

Language-Based Scene Understanding for Mobile Robot

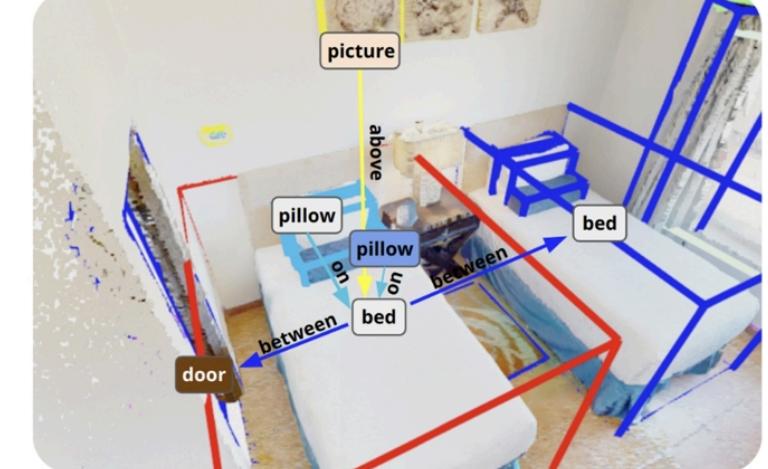
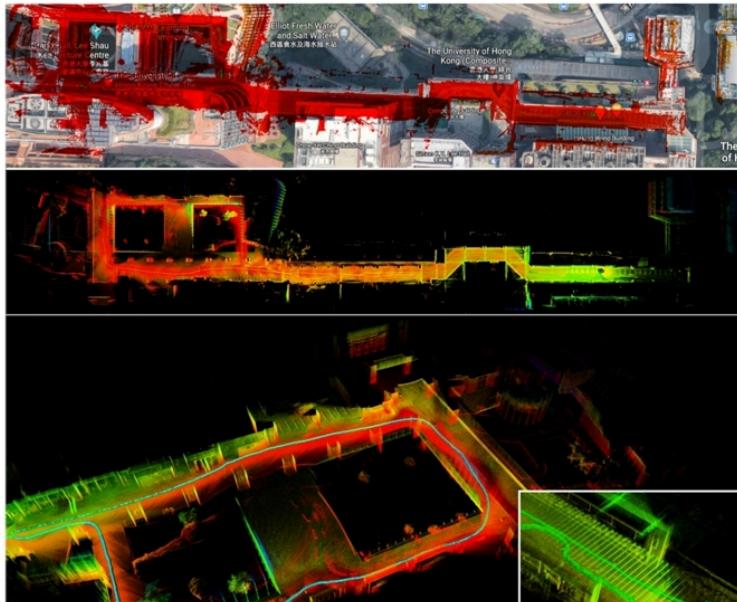
ROAS6800H Seminar

**Yiming ZHU, MPhil's Student
ROAS, Systen Hub
2025/04/01**

Contents

1. Background
2. **SceneScript**: Reconstructing Scenes With An Autoregressive Structured Language Model
What it is, how to train, experiments, why we need?
3. Application

