

World2Minecraft: Occupancy-Driven Simulated Scenes Construction



Lechao Zhang¹, Haoran Xu¹, Jingyu Gong¹, Xuhong Wang², Yuan Xie^{1,3}, Xin Tan^{1,2,*}

¹School of Computer Science and Technology, East China Normal University

²Shanghai Artificial Intelligence Laboratory

³Shanghai Innovation Institute

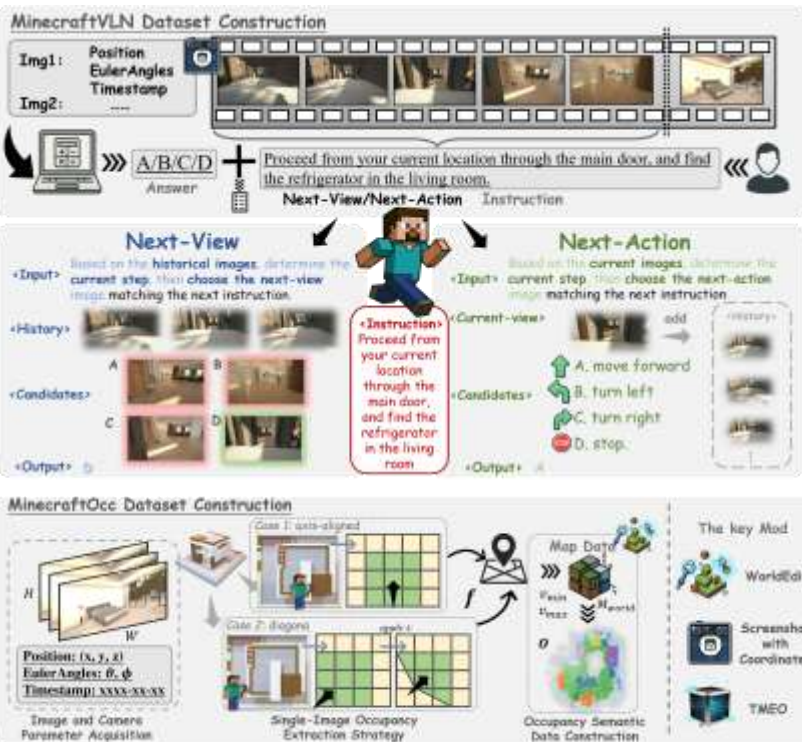


ICLR

Motivation

Embodied AI needs worlds that are both **realistic** and **editable**. Existing methods usually offer **only one side**: real-scene simulators are hard to edit, while Minecraft is editable but far from reality. Moreover, reliable transfer relies on 3D semantic occupancy prediction, which is limited by costly and scarce data. This motivates us to **bridge reality and Minecraft**.

Dataset Construction



Method

Framework of World2Minecraft: real-world scenes are reconstructed into Minecraft via occupancy prediction, enabling VLN tasks such as Next-View and Next-Action in the generated environments.



Overview of World2Minecraft: multi-view observations are converted into fused 3D occupancy and then translated into executable Minecraft commands for scene reconstruction.

Algorithm 1 World2Minecraft: Reality-to-Virtual Transfer

Require: Input: Image set $\mathcal{I} = \{I_1, \dots, I_N\}$

- 1: Camera intrinsic parameters \mathcal{K}
- 2: Camera extrinsic parameters \mathcal{E}
- 3: Pretrained models \mathcal{F}_{mono} , \mathcal{F}_{emb}

