETAS GmbH, Stuttgart

<sup>3</sup>Google AI Berlin

### DeepSDF: Learning Continuous Signed Distance Functions for Shape Representation

Jeong Joon Park<sup>1,3†</sup> Peter Florence<sup>2,3†</sup> Julian Straub<sup>3</sup> Richard Newcombe<sup>3</sup> Steven Lovegrove<sup>3</sup>

<sup>1</sup>University of Washington

<sup>2</sup>Massachusetts Institute of Technology

<sup>3</sup>Facebook Reality Labs

### Learning Implicit Fields for Generative Shape Modeling

Zhiqin Chen

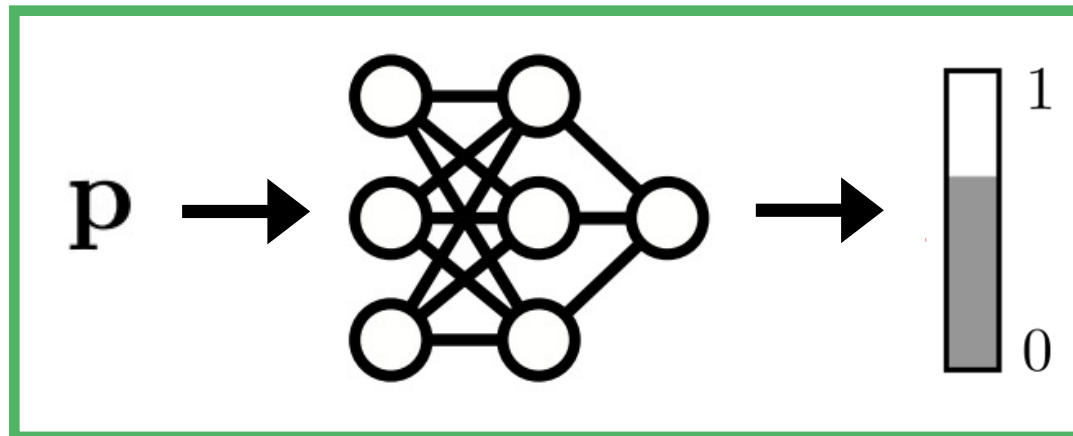
Simon Fraser University

zhiqinc@sfu.ca

Hao Zhang

Simon Fraser University

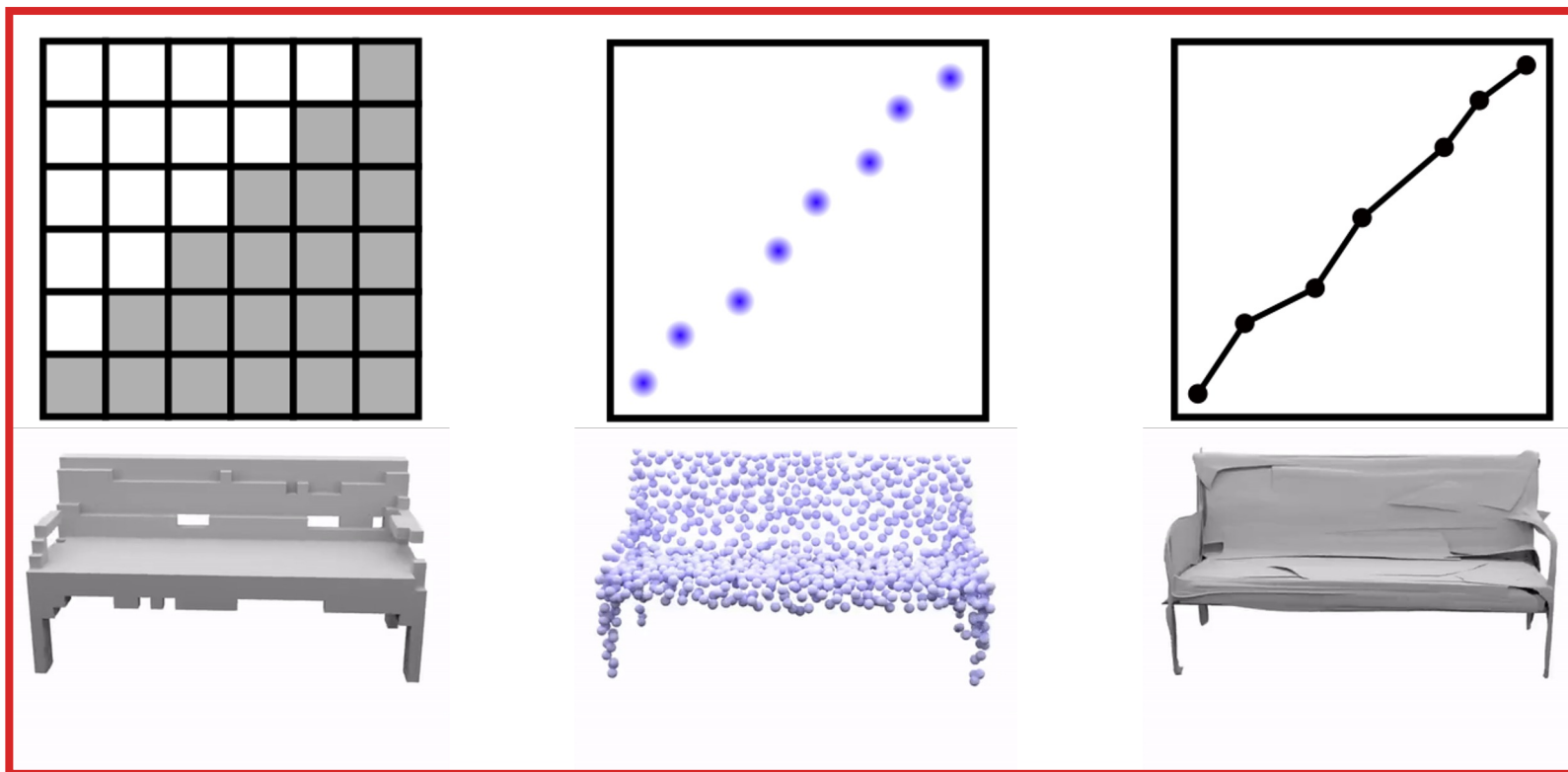
haoz@sfu.ca



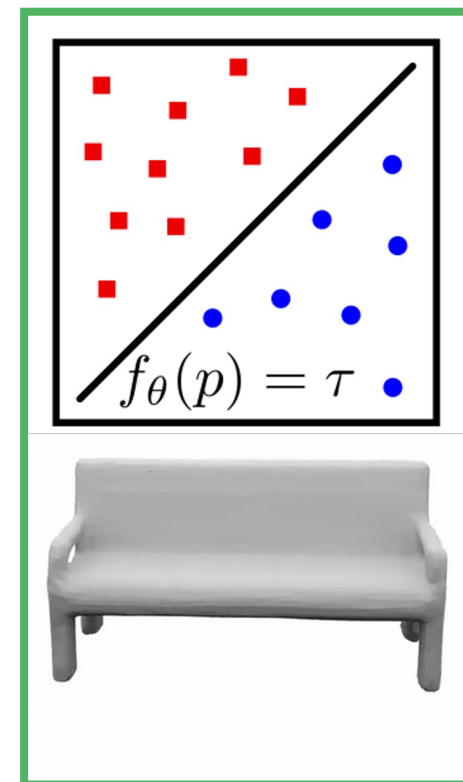


# 3D Representations

## Neural Implicit Representations



**Discretization**

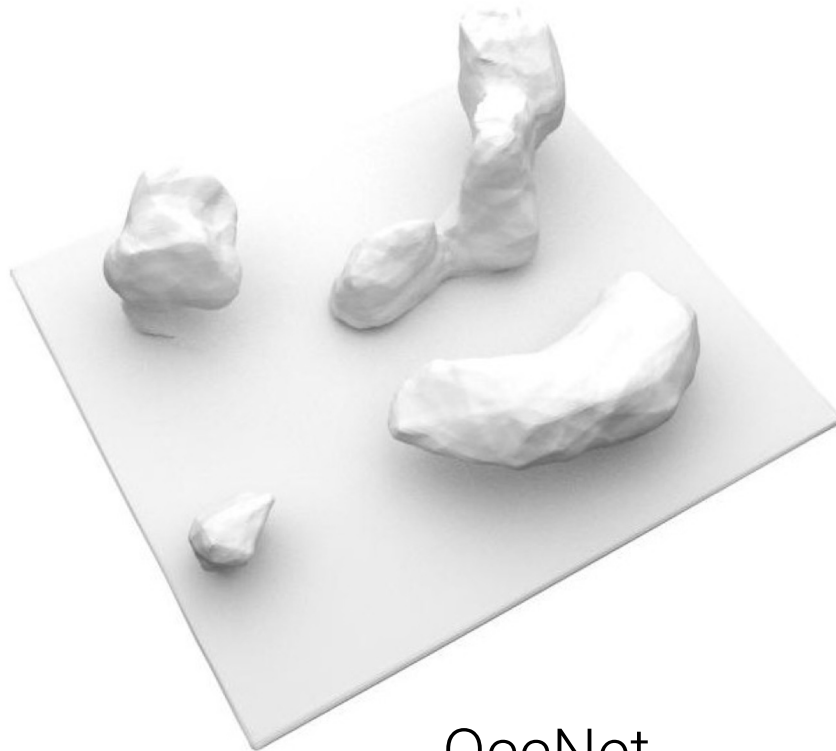


**Continuous**

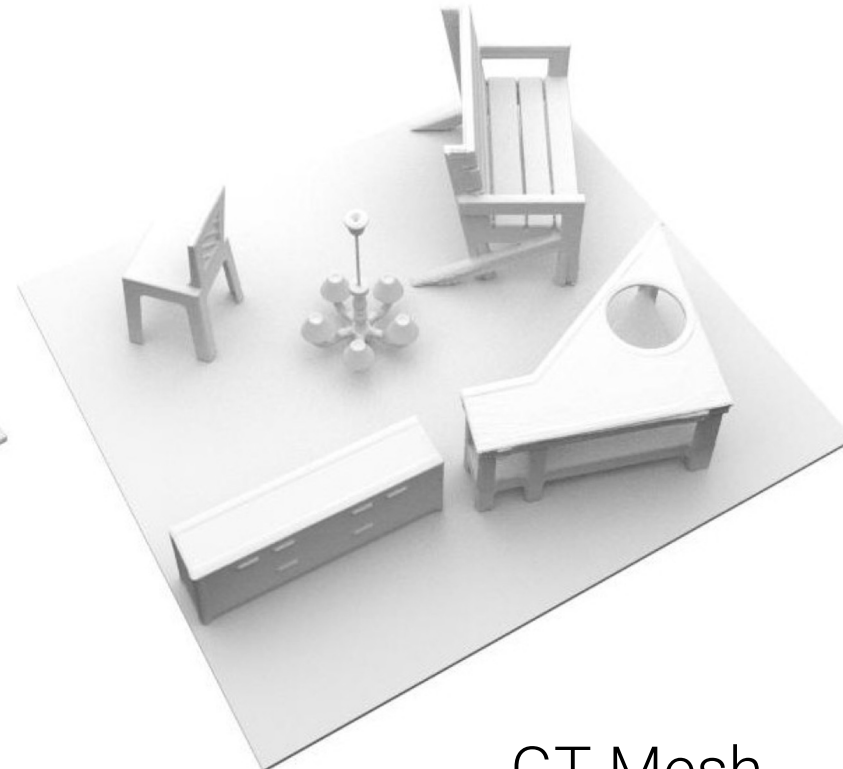
# Limitations

## Neural Implicit Representations

Works well for **simple objects**, but poorly on **complex scenes**



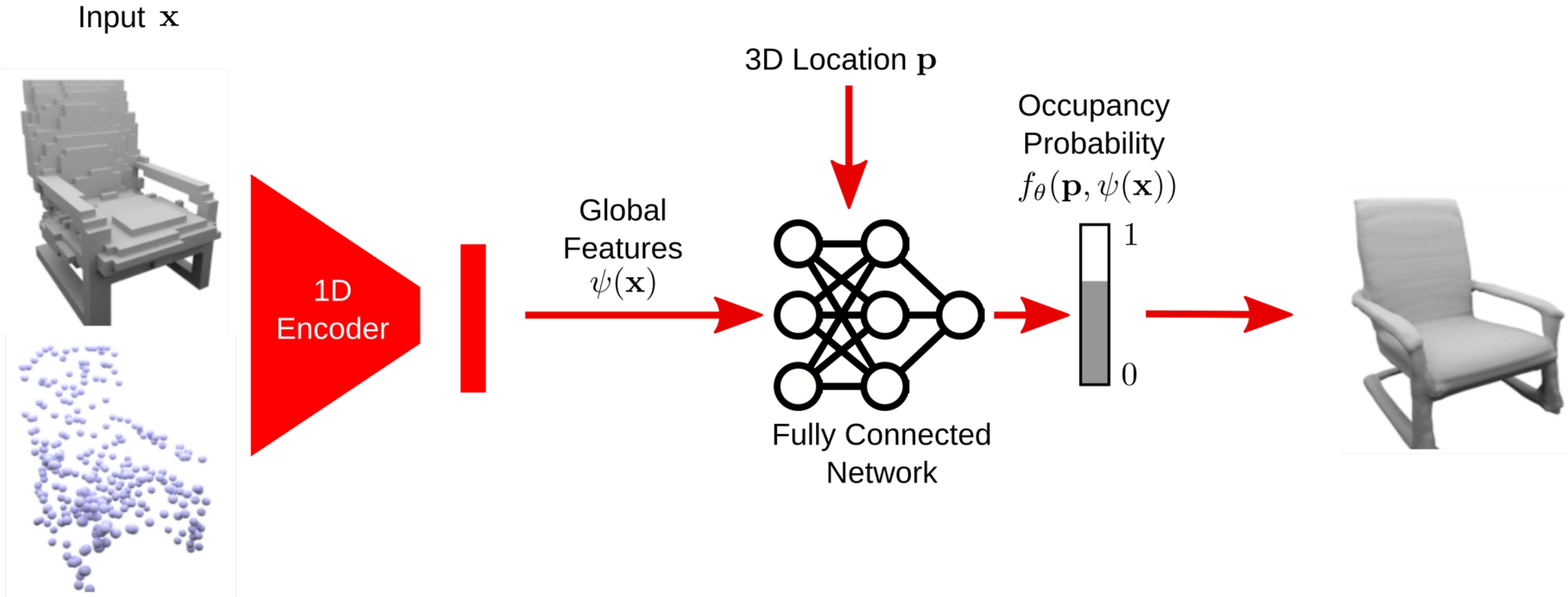
OccNet



GT Mesh

