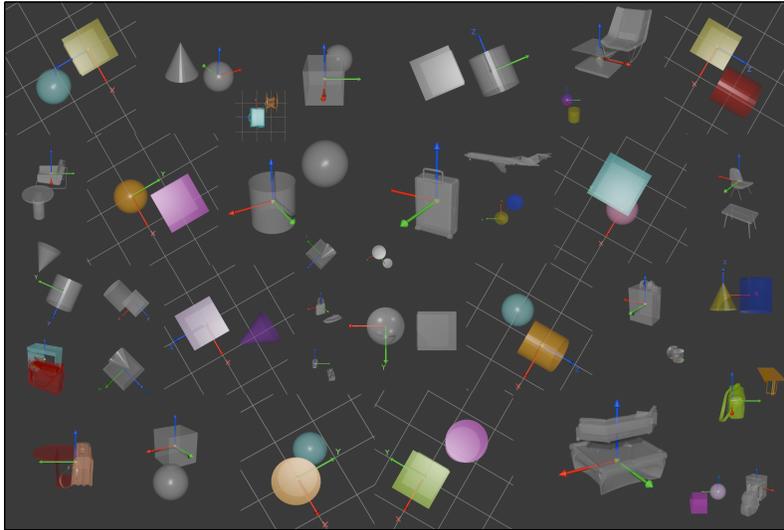

CAN VLMS REASON THROUGH MULTIPLE VIEWS?

Anonymous authors

Paper under double-blind review



ABSTRACT

Recent advances in Vision–Language Models (VLMs) have opened new possibilities for complex spatial reasoning. Benchmarks for VLMs largely assess single- or limited-view perception, leaving untested the core ability to *integrate* observations across viewpoints into a coherent 3D understanding. We introduce **MultiView-Bench**, a benchmark expressly designed to evaluate multi-view integration for holistic 3D scene comprehension. MultiView-Bench is paired with a highly *extensible* data-generation pipeline that supports plug-and-play 3D assets (synthetic or real), configurable distractors, and flexible camera positions and orientations, enabling researchers to readily instantiate new datasets by swapping assets or altering viewpoint configurations. Beyond benchmarking, MultiView-Bench serves as a *fundamental diagnostic* that VLMs should pass before being deployed as agents operating 3D software for downstream tasks such as part assembly for mechanical engineering. We evaluate a broad set of frontier VLMs and uncover consistent failure modes: strong performance on 2D planar relations from a single image, but marked difficulty with 3D spatial relations and with aggregating information across views. We further identify biases in VLMs, including handling unconventional axis directions and sensitivity to object colorways and texture variations. Acknowledging these limitations, we propose **ViewNavigator**, a multi-agent framework that actively selects informative viewpoints, perceive, and fuses multi-view evidence through belief-updating. ViewNavigator improves the performances of diverse base models on MultiView-Bench by more than 50%. MultiView-Bench and its extensible pipeline are designed to equip researchers with a principled testbed for strengthening VLMs’ 3D scene understanding, paving the way for more capable VLM-based agents that can support a wide range of downstream 3D tasks. Project webpage: <https://hantaozhangrichard.github.io/MultiView-Bench/>

1 INTRODUCTION

Recent advances in Large Language Models (LLMs) (Brown et al., 2020; OpenAI, 2023; 2024; Touvron et al., 2023) and Vision-Language Models (VLMs) (Radford et al., 2021; Li et al., 2022; Google, 2023; Dai et al., 2023) have demonstrated remarkable progress in complex perceptual and reasoning tasks, including spatial navigation (Mirowski et al., 2018; Du et al., 2023; Yamada et al., 2023) and image understanding (Dosovitskiy et al., 2021). Their strong generalization capabilities, coupled with emergent reasoning skills, make them compelling candidates for *cognitive systems* that integrate perception and strategic planning (Bubeck et al., 2023). When equipped with appropriate tools and scaffolding, such systems have shown promise in robotics control (Brohan et al., 2023; Liu et al., 2023), 3D modeling (Hu et al., 2024; Gu et al., 2025), and image editing (Huang et al., 2024).

However, effectively solving many of these tasks fundamentally depends on the ability to *perceive and reason about scenes from multiple viewpoints* (Edelman, 1998; Bühlhoff & Edelman, 1992). Humans naturally perform multi-angle observations to construct coherent mental models of objects, resolving perceptual ambiguities that arise from single viewpoints (Shepard & Metzler, 1971; Palmer, 1999). This ability is crucial when assembling complex objects, where each component must be rotated and inspected from multiple viewpoints to determine how it connects with others. In contrast, a single static image often fails to convey critical structural or relational details necessary for accurate reasoning and manipulation, underscoring the importance of *multi-view perception* in spatial cognition (Marr, 2010; Kosslyn, 1994).

A possible workaround involves geometric representations such as point clouds, meshes, or voxels (Qi et al., 2017a; Wu et al., 2015; Mescheder et al., 2019), which encode precise 3D coordinates and shapes. However, processing such low-level geometric data typically requires specialized encoders (Qi et al., 2017b; Wang et al., 2019) and lacks the broad generalization of LLM/VLM-based approaches. Moreover, these representations diverge from the modality of human visual input, limiting their interpretability for *human-readable reasoning* (Tarr & Bühlhoff, 1998).

Current multi-view spatial reasoning benchmarks (Wang et al.; Daxberger et al.; Yin et al., 2025; Zhang et al., 2025) primarily assess egocentric spatial reasoning such as perspective-taking, view-dependent navigation, or dynamic viewpoint transformation. These benchmarks evaluate a model’s ability to map pixels from one 2D plane to another or to predict the visual consequences of camera motion without targeting view-invariant (world-centric) spatial understanding and holistic 3D grounding. Such world-centric reasoning ability is a prerequisite for mechanical engineering and operating 3D modeling software (e.g. Blender) where the agent needs to manipulate in global space, independent of where the user is looking.

In this work, we fill this gap by introducing **MultiView-Bench**, a multi-view spatial reasoning benchmark that explicitly requires the model to integrate multiple viewpoints to reconstruct a static, global mental model of the scene. We ground this by providing a visible, fixed global coordinate system, compelling the model to decouple object positioning from the camera’s transient perspective. Our contributions are as follows:

- We introduce **MultiView-Bench**, a diagnostic benchmark for evaluating VLMs’ ability to understand the 3D global coordinate system in 3D modeling software and integrate multi-view observations into a coherent world-centric 3D scene understanding. Our benchmark serves as a prerequisite test VLMs should pass before deployment to mechanical engineering and 3D modeling. The main benchmark contains 5 task variants each with 100 tasks where we conduct a systematic evaluation of state-of-the-art VLMs.
- Alongside the benchmark, we provide a flexible **data generation pipeline** that allows researchers to easily extend the dataset with new 3D assets, task variants, and viewpoint configurations. We created 20 more different task variants each with 100 tasks and revealed key failure modes, biases, and limitations in VLMs’ multi-view spatial reasoning capabilities through controlled experiments.
- We propose **ViewNavigator**, a multi-agent framework that models perception, planning, and belief-updating. ViewNavigator consistently and significantly enhances the performance of underlying VLMs on MultiView-Bench, demonstrating its potential as a plug-and-play reasoning scaffold for future vision–language systems.

2 RELATED WORK

VLM Benchmarks. Numerous benchmarks have been developed to evaluate VLM capabilities. Foundational datasets such as ShapeNet (Chang et al., 2015) and ModelNet (Wu et al., 2015) focus on 3D object recognition and reconstruction. CLEVR (Johnson et al., 2017) targets compositional reasoning in synthetic and real-world images. More recent work has extended to spatial reasoning: SpatialRGPT (Cheng et al., 2024) and OmniSpatial (Jia et al., 2025) incorporate perspective-taking and dynamic reasoning, but remain primarily *single-view* in nature.

While multi-view benchmarks exist, they often prioritize different cognitive tasks. InternSpatial (Deng et al., 2025) includes is restricted to rotation estimation. ViewSpatial-Bench (Li et al., 2025), SITE (Wang et al.), SPAR (Zhang et al., 2025) and MindCube (Yin et al., 2025) focus on egocentric–allocentric viewpoints transformations and navigation dynamics that requires sequential feature matching Crucially, these benchmarks do not require the model to integrate multiple concurrent views into a single, coherent global representation. Finally, benchmarks like IR3DBench (Liu et al., 2025) and MM-Spatial (Daxberger et al.) test precise 3D coordinate estimation. However, IR3DBench is restricted to single views which allows multiple plausible configurations to produce the same image and MM–Spatial (Daxberger et al.) relies on explicit depth maps, which hinders scalability compared to pure RGB image reasoning. We summarize the distinctions to those image-only 3D multi-view VLM benchmarks in Table 1.

3D Spatial Reasoning with LLM/VLM Agents. Agentic systems leveraging LLMs and VLMs have recently demonstrated impressive performance on 3D tasks, including open-world gaming (Wang et al., 2023; Yao et al., 2023), procedural scene generation (Hu et al., 2024; Huang et al., 2024), and LEGO assembly (Yamada et al., 2024; Pun et al., 2025). In such systems, LLMs often act as *planners*, akin to the prefrontal cortex in the brain (Stokes et al., 2021), while VLMs serve as perceptual modules that transform raw visual inputs into structured descriptions. These pipelines work well when single-image perception suffices (e.g., block-based abstractions in Minecraft), but break down when tasks require precise geometric reasoning over multiple views (Chen et al., 2024; Hong et al., 2023).

While some works attempt to enhance VLM 3D reasoning (Cheng et al., 2024; Chen et al., 2024; Hong et al., 2023), they typically focus on VQA-style setups without extending to real-world applications that require integrated 3D perception and planning. MultiView-Bench is designed precisely to call for the awareness of VLMs’ limitations in multi-view integration and to serve as a selection criterion when building VLM-based agents for real-world embodied intelligence as well as for operating 3D software in 3D asset generation and mechanical engineering.

3 MOTIVATION

To illustrate the necessity and practical importance of our benchmark, we motivate our study through a real-world furniture part assembly task. In this setting, a collection of labeled components (e.g., legs, table tops, backrests) must be connected and arranged to form a functional piece of furniture. Solving this task naturally demands multi-view perception, 3D spatial reasoning, and common-sense knowledge about how objects are typically used and combined. At a minimum, three core abilities are required:

1. **Part Identification and Semantic Reasoning:** Identifying distinct furniture components and inferring their functionality using combined visual and semantic reasoning.
2. **Assembly Planning:** Formulating a coherent and executable plan by determining the correct assembly sequence and how components interconnect.
3. **Assembly Execution:** Precisely placing each component and executing physical assembly actions.

