

---

# LSM-VLM: Long-Short Memory VLM with World-Model for 3D Spatial Reasoning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1       Recent vision-language models (VLMs) have demonstrated strong perceptual and  
2       reasoning capabilities, but their ability to answer 3D spatial questions remains limited,  
3       especially in long-horizon reasoning and cross-view correspondence. Existing  
4       approaches partially address this challenge by introducing structured spatial memory  
5       or world-model-based view synthesis; however, they either fail to capture both  
6       global and local spatial details, or introduce hallucinated evidence. To address these  
7       challenges, we present LSM-VLM, a unified framework that connects perception,  
8       imagination, and reasoning through structured cognitive memory. Specifically, we  
9       construct a long-short memory module that integrates a scene graph and a bird’s-eye-view  
10      (BEV) map, enabling the VLM to capture both object-level relational  
11      information and global spatial layouts. To recover missing spatial evidence, we  
12      leverage a world model to synthesize question-relevant unseen views conditioned  
13      on the current memory state. Since synthesized views may contain unreliable  
14      content, we further introduce a confidence-aware memory update mechanism that  
15      selectively incorporates trustworthy evidence into memory. Experiments show that  
16      LSM-VLM achieves state-of-the-art performance, reaching 69.2% on VSI-Bench  
17      and 63.6% on SQA3D, with consistent gains on long-horizon and cross-view  
18      reasoning tasks such as relative direction reasoning and route planning.<sup>1</sup>

## 19   1 Introduction

20   3D spatial understanding is a fundamental capability for embodied agents to perceive, navigate,  
21   and make decisions in real-world environments [34, 23, 25, 53]. Given an egocentric video and a  
22   question grounded in an embodied scene, an agent must reason not only about objects visible from  
23   the current camera view, but also about the scene layout, object relationships, occluded regions, and  
24   correspondences across different viewpoints. Such reasoning naturally calls for a *cognitive map*: an  
25   internal representation that organizes spatial structures, object locations, and semantic relations into a  
26   form suitable for reasoning and decision-making [36, 53, 18].

27   In recent years, vision-language models (VLMs) have demonstrated strong capabilities in open-  
28   vocabulary perception and visual question answering. However, their spatial reasoning is still  
29   largely driven by image or video tokens, which introduces several potential limitations. First,  
30   when the input video is long, such implicit representations often become insufficient for retaining  
31   all task-relevant evidence and become fragile for 3D spatial VQA [53, 63, 62, 50]. Second, for  
32   questions that require reasoning beyond the current view, such as questions about route planning or  
33   unobserved regions, implicit representations struggle to support effective scene understanding and  
34   reasoning [11, 12, 8, 59].

35   To mitigate long-context limitations, many methods attempt to maintain and update a compact,  
36   structured spatial memory. Some approaches use bird’s-eye-view (BEV) maps as memory [39, 6, 26],

---

<sup>1</sup>An anonymized demo is available at <https://lsm-vlm.github.io/demo/>.

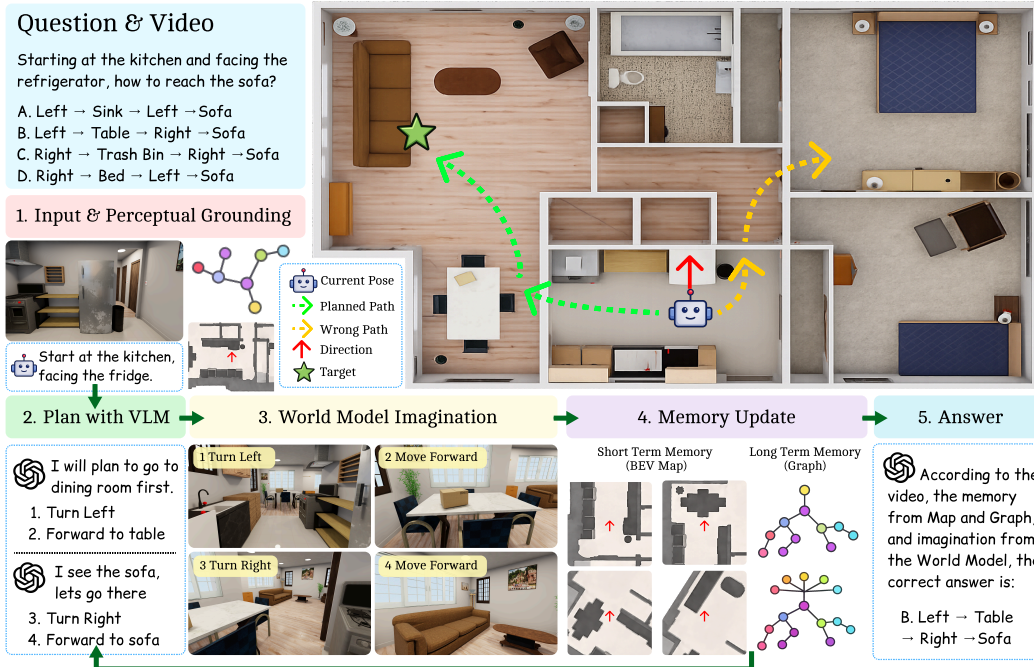


Figure 1: LSM-VLM. Given a video and a question, our framework first grounds the query in the observed scene, then predicts a question-relevant trajectory for exploration. A world model imagines observations along this trajectory, whose evidence is used to selectively update scene graph and BEV map. The enriched spatial memories are finally integrated into the VLM to produce the answer.

37 some use scene graphs [17, 58, 41], and others use latent or recurrent memory tokens [32, 56, 49].  
 38 However, large standalone BEV maps can be difficult for VLMs to interpret semantically, while  
 39 graphs and latent memory tokens are often too abstract to preserve fine-grained local details. For the  
 40 viewpoint-coverage limitation, world models offer a promising way to figure out incomplete visual  
 41 evidence. By generating plausible observations from queried viewpoints, a world model provides a  
 42 form of *imagination*, allowing an agent to inspect what it might see after changing pose [19, 20, 7, 38].  
 43 This capability is particularly useful for spatial VQA, where the answer often depends on evidence  
 44 beyond the current frame. Recent camera-conditioned diffusion and novel-view synthesis models can  
 45 provide such queried views [31, 40, 16, 47, 64], and test-time imagination methods have begun to use  
 46 them for spatial reasoning [55, 27]. However, this strategy is computationally expensive, increases  
 47 visual-context redundancy, and may introduce hallucinated objects or irrelevant details.

48 To address these challenges, we propose LSM-VLM, a unified framework for 3D spatial VQA  
 49 that uses structured spatial memory as an interface between perception, imagination, and language  
 50 reasoning. As shown in Fig. 1, the framework integrates perceptual alignment, dual cognitive memory,  
 51 question-conditioned world-model imagination, and confidence-aware recurrent memory updating  
 52 into a single closed loop. Given an input video, LSM-VLM progressively detects objects and  
 53 constructs a dual cognitive memory, consisting of a long-term scene graph and a short-term BEV  
 54 map with occupancy and semantic channels. The VLM uses this memory to identify informative  
 55 target viewpoints and predict an action sequence to guide the world model. The world model then  
 56 imagines the corresponding observations, providing additional visual evidence beyond the original  
 57 trajectory. A lightweight perception module aligns these imagined observations and converts them  
 58 into structured candidate updates, which are then integrated into memory through a recurrent update  
 59 mechanism with LSTM-like gates [22]. Through this closed loop of planning, imagination, and  
 60 memory updating, LSM-VLM enables the VLM to answer questions based on original observations,  
 61 structured memory, and selectively integrated imagined evidence.

62 This design is well aligned with the requirements of modern spatial reasoning benchmarks. On  
 63 VSI-Bench [53], questions typically require long-horizon evidence accumulation, cross-view corre-  
 64 spondence, relative direction reasoning, and route planning. On SQA3D [34], models need to answer  
 65 situated questions in reconstructed 3D scenes, which requires learning generalizable spatial represen-  
 66 tations. Built on the Qwen3-VL-8B backbone [4], LSM-VLM achieves 69.2% average accuracy

67 on VSI-Bench and 63.6% on SQA3D. Further ablation studies show that performance gains come  
68 from the synergy among dual cognitive memory, world-model imagination, and confidence-aware  
69 recurrent memory integration.

70 Our contributions are summarized as follows:

- 71 • We propose LSM-VLM, a closed-loop framework for 3D spatial VQA that connects  
72 perception, memory, imagination, and language reasoning through a structured spatial-  
73 memory interface, avoiding direct reasoning over long videos or raw imagined images.
- 74 • We design a dual cognitive memory with confidence-aware recurrent updating, which  
75 combines scene graphs for long-term semantic relations and local BEV maps for short-term  
76 spatial reasoning, and selectively integrates evidence imagined by the world model.
- 77 • Experiments on VSI-Bench and SQA3D show promising results, especially on long-horizon  
78 and cross-view reasoning tasks such as relative direction reasoning and route planning.

## 79 **2 Related Work**

### 80 **2.1 VLMs for Spatial VQA**

81 Spatial VQA has evolved from 3D-input benchmarks and 3D-aware models [3, 34, 65, 23] toward  
82 methods that inject 3D awareness into 2D VLMs through spatial QA synthesis, depth modules, or  
83 RGB-D conditioning [11, 12, 8]. Recent video-based studies have further exposed the difficulty  
84 of metric and configurational reasoning for modern MLLMs [53], leading to approaches based on  
85 position-aware tokens, implicit 3D priors, reconstruction-based tuning, and BEV representations [63,  
86 62, 15, 39]. However, existing methods largely reason over the observed video alone, without actively  
87 completing missing evidence from unobserved viewpoints or maintaining a structured spatial memory,  
88 which limits their ability to handle layout, visibility, and cross-view reasoning.

### 89 **2.2 World Models**

90 World models learn predictive representations of environment dynamics [19, 20], with recent genera-  
91 tive video models extending this idea to action- or pose-conditioned rollouts [7, 38, 2]. For spatial  
92 reasoning, a relevant branch is camera-controlled novel view synthesis, where diffusion models  
93 generate object- or scene-level views from specified camera poses [31, 40, 16, 21, 48]. We adopt  
94 Stable Virtual Camera [64], which produces 3D-consistent novel views from arbitrary input images  
95 and target cameras, as our goal-conditioned world model. Unlike recent methods that use imagined  
96 views or distilled world-model priors for spatial reasoning [55, 27, 60, 28], our approach performs  
97 question-conditioned, memory-grounded imagination and consolidates the synthesized evidence into  
98 a dual memory for subsequent reasoning.

### 99 **2.3 Memory Mechanisms for Spatial Reasoning**

100 Spatial reasoning benefits from memory mechanisms that aggregate evidence across time and view-  
101 points into a coherent representation. Prior work has explored cognitive maps for MLLMs [53, 18,  
102 26, 42], structured 3D scene memories [17, 58, 41, 56], and long-horizon video memories such as  
103 working/episodic memory, memory folding, or BEV reconstruction [32, 49, 39]. However, most  
104 methods commit to a single memory modality and do not explicitly address how to absorb noisy  
105 evidence from generative world models without error accumulation. Our approach instead maintains a  
106 persistent dual memory that selectively consolidates simulated spatial evidence for robust cross-view  
107 reasoning.

## 108 **3 Method**

### 109 **3.1 Overview**

110 We propose LSM-VLM, a unified framework that integrates dual structured explicit spatial memory,  
111 question-conditioned world model imagination, and vision-language reasoning. Given a video  
112 sequence  $\mathcal{V} = \{I_1, \dots, I_T\}$  and a question  $q$ , the goal is to output an answer  $a$  that may depend on  
113 scene layout, relative direction, object visibility, and cross-view spatial relations. An overview of our

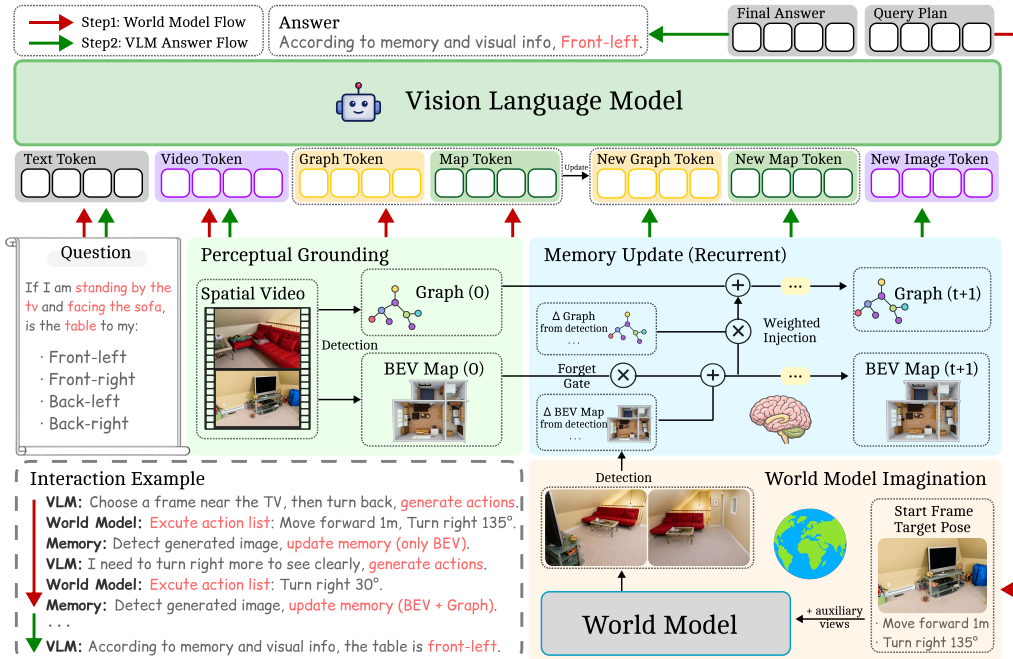


Figure 2: Overview of LSM-VLM: Given a question and an egocentric video, the *Perceptual Grounding* stage initializes a dual cognitive memory consisting of a scene graph  $G_0$  and a BEV map  $B_0$ . The VLM then generates a target pose with an action list, which the *World Model* uses to imagine a new observation. A detection module converts the imagined view into candidate updates ( $\Delta G$ ,  $\Delta B$ ), which are integrated into memory through a forget gate and weighted injection. This updated information will be fed back to the VLM. The loop iterates until the VLM emits the final answer.

114 method is shown in Figure 2. The overall framework consists of two alternating closed loops: the  
 115 World Model Flow for finding spatial evidence and updating memory, and the VLM Answer Flow  
 116 for reading memory and answering. Our framework consists of three stages: *Perceptual Grounding*,  
 117 *World Model Imagination*, and *Memory Update*, as illustrated in Figure 2. In *Perceptual Grounding*,  
 118 we extract relevant information and build a dual cognitive map (3.2). VLM then processes these  
 119 inputs and gives instructions to the world model. In *World Model Imagination*, the world model  
 120 synthesizes additional observations, converts them into a candidate graph, and performs BEV updates.  
 121 (3.3). In *Memory Update*, candidate updates are selectively integrated into the current memory via  
 122 a forget gate and weighted injection, yielding the updated representation for the VLM. (3.4). This  
 123 imagination-and-update loop can be repeated for multiple rounds until the VLM produces the final  
 124 answer. To better align the VLM with the scene graph and BEV map, we further adopt a two-stage  
 125 fine-tuning strategy. The corresponding training details are provided in the experiments section.

### 126 3.2 Perceptual Grounding with Dual Cognitive Memory

127 **Structured spatial memory.** To support long-horizon spatial reasoning, we equip the VLM with  
 128 an effective memory mechanism. Existing methods typically rely on a single intermediate memory,  
 129 such as a BEV map [39, 6, 26], a scene graph [17, 58, 41], or latent memory tokens [32, 56, 49], each  
 130 of which only captures a subset of information needed for spatial reasoning. Graphs capture global  
 131 semantic relations but are often too coarse for fine-grained spatial reasoning, whereas BEV maps  
 132 encode geometry effectively but lack rich object semantics; latent memories are flexible but tend  
 133 to be less interpretable and empirically less stable. Motivated by these observations, LSM-VLM  
 134 maintains dual cognitive memory by using a scene graph as long-term semantic memory and a BEV  
 135 map as short-term spatial memory. After reconstructing with VGGT [43] and semantically parsing  
 136 the input video with SAM3 [10], the system initializes a dual cognitive map  $\mathcal{M}_0 = \{G_0, B_0\}$ , where  
 137  $G_0$  denotes the graph and  $B_0$  denotes the BEV map, and the subscript 0 denotes the memory update  
 138 round in pipeline.

139 **BEV map.** VLMs often struggle to perform stable and accurate spatial reasoning over global BEV  
 140 representations. Therefore, we intentionally keep the BEV representation as simple as possible. As  
 141 shown in 2, the BEV map always uses a local coordinate frame centered on the agent, with the agent  
 142 orientation aligned upward. The BEV map is defined as a two-channel local grid  $B_0 \in \mathbb{R}^{H \times W \times 2}$ ,  
 143 where  $B_0^{\text{occ}} \in \{0, 1\}^{H \times W}$  encodes occupancy and  $B_0^{\text{sem}} \in \{1, \dots, C\}^{H \times W}$  stores semantic labels  
 144 as category indices under the NYU40 [35] taxonomy. This localized design avoids the redundancy  
 145 and reasoning burden introduced by a global map while preserving direction-sensitive geometric  
 146 structure. In this sense, the BEV map mainly serves as a short-term, local, and viewpoint-sensitive  
 147 geometric memory.

148 **Scene graph.** The scene graph  $G_0$  is stored as a structured JSON memory and represents global  
 149 object instances, the agent state, and their spatial relations. We represent the scene graph as  $G_0 =$   
 150  $(\mathcal{O}_0, \mathcal{R}_0)$ , where  $\mathcal{O}_0 = \{o_i\}_{i=1}^{N_0}$  denotes object nodes and  $\mathcal{R}_0 = \{r_{ij}\}$  denotes spatial relations  
 151 between nodes as defined in 3DSSG [51]. Each object node is represented as  $o_i = (c_i, \mathbf{p}_i^w, \mathbf{e}_i^w, s_i, \mathbf{z}_i)$ ,  
 152 where  $c_i$  is the object name,  $\mathbf{p}_i^w = (x_i, y_i)$  is the object center in the world frame,  $\mathbf{e}_i^w = (w_i, h_i, \psi_i)$   
 153 denotes its oriented spatial extent,  $s_i$  is the confidence score, and  $\mathbf{z}_i$  is the semantic encoding under  
 154 the NYU40 taxonomy. This representation provides a compact abstraction of the scene structure  
 155 and facilitates high-level relational reasoning across long temporal horizons. Moreover, the shared  
 156 semantic taxonomy and explicit coordinate attributes align the scene graph with the BEV map,  
 157 forming a unified spatial memory for the VLM.

