

## EXTENDED ABSTRACT

## World2Mind: Cognition Toolkit for Allocentric Spatial Reasoning in Foundation Models

Anonymous CVPR submission

Paper ID \*\*\*\*

## Abstract

001 *Achieving robust spatial reasoning remains a fundamen-*  
 002 *tal challenge for current Multimodal Foundation Models*  
 003 *(MFMs). Existing methods either overfit statistical short-*  
 004 *cuts via 3D grounding data or remain confined to 2D vi-*  
 005 *sual perception, limiting both spatial reasoning accuracy*  
 006 *and generalization in unseen scenarios. Inspired by the*  
 007 *spatial cognitive mapping mechanisms of biological intelli-*  
 008 *gence, we propose **World2Mind**, a training-free spatial in-*  
 009 *telligence toolkit. At its core, World2Mind leverages 3D re-*  
 010 *construction and instance segmentation models to construct*  
 011 *structured spatial cognitive maps, **empowering MFMs to***  
 012 ***proactively acquire targeted spatial knowledge regarding***  
 013 ***interested landmarks and routes of interest.** To provide ro-*  
 014  *bust geometric-topological priors, World2Mind synthesizes*  
 015 *an **Allocentric-Spatial Tree (AST)** that uses elliptical pa-*  
 016 *rameters to model the top-down layout of landmarks ac-*  
 017 *curately. To mitigate the inherent inaccuracies of 3D re-*  
 018 *construction, we introduce a three-stage reasoning chain*  
 019 *comprising tool invocation assessment, modality-decoupled*  
 020 *cue collection, and geometry-semantics interwoven reason-*  
 021 *ing. Extensive experiments demonstrate that World2Mind*  
 022 *boosts the performance of frontier models, such as GPT-*  
 023 *5.2, by 5%~18%. Astonishingly, relying solely on the AST-*  
 024 *structured text, purely text-only foundation models can per-*  
 025 *form complex 3D spatial reasoning, achieving performance*  
 026 *approaching that of advanced multimodal models.*

027 **1. Introduction**

028 Although multimodal foundation models (MFMs) [1, 12,  
 029 24, 28] excel in general visual understanding and cross-  
 030 modal reasoning [29], they struggle significantly in em-  
 031 bodied AI and complex spatial reasoning tasks requir-  
 032 ing physical interaction [15, 17, 18, 26, 34]. This defi-  
 033 ciency stems from their over-reliance on egocentric obser-  
 034 vations and lacking the capacity to abstract global spatial

topology [18, 34], trapping MFMs in an insurmountable  
 “semantic-geometry gap” in tasks like distance estimation,  
 viewpoint transformation, and path planning.

Current efforts to enhance MFMs’ spatial reasoning  
 primarily follow two paradigms. **Training-based meth-**  
**ods** [6, 9, 20] fine-tune models on massive 3D-grounded  
 QA pairs. However, this forces models to overfit statisti-  
 cal shortcuts [14, 25] rather than acquiring genuine spa-  
 tial cognition, leading to poor generalization in out-of-  
 distribution scenarios [33]. Alternatively, introducing ex-  
 plicit 3D modalities [7, 10, 21, 31] exacerbates inter-modal  
 alignment challenges [35]. Meanwhile, recent **tool-based**  
**methods** [8, 19, 36] rely on active rendering under large-  
 scale 3D reconstruction. These are severely bottlenecked by  
 reconstruction quality and remain tethered to low-level vi-  
 sual perception, failing to abstract geometric data into struc-  
 tured semantics for high-level logical reasoning.

Biological Intelligence (BI) offers an ideal blueprint to  
 break the shackles of egocentric observation. Rather than  
 reacting passively to transient visual inputs, the biological  
 brain intrinsically transforms egocentric views into an al-  
 locentric perspective [4]. Supported by place cells in the  
 hippocampus [3, 11, 22] and grid cells in the entorhinal cor-  
 tex [13], mammals construct a global **cognitive map** entirely  
 independent of their egocentric viewpoint [23, 27]. This  
 map forms the cornerstone for strategic mental simulation  
 and advanced reasoning [2].

