



**ICML**  
International Conference  
On Machine Learning

# Active Exploring like a Pigeon: Reinforcing Spatial Reasoning via Agentic Vision-Language Models

Wei Deng<sup>12</sup>, Xianlin Zhang<sup>23</sup>, Mengshi Qi<sup>\*12</sup>

43<sup>nd</sup> International Conference on Machine Learning (ICML), 2026

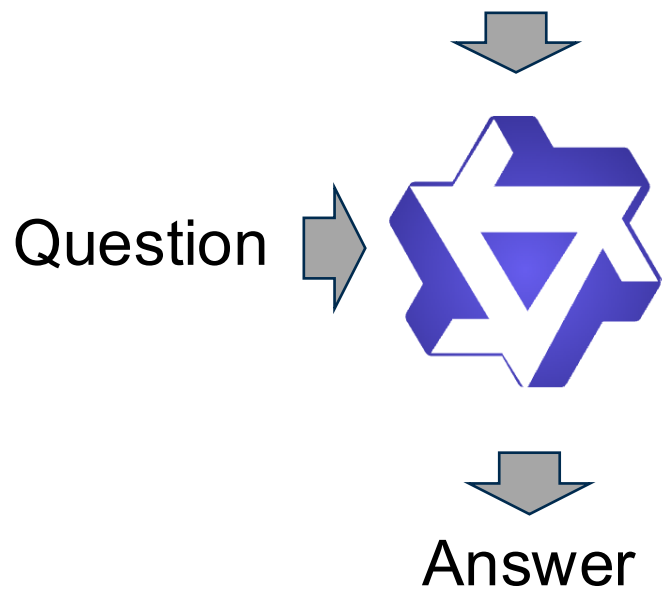
<sup>1</sup>State Key Laboratory of Networking and Switching Technology

<sup>2</sup>School of Digital Media & Design Arts

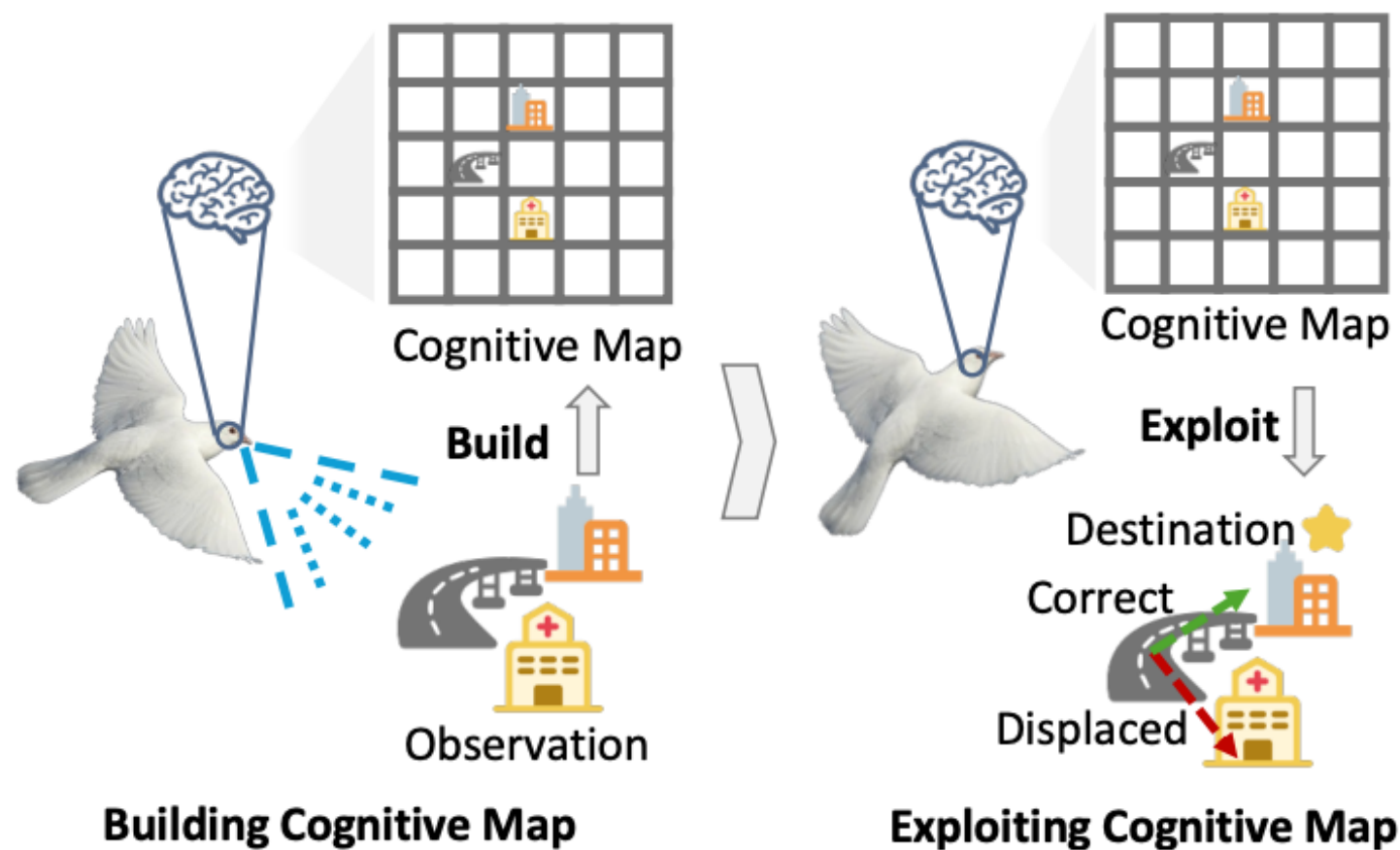
<sup>3</sup>Beijing University of Posts and Telecommunications, China