### 158 3.3 World Model for Spatial Evidence Completion

159 **Target pose prediction.** LSM-VLM augments perception with a question-conditioned world  
 160 imagination module that actively completes missing spatial evidence. At reasoning round  $t$ , the VLM  
 161 reads the question  $q$ , the original observed video  $\mathcal{V}$ , the previously imagined observations  $\tilde{\mathcal{I}}_{<t}$ , where  
 162 the tilde notation denotes imagined content, and the current dual memory  $\mathcal{M}_t = \{G_t, B_t\}$ , and  
 163 predicts a set of informative camera poses together with a target pose:

$$(\mathcal{P}_t, p_t^*) = \text{VLM}_{\theta}^{\text{pose}}(q, \mathcal{V}, \tilde{\mathcal{I}}_{<t}, \mathcal{M}_t), \quad \mathcal{P}_t = \{p_t^0, p_t^1, p_t^2, p_t^3\}.$$

164 Here, each candidate pose  $p_t^j = (x_t^j, y_t^j, \phi_t^j)$  denotes an informative camera pose in the local BEV  
 165 coordinate frame. The primary pose  $p_t^0$  is used as the start pose of the imagined trajectory, while  
 166 the remaining poses  $\{p_t^1, p_t^2, p_t^3\}$  provide auxiliary spatial context, such as nearby or complementary  
 167 viewpoints. The target pose  $p_t^* = (x_t^*, y_t^*, \phi_t^*)$ , also predicted by the VLM, specifies the endpoint  
 168 of the trajectory and indicates where the imagination module should navigate to acquire question-  
 169 relevant missing evidence. Given the current BEV map  $B_t$ , the start pose  $p_t^0$ , and the VLM-predicted  
 170 target pose  $p_t^*$ , we derive a discrete action list

$$\mathcal{A}_t = \{a_t^1, \dots, a_t^{L_t}\} = \text{Plan}(B_t, p_t^{\text{cur}}, p_t^*),$$

171 where each action  $a_t^{\ell}$  is selected from the navigation action space, and  $L_t$  denotes the number of  
 172 planned steps required.  $\text{Plan}(\cdot)$  computes a shortest feasible path on the BEV map and converts it  
 173 into a discrete action sequence. The action list guides the world model to synthesize observations  
 174 along the planned trajectory toward the target viewpoint.

175 **World model imagination and detection.** Conditioned on the observations associated with the  
 176 queried viewpoints and the action list, the world model acts as a goal-conditioned simulator that  
 177 synthesizes a short rollout toward the target pose:

$$\tilde{\mathcal{V}}_t = \mathcal{W}_{\psi}(\mathcal{I}_{\mathcal{P}_t}, \mathcal{A}_t),$$

178 where  $\mathcal{I}_{\mathcal{P}_t} = \{I(p_t^0), I(p_t^1), I(p_t^2), I(p_t^3)\}$  denotes the observations corresponding to the queried  
 179 poses  $\mathcal{P}_t$ , and  $\tilde{\mathcal{V}}_t = \{\tilde{I}_t^1, \dots, \tilde{I}_t^{L_t}\}$  denotes the imagined visual observations along the planned  
 180 trajectory. We then apply the same perception module used in perceptual grounding to the synthesized  
 181 rollout. Specifically, VGGT and SAM3 extract object masks, geometry, camera-aware spatial cues,  
 182 and semantic labels from  $\tilde{\mathcal{V}}_t$ . We denote this perception-and-reconstruction pipeline as  $\Phi_{\text{det}}$ , where  
 183 VGGT provides camera-aware geometric reconstruction and spatial cues, while SAM3 produces  
 184 object-level masks and semantic detections. These predictions are then converted into candidate  
 185 updates for the dual cognitive memory:

$$(\Delta B_t, \Delta G_t, \rho_t) = \Phi_{\text{det}}(\tilde{\mathcal{V}}_t),$$

186 where  $\Delta B_t$  and  $\Delta G_t$  are the proposed updates to the BEV memory and graph memory, respectively,  
 187 and  $\rho_t$  denotes the confidence or reliability of the detected-and-reconstructed evidence. These  
 188 candidate updates are not directly inserted into memory; They are passed to the confidence-aware  
 189 memory consolidation module described below.

### 190 3.4 Confidence-Aware Memory Consolidation and Cross-View Reasoning

191 Although world-model generated observations provide useful complementary evidence, they may  
 192 also introduce synthesis artifacts. Therefore, we draw inspiration from the gating mechanism in  
 193 LSTM and design a dual confidence-aware memory consolidation module.

194 **Memory update.** Given the current memory  $\mathcal{M}_t = \{G_t, B_t\}$  and candidate updates  
 195  $(\Delta G_t, \Delta B_t, \rho_t)$  extracted from the imagined rollout, we update the BEV map and the scene graph  
 196 asynchronously. For the BEV map, the forget operation is applied only to the previous BEV memory.  
 197 It is independent of the current candidate update  $\Delta B_t$ . Specifically, we update the BEV map as

$$B_{t+1} = f_B(B_t, \rho_{t-1}) + i_B(\Delta B_t, \rho_t),$$

198 where  $\rho_{t-1}$  denotes the historical confidence associated with the previous BEV memory  $B_t$ . The  
 199 function  $f_B(\cdot)$  selectively preserves reliable cells and forgets low-confidence cells from the last BEV  
 200 map. The current update  $\Delta B_t$  is then confidence-weighted by  $i_B(\cdot)$  before being injected into the  
 201 retained BEV memory. Before fusion, candidate BEV updates are transformed into the current local  
 202 BEV frame using the estimated relative pose between the imagined camera and the memory frame.

203 For the scene graph, since it serves as long-term semantic memory, candidate graph updates are not  
 204 directly committed. We compute a graph input gate only for candidates that have been confirmed by  
 205 the BEV memory:

$$G_{t+1} = G_t + i_G(G_t, B_{t+1}, \Delta G_t, \rho_t).$$

206 Here,  $B_{t+1}$  is not a single-step observation but the consolidated BEV memory after recurrent updates.  
 207 Thus, using  $B_{t+1}$  allows the weighted injection function  $i_G(\cdot)$  to verify candidate graph updates  
 208 against accumulated multi-step spatial evidence before committing them to the scene graph.

209 **Cross-view reasoning.** After memory consolidation, the VLM predicts the final answer by jointly  
 210 attending to the original observation, imagined evidence, and the updated dual memory:

$$a = \text{VLM}_\theta^{\text{ans}}(q, \mathcal{V}, \tilde{\mathcal{I}}_{\leq t}, \mathcal{M}_{t+1}).$$

211 This design allows LSM-VLM to use imagined observations for spatial evidence completion while  
 212 reducing the risk of contaminating long-term memory with hallucinated content. As a result, the  
 213 model can perform robust cross-view spatial reasoning under partial observation.

## 214 4 Experiments

### 215 4.1 Implementation Details

216 **Two-Stage Training Strategy.** We build our framework on Qwen3-VL-8B-Instruct and adopt  
 217 Stable Virtual Camera as the world model for view generation. To incorporate BEV maps and graph,  
 218 we employ a two-stage training strategy trained on ScanNet [13] and ScanNet++ [57] training splits.

219 In Stage 1, we introduce a map encoder and a graph encoder, together with their corresponding  
 220 projectors, to align the newly introduced modalities with the pretrained language space. The map  
 221 encoder is initialized from ViT-B/16 [14], while the graph encoder is initialized from MiniLM-L6-v2  
 222 [46]. During this stage, the original vision tower, language model, and multimodal merger are frozen,  
 223 and only the newly introduced encoders and projectors are optimized.

224 In Stage 2, we continue training from the Stage 1 checkpoint and perform task-level adaptation  
 225 with LoRA [24]. The map and graph encoders, together with their projectors, remain trainable,  
 226 while LoRA is applied to the selected language modeling modules. The visual backbone and token  
 227 embedding layers remain frozen for stable and efficient adaptation. More details are in Appendix B.

Table 1: Main results on VSI-Bench and SQA3D. All values are accuracies (%) for VSI-Bench and Avg. EM-1 for SQA3D. Higher is better, and - means unavailable.

Method	Avg.	Obj. Cnt.	Abs. Dist.	Obj. Size	Room Size	Rel. Dist.	Rel. Dir.	Route Plan	Appr. Order	SQA3D
<i>Proprietary API Models</i>										
GPT-4o [37, 39]	34.0	46.2	5.3	43.8	38.2	37.0	41.3	31.5	28.5	42.0
Claude-3.7-Sonnet [1, 45]	47.0	-	-	-	-	-	-	-	-	-
Gemini-2.5 Pro [59]	52.7	48.2	35.6	71.3	51.7	58.9	42.4	46.8	66.7	-
<i>Open-source General VLMs</i>										
Qwen3-VL-4B [4, 29]	59.3	-	-	-	-	-	-	-	-	49.9
Qwen3-VL-8B [9, 29]	57.9	67.6	47.0	76.3	61.9	58.0	51.0	35.1	66.3	50.2
Qwen3-VL-30B-A3B [44]	59.5	71.5	42.5	<b>77.4</b>	65.2	54.9	58.8	41.8	64.4	-
LLaVA-OneVision-72B [53]	40.2	43.5	23.9	57.6	37.5	42.5	39.9	32.5	44.6	-
LLaVA-NeXT-Video-72B [62]	40.9	48.9	22.8	57.4	35.3	42.4	36.7	35.0	48.6	-
InternVL3-8B [9]	42.1	66.1	34.9	43.6	47.5	48.0	39.3	26.3	31.4	-
InternVL3.5-8B [44]	49.4	61.3	34.4	61.2	55.4	50.4	44.2	34.0	54.4	-
<i>Spatial-Enhanced Models</i>										
VG-LLM-4B [62]	47.3	66.0	37.8	55.2	59.2	44.6	45.6	33.5	36.4	-
VG-LLM-8B [62]	50.7	67.9	37.7	58.6	62.0	46.6	40.7	32.4	59.2	-
Spatial-MLLM-4B [50]	48.4	65.3	34.8	63.1	45.1	41.3	46.2	33.5	46.3	55.9
VLM-3R-7B [15]	60.9	70.2	49.4	69.2	67.1	65.4	80.5	45.4	40.1	60.7
VST-7B [54]	61.2	71.6	43.8	75.5	69.2	60.0	55.6	44.3	69.2	-
VLM <sup>2</sup> -7B [33]	68.8	72.5	<b>59.6</b>	70.8	69.9	<b>69.0</b>	<b>87.8</b>	52.6	68.3	60.4
LSM-VLM w/o Memory	67.1	72.8	52.3	75.6	68.8	64.2	74.1	51.0	<b>78.3</b>	48.9
LSM-VLM (Ours)	<b>69.2</b>	<b>73.1</b>	53.8	75.8	<b>70.0</b>	65.2	82.8	<b>54.9</b>	<b>78.3</b>	<b>63.6</b>

## 228 4.2 General Spatial Reasoning Benchmark

229 **Benchmark & Metrics** We evaluate general spatial reasoning on two complementary benchmarks,  
 230 **VSI-Bench** [53] and **SQA3D** [34]. VSI-Bench is a video-based benchmark designed for long-horizon  
 231 visual-spatial intelligence from egocentric indoor trajectories. It contains over 5K questions derived  
 232 from ScanNet, ScanNet++ and ARKitScenes [5], and covers eight sub-tasks spanning measurement  
 233 estimation and spatiotemporal reasoning. For multiple-choice or categorical questions, performance is  
 234 measured by answer accuracy, while numerical estimation tasks are evaluated with the mean relative  
 235 accuracy metric; the final Avg. score is the average over all eight sub-tasks.

236 We also report results on SQA3D [34], a standard 3D situated question answering benchmark built  
 237 on ScanNet scenes. We follow the common SQA3D setting and use the average exact-match score on  
 238 the test split. Overall, VSI-Bench evaluates long-horizon spatial reasoning from video observations,  
 239 while SQA3D tests whether the learned spatial representations generalize to standard 3D QA.

240 **Baselines** We compare LSM-VLM against three groups of baselines in Table 1. The first group  
 241 contains **proprietary API models**, including GPT-4o [37], Claude-3.7-Sonnet [1], and Gemini-2.5  
 242 Pro, which represent strong closed-source multimodal systems with large-scale video and image  
 243 understanding capabilities. The second group contains **open-source general VLMs**, including the  
 244 Qwen3-VL series [52], LLaVA-OneVision [30], LLaVA-NeXT-Video [61] and InternVL3/3.5 [45].  
 245 These models are not explicitly designed for spatial reasoning, but they provide strong and widely  
 246 used general-purpose video-language baselines. The third group contains **spatial-enhanced models**,  
 247 such as VG-LLM [62], Spatial-MLLM [50], VLM-3R [15], VST [54], and VLM<sup>2</sup> [33], which inject  
 248 3D priors, reconstruction signals or external memory mechanisms to improve spatial understanding.  
 249 This comparison allows us to assess whether the proposed workflow of dual memory mechanism and  
 250 world-model imagination remains competitive against recent methods.

251 **Results Analysis** Table 1 shows that LSM-VLM achieves the best overall performance on both  
 252 benchmarks. On VSI-Bench, our full model reaches **69.2** average accuracy, outperforming the  
 253 strongest previous result, VLM<sup>2</sup>-7B, by **0.4** points and exceeding the strongest open-source general  
 254 VLMs by a large margin. In particular, compared with the same Qwen3-VL family backbone, our  
 255 method improves over Qwen3-VL-8B from 57.9 to 69.2 (**+11.3**), indicating that the gain does not  
 256 come from model scale alone, but from explicitly introducing structured spatial representations. We  
 257 also obtain the best result on SQA3D, reaching **63.6**, which is **+2.9** above VLM-3R and **+3.2** above  
 258 VLM<sup>2</sup>. This cross-benchmark improvement suggests that our approach learns transferable spatial  
 259 reasoning abilities rather than overfitting to a single video benchmark.

260 Looking into the sub-task breakdown, the main advantage of LSM-VLM lies in **long-horizon**  
 261 **relational and temporal reasoning**. Our model achieves the best performance on *route planning*

BEV	Graph	World Model	Mem.	Avg.	Obj Cnt	Abs Dist	Obj Size	Room Size	Rel Dist	Rel Dir	Route Plan	App. Order
×	×	×	×	61.5	68.8	51.4	76.3	66.2	63.8	55.4	36.1	74.3
✓	×	×	×	61.6	70.0	51.8	75.6	66.9	62.7	56.3	33.5	75.7
×	✓	×	×	61.4	70.1	51.7	75.8	66.5	62.5	55.5	32.9	75.8
✓	✓	×	×	64.3	72.8	52.3	75.5	68.7	64.7	64.7	37.1	<b>78.5</b>
×	×	✓	×	62.6	70.4	52.3	75.4	68.8	63.6	54.7	37.1	78.4
✓	✓	✓	×	67.1	72.8	52.3	75.6	68.8	64.2	74.1	51.0	78.3
✓	✓	✓	✓	<b>69.2</b>	<b>73.1</b>	<b>53.8</b>	<b>75.8</b>	<b>70.0</b>	<b>65.2</b>	<b>82.8</b>	<b>54.9</b>	78.3

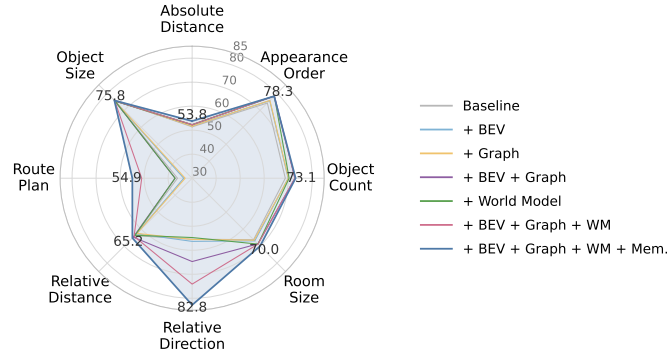
(a) Component ablation.

Stage-1	Stage-2	4B	8B
0	0	59.3	57.9
2K	0	57.7	56.2
2K	2K	59.1	59.8
2K	4K	60.9	61.3
2K	6K	59.2	61.5

(b) Training data scale.

BEV	Graph	8B
1	2	68.9
1	4	69.2
2	2	67.9
2	4	68.5

(c) Memory configuration.



(d) Qualitative visualization.

Figure 3: Ablation study and qualitative analysis on VSI-Bench. (a) Component ablation over BEV memory, graph memory, world model, and recursive memory update. (b) Training ablation under different Stage-1 and Stage-2 training steps. (c) Ablation of BEV and graph memory update intervals using the 8B model. (d) Qualitative visualization of the proposed memory representation.

(54.9) and *appearance order* (78.3). It also remains highly competitive on *object counting*, *room size*, *relative distance*, and *relative direction*, showing that the combination of BEV grounding, graph structure, and world-model completion helps the model maintain a globally coherent scene representation over time. At the same time, our method is not the best on every metric and we view this pattern as meaningful: our design brings the largest benefits to questions that require integrating evidence across viewpoints, preserving temporal order, and reasoning over scene structure, while fine-grained metric estimation remains a challenging direction for future improvement.

### 4.3 Ablation Studies

To better understand where the gains of LSM-VLM come from, we conduct ablation studies on VSI-Bench from four perspectives: component contribution, VLM training strategy, memory update intervals configuration, and computational efficiency. In particular, we examine the effect of different structured inputs and functional modules, analyze how the two-stage VLM adaptation process affects performance, compare different BEV and graph update settings, and finally quantify the runtime overhead introduced by world-model-based spatial evidence completion.

**Component Ablation** We first study the contribution of each major component in Figure 3(a). Starting from the fine-tuned VLM without any additional input, the model reaches 61.5 average accuracy on VSI-Bench, showing that fine-tuning itself does not simply bias the model toward the answer distribution. Adding only BEV memory or only graph memory brings almost no gain, whereas using the two structured memories together improves accuracy by 2.8 points over the VLM-only baseline, confirming that the two memories are complementary rather than redundant.

We then evaluate world-model imagination. Using the world model alone increases the average score from 61.5 to 62.6, suggesting that synthesized views already provide useful supplementary evidence even without explicit memory. However, the gain becomes much larger when imagination is grounded by structured memory: combining BEV, graph, and world model lifts the score to 67.1, with strong improvements on relative direction (74.1) and route planning (51.0). Finally, enabling recursive memory update yields the full model at 69.2, adding another 2.1 points and further improving relative distance, relative direction, and route planning. This result shows that imagined observations are most effective when converted into structured updates and consolidated back into persistent memory.