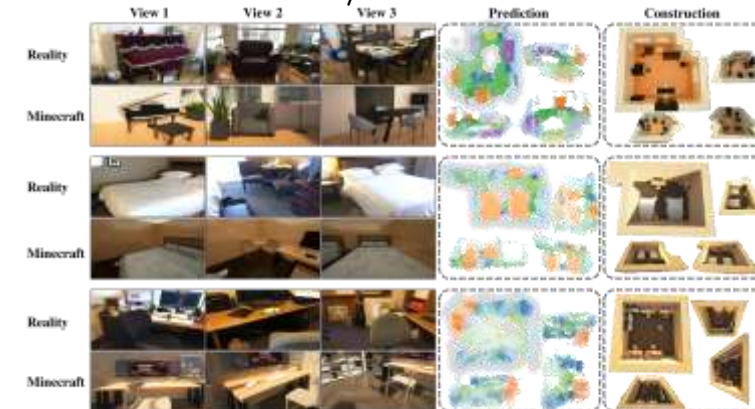
Ensure: Output: Reconstructed Minecraft scene

procedure RECONSTRUCTSCENE($\mathcal{I}, \mathcal{K}, \mathcal{E}, \mathcal{F}_{mono}, \mathcal{F}_{emb}$)

- 4: $\mathcal{O}_{mono} \leftarrow \emptyset$ ▷ Initialize monocular predictions set
- 5: **for** each image $I_i \in \mathcal{I}$ **do** ▷ Process each view
- 6: $\mathcal{O}_{mono} \leftarrow \mathcal{F}_{mono}(I_i, \mathcal{K})$
- 7: $\mathcal{O}_{mono} \leftarrow \mathcal{O}_{mono} \cup \{\mathcal{O}_{mono}^i\}$ ▷ Generate per-view occupancy
- 8: $\hat{\mathcal{O}}_{scene} \leftarrow \mathcal{F}_{embodist}(\mathcal{O}_{mono}, \mathcal{K}, \mathcal{E})$ ▷ Fuse multi-view predictions
- 9: $\mathcal{D} \leftarrow \mathcal{K} * \hat{\mathcal{O}}_{scene}$ ▷ Compute density map via 3D convolution
- 10: $\mathcal{C} \leftarrow \{v \in \hat{\mathcal{O}}_{scene} \mid \mathcal{D}(v) \geq \tau\}$ ▷ Extract centers above threshold τ
- 11: $\mathcal{C}' \leftarrow \text{Cluster}(\mathcal{C}, \eta)$ ▷ Merge centers within distance η
- 12: $\mathcal{M} \leftarrow \text{TranslateToMinecraft}(\mathcal{C}')$ ▷ Generate Minecraft building commands
- 13: ExecuteCommands(\mathcal{M}) ▷ Render scene in Minecraft
- 14: **return** MinecraftScene ▷ Return reconstructed virtual scene

Result

Reality-to-Minecraft reconstruction results. Across different viewpoints, the reconstructed scenes stay consistent with real-world observations



VLN results in reconstructed scene. Guided by Gemini-2.5-Pro, the agent follows instruction "Go to the piano" and successfully reaches the target.



Comparative Experimental Results

Method	OOB ↓	Collision ↓	Semantic ↑	Visual ↑	Complete ↑	Aesthetic ↑
LayoutGPT	0.279	4.5	0.689	5.000	3.856	4.582
I-Design	0.423	0	0.884	6.001	4.734	5.352
LayoutVLM	0	0.9	0.348	3.625	2.270	2.708
World2Minecraft (Ours)	0.024	0.2	0.913	6.145	5.186	6.022

Dataset Composition	Task	Qwen2.5-VL-3B			Qwen2.5-VL-7B		
		No Train	SFT	RFT	No Train	SFT	RFT
Base	Next-View	0.2195	0.5610	0.2927	0.3905	0.5854	0.4390
	Next-Action	0.1943	0.7200	0.6343	0.3829	0.8000	0.6343
Extend	Next-View	0.2261	0.7087	0.3043	0.2913	0.6826	0.6043
	Next-Action	0.3657	0.5437	0.6667	0.3786	0.6019	0.6343
Combined	Next-View	0.2288	0.5609	0.3137	0.2878	0.6642	0.6753
	Next-Action	0.3037	0.4835	0.6570	0.3760	0.6281	0.6219