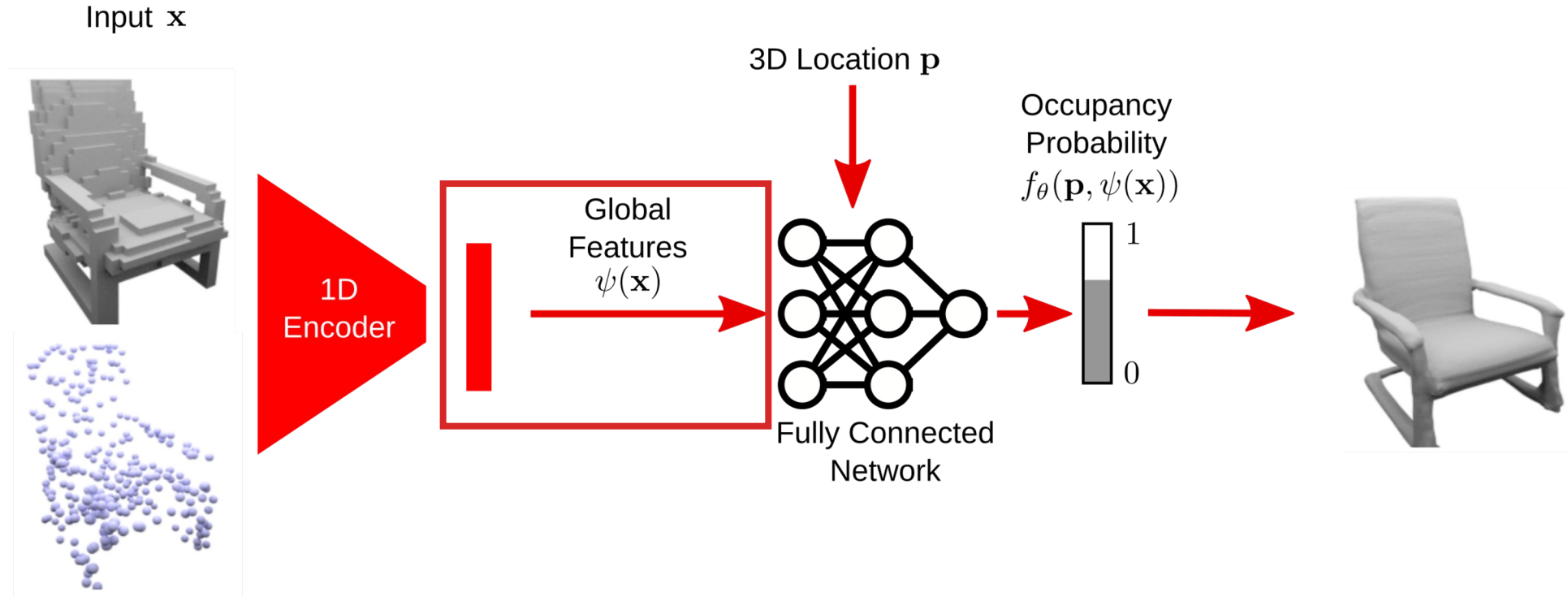
# Limitations

## Neural Implicit Representations



# Limitations

## Neural Implicit Representations

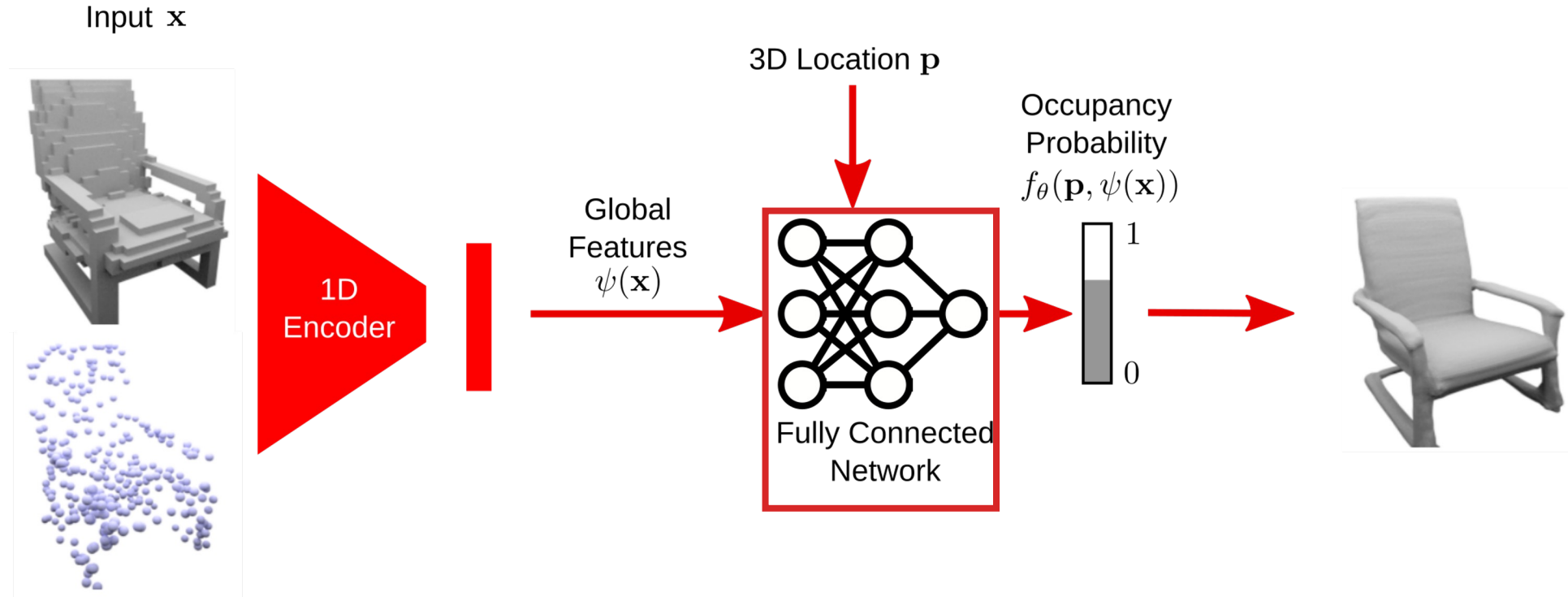


- Global latent code  $\Rightarrow$  **overly smooth geometry**



# Limitations

## Neural Implicit Representations

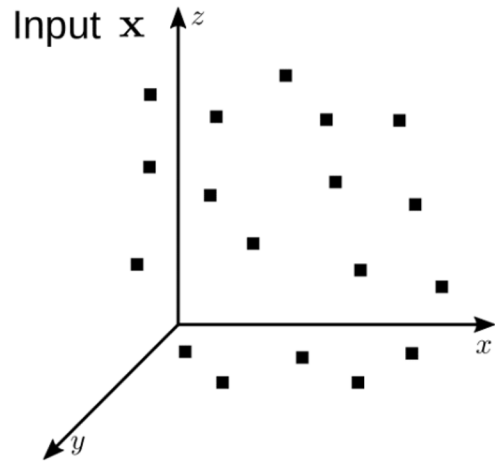


- Global latent code  $\Rightarrow$  **overly smooth geometry**
- Fully-connected architecture  $\Rightarrow$  **no translation equivariance**

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

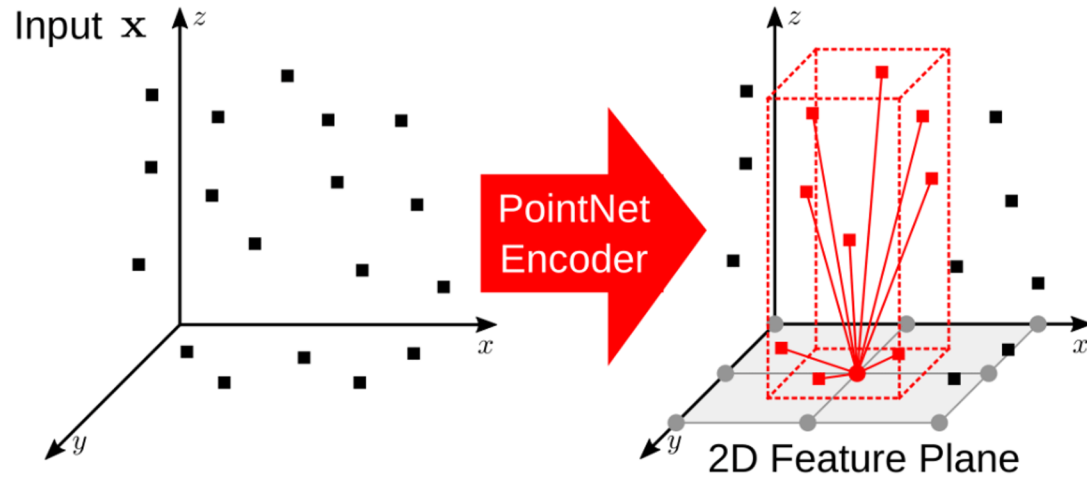
# Main Idea

## Convolutional Occupancy Networks



# Main Idea

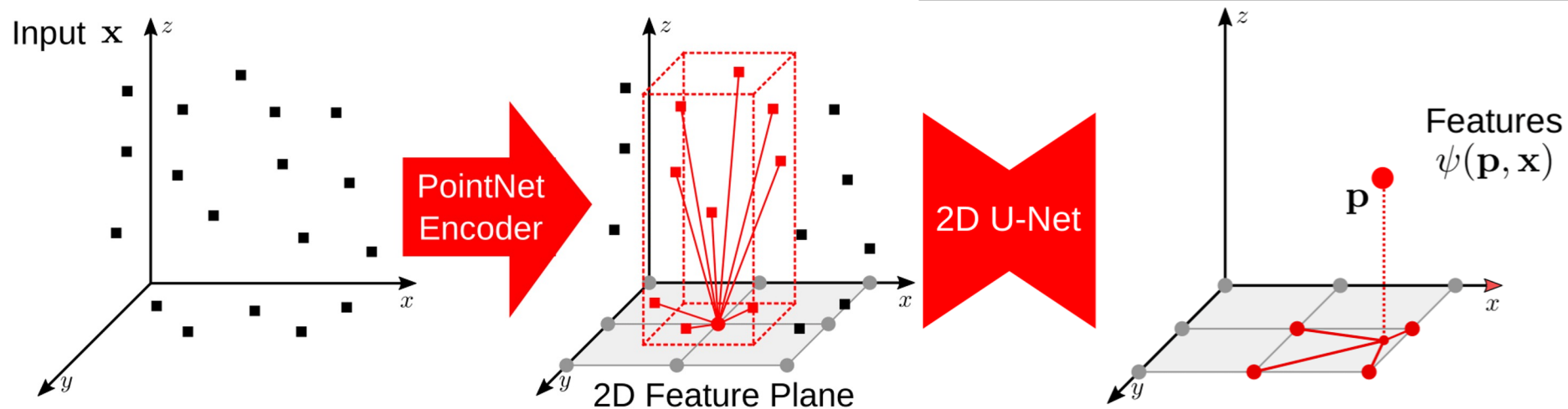
## Convolutional Occupancy Networks



- **2D Plane Encoder:** Project point features onto the canonical plane

# Main Idea

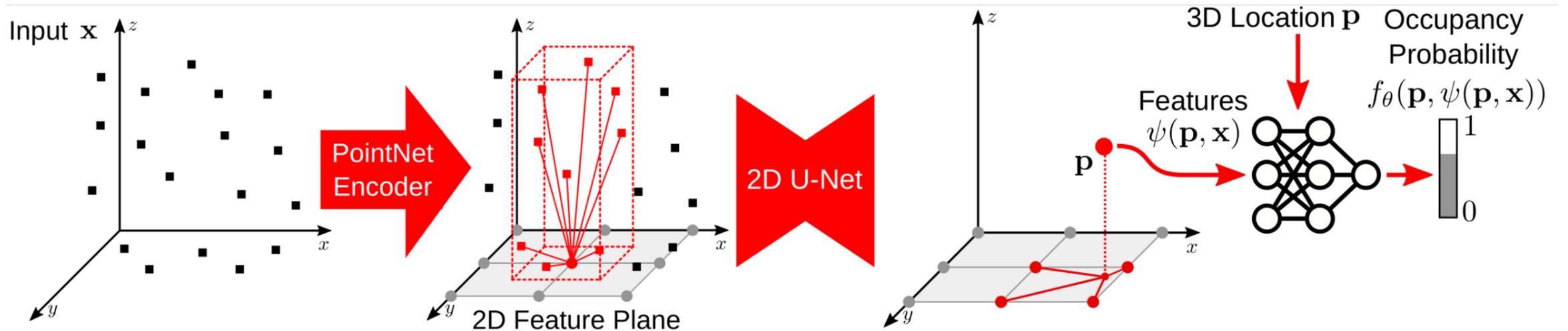
## Convolutional Occupancy Networks



- **2D Plane Encoder**: Project point features onto the canonical plane
- **2D Plane Decoder**: Processed by UNet, query features via bilinear interpolation