From experimentation, VLMs have considerable potential to tackle the first two stages, effectively operating as the cognitive system or *brain* of a robotic agent. The subsequent execution can then rely on specialized robotic control modules that function as the robot’s *motor system*.

Table 1: Distinctions between existing image-only multi-view VLM benchmarks (Wang et al.; Yin et al., 2025; Zhang et al., 2025) and MultiView-Bench.

Feature	Existing Multi-View Benchmarks	MultiView-Bench
Environment	Photorealistic / Real-World Scenes High visual noise, complex textures, messy lighting.	Synthetic / Blender-GUI Clean, schematic geometry with visible global axes/grids.
Reference Frame	View-dependent Relative to the camera (Left, Right, Forward).	View-invariant Relative to a fixed global coordinate system.
Reasoning Type	Transformation $View_1 \rightarrow \text{Transformation} \rightarrow View_2$	Intersection $View_1 + View_2 + \dots \rightarrow \text{Global State}$
Primary Goal	End-to-End Performance Testing general robustness in noisy, real-world environments.	Diagnostic Prerequisite "Unit testing" precise spatial logic before deployment.
Domain	Navigation & Embodied AI "Where do I go next?" "What will I see next?"	Engineering & Manipulation "Where is this part located exactly?" "Where should I put the next block?"

While recent LLM-based approaches have shown strong capabilities in environments such as Minecraft or Blender (Yamada et al., 2024; Pun et al., 2025; Wang et al., 2023), where assemblies involve standardized, uniform blocks. Such blocks lend themselves to lossless bounding-box (length, width, height) representation and these LLM-based agents' success largely depends on simple geometric descriptions and mathematical reasoning. However, real-world furniture assembly poses significantly greater complexity due to irregular, non-convex shapes that defy concise, lossless linguistic descriptions. Hence, purely semantic or bounding-box representations are inadequate for precise assembly tasks involving intricate real-world parts.

Figure 1 illustrates typical furniture assembly tasks that explicitly require multi-view perception to comprehend and accurately reason about the spatial configuration of parts. These examples underscore the critical need for robust VLMs capable of integrating information across multiple visual perspectives to build accurate internal 3D understanding.

4 MULTIVIEW-BENCH

In this section, we introduce **MultiView-Bench** (Multi-view Benchmark), a foundational evaluation designed to test VLMs' multi-view spatial reasoning capabilities, preparing them for complex real-world tasks like mechanical engineering or 3D scene reconstruction.

4.1 SETUP

The core task of MultiView-Bench assesses a VLM's ability to reason about the relative positions of objects within a 3D scene. VLMs must observe scenes from multiple viewpoints to infer spatial relationships accurately (Figure 2). To ensure consistency, we introduce a fixed global coordinate system with clearly marked axes—X (red), Y (green), and Z (blue)—providing a viewpoint-independent frame of reference for spatial descriptions. This coordinate system is created as fixed 3D meshes in Blender so it keeps invariant upon camera change. We added appropriate opacity in objects' materials so that axes would be clearly seen. This is analogous to Blender GUI and this benchmark can be seen as a fundamental test to pass before building any VLM-based agents for 3D model design for manufacturing as it requires VLMs to visually understand 3D coordinate system that is crucial in any 3D modeling software.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

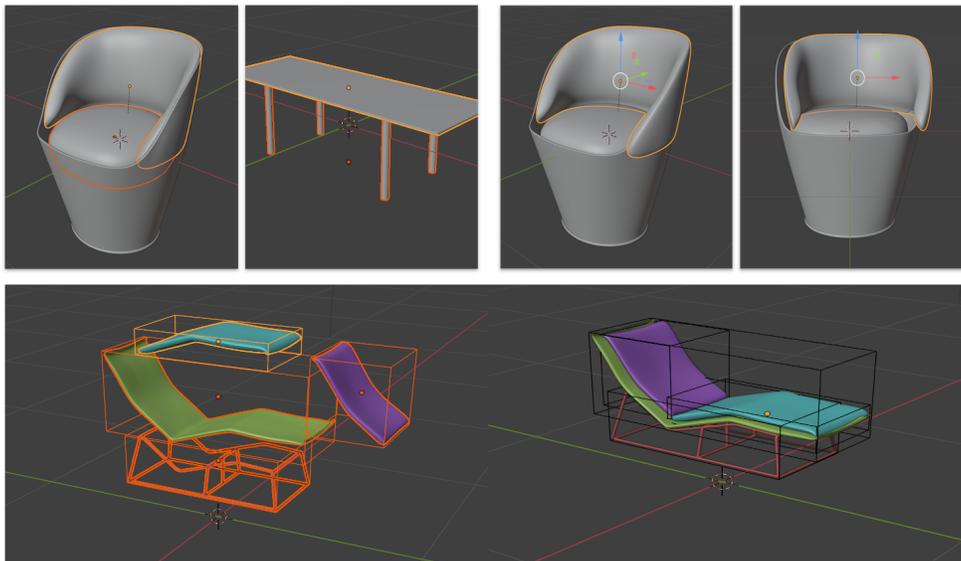


Figure 1: Furniture Assembly Example. Top Left: Many real-world objects do not lend themselves to simple natural language description. The table on the right can be described using fundamental convex shapes and their bounding boxes but the chair on the left has non-convex parts without an analytical expression. Thus it is preferable to include visual information. Top Right: Often times single-view observation leads to visual misconception and does not reveal some alignment issues. The chair looks well-assembled in the view shown on the left but when it turns to the view shown on the right we see the backrest is slightly misplaced in the X-axis. Bottom: Only using the bounding box dimensions, we are unable to assemble furniture that have non-convex parts. In the left instance, the bounding boxes are perfectly aligned but the assembly is problematic. A good assembly example is shown on the right.

After multi-view observations, VLMs must describe the object’s relative position to the central object along each axis using the format $(\pm X/0, \pm Y/0, \pm Z/0)$, ensuring precise and parsable responses to support large-scale evaluations.

4.2 DATASET CREATION

To construct our benchmark dataset, we design a modular pipeline that procedurally generates diverse 3D scenes with controlled variations (Figure 2).

3D assets. A wide variety of 3D models are readily available from online repositories. In MultiView-Bench, we used the fundamental geometric object randomly sampled from cube, sphere, cylinder and cone for the synthetic tasks. For real-world objects (3D Real World), we use 3DCoMPaT++ dataset (Slim et al. (2023)) which consists of thousands of real-world objects from different categories like table, chair, airplane and so on. To minimize visual ambiguity in spatial comparisons, all objects are rescaled to share a common bounding box (length, width, height) so that the relative position can be easily inferred from comparison for arbitrary edges or vertices. Beyond the main objects of interest, additional distractor objects can be introduced as confounders to increase scene complexity. The modularity of our pipeline also makes it easy to substitute alternative 3D assets and construct new datasets tailored to specific domains.

Object placement. We first fix a central object at the origin and then randomize other objects’ positions within the scene while enforcing a minimum and maximum separation distance. This prevents overlaps and ensures objects remain in close proximity. We set a threshold margin for 0 relation by pushing any smaller deviations along an axis to exactly 0. To better analyze model limitations, we also construct controlled task variants where target objects are restricted to simplified spatial layouts. In the DoF=1 variant, objects are placed along a single axis as the central object,

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
321
322
323

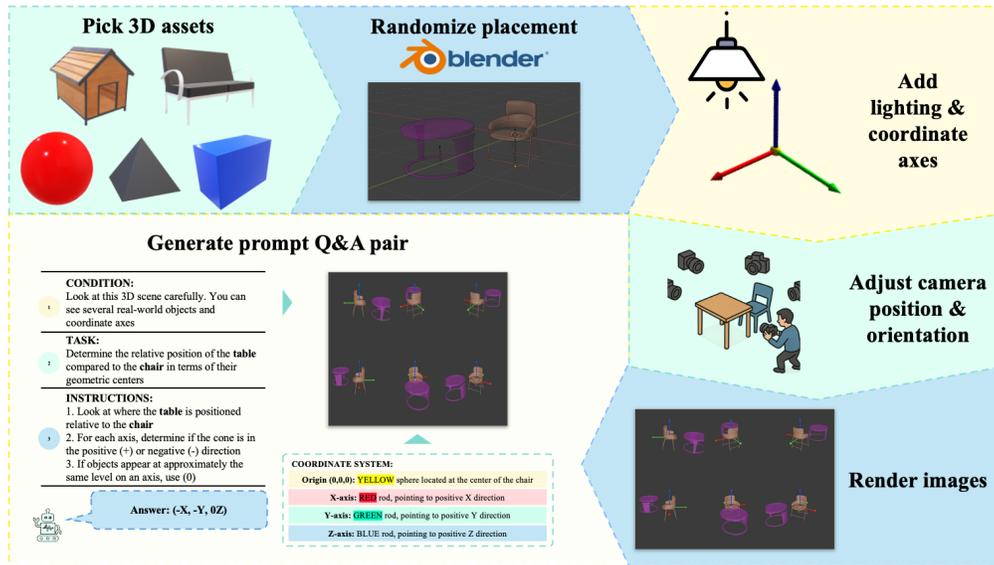


Figure 2: Data generation pipeline of MultiView-Bench: this pipeline allows us to choose 3D assets, manipulate placements, and adjust camera configurations to test different spatial cognitive abilities.

reducing the task to detecting 1D relative relationships. In the DoF=2 variant, objects lie on the same 2D plane as the central object and DoF=3 refers to no constraints of placement in the 3D space.

Camera viewpoints. In the main benchmark, we render six viewpoints with uniformly distributed azimuth angles and slight elevations. This configuration guarantees visibility of all three axes. Importantly, the pipeline is not restricted to this setup: arbitrary viewpoint configurations can be specified, enabling analysis of inductive biases in VLMs and the creation of specialized tasks such as egocentric-to-allothetic transformations or spatial navigation through 3D environments. Examples of such extensions to tasks proposed in Jia et al. (2025) and Yin et al. (2025) are provided in the Appendix A.5.

