From Words to Worlds: Spatial Intelligence is AI's Next Frontier

# Huge Performance Gap

GPT-5-thinking

**MM1**

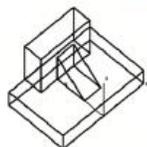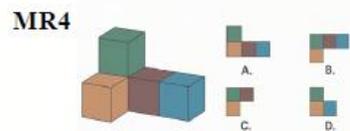**Question:** What is the height of region 1 in meters? **GT:** 2.7m.

**Answer:** ✅ 2m.

**MR2**

**Question:** Given the front, side and top-down view of a 3D object, analyze its structure and reconstruct it in 3D axis.

**Answer:** ✅

**MR3**

**Question:** Generate a 90 degrees top-down view of this scene.

**Answer:** ✅

**MR4**

**Question:** Which option is the correct top-down view of the object? **GT:** B.

**Answer:** ❌ A.

**SR5**

**Question:** Which object is higher in the 3D world space, the clock or the house in the back? **GT:** The house in the back.
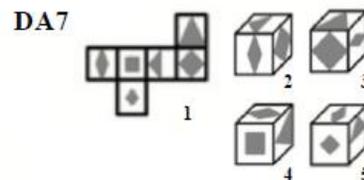
**Answer:** ❌ Clock.

GPT-5-thinking

**PT6**

**Question:** The images are frames from a video. The first image is from the beginning of the video and the second image is from the end. Is the camera moving left or right when shooting the video? **GT:** Left.
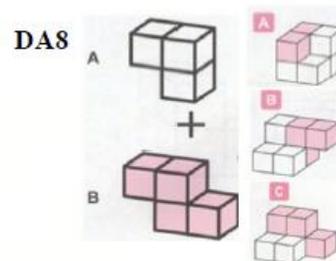
**Answer:** ❌ Right.

**DA7**

**Question:** Flip the shape in image 1 to form a 3D cube. Which of the image 2, 3, 4, 5 is a possible view of the formed cube? **GT:** Image 4.

**Answer:** ❌ Image 2.

**DA8**

**Question:** Which of A, B, C is possible to be built when rotating and combining the two 3D structure in image 1? **GT:** C.

**Answer:** ❌ A and B.

**CR9**

**Question:** How many 3D blocks in the image? **GT:** 8.

**Answer:** ❌ 9.

| Models | VSI [66] | SITE [57] | MMSI [68] | OmniSpatial [23] | MindCube* [69] | STARE [32] | CoreCognition [33] | SpatialViz [55] |
|---|---|---|---|---|---|---|---|---|
| Metric | MRA, Acc | CAA | Acc | Acc | Acc | Acc, F1 | Acc | Acc |
| **Random Choice** | 34.00 | 0.0 | 25.00 | 24.98 | 32.35 | 34.80 | 33.93 | 25.08 |
| **Proprietary Models** | | | | | | | | |
| Seed-1.6-2025-06-15 [51] | 49.91 | 54.61 | 38.30 | 49.32 | 48.75 | 46.06 | 77.17 | 34.58 |
| Gemini-2.5-pro-2025-06 [52] | 53.57 | 57.06 | 38.00 | 55.38 | 57.60 | 49.14 | 76.70 | 42.71 |
| Grok-4-2025-07-09 [62] | 47.92 | 47.01 | 37.80 | 46.84 | **63.56** | 26.90 | 79.27 | 19.40[†] |
| GPT-5-nano-2025-08-07 [45] | 43.22 | 35.81 | 28.90 | 47.81 | 41.48 | 46.05 | 67.92 | 35.59 |
| GPT-5-mini-2025-08-07 [45] | 48.67 | 52.47 | 34.10 | 55.52 | 56.69 | 52.51 | 77.77 | 44.66 |
| GPT-5-2025-08-07 [45] | **55.03** | **61.88** | **41.80** | **59.90** | 56.30 | **54.59** | **84.37** | **51.27** |
| **Open-source Models** | | | | | | | | |
| Qwen2.5-VL-3B-Instruct [1] | 27.00 | 33.14 | 28.60 | 42.47 | 37.60 | 37.83 | 60.19 | 21.86 |
| Qwen2.5-VL-7B-Instruct [1] | 32.30 | 37.64 | 26.80 | 39.07 | 36.05 | 35.03 | 62.16 | 26.78 |
| Qwen2.5-VL-72B-Instruct [1] | 35.77 | 47.41 | **32.50** | 47.81 | 42.40 | 38.37 | 69.22 | **32.54** |
| InternVL3-8B [79] | 42.14 | 41.15 | 28.00 | 46.25 | 41.54 | 41.36 | 60.92 | 30.00 |
| InternVL3-78B [79] | 47.55 | **52.72** | 30.50 | **50.95** | **49.52** | **42.00** | **71.16** | 31.10 |
| InternVL3.5-8B [56] | 56.05 | 43.79 | 27.30 | 46.71 | 42.50 | 40.18 | 66.40 | 23.98 |
| Qwen3-8B-Instruct [65] | **57.90** | 45.83 | 31.10 | 45.73 | 29.42 | 39.76 | 69.67 | 17.54[†] |
| **Human Evaluation** | | | | | | | | |
| Δ(Best Model,Human) | -21.3 | -5.62 | -55.40 | -32.73 | -30.99 | -42.06 | -2.61 | -31.19 |
| **Human** | **79.2** | **67.5** | **97.2** | **92.63** | **94.55** | **96.50** | **86.98** | **82.46** |