**Training Strategy Ablation for VLM.** Figure 3(b) analyzes the two-stage adaptation strategy for the VLM. Without any additional training, the 4B and 8B backbones achieve 59.3 and 57.9,

292 respectively. Applying only Stage 1 instead lowers the scores, since new tokens disrupt the original  
 293 concatenation and interpretation. Once Stage 2 tuning is introduced, performance consistently  
 294 recovers and then surpasses the original backbone. For the 8B model, extending Stage 2 to 6K steps  
 295 gives the best result of 61.5. In contrast, the smaller 4B model saturates earlier and is more sensitive  
 296 to overfitting. Overall, these results validate the necessity of the proposed two-stage training strategy.

297 **Memory Update Ablation** Because the BEV map and the scene graph play different short-term  
 298 and long-term memory roles in our framework, their updates are not synchronized. Figure 3(c)  
 299 compares four BEV&graph update interval configurations. Among them, the best result is obtained  
 300 with the asymmetric setting of BEV=1 and Graph=4, which reaches 69.2. The overall pattern  
 301 suggests that the two memories should not be treated identically. This finding supports our design  
 302 choice of maintaining dual memory with asynchronous update behavior instead of collapsing both  
 303 representations into a single shared memory process.

Table 2: Runtime breakdown by component.

Component	Time / Call	Call Times	Total Time
VLM (128 frames)	8.2s	3.4	27.88s
World Model	13.8s	1.7	23.46s
VGGT (1 frame)	0.1s	152.6	15.26s
SAM3 (1 frame)	0.1s	152.6	15.26s
Memory Update	0.1s	2.4	0.24s

304 **Efficiency Analysis** We finally analyze the computational overhead of each component in Table 2.  
 305 When the target view is covered previously, we skip the world model and reuse the nearest observed  
 306 view. The total inference time is about 82.1 seconds per sample on average. Among all modules, the  
 307 VLM contributes the largest cumulative cost, requiring 8.2 seconds per call and 27.88 seconds in  
 308 total because it is invoked 3.4 times on average. The world model is the most expensive single call at  
 309 13.8 seconds, but it is queried 1.7 times. VGGT and SAM3 are both lightweight on a per-call basis  
 310 (0.1 second), but they are applied densely over many frames. The memory update module itself is  
 311 lightweight. Importantly, the world model does not dominate the full runtime and it is affordable.

## 312 5 Conclusion

313 In this paper, we presented LSM-VLM, a closed-loop framework for 3D spatial visual question  
 314 answering motivated by a central challenge raised in the introduction: spatial reasoning in embodied  
 315 scenes cannot rely only on implicit video tokens, especially when evidence is distributed across  
 316 long trajectories and missing viewpoints. Instead of asking the VLM to directly reason over long  
 317 videos or raw generated images, LSM-VLM uses structured spatial memory as the interface between  
 318 perception, imagination, and language reasoning. Specifically, it maintains a dual cognitive map  
 319 with a scene graph for long-term semantic relations and a local BEV map for short-term geometric  
 320 grounding, queries a question-conditioned world model to complete missing evidence, and selectively  
 321 consolidates imagined observations through confidence-aware gated memory updates. This design  
 322 enables the model to reason jointly over what exists in the scene, how entities are spatially arranged,  
 323 and what may be observed from unvisited viewpoints.

324 Experiments on VSI-Bench and SQA3D show that this combination leads to strong and consistent  
 325 gains, improving the Qwen3-VL-8B backbone from 57.9 to 69.2 on VSI-Bench and achieving 63.6  
 326 on SQA3D. Ablation studies further confirm that the improvements come from the synergy among  
 327 dual memory, world-model imagination, and recurrent memory integration, with especially clear  
 328 benefits on long-horizon and cross-view tasks such as relative direction reasoning, route planning,  
 329 and appearance order.

330 **Limitations and Future Work.** Our current study is primarily built on the Qwen3-VL 4B and 8B  
 331 backbones, and we have not yet systematically evaluated how well the proposed framework transfers  
 332 to other VLM families. Extending the method to a broader range of backbones is therefore an  
 333 important direction for future work. In particular, the overall pipeline is compatible with asynchronous  
 334 execution, and future systems could further reduce latency by parallelizing parts of the VLM and  
 335 world-model inference process. Finally, the current memory is used mainly for VLM-side reasoning  
 336 and update, but it is not yet incorporated as an explicit condition for the world model itself. Allowing  
 337 the world model to directly access structured memory may improve the quality, consistency, and task  
 338 relevance of imagined views.

## 339 References

- 340 [1] Claude 3.7 sonnet system card. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:276612236)  
341 276612236.
- 342 [2] Mido Assran, Adrien Bardes, David Fan, Quentin Garrido, Russell Howes, Mojtaba, Komeili,  
343 Matthew Muckley, Ammar Rizvi, Claire Roberts, Koustuv Sinha, Artem Zholus, Sergio Arnaud,  
344 Abha Gejji, Ada Martin, Francois Robert Hogan, Daniel Dugas, Piotr Bojanowski, Vasil  
345 Khalidov, Patrick Labatut, Francisco Massa, Marc Szafraniec, Kapil Krishnakumar, Yong Li,  
346 Xiaodong Ma, Sarath Chandar, Franziska Meier, Yann LeCun, Michael Rabbat, and Nicolas  
347 Ballas. V-jepa 2: Self-supervised video models enable understanding, prediction and planning,  
348 2025. URL <https://arxiv.org/abs/2506.09985>.
- 349 [3] Daichi Azuma, Taiki Miyanishi, Shuhei Kurita, and Motoaki Kawanabe. Scanqa: 3d question  
350 answering for spatial scene understanding. In *2022 IEEE/CVF Conference on Computer Vision*  
351 *and Pattern Recognition (CVPR)*, pages 19107–19117, 2022. doi: 10.1109/CVPR52688.2022.  
352 01854.
- 353 [4] Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao  
354 Deng, Wei Ding, Chang Gao, Chunjiang Ge, Wenbin Ge, Zhifang Guo, Qidong Huang, Jie  
355 Huang, Fei Huang, Binyuan Hui, Shutong Jiang, Zhaohai Li, Mingsheng Li, Mei Li, Kaixin  
356 Li, Zicheng Lin, Junyang Lin, Xuejing Liu, Jiawei Liu, Chenglong Liu, Yang Liu, Dayiheng  
357 Liu, Shixuan Liu, Dunjie Lu, Ruilin Luo, Chenxu Lv, Rui Men, Lingchen Meng, Xuancheng  
358 Ren, Xingzhang Ren, Sibao Song, Yuchong Sun, Jun Tang, Jianhong Tu, Jianqiang Wan, Peng  
359 Wang, Pengfei Wang, Qiuyue Wang, Yuxuan Wang, Tianbao Xie, Yiheng Xu, Haiyang Xu, Jin  
360 Xu, Zhibo Yang, Mingkun Yang, Jianxin Yang, An Yang, Bowen Yu, Fei Zhang, Hang Zhang,  
361 Xi Zhang, Bo Zheng, Humen Zhong, Jingren Zhou, Fan Zhou, Jing Zhou, Yuanzhi Zhu, and  
362 Ke Zhu. Qwen3-vl technical report, 2025. URL <https://arxiv.org/abs/2511.21631>.
- 363 [5] Gilad Baruch, Zhuoyuan Chen, Afshin Dehghan, Tal Dimry, Yuri Feigin, Peter Fu, Thomas  
364 Gebauer, Brandon Joffe, Daniel Kurz, Arik Schwartz, and Elad Shulman. Arkitscenes: A  
365 diverse real-world dataset for 3d indoor scene understanding using mobile rgb-d data, 2022.  
366 URL <https://arxiv.org/abs/2111.08897>.
- 367 [6] Felix Brandstätter, Erik Schütz, Katharina Winter, and Fabian B. Flohr. Bev-llm: Leveraging  
368 multimodal bev maps for scene captioning in autonomous driving. In *2025 IEEE Intelligent*  
369 *Vehicles Symposium (IV)*, pages 345–350, 2025. doi: 10.1109/IV64158.2025.11097781.
- 370 [7] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr,  
371 Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh.  
372 Video generation models as world simulators. 2024. URL [https://openai.com/research/](https://openai.com/research/video-generation-models-as-world-simulators)  
373 [video-generation-models-as-world-simulators](https://openai.com/research/video-generation-models-as-world-simulators).
- 374 [8] Wenxiao Cai, Iaroslav Ponomarenko, Jianhao Yuan, Xiaoqi Li, Wankou Yang, Hao Dong, and  
375 Bo Zhao. Spatialbot: Precise spatial understanding with vision language models. In *2025 IEEE*  
376 *International Conference on Robotics and Automation (ICRA)*, pages 9490–9498, 2025. doi:  
377 10.1109/ICRA55743.2025.11128671.
- 378 [9] Zhongang Cai, Ruisi Wang, Chenyang Gu, Fanyi Pu, Junxiang Xu, Yubo Wang, Wanqi Yin,  
379 Zhitao Yang, Chen Wei, Qingping Sun, Tongxi Zhou, Jiaqi Li, Hui En Pang, Oscar Qian,  
380 Yukun Wei, Zhiqian Lin, Xuanke Shi, Kewang Deng, Xiaoyang Han, Zukai Chen, Xiangyu  
381 Fan, Hanming Deng, Lewei Lu, Liang Pan, Bo Li, Ziwei Liu, Quan Wang, Dahua Lin, and  
382 Lei Yang. Scaling spatial intelligence with multimodal foundation models, 2026. URL <https://arxiv.org/abs/2511.13719>.  
383
- 384 [10] Nicolas Carion, Laura Gustafson, Yuan-Ting Hu, Shoubhik Debnath, Ronghang Hu, Didac Suris,  
385 Chaitanya Ryalí, Kalyan Vasudev Alwala, Haitham Khedr, Andrew Huang, Jie Lei, Tengyu  
386 Ma, Baishan Guo, Arpit Kalla, Markus Marks, Joseph Greer, Meng Wang, Peize Sun, Roman  
387 Rädle, Triantafyllos Afouras, Effrosyni Mavroudi, Katherine Xu, Tsung-Han Wu, Yu Zhou,  
388 Liliane Momeni, Rishi Hazra, Shuangrui Ding, Sagar Vaze, Francois Porcher, Feng Li, Siyuan  
389 Li, Aishwarya Kamath, Ho Kei Cheng, Piotr Dollár, Nikhila Ravi, Kate Saenko, Pengchuan  
390 Zhang, and Christoph Feichtenhofer. Sam 3: Segment anything with concepts, 2026. URL  
391 <https://arxiv.org/abs/2511.16719>.

- 392 [11] Boyuan Chen, Zhuo Xu, Sean Kirmani, Brain Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei  
393 Xia. SpatialVlm: Endowing vision-language models with spatial reasoning capabilities. In  
394 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,  
395 pages 14455–14465, June 2024.
- 396 [12] An-Chieh Cheng, Hongxu Yin, Yang Fu, Qiushan Guo, Ruihan Yang, Jan Kautz, Xiaolong  
397 Wang, and Sifei Liu. SpatialRGPT: Grounded spatial reasoning in vision-language models. In  
398 *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL  
399 <https://openreview.net/forum?id=JKEIYQUSUc>.
- 400 [13] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias  
401 Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes, 2017. URL <https://arxiv.org/abs/1702.04405>.  
402
- 403 [14] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
404 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,  
405 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image  
406 recognition at scale, 2021. URL <https://arxiv.org/abs/2010.11929>.
- 407 [15] Zhiwen Fan, Jian Zhang, Renjie Li, Junge Zhang, Runjin Chen, Hezhen Hu, Kevin Wang,  
408 Huaizhi Qu, Shijie Zhou, Dilin Wang, Zhicheng Yan, Hongyu Xu, Justin Theiss, Tianlong  
409 Chen, Jiachen Li, Zhengzhong Tu, Zhangyang Wang, and Rakesh Ranjan. Vlm-3r: Vision-  
410 language models augmented with instruction-aligned 3d reconstruction, 2026. URL <https://arxiv.org/abs/2505.20279>.  
411
- 412 [16] Ruiqi Gao, Aleksander Hołyński, Philipp Henzler, Arthur Brussee, Ricardo Martin-Brualla,  
413 Pratul Srinivasan, Jonathan T. Barron, and Ben Poole. Cat3d: create anything in 3d with  
414 multi-view diffusion models. In *Proceedings of the 38th International Conference on Neural  
415 Information Processing Systems, NIPS '24*, Red Hook, NY, USA, 2024. Curran Associates Inc.  
416 ISBN 9798331314385.
- 417 [17] Qiao Gu, Ali Kuwajerwala, Sacha Morin, Krishna Murthy Jatavallabhula, Bipasha Sen,  
418 Aditya Agarwal, Corban Rivera, William Paul, Kirsty Ellis, Rama Chellappa, Chuang Gan,  
419 Celso Miguel de Melo, Joshua B. Tenenbaum, Antonio Torralba, Florian Shkurti, and Liam  
420 Paull. Conceptgraphs: Open-vocabulary 3d scene graphs for perception and planning. In *2024  
421 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5021–5028, 2024.  
422 doi: 10.1109/ICRA57147.2024.10610243.
- 423 [18] Chanyoung Gwak, Yoonwoo Jeong, Byungwoo Jeon, Hyunseok Lee, Jinwoo Shin, and Minsu  
424 Cho. Cog3dmap: Multi-view vision-language reasoning with 3d cognitive maps, 2026. URL  
425 <https://arxiv.org/abs/2603.23023>.
- 426 [19] David Ha and Jürgen Schmidhuber. World models. 2018. doi: 10.5281/ZENODO.1207631.  
427 URL <https://zenodo.org/record/1207631>.
- 428 [20] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains  
429 through world models, 2024. URL <https://arxiv.org/abs/2301.04104>.
- 430 [21] Hao He, Yinghao Xu, Yuwei Guo, Gordon Wetzstein, Bo Dai, Hongsheng Li, and Ceyuan  
431 Yang. Cameractrl: Enabling camera control for text-to-video generation, 2025. URL <https://arxiv.org/abs/2404.02101>.  
432
- 433 [22] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9  
434 (8):1735–1780, 1997. doi: 10.1162/neco.1997.9.8.1735.
- 435 [23] Yining Hong, Haoyu Zhen, Peihao Chen, Shuhong Zheng, Yilun Du, Zhenfang Chen, and  
436 Chuang Gan. 3d-llm: injecting the 3d world into large language models. In *Proceedings of the  
437 37th International Conference on Neural Information Processing Systems, NIPS '23*, Red Hook,  
438 NY, USA, 2023. Curran Associates Inc.
- 439 [24] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang,  
440 Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. URL  
441 <https://arxiv.org/abs/2106.09685>.

- 442 [25] Jiagyong Huang, Silong Yong, Xiaojian Ma, Xiongkun Linghu, Puhao Li, Yan Wang, Qing  
443 Li, Song-Chun Zhu, Baoxiong Jia, and Siyuan Huang. An embodied generalist agent in 3d  
444 world. In *Proceedings of the 41st International Conference on Machine Learning, ICML'24*.  
445 JMLR.org, 2024.
- 446 [26] Yibin Huang, Wang Xu, Wanyue Zhang, Helu Zhi, Jingjing Huang, Yangbin Xu, Yangang  
447 Sun, Conghui Zhu, and Tiejun Zhao. Video2layout: Recall and reconstruct metric-grounded  
448 cognitive map for spatial reasoning, 2026. URL <https://arxiv.org/abs/2511.16160>.
- 449 [27] Saurav Jha, M. Jehanzeb Mirza, Wei Lin, Shiqi Yang, and Sarath Chandar. Probing the  
450 effectiveness of world models for spatial reasoning through test-time scaling, 2025. URL  
451 <https://arxiv.org/abs/2512.05809>.
- 452 [28] Haoyi Jiang, Liu Liu, Xinjie Wang, Yonghao He, Wei Sui, Zhizhong Su, Wenyu Liu, and  
453 Xinggong Wang. Spa3r: Predictive spatial field modeling for 3d visual reasoning, 2026. URL  
454 <https://arxiv.org/abs/2602.21186>.
- 455 [29] Ruei-Chi Lai, Bolivar Enrique Solarte, Chin-Hsuan Wu, Yi-Hsuan Tsai, and Min Sun. Seeing  
456 once is enough? online geometry-aware token pruning for 3d question answering. In *The First*  
457 *Workshop on Efficient Spatial Reasoning*, 2026. URL [https://openreview.net/forum?](https://openreview.net/forum?id=jnDbE6cV2D)  
458 [id=jnDbE6cV2D](https://openreview.net/forum?id=jnDbE6cV2D).
- 459 [30] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan  
460 Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer,  
461 2024. URL <https://arxiv.org/abs/2408.03326>.
- 462 [31] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl  
463 Vondrick. Zero-1-to-3: Zero-shot one image to 3d object. In *2023 IEEE/CVF International*  
464 *Conference on Computer Vision (ICCV)*, pages 9264–9275, 2023. doi: 10.1109/ICCV51070.  
465 2023.00853.
- 466 [32] Zuntao Liu, Yi Du, Taimeng Fu, Shaoshu Su, Cherie Ho, and Chen Wang. Vision-language  
467 memory for spatial reasoning, 2025. URL <https://arxiv.org/abs/2511.20644>.
- 468 [33] Zuntao Liu, Yi Du, Taimeng Fu, Shaoshu Su, Cherie Ho, and Chen Wang. Vision-language  
469 memory for spatial reasoning, 2025. URL <https://arxiv.org/abs/2511.20644>.
- 470 [34] Xiaojian Ma, Silong Yong, Zilong Zheng, Qing Li, Yitao Liang, Song-Chun Zhu, and Siyuan  
471 Huang. Sqa3d: Situated question answering in 3d scenes, 2023. URL <https://arxiv.org/abs/2210.07474>.
- 473 [35] Pushmeet Kohli Nathan Silberman, Derek Hoiem and Rob Fergus. Indoor segmentation and  
474 support inference from rgbd images. In *ECCV*, 2012.
- 475 [36] J. O’Keefe and L. Nadel. *The Hippocampus as a Cognitive Map*. Clarendon Press, 1978. ISBN  
476 9780198572060. URL <https://books.google.com.hk/books?id=trR-AAAAIAAJ>.
- 477 [37] OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, and et.al. Gpt-4o system card, 2024.  
478 URL <https://arxiv.org/abs/2410.21276>.
- 479 [38] Jack Parker-Holder, Philip Ball, Jake Bruce, Vibhavari Dasagi, Kristian Holsheimer,  
480 Christos Kaplanis, Alexandre Moufarek, Guy Scully, Jeremy Shar, Jimmy Shi, Stephen  
481 Spencer, Jessica Yung, Michael Dennis, Sultan Kenjeyev, Shangbang Long, Vlad Mnih,  
482 Harris Chan, Maxime Gazeau, Bonnie Li, Fabio Pardo, Luyu Wang, Lei Zhang, Frederic  
483 Besse, Tim Harley, Anna Mitenkova, Jane Wang, Jeff Clune, Demis Hassabis, Raia  
484 Hadsell, Adrian Bolton, Satinder Singh, and Tim Rocktäschel. Genie 2: A large-scale  
485 foundation world model. 2024. URL [https://deepmind.google/discover/blog/](https://deepmind.google/discover/blog/genie-2-a-large-scale-foundation-world-model/)  
486 [genie-2-a-large-scale-foundation-world-model/](https://deepmind.google/discover/blog/genie-2-a-large-scale-foundation-world-model/).
- 487 [39] Zhangyang Qi, Zhixiong Zhang, Ye Fang, Jiaqi Wang, and Hengshuang Zhao. Gpt4scene:  
488 Understand 3d scenes from videos with vision-language models, 2025. URL <https://arxiv.org/abs/2501.01428>.
- 489