To bridge this gap between foundation models and  
 BI in spatial representation and reasoning, we propose  
**World2Mind**, a **plug-and-play spatial cognition toolkit**  
 that equips models with human-like mental simulation ca-  
 pabilities. World2Mind integrates an efficient geometry-  
 semantics alignment pipeline, leveraging pre-trained visual  
 geometry [16, 30, 32] and instance segmentation models [5]  
 to extract semantic voxel grids. From these, it constructs  
 two core representations: 1) a Route Cognitive Map for  
 passability prediction, and 2) a Landmark Cognitive Map  
 for object topology. Encapsulated as an accessible toolset,

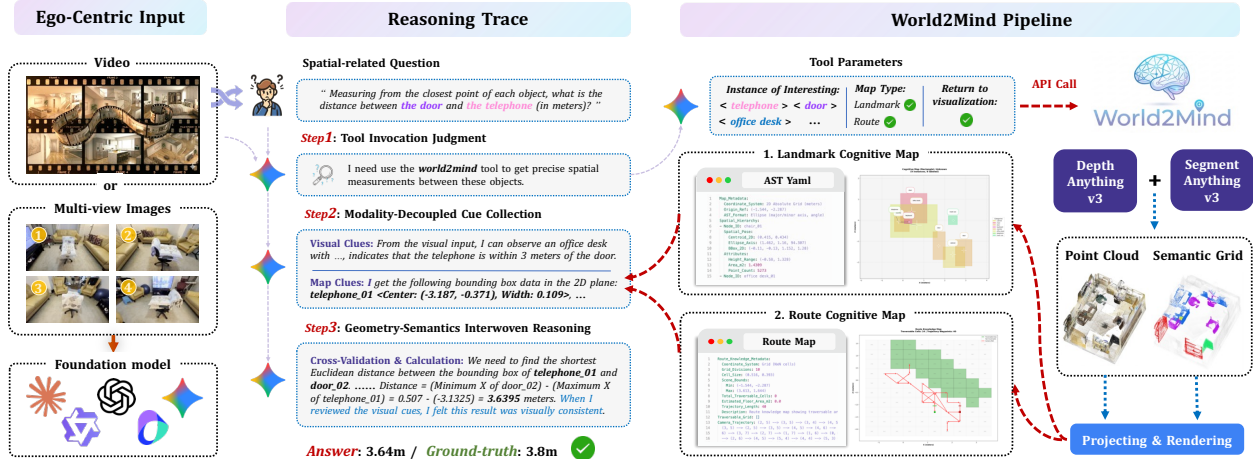


Figure 1. **Overview of foundation models performing allocentric spatial reasoning via the proposed World2Mind toolkit.** Given ego-centric video or multi-view observations, the model first assesses the necessity of tool invocation and subsequently passes key parameters (e.g., instances of interest) to World2Mind to drive the generation of spatial cognitive maps. World2Mind integrates an efficient pipeline for 3D reconstruction and semantic-geometry alignment, returning the required structured spatial knowledge through targeted projection and rendering mechanisms. Furthermore, the model conducts geometry-semantics interwoven reasoning based on both the raw visual observations and the geometric cues provided by World2Mind, ultimately yielding highly reliable answers.

073 World2Mind enables models to dynamically specify parameters (e.g., instances of interest, required spatial knowledge, and map visualizations) to proactively acquire targeted allocentric spatial knowledge on demand.

077 To provide robust geometric-topological priors, we formally define the **Allocentric-Spatial Tree (AST)** as the core spatial representation in World2Mind. The AST is a directed acyclic graph utilizing geometrically stable landmarks (e.g., beds, tables) as core nodes to hierarchically associate surrounding smaller instances. Crucially, to approximate the fuzzy nature of human cognition, the AST models spatial footprints using rectangle-elliptical parameters (bounding boxes, major/minor axes, eccentricity, and rotation angles). These designs equip models with robust, dense, and highly actionable geometric-topological priors.

088 However, merely offering spatial representation is insufficient to guarantee robust reasoning. In complex physical scenarios, reconstruction quality often suffers severe corruption due to occlusions or restricted viewpoints, leading to conflicts with objective geometric laws and raw visual observations. To mitigate this risk, we integrate a rigorous **spatial reasoning chain** into World2Mind: **1) Difficulty Assessment and Tool Invocation**, preventing over-computation on simple superficial queries; **2) Modality-Decoupled Cue Collection**, independently extracting information from ego-centric vision, AST structured text, and map visualizations; and **3) Geometry-Semantics Interwoven Reasoning**, guiding the model to resolve cross-modal conflicts proactively and ultimately yield reliable spatial decisions.