Rendering. We render images from each viewpoint. While rendering complex scenes with many textured objects can be computationally expensive, the process is parallelizable across CPU cores, allowing multiple viewpoints to be rendered simultaneously. This design makes large-scale dataset creation both efficient and scalable.

Q&A generation. Finally, we generate question-answer pairs automatically. Relative spatial relations are computed directly from Blender’s intrinsic coordinate system, ensuring reliable ground-truth supervision.

Overall, this dataset provides a rigorous test of VLMs’ ability to integrate multi-view information and reason spatially. Moreover, the extensibility of the pipeline makes it a versatile testbed for probing the limits of 3D reasoning, analyzing inductive biases, and generating tailored training datasets for downstream applications.

5 FAILURE AND BIAS ANALYSIS

In this section, we investigate the limitations and biases of off-the-shelf VLMs by conducting detailed case studies on different VLMs.

5.1 FAILURE PATTERNS

First we notice that providing only a single isometric view (clearly displaying all three axes) significantly reduced accuracy for both models (Figure 3, 3D Single View), underscoring the essential

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

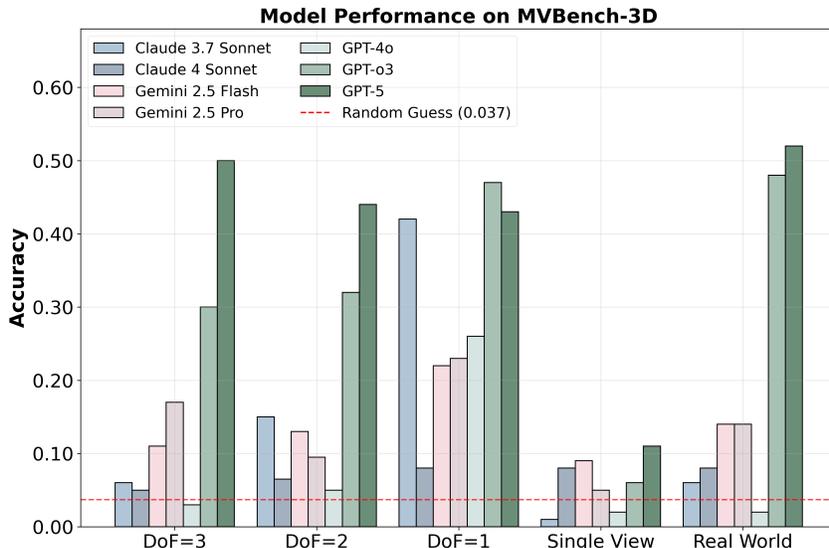


Figure 3: Model Performance on MultiView-Bench: we systematically evaluated the performance of 7 leading VLMs. We found that GPT-5 has the best overall performance while models like Claude series and GPT-4o can hardly beat random chance (red dashed line) in harder tasks 3D DoF=3 and 3D Real World. As we decrease the DoF, models tend to have higher accuracy with a great leap as we decrease the DoF to 1. Surprisingly, Claude 3.7 Sonnet beat its successor significantly in low DoF synthetic tasks and reached near GPT-5 level when DoF=1. We also show that single isometric view is not enough for solving this task as all models’ performances are around random chance. Full evaluation results are presented in Appendix A.1. The random chance is computed by $\frac{1}{3^3}$ to account for the 3 choices for each of the 3 axes.

role of multi-view observations. Single-view conditions resulted in substantial information loss, demonstrating that accurate spatial reasoning fundamentally requires multi-angle perception.

To better understand the sources of these failures, we decomposed the spatial reasoning task into intermediate sub-steps analogous to human reasoning processes:

- **Step 1: Object Identification:** Recognizing and distinguishing individual objects.
- **Step 2: 2D Spatial Relation:** Articulating the relative object positions within a given 2D viewpoint.
- **Step 3: Axis Direction Identification:** Interpreting and clearly describing the directionality of each axis within a 3D scene.
- **Step 4: 3D Position Translation:** Translating observed 2D spatial relationships into explicit 3D positional descriptions.

We performed some web trials and revealed that VLMs primarily struggle with accurately identifying and articulating the directionality of axes within the 3D context (Step 3), as illustrated in Figure 4. Models consistently exhibited difficulty in expressing axis directions using clear and unambiguous 3D spatial language.

Motivated by this observation, we further explored whether VLM performance could be enhanced by decomposing the task into simpler 2D views. We designed a variant of the MultiView-Bench task (2D) utilizing three canonical (front, side, top) views, as shown in Figure 8, each emphasizing one 2D plane (XZ, YZ, XY) and clearly displaying only two axes per view. Models were tested under two configurations: first, providing all three canonical views simultaneously to produce a single integrated 3D answer (Single-agent), and second, utilizing a multi-agent approach wherein each agent independently assessed one canonical view, with the final 3D answer obtained by straightforward integration of individual responses (resolving inconsistencies by random selection).

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397

Which object is at the center and which one is next to it?

At the center (origin 0,0,0): The brown/tan colored cube is positioned at the center of the coordinate system. This is where the three coordinate axes (red X-axis, green Y-axis, and blue Z-axis) intersect, and there's a small yellow sphere marking the origin point.

Next to it: The purple sphere is positioned above and slightly offset from the center cube. It appears to be floating in the positive Z direction (upward along the blue axis) relative to the central cube.

What is the relative direction of the purple sphere to the brown cube on this 2D image?

Looking at this 2D image, the purple sphere is positioned in the upper left relative to the brown cube.

Where does each axis point to in this 3D scene?

Looking at this 3D scene, I can identify the coordinate system based on the colored rods extending from the yellow sphere (origin):
Origin (0,0,0): Located at the yellow sphere in the center
X-axis (RED rod): Points to the left in the image – **negative X direction from our viewing perspective**
Y-axis (GREEN rod): Points forward/away from us into the scene – **positive Y direction**
Z-axis (BLUE rod): Points upward – **positive Z direction**

This appears to be a right-handed coordinate system where: The red X-axis extends **horizontally to the left**; The green Y-axis **extends into the depth of the scene**; The blue Z-axis extends vertically upward; The brown cube sits on the grid plane **below the origin**; The purple sphere floats above the scene.

Figure 4: Example conversation with VLM: VLM did well in 2D reasoning in the first 2 steps but failed to reason about the 3D scene in Step 3. Wrong reasoning traces are highlighted in yellow. Signs of inductive bias on the 3D coordinate system are highlighted in cyan.

401
402
403
404
405
406

Results from these experiments, depicted in Figure 5, demonstrate notable improvements by simplifying 3D reasoning tasks into 2D multi-view reasoning tasks and even further improvements by decomposing multi-view 2D tasks into single-view sub-tasks. These results validate that VLMs' struggle both at 3D perception and multi-view integration.

407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422

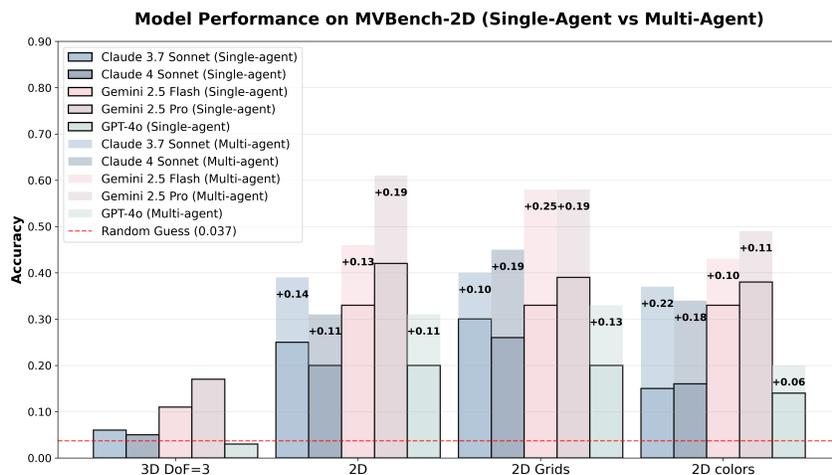


Figure 5: By simplifying the task into 2D multi-view task, VLMs demonstrate significantly higher performance. Decomposing 2D multi-view task into a multi-agent 2D single-view task can further improve VLMs' performances by a great margin. Adding visual aids like grids and distinct colors enhances the performance of Claude series but shows limited help to Gemini series and GPT-4o.

423
424
425
426
427

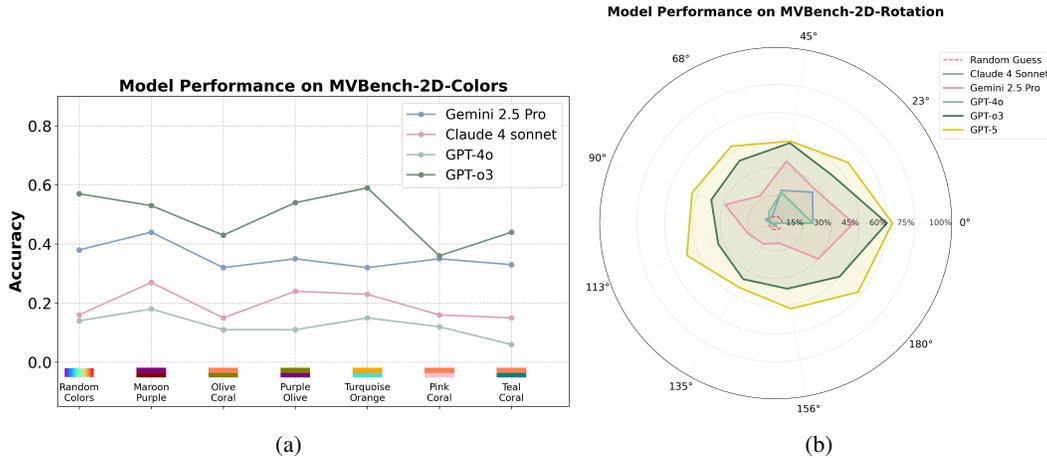
5.2 BIAS DISCOVERY

428
429
430
431

Visual Enhancements. We further explore how visual enhancements, such as distinct color schemes and grids, influence VLMs' 3D perception and spatial reasoning performance. Specifically, we evaluated VLMs under conditions of randomized color assignment, fixed color combinations, 3D

432 scenes with integrated 2D grid planes, and purely 2D views with grid overlays. Previous research
 433 indicates that structured visual aids can improve VLM performance on visual reasoning tasks like
 434 counting and scene comprehension (Izadi et al., 2025). Our experiments confirm that visual struc-
 435 tures indeed boost performance for models like Claude series; however, surprisingly, Gemini 2.5
 436 Pro’s performance declines under these conditions. Additionally, we observed distinct color biases
 437 among different models, with each showing preferential responses to particular color combinations,
 438 underscoring inherent perceptual biases (Figure 6).