- 490 [40] Kyle Sargent, Zizhang Li, Tanmay Shah, Charles Herrmann, Hong-Xing Yu, Yunzhi Zhang,  
491 Eric Ryan Chan, Dmitry Lagun, Li Fei-Fei, Deqing Sun, and Jiajun Wu. Zeronvs: Zero-shot  
492 360-degree view synthesis from a single image, 2024. URL [https://arxiv.org/abs/2310.](https://arxiv.org/abs/2310.17994)  
493 17994.
- 494 [41] Saumya Saxena, Blake Buchanan, Chris Paxton, Peiqi Liu, Bingqing Chen, Narunas Vaskevicius,  
495 Luigi Palmieri, Jonathan Francis, and Oliver Kroemer. Grapheqa: Using 3d semantic scene  
496 graphs for real-time embodied question answering, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2412.14480)  
497 2412.14480.
- 498 [42] Jinzhou Tang, Jusheng zhang, Sidi Liu, Waikit Xiu, Qinhan Lv, and Xiying Li. Beyond  
499 pixels: Introducing geometric-semantic world priors for video-based embodied models via  
500 spatio-temporal alignment, 2025. URL <https://arxiv.org/abs/2509.00210>.
- 501 [43] Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea Vedaldi, Christian Rupprecht, and David  
502 Novotny. Vggt: Visual geometry grounded transformer, 2025. URL [https://arxiv.org/](https://arxiv.org/abs/2503.11651)  
503 [abs/2503.11651](https://arxiv.org/abs/2503.11651).
- 504 [44] Pan Wang, Yang Liu, Guile Wu, Eduardo R. Corral-Soto, Chengjie Huang, Binbin Xu, Dongfeng  
505 Bai, Xu Yan, Yuan Ren, Xingxin Chen, Yizhe Wu, Tao Huang, Wenjun Wan, Xin Wu, Pei Zhou,  
506 Xuyang Dai, Kangbo Lv, Hongbo Zhang, Yosef Fried, Aixue Ye, Bailan Feng, Zhenyu Chen,  
507 Zhen Li, Yingcong Chen, Yiyi Liao, and Bingbing Liu. Spatial4d-bench: A versatile 4d spatial  
508 intelligence benchmark, 2026. URL <https://arxiv.org/abs/2601.00092>.
- 509 [45] Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long Cui, Xingguang Wei, Zhaoyang  
510 Liu, Linglin Jing, Shenglong Ye, Jie Shao, Zhaokai Wang, Zhe Chen, Hongjie Zhang, Ganlin  
511 Yang, Haomin Wang, Qi Wei, Jinhui Yin, Wenhao Li, Erfei Cui, Guanzhou Chen, Zichen Ding,  
512 Changyao Tian, Zhenyu Wu, Jingjing Xie, Zehao Li, Bowen Yang, Yuchen Duan, Xuehui Wang,  
513 Zhi Hou, Haoran Hao, Tianyi Zhang, Songze Li, Xiangyu Zhao, Haodong Duan, Nianchen Deng,  
514 Bin Fu, Yanan He, Yi Wang, Conghui He, Botian Shi, Junjun He, Yingtong Xiong, Han Lv, Lijun  
515 Wu, Wenqi Shao, Kaipeng Zhang, Huipeng Deng, Biqing Qi, Jiaye Ge, Qipeng Guo, Wenwei  
516 Zhang, Songyang Zhang, Maosong Cao, Junyao Lin, Kexian Tang, Jianfei Gao, Haihan Huang,  
517 Yuzhe Gu, Chengqi Lyu, Huanze Tang, Rui Wang, Haijun Lv, Wanli Ouyang, Limin Wang,  
518 Min Dou, Xizhou Zhu, Tong Lu, Dahua Lin, Jifeng Dai, Weijie Su, Bowen Zhou, Kai Chen,  
519 Yu Qiao, Wenhao Wang, and Gen Luo. Internvl3.5: Advancing open-source multimodal models  
520 in versatility, reasoning, and efficiency, 2025. URL <https://arxiv.org/abs/2508.18265>.
- 521 [46] Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. Minilm: Deep  
522 self-attention distillation for task-agnostic compression of pre-trained transformers, 2020. URL  
523 <https://arxiv.org/abs/2002.10957>.
- 524 [47] Zhouxia Wang, Ziyang Yuan, Xintao Wang, Yaowei Li, Tianshui Chen, Menghan Xia, Ping  
525 Luo, and Ying Shan. Motionctrl: A unified and flexible motion controller for video generation.  
526 In *ACM SIGGRAPH 2024 Conference Papers*, SIGGRAPH '24, New York, NY, USA, 2024.  
527 Association for Computing Machinery. ISBN 9798400705250. doi: 10.1145/3641519.3657518.  
528 URL <https://doi.org/10.1145/3641519.3657518>.
- 529 [48] Daniel Watson, Saurabh Saxena, Lala Li, Andrea Tagliasacchi, and David J. Fleet. Controlling  
530 space and time with diffusion models, 2025. URL <https://arxiv.org/abs/2407.07860>.
- 531 [49] Yuxi Wei, Wei Huang, Qirui Chen, Lu Hou, and Xiaojuan Qi. See, remember, explore: A  
532 benchmark and baselines for streaming spatial reasoning, 2026. URL [https://arxiv.org/](https://arxiv.org/abs/2603.23864)  
533 [abs/2603.23864](https://arxiv.org/abs/2603.23864).
- 534 [50] Diankun Wu, Fangfu Liu, Yi-Hsin Hung, and Yueqi Duan. Spatial-mlm: Boosting mllm  
535 capabilities in visual-based spatial intelligence, 2025. URL [https://arxiv.org/abs/2505.](https://arxiv.org/abs/2505.23747)  
536 23747.
- 537 [51] Shun-Cheng Wu, Keisuke Tateno, Nassir Navab, and Federico Tombari. Incremental 3d semantic  
538 scene graph prediction from rgb sequences. In *Proceedings of the IEEE/CVF Conference on*  
539 *Computer Vision and Pattern Recognition*, pages 5064–5074, 2023.

- 540 [52] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang  
541 Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu,  
542 Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin  
543 Yang, Jiayi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang,  
544 Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui  
545 Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang  
546 Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger  
547 Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan  
548 Qiu. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.
- 549 [53] Jihan Yang, Shusheng Yang, Anjali W. Gupta, Rilyn Han, Li Fei-Fei, and Saining Xie. Thinking  
550 in space: How multimodal large language models see, remember, and recall spaces, 2025. URL  
551 <https://arxiv.org/abs/2412.14171>.
- 552 [54] Rui Yang, Ziyu Zhu, Yanwei Li, Jingjia Huang, Shen Yan, Siyuan Zhou, Zhe Liu, Xiangtai Li,  
553 Shuangye Li, Wenqian Wang, Yi Lin, and Hengshuang Zhao. Visual spatial tuning, 2025. URL  
554 <https://arxiv.org/abs/2511.05491>.
- 555 [55] Yuncong Yang, Jiageng Liu, Zheyuan Zhang, Siyuan Zhou, Reuben Tan, Jianwei Yang, Yilun  
556 Du, and Chuang Gan. Mindjourney: Test-time scaling with world models for spatial reasoning,  
557 2025. URL <https://arxiv.org/abs/2507.12508>.
- 558 [56] Yuncong Yang, Han Yang, Jiachen Zhou, Peihao Chen, Hongxin Zhang, Yilun Du, and Chuang  
559 Gan. 3d-mem: 3d scene memory for embodied exploration and reasoning, 2025. URL  
560 <https://arxiv.org/abs/2411.17735>.
- 561 [57] Chandan Yeshwanth, Yueh-Cheng Liu, Matthias Nießner, and Angela Dai. Scannet++: A high-  
562 fidelity dataset of 3d indoor scenes, 2023. URL <https://arxiv.org/abs/2308.11417>.
- 563 [58] Tatiana Zemsikova and Dmitry Yudin. 3dgraphllm: Combining semantic graphs and large  
564 language models for 3d scene understanding, 2025. URL <https://arxiv.org/abs/2412.18450>.
- 565
- 566 [59] Gongjie Zhang, Wenhao Li, Quan hao Qian, Juniu Wang, Deli Zhao, Shijian Lu, and Ran  
567 Xu. On the generalization capacities of mllms for spatial intelligence, 2026. URL <https://arxiv.org/abs/2603.06704>.
- 568
- 569 [60] Kevin Zhang, Kuangzhi Ge, Xiaowei Chi, Renrui Zhang, Shaojun Shi, Zhen Dong, Sirui Han,  
570 and Shanghang Zhang. Can world models benefit vlms for world dynamics?, 2025. URL  
571 <https://arxiv.org/abs/2510.00855>.
- 572 [61] Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. Llava-  
573 video: Video instruction tuning with synthetic data, 2025. URL <https://arxiv.org/abs/2410.02713>.
- 574
- 575 [62] Duo Zheng, Shijia Huang, Yanyang Li, and Liwei Wang. Learning from videos for 3d world:  
576 Enhancing mllms with 3d vision geometry priors, 2025. URL <https://arxiv.org/abs/2505.24625>.
- 577
- 578 [63] Duo Zheng, Shijia Huang, and Liwei Wang. Video-3d llm: Learning position-aware video  
579 representation for 3d scene understanding. In *2025 IEEE/CVF Conference on Computer Vision  
580 and Pattern Recognition (CVPR)*, pages 8995–9006, 2025. doi: 10.1109/CVPR52734.2025.  
581 00841.
- 582 [64] Jensen Zhou, Hang Gao, Vikram Voleti, Aaryaman Vasishta, Chun-Han Yao, Mark Boss, Philip  
583 Torr, Christian Rupprecht, and Varun Jampani. Stable virtual camera: Generative view synthesis  
584 with diffusion models, 2025. URL <https://arxiv.org/abs/2503.14489>.
- 585 [65] Ziyu Zhu, Xiaojuan Ma, Yixin Chen, Zhidong Deng, Siyuan Huang, and Qing Li. 3d-vista:  
586 Pre-trained transformer for 3d vision and text alignment. In *2023 IEEE/CVF International  
587 Conference on Computer Vision (ICCV)*, pages 2899–2909, 2023. doi: 10.1109/ICCV51070.  
588 2023.00272.

589 **Appendix**

590 **A Pseudocode For Inference** **15**

591 **B Training Setup Details** **17**

592 **C Inference setup details** **19**

593   C.1 Instantiation with the Trained VLM and World Model . . . . . 20

594   C.2 Inference Configurations Used in Practice . . . . . 22

595 **D BEV Map and Graph Examples** **22**

596 **E Case Study** **25**

597   E.1 Positive Case 1: Relative Direction Reasoning . . . . . 26

598   E.2 Negative Case 1: Hallucinated Spatial Evidence . . . . . 28

599 **F Statistical Significance** **31**

600 **G Broader Impact** **31**

601 **H Assets and License** **32**

602 **I Prompt Examples and Token Examples** **33**

603   I.1 Prompt Examples . . . . . 33

604   I.2 Token Examples . . . . . 34

605 **A Pseudocode For Inference**

606 We summarize the inference procedure of LSM-VLM in Algorithms 1–3. Algorithm 1 gives the  
 607 overall closed-loop inference process. Given an egocentric video  $\mathcal{V}$  and a question  $q$ , the model first  
 608 constructs an initial dual memory  $\mathcal{M}_0 = \{G_0, B_0\}$ , where  $G_0$  is the scene graph and  $B_0$  is the local  
 609 BEV map. The detector  $\Phi_{\text{det}}(\cdot)$  extracts visual evidence and returns the initial memory together with  
 610 confidence  $\rho_0$ . At each round, the model imagines missing spatial evidence, updates the memory,  
 611 and then feeds the last imagined frame  $\tilde{I}_t^{\text{last}}$  together with the updated memory  $\mathcal{M}_{t+1}$  to the VLM  
 612 for answering. The loop stops when the VLM outputs a final answer signal or when the maximum  
 613 number of rounds is reached. In our implementation, we set the maximum number of rounds as 4.

614 Algorithm 2 describes the world-model imagination step. At round  $t$ , the current memory is  $\mathcal{M}_t =$   
 615  $\{G_t, B_t\}$ . The pose predictor  $\text{VLM}_\theta^{\text{pose}}$  predicts informative viewpoints  $\mathcal{P}_t = \{p_t^0, p_t^1, p_t^2, p_t^3\}$  based  
 616 on the question, observed video, imagined history, and current memory. The planner  $\text{Plan}(\cdot)$  then  
 617 computes an action sequence  $\mathcal{A}_t$  from the current pose  $p_t^{\text{cur}}$  to the target pose  $p_t^*$ . The world model  $\mathcal{W}_\psi$   
 618 generates an imagined rollout  $\tilde{\mathcal{V}}_t$ , which is parsed into candidate memory updates  $(\Delta B_t, \Delta G_t, \rho_t)$ .

619 Algorithm 3 shows how candidate updates are fused into memory. For the BEV map,  $f_B(B_t, \rho_{t-1})$   
 620 keeps reliable previous cells, while  $i_B(\Delta B_t, \rho_t)$  injects new confidence-weighted BEV evidence.  
 621 Since the newly generated BEV update  $\Delta B_t$  and the previous BEV memory  $B_t$  are not necessarily  
 622 defined in the same coordinate frame, we always use an agent-centric local frame, where the agent  
 623 position is placed at the center and its heading direction is aligned upward. Before memory fusion, all  
 624 spatial coordinates are transformed into this current local frame, including the local object coordinates  
 625 stored in the scene graph. The updated BEV memory is denoted as  $B_{t+1}$ . For the scene graph,  
 626  $i_G(G_t, B_{t+1}, \Delta G_t, \rho_t)$  verifies candidate graph updates using the consolidated BEV map before  
 627 adding them to  $G_t$ . The final updated memory is  $\mathcal{M}_{t+1} = \{G_{t+1}, B_{t+1}\}$ .

---

**Algorithm 1** Overall LSM-VLM Inference

---

**Require:** Egocentric video  $\mathcal{V}$ , question  $q$ , maximum rounds  $T_{\max}$

**Ensure:** Final answer  $a$

▷ **Perceptual Grounding**

- 1: Extract geometry, object masks, and semantic labels from  $\mathcal{V}$  using VGGT and SAM3.

$$(B_0, G_0, \rho_0) \leftarrow \Phi_{\text{det}}(\mathcal{V}).$$

- 2: Build initial scene graph  $G_0$  and local BEV map  $B_0$ .

- 3: Initialize dual memory  $\mathcal{M}_0 \leftarrow \{G_0, B_0\}$ .

- 4: Initialize imagined buffer  $\tilde{\mathcal{I}}_{<0} \leftarrow \emptyset$  and confidence  $\rho_0$ .

- 5: **for**  $t = 0$  **to**  $T_{\max} - 1$  **do**

▷ **Evidence Completion Loop (World Model Flow)**

- 6:  $(\tilde{\mathcal{V}}_t, \Delta B_t, \Delta G_t, \rho_t) \leftarrow \text{IMAGINEEVIDENCE}(q, \mathcal{V}, \tilde{\mathcal{I}}_{<t}, \mathcal{M}_t)$ .

- 7: Extract the last imagined frame & Update imagined evidence history:

$$\tilde{I}_t^{\text{last}} \leftarrow \text{LastFrame}(\tilde{\mathcal{V}}_t), \quad \tilde{\mathcal{I}}_{<t+1} \leftarrow \tilde{\mathcal{I}}_{<t} \cup \{\tilde{I}_t^{\text{last}}\}.$$

- 8:  $\mathcal{M}_{t+1} \leftarrow \text{UPDATEMEMORY}(\mathcal{M}_t, \Delta B_t, \Delta G_t, \rho_{t-1}, \rho_t)$ .

▷ **VLM Answer Flow**

- 9: Feed the last imagined frame and updated memory to the VLM:

$$(y_t, s_t) \leftarrow \text{VLM}_{\theta}^{\text{ans}}(q, \mathcal{V}, \tilde{\mathcal{I}}_{<t+1}, \mathcal{M}_{t+1}).$$

- 10: **if**  $s_t = \text{FINALANSWER}$  **or**  $t = T_{\max} - 1$  **then**

- 11: **return**  $y_t$ .

- 12: **end if**

- 13: **end for**
- 

---

**Algorithm 2** IMAGINEEVIDENCE: Question-Conditioned World Imagination

---

**Require:** Question  $q$ , video  $\mathcal{V}$ , imagined buffer  $\tilde{\mathcal{I}}_{<t}$ , memory  $\mathcal{M}_t = \{G_t, B_t\}$

**Ensure:** Imagined rollout  $\tilde{\mathcal{V}}_t$ , candidate updates  $\Delta B_t, \Delta G_t$ , confidence  $\rho_t$

▷ **Target Pose Prediction**

- 1: Predict informative viewpoints:

$$\mathcal{P}_t = \{p_t^0, p_t^1, p_t^2, p_t^3\} \leftarrow \text{VLM}_{\theta}^{\text{pose}}(q, \mathcal{V}, \tilde{\mathcal{I}}_{<t}, \mathcal{M}_t).$$

- 2: Select target pose  $p_t^*$  according to  $q$  and  $\mathcal{M}_t$ .

- 3: Plan actions on the current BEV map:

$$\mathcal{A}_t \leftarrow \text{Plan}(B_t, p_t^{\text{cur}}, p_t^*).$$

▷ **World Model Rollout**

- 4: Retrieve observations at queried viewpoints:

$$\mathcal{I}_{\mathcal{P}_t} \leftarrow \{I(p_t^0), I(p_t^1), I(p_t^2), I(p_t^3)\}.$$

- 5: Generate imagined rollout:

$$\tilde{\mathcal{V}}_t \leftarrow \mathcal{W}_{\psi}(\mathcal{I}_{\mathcal{P}_t}, \mathcal{A}_t).$$

▷ **Candidate Evidence Extraction**

- 6: Apply VGGT and SAM3 to  $\tilde{\mathcal{V}}_t$ .