102 Extensive evaluations across various spatial reasoning

benchmarks demonstrate that World2Mind yields stable performance improvements of **6%–18%** for frontier models like GPT-5.2, while maintaining exceptional efficiency and reasoning interpretability. Astonishingly, leveraging the pure, high-density allocentric priors provided by the AST, text-only foundation models can execute complex 3D reasoning directly within their parameter space simply by reading the AST representation, approaching the performance of advanced multimodal models. Our findings offer a highly promising pathway to overcome the spatial cognition bottleneck in foundation models.

## 2. Method Overview

This section details the technical overview of the proposed world2mind, as illustrated in Fig. 1

### 2.1. Geometry-Semantic Alignment Pipeline

Given an ego-centric video sequence or multi-view image set  $\{I_t\}_{t=1}^T$ , our primary objective is to transcend the limitations of 2D vision and construct a robust 3D semantic representation of the physical world.

**1) Depth Estimation & Semantic Extraction.** We employ **Depth Anything V3** [16] for monocular depth estimation, obtaining the depth map  $D_t \in \mathbb{R}^{H \times W}$  and camera pose  $T_t \in SE(3)$ . Concurrently, we utilize **SAM3** [5] to extract open-vocabulary semantic masks  $M_t$  based on a user-specified category list  $\mathcal{C}$ . To suppress the accumulation of long-tail errors inherent in depth estimation, we introduce a dual-level filtering mechanism based on the predicted confidence map  $C_t \in [0, 1]^{H \times W}$ . Specifically, we formulate a

Table 1. Main results on the VSI-Bench [34] benchmark (Tiny subset).

Models	Avg.	Numerical Answer (%)				Multiple-Choice Answer (%)			
		Obj. Count	Abs. Dist.	Obj. Size	Room Size	Rel. Dist.	Rel. Dir.	Route Plan	Appr. Order
<i>Frontier models w/o. World2Mind</i>									
🌀 GPT-5.2	46.7	52.5	34.9	67.5	50.6	42.0	40.7	34.7	51.0
🌟 Claude-4.6-Opus	38.4	46.9	18.5	62.1	26.8	40.0	47.2	34.7	30.6
🔹 Gemini-3-Pro	55.2	47.8	32.1	<b>71.3</b>	55.0	54.0	44.8	57.1	79.6
<i>Frontier models w./ World2Mind</i>									
🌀 GPT-5.2	54.0 (↑7.3)	47.4 (↓5.1)	33.4 (↓1.5)	63.3 (↓4.2)	52.4 (↑1.8)	<b>64.0</b> (↑22.0)	41.1 (↑0.4)	51.0 (↑16.3)	79.6 (↑28.6)
🌟 Claude-4.6-Opus	56.0 (↑17.7)	<b>59.0</b> (↑12.0)	34.3 (↑15.8)	67.3 (↑5.2)	54.8 (↑28.0)	<b>64.0</b> (↑24.0)	62.7 (↑15.6)	<b>65.3</b> (↑30.6)	40.8 (↑10.2)
🔹 Gemini-3-Pro	<b>61.0</b> (↑5.8)	51.8 (↑4.1)	<b>36.8</b> (↑4.7)	57.7 (↓13.5)	<b>62.6</b> (↑7.6)	62.0 (↑8.0)	<b>67.7</b> (↑22.9)	<b>65.3</b> (↑8.2)	<b>83.7</b> (↑4.1)

Table 2. Results on the MindCube-Tiny [18] benchmark.