\*Corresponding author

## All Frames

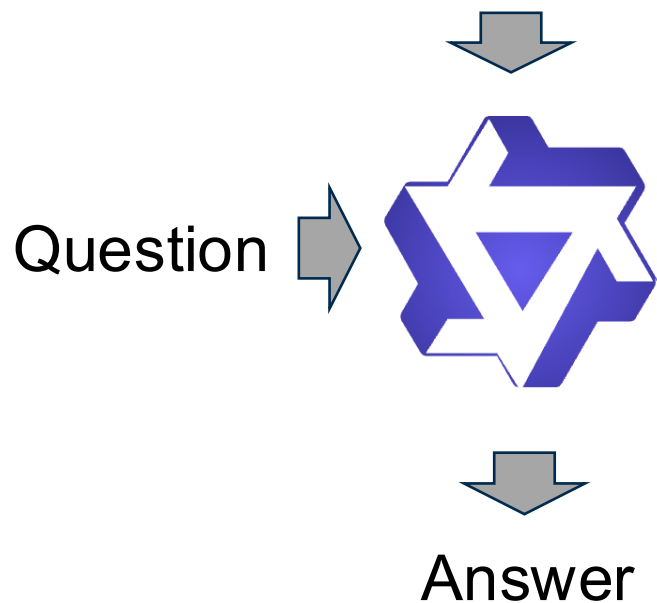


**VLM serves as a passive observer for spatial reasoning**

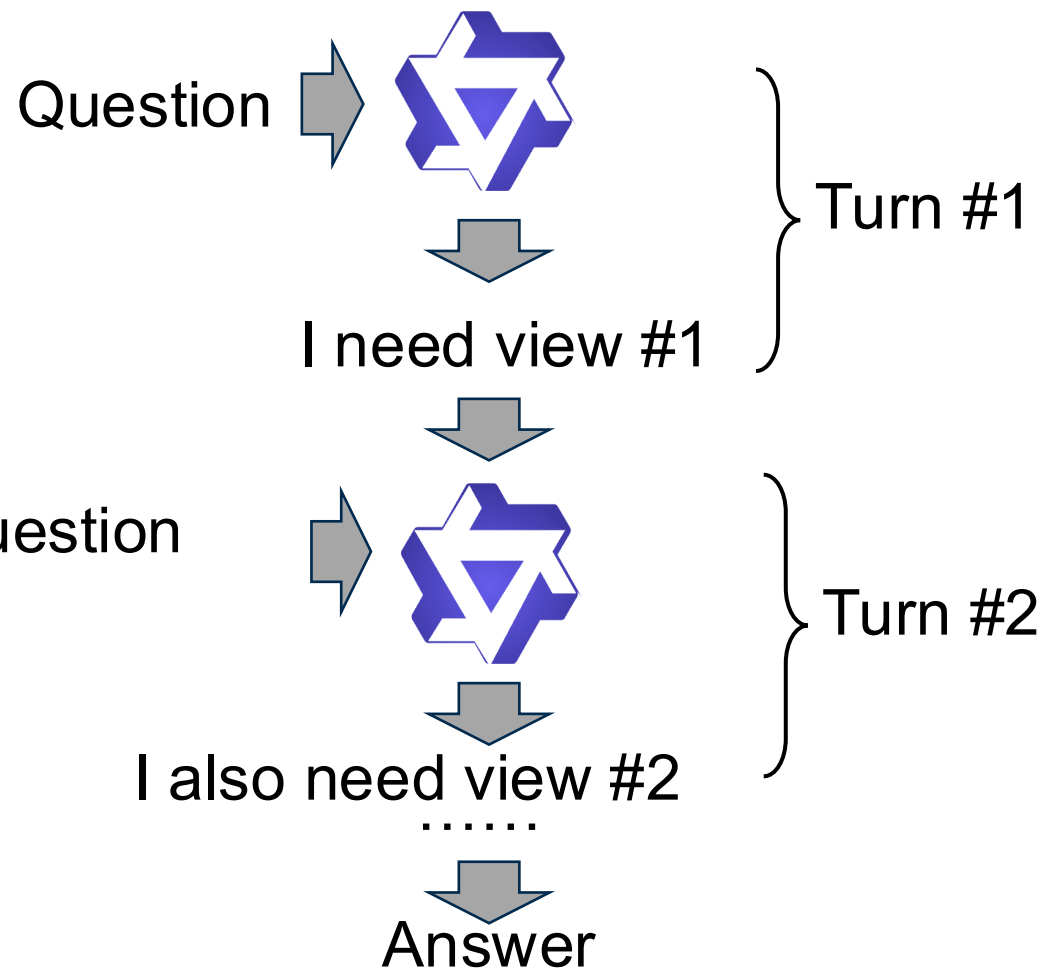


**Homing pigeons build and exploiting cognitive maps for navigation**

All Frames

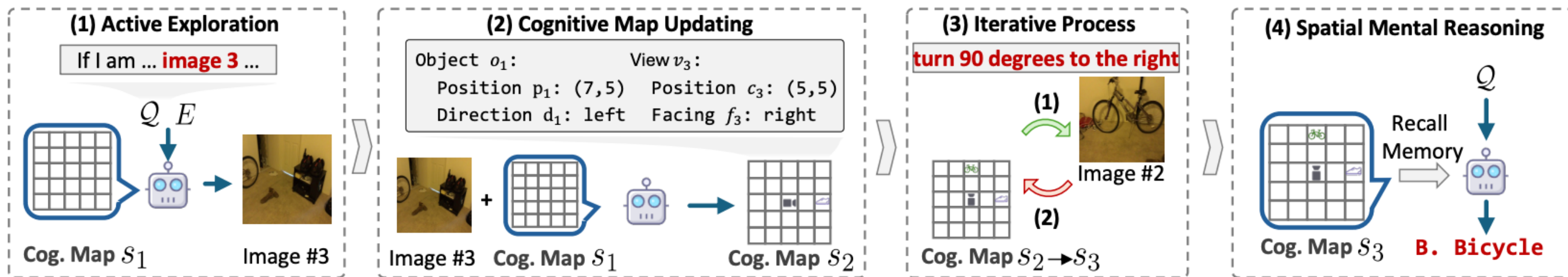


**Previous: Passive perception**



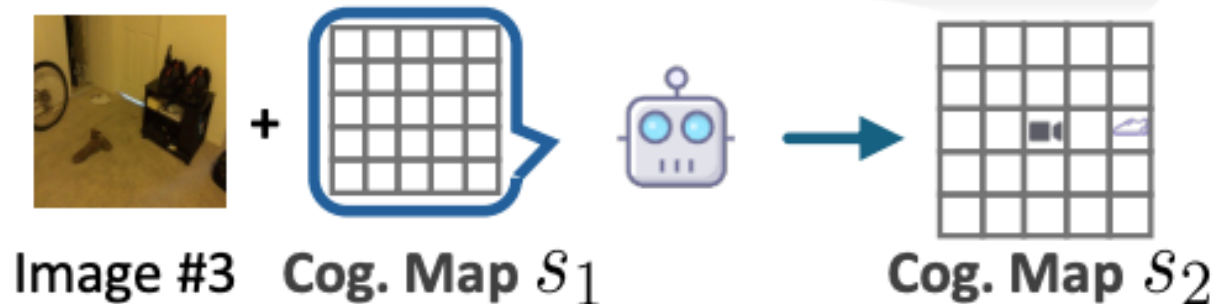