439
 440 **Coordinate Rotation.** Upon closer examination of model responses, we identified a notable pat-
 441 tern: models frequently disregarded explicitly depicted coordinate directions, defaulting instead to
 442 reasoning based on the conventional right-handed coordinate system (Figure 4). We hypothesize that
 443 this behavior reflects a strong inductive bias derived from extensive exposure to standard coordinate
 444 conventions during pretraining. To rigorously investigate this bias, we conducted rotation experi-
 445 ments using the canonical 2D view task. In these experiments, axes directions were deliberately
 446 rotated away from conventional orthogonal orientations (such as 90° , 180° , etc.) to non-standard an-
 447 gles. Results shown in Figure 6 clearly showed significant performance degradation under these un-
 448 conventional orientations, confirming that VLMs heavily rely on learned coordinate-system priors.
 449 These findings highlight the critical need to address such inherent biases to enhance the robustness
 450 and generalization capabilities of VLMs.



464
 465
 466 Figure 6: (a) VLMs exhibit fluctuating performance across different color combinations. Notably,
 467 these fluctuations diverge from human perception: the models perform better on the less distinguish-
 468 able Maroon/Purple pair, yet worse on the more distinct Olive/Coral pair. (b) VLMs’ performances
 469 degrade under unconventional axis orientations.

471 6 VIEWNAVIGATOR

472 In this section, we introduce ViewNavigator, a multi-agent system designed to actively reason about
 473 spatial relationships between objects within a 3D environment.

474 Our agent architecture integrates a VLM and a LLM in a closed-loop manner without requiring
 475 post-training or external geometry-based image analysis. The LLM strategically plans the next
 476 move, deciding the viewpoint to look at. The VLM processes visual inputs from one viewpoint
 477 and its jittered viewpoints each time. The probabilistic belief module (details in Appendix A.3)
 478 integrates feedback from VLM to maintain a memory of the trajectory and belief state, which the
 479 LLM retrieves to guide future actions. The LLM emits the final answer if it is confident enough.

482 6.1 BELIEF MODULE

483 The agent maintains a probabilistic belief over spatial directions for each axis using a Dirichlet
 484 distribution. After sampling multiple jittered views, the VLM’s categorical votes are aggregated into
 485 smoothed proportions. To guard against unstable or noisy predictions, we discount the evidence size

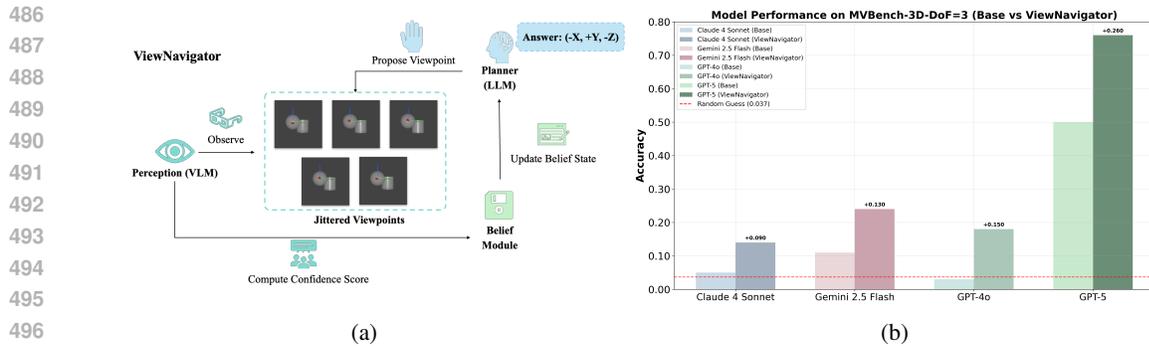


Figure 7: (a) ViewNavigator workflow. (b) ViewNavigator framework significantly enhances various base models’ performances on the 3D DoF=3 tasks.

with a confidence score, computed either from a Wilson lower bound on the majority vote or from the entropy of the vote distribution. The resulting effective evidence determines how strongly the current observations update the belief state. This iterative procedure ensures that stable, consistent clusters strengthen the belief, while ambiguous clusters contribute more cautiously. Full details can be referred to A.3.

6.2 ACTIVE VIEW SELECTION AND AGGREGATION

The LLM planner actively proposes the next best camera viewpoint based on the current belief state and previous view history, aiming to maximize information gain and reduce uncertainty. After each belief update, the agent checks if the posterior probability of the dominant class on each axis exceeds a confidence threshold τ (e.g., $\tau = 0.9$) and if sufficient evidence concentration is reached (e.g., total evidence $\sum_s \alpha_{A,s} \geq \kappa_{\min}$). If these criteria are met for all axes, the agent terminates the viewpoint exploration and outputs the final prediction:

$$\text{prediction} = \arg \max_{s \in \{+, 0, -\}} \frac{\alpha_{A,s}}{\sum_t \alpha_{A,t}}, \quad \forall A \in \{X, Y, Z\}. \quad (1)$$

ViewNavigator significantly enhances the performances across diverse base models by a large margin (Figure 7). Detailed configurations and prompts for ViewNavigator is presented in Appendix A.4

7 CONCLUSION

In this work, we presented **MultiView-Bench**, a diagnostic benchmark specifically designed to test the ability of VLMs to understand 3D global coordinate system and integrate multi-view observations into a coherent 3D scene understanding. Alongside the benchmark, we introduced an extensible **data-generation pipeline** that allows researchers to readily construct new datasets and a multi-agent framework **ViewNavigator** that significantly improves the performance of diverse base models on MultiView-Bench. Our systematic evaluation of leading VLMs on MultiView-Bench revealed fundamental limitations: while these models excel at recognizing 2D planar relations from single images, they struggle with integrating information across multiple views, interpreting 3D spatial relations, and generalizing under unconventional axes or texture variations. Taken together, MultiView-Bench, its extensible pipeline, and ViewNavigator form both a diagnostic tool and a stepping stone toward more powerful VLM-based agents. This benchmark is designed to raise awareness of the limitations of current VLMs in multi-view integration and horizontally benchmark and track improvements of VLMs. Our benchmark also serves as a selection standard for base model when building VLM-based 3D-reasoning agents for mechanical engineering or operating 3D modeling software GUI. We hope that future research builds on this foundation to equip VLMs with the spatial understanding necessary for diverse downstream 3D tasks such as part assembly, scene editing, and 3D assets creation.

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

ACKNOWLEDGMENTS

I would like to express my sincere gratitude to Yujing Zhang for her invaluable contribution in stylizing the graphics used in this paper.

REPRODUCIBILITY STATEMENT

We ensure reproducibility of our experiments by reporting the exact hyperparameters, experimental setup, and prompts. For all VLM and LLM API calls (except GPT-o3 and GPT-5 where the temperature is fixed), we use `temperature = 1.0` and `max.tokens = 4096`. For model evaluations, we used 100 tasks per task variant and repeated each experiment three times, reporting the averaged results. For ViewNavigator, due to the cost and time, we use 50 tasks per task variant and repeated three times as well. Full details of the experiments setup and the exact prompts for each task are provided in Appendix A.4.

LLM USAGE STATEMENT

LLM is used to polish a few sentences in the paper and format some mathematical expressions in LaTeX.

REFERENCES

- Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Chormanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, Pete Florence, Chuyuan Fu, Montse Gonzalez Arenas, Keerthana Gopalakrishnan, Kehang Han, Karol Hausman, Alexander Herzog, Jasmine Hsu, Brian Ichter, Alex Irpan, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Lisa Lee, Tsang-Wei Edward Lee, Sergey Levine, Yao Lu, Henryk Michalewski, Igor Mordatch, Karl Pertsch, Kanishka Rao, Krista Reymann, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Pierre Sermanet, Jaspier Singh, Anikait Singh, Radu Soricut, Huong Tran, Vincent Vanhoucke, Quan Vuong, Ayzaan Wahid, Stefan Welker, Paul Wohlhart, Jialin Wu, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-2: Vision-language-action models transfer web knowledge to robotic control, 07 2023. URL <https://arxiv.org/abs/2307.15818>.
- Tom B Brown, Benjamin Mann, Nick Ryder, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems (NeurIPS)*, 33:1877–1901, 2020.
- Sébastien Bubeck et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- Heinrich H. Bülthoff and Shimon Edelman. Psychophysical support for a two-dimensional view interpolation theory of object recognition. *Proceedings of the National Academy of Sciences of the United States of America*, 89(1):60–64, 1992. ISSN 00278424, 10916490. URL <http://www.jstor.org/stable/2358475>.
- Angel X Chang, Thomas Funkhouser, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015.
- Boyuan Chen, Zhuo Xu, Sean Kirmani, Brian Ichter, Danny Driess, Pete Florence, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. Spatialvlm: Endowing vision-language models with spatial reasoning capabilities, 2024. URL <https://arxiv.org/abs/2401.12168>.
- 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, 2024. URL <https://arxiv.org/abs/2406.01584>.
- Wenliang Dai et al. Instructblip: Towards general-purpose vision-language models with instruction tuning. In *NeurIPS*, 2023.
- Erik Daxberger, Nina Wenzel, David Griffiths, Haiming Gang, Justin Lazarow, Gefen Kohavi, Kai Kang, Marcin Eichner, Yinfei Yang, Afshin Dehghan, and Peter Grasch. Mm-spatial: Exploring

594 3d spatial understanding in multimodal llms. URL https://openaccess.thecvf.com/content/ICCV2025/papers/Daxberger_MM-Spatial_Exploring_3D_Spatial_Understanding_in_Multimodal_LLMs_ICCV_2025_paper.pdf.

595

596

597

598 Nianchen Deng, Lixin Gu, Shenglong Ye, Yinan He, and Wenhai Wang. Internspatial: A comprehensive dataset for spatial reasoning in vision-language models, 06

599 2025. URL https://www.researchgate.net/publication/392942436_InternSpatial_A_Comprehensive_Dataset_for_Spatial_Reasoning_in_Vision-Language_Models.

