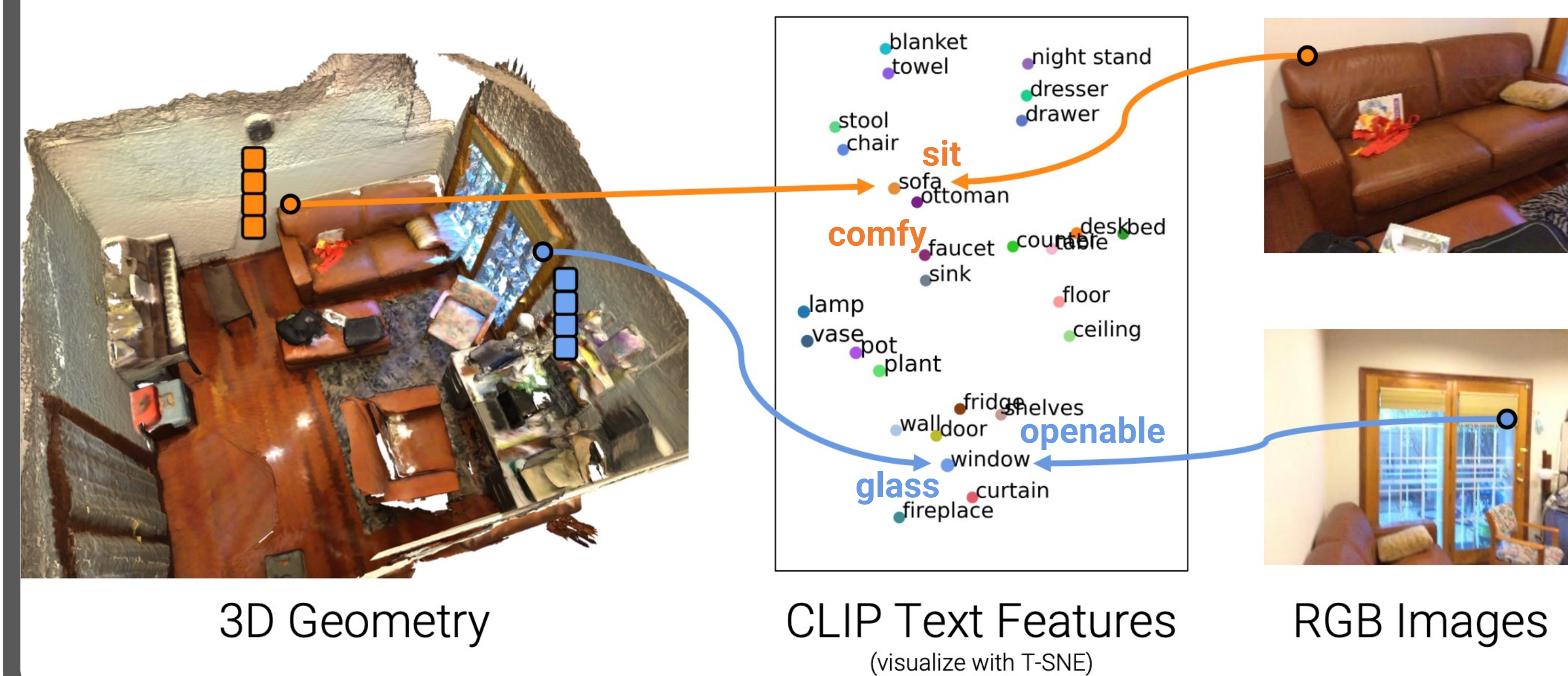


### 1. Introduction

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