# Main Idea

## Convolutional Occupancy Networks

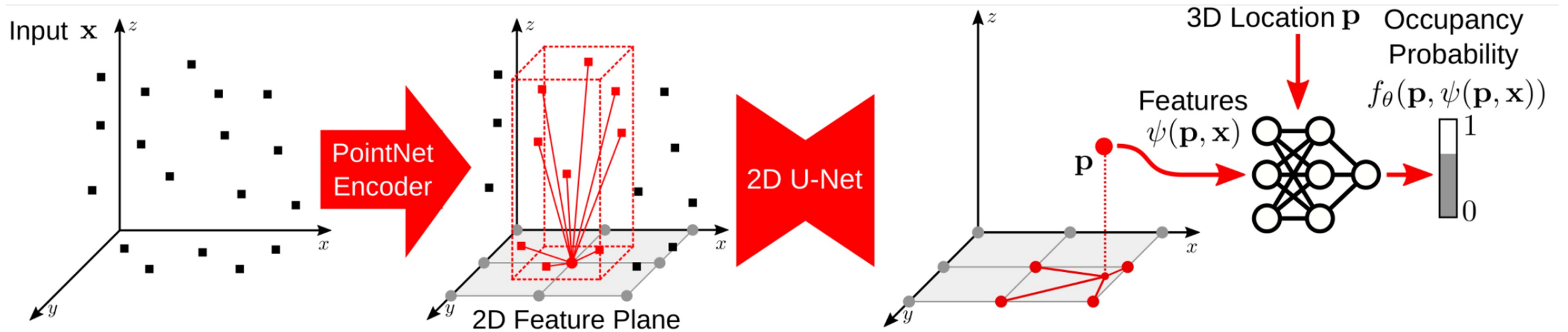


- **2D Plane Encoder:** Project point features onto the canonical plane
- **2D Plane Decoder:** Processed by UNet, query features via bilinear interpolation
- **Occupancy Net:** Shallow MLP  $f_{\theta}(\cdot)$



# Main Idea

## Convolutional Occupancy Networks

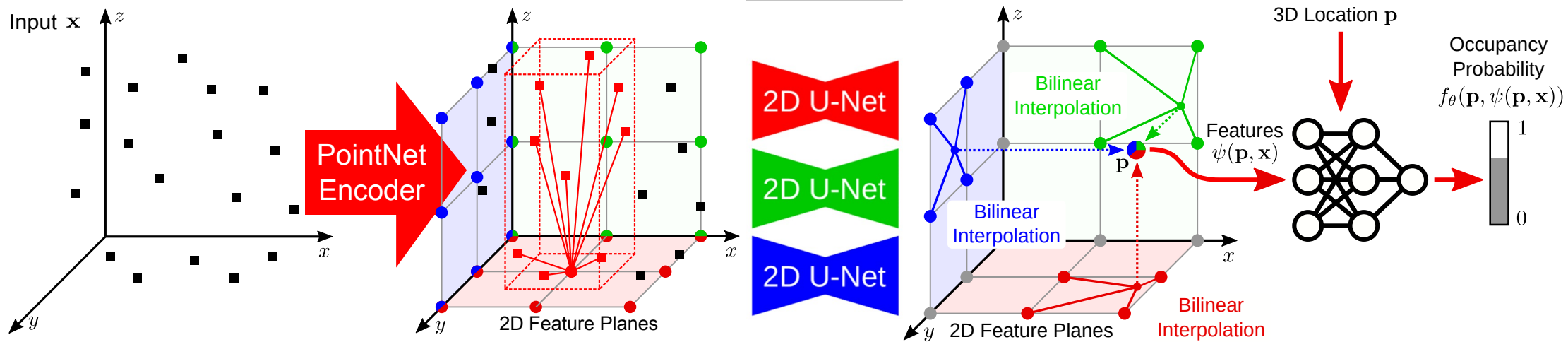


**Can we do better?**

For now, features only on the **ground plane**...

# Main Idea – “Tri-plane”

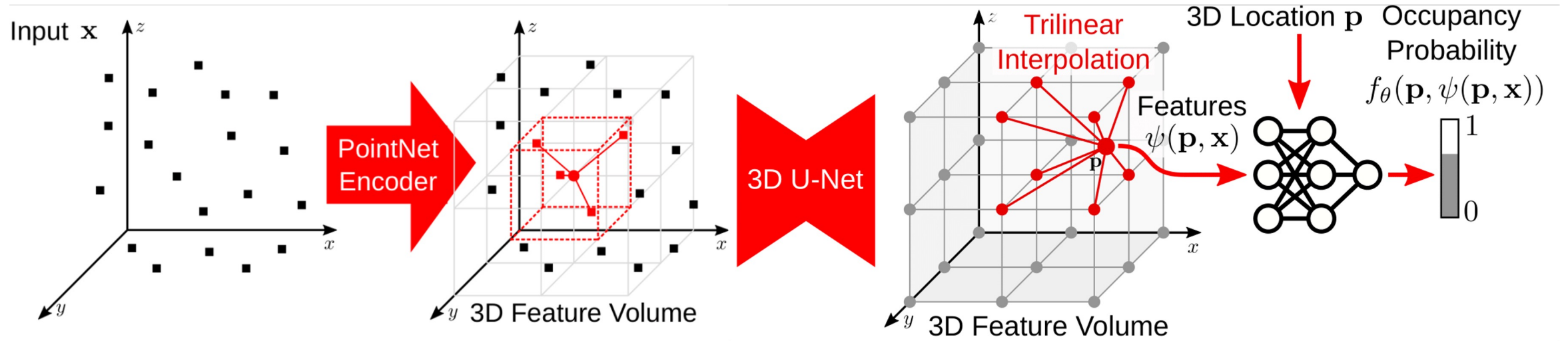
## Convolutional Occupancy Networks



Project features on **X, Y, Z canonical planes**

# Main Idea – 3D Volume

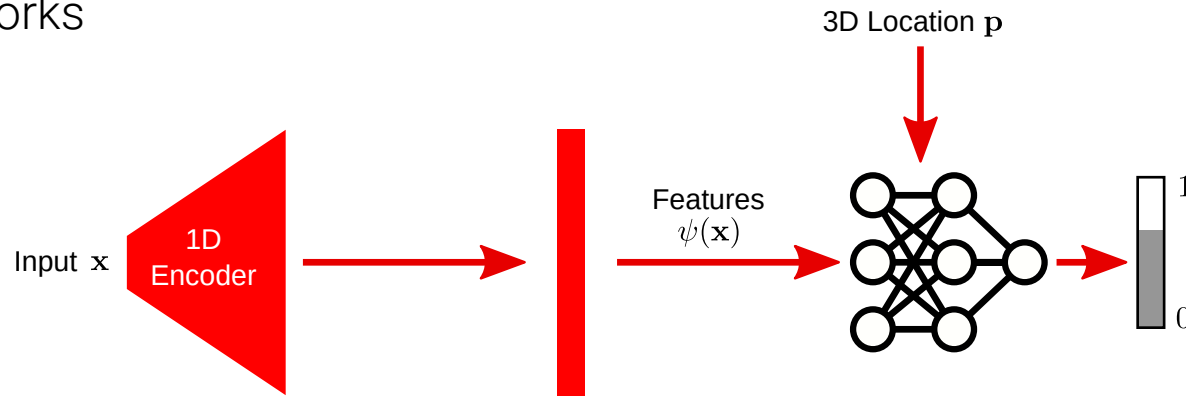
## Convolutional Occupancy Networks



Encode local information into a **3D feature volume**

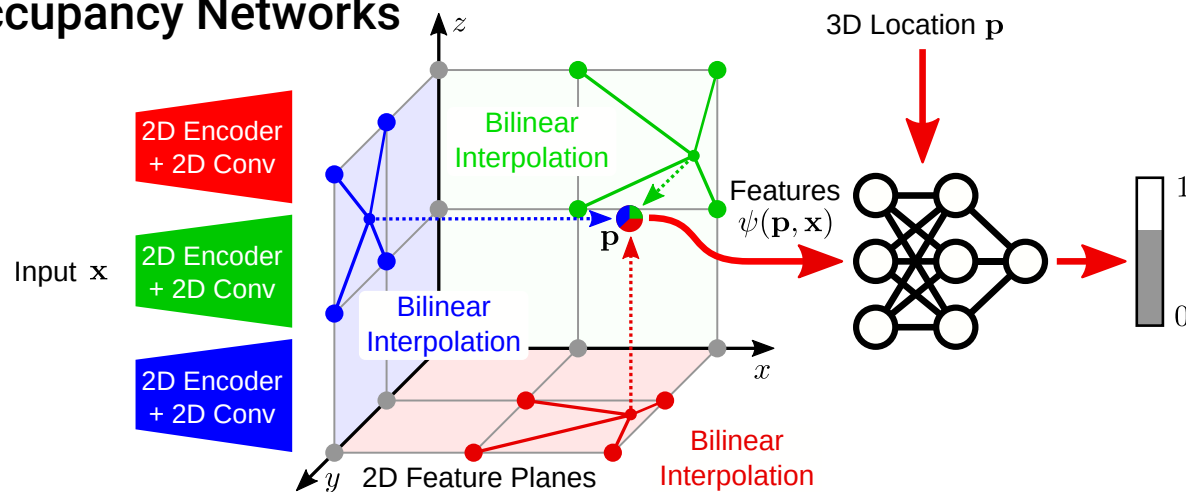
# Comparison

## Occupancy Networks



- global feature
- heavy FC network
- no translation equivariance

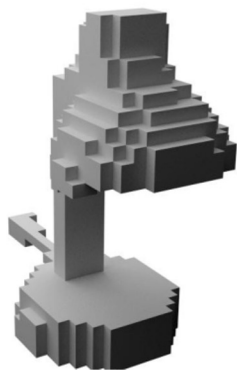
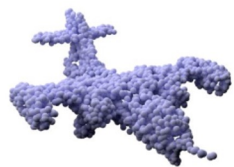
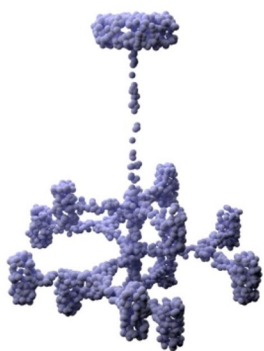
## Convolutional Occupancy Networks



- + local feature
- + shallow FC network
- + translation equivariance

# Results

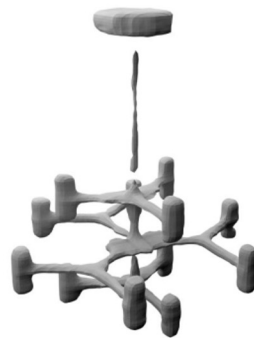
# Object-Level Reconstruction



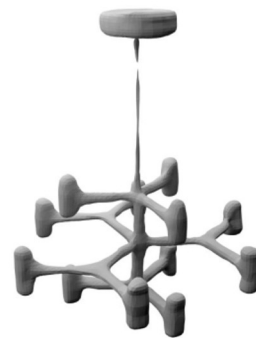
Input



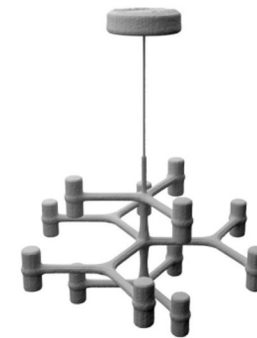
OccNet



Ours  
(Tri-plane)



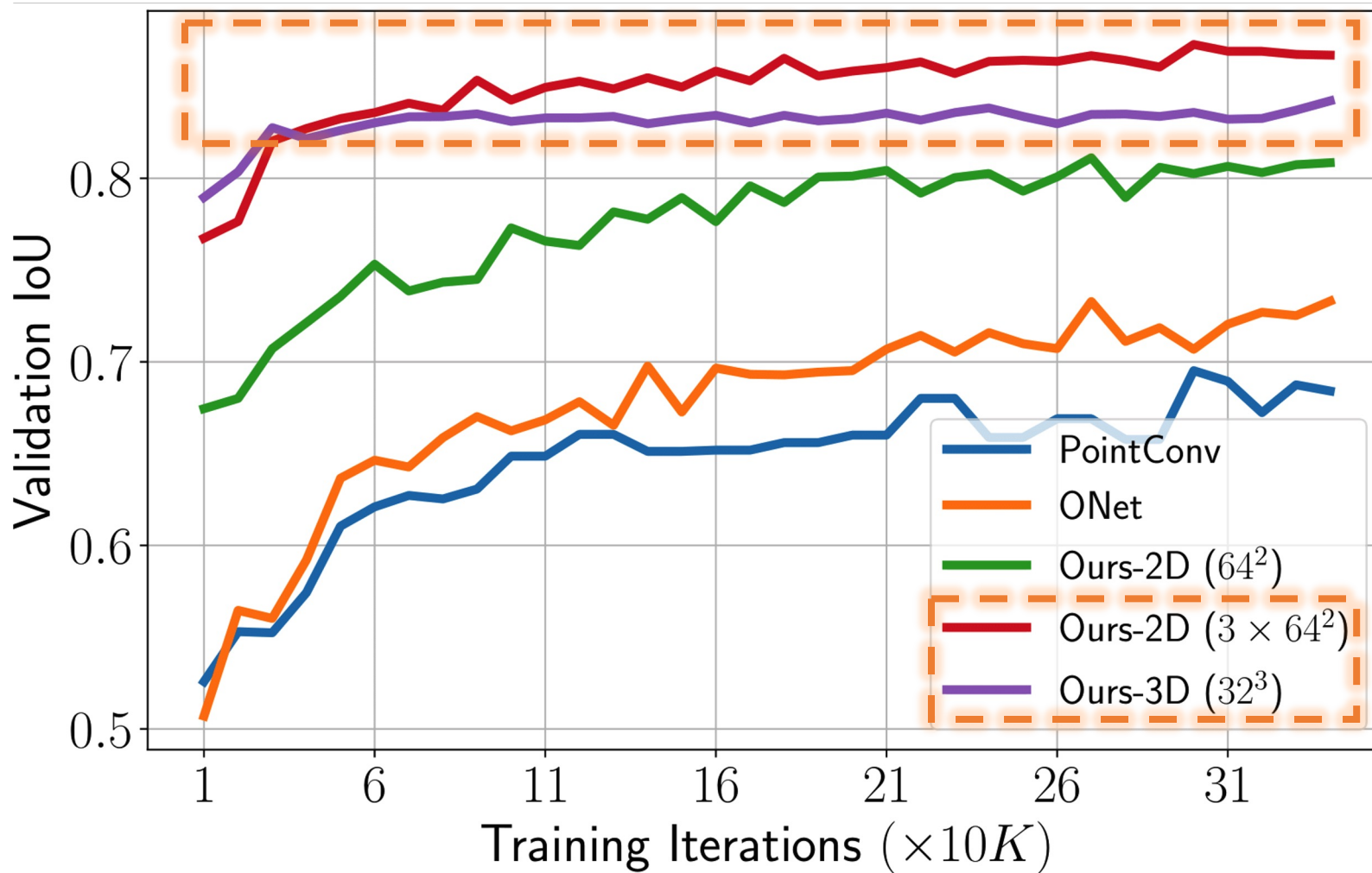
Ours  
(3D Volume)