Models	Avg.	Around	Among	Rotation
<i>Frontier models w/o. World2Mind</i>				
🌀 GPT-5.2	49.9	62.4	45.2	48.5
🌟 Claude-4.6-Opus	48.5	58.8	50.7	29.0
🔹 Gemini-3-Pro	75.1	77.2	68.2	93.0
<i>Frontier models w./ World2Mind</i>				
🌀 GPT-5.2	54.6 (↑4.7)	60.4 (↓2.0)	47.7 (↑2.5)	68.0 (↑19.5)
🌟 Claude-4.6-Opus	62.9 (↑14.4)	82.4 (↑23.6)	60.8 (↑10.1)	45.0 (↑16.0)
🔹 Gemini-3-Pro	<b>81.6</b> (↑6.5)	<b>86.0</b> (↑8.8)	<b>75.8</b> (↑7.6)	<b>93.5</b> (↑0.5)

131 binary validity mask  $\mathcal{V}_t \in \{0, 1\}^{H \times W}$  as follows:

132 
$$\mathcal{V}_t(u, v) = \mathbb{I}(C_t(u, v) > \tau_{\text{pixel}}) \cdot \mathbb{I}(\mu_t > \tau_{\text{frame}}), \quad (1)$$

133 where  $\mu_t = \frac{1}{HW} \sum_{x=1}^H \sum_{y=1}^W C_t(x, y)$  denotes the global  
 134 spatial confidence of frame  $t$ , and  $\mathbb{I}(\cdot)$  is the indicator func-  
 135 tion that returns 1 if the condition is met and 0 otherwise.  
 136 The variables  $\tau_{\text{pixel}}$  and  $\tau_{\text{frame}}$  represent the pixel-level and  
 137 frame-level thresholds, respectively. A pixel is incorporated  
 138 into the subsequent reconstruction only when  $\mathcal{V}_t(u, v) = 1$ .

139 **② Point Cloud Mapping & Density Filtering.** Qualifying  
 140 2D pixels are back-projected into the world coordinate sys-  
 141 tem via the camera intrinsic matrix  $K$ , generating a global  
 142 point cloud  $\mathcal{P} = \{(\mathbf{p}_i, s_i, \mathbf{rgb}_i)\}_{i=1}^M$  carrying semantic la-  
 143 bels  $s_i \in \mathcal{C}$ . Addressing the boundary outliers inherent in  
 144 depth estimation, we propose a core region extraction strat-  
 145 egy: for each point, we calculate its  $K$ -nearest neighbor lo-  
 146 cal density  $\rho_i = \frac{1}{K} \sum_{j \in \mathcal{N}_K(i)} \|\mathbf{p}_i - \mathbf{p}_j\|^{-1}$ , and eliminate  
 147 low-density "tail" points based on density percentiles. This  
 148 yields an exceptionally pure geometry-semantic substrate.

## 149 2.2. Allocentric Cognitive Mapping

150 Inspired by the spatial mapping mechanisms of BI, we dis-  
 151 till the unstructured point cloud into two highly abstract  
 152 cognitive maps, enabling the model to proactively acquire  
 153 spatial knowledge on demand via tool invocation.

154 **① Landmark Cognitive Mapping.** Traditional methods  
 155 rely on ambiguous relative relations or simplified grid repre-  
 156 sentation [18, 34]. To overcome this, we formally define the  
 157 *Allocentric-Spatial Tree (AST)*, which reorganizes spatial  
 158 entities as a directed acyclic graph within an absolute coordi-  
 159 nate system. Specifically, we perform adaptive DBSCAN  
 160 clustering on each semantic category within the point cloud  
 161 to separate distinct instances. For each instance node, the  
 162 AST discards traditional bounding boxes and instead fits  
 163 a minimum bounding ellipse in the top-down view (X-Z  
 164 plane), extracting the centroid  $(x_c, z_c)$ , major and minor  
 165 axes  $a$  and  $b$ , and rotation angle  $\theta$ . This parameterization:  
 166 **1)** significantly enhances robustness against reconstruction  
 167 boundary noise; **2)** perfectly aligns with the fuzzy probabili-  
 168 ty nature of human spatial footprint perception. Output as  
 169 dense structured text (e.g., YAML), the AST explicitly en-  
 170 codes hierarchical containment relationships among entities  
 171 along with multi-dimensional geometric attributes.

172 **② Route Cognitive Mapping.** For navigation-oriented  
 173 tasks, World2Mind also enables extracting the masks of  
 174 traversable categories (e.g., floors), back-projects and vox-  
 175 elizes them, and subsequently partitions them into an  $N \times N$   
 176 grid map on the top-down plane. Combined with the map-  
 177 ping of the camera trajectory sequence  $\{T_t\}$ , this route map  
 178 provides the model with explicit priors regarding passability  
 179 and the human observer’s motion trajectory.