- 7: Convert detections into candidate memory updates:

$$(\Delta B_t, \Delta G_t, \rho_t) \leftarrow \Phi_{\text{det}}(\tilde{\mathcal{V}}_t).$$

- 8: Transform  $\Delta B_t$  into the current local BEV frame using the estimated relative pose.

- 9: **return**  $(\tilde{\mathcal{V}}_t, \Delta B_t, \Delta G_t, \rho_t)$ .
-

---

**Algorithm 3** UPDATEMEMORY: Confidence-Aware Memory Consolidation

---

**Require:** Current memory  $\mathcal{M}_t = \{G_t, B_t\}$ , candidate updates  $\Delta B_t, \Delta G_t$ , historical confidence  $\rho_{t-1}$ , current confidence  $\rho_t$

**Ensure:** Updated memory  $\mathcal{M}_{t+1}$

▷ **Coordinate Alignment**

- 1: Align BEV cells and graph coordinates to the current agent-centric local frame:

$$(B_t^{\text{loc}}, G_t^{\text{loc}}, \Delta B_t^{\text{loc}}, \Delta G_t^{\text{loc}}) \leftarrow \text{AlignToAgentFrame}(B_t, G_t, \Delta B_t, \Delta G_t).$$

▷ **BEV Memory Consolidation**

- 2: Preserve reliable cells and inject new confidence-weighted BEV evidence:

$$B_t^{\text{keep}} \leftarrow f_B(B_t^{\text{loc}}, \rho_{t-1}), \quad B_t^{\text{new}} \leftarrow i_B(\Delta B_t^{\text{loc}}, \rho_t).$$

- 3: Update BEV memory:

$$B_{t+1} \leftarrow B_t^{\text{keep}} + B_t^{\text{new}}.$$

▷ **Scene Graph Consolidation**

- 4: Verify candidate graph updates and update long-term semantic memory:

$$\Delta G_t^{\text{valid}} \leftarrow i_G(G_t^{\text{loc}}, B_{t+1}, \Delta G_t^{\text{loc}}, \rho_t), \quad G_{t+1} \leftarrow G_t^{\text{loc}} + \Delta G_t^{\text{valid}}.$$

▷ **Dual Memory Output**

- 5:  $\mathcal{M}_{t+1} \leftarrow \{G_{t+1}, B_{t+1}\}$ .
  - 6: **return**  $\mathcal{M}_{t+1}$ .
- 

## 628 B Training Setup Details

629 Our training pipeline is implemented in a Hugging Face and DeepSpeed codebase built around  
630 transformers, peft, and trl, with custom modules under src/train, src/dataset, and  
631 src/trainer. The default backbone used in both stages is Qwen/Qwen3-VL-8B-Instruct. The  
632 implementation extends the standard Qwen-VL supervised fine-tuning pipeline by adding two extra  
633 modality branches, namely a *map branch* and a *graph branch*, and by monkey-patching the Qwen  
634 forward function so that these additional modality features can be injected into the language token  
635 stream during training.

636 **Additional modality branches.** The map branch uses google/vit-base-patch16-224-in21k  
637 as the map encoder, while the graph branch uses sentence-transformers/all-MiniLM-L6-v2  
638 as the graph encoder. Two learnable linear projectors are newly introduced to align the map and graph  
639 features to the hidden size of the Qwen language model. Concretely, the model adds map\_encoder,  
640 graph\_encoder, map\_projector, and graph\_projector on top of the original Qwen-VL archi-  
641 tecture. Two new special tokens, <|map\_pad|> and <|graph\_pad|>, are added to the tokenizer  
642 vocabulary and are used as placeholders for the injected map and graph features.

643 **Map representation.** Each training sample may contain a video, a top-down occupancy map  
644 sequence, and a graph file. For maps, the code first collects the map frames associated with the  
645 sample. When training on video data, the map frames are aligned to the sampled video frame indices;  
646 in the default implementation, this is done by selecting files such as occ\_{frame\_idx}.png from  
647 the occupancy-map directory. The selected map frames are tiled into a single 2D canvas arranged  
648 on a near-square grid, and the resulting canvas is processed by the map image processor. When  
649 semantic supervision is enabled (semantic=True in both training scripts), an additional semantic  
650 tensor is loaded from the corresponding .npz files and concatenated channel-wise with the RGB map  
651 canvas. As a result, the input to the map encoder has 43 channels (3 RGB channels plus 40 semantic  
652 channels in the default setup). Since the pretrained ViT map encoder expects 3-channel input, its  
653 patch embedding convolution is explicitly adapted to 43 input channels by copying the original RGB  
654 weights and initializing the extra channels with the mean of the original filters.

655 **Graph representation.** For graph inputs, the code loads the graph file (typically  
656 graph/graph.json) as raw text and tokenizes it with the graph encoder tokenizer. Graph se-

657 quences are truncated or padded to a maximum length of 256 tokens. If a graph exists but the human  
658 prompt does not explicitly contain a graph placeholder, the loader automatically inserts `<graph>`  
659 into the prompt so that graph features are always aligned with an explicit placeholder position in the  
660 language sequence.

661 **Feature injection into Qwen.** The Qwen forward pass is modified to support mixed-modality  
662 inputs. During training, the visual encoder processes the image or video tokens as usual, while the  
663 new map and graph encoders produce separate hidden-state sequences. These hidden states are then  
664 compressed into a fixed number of placeholder-aligned tokens by chunk-wise average pooling after  
665 dropping the first encoder token. In the default configuration, each map input is compressed into  
666 16 tokens and each graph input is compressed into 8 tokens. Accordingly, in the text preprocessing  
667 stage, `<map>` is expanded into 16 repeated `<|map_pad|>` tokens and `<graph>` is expanded into 8  
668 repeated `<|graph_pad|>` tokens. After projection to the Qwen hidden dimension, the resulting map  
669 and graph features replace the embeddings at those placeholder positions, so the language model  
670 receives aligned map-aware and graph-aware token embeddings inside the normal autoregressive  
671 sequence.

672 **Data formatting and supervision.** The training data are stored in JSON format and contain  
673 conversation pairs together with multimodal references, e.g., `video`, `map`, and `graph`. A typical  
674 prompt contains `<video>`, `<map>`, and `<graph>` markers followed by the task instruction. The data  
675 loader converts the original LLaVA-style placeholders into Qwen-compatible tokens, assembles the  
676 multimodal tensors, and constructs next-token prediction targets in the standard supervised fine-tuning  
677 format. Only assistant responses contribute to the loss; all user-side prompt tokens are masked with  
678 `IGNORE_INDEX`. Since the default backbone in our scripts is Qwen3-VL, no extra system prompt is  
679 prepended by the data loader.

680 **Optimization framework.** Both stages use DeepSpeed ZeRO-2, `adamw_torch`, `bfloat16` train-  
681 ing, TF32 matmul, gradient checkpointing, cosine learning-rate decay, and a warmup ratio of 0.03.  
682 The scripts enable FlashAttention-2 (`disable_flash_attn2=False`) and disable Liger kernels  
683 (`use_liger_kernel=False`). Weight decay is 0.1. Training logs are written every step and check-  
684 points are saved every 2000 steps with a maximum of 5 retained checkpoints. The scripts use  
685 on-the-fly preprocessing (`lazy_preprocess=True`) with `data_loader_num_workers=0`. Images  
686 are constrained to a pixel budget between  $512 \times 28^2$  and  $1280 \times 28^2$ , and videos are sampled with  
687 128 frames by default.

688 **Stage 1: modality-branch pretraining.** Stage 1 is designed to learn the newly added map and  
689 graph components while keeping the original Qwen-VL backbone fixed. Specifically, LoRA is  
690 disabled, the Qwen language model is frozen, the original visual tower is frozen, and the visual  
691 merger is frozen. In contrast, both the map encoder and graph encoder are trainable, and both  
692 projectors are trainable. The default learning rates are  $2 \times 10^{-6}$  for the map encoder,  $2 \times 10^{-6}$   
693 for the graph encoder,  $1 \times 10^{-5}$  for the map projector, and  $1 \times 10^{-5}$  for the graph projector. The  
694 generic learning rate argument is  $5 \times 10^{-5}$ , but in practice the trainable parameters in this stage are  
695 dominated by the modality-specific encoders and projectors because the base Qwen-VL modules  
696 remain frozen. The default script runs for 10 epochs on 8 GPUs with per-device batch size 1, giving a  
697 default global batch size of 8 and gradient accumulation of 1.

698 **Stage 2: LoRA-based alignment.** Stage 2 starts from a Stage-1 checkpoint (by default  
699 `output/4_09_stage1_map_graph_only/checkpoint-2000`) and performs a weights-only con-  
700 tinuation rather than resuming the optimizer or scheduler states. LoRA is enabled with rank 16,  
701 scaling factor 32, and dropout 0.05. The code applies LoRA to all eligible linear or embedding  
702 modules except those explicitly excluded by name. By default, LoRA is excluded from the visual  
703 backbone, the map encoder, the graph encoder, the two projectors, the language-model embedding  
704 layer, and the LM head. Therefore, Stage 2 mainly aligns the language-model internals to the already  
705 learned map/graph representations while preserving the newly introduced modality branches as  
706 separate trainable components. The Qwen language model, vision tower, and visual merger remain  
707 frozen in the standard full-parameter sense, but LoRA adapters on the allowed language-model  
708 modules are trainable. In the provided script, both the map encoder and graph encoder also remain  
709 unfrozen by default (`FREEZE_MAP_ENCODER=False`, `FREEZE_GRAPH_ENCODER=False`), so Stage  
710 2 jointly updates the LoRA adapters, the map encoder, the graph encoder, and both projectors. The

711 default script runs for 10 epochs on 8 GPUs with per-device batch size 1, corresponding to a default  
712 global batch size of 8 and gradient accumulation of 1.

713 **Learning-rate grouping.** The trainer uses explicit parameter groups with separate learning rates  
714 for the visual tower, merger, map encoder, graph encoder, map projector, and graph projector. In our  
715 two-stage setting, this design is important because it allows the newly introduced modality branches  
716 to be trained conservatively while LoRA parameters in the language model use the base learning rate.  
717 This decoupled optimization is one of the key implementation choices that stabilizes training after  
718 adding heterogeneous map and graph inputs to a large pretrained vision-language backbone.

719 **Checkpointing and reproducibility details.** When loading from a local checkpoint, the code  
720 first instantiates the backbone, then attaches the custom map/graph modules, and finally reloads  
721 the checkpoint weights with non-strict matching so that the extra modality parameters are correctly  
722 restored. If LoRA is enabled, checkpoints additionally save the non-LoRA trainable parameters  
723 separately, which is necessary because the map/graph modules are not part of the LoRA adapters  
724 themselves. The training code also logs trainable-parameter summaries and gradient norms for the  
725 map encoder, graph encoder, and both projectors, which we used as a sanity check to ensure that  
726 gradients propagated correctly through the newly added branches.

## 727 C Inference setup details

728 We implement a closed-loop inference pipeline that explicitly separates perceptual grounding, plan-  
729 ning, future-view imagination, memory consolidation, and final answering. The lightweight unit  
730 test in `tests/test_pipeline.py` validates the contract of this pipeline, while the full implemen-  
731 tation in `origin_code/inference` instantiates the heavy-weight vision-language model (VLM),  
732 map/graph-conditioned prompting, and Stable Virtual Camera (SVC) world-model generation.

733 **Pipeline contract.** The pipeline interface is defined by `LSMVLMPipeline` with con-  
734 figuration `PipelineConfig(rounds, bev_height, bev_width, bev_resolution,`  
735 `graph_update_interval, manual_actions)`. Given a question, a sequence of observa-  
736 tions, and optional initial graph/BEV files, the pipeline executes the following stages:

- 737 1. load the initial cognitive memory from a scene graph and a BEV map;
- 738 2. ground the input observations into an initial graph and BEV update;
- 739 3. initialize the recurrent memory;
- 740 4. predict a target pose or action sequence conditioned on the current memory and the question;
- 741 5. imagine future observations along the planned trajectory using a world model;
- 742 6. detect candidate objects and BEV occupancy from the imagined views;
- 743 7. update the memory using a confidence-aware recurrent rule;
- 744 8. answer the question from the updated memory.

745 The unit test additionally verifies that the pipeline serializes intermediate memory states to disk  
746 (e.g., `memory_round_000.json`, `memory_round_001.json`, and `result.json`) and that action  
747 parsing and execution are deterministic.

748 **Dual-memory representation.** The recurrent memory is represented as a pair of structured states:

- 749 • a *scene graph* containing object nodes, relations, and the agent pose;
- 750 • a *BEV map* containing occupancy, semantic labels, confidence, origin, and metric resolution.

751 Each object node stores a node ID, category name, planar position, extent, confidence, semantic ID,  
752 and auxiliary attributes. The BEV memory stores three aligned grids: occupancy, semantic labels,  
753 and confidence. Coordinates are converted between world and grid frames using the map origin  
754 and resolution. This design matches the role split used throughout the codebase: the graph provides  
755 sparse symbolic structure, while the BEV map provides dense spatial support.

756 **Confidence-aware recurrent update.** Memory updates are consolidated by  
757 `ConfidenceAwareMemoryUpdater`. Let the previous BEV confidence be decayed by a fac-  
758 tor of 0.92, and let cells whose decayed confidence falls below 0.25 be forgotten. New candidate  
759 occupancy and semantic evidence are injected wherever the new confidence exceeds the retained  
760 confidence. For graph updates, candidate nodes are added only if they are confirmed by the BEV  
761 map, i.e., the confidence at the node’s projected BEV location exceeds the graph confirmation  
762 threshold (0.35 by default). Graph nodes with the same class name are merged when they fall within  
763 a Euclidean distance of 0.75 m. After each update, pairwise spatial relations such as `left_of`,  
764 `right_of`, `in_front_of`, and `behind` are recomputed for nearby objects.

765 **Trajectory parameterization.** Planned actions are represented as short symbolic commands such  
766 as `move-forward {meters}`, `turn-left {degrees}`, and `turn-right {degrees}`. The action  
767 parser accepts either a JSON array or a comma-separated textual list. Actions are deterministically  
768 applied to a planar pose  $(x, y, \theta)$ , and the resulting waypoints are converted to camera extrinsics for  
769 the world model. The unit test explicitly checks that parsing `["turn-left 90", "move-forward  
770 1"]` from the default origin pose produces the expected translated pose.

## 771 C.1 Instantiation with the Trained VLM and World Model

772 The abstract pipeline above is instantiated in `origin_code/inference` using the trained map/graph-  
773 conditioned Qwen-VL model and a future-view generator. The code supports both a direct-answer  
774 baseline and a closed-loop world-model setting.

775 **VLM loading and modality restoration.** At inference time, the code reconstructs the multimodal  
776 model in the same way as during training. It first reads `training_args.bin` from the checkpoint  
777 directory to recover the base model ID, map encoder ID, graph encoder ID, map token count, graph  
778 token count, semantic-map flag, and numerical precision. It then:

- 779 1. loads the Qwen-VL base model;
- 780 2. monkey-patches the forward function to accept `map_pixel_values` and  
781 `graph_input_ids`;
- 782 3. re-attaches `map_encoder`, `graph_encoder`, `map_projector`, and `graph_projector`;
- 783 4. restores either a full checkpoint or a LoRA checkpoint;
- 784 5. if LoRA is used, separately loads `non_lora_state_dict.bin` before merging the adapter  
785 into the base model.

786 This is important because the map and graph branches are not native parts of the original Qwen-VL  
787 architecture and therefore must be explicitly re-created before checkpoint restoration.

788 **Prompt and modality preparation.** The inference code preserves the same placeholder mechanism  
789 used in training. Human prompts are normalized so that `<map>` and `<graph>` are expanded to repeated  
790 `<|map_pad|>` and `<|graph_pad|>` spans according to the checkpoint-specific token counts. When  
791 the graph is available but the prompt omits an explicit graph placeholder, the code automatically  
792 inserts one. For video inputs, the model samples frames using the same Qwen video preprocessing  
793 utilities as in training; by default, the scripts use up to 128 frames. The map input is constructed  
794 by aligning occupancy-map frames with the sampled video frame indices and then building a tiled  
795 map canvas. When semantic maps are enabled, semantic one-hot channels are concatenated to the  
796 RGB map channels before being fed to the adapted ViT map encoder. Graph inputs are read from  
797 `graph/graph.json` and tokenized to a maximum length of 256.

798 **Two-stage closed-loop reasoning.** The full inference pipeline uses the VLM twice. In the first  
799 stage, the VLM acts as a planner. In the second stage, after the world model synthesizes future  
800 observations, the VLM acts as an observer-answerer.

801 There are two planner modes:

- 802 • `actions`: the model directly predicts a short action sequence in JSON form;

803       • `navigation_target`: the model predicts a graph node ID and object name, and the code  
804       converts this symbolic target into a metric goal pose and a derived action sequence.

805   When graph input is unavailable for the current sample, the code automatically falls back from  
806   `navigation_target` mode to `actions` mode.

807   **Question-aware planning prompts.** Planner prompts are specialized by task type. The code  
808   contains separate prompt templates for generic action planning, navigation-target selection, SQA3D-  
809   style situation-grounded navigation, and route-planning tasks. For object-relative-direction questions,  
810   the planner prompt is augmented with an explicit explanation that the “standing by” object defines the  
811   starting position and the “facing” object defines the initial heading. For route-planning questions, the  
812   planner is encouraged to select the final destination object node rather than an intermediate landmark.  
813   For SQA3D, the planner is instructed to infer the start pose from the situation description before  
814   choosing a target that is useful for answering the question.

815   **Source-view selection before world-model rollout.** A distinctive implementation detail is that  
816   the planner is not always anchored to the last video frame. The code resolves ground-truth camera  
817   poses for all video frames and, when the question text provides a reliable start anchor, it selects a  
818   source frame whose pose best matches that textual anchor. This is used for route-planning, SQA3D,  
819   and object-relative-direction questions. If the selected planner-start pose differs too much from the  
820   current visual anchor (more than 1.0 m position error or more than 45° yaw error), the code inserts  
821   an explicit relocation segment before the task-specific actions. This makes the generated trajectory  
822   consistent with the semantics of the question rather than merely with the temporal end of the source  
823   video.