600

601

602

603 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, et al. An image is worth 16x16 words:

604 Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2021.

605

606

607 Hao Du, Yiyuan Zhao, et al. Navgpt: Explicit reasoning in vision-and-language navigation with

608 large language models. *arXiv preprint arXiv:2302.09170*, 2023.

609

610 Shimon Edelman. Representation is representation of similarities. *Behavioral and Brain Sciences*,

611 21(4):449–467, 1998.

612

613 Google. Gemini: A family of highly capable multimodal models, 12 2023. URL <https://arxiv.org/abs/2312.11805>.

614

615 Yunqi Gu, Ian Huang, Jihyeon Je, Guandao Yang, and Leonidas Guibas. Blendergym: Bench-

616 marking foundational model systems for graphics editing, 2025. URL <https://arxiv.org/abs/2504.01786>.

617

618 Yining Hong, Chunru Lin, Yilun Du, Zhenfang Chen, Joshua B Tenenbaum, and Chuang Gan. 3d

619 concept learning and reasoning from multi-view images, 2023. URL <https://arxiv.org/abs/2303.11327>.

620

621 Ziniu Hu, Ahmet Iscen, Aashi Jain, Thomas Kipf, Yisong Yue, David A Ross, Cordelia Schmid, and

622 Alireza Fathi. Scenecraft: An llm agent for synthesizing 3d scene as blender code, 2024. URL

623 <https://arxiv.org/abs/2403.01248>.

624

625 Ian Huang, Guandao Yang, and Leonidas Guibas. Blenderalchemy: Editing 3d graphics with vision-

626 language models, 2024. URL <https://arxiv.org/abs/2404.17672>.

627

628 Amirmohammad Izadi, Mohammad Ali Banayeeanzade, Fatemeh Askari, Ali Rahimiakbar, Mo-

629 hammad Mahdi Vahedi, Hosein Hasani, and Baghshah Mahdieh Soleymani. Visual struc-

630 tures helps visual reasoning: Addressing the binding problem in vlms, 2025. URL <https://arxiv.org/abs/2506.22146>.

631

632 Mengdi Jia, Zekun Qi, Shaochen Zhang, Wenyao Zhang, and Li Yi. Omnispatial: Towards

633 comprehensive spatial reasoning benchmark for vision language models, 06 2025. URL

634 https://www.researchgate.net/publication/392372432_Omnispatial_Towards_Comprehensive_Spatial_Reasoning_Benchmark_for_Vision_Language_Models.

635

636

637 Justin Johnson et al. Clevr: A diagnostic dataset for compositional language and elementary visual

638 reasoning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2901–2910, 2017.

639

640 Stephen M Kosslyn. *Image and Brain: The Resolution of the Imagery Debate*. MIT press, 1994.

641

642 Dingming Li, Hongxing Li, Zixuan Wang, Yuchen Yan, Hang Zhang, Siqi Chen, Guiyang Hou,

643 Shengpei Jiang, Wenqi Zhang, Yongliang Shen, Weiming Lu, and Yueting Zhuang. Viewspatial-

644 bench: Evaluating multi-perspective spatial localization in vision-language models, 2025. URL

645 <https://arxiv.org/abs/2505.21500>.

646

647 Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-

training for unified vision-language understanding and generation. In *International Conference on Machine Learning (ICML)*, pp. 12888–12900, 2022.

648 Fangchen Liu et al. Llm+p: Empowering large language models with optimal planning proficiency.
649 In *ICLR*, 2023.
650

651 Parker Liu, Chenxin Li, Zhengxin Li, Yipeng Wu, Wuyang Li, Zhiqin Yang, Zhenyuan Zhang,
652 Yunlong Lin, Sirui Han, and Brandon Y Feng. Ir3d-bench: Evaluating vision-language model
653 scene understanding as agentic inverse rendering, 2025. URL [https://arxiv.org/abs/
654 2506.23329](https://arxiv.org/abs/2506.23329).

655 David Marr. *Vision: A Computational Investigation into the Human Representation and*
656 *Processing of Visual Information*. The MIT Press, 07 2010. ISBN 9780262514620.
657 doi: 10.7551/mitpress/9780262514620.001.0001. URL [https://doi.org/10.7551/
659 mitpress/9780262514620.001.0001](https://doi.org/10.7551/
658 mitpress/9780262514620.001.0001).

660 Lars Mescheder, Michael Oechsle, et al. Occupancy networks: Learning 3d reconstruction in func-
661 tion space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recog-
662 nition (CVPR)*, pp. 4460–4470, 2019.

663 Piotr Mirowski, Razvan Pascanu, et al. Learning to navigate in cities without a map. *Advances in*
664 *Neural Information Processing Systems (NeurIPS)*, pp. 2419–2430, 2018.
665

666 OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
667

668 OpenAI. Gpt-4o system card, 2024. URL <https://arxiv.org/abs/2410.21276>.

669 Stephen E. Palmer. *Vision science: Photons to phenomenology*. MIT Press, 1999.
670

671 Ava Pun, Kangle Deng, Ruixuan Liu, Deva Ramanan, Changliu Liu, and Jun-Yan Zhu. Generating
672 physically stable and buildable brick structures from text, 2025. URL [https://arxiv.org/
674 abs/2505.05469](https://arxiv.org/
673 abs/2505.05469).

675 Charles R Qi, Hao Su, et al. Pointnet: Deep learning on point sets for 3d classification and segmen-
676 tation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*
677 *(CVPR)*, pp. 652–660, 2017a.

678 Charles R Qi et al. Pointnet++: Deep hierarchical feature learning on point sets in a metric space.
679 In *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 5099–5108, 2017b.
680

681 Alec Radford, Jong Wook Kim, Chris Hallacy, et al. Learning transferable visual models from
682 natural language supervision. In *International Conference on Machine Learning (ICML)*, pp.
683 8748–8763, 2021.

684 Roger N Shepard and Jacqueline Metzler. Mental rotation of three-dimensional objects. *Science*,
685 171(3972):701–703, 1971.
686

687 Habib Slim, Xiang Li, Yuchen Li, Mahmoud Ahmed, Mohamed Ayman, Ujjwal Upadhyay,
688 Ahmed Abdelreheem, Arpit Prajapati, Suhail Pothigara, Peter Wonka, and Mohamed Elhoseiny.
689 3dcompat⁺⁺: An improved large-scale 3d vision dataset for compositional recognition, 2023.
690 URL <https://arxiv.org/abs/2310.18511>.

691 Mark G. Stokes et al. The prefrontal cortex and cognitive control. *Annual Review of Neuroscience*,
692 44:403–423, 2021.
693

694 Michael J Tarr and Heinrich H Bülthoff. Image-based object recognition in man, monkey and
695 machine. *Cognition*, 67(1-2):1–20, 1998.

696 Hugo Touvron et al. Llama: Open and efficient foundation language models. *arXiv preprint*
697 *arXiv:2302.13971*, 2023.
698

699 Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan,
700 and Anima Anandkumar. Voyager: An open-ended embodied agent with large language mod-
701 els, 05 2023. URL [https://arxiv.org/abs/2305.16291#:~:text=Voyager%
20interacts%20with%20GPT-4%20via%20blackbox%20queries%2C%20which](https://arxiv.org/abs/2305.16291#:~:text=Voyager%20interacts%20with%20GPT-4%20via%20blackbox%20queries%2C%20which).

702 Wenqi Wang, Reuben Tan, Pengyue Zhu, Jianwei Yang, Zhengyuan Yang, Lijuan Wang, An-
703 drey Kolobov, Jianfeng Gao, and Boqing Gong. Site: towards spatial intelligence thorough
704 evaluation site-project-page. URL [https://openaccess.thecvf.com/content/
705 ICCV2025/papers/Wang_SITE_towards_Spatial_Intelligence_Thorough_
706 Evaluation_ICCV_2025_paper.pdf](https://openaccess.thecvf.com/content/ICCV2025/papers/Wang_SITE_towards_Spatial_Intelligence_Thorough_Evaluation_ICCV_2025_paper.pdf).

707 Yue Wang et al. Dynamic graph cnn for learning on point clouds. In *ACM Transactions on Graphics*
708 (*TOG*), 2019.

709

710 Zhirong Wu, Shuran Song, et al. 3d shapenets: A deep representation for volumetric shapes. In
711 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp.
712 1912–1920, 2015.

713 Yutaro Yamada, Yihan Bao, Andrew K Lampinen, Jungo Kasai, and Ilker Yildirim. Evaluating
714 spatial understanding of large language models, 2023. URL [https://arxiv.org/abs/
715 2310.14540](https://arxiv.org/abs/2310.14540).

716

717 Yutaro Yamada, Khyathi Chandu, Yuchen Lin, Jack Hessel, Ilker Yildirim, and Yejin Choi. L3go:
718 Language agents with chain-of-3d-thoughts for generating unconventional objects, 2024. URL
719 <https://arxiv.org/abs/2402.09052>.

720 Shunyu Yao et al. React: Synergizing reasoning and acting in language models. In *ICLR*, 2023.

721

722 Baiqiao Yin, Qineng Wang, Pingyue Zhang, Jianshu Zhang, Kangrui Wang, Zihan Wang, Jieyu
723 Zhang, Keshigeyan Chandrasegaran, Han Liu, Ranjay Krishna, Saining Xie, Manling Li, Jia-
724 jun Wu, and Li Fei-Fei. Spatial mental modeling from limited views, 2025. URL [https:
725 //arxiv.org/abs/2506.21458](https://arxiv.org/abs/2506.21458).

726 Jiahui Zhang, Yurui Chen, Yanpeng Zhou, Yueming Xu, Ze Huang, Jilin Mei, Junhui Chen, Yu-
727 Jie Yuan, Xinyue Cai, Guowei Huang, Xingyue Quan, Hang Xu, and Li Zhang. From flatland
728 to space: Teaching vision-language models to perceive and reason in 3d, 2025. URL [https:
729 //arxiv.org/abs/2503.22976](https://arxiv.org/abs/2503.22976).

