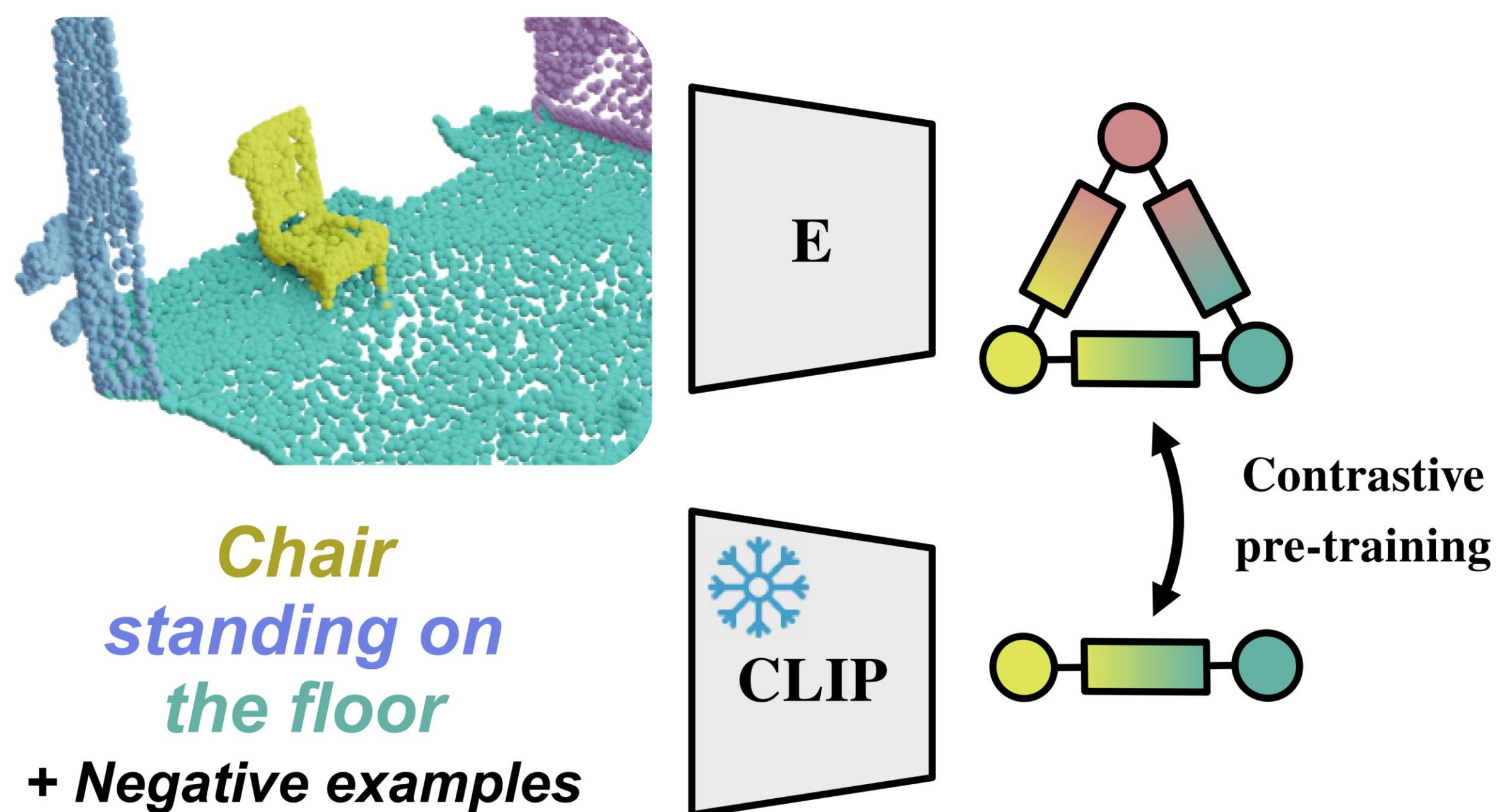


### 1. Introduction

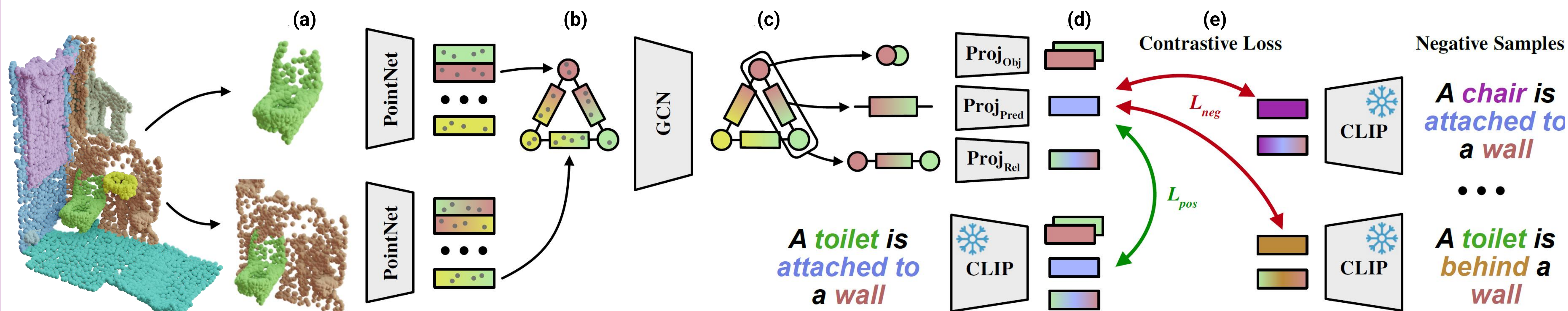
**Key Insight:** Language and Scene Graphs are very similar



+ **Negative examples**

**Goal:** Improved 3D scene graph prediction by language supervised pre-training

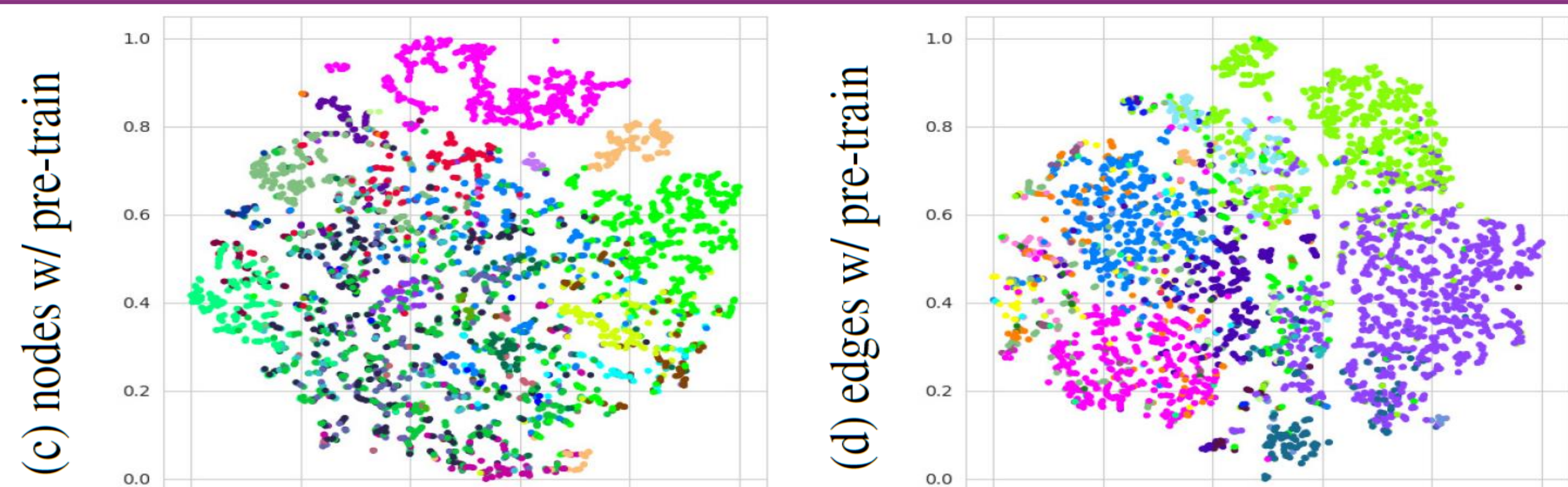
### 2. Method



Construct a latent scene graph from a 3D point cloud. Then align the 3D scene graph embedding with the well-structured language model embedding (CLIP/BERT) using a contrastive loss.

$$\mathcal{L}_{pos} = \sum_{i=1}^N \frac{1}{|K|} \sum_{j \in K} 1 - \cos(f_i, f_{h(j)}^t) \quad \mathcal{L}_{neg} = \sum_{i=1}^N \frac{1}{|M|} \sum_{j \in M} \max(0, \cos(f_i, f_{h(j)}^t) - \tau)$$