GT Mesh



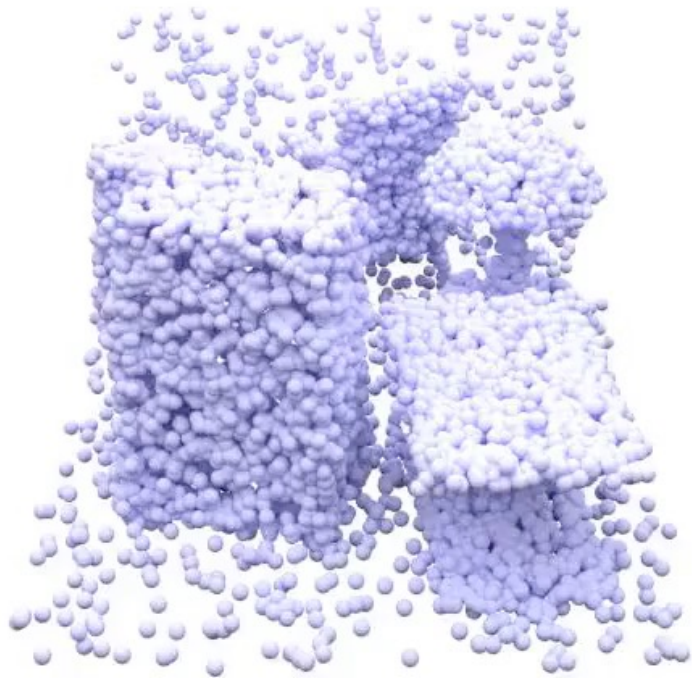
# Training Speed





# Scene-Level Reconstruction

Train and evaluate on synthetic room



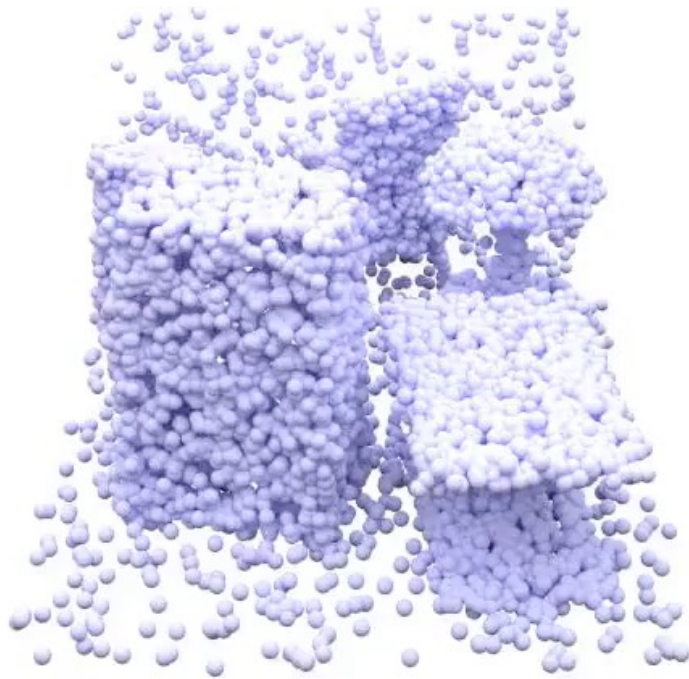
Input



GT Mesh

# Scene-Level Reconstruction

OccNet **fails** on room-level reconstruction



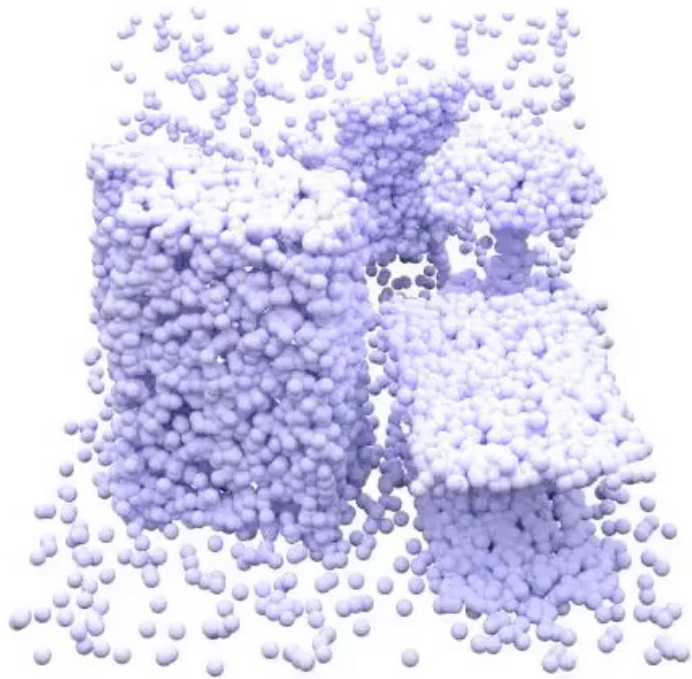
Input



OccNet

# Scene-Level Reconstruction

SPSR requires surface normal, output is **noisy**



Input

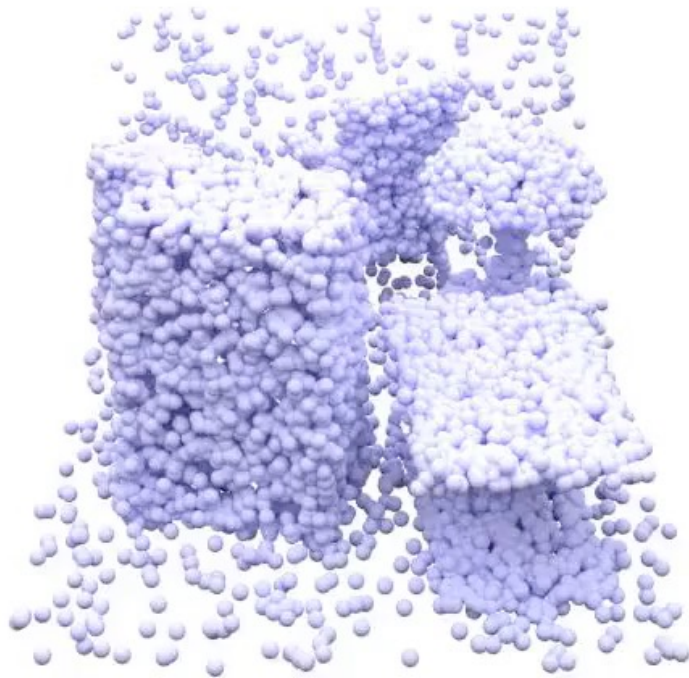


SPSR

(Screened Poisson Surface Reconstruction)

# Scene-Level Reconstruction

Ours preserves better details



Input



Ours

# Scene-Level Reconstruction



OccNet



SPSR



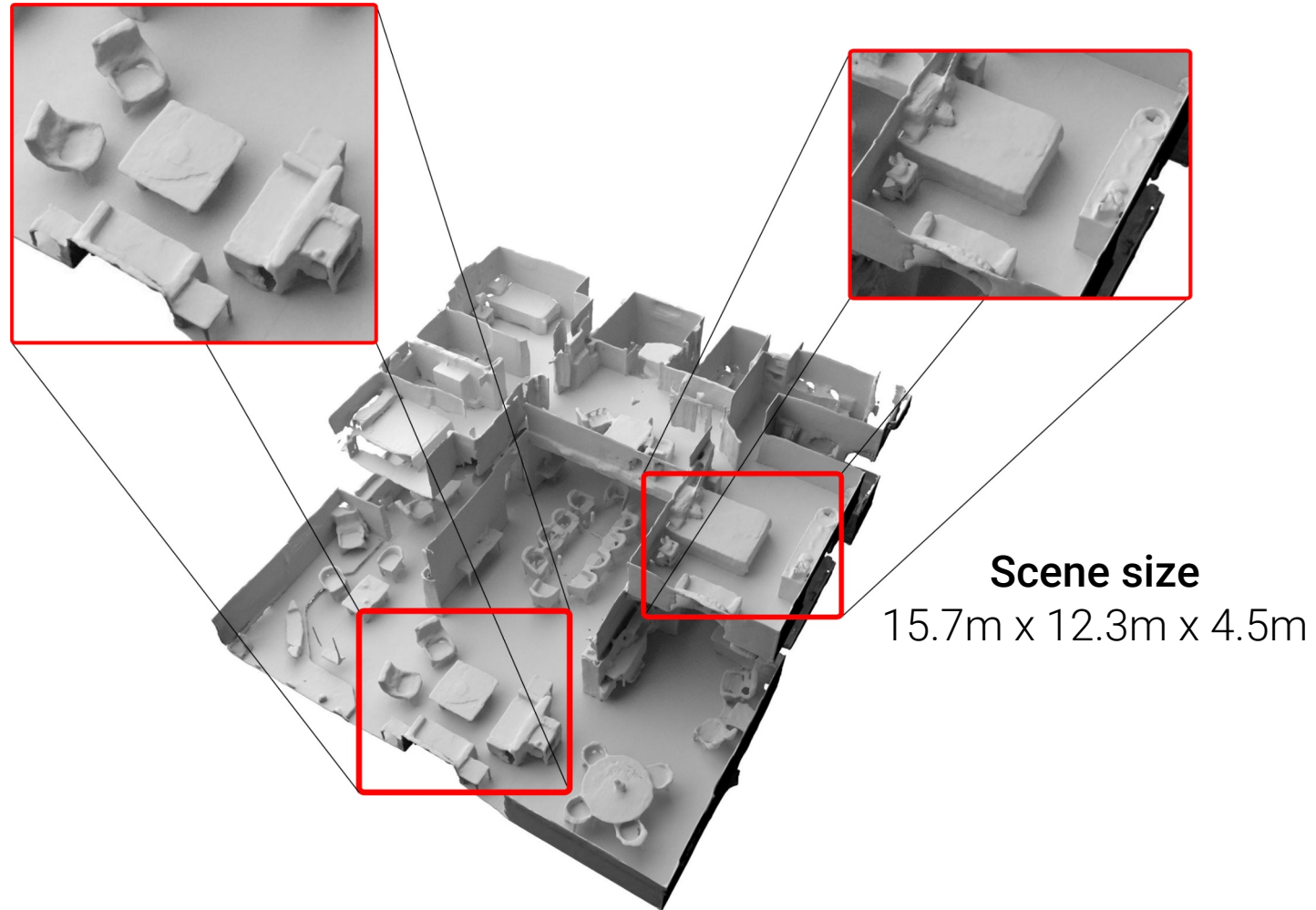
**Ours**



# Large-Scale Reconstruction

Reconstruct a big house in Matterport3D

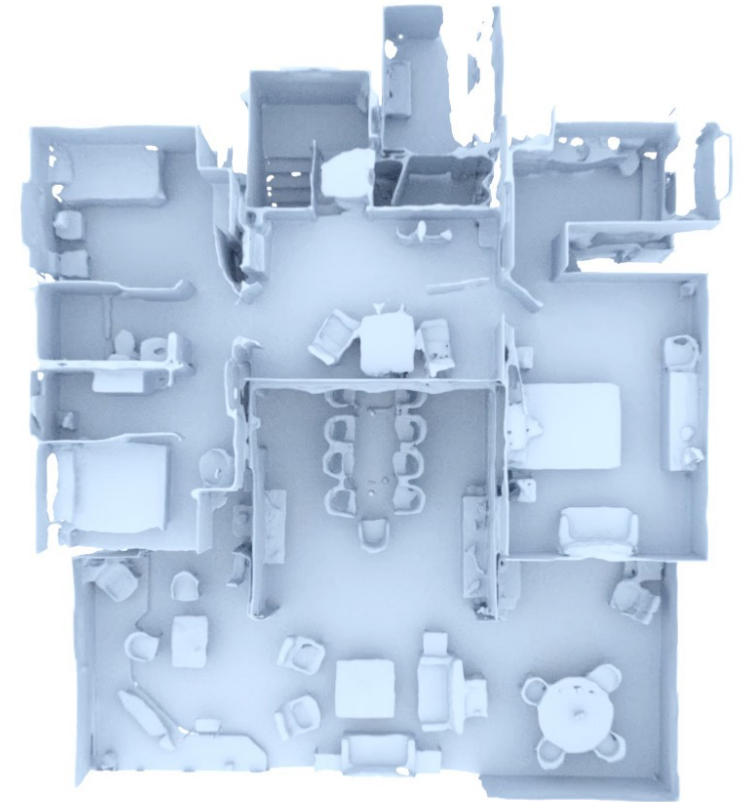
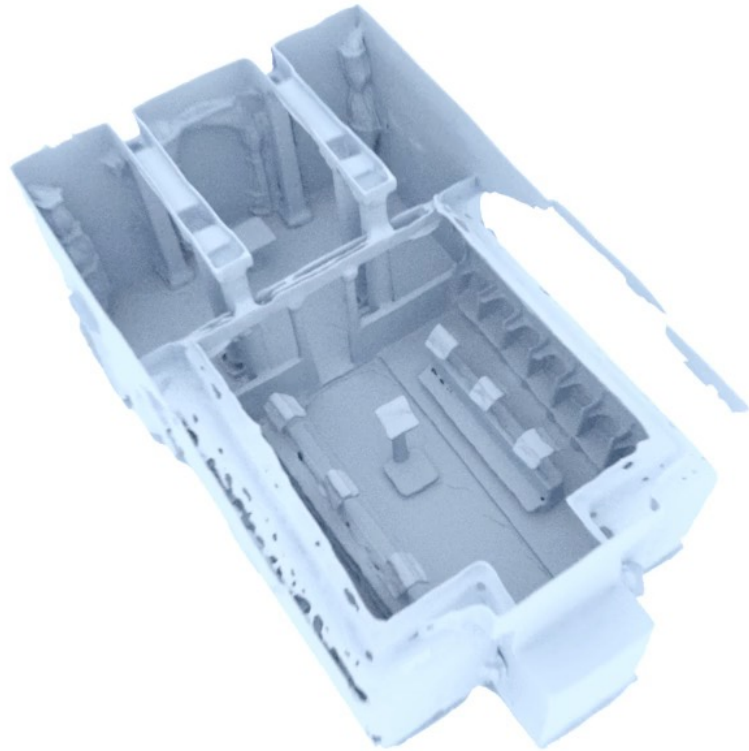
- Fully convolutional model
  - Sliding-window evaluation
  - Scale to any size
- Trained on **synthetic crops**