**Ours: Active perception**

1. Active exploration
2. Cognitive map updating
3. Iterative process
4. Spatial mental reasoning

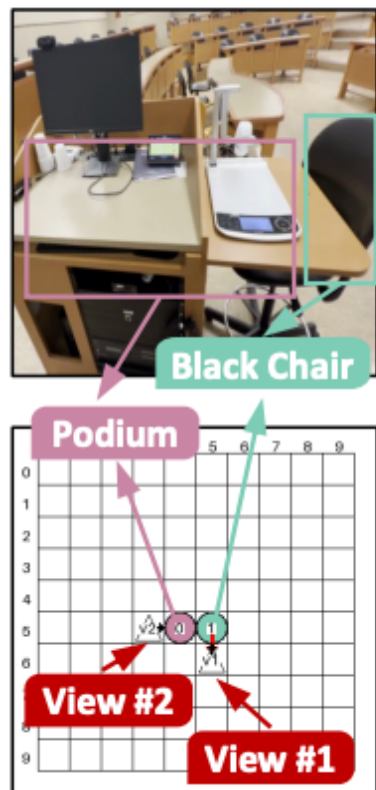


- A unified top-down coordinate system
- Storing objects with their positions and orientations, and camera viewpoints
- A persistent memory

Object $o_1$ :	View $v_3$ :
Position $p_1$ : (7,5)	Position $c_3$ : (5,5)
Direction $d_1$ : left	Facing $f_3$ : right