Background



(a) Robust Odometry¹ (b) Complex Scenes² (c) Rich Annotations³

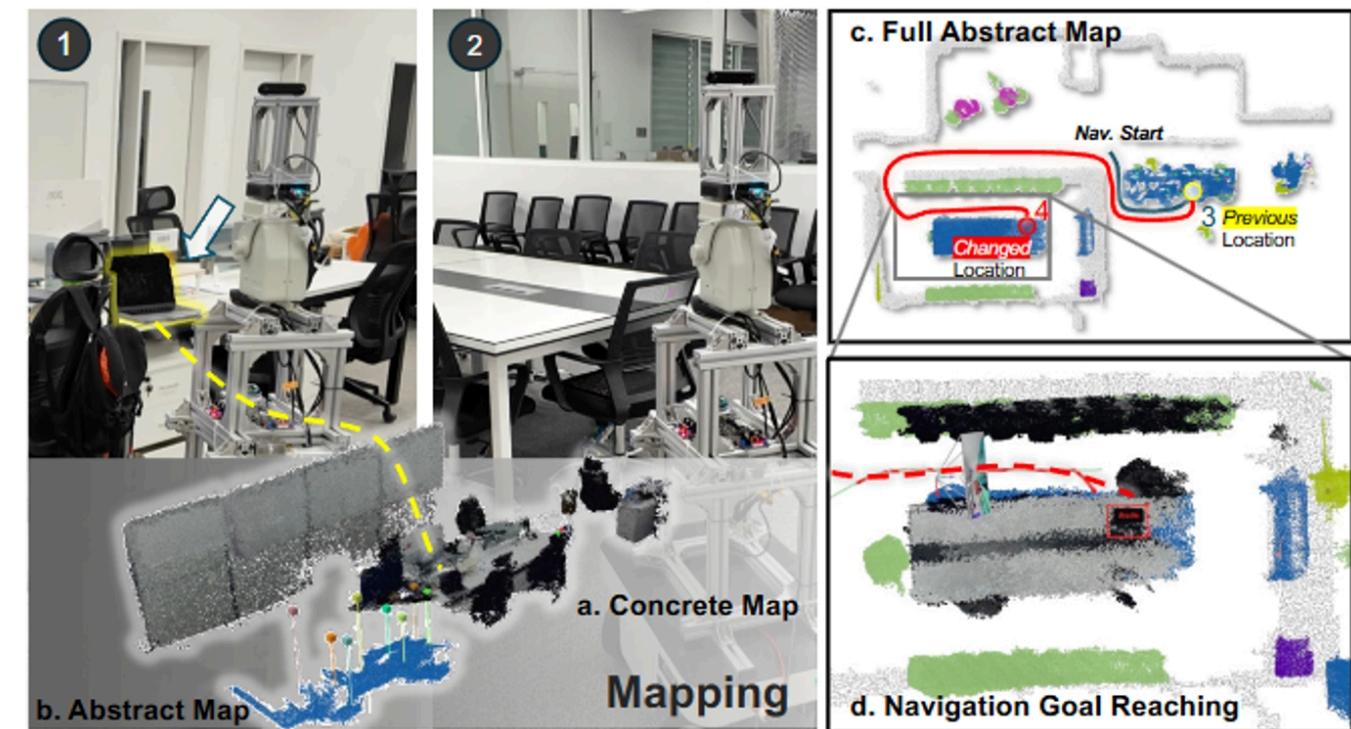
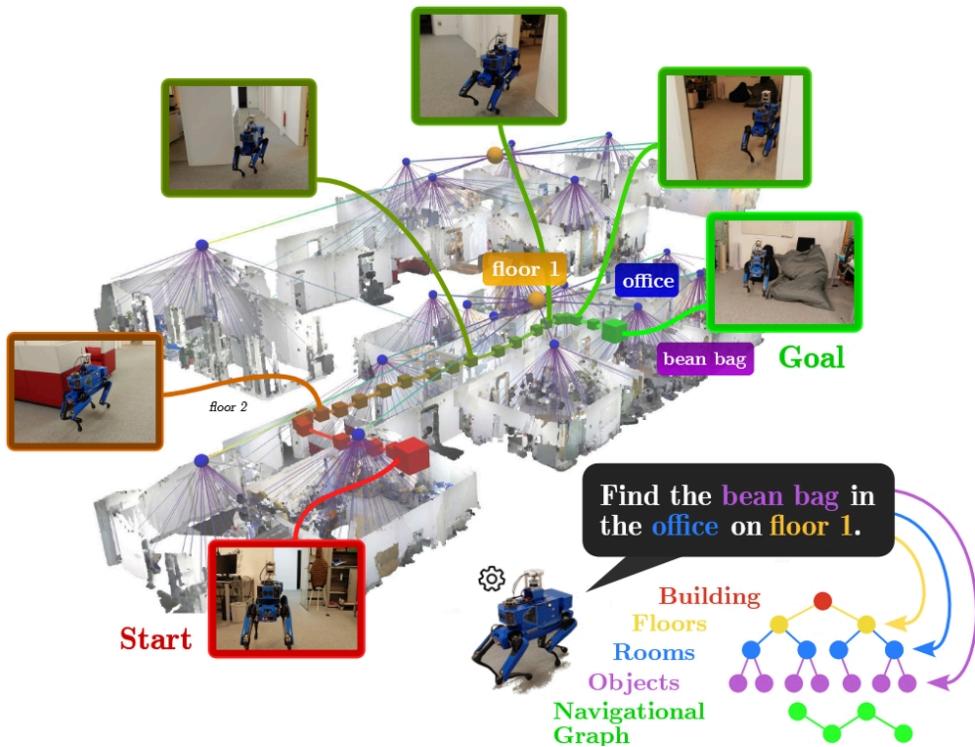
1: Fast-lio <https://ieeexplore.ieee.org/abstract/document/9372856/>

2: Habitat3.0 <https://arxiv.org/abs/2310.13724>

3: VLA-3D <https://arxiv.org/abs/2411.03540>

Background

- Mainstream methods¹ use scene graph to represent a scene.
- Our work DualMap² aims at real-time language navigation.