## 180 2.3. Geometry-Semantics Interwoven Reasoning

181 In physical scenarios, 2D visual observations are suscepti-  
 182 ble to occlusions and adverse viewpoints, while 3D recon-  
 183 struction information may contain local errors. To resolve  
 184 potential contradictions between these two modalities, we  
 185 design a rigorous three-stage interwoven reasoning chain.

186 **Stage 1: Tool Invocation Judgement.** To reduce unnec-  
 187 essary computational overhead, the model must first evalu-  
 188 ate the spatial relevance of the query. The model should  
 189 proactively invoke World2Mind only when the task explic-  
 190 itly involves spatial reasoning, such as occlusion inference,  
 191 distance estimation, or path planning.

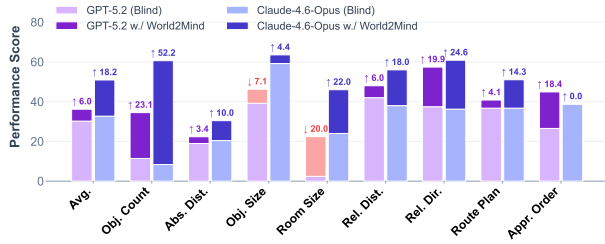


Figure 2. **Performance comparison under the text-only model (“blind”) setting.** we report the performance gap on the VSI-Bench (Tiny subset) between foundation models relying solely on commonsense reasoning and those leveraging world2mind to acquire structured spatial knowledge for allocentric reasoning.

192 **Stage 2: Modality-Decoupled Cue Collection.** We force  
193 the model to extract information independently to prevent  
194 early modality bias. The model must simultaneously gather  
195 corroborating evidence from three independent sources:  
196 egocentric vision, the AST text returned by World2Mind,  
197 and optional 2D top-down map visualizations.

198 **Stage 3: Conflict Resolution and Cross-Validation.** This is  
199 the crux of the reasoning chain. The model needs to proac-  
200 tively coordinate evidence across different modalities and  
201 identify cross-modal conflicts, such as missing objects due  
202 to visual truncation or coordinate drift caused by depth er-  
203 rors. By cross-validating the objective geometric param-  
204 eters of the AST against subjective visual appearances, the  
205 model can dynamically weigh the credibility between visual  
206 illusions and reconstruction artifacts, ultimately outputting  
207 highly reliable and logically interpretable decisions.

## 208 3. Experiment

### 209 3.1. Experimental Setup

210 We conduct evaluations on two challenging spatial reason-  
211 ing benchmarks: VSI-Bench [34], which focuses on video-  
212 based reasoning in real-world physical scenes, and Mind-  
213 Cube [18], which emphasizes multi-view cognitive map-  
214 ping and mental simulation. Our evaluation primarily tar-  
215 gets the most frontier multimodal foundation models, in-  
216 cluding GPT, Claude, and Gemini, given their exceptional  
217 proficiency in tool invocation and instruction-following.

### 218 3.2. Main Results and Analysis

219 **Evaluation on VSI-Bench.** As shown in Tab. 1, the seam-  
220 less integration of World2Mind yields a substantial leap in  
221 the average performance (Avg.) across all frontier mod-  
222 els. Specifically, GPT-5.2 achieves a 7.3% improvement,  
223 while Claude-4.6-Opus achieves a staggering 17.7% im-  
224 provement. In-depth analysis of the subtasks reveals that the  
225 performance gains are most striking in task categories that  
226 rely heavily on allocentric priors. For instance, on Relative  
227 Direction, Route Planning, and Relative Distance, the per-

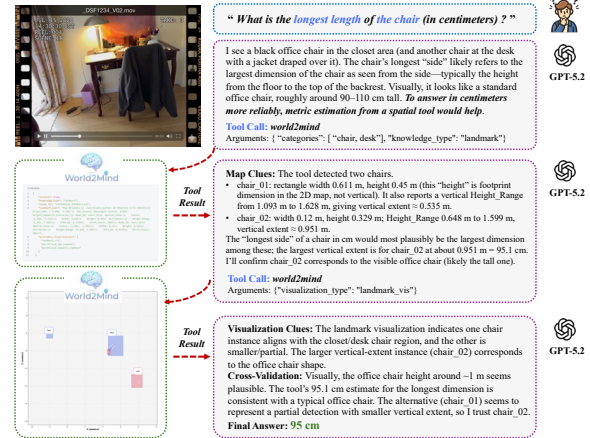


Figure 3. Complete reasoning trace under World2Mind.