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756 A APPENDIX

757
758 A.1 MODEL EVALUATIONS IN ALL MULTIVIEW-BENCH TASKS

759
760 Table 2: Model Performance on MultiView-Bench

761

762 Tasks/Models	Claude 3.7 Sonnet	Claude 4 Sonnet	Gemini 2.5 Flash	Gemini 2.5 Pro	GPT-4o	GPT-5	GPT-o3
763 3D DoF=3	0.060	0.050	0.110	0.170	0.030	0.500	0.300
764 3D DoF=2	0.150	0.065	0.130	0.095	0.050	0.440	0.320
765 3D DoF=1	0.420	0.080	0.220	0.230	0.260	0.430	0.470
766 3D Single View	0.010	0.080	0.090	0.050	0.020	0.110	0.060
767 3D Real World	0.060	0.080	0.140	0.140	0.020	0.520	0.480
768 2D Three Views (Rotation 0)	0.250	0.200	0.330	0.420	0.200	0.630	0.550
769 2D Three Views multiagent (Rotation 0)	0.390	0.310	0.460	0.610	0.310	0.660	0.560
770 2D Three Views Grids	0.300	0.260	0.330	0.390	0.200	0.690	0.540
771 2D Three Views Grids multiagent	0.400	0.450	0.580	0.580	0.330	0.580	0.570
772 2D Three Views Colors (Random)	0.150	0.160	0.330	0.380	0.140	0.640	0.560
773 2D Three Views Colors multiagent (Random)	0.370	0.340	0.430	0.490	0.200	0.700	0.420
774 2D Three Views Colors (Maroon+Purple)	0.190	0.270	0.330	0.440	0.180	0.680	0.590
775 2D Three Views Colors (Turquoise+Orange)	0.210	0.230	0.300	0.320	0.150	0.710	0.530
776 2D Three Views Colors (Teal+Olive)	0.230	0.240	0.320	0.350	0.110	0.660	0.540
777 2D Three Views Colors (Teal+Coral)	0.150	0.150	0.200	0.330	0.060	0.530	0.430
778 2D Three Views Colors (Olive+Coral)	0.200	0.150	0.200	0.320	0.110	0.600	0.440
779 2D Three Views Colors (Pink+Coral)	0.180	0.160	0.130	0.350	0.120	0.580	0.360
780 2D Three Views Rotation (23°)	0.210	0.260	0.310	0.280	0.140	0.510	0.400
781 2D Three Views Rotation (45°)	0.240	0.180	0.310	0.340	0.170	0.450	0.440
782 2D Three Views Rotation (68°)	0.090	0.040	0.290	0.170	0.070	0.480	0.390
783 2D Three Views Rotation (90°)	0.060	0.060	0.130	0.290	0.040	0.480	0.370
784 2D Three Views Rotation (113°)	0.030	0.020	0.200	0.160	0.010	0.510	0.330
785 2D Three Views Rotation (135°)	0.070	0.000	0.160	0.130	0.020	0.400	0.350
786 2D Three Views Rotation (156°)	0.030	0.020	0.170	0.110	0.010	0.470	0.360
787 2D Three Views Rotation (180°)	0.000	0.000	0.120	0.300	0.000	0.580	0.450

788
789 A.2 MORE EXAMPLE TASKS

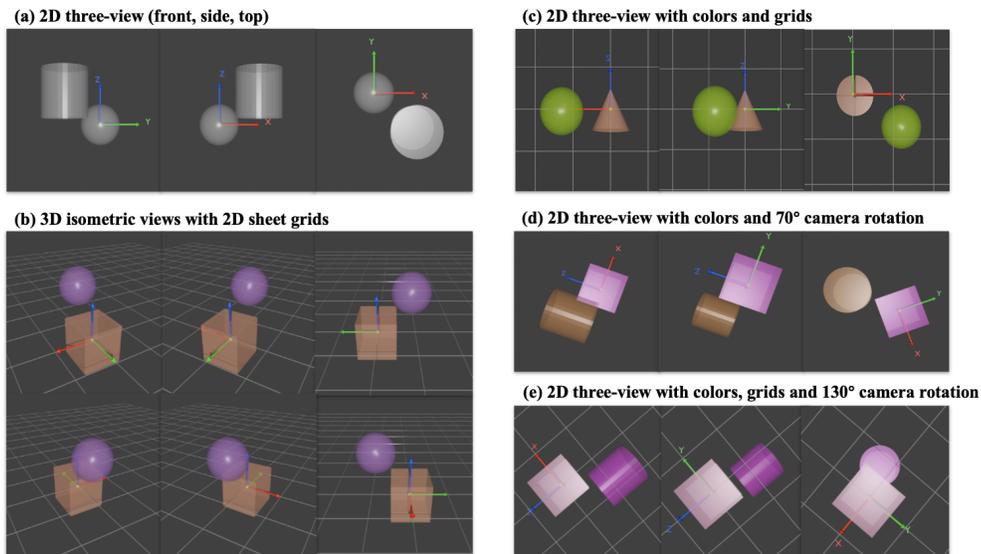


Figure 8: More Example Tasks

810 A.3 BELIEF STATE AND UPDATE

811
812 The agent maintains a probabilistic belief state over spatial directions (+, 0, -) for each axis inde-
813 pendently. Specifically, we model this belief as a Dirichlet distribution parameterized by vector
814 $\alpha_A = [\alpha_{A,+}, \alpha_{A,0}, \alpha_{A,-}]$ for each axis $A \in \{X, Y, Z\}$. Initially, each axis is given an uniform
815 prior: $\alpha_A = [1, 1, 1]$.

816 Upon selecting a camera viewpoint, the agent captures multiple images using micro-jitters (small
817 perturbations around a base viewpoint) to assess stability in the VLM’s answers, which is used as
818 a confidence score. For each jittered viewpoint, the VLM returns a categorical judgment (+, 0, -)
819 independently for each axis, resulting in vote counts $k_{A,+}, k_{A,0}, k_{A,-}$ from a set of n images.

820 To update the belief, we first compute smoothed proportions:

$$821 \hat{p}_{A,s} = \frac{k_{A,s} + \lambda}{n + 3\lambda}, \quad s \in \{+, 0, -\} \quad (2)$$

822 where λ is a smoothing constant (default $\lambda = 1$) to avoid over-confidence for small n .

823 These proportions represent the directional preference of the cluster, while the confidence score
824 discounts clusters that show large variability under micro-jitters.

825 We propose two methods to compute the confidence score: **Wilson Lower Bound Score** and **Rela-**
826 **tive Entropy Score**, both of which achieved comparable performance.

827
828 **Wilson Lower Bound Score.** Given that the majority label among the n answers occurs k_{\max}
829 times, the empirical majority proportion is $\hat{p} = k_{\max}/n$. The Wilson score interval offers a conser-
830 vative estimate of the true binomial proportion, particularly robust for small n or when proportions
831 are near 0 or 1. The 95% Wilson lower bound is computed as:

$$832 \text{LB} = \frac{\hat{p} + z^2/(2n) - z\sqrt{\frac{\hat{p}(1-\hat{p})}{n} + \frac{z^2}{4n^2}}}{1 + z^2/n}, \quad z = 1.96. \quad (3)$$

833 We map this to a conservative confidence score relative to a random baseline (uniform guess = 1/3):

$$834 \omega_A = \left(\frac{\max(\text{LB}, 1/3) - 1/3}{2/3} \right)^\gamma, \quad (4)$$

835 where $\gamma \in [1, 2]$ controls sensitivity.

836
837 **Relative Entropy Score.** Let $H(\hat{\mathbf{p}}_A) = -\sum_s \hat{p}_{A,s} \log \hat{p}_{A,s}$ be the entropy of the smoothed vote
838 distribution and $H_{\max} = \log 3$ its maximum for three equally likely outcomes. The normalized
839 entropy gap from uniform is:

$$840 \omega_A = \left(1 - \frac{H(\hat{\mathbf{p}}_A)}{H_{\max}} \right)^\gamma, \quad (5)$$

841 with γ again controlling sensitivity.

842 These two methods prevent overconfidence when a cluster’s votes are unstable, rewarding highly
843 peaked vote distributions and penalizing near-uniform ones.

844 The effective evidence size is then:

$$845 n_{\text{eff},A} = n \cdot \omega_A. \quad (6)$$

846
847 **Belief Update.** The smoothed proportions $\hat{p}_{A,s}$ are scaled by $n_{\text{eff},A}$ to yield soft counts:

$$848 \Delta\alpha_{A,s} = n_{\text{eff},A} \cdot \hat{p}_{A,s}, \quad \forall s \in \{+, 0, -\} \quad (7)$$

849 These are added to the Dirichlet parameters to yield the new belief:

$$850 \alpha_{A,s} \leftarrow \alpha_{A,s} + \Delta\alpha_{A,s} \quad (8)$$

851 This belief is updated iteratively over successive jittered view clusters.

864 **Active View Selection and Aggregation** The LLM planner actively proposes the next best camera
865 viewpoint based on the current belief state and previous view history, aiming to maximize in-
866 formation gain and reduce uncertainty. After each belief update, the agent checks if the posterior
867 probability of the dominant class on each axis exceeds a confidence threshold τ and if sufficient
868 evidence concentration is reached (e.g., total evidence $\sum_s \alpha_{A,s} \geq \kappa_{\min}$). If these criteria are met for
869 all axes, the agent terminates the viewpoint exploration and outputs the final prediction:

$$870 \text{ prediction} = \arg \max_{s \in \{+,0,-\}} \frac{\alpha_{A,s}}{\sum_t \alpha_{A,t}}, \quad \forall A \in \{X, Y, Z\}. \quad (9)$$

873 A.4 EXPERIMENTS SETUP

874
875 All experiments are conducted using our proposed **MultiView-Bench** benchmark. The 3D scenes
876 and corresponding multi-view images are procedurally generated using Blender. Each scene is con-
877 structed with a fixed global coordinate system, represented by colored axes (X: red, Y: green, Z:
878 blue), to provide a consistent frame of reference across all viewpoints. The generation pipeline
879 allows for the use of various 3D assets, randomized object placements, and configurable camera
880 positions, as detailed in Section 4 of the main paper.