### 3. Object & Predicate Embedding



Structured latent space for improved downstream prediction

### 4. Quantitative Results

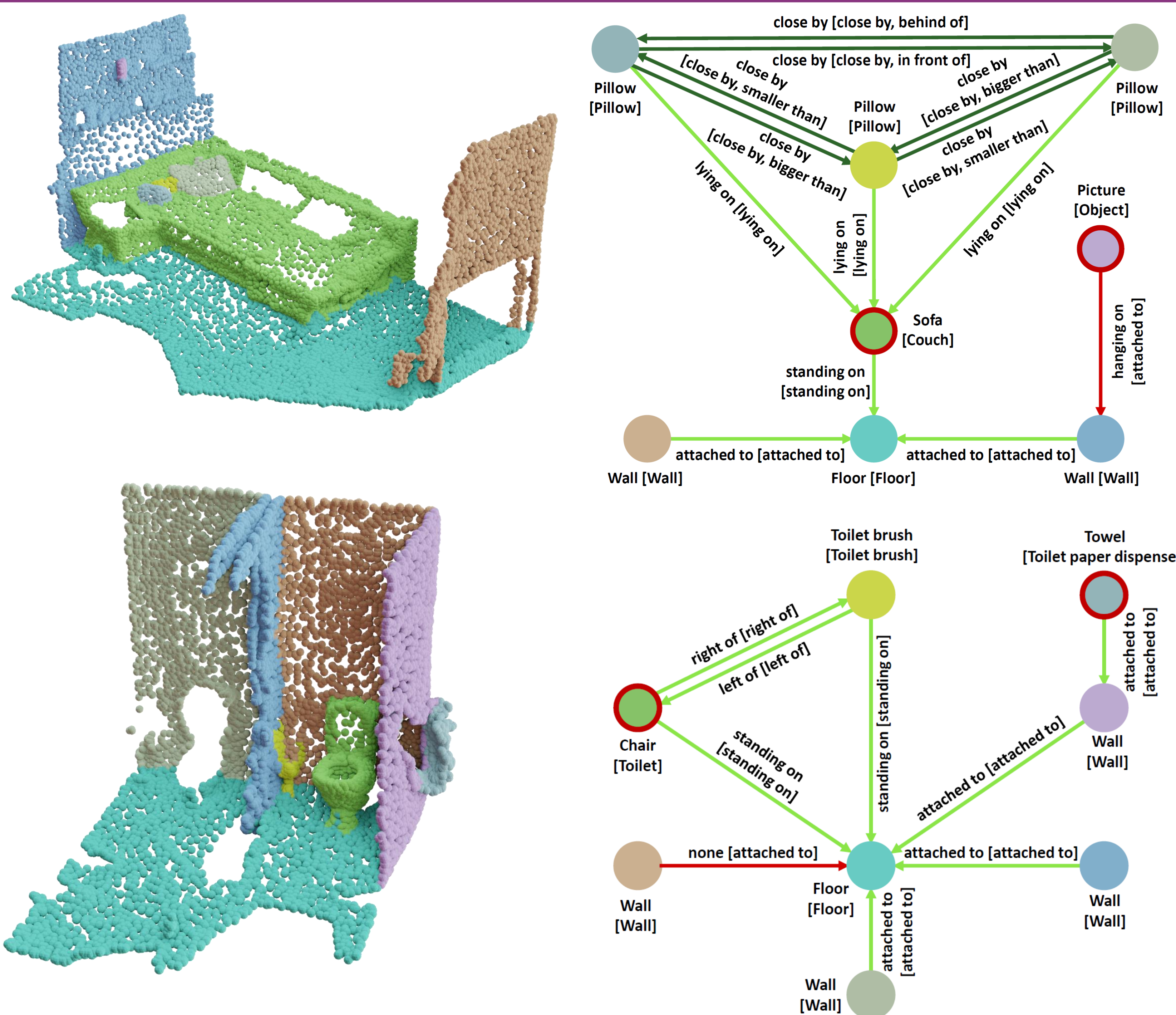
	Object		Predicate		Relationships	
	R@5	R@10	R@3	R@5	R@50	R@100
3DSSG	0.68	0.78	0.89	0.93	0.40	0.66
SGFN	0.70	0.80	<b>0.97</b>	<b>0.99</b>	0.85	0.87
<b>Lang3DSG</b>	<b>0.77</b>	<b>0.84</b>	<b>0.97</b>	<b>0.99</b>	<b>0.87</b>	<b>0.89</b>

	Pre-train		Object		Predicate	
	PCL	SG	R@5	mR@5	R@3	mR@3
Ours (no pre-train)			0.63	0.30	0.94	0.57
STRL	✓		0.75	0.35	0.94	0.50
DepthContrast	✓		0.77	0.36	0.94	0.51
<b>Ours (no GCN)</b>		✓	0.74	0.37	0.94	0.60
<b>Ours</b>		✓	<b>0.77</b>	<b>0.43</b>	<b>0.96</b>	<b>0.67</b>

Scene Graph predictions requires graph-based pre-training

### 5. 3D Semantic Scene Graph Predictions



### 6. Zero-shot room type classification