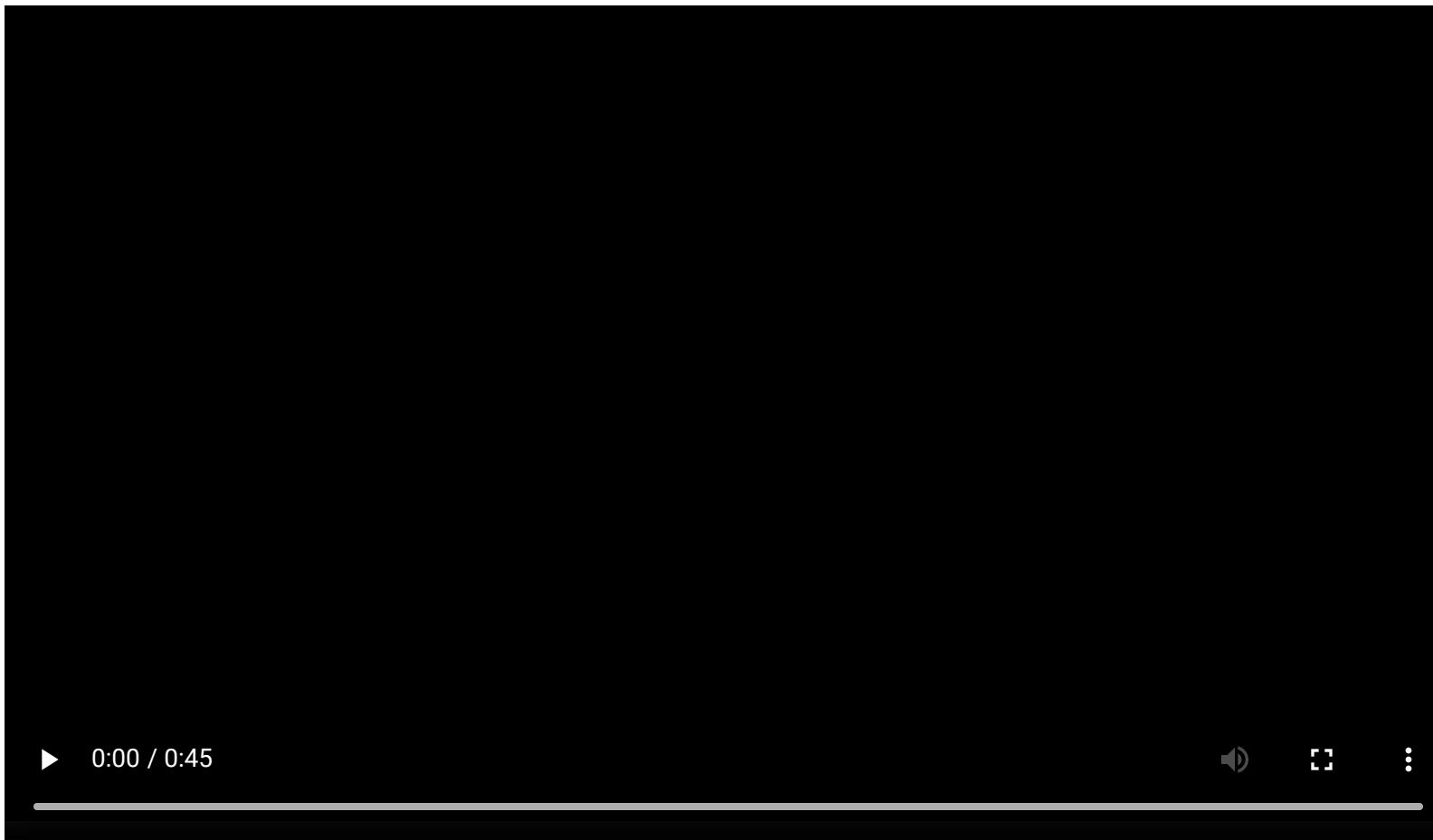


1: HOV-SG <https://arxiv.org/abs/2403.17846>

2: under review

What is SceneScript?

SceneScript (Meta Reality Labs, ECCV24¹) is a novel method of **end-to-end** representing scenes using **language**.