881 A.4.1 HYPERPARAMETERS

882 **VLM API Calls (Single-Agent & Multi-Agent)** For models other than GPT-o3 and GPT-5:

- 883 • `temperature` = 1.0
- 884 • `max_tokens` = 4096

885
886 **ViewNavigator Framework** The ViewNavigator agentic framework was configured as follows:

- 887 • `max_steps` = 10 (Maximum number of viewpoints the agent can select)
 - 888 • `r_az` = 5 (Radius in degrees for azimuthal jitter)
 - 889 • `r_el` = 5 (Radius in degrees for elevation jitter)
 - 890 • `tau` = 0.6 (Confidence threshold τ for the belief state using the Wilson Lower Bound
891 Score)
 - 892 • `jitter_size` = 5 (Number of jittered images per viewpoint)
- 893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

A.4.2 PROMPTS

The exact prompts used in our experiments are provided below. Placeholders such as {central_obj_type} are filled dynamically during data generation. Labels such as % VLM SYSTEM PROMPT are included only for readability in the paper.

Prompt for MultiView-Bench-3D and MultiView-Bench-2D Tasks (Single-agent):

Look at this 3D scene carefully from different viewpoints. You can see several geometric objects and coordinate axes.

COORDINATE SYSTEM:

- X-axis: RED rod, pointing to positive X direction
- Y-axis: GREEN rod, pointing to positive Y direction
- Z-axis: BLUE rod, pointing to positive Z direction
- Origin (0,0,0): YELLOW sphere, located at the center of the {central_obj_type}

TASK:

Determine the relative position of the {sampled_obj_type} compared to the {central_obj_type} in terms of their geometric centers.

INSTRUCTIONS:

1. Look at where the {sampled_obj_type} is positioned relative to the {central_obj_type}
2. For each axis, determine if the {sampled_obj_type} is in the positive (+) or negative (-) direction using the coordinate system shown in the images.
3. If objects appear at approximately the same level on an axis, use (0)

ANSWER FORMAT:

Respond with exactly this format: <answer>(+X, +Y, +Z)</answer>
Examples: <answer>(+X, -Y, +Z)</answer> or <answer>(-X, 0Y, -Z)</answer>
or <answer>(0X, +Y, 0Z)</answer>

What is the relative position of the {sampled_obj_type} to the {central_obj_type}?

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

Prompt for MultiView-Bench-2D Tasks (Multi-agent):

```
% Prompt for Front View (XZ plane)
Look at this Front View (XZ plane) carefully. You can see several
geometric objects and coordinate axes.

VIEW DESCRIPTION:
This is the Front View (XZ plane), looking along the Y-axis.

COORDINATE SYSTEM:
- X-axis: RED rod, pointing to positive X direction
- Z-axis: BLUE rod, pointing to positive Z direction
- Origin (0,0,0): YELLOW sphere, located at the center of the
{central_obj_type}

TASK:
Determine the relative position of the {sampled_obj_type} compared
to the {central_obj_type} in terms of their geometric centers,
focusing only on the X and Z axes visible in this view.

INSTRUCTIONS:
1. Look at where the {sampled_obj_type} is positioned relative to
the {central_obj_type}
2. For each visible axis (X, Z), determine if the
{sampled_obj_type} is in the positive (+) or negative (-) direction
using the coordinate system shown in the image.
3. If objects appear at approximately the same level on an axis,
use (0)

ANSWER FORMAT:
Respond with exactly this format for the X and Z axes: <answer>(±X,
±Z)</answer>
Examples: <answer>(+X, -Z)</answer> or <answer>(0X, +Z)</answer>

What is the relative position of the {sampled_obj_type} to the
{central_obj_type} in the X and Z axes?
```

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

```
% Prompt for Side View (YZ plane)
Look at this Side View (YZ plane) carefully. You can see several
geometric objects and coordinate axes.

VIEW DESCRIPTION:
This is the Side View (YZ plane), looking along the X-axis.

COORDINATE SYSTEM:
- Y-axis: GREEN rod, pointing to positive Y direction
- Z-axis: BLUE rod, pointing to positive Z direction
- Origin (0,0,0): YELLOW sphere, located at the center of the
{central_obj_type}

TASK:
Determine the relative position of the {sampled_obj_type} compared
to the {central_obj_type} in terms of their geometric centers,
focusing only on the Y and Z axes visible in this view.

INSTRUCTIONS:
1. Look at where the {sampled_obj_type} is positioned relative to
the {central_obj_type}
2. For each visible axis (Y, Z), determine if the
{sampled_obj_type} is in the positive (+) or negative (-) direction
using the coordinate system shown in the image.
3. If objects appear at approximately the same level on an axis,
use (0)

ANSWER FORMAT:
Respond with exactly this format for the Y and Z axes: <answer>(<math>\pm Y,</math>
<math>\pm Z</math>)/<answer>
Examples: <answer>(<math>+Y,</math> <math>-Z</math>)/<answer> or <answer>(<math>0Y,</math> <math>+Z</math>)/<answer>

What is the relative position of the {sampled_obj_type} to the
{central_obj_type} in the Y and Z axes?"
```

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

```
% Prompt for Top View (XY plane)
Look at this Top View (XY plane) carefully. You can see several
geometric objects and coordinate axes.

VIEW DESCRIPTION:
This is the Top View (XY plane), looking along the Z-axis from
above.

COORDINATE SYSTEM:
- X-axis: RED rod, pointing to positive X direction
- Y-axis: GREEN rod, pointing to positive Y direction
- Origin (0,0,0): YELLOW sphere, located at the center of the
{central_obj_type}

TASK:
Determine the relative position of the {sampled_obj_type} compared
to the {central_obj_type} in terms of their geometric centers,
focusing only on the X and Y axes visible in this view.

INSTRUCTIONS:
1. Look at where the {sampled_obj_type} is positioned relative to
the {central_obj_type}
2. For each visible axis (X, Y), determine if the
{sampled_obj_type} is in the positive (+) or negative (-) direction
using the coordinate system shown in the image.
3. If objects appear at approximately the same level on an axis,
use (0)

ANSWER FORMAT:
Respond with exactly this format for the X and Y axes: <answer>(+X,
±Y)</answer>
Examples: <answer>(+X, -Y)</answer> or <answer>(0X, +Y)</answer>

What is the relative position of the {sampled_obj_type} to the
{central_obj_type} in the X and Y axes?
```

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

Prompts for ViewNavigator:

VLM PERCEPTION MODULE PROMPTS

```
% VLM SYSTEM PROMPT
You are a precise vision judge. The image shows colored world axes:

COORDINATE SYSTEM:
- X-axis: RED rod, pointing to positive X direction
- Y-axis: GREEN rod, pointing to positive Y direction
- Z-axis: BLUE rod, pointing to position Z direction
- Origin (0,0,0): YELLOW sphere, located at the center of the
{central_object}
- Both CENTRAL and TARGET objects have the same scale in X, Y, Z
dimensions

TASK:
Determine the relative position of the {target_object} compared to
the
{central_object} in terms of their geometric centers.

INSTRUCTIONS:
1. Only focus on axis {axis} for this view and only give answer for
these axes.
2. Compare the TARGET center to the CENTRAL center along each of
{axis}:
  . "+" if TARGET lies in the positive direction
  . "-" if in the negative direction
  . "0" if approximately equal (centers aligned along that axis)
3. Wrap your full step-by-step reasoning in <think>...</think>.
4. Then emit exactly one line, wrapped in <answer>...</answer>,
listing only
   axes {axis} with their sign or 0.
Do not include any extra text or prose.
```

```
% VLM USER PROMPT
CENTRAL: {central}
TARGET: {target}

Return exactly:
<think>...step-by-step reasoning...</think>
<answer>(+X, +Y, +Z)</answer>

Valid examples:
  <think>I see red and green axes...</think>
  <answer>(+X, -Y)</answer>

  <think>Blue axis only is clear...</think>
  <answer>(-Z)</answer>
```

1188 LLM PLANNER MODULE PROMPTS

```
1189
1190 % LLM SYSTEM PROMPT
1191 You control a camera in a 3D scene. Your goal is to decide the
1192 signs (+,0,-)
1193 of TARGET relative to CENTRAL on axes X,Y,Z by choosing successive
1194 viewpoints.
1195
1196 **Camera Coordinate System:**
1197 - Azimuth 0°: X-axis points towards viewer, Y-axis points right
1198 - As azimuth increases (clockwise rotation):
1199   - Azimuth 90°: Y-axis points towards viewer, X-axis points left
1200   - Azimuth 180°: X-axis points away from viewer, Y-axis points
1201     left
1202   - Azimuth 270°: Y-axis points away from viewer, X-axis points
1203     right
1204 - Elevation 0°: Camera views from directly above (top-down)
1205 - Elevation 90°: Camera views from horizontal level
1206 - Elevation 180°: Camera views from directly below (bottom-up)
1207
1208 On every turn you will receive:
1209 - threshold \tau (a float in [0,1])
1210 - belief_state:
1211   {
1212     "X": {"+": p_plus, "0": p_zero, "-": p_minus},
1213     "Y": {...},
1214     "Z": {...}
1215   }
1216 - history: a list of previously checked views, each entry:
1217   {
1218     "view": {"az": az_deg, "el": el_deg},
1219     "answer": "(±X, ±Y, ±Z)" or shorter,
1220     "confidence": {"X":cX, "Y":cY, "Z":cZ}
1221   }
1222
1223 If **all** axes have max(belief)  $\geq$   $\tau$ , you should stop.
1224 Otherwise choose the
1225 next best view. Note that you can revisit some views to strengthen
1226 your belief.
1227
1228 You should also decide which axes you want to focus on in a view.
1229 For example,
1230 if you choose a view that shows the XY plane, then you should focus
1231 on only
1232 the X axis and Y axis or even just focus on X or Y axis.
1233
1234 Note that the confidence score represents the reliability of the
1235 answer got from
1236 that view. Zero confidence score for a view indicates that the
1237 relative
1238 position is not clear revealed through that view.
1239
1240 **Rules**
1241 1. Wrap your internal reasoning in <think>...</think>
1242 2. Then emit exactly one <answer>...</answer> containing **only**
1243 this JSON:
1244
1245 {
1246   "action": "CAPTURE"|"STOP",
1247   "view": {"az": <number>, "el": <number>} | null,
1248   "axis": ["X", "Y"]
1249 }
1250 No extra text or fields.
```