**Our reconstruction output**

# Large-Scale Reconstruction

Reconstruct large-scale scenes in Matterport3D



# ConvOccNet - TL;DR

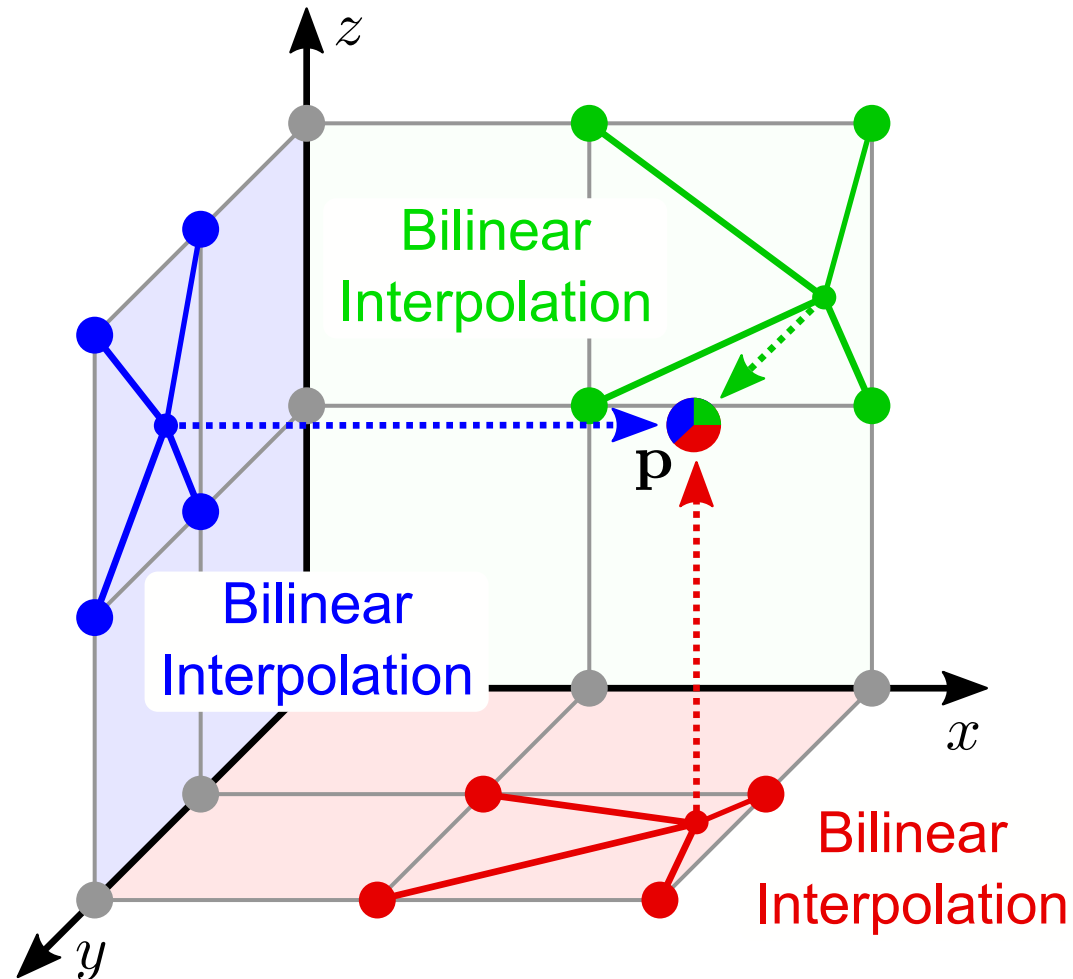
+ Three **hybrid representations** for neural fields

a) Ground plane    b) **Tri-plane**    c) 3D volume

+ CNN's **translation equivariance** rocks

+ **Synthetic-to-real** generalization

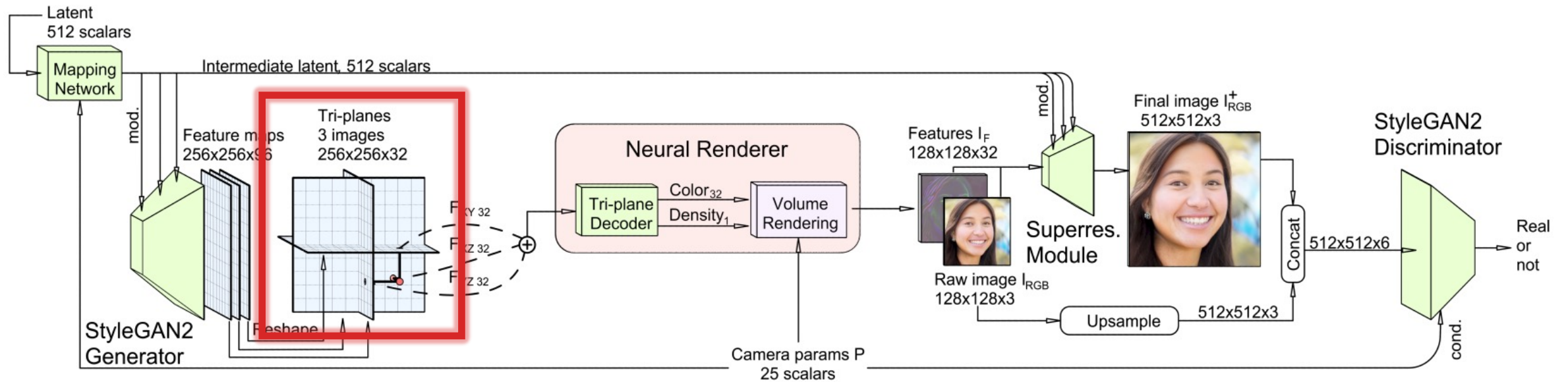
# “Tri-plane” Representations



**Reviewer 2:** *“What is the point of having that 3-plane representation?”*

# “Tri-plane” Representations

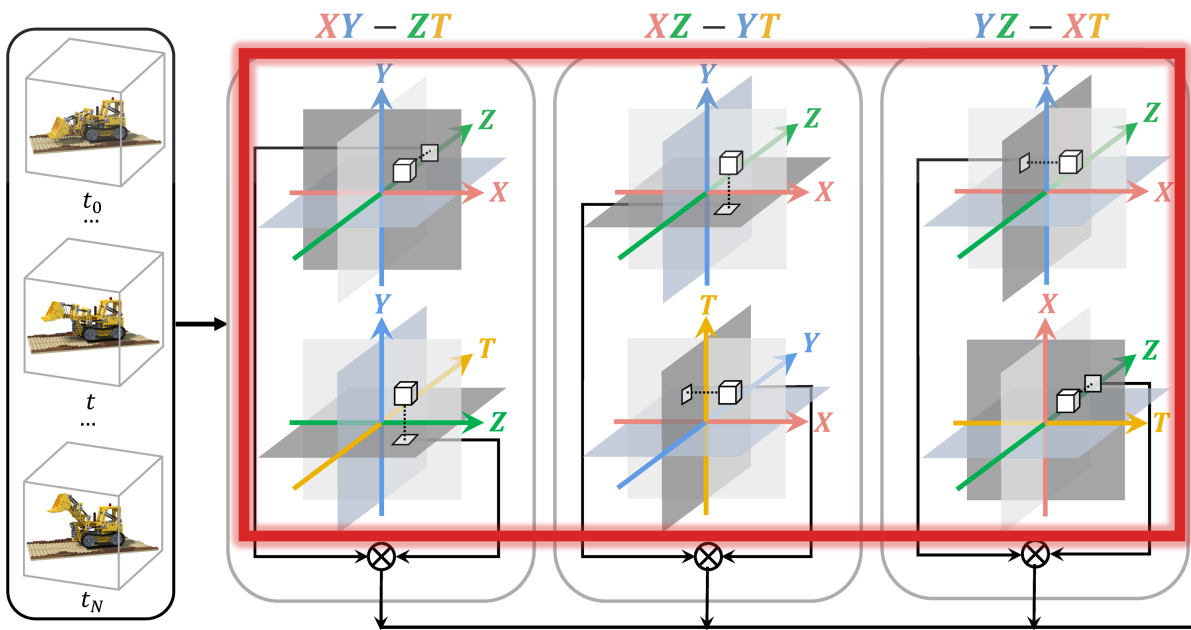
## High Fidelity 3D-Aware View Synthesis





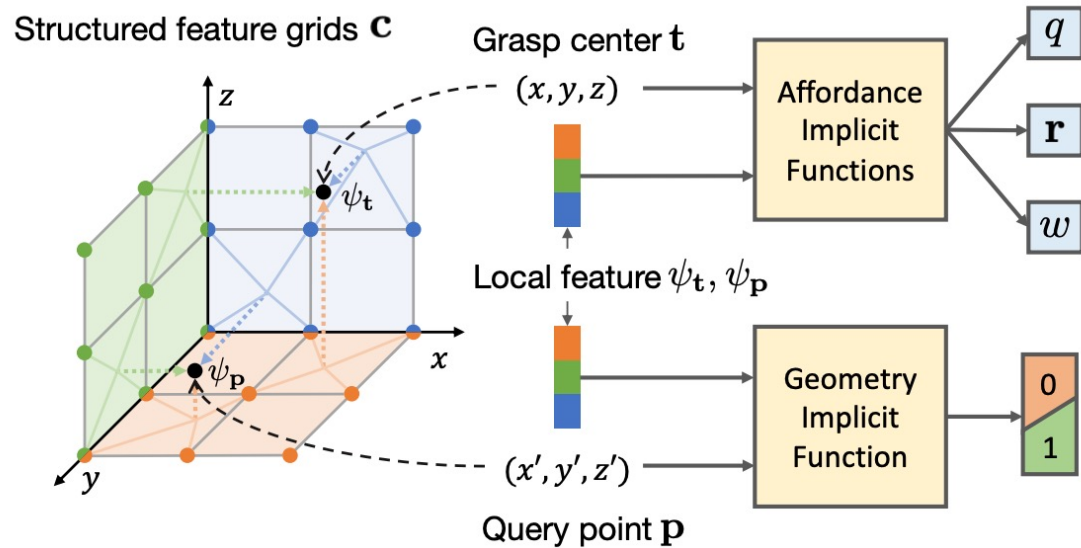
# “Tri-plane” Representations

## Efficient 4D View Synthesis



# “Tri-plane” Representations

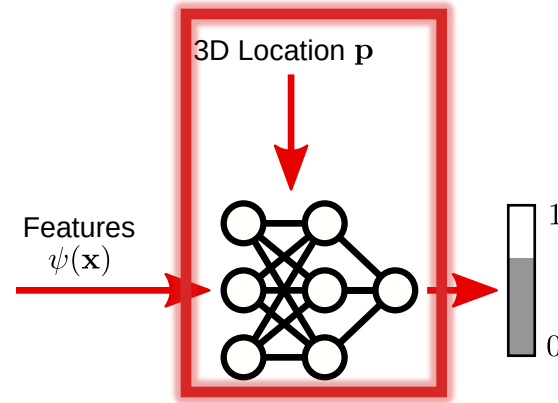
## Robot Grasping



# ConvOccNet - Limitations

## — Very slow inference

For a grid of **128<sup>3</sup>**, > 2 million MLP forward passes !



---

## Shape As Points: A Differentiable Poisson Solver

---

Songyou Peng<sup>1,2</sup>   Chiyu “Max” Jiang\*<sup>†</sup>   Yiyi Liao<sup>2,3†</sup>   Michael Niemeyer<sup>2,3</sup>

Marc Pollefeys<sup>1,4</sup>   Andreas Geiger<sup>2,3</sup>

<sup>1</sup>ETH Zurich

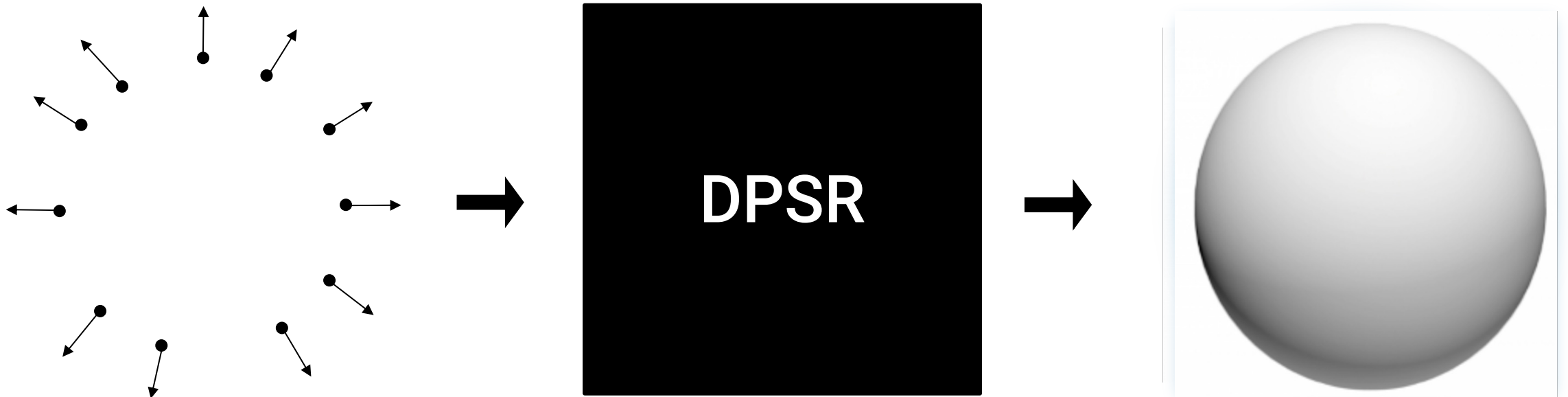
<sup>2</sup>Max Planck Institute for Intelligent Systems, Tübingen

<sup>3</sup>University of Tübingen

<sup>4</sup>Microsoft

# Shape As Points

A Differentiable Point-to-Mesh Layer



**No network evaluation, **fast!****



Inputs



GT Mesh



ConvOccNet

327 ms




**SAP**

12 ms

# ConvOccNet - Limitations

- Very slow inference

 We have **SAP**

- Only reconstruct from 3D noisy point clouds

Can we **online** reconstruct purely **from 2D observations**?