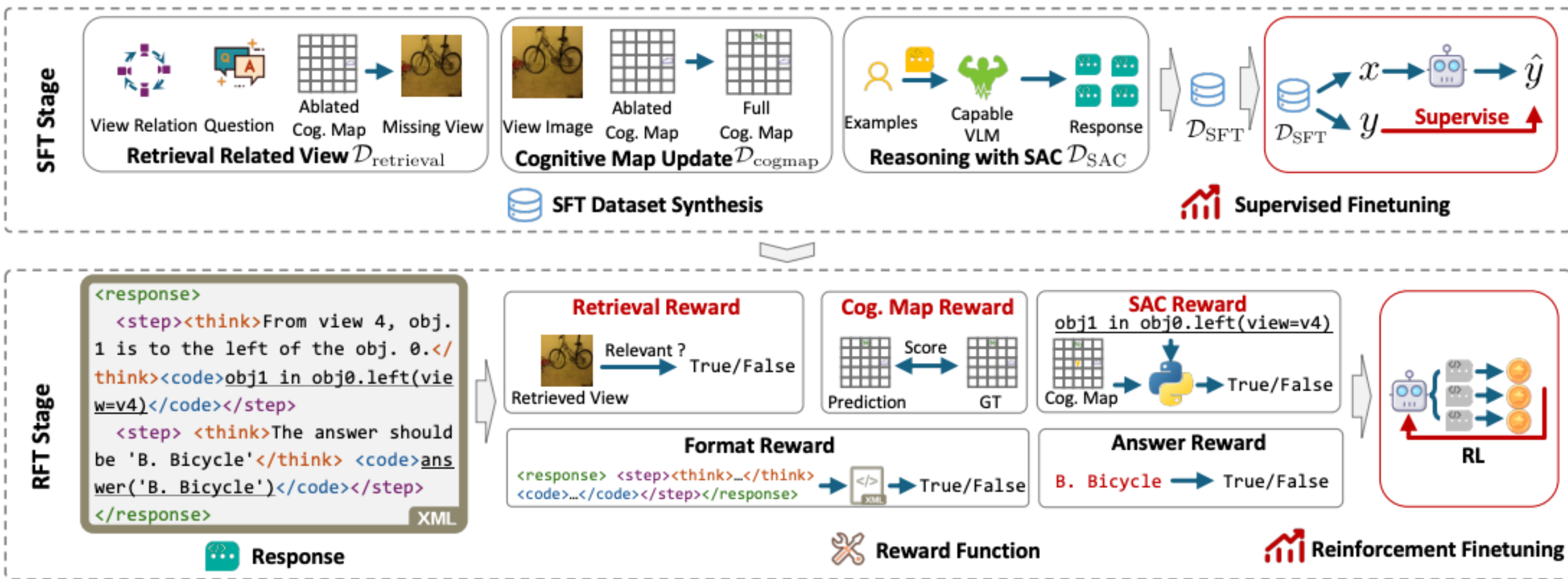


- SAC is an executable Python expression
- SAC describe spatial relationships
- Example: object\_1 on the left of object\_2 from the perspective of view 4 -> obj1 in obj0.left(view=v4)



```
<step>
  <think>According to the cognitive map, the first view
  (v1) is at `(5, 6)` and the second view (v2) is at `(
  3, 5)`. The "podium" (obj0) is at `(4, 5)` and the "b
  lack chair" (obj1) is at `(5, 5)`. The second view is
  in the forward direction relative to the first view.<
  /think>
  <code>v2 in v1.front(view=v1) and v2 in v1.left(view=
  v1)</code>
</step>
<step>
  <think>So, the answer is B.</think>
  <code>agent.answer('B')</code>
</step>
```

- SFT: training with view retrieval, cognitive map updating, and spatial reasoning with SAC
- RFT: Finetuning with GRPO and dense reward functions