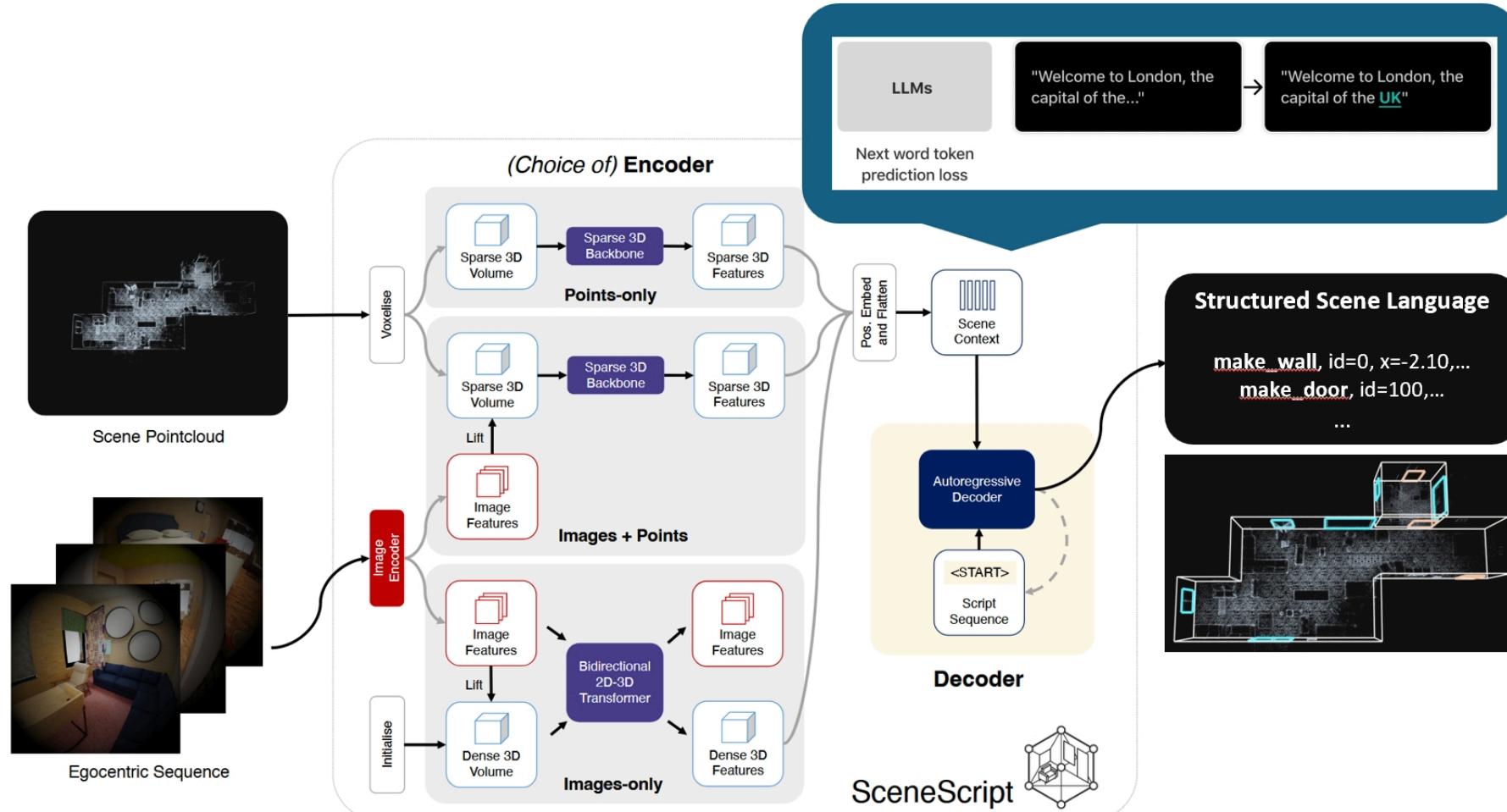


```
1 make_wall,
2   (id = 0),
3   (a_x = -5.518),
4   (a_y = 7.228),
5   (a_z = 0.0),
6   (b_x = -0.91),
7   (b_y = 7.228),
8   (height = 3.089)
9 ...
10 make_door,
11   (id = 1000),
12   (wall_id = 0),
13   (position_x = -3.428),
14   (position_y = 7.228),
15   (size_x = 2.533),
16   (size_y = 1.991)
17 make_window,
18   (id = 2001),
19   (wall_id = 2),
20   (position_x = -3.348),
21   (position_z = 1.514),
22   (size_y = 2.735)
```