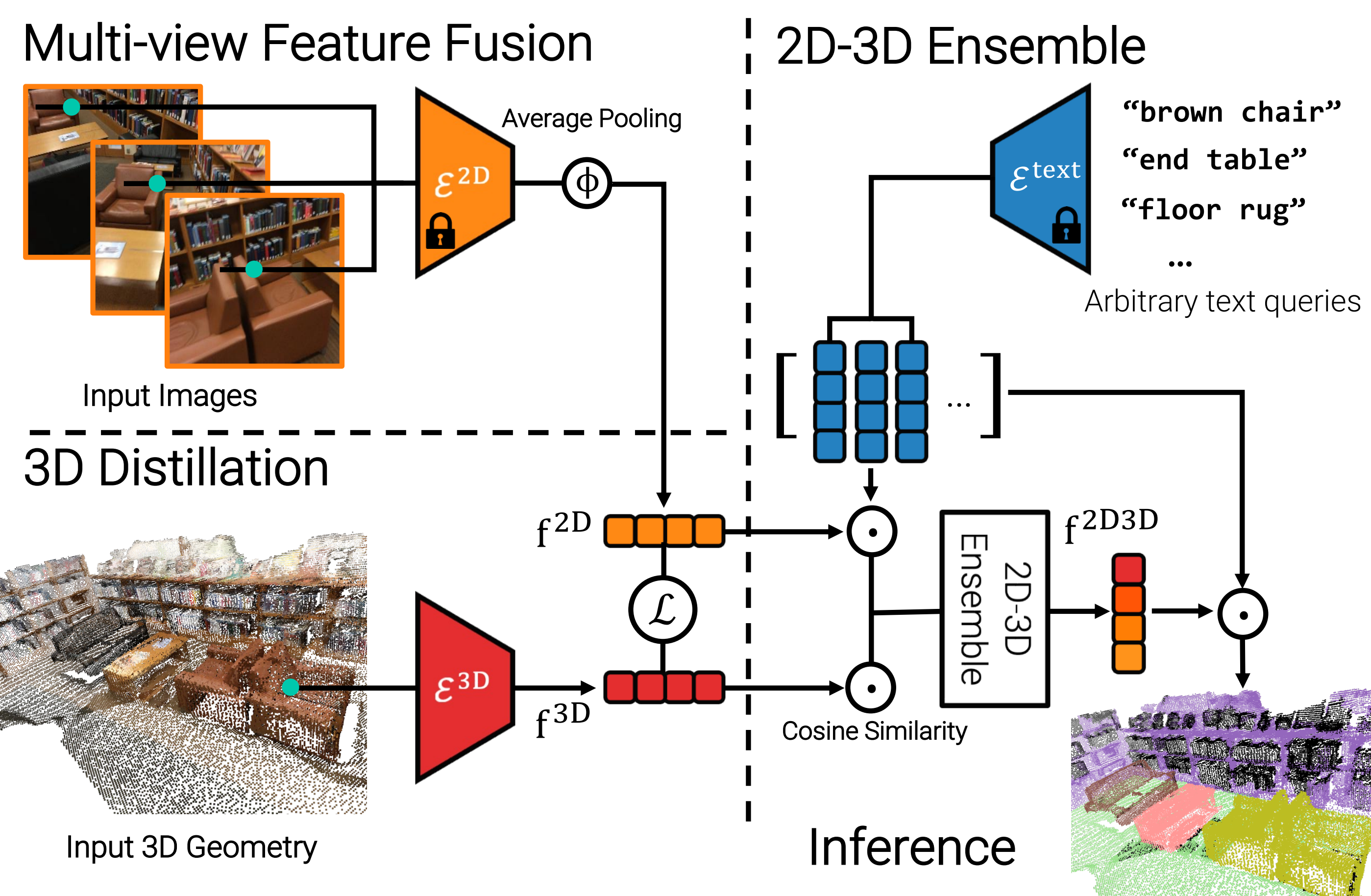
**Goal:** A **zero-shot** approach to perform novel 3D scene understanding tasks **w/o annotation labels**