# ConvOccNet - Limitations

— Very slow inference



We have **SAP**

— Only reconstruct from 3D noisy point clouds



**NICE-SLAM: Neural Implicit Scalable Encoding for SLAM**

Zihan Zhu<sup>1,2\*</sup>

Songyou Peng<sup>2,4\*</sup>

Viktor Larsson<sup>3</sup>

Weiwei Xu<sup>1</sup>

Hujun Bao<sup>1</sup>

Zhaopeng Cui<sup>1†</sup>

Martin R. Oswald<sup>2,5</sup>

Marc Pollefeys<sup>2,6</sup>

<sup>1</sup>State Key Lab of CAD&CG, Zhejiang University

<sup>2</sup>ETH Zurich

<sup>3</sup>Lund University

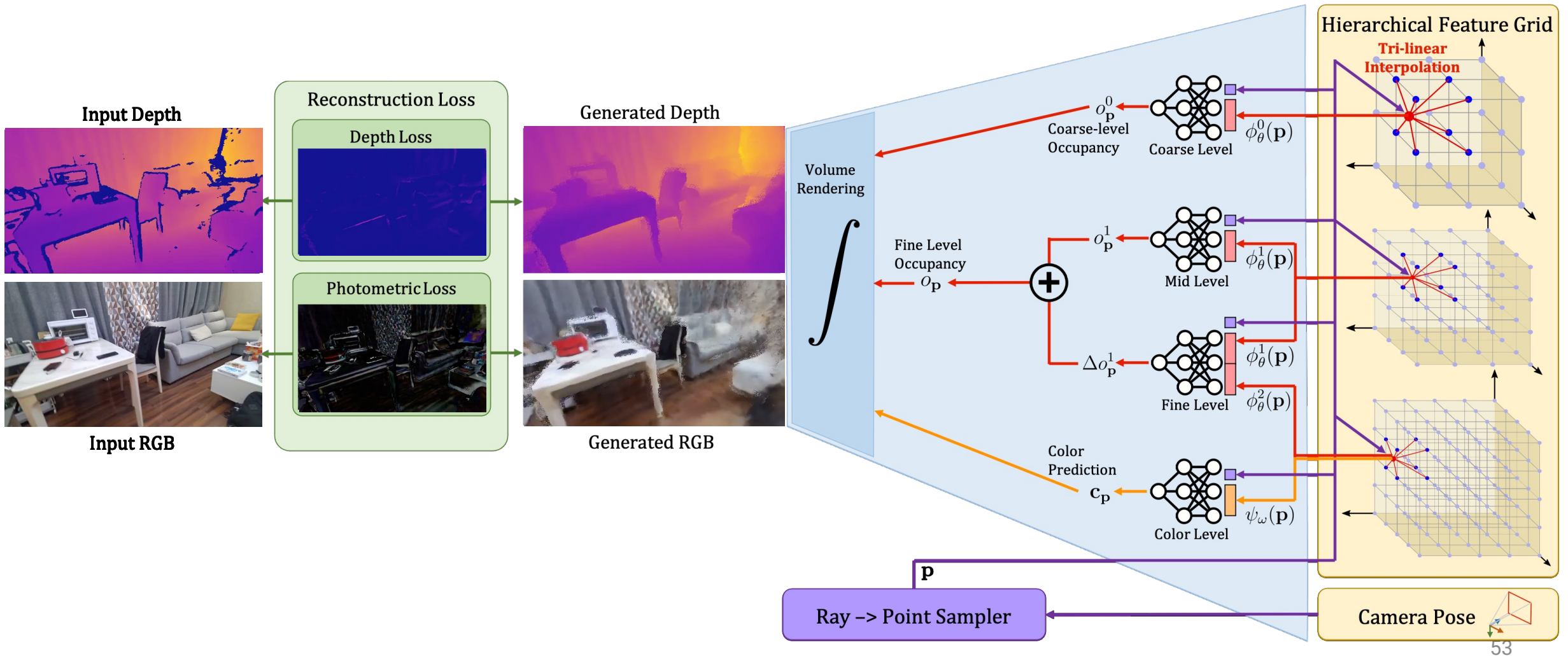
<sup>4</sup>MPI for Intelligent Systems, Tübingen

<sup>5</sup>University of Amsterdam

<sup>6</sup>Microsoft

# NICE-SLAM

## Neural Implicit Scalable Encoding for SLAM

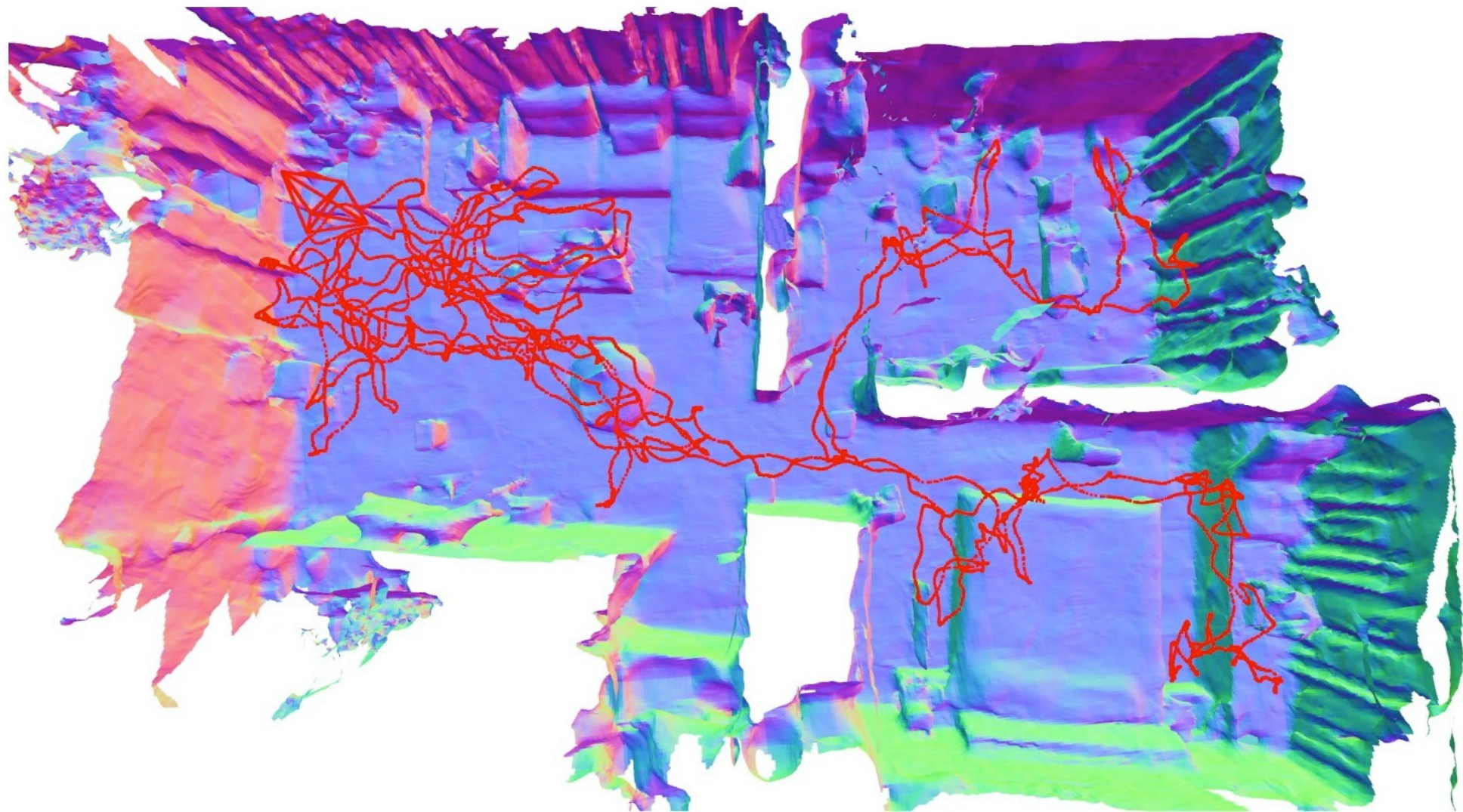




## RGB-D Sequences



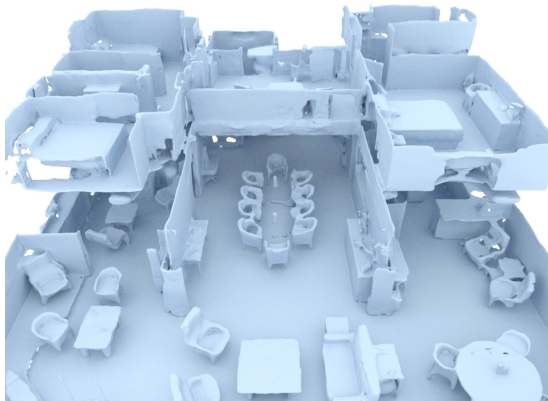
40x Speed



# This Thesis

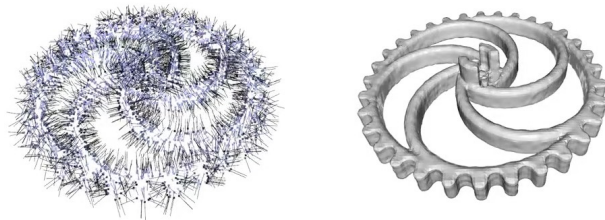
Develop 3D Neural Scene Representations  
for **3D Reconstruction** and **3D Scene Understanding**