228 formance of Claude-4.6-Opus skyrockets by 15.6%,  
229 and 24.0%, respectively. This compellingly demonstrates  
230 the critical role of the allocentric spatial knowledge pro-  
231 vided by World2Mind in bridging the spatial reasoning gap.  
232 **Evaluation on MindCube.** The results in Tab. 1 further  
233 corroborate the universality and robustness of our frame-  
234 work in sparse multi-view inputs. Even for Gemini-3-Pro,  
235 whose native spatial reasoning capability is already top-tier  
236 (with a baseline Avg. of 75.1%), World2Mind successfully  
237 shatters its performance ceiling, pushing its average accu-  
238 racy to 81.6% (+6.5%). Notably, in tasks like “Rotation”  
239 that severely test 3D spatial imagination, the model achieves  
240 a remarkable performance breakthrough (GPT-5.2 improves  
241 by 19.5%) due to its ability to perform logical deduction  
242 grounded in the AST.

### 243 3.3. Ablation and Case Study

244 To explore the limits of how structured text of AST em-  
245 powers the spatial cognition of large language models,  
246 we follow [34] to conduct ablation studies under the text-  
247 only (“blind”) setting (see Fig. 2). When visual image in-  
248 puts are completely stripped away, foundation models that  
249 rely solely on commonsense priors degrade to near-random  
250 guessing on spatial tasks. Astonishingly, however, when  
251 equipped with World2Mind, both GPT-5.2 and Claude-  
252 4.6-Opus exhibit a remarkable performance rebound in the  
253 “blind” state. On core reasoning tasks such as Object Size  
254 and Route Planning, their scores closely approach those  
255 achieved with full visual inputs. This profound finding *indi-*  
256 *cates that pure, high-quality allocentric geometric priors*  
257 *are entirely sufficient to ignite powerful 3D mental recon-*  
258 *struction and simulation capabilities of foundation models*  
259 *under the text-based reasoning.* Furthermore, we visualize  
260 the complete reasoning traces powered by World2Mind in  
261 Fig. 3, which clearly demonstrate that the interwoven rea-  
262 soning process exhibits exceptional robustness and logical  
263 interpretability when resolving cross-modal conflicts.