**Key idea:** Co-embed 3D features with CLIP image features → naturally also with CLIP text features

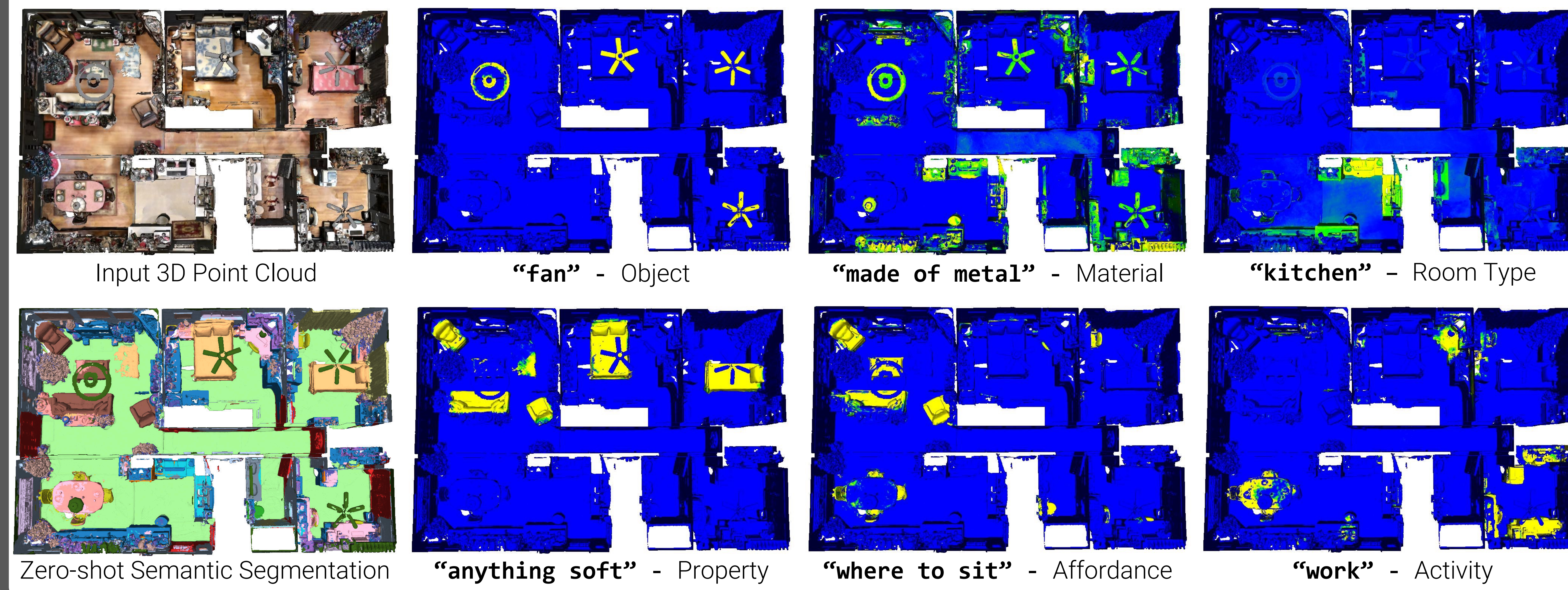


### 2. Method

How to produce text-image-3D co-embedding?

