

# 3D Scene Understanding with Open Vocabularies

Songyou Peng\*

Kyle Genova

Chiyu "Max" Jiang

Andrea Tagliasacchi

Marc Pollefeys

Thomas Funkhouser

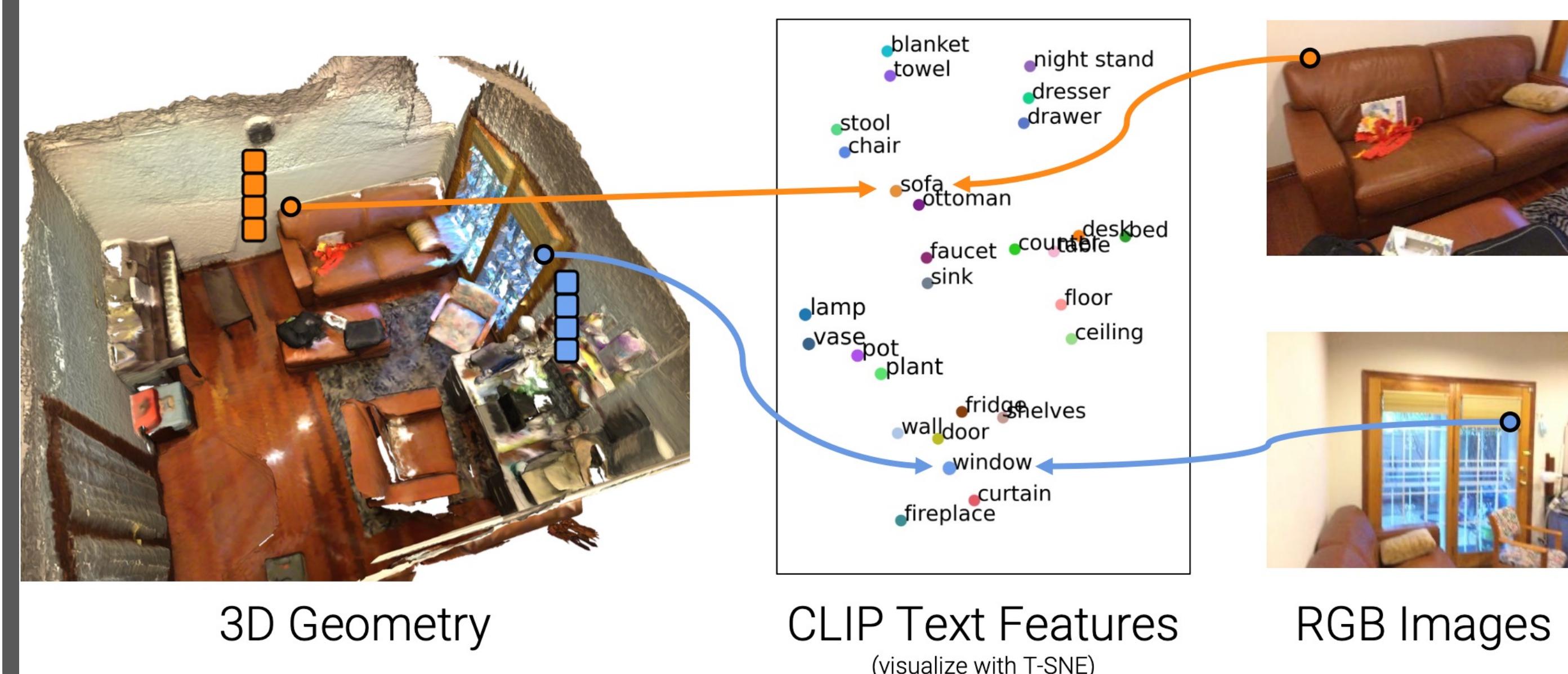
\* Work done during an internship at Google Research