Structured scene language

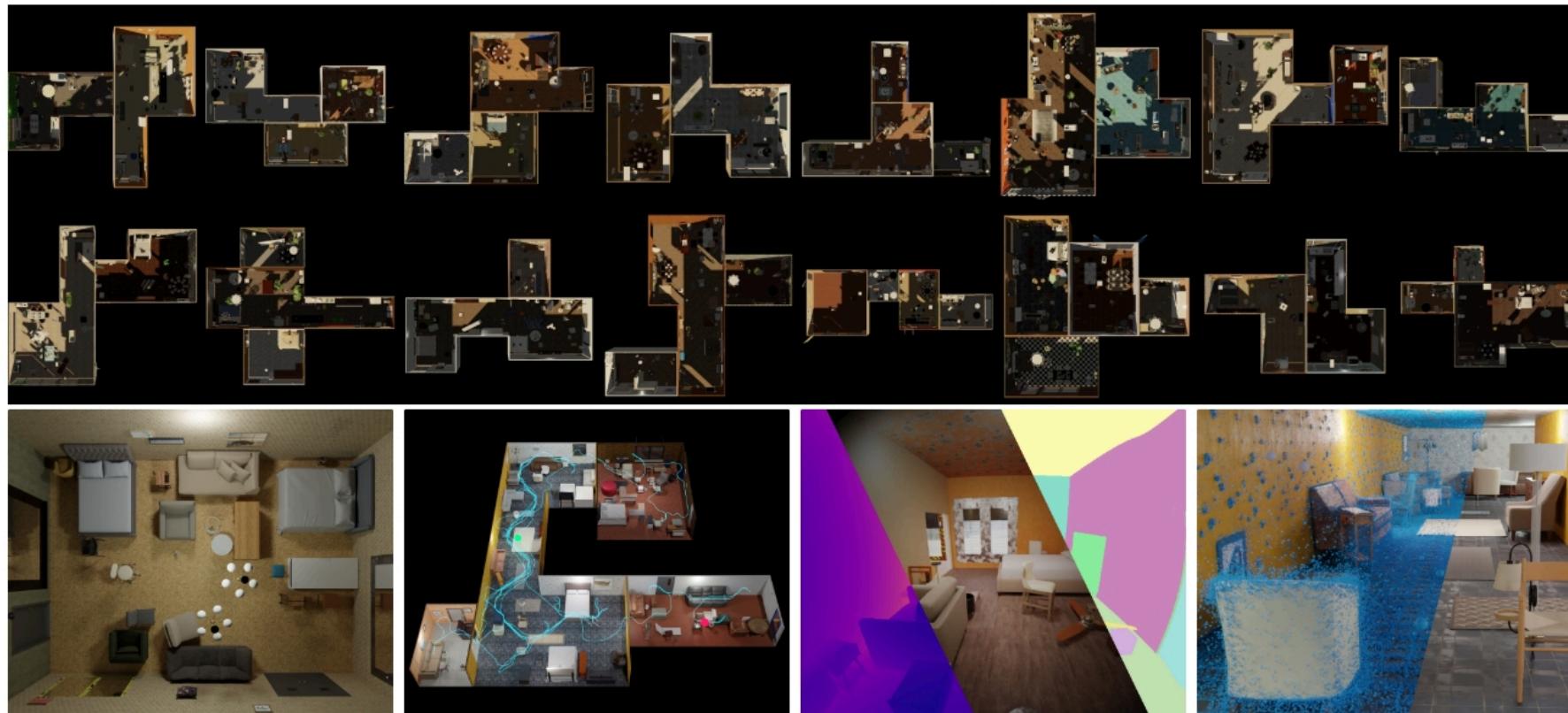
How is SceneScript trained?

Similar to next token prediction in LLMs.



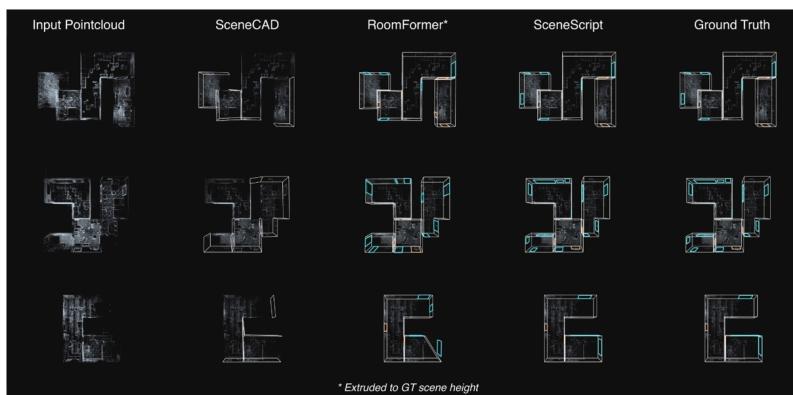
How is SceneScript trained?

- 3D scene data is not as accessible as web texts for LLM training.
- Synthetic dataset [Aria Synthetic Environments](#) with 100K scenes.

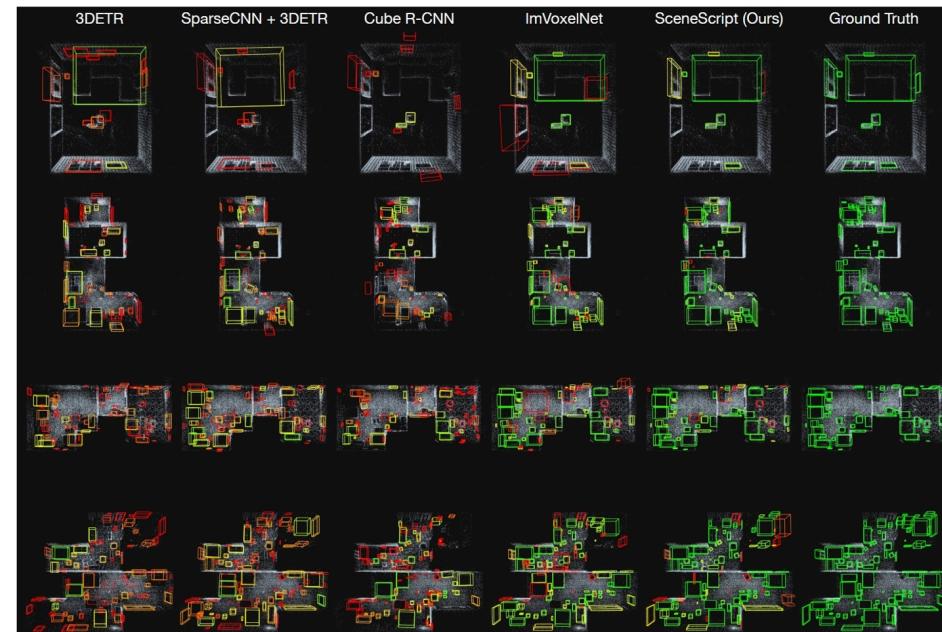


Experimental Results

The authors tested the performance of SceneScripts on layout estimation & object detection.