824   **World-model conditioning and rollout.** The world model is instantiated through  
825   `generate_video_from_image_camera_poses` and is configured in the provided scripts  
826   with `stabilityai/stable-virtual-camera` as the default backend. The default configuration in  
827   `inference_world_model_plan.sh` uses:

828       • model version 1.1;  
829       • spatial resolution  $576 \times 576$ ;  
830       • shortest side 576;  
831       • context frames 21, 21;  
832       • classifier-free guidance 4.0, 2.0;  
833       • guider types 1, 2;  
834       • minimum CFG 1.2;  
835       • camera scale 2.0;  
836       • 20 denoising steps;  
837       • trajectory prior enabled;  
838       • crop-based transform target with scale 1.0;  
839       • frame interval 3;  
840       • field of view 90°;  
841       • camera pitch 20°;  
842       • output frame rate 8 FPS.

843   The model is conditioned on multiple source views extracted from the original video. By default,  
844   three matched conditioning frames plus the anchor frame are selected so that the conditioning set is  
845   geometrically close to the planned future path.

846   **Observer-view construction.** After planning, the code converts the symbolic actions into dense  
847   camera poses and writes them to `camera_poses.json`. It then concatenates: (i) the source condi-  
848   tioning poses, (ii) the planned forward trajectory, and, for some question types, (iii) additional orbit  
849   views around the target pose. For object-relative-direction tasks, the code renders a dense 15° orbit  
850   around the final pose and feeds a sparse 45° subset of up to seven views back into the VLM, so the  
851   observer receives multiple imagined perspectives from the same target location.

852 **Follow-up answering.** Once the imagined video is generated, the code either extracts the final  
853 predicted frame or a selected subset of observer views and asks the same VLM to answer the original  
854 question. The follow-up prompt explicitly states that the image(s) are predicted future views after  
855 executing the planned actions or after navigating to the selected target node. In navigation-target  
856 mode, the follow-up prompt also injects the metric goal pose  $(x, y, \text{yaw})$  in map coordinates. If both  
857 map and graph are disabled, the code skips the world-model stage entirely and directly applies the  
858 trained VLM to the original video.

## 859 C.2 Inference Configurations Used in Practice

860 The provided shell scripts define the main evaluation settings.

861 **Direct trained-model inference.** The script `inference_trained.sh` evaluates the trained check-  
862 point from `output/4_10_stage2_lora_align/checkpoint-6000` by default. It enables seman-  
863 tic maps by default, resumes unfinished runs, and launches three modality ablations:

- 864 • map + graph,
- 865 • map only,
- 866 • graph only.

867 The default launcher uses GPUs 0, 1, 2, 3 with 3 workers per GPU.

868 **Closed-loop world-model inference.** The script `inference_world_model_plan.sh` evaluates  
869 the full planner–world-model–observer pipeline. By default it uses:

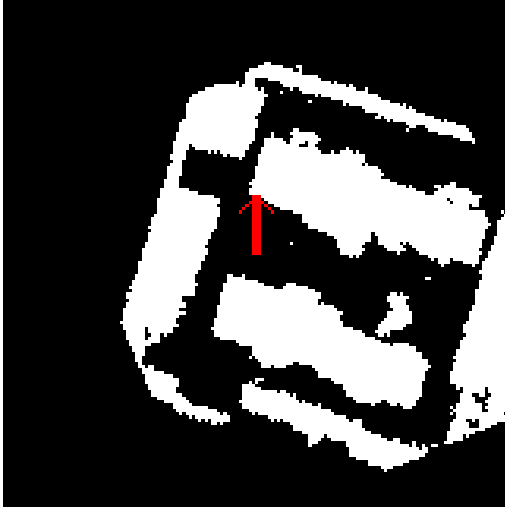
- 870 • the same stage-2 checkpoint at step 6000;
- 871 • 2 workers per GPU;
- 872 • route-planning validation data from `dataset/data_val_route_planning_scannet_scannetpp.json`;
- 873 • planner output mode `actions`;
- 874 • map directory `obs_bbox_id/occupancy`;
- 875 • graph file `graph/graph.json`;
- 876 • maximum VLM generation length 256, temperature 0, and top- $p$  0.9;
- 877 • VLM video sampling with up to 128 frames.

878 **Outputs for reproducibility.** For each sample, the code writes a detailed result package  
879 containing the resolved input paths, selected source frames, source and target poses, plan-  
880 ner prompt, planner raw output, parsed actions or navigation target, world-model camera  
881 poses, rendered future view(s), observer prompt, and final answer. In the heavy-weight  
882 world-model setting, the output directory also contains `world_model.mp4`, `world_model.json`,  
883 `camera_poses.json`, `world_model_camera_poses.json`, and per-sample result summaries  
884 such as `vlm_world_model_result.json`. This serialization mirrors the test-pipeline design goal:  
885 every major intermediate state is materialized to disk for debugging and analysis.

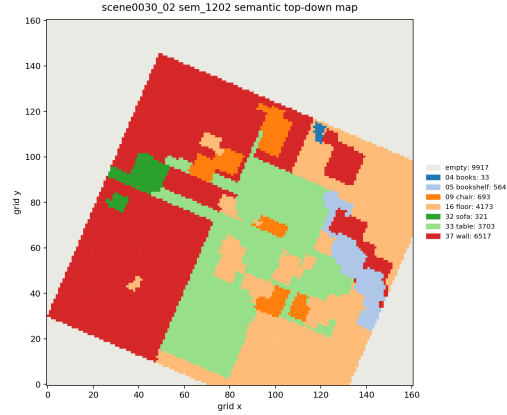
886 **Interpretation.** Conceptually, `tests/test_pipeline.py` specifies the algorithmic skeleton of  
887 LSM-VLM, namely recurrent memory initialization, plan generation, imagination, update, and answer.  
888 The code under `origin_code/inference` provides the concrete realization used in experiments:  
889 a map/graph-augmented Qwen-VL planner-observer, question-aware source-pose selection, Stable  
890 Virtual Camera rollout, and confidence-aware closed-loop reasoning over future predicted views.

## 891 D BEV Map and Graph Examples

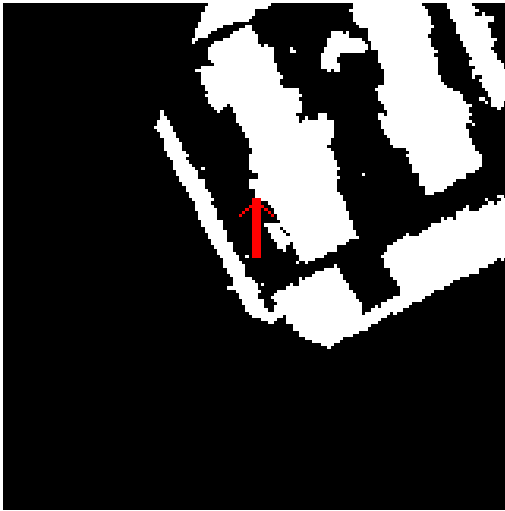
892 Here, the short-term memory represented by our BEV map does not refer to a memory that is  
893 immediately forgotten after observation. Instead, similar to the hidden state in an LSTM, this short-  
894 term memory is continuously maintained and updated over time. It is termed “short-term” because  
895 it is updated at a higher frequency and may gradually forget outdated information. In contrast,  
896 long-term memory is more stable and persistent, and is not subject to forgetting.



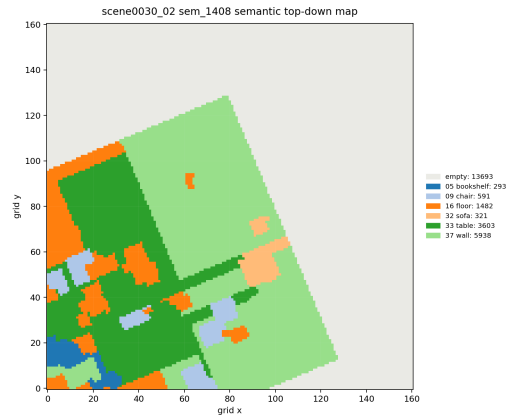
(a) Occupancy Channel from the first viewpoint.



(b) Semantic Channel from the first viewpoint.



(c) Occupancy Channel from another viewpoint.



(d) Semantic Channel from another viewpoint.

Figure 4: Examples of the BEV map representation from two viewpoints. Each viewpoint includes the Occupancy Channel and the Semantic Channel. The BEV map is agent-centric, with the agent located at the map center and its heading aligned upward.

897 In this section, we provide examples of the Occupancy Channel and Semantic Channel in our BEV  
 898 map. The BEV map has a resolution of  $161 \times 161$  grid cells, where each cell corresponds to 0.05  
 899 meters in the real world. The agent is always placed at the center of the map and is indicated by an  
 900 upward-facing arrow. When the agent rotates, the entire map is rotated accordingly around the agent  
 901 as the anchor, ensuring that the map remains expressed in the current agent-centric coordinate frame.

902 We also visualize the Semantic Channel of the BEV map. In our implementation, each grid cell in the  
 903 Semantic Channel stores a 40-dimensional one-hot semantic encoding. For visualization purposes,  
 904 we assign a distinct color to each semantic code and render the resulting semantic map as a colored  
 905 grid.

906 The other two images below show the Occupancy Channel and Semantic Channel of the BEV map  
 907 from another viewpoint.

908 We also present an example of the graph representation and provide a detailed explanation of each  
 909 field. The schema below follows the generated `graph/graph.json` files and the fields consumed

910 by CognitiveMap for object lookup, viewpoint planning, memory update, and VLM-readable graph  
911 serialization.

```
912 {
913   "schema_version": "1.1",           // Graph format version.
914   "frame_id": 0,                     // Source frame used to initialize/update the graph.
915   "timestamp": 0,                    // Optional time index; kept for temporal extensions.
916
917   "frames": {
918     "world": "map",                  // World frame used by pose_w and rel_world.
919     "agent": "agent"                 // Agent-local frame used by rel_agent.
920   },
921
922   "conventions": {
923     "units": {
924       "distance": "m",               // Distances are measured in meters.
925       "angle": "deg"                 // Angles are measured in degrees.
926     },
927     "agent_axes": "x_fwd, y_left",   // Agent-relative x is forward, y is left.
928     "bearing_deg": "positive_to_right_of_agent_heading",
929                                     // Positive bearing is to the agent's right.
930     "yaw_deg": "counter_clockwise_from_world_x"
931                                     // Agent/object yaw convention in the map frame.
932   },
933
934   "nodes": [
935     {
936       "node_id": "A0",               // Unique agent node id; CognitiveMap updates this pose.
937       "pose_w": {
938         "x": 2.51,                   // Agent x coordinate in the map/world frame.
939         "y": 4.87,                   // Agent y coordinate in the map/world frame.
940         "yaw": -30.79                // Agent heading; used for bearing and FOV tests.
941       },
942       "fov": {
943         "half_angle_deg": 45.0,      // Half field-of-view angle for visibility labels.
944         "range_m": 4.0               // FOV range for in_fov computation.
945       }
946     },
947     {
948       "node_id": "000",              // Unique object node id used by graph edges/planning.
949       "name": "chair",               // Detector/raw object name; indexed for text matching.
950       "type": "chair",              // Canonical semantic class; also indexed for lookup.
951       "pose_w": {
952         "x": 1.06,                   // Object center x in the map/world frame.
953         "y": 4.60                     // Object center y; used for navigation and merging.
954       },
955       "extent_w": {
956         "w": 0.95,                   // Object footprint width in the world/map frame.
957         "h": 0.85,                   // Object footprint height/depth in the world/map frame.
958         "yaw": 74.78                 // Object footprint orientation.
959       },
960       "score": 0.80,                 // Detection/memory confidence for this object.
961       "semantic_one_hot": [0, 0, 0] // NYU40-style semantic one-hot vector, shortened here.
962     }
963   ],
964
965   "edges": [
966     {
967       "src": "A0",                   // Source node id.
```

```

968     "dst": "000",           // Destination object node id.
969     "type": "agent_object", // Agent-to-object relation.
970     "rel_agent": {
971         "range_m": 1.48,     // Agent-object distance.
972         "bearing_deg": 138.99, // Bearing in the agent frame.
973         "rel_xy": {
974             "x_fwd": -1.12,   // Object x in agent frame; forward is positive.
975             "y_left": -0.97  // Object y in agent frame; left is positive.
976         },
977         "in_fov": false      // Whether object lies inside the agent FOV.
978     },
979     "qualitative": {
980         "sector": "behind",   // Discrete bearing/FOV bucket for text serialization.
981         "proximity": "near"  // Discrete range bucket.
982     },
983     "confidence": 0.80,      // Relation confidence.
984     "last_updated_frame": 0  // Last frame that produced/updated this relation.
985 },
986 {
987     "src": "000",
988     "dst": "001",
989     "type": "oo_spatial",   // Object-object spatial relation.
990     "rel_world": {
991         "delta_xy": {
992             "dx": 2.23,       // dst.x - src.x in the world/map frame.
993             "dy": -0.84      // dst.y - src.y in the world/map frame.
994         },
995         "distance_m": 2.38   // Euclidean object-object distance.
996     },
997     "relation": {
998         "primary": "right_of" // Main spatial predicate used by graph text/reasoning.
999     },
1000    "qualitative": {
1001        "left_right_world": "right", // Coarse world-frame lateral relation.
1002        "near_level": "far"        // Coarse distance bucket.
1003    },
1004    "confidence": 0.52,
1005    "last_updated_frame": 0
1006 }
1007 ],
1008
1009 "views_for_vlm": {
1010     "node_table": [
1011         "A0 agent pose_w=(2.51,4.87,yaw=-31) fov=+/-45deg range=4m",
1012         "000 chair type=chair pose_w=(1.06,4.60) score=0.80"
1013     ],
1014     "edge_table": [
1015         "A0 -> 000 sector=behind range=1.48 bearing=+139 in_fov=0",
1016         "000 <-> 001 far d=2.38 right_world"
1017     ]
1018     // Compact textual graph view tokenized by the VLM.
1019 }
1020 }

```

## 1021 E Case Study

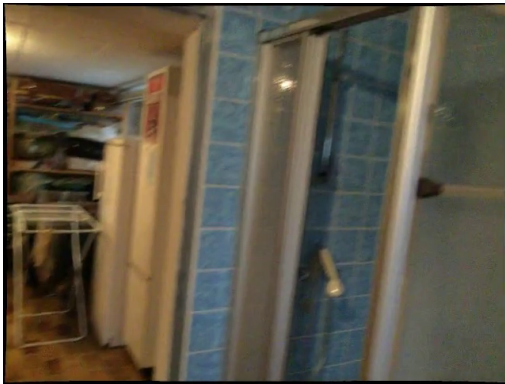
1022 In this section, we provide qualitative case studies to better illustrate how the proposed cognitive-map-  
1023 based interaction pattern works in practice. Specifically, we include one successful example and one



(a) Case1-View1.



(b) Case1-View2.



(c) Case1-View3.



(d) Case1-View4.

Figure 5: Input Video Views.

1024 failure example. The successful examples demonstrate how the world model and memory module  
1025 contribute to spatial reasoning and long-horizon embodied understanding. The failure examples  
1026 highlight current limitations, including unverifiable hallucinations in the world model and cases  
1027 where our interaction pattern does not provide sufficient benefit.

### 1028 E.1 Positive Case 1: Relative Direction Reasoning

1029 The first positive example focuses on relative direction reasoning. This type of question requires the  
1030 agent not only to recognize objects in the scene, but also to reason about their spatial relationships  
1031 from an egocentric point of view. In particular, the question asks the model to infer the relative  
1032 quadrant of one object with respect to the user’s position and facing direction.

#### 1033 Question.

1034 If I am standing by the chair and facing the refrigerator, is the washing machine to  
1035 my front-left, front-right, back-left, or back-right?

1036 The directions refer to the quadrants of a Cartesian plane, where I am standing at  
1037 the origin and facing along the positive  $y$ -axis.

1038 **Input Video.** Here we show four images in the input video. The object chair(left-down side  
1039 of Case1-View1), washing machine(left-down side of Case1-View2) and refrigerator(left side of  
1040 Case1-View3) have been detected during the detection video part. So right now the input memory  
1041 include these three objects. But the problem is that VLM is not good at direction change, so right  
1042 now if you let it answer directly, you will get a wrong answer.

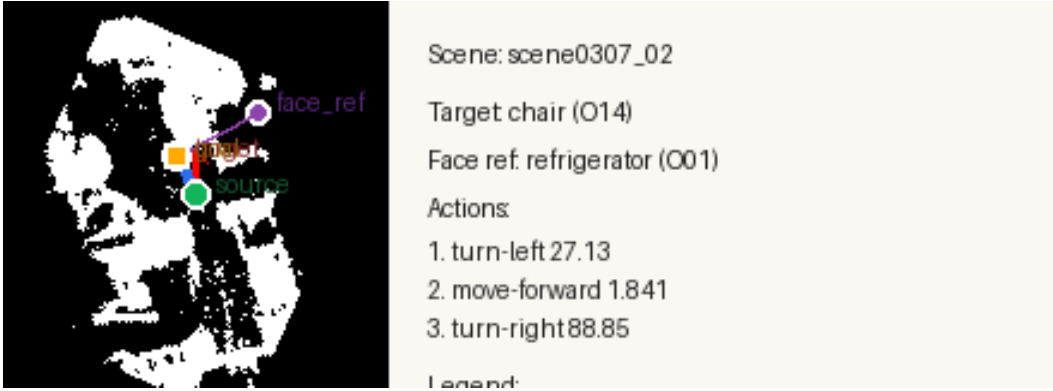


Figure 6: Action list plan visualized on the BEV map.

1043 **VLM-generated action list.** As shown in Fig 6, given the question and the scene graph, the VLM  
 1044 first identifies the three objects explicitly mentioned in the query: the chair, the refrigerator, and the  
 1045 washing machine. These objects are treated as the primary reasoning anchors. In the scene graph, the  
 1046 chair is grounded to node O14, the refrigerator is grounded to node O01, and the washing machine is  
 1047 grounded to node O17. The chair defines the agent’s target position, because the question asks the  
 1048 agent to imagine standing by the chair. The refrigerator defines the orientation reference, because the  
 1049 agent should face the refrigerator after reaching the chair. The washing machine is not selected as  
 1050 the navigation target; instead, it is reserved as the query object whose relative direction should be  
 1051 evaluated after the agent reaches the correct egocentric pose.

1052 The VLM then converts this spatial reasoning objective into an executable action list for the world  
 1053 model. Starting from the last observed camera pose, the system plans a short navigation trajectory  
 1054 toward the chair and then rotates the camera to face the refrigerator. The resulting action list is:

```
1055     turn-left 27.13;
1056     move-forward 1.841;
1057     turn-right 88.85.
```

1058 **World model imagination.** Following the VLM-generated action list, the world model simulates  
 1059 the agent’s future observations step by step. As shown in Fig. ??, the agent first turns left from the  
 1060 last observed pose, then moves forward toward the chair, and finally turns right to align its view  
 1061 with the refrigerator. The imagined observations indicate that the agent successfully reaches the  
 1062 neighborhood of the chair and obtains the intended egocentric pose facing the refrigerator. Under this  
 1063 simulated viewpoint, the washing machine becomes the query object to be localized relative to the  
 1064 agent’s final orientation, enabling the system to infer its direction from the imagined scene rather  
 1065 than from the original observation alone. From the World Model View4, we know the world model  
 1066 now have already generate a good view point.