## 1. Introduction

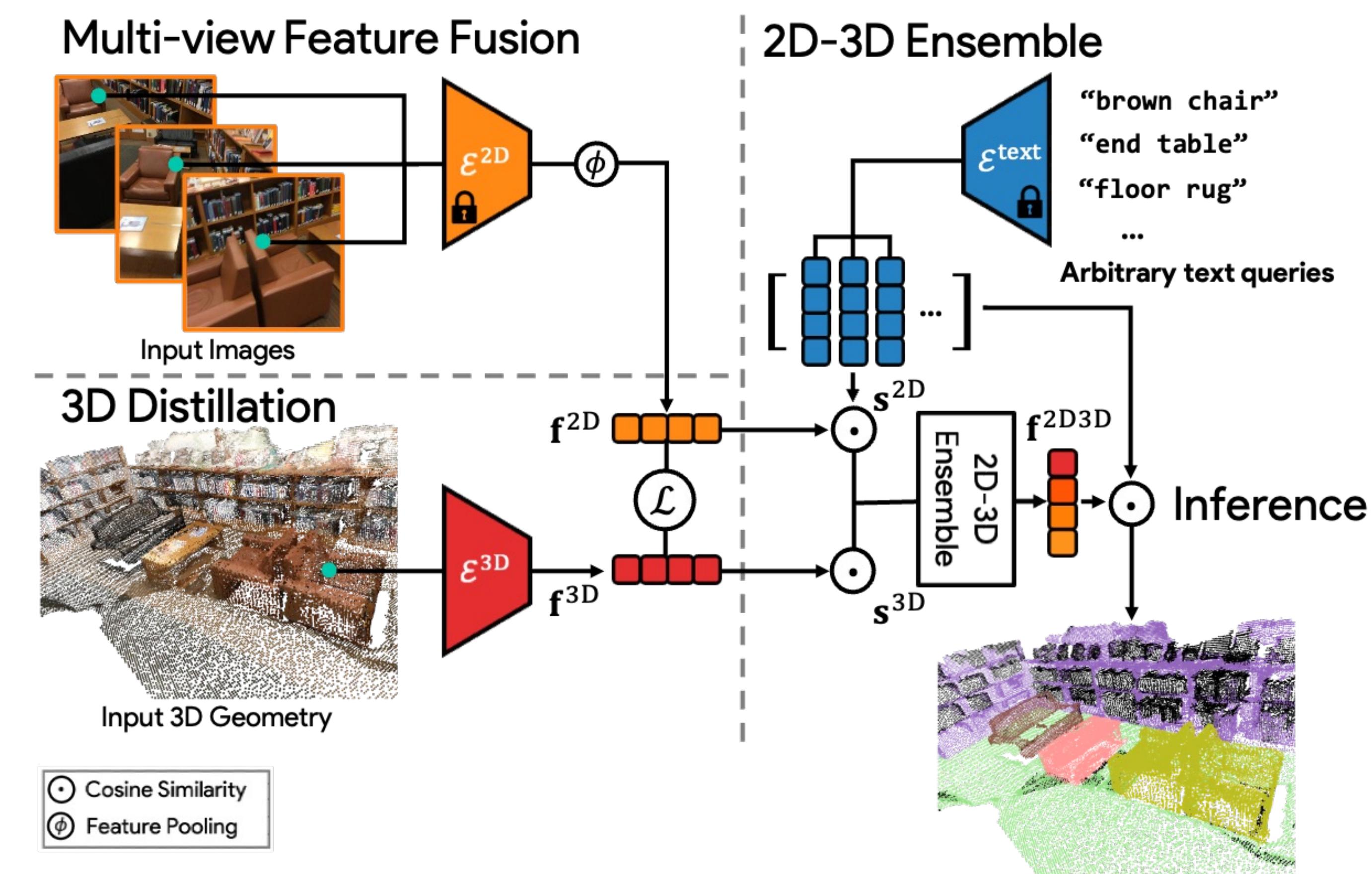
**Problem:** Traditional 3D scene understanding only train and test on some fixed common classes