Method	F1 @5cm				Avg F1			
	mean	wall	door	window	mean	wall	door	window
SceneCAD '20 [2]	-	0.048	-	-	-	0.275	-	-
RoomFormer '23 [52]	0.139	0.159	0.148	0.110	0.464	0.505	0.481	0.407
Ours (Point cloud)	0.848	0.930	0.922	0.692	0.784	0.816	0.811	0.724
Ours (Lifted features)	0.903	0.943	0.959	0.806	0.801	0.818	0.822	0.764
Ours (Image-only)	0.661	0.687	0.798	0.497	0.719	0.727	0.772	0.658



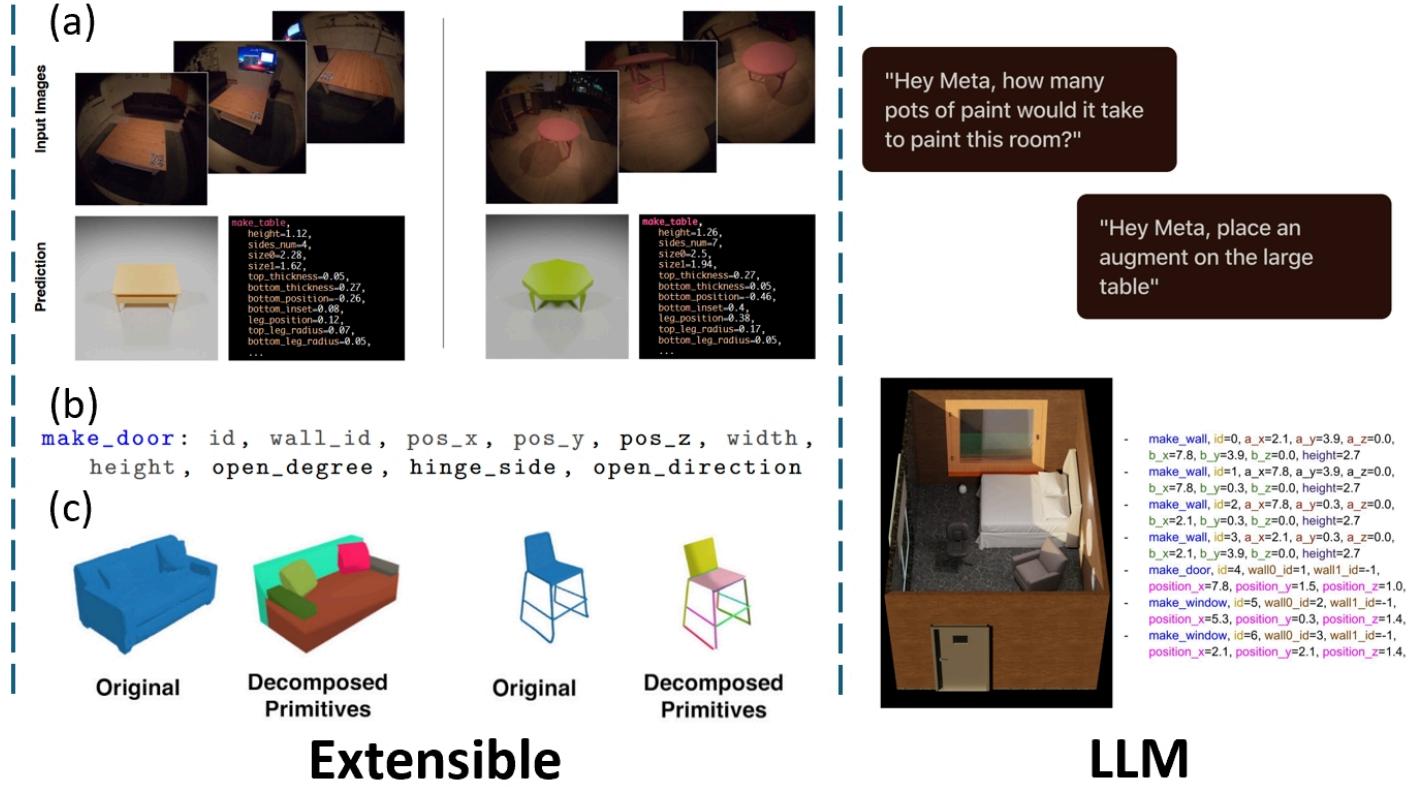
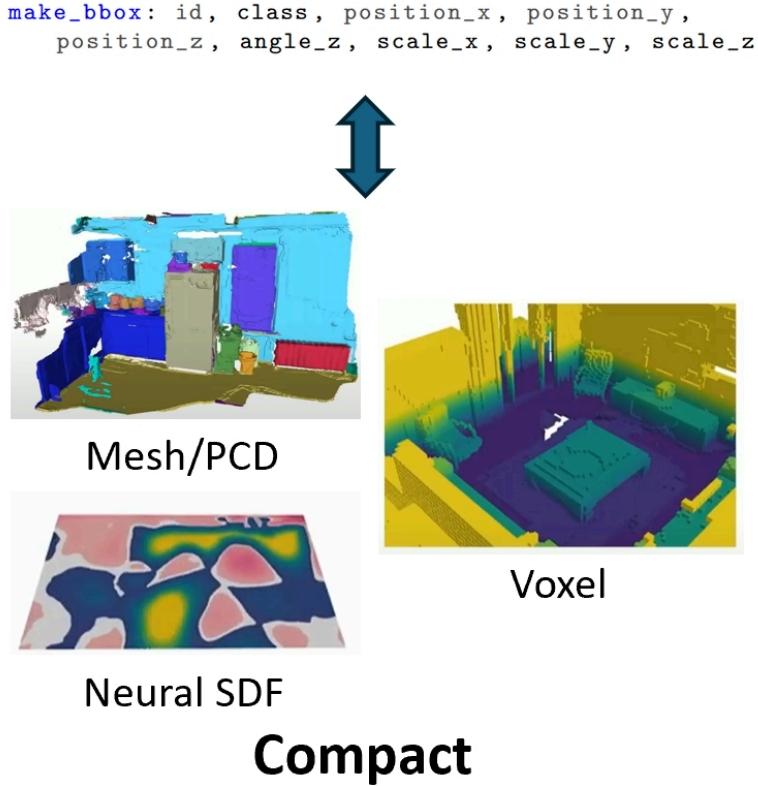
(a) *Aria Synthetic Environments*