1067 **Memory Update and VLM answering.** After the world model completes the imagined action  
 1068 sequence, the system updates both the BEV map and the scene graph according to the agent’s  
 1069 simulated final pose. At this stage, the agent’s position and orientation are aligned with the condition  
 1070 specified in the question: the agent is standing near the chair and facing the refrigerator. This  
 1071 pose normalization is important because the original video observations are collected from different  
 1072 viewpoints, whereas the question requires reasoning in a specific egocentric coordinate frame.

1073 In the updated BEV map, the chair is treated as the origin of the local coordinate system, and  
 1074 the direction from the chair to the refrigerator defines the positive  $y$ -axis. The relative position of  
 1075 the washing machine is then evaluated under this newly established egocentric frame. Since the  
 1076 washing machine lies behind the agent and to the agent’s right with respect to the refrigerator-facing  
 1077 orientation, it falls into the back-right quadrant of the Cartesian plane described in the question.

1078 Therefore, after incorporating the world-model imagination result into memory, the VLM determines  
 1079 that the available spatial information is sufficient to answer the question reliably. The final answer is:



(a) World-Model View1.



(b) World-Model View2.



(c) World-Model View3.



(d) World-Model View4.

Figure 7: Input Video Views.

1080 **Back-Right.**

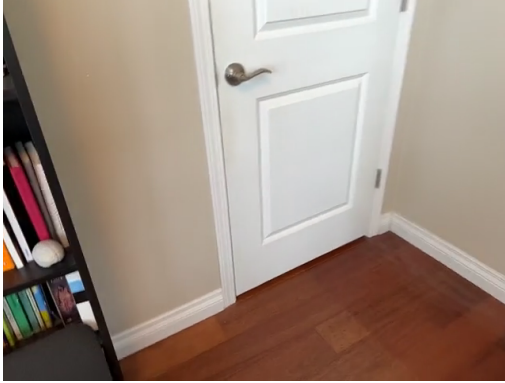
1081 **E.2 Negative Case 1: Hallucinated Spatial Evidence**

1082 This negative example focuses on hallucinated spatial evidence produced during question-conditioned  
 1083 imagination. The question requires the agent to reason about an object or relation that is only  
 1084 partially visible in the input video. The system therefore asks the world model to imagine a missing  
 1085 viewpoint. However, the imagined rollout introduces an unsupported object appearance or an incorrect  
 1086 object-object relation, and this hallucinated evidence propagates into the final VLM answer.

1087 **Question.**

1088 You are a robot beginning at the black desk chair and facing the bookshelf. You  
 1089 want to navigate to the red desk chair. You will perform the following actions  
 1090 (Note: for each [please fill in], choose either 'turn back,' 'turn left,' or 'turn right.):  
 1091 1. [please fill in] 2. Go forward until the red desk chair. You have reached the final  
 1092 destination.

1093 **Input Video.** The input video contains the reliable visual evidence used to initialize the memory.  
 1094 In this case, the relevant anchors are *black desk chair*, *bookshelf*, and *red desk chair*. These objects  
 1095 are either directly detected in the observed frames or inferred from high-confidence memory entries.  
 1096 Importantly, the hallucinated entity or relation that later appears in the imagined rollout is not  
 1097 supported by these observed frames.



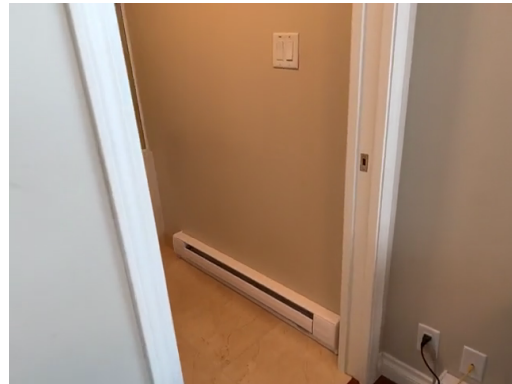
(a) Case2-View1.



(b) Case2-View2.



(c) Case2-View3.



(d) Case2-View4.

Figure 8: Input Video Views.

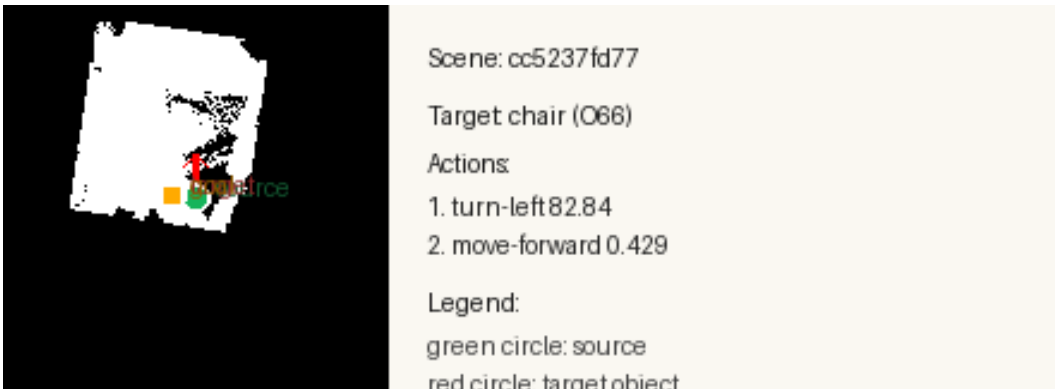
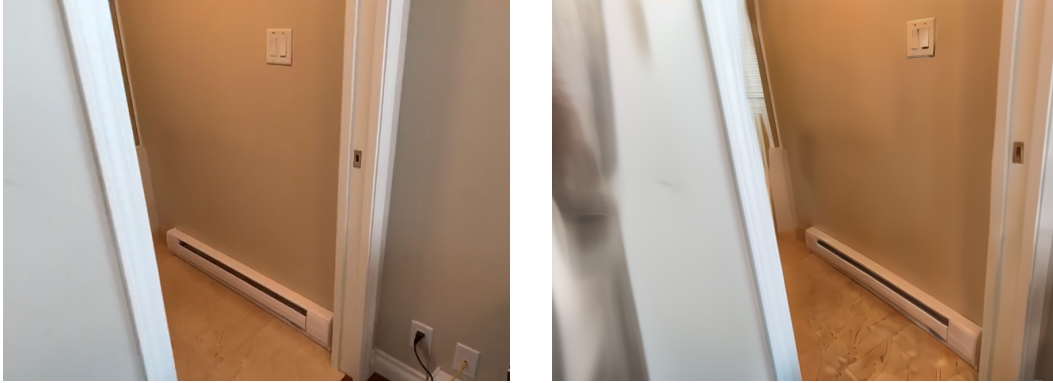


Figure 9: Action list plan visualized on the BEV map.

1098 **VLM-generated action list.** Given the question and the current scene graph, the VLM identifies  
 1099 the objects and spatial relations that are necessary for answering. It then selects an informative target  
 1100 pose intended to reveal the missing evidence. The resulting action list can be summarized as:

1101           turn-left 82.84;  
 1102           move-forward 0.429;

1103 The planned trajectory is reasonable under the current memory: it moves the agent toward *red desk*  
 1104 *chair* and rotates the view toward the region where the relative position of red desk chair with respect



(a) World-Model View1.

(b) World-Model View2.



(c) World-Model View3.



(d) World-Model View4.

Figure 10: Negative Case World Model Imagine Results.

1105 to the bookshelf is expected. The failure does not come from the action list itself, but from the  
 1106 imagined observation generated after executing this plan.

1107 **World model imagination.** Following the action list, the world model generates a future observation  
 1108 sequence. In the early imagined frames, the scene remains consistent with the observed video and the  
 1109 memory. However, after the viewpoint changes to the unobserved region, the world model introduces  
 1110 hallucinated object. This evidence is visually plausible, but it is not grounded in the input video or in  
 1111 the high-confidence scene graph.

1112 **Memory Update and VLM answering.** After the imagined rollout, the candidate update is passed  
 1113 to the memory module. Ideally, confidence-aware memory update should suppress hallucinated  
 1114 content when it conflicts with reliable prior evidence. In this case, however, the hallucinated evidence  
 1115 is difficult to reject because it appears in an unobserved region and does not directly contradict a  
 1116 high-confidence existing entry. As a result, the VLM treats the imagined evidence as if it were valid  
 1117 scene evidence.

1118 The final answer is therefore biased toward the hallucinated observation:

1119 **Turn Right**

1120 The correct answer should be:

1121 **Turn Back**

1122 **Failure analysis.** This case illustrates a hallucination pattern with three stages. First, the question  
 1123 requires evidence from a viewpoint that is weakly covered or not covered by the input video. Second,  
 1124 the world model fills this missing region with plausible but unsupported content. Third, because the

Table 3: Repeated-seed results on VSI-Bench. We report mean  $\pm$  standard deviation over  $N = 3$  world-model seeds. The dataset split, model checkpoint, prompts, and evaluation protocol are fixed across all runs.

Method	Avg.	Rel. Dir.	Route Plan	App. Order
LSM-VLM w/o Mem.	$67.1 \pm 0.3$	$74.1 \pm 0.3$	$51.0 \pm 1.4$	$78.3 \pm 0.1$
LSM-VLM	$69.2 \pm 0.4$	$82.8 \pm 0.3$	$54.9 \pm 1.7$	$78.3 \pm 0.1$

hallucinated content is not directly contradicted by existing memory, it can pass through the update stage and influence the VLM’s final reasoning. This failure suggests that future systems should add stronger verification for imagined evidence, for example by requiring multi-view consistency, checking object persistence across imagined frames, or explicitly marking low-support imagined content as uncertain rather than treating it as observed evidence.

## F Statistical Significance

**Sources of randomness.** We inspect the sources of randomness in our inference pipeline. The VLM evaluation uses greedy decoding with temperature set to 0 by default, and therefore does not introduce active sampling randomness. The world-model video generation is the main stochastic component and uses a fixed random seed by default, set to 42 in our standard evaluation. Therefore, once the dataset split, model checkpoint, prompts, preprocessing pipeline, and evaluation protocol are fixed, the remaining evaluation variability mainly comes from the random seed used by the world model. To quantify this variability, we repeat evaluation by varying only the world-model seed.

**Repeated-seed evaluation.** We evaluate the full LSM-VLM model and the strongest internal ablation, LSM-VLM without confidence-aware recurrent memory update, under multiple world-model seeds. Specifically, we use  $N = 3$  seeds,  $\{0, 42, 3407\}$ , for the world-model generation module while keeping all other factors unchanged. The main paper reports the mean performance over these three seeds, and this appendix provides the corresponding standard deviations for the key metrics supporting our main claims.

For each metric, given scores  $\{x_i\}_{i=1}^N$ , we compute the mean and sample standard deviation as

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i, \quad s = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2}.$$

Unless otherwise stated, all error bars denote one standard deviation across world-model seeds. These error bars capture inference-time variability induced by world-model stochasticity, rather than training-time variability.

**Results and uncertainty analysis.** As shown in Table 3, LSM-VLM consistently outperforms the strongest internal ablation under repeated world-model seeds. The full model improves the overall VSI-Bench average from  $67.1 \pm 0.3$  to  $69.2 \pm 0.4$ , corresponding to a gain of +2.1 points. The improvement is especially large on relative direction reasoning, where LSM-VLM improves from  $74.1 \pm 0.3$  to  $82.8 \pm 0.3$ , and on route planning, where it improves from  $51.0 \pm 1.4$  to  $54.9 \pm 1.7$ . Route planning exhibits larger variance than the other metrics, which is expected because it depends more directly on the stochastic imagined rollouts generated by the world model. In contrast, appearance order remains nearly unchanged across seeds, suggesting that this task is mainly determined by observed temporal evidence rather than world-model stochasticity.

## G Broader Impact

This work studies 3D spatial reasoning for vision-language models in embodied scenes. Potential positive impacts include improving the spatial understanding capabilities of embodied AI systems, such as assistive robots, navigation agents, household service robots, and systems that help users understand complex indoor environments. By explicitly maintaining structured spatial memories and

1162 completing missing visual evidence through question-conditioned imagination, the proposed frame-  
1163 work may make embodied agents more reliable in long-horizon reasoning, cross-view correspondence,  
1164 route planning, and situated question answering.

1165 At the same time, the proposed techniques also have potential negative societal impacts. First, stronger  
1166 spatial reasoning and scene-memory capabilities could be misused for surveillance, unauthorized  
1167 monitoring, or privacy-invasive reconstruction of indoor environments, especially when applied to  
1168 videos of private spaces. Second, because the method uses a generative world model to imagine  
1169 unobserved viewpoints, it may produce hallucinated objects, incorrect spatial relations, or misleading  
1170 scene layouts. If such outputs are used in downstream embodied systems, incorrect reasoning could  
1171 lead to unsafe navigation or inappropriate actions. Third, deployment in real-world assistive or robotic  
1172 systems may amplify dataset biases or perception errors, for example if the underlying detectors,  
1173 VLM backbone, or training data perform unevenly across environments, object categories, or user  
1174 groups.

1175 We discuss several technical design choices that partially mitigate these risks. In particular, LSM-  
1176 VLM does not directly rely on raw imagined images as final evidence; instead, imagined observations  
1177 are converted into structured candidate updates and selectively consolidated through confidence-aware  
1178 memory updates. This design is intended to reduce the risk of contaminating long-term memory  
1179 with hallucinated or low-confidence content. For real-world deployment, additional safeguards  
1180 would be necessary, including privacy-preserving data collection, user consent for indoor video use,  
1181 restricted or gated release in sensitive settings, monitoring for hallucinated spatial evidence, and  
1182 human oversight when the system is used for safety-critical navigation or assistive decision-making.

## 1183 H Assets and License

1184 Our work builds on several publicly available research assets, including pretrained models, perception  
1185 modules, datasets, and benchmarks. We use Qwen3-VL-8B-Instruct as the vision-language backbone,  
1186 Stable Virtual Camera as the camera-conditioned world model for novel-view generation, and VGGT  
1187 together with SAM3 for geometry estimation, segmentation, and perceptual grounding. We cite the  
1188 original papers for all these components in the main paper.

1189 For training and evaluation, we use existing 3D indoor-scene datasets and spatial reasoning bench-  
1190 marks. Specifically, we train the proposed map and graph encoders using the training splits of ScanNet  
1191 and ScanNet++, and evaluate the method on VSI-Bench and SQA3D. ScanNet and ScanNet++ pro-  
1192 vide annotated indoor 3D scenes, while VSI-Bench and SQA3D are used only for benchmarking  
1193 spatial reasoning and situated question answering. We do not introduce a new dataset containing  
1194 human subjects, private user data, or newly collected indoor videos.

1195 All existing assets are used for research purposes and are credited through citations to their original  
1196 papers. We follow the intended research-use setting of these assets and do not redistribute the original  
1197 datasets, pretrained model weights, or third-party code as part of this submission. Users who wish  
1198 to reproduce our experiments should obtain each asset from its official source and comply with the  
1199 corresponding license, terms of use, and data-access requirements. Any released code from our side  
1200 will contain instructions for installing dependencies and preparing these external assets, but will not  
1201 repack third-party data or model weights unless their licenses explicitly permit redistribution.

1202 The main external assets used in this work are:

- 1203 • **Qwen3-VL-8B-Instruct**: used as the pretrained VLM backbone. We cite the original  
1204 Qwen3-VL technical report and follow its model usage terms.
- 1205 • **Stable Virtual Camera**: used as the camera-conditioned world model for generating  
1206 imagined observations from target viewpoints. We cite the original work and use it under its  
1207 released research license and terms.
- 1208 • **VGGT**: used for camera-aware geometry and visual grounding. We cite the original work  
1209 and follow its release terms.
- 1210 • **SAM3**: used for object segmentation and mask extraction. We cite the original work and  
1211 follow its release terms.

- 1212 • **ScanNet and ScanNet++**: used for training the newly introduced map and graph encoders.  
1213 We do not redistribute the datasets and require users to obtain them from the official dataset  
1214 sources under their respective licenses and data-use agreements.
- 1215 • **VSI-Bench and SQA3D**: used as evaluation benchmarks for video-based spatial reasoning  
1216 and 3D situated question answering. We cite the original benchmark papers and follow their  
1217 official usage terms.

1218 To the best of our knowledge, our use of these assets is consistent with their intended academic  
1219 research purposes. We will include license and access information in the released code documentation  
1220 where applicable, and we will ask users to follow the licenses and terms of all upstream assets when  
1221 reproducing or extending our work.

## 1222 I Prompt Examples and Token Examples

### 1223 I.1 Prompt Examples

1224 Our inference pipeline uses a two-stage prompting strategy. In the first stage, the VLM receives an  
1225 egocentric video together with an occupancy map and a scene graph, and predicts an informative  
1226 navigation target for world-model imagination. In the second stage, after the world model renders the  
1227 target view and the memory module updates the scene representation, the VLM answers the question  
1228 using both the imagined observation and the structured memory.

#### 1229 Stage 1: World Model Planning Prompt.

1230 You are given an egocentric video, an occupancy map, and a scene  
1231 graph from the same indoor scene.  
1232 Choose one target object node from the graph that defines where the  
1233 agent should navigate in order to answer the question.  
1234 Do not answer the question directly. Prefer the object that  
1235 provides the most informative viewpoint for reasoning.  
1236  
1237 Question:  
1238 If I am standing by the backpack and facing the door, is the trash  
1239 bin to my front-left, front-right, back-left, or back-right?  
1240  
1241 Options:  
1242 A. front-right  
1243 B. back-left  
1244 C. back-right  
1245 D. front-left  
1246  
1247 Question-grounded anchors:  
1248 - standing object: backpack (start-position anchor)  
1249 - facing object: door (heading anchor)  
1250 - queried object: trash bin  
1251  
1252 Return JSON only:  
1253 {"target\_node\_id": "<graph node id>", "target\_name": "<object name>",  
1254 "reason": "<short reason>"}

#### 1255 Stage 2: Answer Generation Prompt.

1256 This image is the predicted future view after navigating to target  
1257 trash bin (017).  
1258 The navigation goal pose in map coordinates is x=2.14, y=1.36,  
1259 yaw=1.57 rad.  
1260 Answer the question below using this image together with the updated  
1261 occupancy map and scene graph when helpful.  
1262  
1263 Question:  
1264 If I am standing by the backpack and facing the door, is the trash  
1265 bin to my front-left, front-right, back-left, or back-right?