# Chain-of-Thought?



Saining Xie, Manling Li, Jiajun Wu, Li Fei-Fei, et al, *MindCube*, 2025

| CoT Style | Average # Output Token | VSI-Bench: Obj. Rel. Direction | | | |
|---|---|---|---|---|---|
| | | Overall | Easy | Medium | Hard |
| **InternVL3-8B** | 1 | 39.3 | 48.8 | 47.0 | 21.9 |
| No CoT | 3.4 | 54.9 | 62.2 | 55.8 | 46.6 |
| CoT-GPT-5 | 1070.7 | 40.0 | 41.4 | 43.1 | 36.1 |
| CoT-MindCube-Aug-CGMap | 1490.6 | 39.9 | 45.9 | 42.7 | 33.7 |
| CoT-SenseNova-SI-CGMap | 2262.8 | 47.9 | 53.9 | 51.3 | 41.0 |

**Text-based CoT is <u>not efficient or generalizable</u>!**

Metric Measurement · Mental Reconstruction · Spatial Relations · Perspective-taking · Deformation & Assembly · Comprehensive Reasoning

Hokin Deng *et al., CoreCognition*, 2025

# Scaling Spatial Intelligence

Minimize Redundant Information

Avoid Insufficient Information

Avoid Ambiguity

# Scaling Spatial Intelligence

| Models | VSI-Bench [56] | MMSI-Bench [60] | MindCube* [62] | ViewSpatial [29] | SITE [50] | MMB-EN [33] |
|---|---|---|---|---|---|---|
| Metric | MRA, Acc | Acc | Acc | Acc | CAA | Acc |
| **Human** | **79.2** | **97.2** | **94.5** | - | **67.5** | - |
| Random Choice | 34.0 | 25.0 | 33.0 | 26.3 | 0.0 | 25.0 |
| **Proprietary Models** | | | | | | |
| Seed-1.6-2025-06-15 [42] | 49.9 | 38.3 | 48.7 | 43.8 | 54.6 | 87.5 |
| Gemini-2.5-Pro-2025-06 [43] | 53.5 | 38.0 | 57.6 | 46.0 | 57.0 | **90.1** |
| Grok-4-2025-07-09 [54] | 47.9 | 37.8 | 63.5 | 43.2 | 47.0 | 86.3 |
| GPT-5-2025-08-07 [37] | **55.0** | 41.8 | 56.3 | 45.5 | 61.8 | 85.2 |
| Gemini-3-Pro-Preview [19] | 52.5 | **45.2** | **70.8** | **50.3** | **62.2** | - |
| **Open-source General Models** | | | | | | |
| Bagel-7B-MoT [15] | 31.4 | 31.0 | 34.7 | 41.3 | 37.0 | 82.8 |
| Qwen2.5-VL-3B-Instruct [3] | 27.0 | 28.6 | 37.6 | 31.9 | 33.1 | 77.4 |
| Qwen2.5-VL-7B-Instruct [3] | 32.3 | 26.8 | 36.0 | 36.8 | 37.6 | 82.6 |
| Qwen3-VL-2B-Instruct [13] | 50.3 | 28.9 | 34.5 | 36.9 | 35.6 | 75.1 |
| Qwen3-VL-8B-Instruct [13] | 57.9 | 31.1 | 29.4 | 42.2 | **45.8** | **84.6** |
| InternVL3-2B [65] | 32.9 | 26.5 | 37.5 | 32.5 | 30.0 | 79.7 |
| InternVL3-8B [65] | 42.1 | 28.0 | 41.5 | 38.6 | 41.1 | 81.7 |
| **Open-source Spatial Intelligence Models** | | | | | | |
| MindCube-3B-RawQA-SFT [62] | 17.2 | 1.7 | **51.7** | 24.1 | 6.3 | 32.3 |
| SpatialLadder-3B [30] | 44.8 | 27.4 | 43.4 | 39.8 | 27.9 | 72.5 |
| Spatial-MLLM-4B [52] | 46.3 | 26.1 | 33.4 | 34.6 | 18.0 | 64.5 |
| SpaceR-7B [38] | 41.5 | 27.4 | 37.9 | 35.8 | 34.2 | 75.4 |
| ViLaSR-7B [53] | 44.6 | 30.2 | 35.1 | 35.7 | 38.7 | 81.1 |
| VST-3B-SFT [58] | 57.9[†] | 30.2[†] | 35.9 | **52.8** | 35.8 | 80.9[†] |
| VST-7B-SFT [58] | 60.6[†] | **32.0**[†] | 39.7 | 50.5 | 39.6 | 83.3[†] |
| Cambrian-S-3B [59] | 57.3[†] | 25.2 | 32.5 | 39.0 | 28.3 | 76.0[†] |
| Cambrian-S-7B [59] | **67.5**[†] | 25.8 | 39.6 | 40.9 | 33.0 | 80.4[†] |
| **Ours** | | | | | | |
| SenseNova-SI Bagel-7B-MoT | 41.6(+32.5%) | 36.2(+16.8%) | 50.8(+46.4%) | 50.3(+21.8%) | 41.6(+12.4%) | 83.4(+0.72%) |
| SenseNova-SI Qwen3-VL-8B | 62.9(+8.6%) | 37.5(+20.6%) | 74.8(+154.4%) | 48.4(+14.7%) | 50.1(+9.3%) | 83.5(-1.30%) |
| SenseNova-SI InternVL3-2B | 63.7(+93.6%) | 34.2(+29.1%) | 41.8(+11.5%) | 52.6(+61.8%) | 36.7(+22.3%) | 78.9(-1.00%) |
| SenseNova-SI InternVL3-8B | **68.7**(+63.2%) | **43.3**(+54.6%) | **85.6**(+106.3%) | **54.6**(+41.5%) | 47.7(+16.1%) | **84.9**(+3.92%) |

**Best Proprietary**

**Best Expert**

**SenseNova-SI**

[Obj-Orient View (Ego-Exo4D)]

Which egocentric view image correctly matches the exocentric view?

- B.

[SITE Mov&Nav-Maze]

How many right turns are there in the provided path (marked by Blue) from S (green block) to E (red block)?

— C. 2.

90.4%

[MMSI Pos-Cam-Cam]

Camera was facing the west side of the room when the first picture was taken, which direction is the camera facing in the room when the second picture is taken?

— C. North.

23.8%

[Cam Rotation (MessyTable)]

From image 1 to 2. What is the rotation direction?

- Rotated to the Left.

[MMSI Pos-Cam-Cam]

Camera coordinate system +Y up, -Z forward, right-handed. How can the first image be obtained from the second image?