## 1. Complex Scenes



**ConvOccNet**  
ECCV 2020 (Spotlight)

## 2. Fast Inference



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

## 3. From 2D Observations



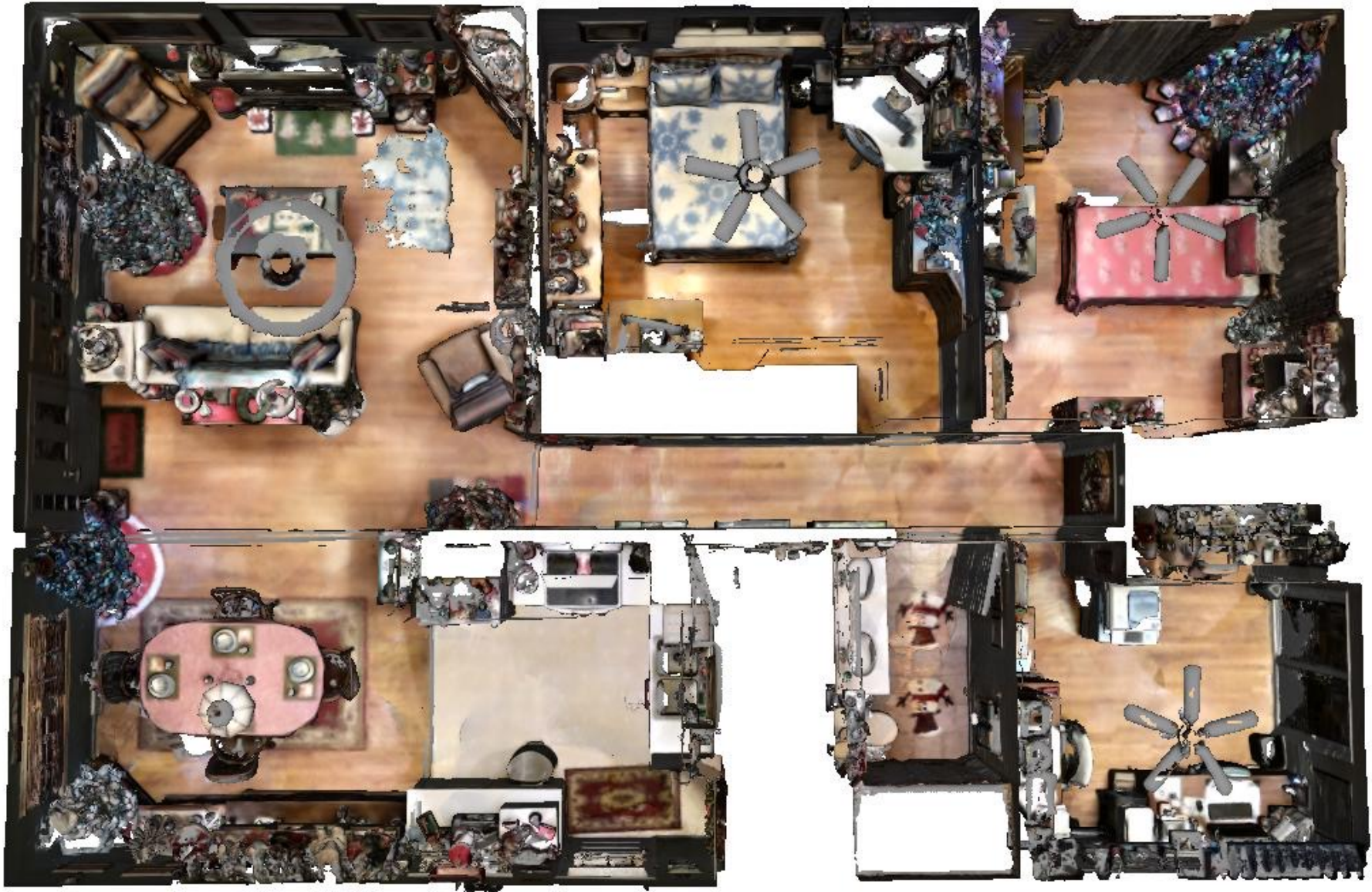
**NICE-SLAM**  
CVPR 2022

## 4. Arbitrary Queries



**OpenScene**  
CVPR 2023





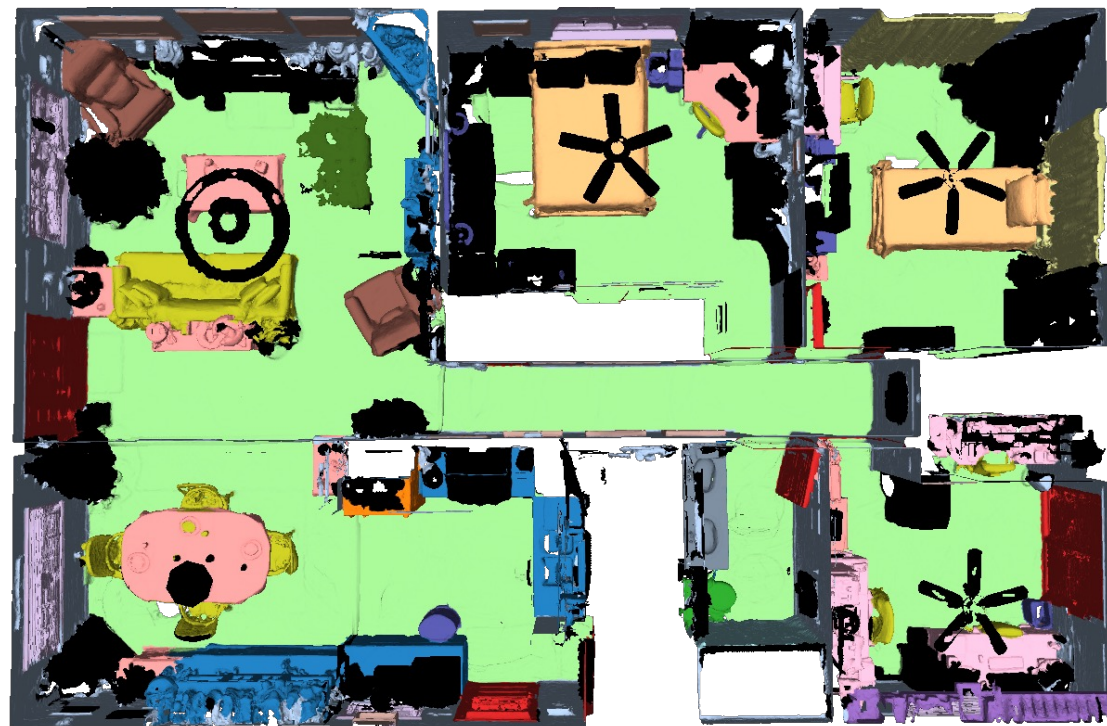
Input 3D Geometry





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 learn a scene representation to handle all these tasks  
**without labeled 3D data?**

# This Thesis

Develop 3D Neural Scene Representations  
for **3D Reconstruction** and **3D Scene Understanding**

## 1. Complex Scenes



**ConvOccNet**  
ECCV 2020 (Spotlight)

## 2. Fast Inference



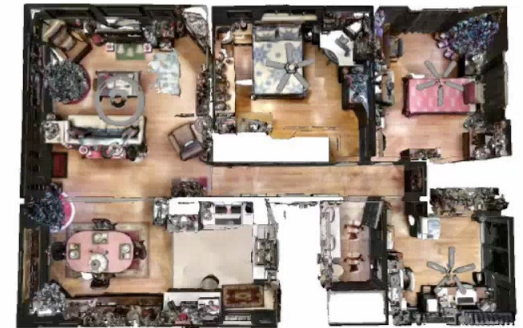
**Shape As Points**  
NeurIPS 2021 (Oral)

## 3. From 2D Observations



**NICE-SLAM**  
CVPR 2022

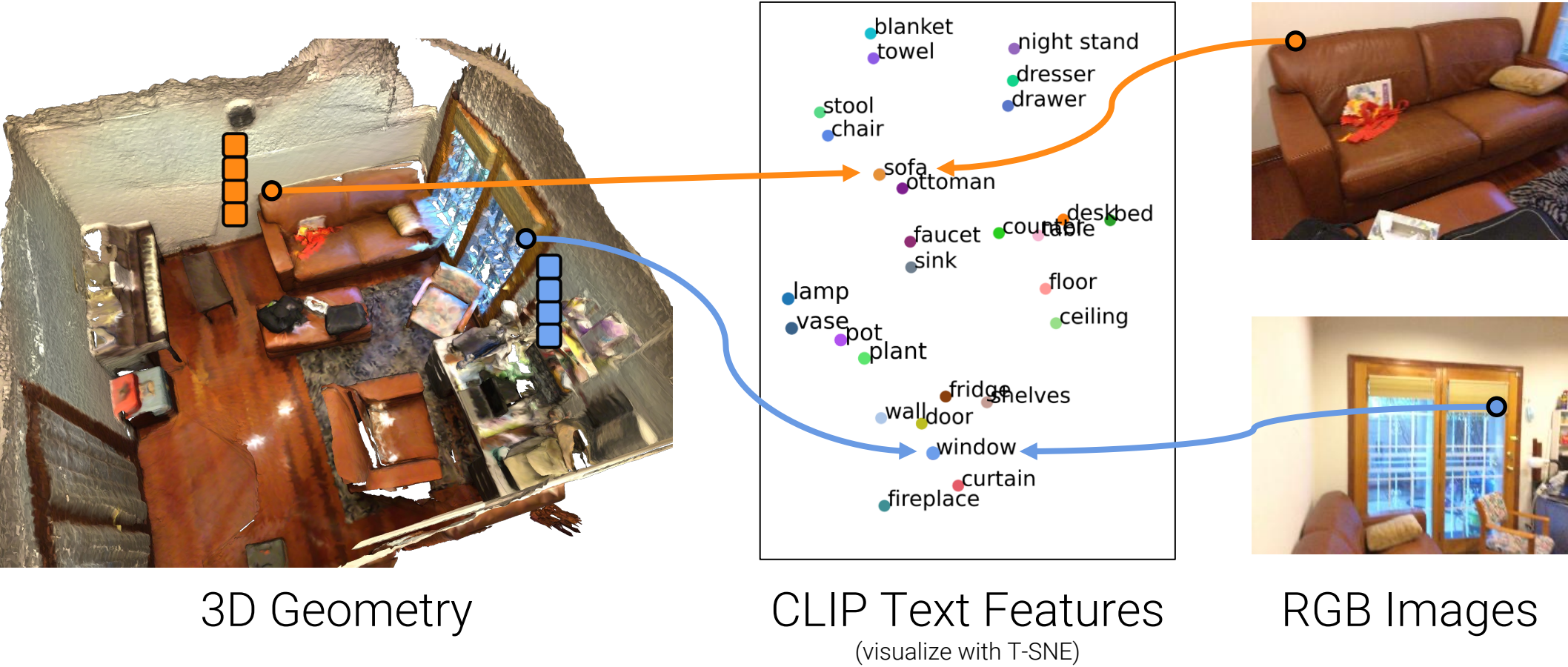
## 4. Arbitrary Queries



**OpenScene**  
CVPR 2023

# Key Idea

## Co-embed 3D Features with CLIP Features



3D Geometry

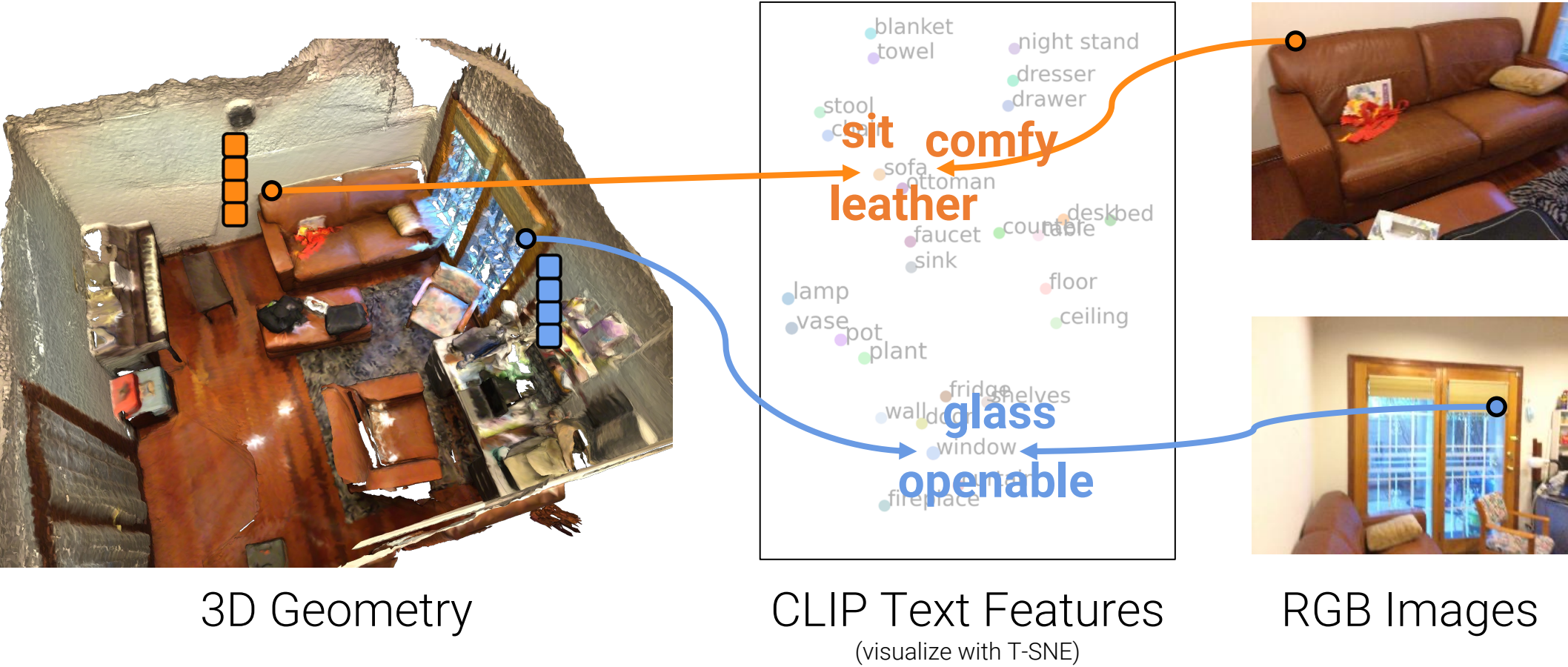
CLIP Text Features  
(visualize with T-SNE)

RGB Images



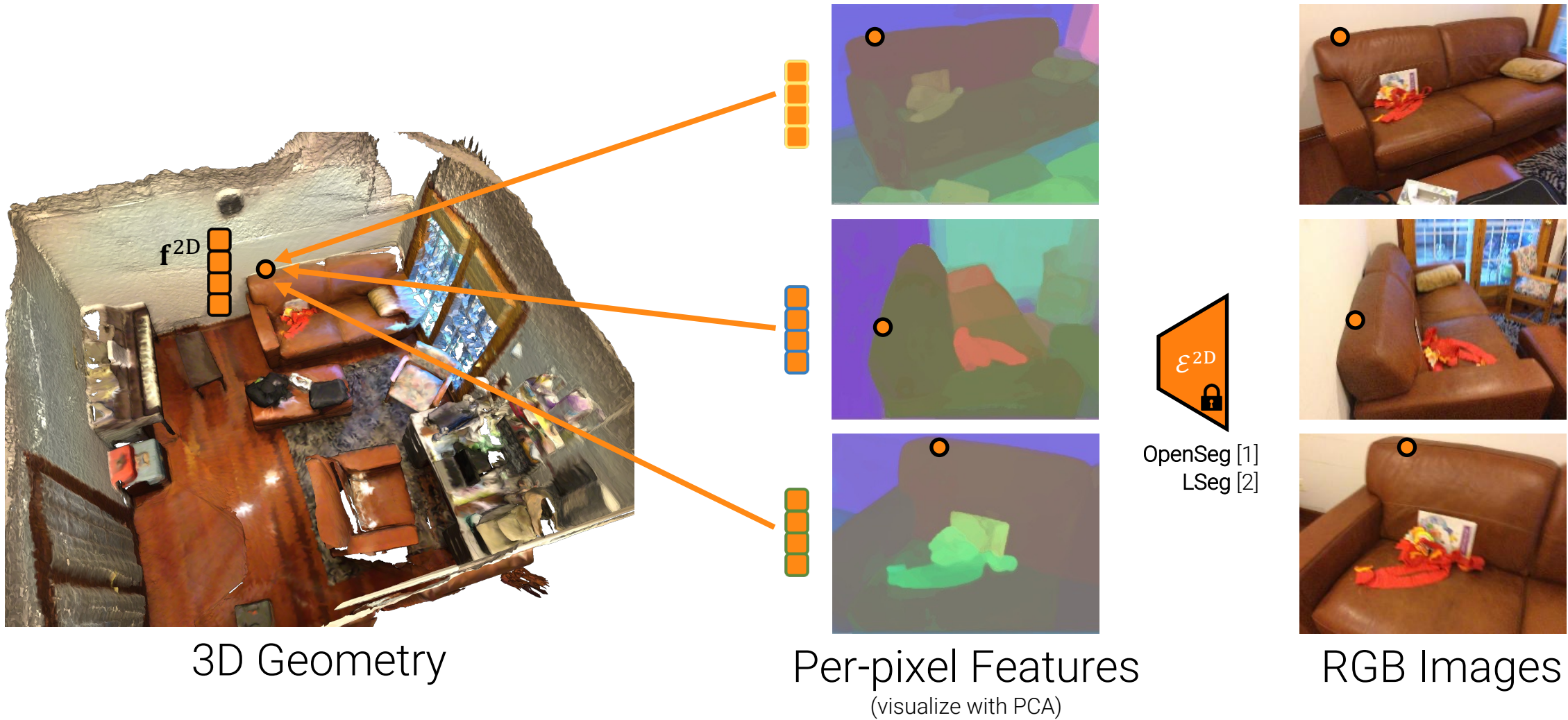
# Key Idea

## 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

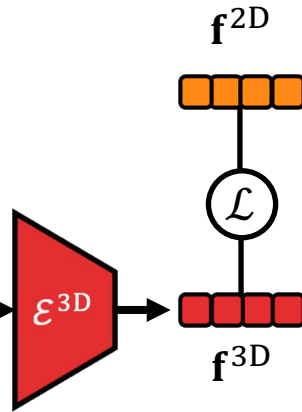
[2] Li, Weinberger, Belongie, Koltun, Ranftl: [Language-driven Semantic Segmentation](#). ICLR 2022