264

## References

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

- [1] Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei Ding, Chang Gao, Chunjiang Ge, et al. Qwen3-vl technical report. *arXiv preprint arXiv:2511.21631*, 2025. 1
- [2] Jacob LS Bellmund, Peter Gärdenfors, Edvard I Moser, and Christian F Doeller. Navigating cognition: Spatial codes for human thinking. *Science*, 362(6415):eaat6766, 2018. 1
- [3] Nicola J Broadbent, Larry R Squire, and Robert E Clark. Spatial memory, recognition memory, and the hippocampus. *Proceedings of the National Academy of Sciences*, 101(40):14515–14520, 2004. 1
- [4] Neil Burgess. Spatial memory: how egocentric and allocentric combine. *Trends in cognitive sciences*, 10(12):551–557, 2006. 1
- [5] Nicolas Carion, Laura Gustafson, Yuan-Ting Hu, Shoubhik Debnath, Ronghang Hu, Didac Suris, Chaitanya Ryal, Kalyan Vasudev Alwala, Haitham Khedr, Andrew Huang, et al. Sam 3: Segment anything with concepts. *arXiv preprint arXiv:2511.16719*, 2025. 1, 2
- [6] Boyuan Chen, Zhuo Xu, Sean Kirmani, Brain Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14455–14465, 2024. 1
- [7] Pingyi Chen, Yujing Lou, Shen Cao, Jinhui Guo, Lubin Fan, Yue Wu, Lin Yang, Lizhuang Ma, and Jieping Ye. Sd-vlm: Spatial measuring and understanding with depth-encoded vision-language models. *arXiv preprint arXiv:2509.17664*, 2025. 1
- [8] Siyi Chen, Mikaela Angelina Uy, Chan Hee Song, Faisal Ladhak, Adithyavairavan Murali, Qing Qu, Stan Birchfield, Valts Blukis, and Jonathan Tremblay. Spacetools: Tool-augmented spatial reasoning via double interactive rl. *arXiv preprint arXiv:2512.04069*, 2025. 1
- [9] An-Chieh Cheng, Hongxu Yin, Yang Fu, Qiushan Guo, Ruihan Yang, Jan Kautz, Xiaolong Wang, and Sifei Liu. Spatialrgpt: Grounded spatial reasoning in vision-language models. *Advances in Neural Information Processing Systems*, 37:135062–135093, 2024. 1
- [10] Erik Daxberger, Nina Wenzel, David Griffiths, Haiming Gang, Justin Lazarow, Gefen Kohavi, Kai Kang, Marcin Eichner, Yinfei Yang, Afshin Dehghan, et al. Mm-spatial: Exploring 3d spatial understanding in multimodal llms. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7395–7408, 2025. 1
- [11] Howard Eichenbaum. The role of the hippocampus in navigation is memory. *Journal of neurophysiology*, 117(4):1785–1796, 2017. 1
- [12] Google. Gemini 3.1 pro: Best for complex tasks and bringing creative concepts to life. <https://deepmind.google/models/gemini/pro/>, 2026. 1
- [13] Torkel Hafting, Marianne Fyhn, Sturla Molden, May-Britt Moser, and Edvard I Moser. Microstructure of a spatial map in the entorhinal cortex. *Nature*, 436(7052):801–806, 2005. 1
- [14] Jiaxin Huang, Ziwen Li, Hanlve Zhang, Runnan Chen, Xiao He, Yandong Guo, Wenping Wang, Tongliang Liu, and Mingming Gong. Surprise3d: A dataset for spatial understanding and reasoning in complex 3d scenes. *arXiv preprint arXiv:2507.07781*, 2025. 1
- [15] Shuai Huang, Wenxuan Zhao, and Jun Gao. Si-bench: Benchmarking social intelligence of large language models in human-to-human conversations. *arXiv preprint arXiv:2510.23182*, 2025. 1
- [16] Haotong Lin, Sili Chen, Junhao Liew, Donny Y Chen, Zhenyu Li, Guang Shi, Jiashi Feng, and Bingyi Kang. Depth anything 3: Recovering the visual space from any views. *arXiv preprint arXiv:2511.10647*, 2025. 1, 2
- [17] Jingli Lin, Runsen Xu, Shaohao Zhu, Sihan Yang, Peizhou Cao, Yunlong Ran, Miao Hu, Chenming Zhu, Yiman Xie, Yilin Long, et al. Mmsi-video-bench: A holistic benchmark for video-based spatial intelligence. *arXiv preprint arXiv:2512.10863*, 2025. 1
- [18] Fangzheng Liu, Don Derek Haddad, and Joe Paradiso. Mindcube: an interactive device for gauging emotions. In *Adjunct Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology*, pages 1–2, 2024. 1, 3, 4
- [19] Zhanpeng Luo, Ce Zhang, Silong Yong, Cunxi Dai, Qianwei Wang, Haoxi Ran, Guanya Shi, Katia Sycara, and Yaqi Xie. pypatial: Generating 3d visual programs for zero-shot spatial reasoning. *arXiv preprint arXiv:2603.00905*, 2026. 1
- [20] Wufei Ma, Yu-Cheng Chou, Qihao Liu, Xingrui Wang, Celso de Melo, Jianwen Xie, and Alan Yuille. Spatialreasoner: Towards explicit and generalizable 3d spatial reasoning. *arXiv preprint arXiv:2504.20024*, 2025. 1
- [21] Zhenhua Ning, Zhuotao Tian, Shaoshuai Shi, Guangming Lu, Daojing He, Wenjie Pei, and Li Jiang. Enhancing spatial reasoning in multimodal large language models through reasoning-based segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7851–7860, 2025. 1
- [22] John O’Keefe and Jonathan Dostrovsky. The hippocampus as a spatial map: preliminary evidence from unit activity in the freely-moving rat. *Brain research*, 1971. 1
- [23] John O’keefe and Lynn Nadel. *The hippocampus as a cognitive map*. Oxford university press, 1978. 1
- [24] OpenAI. Gpt-4v(ision) system card. [https://cdn.openai.com/papers/GPTV\\_System\\_Card.pdf](https://cdn.openai.com/papers/GPTV_System_Card.pdf), 2023. 1
- [25] Jianing Qi, Jiawei Liu, Hao Tang, and Zhigang Zhu. Beyond semantics: Rediscovering spatial awareness in vision-language models. *arXiv preprint arXiv:2503.17349*, 2025. 1
- [26] Santhosh Kumar Ramakrishnan, Erik Wijmans, Philipp Kraehenbuehl, and Vladlen Koltun. Does spatial cognition emerge in frontier models? *arXiv preprint arXiv:2410.06468*, 2024. 1
- [27] Daniela Schiller, Howard Eichenbaum, Elizabeth A Buffalo, Lila Davachi, David J Foster, Stefan Leutgeb, and Charan Ranganath. Memory and space: towards an understanding of the cognitive map. *Journal of Neuroscience*, 35(41):13904–13911, 2015. 1