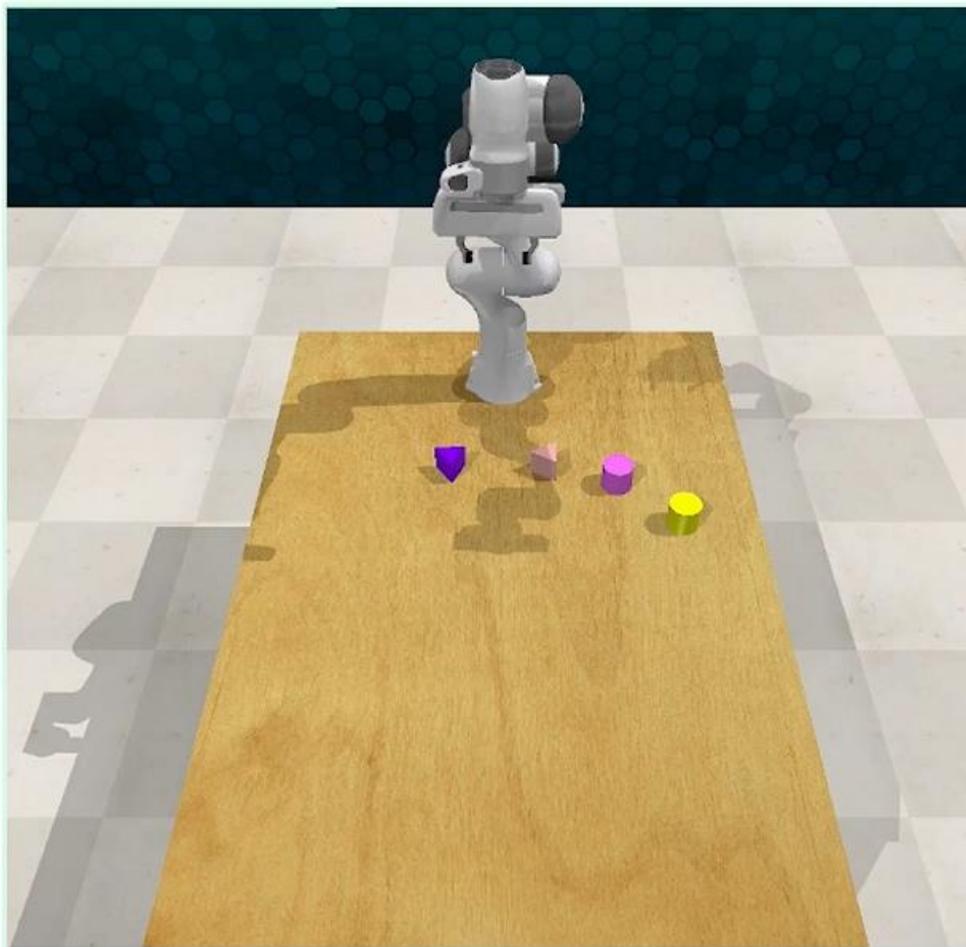—A. Rotate a positive angle around the Y axis.

66.7%

[MMSI Attr-Appr]