**Goal:** A zero-shot approach to perform novel 3D scene understanding tasks with **arbitrary queries**

**Our Key Idea:** Co-embed 3D features with CLIP text and image features

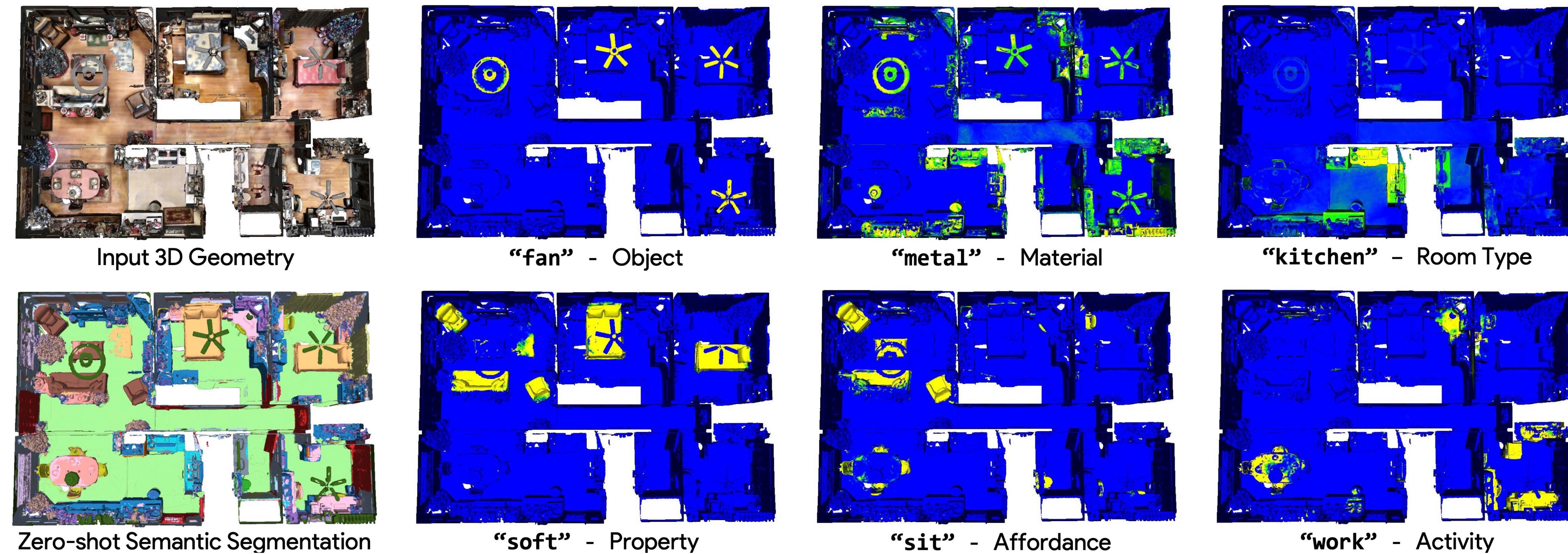


## 2. Method



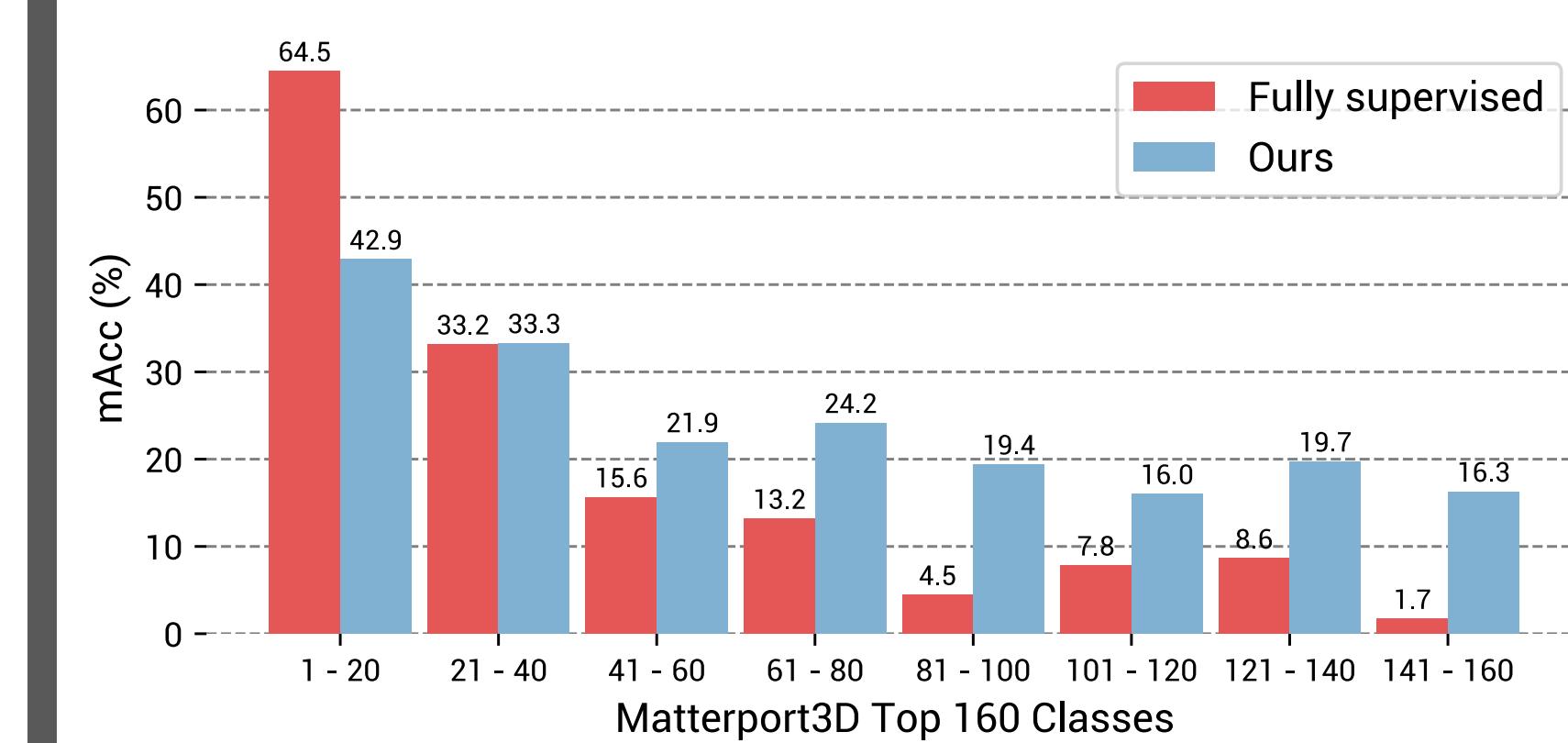
Paper, code, and real-time demo are available:  
[pengsongyou.github.io/openscene](https://pengsongyou.github.io/openscene)