(b) *ScanNet* [10]

Method	Input	F1		Method	Input	F1	
		@.25 IoU	@.50 IoU			@.25 IoU	@.50 IoU
3DETR '21 [26]	Points	0.201	0.078	3DETR '21 [26]	Points	0.480	0.349
Cube R-CNN '23 [4]	RGB	0.394	0.228	Cube R-CNN '21 [26]	Points	0.536	0.407
ImVoxelNet '22 [36]	RGB	0.584	0.516	SoftGroup '22 [46]	RGB Points	0.622	0.573
Ours	Points	0.620	0.577	Ours	RGB Points	0.506	0.406

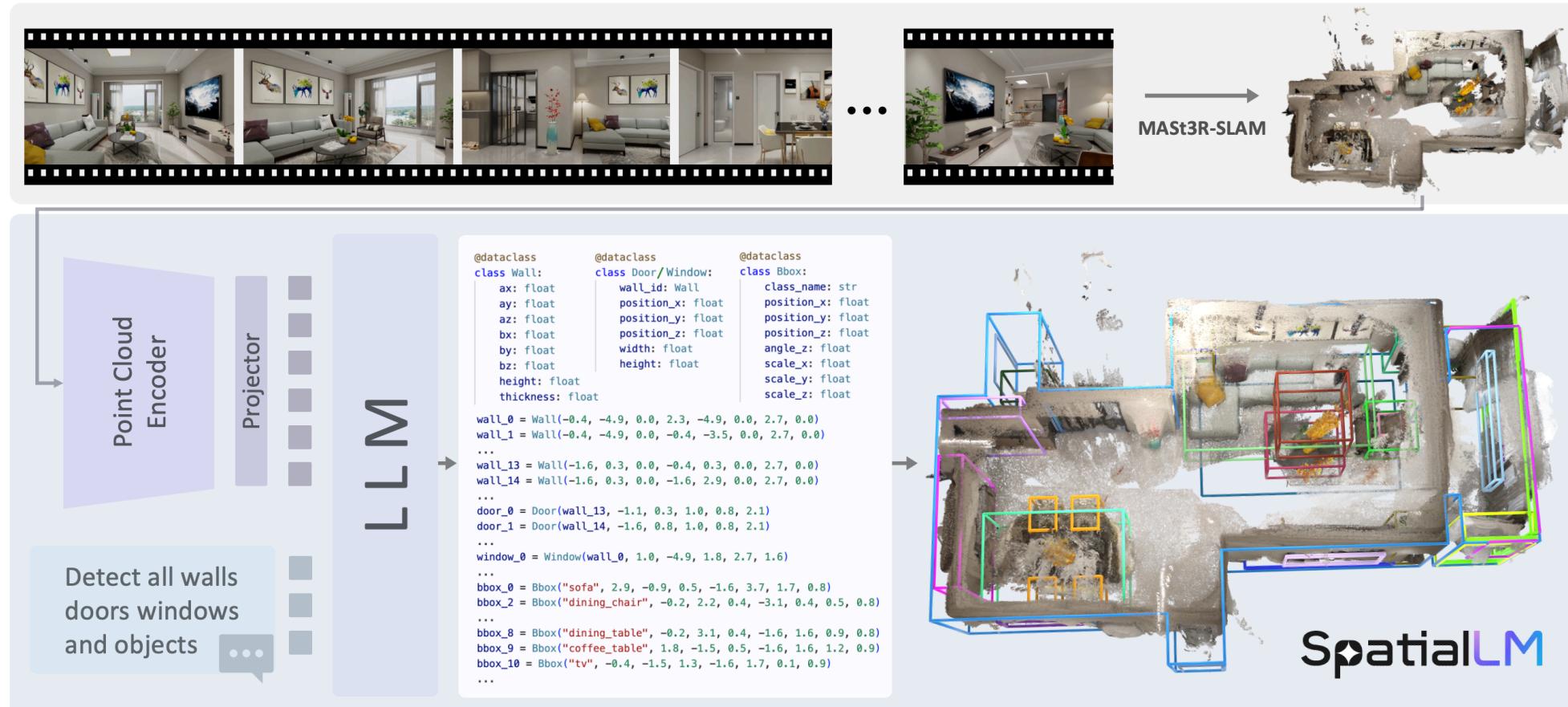
Why SceneScript?