How many different drawers with a width greater than their height appear in total?
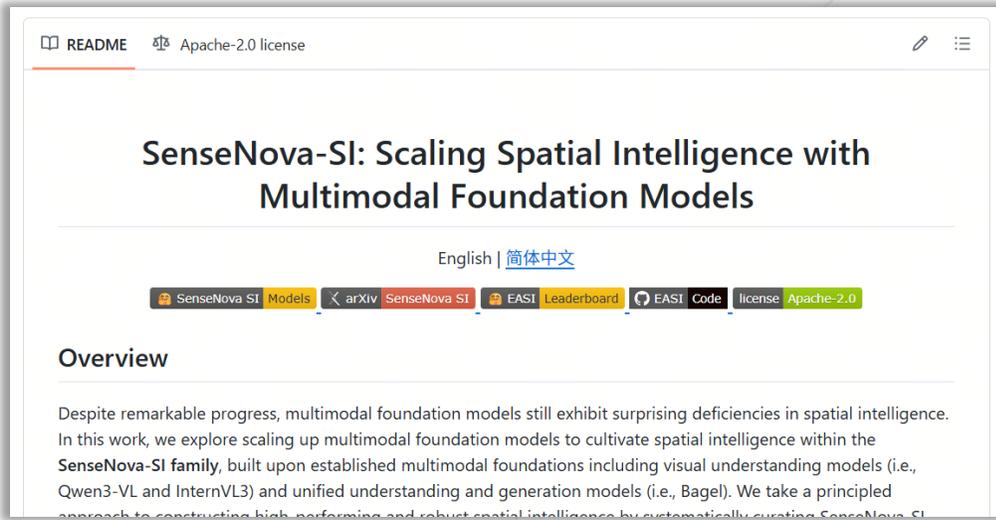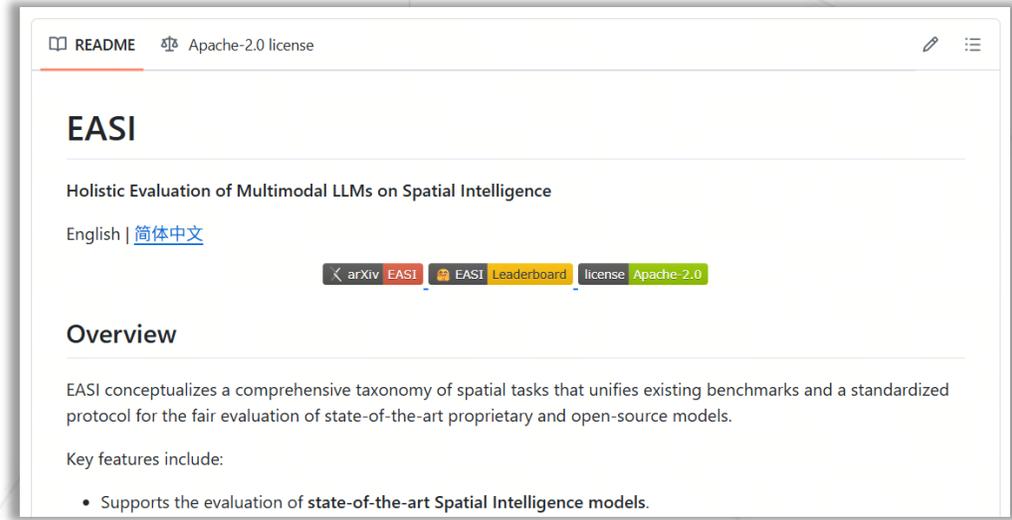