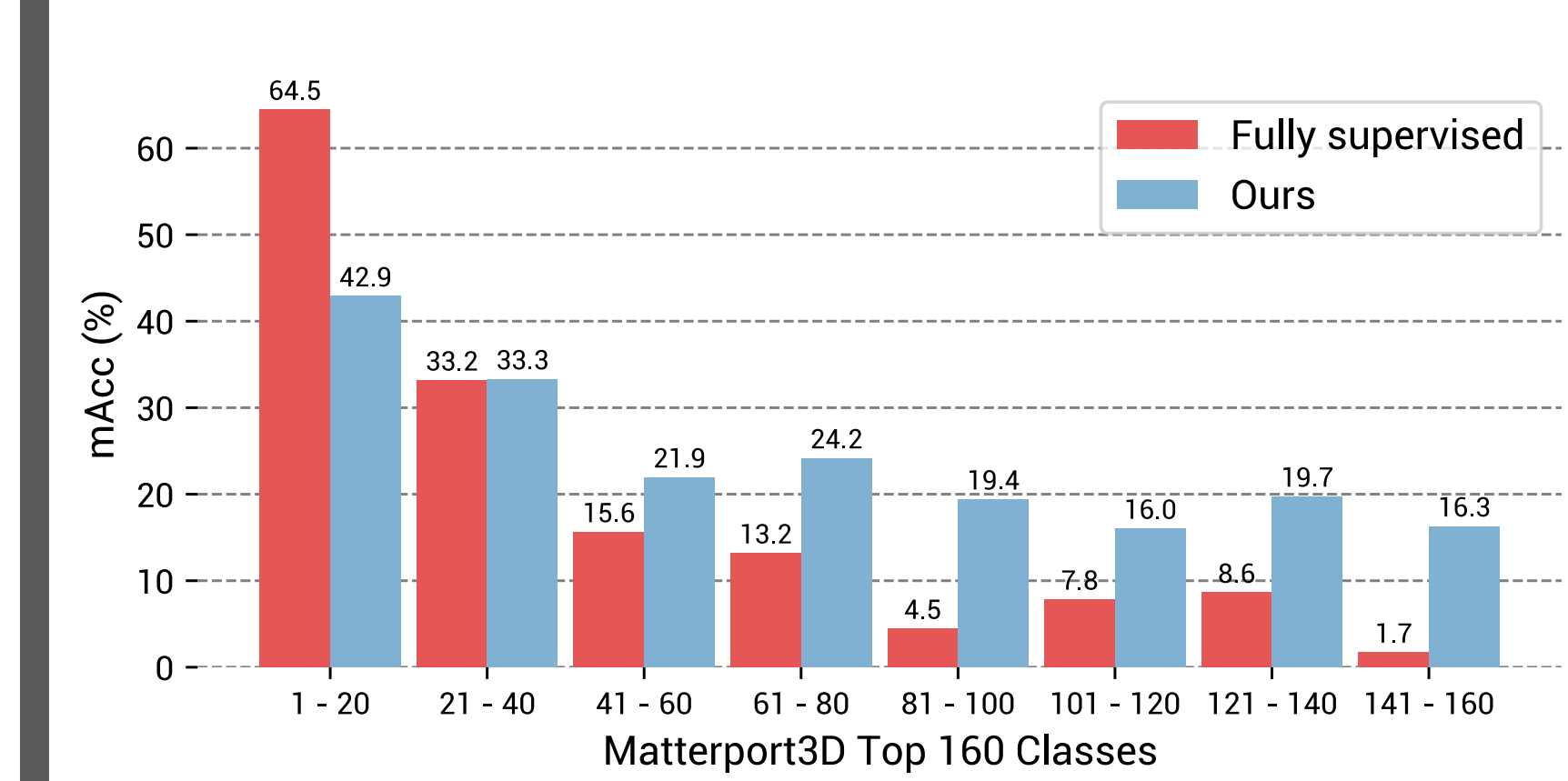


### 3. Zero-shot Open-vocabulary Scene Exploration

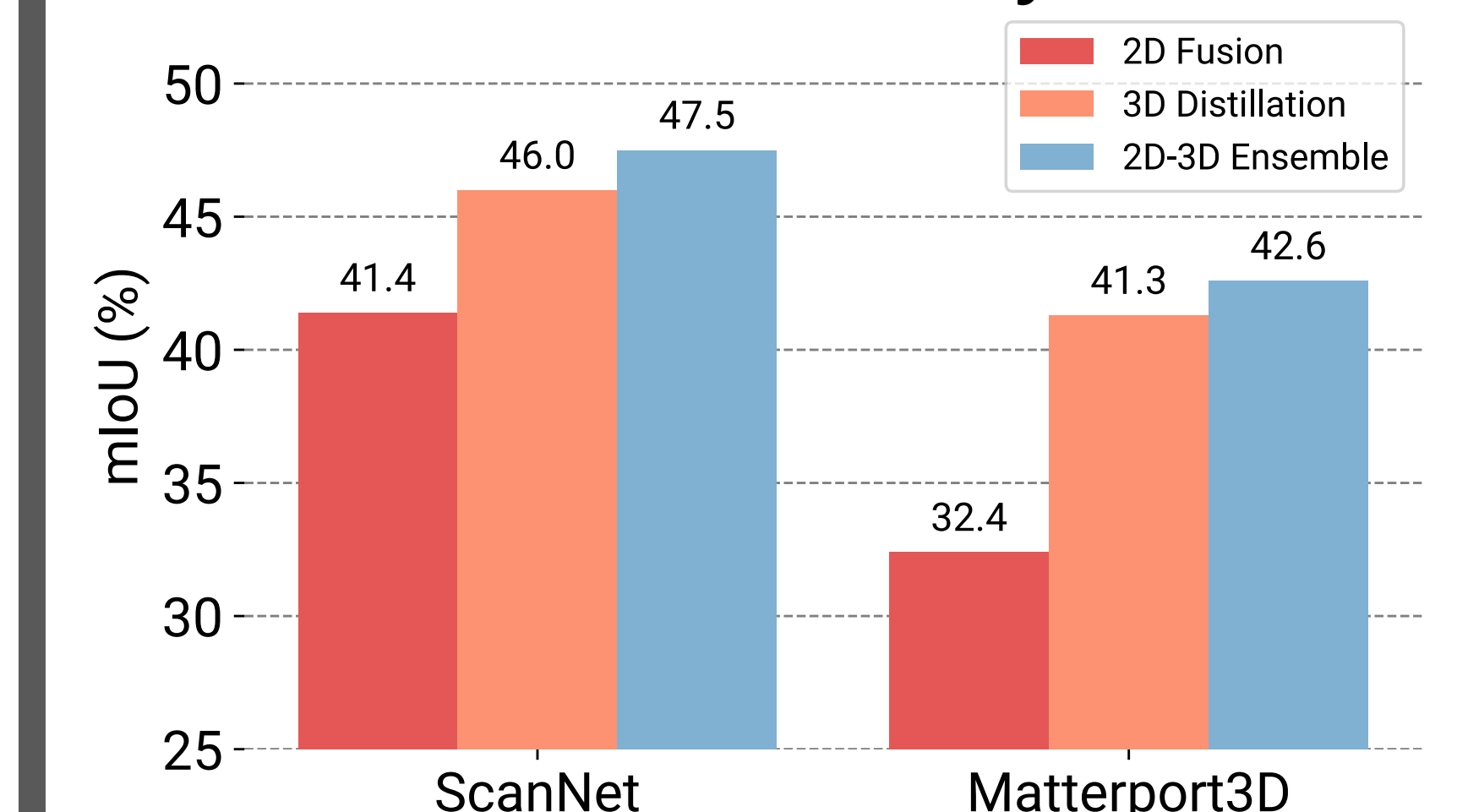


### 5. More Studies

**Robust to Tailed Classes**