## 3. Zero-shot Open-vocabulary Scene Exploration



## 5. More Studies

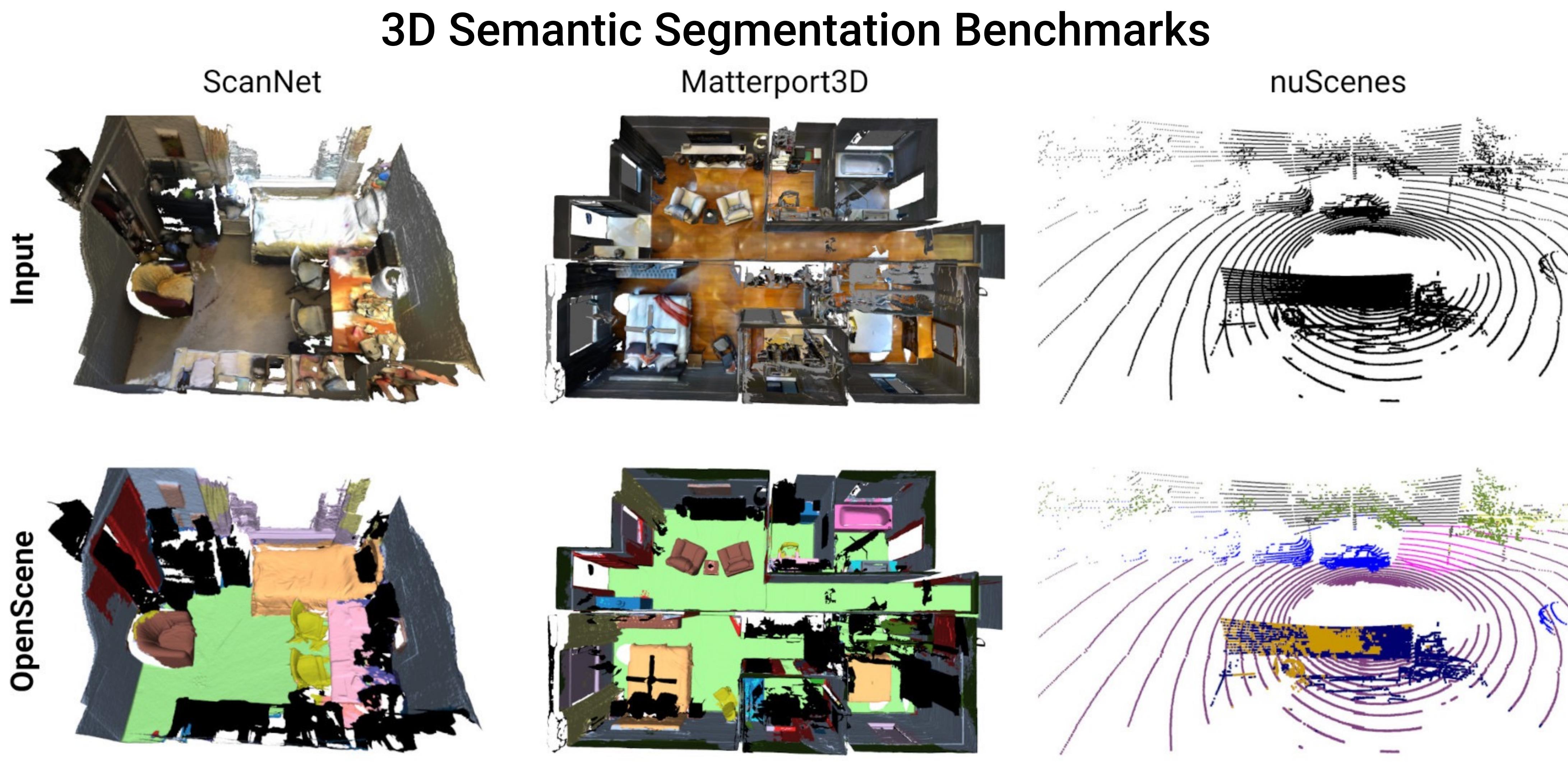
### Robust to Finding Rare Object



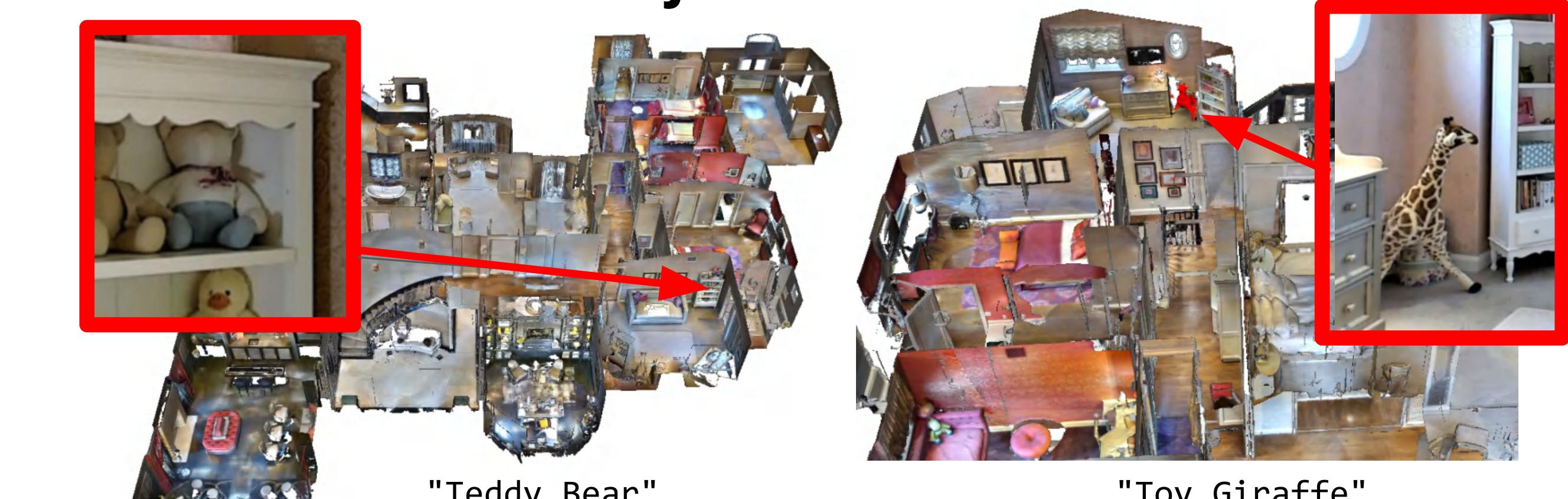
### Ablation Study

	ScanNet [1]		Matterport3D [4]	
	mIoU	mAcc	mIoU	mAcc
<b>Ours</b>	2D Fusion	50.0	62.7	32.3
LSeg	3D Distill	52.9	63.2	41.9
	2D-3D Ens.	<b>54.2</b>	66.6	<b>43.4</b>
<b>Ours</b>	2D Fusion	41.4	63.6	32.4
OpenSeg	3D Distill	46.0	66.3	41.3
	2D-3D Ens.	47.5	<b>70.7</b>	42.6
				<b>59.2</b>

## 4. Additional Applications



## Object Retrieval



## Image-based 3D Object Detection