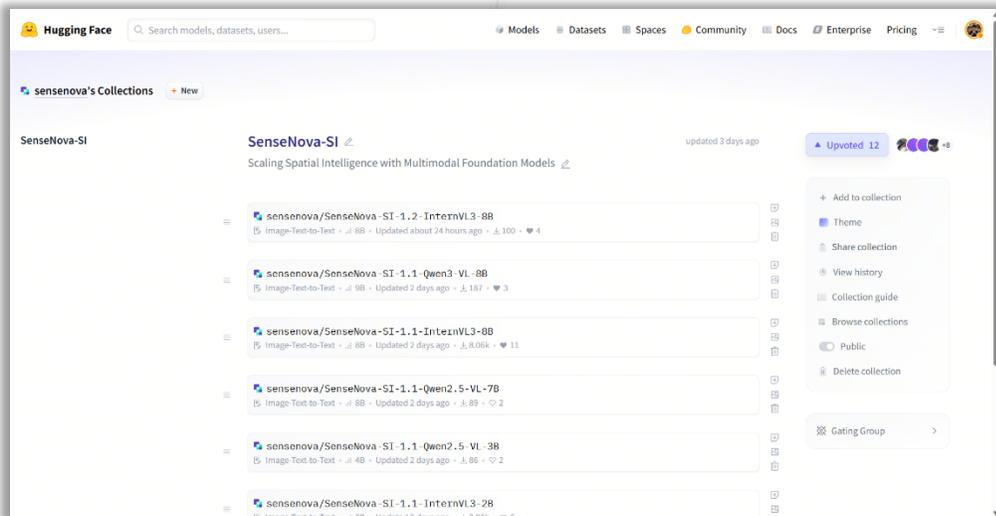
— D. 5.

25.6%

SSI-InternVL3-8B

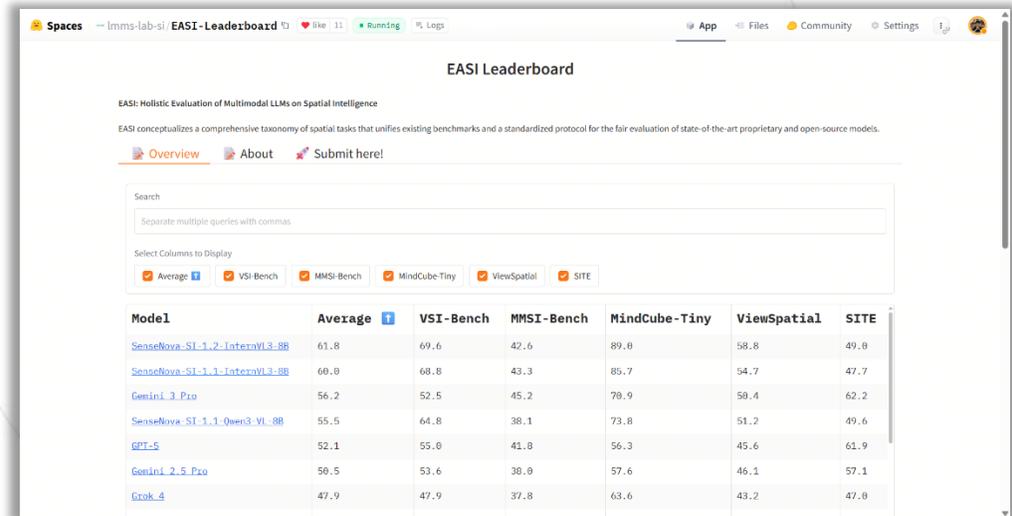Stack the left triangular prism on top of the right cylinder.