**Ablation Study**



### 4. Additional Applications

3D Semantic Segmentation Benchmarks



Rare Object Retrieval

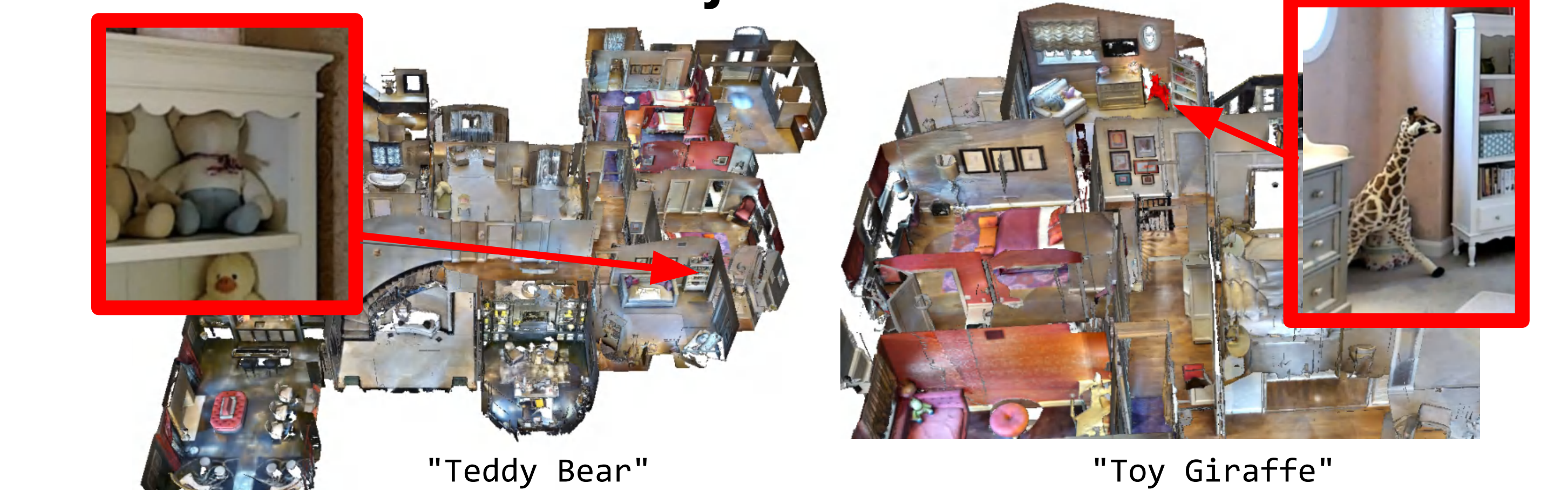


Image-based 3D Object Detection