1266  
 1267 Options:  
 1268 A. front-right  
 1269 B. back-left  
 1270 C. back-right  
 1271 D. front-left  
 1272  
 1273 Return only the answer.

## 1274 I.2 Token Examples

1275 To make the multimodal input explicit, we show the Stage 1 planning prompt in both human-readable  
 1276 form and its serialized token form before being fed into the VLM.

### 1277 Human-readable Stage 1 prompt.

1278 <video>  
 1279 <map>  
 1280 <graph>  
 1281 You are given an egocentric video, an occupancy map, and a scene  
 1282 graph from the same indoor scene.  
 1283 Choose one target object node from the graph that defines where the  
 1284 agent should navigate in order to answer the question.  
 1285 Do not answer the question directly. Prefer the object that  
 1286 provides the most informative viewpoint for reasoning.  
 1287  
 1288 Question:  
 1289 If I am standing by the backpack and facing the door, is the trash  
 1290 bin to my front-left, front-right, back-left, or back-right?  
 1291  
 1292 Options:  
 1293 A. front-right  
 1294 B. back-left  
 1295 C. back-right  
 1296 D. front-left  
 1297  
 1298 Question-grounded anchors:  
 1299 - standing object: backpack (start-position anchor)  
 1300 - facing object: door (heading anchor)  
 1301 - queried object: trash bin  
 1302  
 1303 Return JSON only:  
 1304 {"target\_node\_id": "<graph node id>", "target\_name": "<object name>",  
 1305 "reason": "<short reason>"}

1306 **Serialized multimodal prompt.** Let  $K_m$  and  $K_g$  denote the numbers of map tokens and graph  
 1307 tokens, respectively. After preprocessing, the prompt is transformed into:

1308 <|im\_start|>system  
 1309 You are a helpful assistant.<|im\_end|>  
 1310 <|im\_start|>user  
 1311 <|vision\_start|><|video\_pad|><|vision\_end|>  
 1312 <|map\_pad|> ... <|map\_pad|> ( $K_m$  times)  
 1313 <|graph\_pad|> ... <|graph\_pad|> ( $K_g$  times)  
 1314 You are given an egocentric video, an occupancy map, and a scene  
 1315 graph from the same indoor scene.  
 1316 Choose one target object node from the graph that defines where the  
 1317 agent should navigate in order to answer the question.  
 1318 Do not answer the question directly. Prefer the object that  
 1319 provides the most informative viewpoint for reasoning.  
 1320 Question:  
 1321 If I am standing by the backpack and facing the door, is the trash  
 1322 bin to my front-left, front-right, back-left, or back-right?  
 1323 Options:

1324 A. front-right  
 1325 B. back-left  
 1326 C. back-right  
 1327 D. front-left  
 1328 Question-grounded anchors:  
 1329 - standing object: backpack (start-position anchor)  
 1330 - facing object: door (heading anchor)  
 1331 - queried object: trash bin  
 1332 Return JSON only: {"target\_node\_id": "<graph node id>",  
 1333 "target\_name": "<object name>", "reason": "<short reason>"}<|im\_end|>  
 1334 <|im\_start|>assistant

1335 **Compact notation.** The final VLM input can be summarized as

$$X = [\text{SYS}] \oplus [\text{VID}] \oplus [\text{MAP}]^{K_m} \oplus [\text{GRAPH}]^{K_g} \oplus [\text{TEXT\_plan}],$$

1336 where  $\oplus$  denotes concatenation,

$$[\text{VID}] = \langle | \text{vision\_start} | \rangle \langle | \text{video\_pad} | \rangle \langle | \text{vision\_end} | \rangle,$$

1337

$$[\text{MAP}] = \langle | \text{map\_pad} | \rangle, \quad [\text{GRAPH}] = \langle | \text{graph\_pad} | \rangle.$$

1338

## 1339 **NeurIPS Paper Checklist**

### 1340 **1. Claims**

1341 Question: Do the main claims made in the abstract and introduction accurately reflect the  
1342 paper’s contributions and scope?

1343 Answer: [Yes].

1344 Justification: The abstract and Introduction accurately summarize the scope and contribu-  
1345 tions of LSM-VLM, including the dual-memory design, the world-model-based reasoning  
1346 component, and the empirical improvements reported on VSI-Bench and SQA3D. The  
1347 claims are aligned with the experimental results and are limited to the evaluated embod-  
1348 ied/spatial reasoning settings.

1349 Guidelines:

- 1350 • The answer [N/A] means that the abstract and introduction do not include the claims  
1351 made in the paper.
- 1352 • The abstract and/or introduction should clearly state the claims made, including the  
1353 contributions made in the paper and important assumptions and limitations. A [No] or  
1354 [N/A] answer to this question will not be perceived well by the reviewers.
- 1355 • The claims made should match theoretical and experimental results, and reflect how  
1356 much the results can be expected to generalize to other settings.
- 1357 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
1358 are not attained by the paper.

### 1359 **2. Limitations**

1360 Question: Does the paper discuss the limitations of the work performed by the authors?

1361 Answer: [Yes].

1362 Justification: The Conclusion includes a dedicated “Limitations and Future Work” paragraph  
1363 covering backbone coverage, inference efficiency, and the current lack of explicit memory  
1364 conditioning in the world model.

1365 Guidelines:

- 1366 • The answer [N/A] means that the paper has no limitation while the answer [No] means  
1367 that the paper has limitations, but those are not discussed in the paper.
- 1368 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 1369 • The paper should point out any strong assumptions and how robust the results are to  
1370 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
1371 model well-specification, asymptotic approximations only holding locally). The authors  
1372 should reflect on how these assumptions might be violated in practice and what the  
1373 implications would be.
- 1374 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
1375 only tested on a few datasets or with a few runs. In general, empirical results often  
1376 depend on implicit assumptions, which should be articulated.
- 1377 • The authors should reflect on the factors that influence the performance of the approach.  
1378 For example, a facial recognition algorithm may perform poorly when image resolution  
1379 is low or images are taken in low lighting. Or a speech-to-text system might not be  
1380 used reliably to provide closed captions for online lectures because it fails to handle  
1381 technical jargon.
- 1382 • The authors should discuss the computational efficiency of the proposed algorithms  
1383 and how they scale with dataset size.
- 1384 • If applicable, the authors should discuss possible limitations of their approach to  
1385 address problems of privacy and fairness.
- 1386 • While the authors might fear that complete honesty about limitations might be used by  
1387 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
1388 limitations that aren’t acknowledged in the paper. The authors should use their best  
1389 judgment and recognize that individual actions in favor of transparency play an impor-  
1390 tant role in developing norms that preserve the integrity of the community. Reviewers  
1391 will be specifically instructed to not penalize honesty concerning limitations.

1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403  
1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444

### 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [N/A].

Justification: The paper does not present theoretical theorems, formal claims, or proofs; it is an empirical method paper.

Guidelines:

- The answer [N/A] means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

### 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes].

Justification: The Method section specifies the model components and memory update procedure, while the Experiments section describes datasets, benchmarks, training stages, optimization details, baselines, ablations, and runtime analysis needed to understand and reproduce the main claims.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- If the paper includes experiments, a [No] answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).

1445 (d) We recognize that reproducibility may be tricky in some cases, in which case  
1446 authors are welcome to describe the particular way they provide for reproducibility.  
1447 In the case of closed-source models, it may be that access to the model is limited in  
1448 some way (e.g., to registered users), but it should be possible for other researchers  
1449 to have some path to reproducing or verifying the results.

## 1450 5. Open access to data and code

1451 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
1452 tions to faithfully reproduce the main experimental results, as described in supplemental  
1453 material?

1454 Answer: [Yes].

1455 Justification: We will provide anonymized code with instructions for reproducing the main  
1456 experimental results, including environment setup, data preparation, and evaluation scripts,  
1457 in the supplemental material.

1458 Guidelines:

- 1459 • The answer [N/A] means that paper does not include experiments requiring code.
- 1460 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/  
1461 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1462 • While we encourage the release of code and data, we understand that this might not  
1463 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not  
1464 including code, unless this is central to the contribution (e.g., for a new open-source  
1465 benchmark).
- 1466 • The instructions should contain the exact command and environment needed to run to  
1467 reproduce the results. See the NeurIPS code and data submission guidelines ([https://  
1468 neurips.cc/public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1469 • The authors should provide instructions on data access and preparation, including how  
1470 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 1471 • The authors should provide scripts to reproduce all experimental results for the new  
1472 proposed method and baselines. If only a subset of experiments are reproducible, they  
1473 should state which ones are omitted from the script and why.
- 1474 • At submission time, to preserve anonymity, the authors should release anonymized  
1475 versions (if applicable).
- 1476 • Providing as much information as possible in supplemental material (appended to the  
1477 paper) is recommended, but including URLs to data and code is permitted.

## 1478 6. Experimental setting/details

1479 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-  
1480 rameters, how they were chosen, type of optimizer) necessary to understand the results?

1481 Answer: [Yes].

1482 Justification: The Experiments section specifies the training data sources and splits, two-  
1483 stage training protocol, trainable modules, LoRA adaptation, optimizer, learning rates, batch  
1484 size, training steps, benchmarks, metrics, and baseline categories. More training/evaluation  
1485 details will be in appendix.

1486 Guidelines:

- 1487 • The answer [N/A] means that the paper does not include experiments.
- 1488 • The experimental setting should be presented in the core of the paper to a level of detail  
1489 that is necessary to appreciate the results and make sense of them.
- 1490 • The full details can be provided either with the code, in appendix, or as supplemental  
1491 material.

## 1492 7. Experiment statistical significance

1493 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
1494 information about the statistical significance of the experiments?

1495 Answer: [Yes].

1496 Justification: The paper reports benchmark scores and ablation results in the main text,  
1497 and the supplementary material provides additional repeated-run analyses with uncertainty  
1498 estimates for the key experiments supporting the main claims. Specifically, the supplement  
1499 reports error bars computed over repeated runs under the same evaluation protocols, and  
1500 describes the source of variability and the calculation procedure used for these estimates.

1501 Guidelines:

- 1502 • The answer [N/A] means that the paper does not include experiments.
- 1503 • The authors should answer [Yes] if the results are accompanied by error bars, confidence  
1504 intervals, or statistical significance tests, at least for the experiments that support the  
1505 main claims of the paper.
- 1506 • The factors of variability that the error bars are capturing should be clearly stated (for  
1507 example, train/test split, initialization, random drawing of some parameter, or overall  
1508 run with given experimental conditions).
- 1509 • The method for calculating the error bars should be explained (closed form formula,  
1510 call to a library function, bootstrap, etc.)
- 1511 • The assumptions made should be given (e.g., Normally distributed errors).
- 1512 • It should be clear whether the error bar is the standard deviation or the standard error  
1513 of the mean.
- 1514 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
1515 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
1516 of Normality of errors is not verified.
- 1517 • For asymmetric distributions, the authors should be careful not to show in tables or  
1518 figures symmetric error bars that would yield results that are out of range (e.g., negative  
1519 error rates).
- 1520 • If error bars are reported in tables or plots, the authors should explain in the text how  
1521 they were calculated and reference the corresponding figures or tables in the text.

## 1522 8. Experiments compute resources

1523 Question: For each experiment, does the paper provide sufficient information on the com-  
1524 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
1525 the experiments?

1526 Answer: [Yes].

1527 Justification: The Experiments section states that training is performed on 8 H100 GPUs  
1528 and Table 2 reports per-component inference time, call counts, and total inference time. The  
1529 paper also describes efficiency considerations in the ablation section. More details are in  
1530 appendix.

1531 Guidelines:

- 1532 • The answer [N/A] means that the paper does not include experiments.
- 1533 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
1534 or cloud provider, including relevant memory and storage.
- 1535 • The paper should provide the amount of compute required for each of the individual  
1536 experimental runs as well as estimate the total compute.
- 1537 • The paper should disclose whether the full research project required more compute  
1538 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
1539 didn't make it into the paper).

## 1540 9. Code of ethics

1541 Question: Does the research conducted in the paper conform, in every respect, with the  
1542 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

1543 Answer: [Yes].

1544 Justification: The research uses public benchmarks and existing model components for spa-  
1545 tial VQA research and does not involve human-subject experiments, private data collection,  
1546 or deception. The authors have reviewed the NeurIPS Code of Ethics and are not aware of  
1547 deviations.

1548 Guidelines:

- 1549
- The answer [N/A] means that the authors have not reviewed the NeurIPS Code of Ethics.
- 1550
- If the authors answer [No], they should explain the special circumstances that require a deviation from the Code of Ethics.
- 1551
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).
- 1552
- 1553
- 1554

## 10. Broader impacts

1555

1556 Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

1557

1558 Answer: [Yes].

1559 Justification: The Appendix includes a “Broader Impact” section discussing potential benefits for embodied and assistive systems, as well as privacy, surveillance, hallucination, and safety risks with suggested mitigations.

1560

1561

1562 Guidelines:

- The answer [N/A] means that there is no societal impact of the work performed.
  - If the authors answer [N/A] or [No], they should explain why their work has no societal impact or why the paper does not address societal impact.
  - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
  - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate Deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
  - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
  - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).
- 1563
- 1564
- 1565
- 1566
- 1567
- 1568
- 1569
- 1570
- 1571
- 1572
- 1573
- 1574
- 1575
- 1576
- 1577
- 1578
- 1579
- 1580
- 1581
- 1582
- 1583
- 1584

## 11. Safeguards

1585

1586 Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pre-trained language models, image generators, or scraped datasets)?

1587

1588

1589 Answer: [N/A].

1590 Justification: The paper does not release a new high-risk pretrained language model, image generator, scraped dataset, or deployment-ready system. The work uses existing model components in a research pipeline, and no new broadly deployable asset with elevated misuse risk is released in this submission.

1591

1592

1593

1594 Guidelines:

- The answer [N/A] means that the paper poses no such risks.
  - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
  - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- 1595
- 1596
- 1597
- 1598
- 1599
- 1600
- 1601

- 1602           • We recognize that providing effective safeguards is challenging, and many papers do  
1603           not require this, but we encourage authors to take this into account and make a best  
1604           faith effort.

1605 **12. Licenses for existing assets**

1606 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
1607 the paper, properly credited and are the license and terms of use explicitly mentioned and  
1608 properly respected?

1609 Answer: [Yes].

1610 Justification: The paper properly credits the creators of the existing assets used in the work,  
1611 including VSI-Bench, SQA3D, the underlying model components, and implementation  
1612 libraries. All assets are used in accordance with their stated licenses and terms, and we  
1613 do not repackage or redistribute existing datasets under modified licenses. We cited their  
1614 location in the paper.

1615 Guidelines:

- 1616       • The answer [N/A] means that the paper does not use existing assets.
- 1617       • The authors should cite the original paper that produced the code package or dataset.
- 1618       • The authors should state which version of the asset is used and, if possible, include a  
1619       URL.
- 1620       • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 1621       • For scraped data from a particular source (e.g., website), the copyright and terms of  
1622       service of that source should be provided.
- 1623       • If assets are released, the license, copyright information, and terms of use in the  
1624       package should be provided. For popular datasets, `paperswithcode.com/datasets`  
1625       has curated licenses for some datasets. Their licensing guide can help determine the  
1626       license of a dataset.
- 1627       • For existing datasets that are re-packaged, both the original license and the license of  
1628       the derived asset (if it has changed) should be provided.
- 1629       • If this information is not available online, the authors are encouraged to reach out to  
1630       the asset’s creators.

1631 **13. New assets**

1632 Question: Are new assets introduced in the paper well documented and is the documentation  
1633 provided alongside the assets?

1634 Answer: [Yes].

1635 Justification: The submission will include anonymized code as a new asset, together with  
1636 documentation for installation, data preparation, training/evaluation, and reproduction of  
1637 the main results. Release of model checkpoints is still under consideration and will be  
1638 documented if included.

1639 Guidelines:

- 1640       • The answer [N/A] means that the paper does not release new assets.
- 1641       • Researchers should communicate the details of the dataset/code/model as part of their  
1642       submissions via structured templates. This includes details about training, license,  
1643       limitations, etc.
- 1644       • The paper should discuss whether and how consent was obtained from people whose  
1645       asset is used.
- 1646       • At submission time, remember to anonymize your assets (if applicable). You can either  
1647       create an anonymized URL or include an anonymized zip file.

1648 **14. Crowdsourcing and research with human subjects**

1649 Question: For crowdsourcing experiments and research with human subjects, does the paper  
1650 include the full text of instructions given to participants and screenshots, if applicable, as  
1651 well as details about compensation (if any)?

1652 Answer: [N/A].

1653 Justification: The research does not involve crowdsourcing experiments or research with  
1654 human subjects.

1655 Guidelines:

- 1656 • The answer [N/A] means that the paper does not involve crowdsourcing nor research  
1657 with human subjects.
- 1658 • Including this information in the supplemental material is fine, but if the main contribu-  
1659 tion of the paper involves human subjects, then as much detail as possible should be  
1660 included in the main paper.
- 1661 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
1662 or other labor should be paid at least the minimum wage in the country of the data  
1663 collector.

1664 **15. Institutional review board (IRB) approvals or equivalent for research with human**  
1665 **subjects**

1666 Question: Does the paper describe potential risks incurred by study participants, whether  
1667 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
1668 approvals (or an equivalent approval/review based on the requirements of your country or  
1669 institution) were obtained?

1670 Answer: [N/A].

1671 Justification: The research does not involve crowdsourcing or human-subject studies, so  
1672 IRB or equivalent review is not applicable.

1673 Guidelines:

- 1674 • The answer [N/A] means that the paper does not involve crowdsourcing nor research  
1675 with human subjects.
- 1676 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
1677 may be required for any human subjects research. If you obtained IRB approval, you  
1678 should clearly state this in the paper.
- 1679 • We recognize that the procedures for this may vary significantly between institutions  
1680 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
1681 guidelines for their institution.
- 1682 • For initial submissions, do not include any information that would break anonymity (if  
1683 applicable), such as the institution conducting the review.

1684 **16. Declaration of LLM usage**

1685 Question: Does the paper describe the usage of LLMs if it is an important, original, or  
1686 non-standard component of the core methods in this research? Note that if the LLM is used  
1687 only for writing, editing, or formatting purposes and does *not* impact the core methodology,  
1688 scientific rigor, or originality of the research, declaration is not required.

1689 Answer: [Yes].

1690 Justification: The paper explicitly describes the use of VLMs as a core methodological  
1691 component, including Qwen3-VL-8B as the backbone and the VLM roles in pose prediction,  
1692 memory reading, and answer generation.

1693 Guidelines:

- 1694 • The answer [N/A] means that the core method development in this research does not  
1695 involve LLMs as any important, original, or non-standard components.
- 1696 • Please refer to our LLM policy in the NeurIPS handbook for what should or should not  
1697 be described.