# Step 2: 3D Feature Distillation

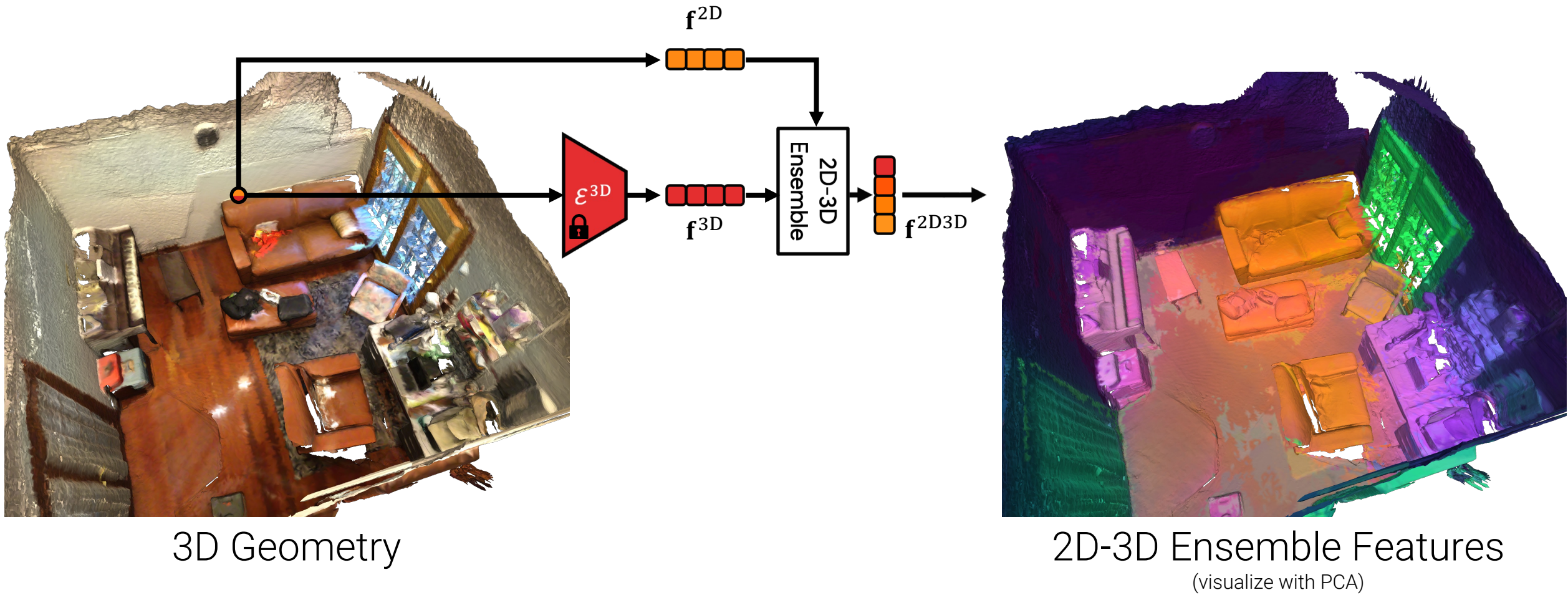


3D Geometry



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

# 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	■ toilet brush	■ computer
■ picture	■ mirror	■ banister	■ counter	■ photo	■ refridgerator	■ purse	■ bin	■ headboard	■ printer	■ sheet	■ rope	■ display case	■ bottle of soap	■ drum	■ knob
■ window	■ towel	■ stairs	■ bench	■ toilet paper	■ bookshelf	■ door way	■ chest	■ telephone	■ telephone	■ bucket	■ ball	■ toilet paper holder	■ water cooler	■ whiteboard	■ computer
■ chair	■ sink	■ stool	■ garbage bin	■ fan	■ wardrobe	■ basket	■ microwave	■ blanket	■ blanket	■ glass	■ exercise equipment	■ tea pot	■ range hood	■ paper	■ projector
■ pillow	■ shelves	■ vase	■ fireplace	■ railing	■ pipe	■ chandelier	■ blinds	■ flower pot	■ handle	■ dishwasher	■ tray	■ stuffed animal	■ candelabra		

# Image-based 3D Scene Query



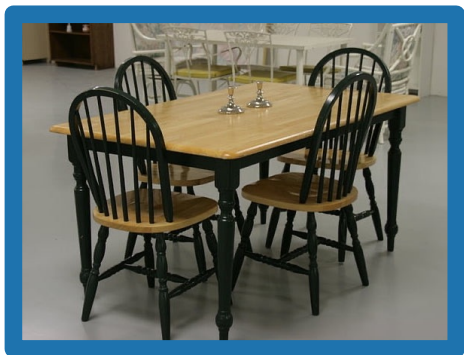


Image Queries

Given 3D Geometry

# **Interactive Demo**

Open-vocabulary 3D Scene Exploration



Text queries:



# OpenScene - TL;DR

- + Open up a **wide range of applications** by leveraging large vision-language models
- + Inspire future works to shift to open-vocabulary tasks
- Currently all power comes from 2D foundation models

# This Thesis

Develop Neural Scene Representations  
for **3D Reconstruction** and **3D Scene Understanding**

## 1. Complex Scenes



ConvOccNet  
ECCV 2020 (Spotlight)

## 2. Fast Inference



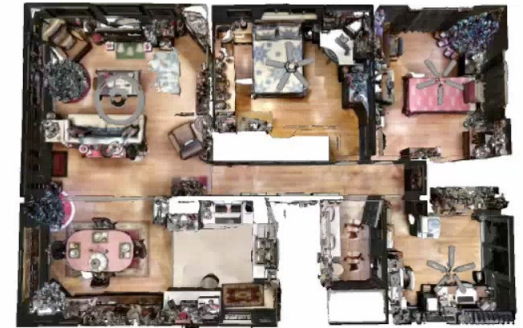
Shape As Points  
NeurIPS 2021 (Oral)

## 3. From 2D Observations



NICE-SLAM  
CVPR 2022

## 4. Arbitrary Queries



OpenScene  
CVPR 2023



# This Thesis

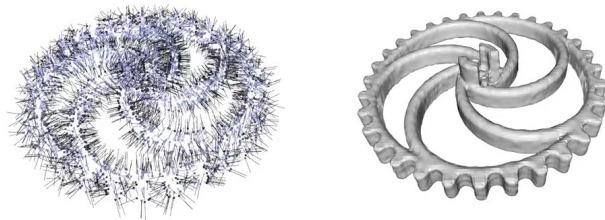
Develop Neural Scene Representations  
for **3D Reconstruction** and **3D Scene Understanding**

## 1. Complex Scenes



**ConvOccNet**  
ECCV 2020 (Spotlight)

## 2. Fast Inference



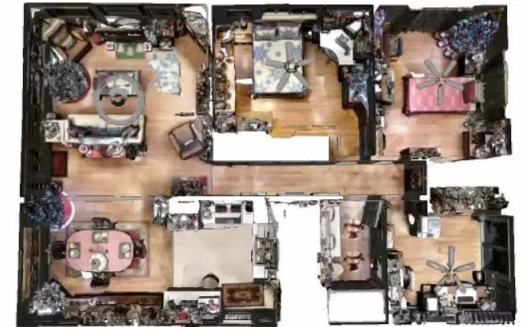
**Shape As Points**  
NeurIPS 2021 (Oral)

## 3. From 2D Observations



**NICE-SLAM**  
CVPR 2022

## 4. Arbitrary Queries



**OpenScene**  
CVPR 2023

**What is Next?**

# Practicality





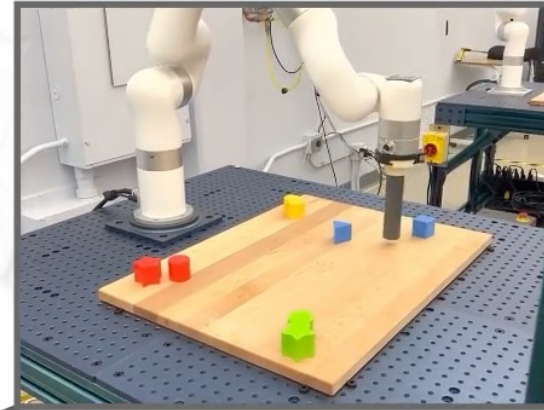
# Leverage Foundation Models for Everything

Robot Mobile Manipulation



Task: give me the chips from the drawer  
Next step: **Pick up the green chip bag**

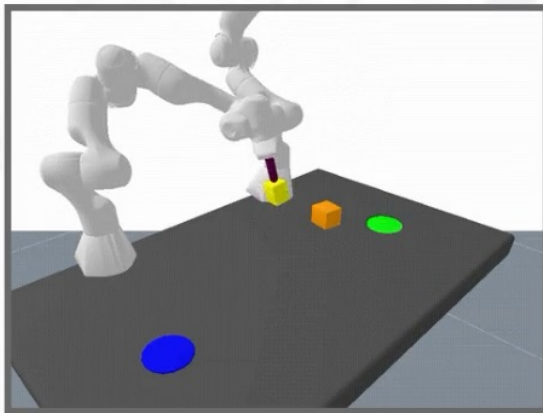
Robot Tabletop Manipulation



Task: sort blocks by colors into corners  
Next step: **Push blue blocks to the right**

**PaLM-E**  
**562B**

Task and Motion Panning



Q: How to put yellow block on blue plate?  
A: **Hand the yellow block to other arm**

Visual Question Answering



Q: What's in the image in emojis?  
A: 🍏🍌🍇🍏🍊🍏🍒

Video Credit:  
PALM-E

# Acknowledgements

## Supervisors



Marc Pollefeys



Andreas Geiger

## External Examiners



Leo Guibas



Vincent Sitzmann

## Collaborators



Michael Niemeyer



Lars Mescheder



Michael Oechsle



Yiyi Liao



Chiyu "Max" Jiang



Christian Reiser



Zihan Zhu



Zhaopeng Cui



Shaohui Liu



Viktor Larsson



Martin R. Oswald



Zehao Yu



Tom Funkhouser



Kyle Genova



Andrea Tagliasacchi



Shengqu Cai





Thank you all for such a wonderful journey!

ICCV23

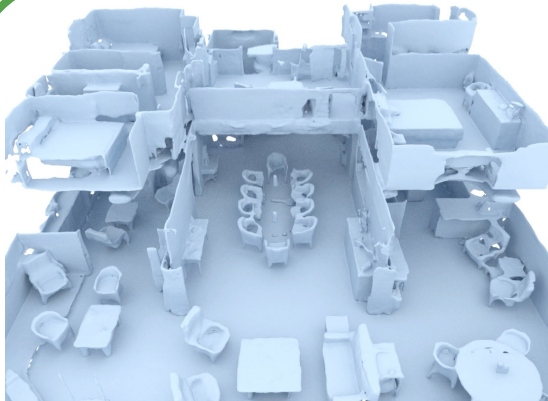
PARIS  
CVF

COMPUTER SOCIETY



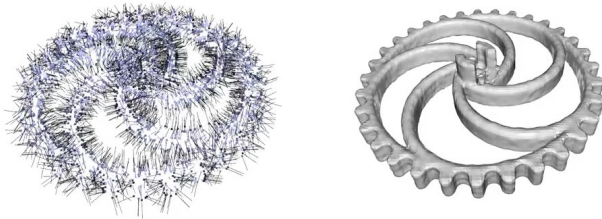
# Neural Scene Representations for 3D Reconstruction and Scene Understanding

Songyou Peng



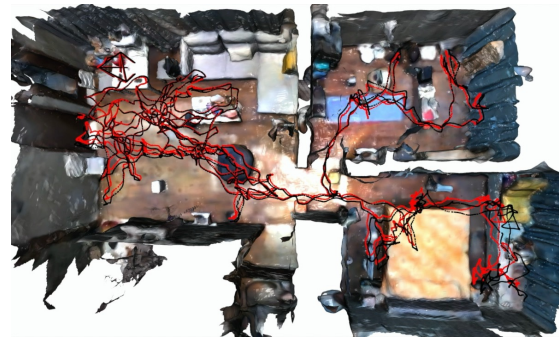
ConvOccNet

[pengsongyou.github.io/conv\\_onet](https://pengsongyou.github.io/conv_onet)



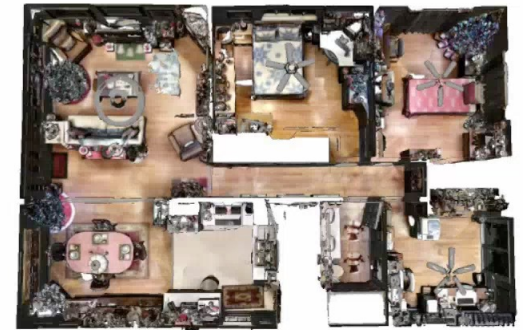
Shape As Points

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



NICE-SLAM

[pengsongyou.github.io/nice-slam](https://pengsongyou.github.io/nice-slam)



OpenScene

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