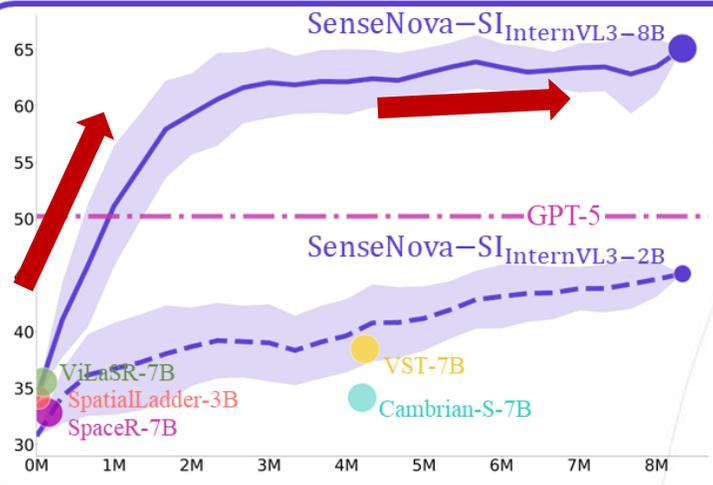# Open-source
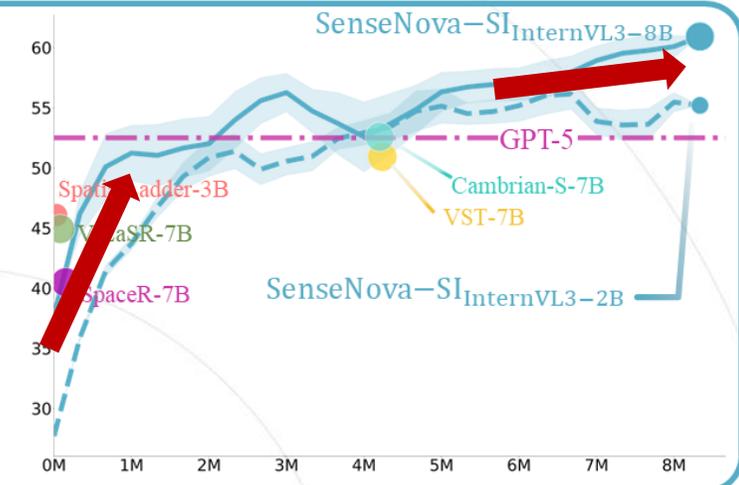

Code for SenseNova-SI


EASI: Evaluation Toolbox


Weights for SenseNova-SI


EASI Leaderboard

# What is Next?



Perspective-taking

SenseNova−SI$_{InternVL3-8B}$
SenseNova−SI$_{InternVL3-2B}$
GPT-5
ViLaSR-7B
SpatialLadder-3B
SpaceR-7B
VST-7B
Cambrian-S-7B

Metric Measurement

1.5m
3m

SenseNova−SI$_{InternVL3-8B}$
GPT-5
SpatialLadder-3B
ViLaSR-7B
SpaceR-7B
Cambrian-S-7B
VST-7B
SenseNova−SI$_{InternVL3-2B}$

VSI: 68.7

MMSI: 43.3

SITE: 47.7

MMBench-EN: 84.9

ViewSpatial: 54.6

MindCube: 85.6

SenseNova-SI$_{InternVL3-8B}$
InternVL3-8B
SpaceR-7B
VST-7B
Qwen3-VL-8B
ViLaSR-7B
Cambrian-S-7B

Spatial Relations

GPT-5
SenseNova−SI$_{InternVL3-8B}$
SenseNova−SI$_{InternVL3-2B}$
VST-7B
Cambrian-S-7B
ViLaSR-7B
SpaceR-7B
SpatialLadder-3B

top
left & behind
right
front

Comprehensive Reasoning

SenseNova−SI$_{InternVL3-8B}$
GPT-5
VST-7B
Cambrian-S-7B
ViLaSR-7B
SpaceR-7B
SpatialLadder-3B
SenseNova−SI$_{InternVL3-2B}$

# Thank you!