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

- 378 [28] Aaditya Singh, Adam Fry, Adam Perelman, Adam Tart, Adi  
379 Ganesh, Ahmed El-Kishky, Aidan McLaughlin, Aiden Low,  
380 AJ Ostrow, Akhila Ananthram, et al. Openai gpt-5 system  
381 card. *arXiv preprint arXiv:2601.03267*, 2025. 1
- 382 [29] Zhaochen Su, Peng Xia, Hangyu Guo, Zhenhua Liu, Yan Ma,  
383 Xiaoye Qu, Jiaqi Liu, Yanshu Li, Kaide Zeng, Zhengyuan  
384 Yang, et al. Thinking with images for multimodal reasoning:  
385 Foundations, methods, and future frontiers. *arXiv preprint*  
386 *arXiv:2506.23918*, 2025. 1
- 387 [30] Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea  
388 Vedaldi, Christian Rupprecht, and David Novotny. Vgg: Vi-  
389 sual geometry grounded transformer. In *Proceedings of the*  
390 *Computer Vision and Pattern Recognition Conference*, pages  
391 5294–5306, 2025. 1
- 392 [31] Yuxin Wang, Lei Ke, Boqiang Zhang, Tianyuan Qu, Hanxun  
393 Yu, Zhenpeng Huang, Meng Yu, Dan Xu, and Dong  
394 Yu. N3d-vlm: Native 3d grounding enables accurate spa-  
395 tial reasoning in vision-language models. *arXiv preprint*  
396 *arXiv:2512.16561*, 2025. 1
- 397 [32] Yifan Wang, Jianjun Zhou, Haoyi Zhu, Wenzheng Chang,  
398 Yang Zhou, Zizun Li, Junyi Chen, Jiangmiao Pang, Chunhua  
399 Shen, and Tong He. Permutation-equivariant visual geome-  
400 try learning. *arXiv preprint arXiv:2507.13347*, 2025. 1
- 401 [33] Mingrui Wu, Zhaozhi Wang, Fangjinhua Wang, Jiaolong  
402 Yang, Marc Pollefeys, and Tong Zhang. From indoor to open  
403 world: Revealing the spatial reasoning gap in mllms. *arXiv*  
404 *preprint arXiv:2512.19683*, 2025. 1
- 405 [34] Jihan Yang, Shusheng Yang, Anjali W Gupta, Rilyn Han,  
406 Li Fei-Fei, and Saining Xie. Thinking in space: How mul-  
407 timodal large language models see, remember, and recall  
408 spaces. In *Proceedings of the Computer Vision and Pattern*  
409 *Recognition Conference*, pages 10632–10643, 2025. 1, 3, 4
- 410 [35] Weichen Zhang, Ruiying Peng, Chen Gao, Jianjie Fang, Xin  
411 Zeng, Kaiyuan Li, Ziyou Wang, Jinqiang Cui, Xin Wang,  
412 Xinlei Chen, et al. The point, the vision and the text: Does  
413 point cloud boost spatial reasoning of large language mod-  
414 els? *arXiv preprint arXiv:2504.04540*, 2025. 1
- 415 [36] Zaibin Zhang, Yuhan Wu, Lianjie Jia, Yifan Wang, Zhongbo  
416 Zhang, Yijiang Li, Binghao Ran, Fuxi Zhang, Zhuohan Sun,  
417 Zhenfei Yin, et al. Think3d: Thinking with space for spatial  
418 reasoning. *arXiv preprint arXiv:2601.13029*, 2026. 1