Low memory, extensibility and integration with language models.



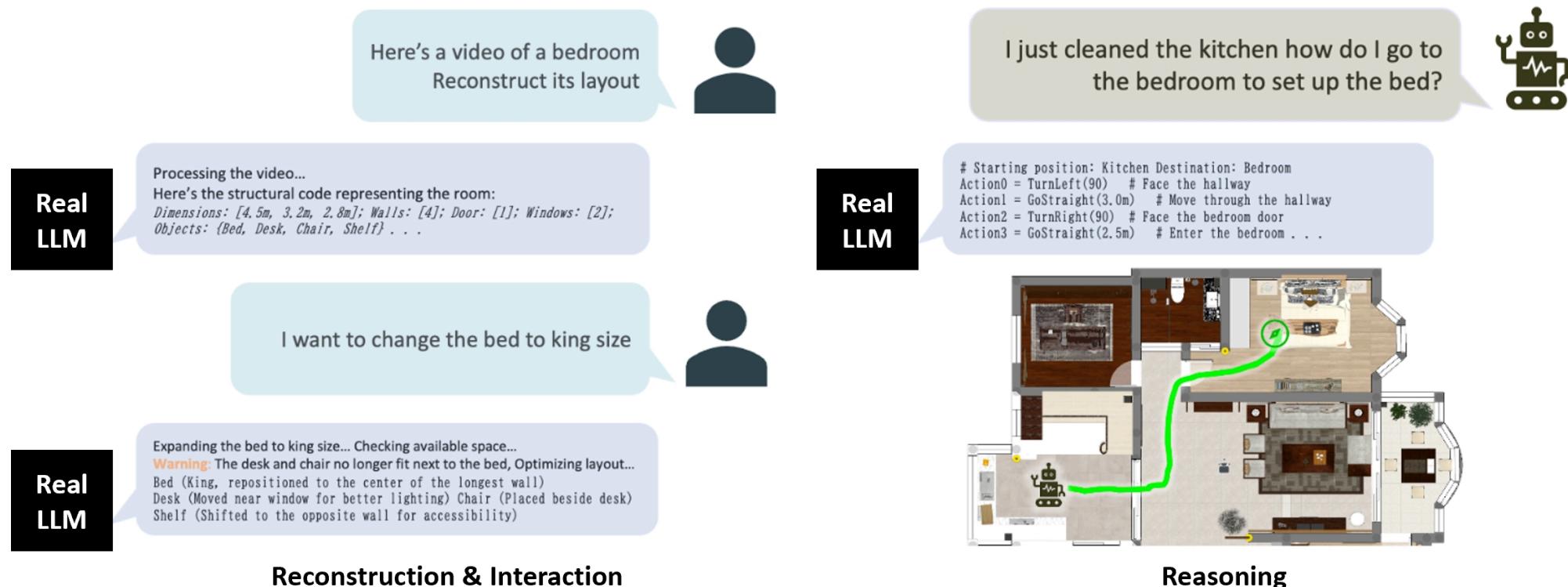
Application

SpatialLM uses the pointcloud encoder of SceneScript. It validates the SceneScript's compatibility with LLM.



Application

- Spatial "LLM" means Large Language Model, which is **over-claimed**.
- But it inspires us about future directions.



Block Test

This is a test for blocks.

Key Point

The next generation of robots will need to build **metric-semantic** representations of environments, enriched with **inter-object relations** and **human-in-the-loop** capabilities.

They looks all set.

Thank you!