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

```
% LLM FIRST TURN PROMPT
# First turn (no belief_state or history)
Task: find ( $\pm X, \pm Y, \pm Z$ ) for TARGET={target} vs CENTRAL={central}.
Threshold \tau = {\tau}.

Propose your initial viewpoint.
Respond with:
<think>...</think>
<answer>{{
  "action": "CAPTURE",
  "view": {{ "az": <num>, "el": <num> }},
  "axis": ["axes to focus on for this view"]
}}</answer>

% LLM INTERMEDIATE TURN PROMPT
# Subsequent turn
Threshold \tau = {\tau}
belief_state = {belief_state}
history      = {history}

Decide whether to STOP or pick another view.
Respond with:
<think>...</think>
<answer>{{
  "action": "CAPTURE"|"STOP",
  "view":   {{ "az": <num>, "el": <num> }} | null,
  "axis": ["axes to focus on for this view"] | null
}}</answer>
```

A.5 EXAMPLES OF MORE SPATIAL TASKS

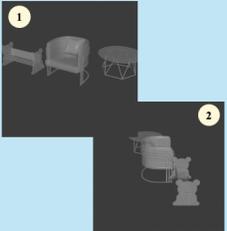
<p>Hypothetical Perspective-Taking</p> <p>Question: If you are sitting on the chair facing the front, the table is on the left or right of you?</p>  <p>Answer: Left side.</p>	<p>Spatial Navigation</p> <p>Question: If you are positioned at this viewpoint, then turn left and move forward, will you get closer to the chair?</p>  <p>Answer: No.</p>
<p>Spatial Mental Modeling</p> <p>Question: If you are at view 1 and move to view 2, which one is the furthest from you?</p>  <p>Answer: Table.</p>	<p>Allocentric Perspective-Taking</p> <p>Question: The table is on the left or right side of the bench?</p>  <p>Answer: Left side.</p>

Figure 9: Examples of more spatial tasks that can be created using our data generation pipeline.

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

A.6 FAILURE CASES

In this section, we present a glimpse of failure cases we observed during the evaluations.

A.6.1 3D DoF=3 (CLAUDE 3.7 SONNET)

In this challenging 3D spatial reasoning task, the model correctly identifies the relationships on the X and Y axes but fails to accurately determine the relative position on the Z-axis.

- **Ground Truth:** $(-X, +Y, +Z)$
- **Model's Answer:** $(-X, +Y, 0Z)$

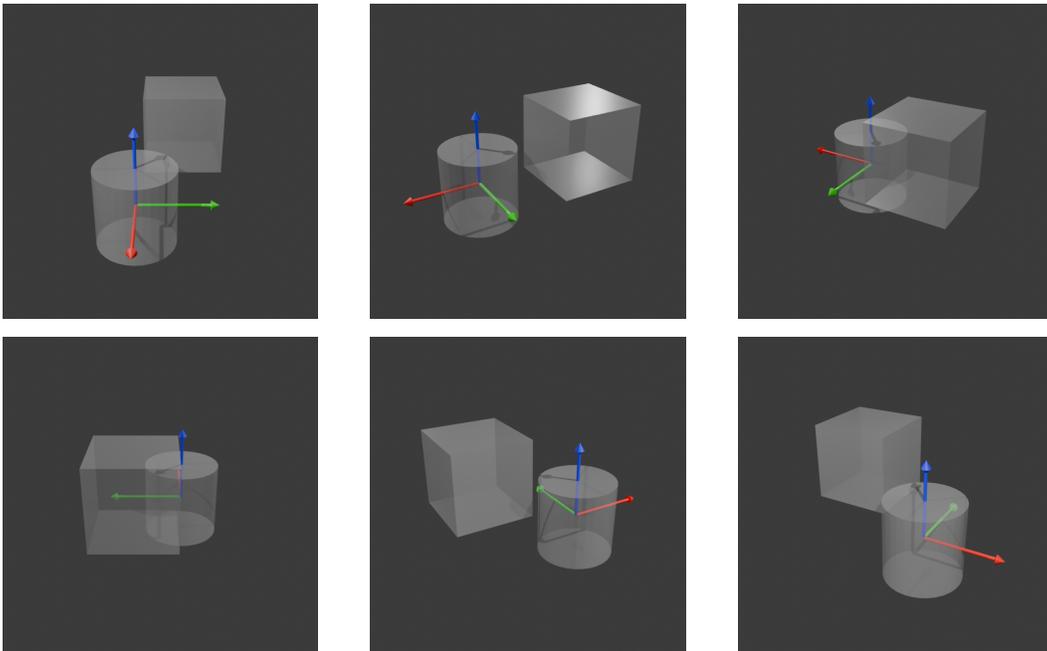


Figure 10: Six isometric views provided to the model for the 3D DoF=3 task.

Reasoning and Answer:

Looking at the different viewpoints of the 3D scene, I can analyze the relative position of the cube compared to the cylinder.

From the images, I can see that:

- X-axis: The cube appears to be positioned mainly in the negative X direction relative to the cylinder (to the left in several views)
- Y-axis: The cube appears to be positioned in the positive Y direction relative to the cylinder (forward in several views)
- Z-axis: The cube and cylinder appear to be at approximately the **same height on the Z-axis**

<answer>(-X, +Y, **0Z**)</answer>

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

A.6.2 2D COLORS (CLAUDE 4 SONNET)

In this task, the model is provided with three canonical 2D views (front, side, and top). The model demonstrates a classic failure in view integration. It correctly processes the x-axis relationship (-X), which is salient in the top and front views. However, when trying to determine depth (Y-axis from the side/top views) and height (Z-axis from the front/side views), it incorrectly concludes the objects are at the "same level."

- **Ground Truth:** (-X, +Y, +Z)
- **Model's Answer:** (-X, 0Y, 0Z)

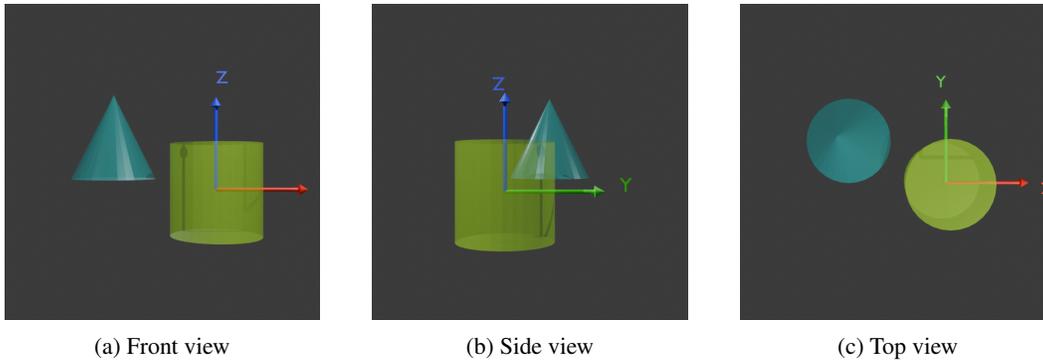


Figure 11: Example task of 2D three view with colors.

Reasoning and Answer:

Based on examining the 3D scene from multiple viewpoints, I can determine the relative position of the teal cone compared to the lime cylinder.
Looking at all three viewpoints:

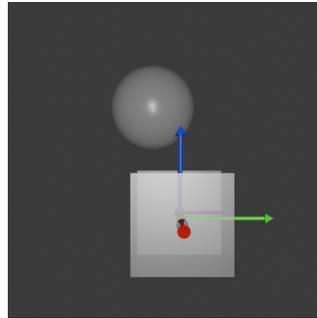
- X-axis: The cone appears to be to the left of the cylinder (negative X direction)
- Y-axis: The cone appears to be at approximately the **same level** as the cylinder on the Y-axis
- Z-axis: The cone appears to be at approximately the **same height** as the cylinder on the Z-axis

<answer>(-X, **0Y**, **0Z**)</answer>

1404 A.6.3 VIEWNAVIGATOR (GPT-5): VLM PERCEPTION FAILURE
1405

1406 This case illustrates a failure where the agent incorrectly determines the position on the Y-axis,
1407 leading to an incorrect final answer. The failure stems from a critical misperception by the VLM,
1408 which then misleads the LLM planner. This single error propagates through the agent's belief state,
1409 causing the LLM planner to develop confirmation bias and eventually accept the wrong conclusion.

- 1410 • **Ground Truth:** $(-X, -Y, +Z)$
- 1411 • **Agent's Final Answer:** $(-X, 0Y, +Z)$



1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424 Figure 12: Image observed by the VLM from a side view ($a_z=0, e_l=90$).
1425

1426 **Reasoning and Answer:**
1427

1428 `<think>Identify axes: green is +Y to the right; blue is +Z upward. Compare centers:
1429 the sphere is directly above the cube with no horizontal (Y) offset, so along Y the centers
1430 align (0). Along Z, the sphere lies above the cube center, so it is in the positive Z direction
1431 (+).</think> <answer>(0Y, +Z)</answer>`
1432

1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

A.6.4 VIEWNAVIGATOR (GEMINI 2.5 FLASH): LLM PLANNING FAILURE

This case demonstrates a fundamental failure in the LLM’s spatial reasoning, where there is a disconnect between its stated goal and its proposed action. It correctly identifies the need for a view of the XZ plane but incorrectly believes that an elevation of 90° will achieve this. A view with an elevation of 90° is a top-down view, which clearly displays the XY plane while making it impossible to accurately judge relative positions along the X-axis. The model’s stated goal is contradicted by its proposed action, highlighting a limited capacity to translate a desired outcome into the correct geometric parameters.

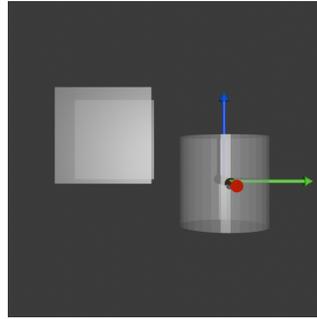


Figure 13: The view generated by the LLM’s proposal ($az=0, el=90$).

Reasoning and Planning:

```
<think> ... A view with azimuth  $0^\circ$  and elevation  $90^\circ$  would provide a side view focusing on the XZ plane, which might help clarify the X axis. This view also gives us a chance to improve our understanding of the Z axis. Therefore, I will choose a view with azimuth  $0^\circ$  and elevation  $90^\circ$ , focusing on the X and Z axes. </think> <answer> "action": "CAPTURE", "view": "az": 0, "el": 90, "axis": ["X", "Z"] </answer>
```