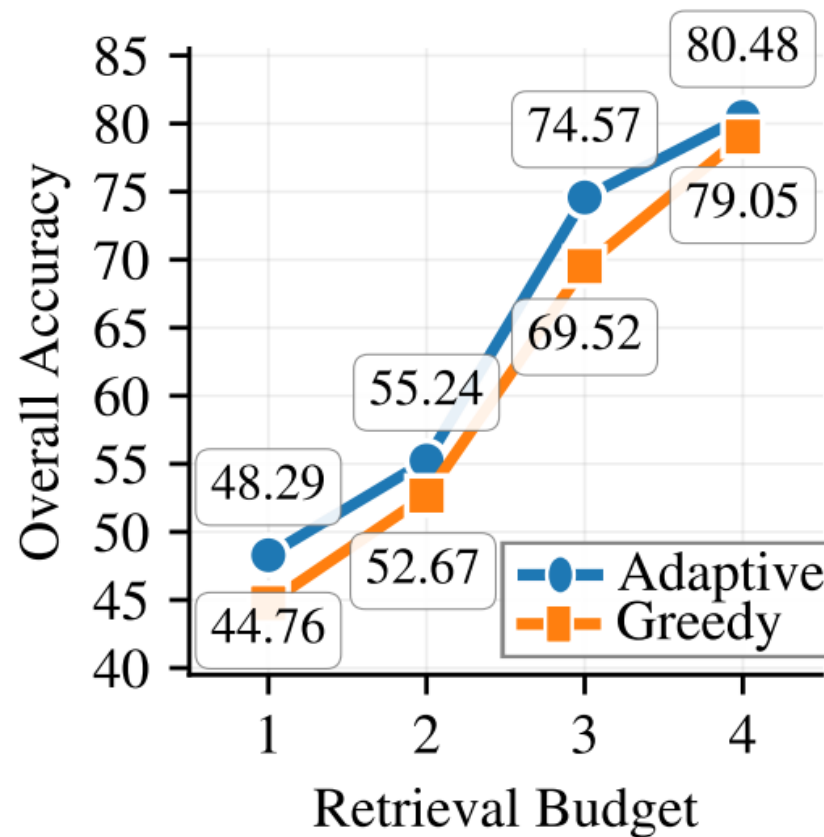
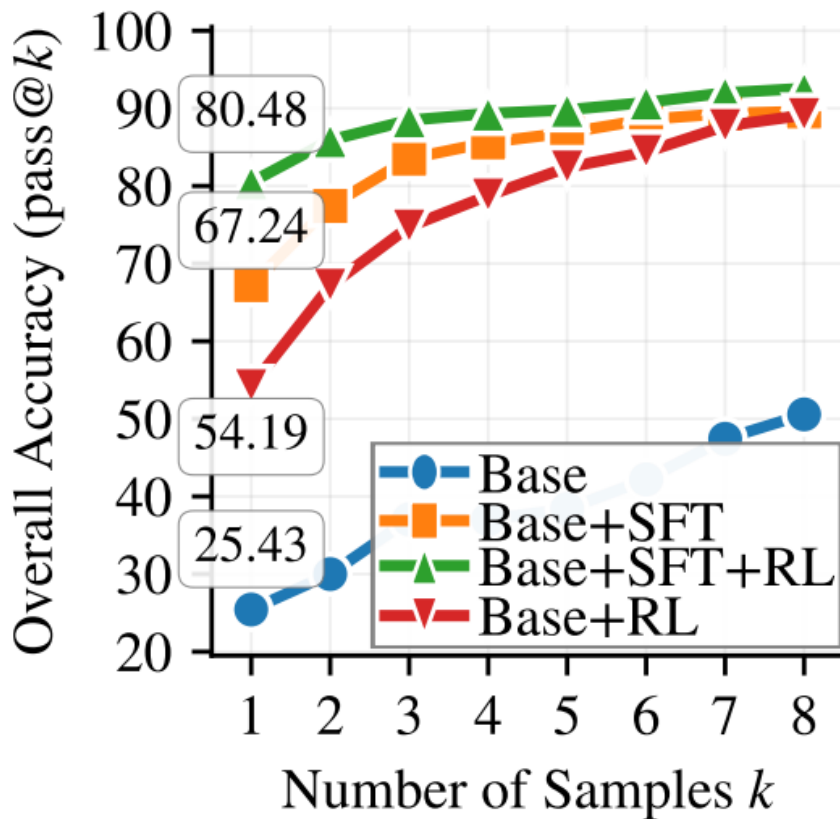
- We achieve the best performance on MindCube
- We achieve a significant improvement on the Rotation subset, which is important for EAI

Methods	MindCube-Tiny Benchmark				
	Features	Overall ↑	Rotation ↑	Among ↑	Around ↑
<i>Baseline</i>					
Random (chance)	–	32.35	36.36	32.29	30.66
Random (frequency)	–	33.02	38.30	32.66	35.79
<i>Open-Source VLMs</i>					
Qwen2.5-VL-7B-Instruct (Bai et al., 2025)	Passive	29.26	38.76	29.50	21.35
Qwen2.5-VL-3B-Instruct (Bai et al., 2025)	Passive	33.21	37.37	33.26	30.34
<i>Proprietary Models</i>					
GPT-4o (OpenAI, 2024)	Passive	38.81	32.65	40.17	29.16
Claude-4-Sonnet-20250514 (Anthropic, 2024)	Passive	44.75	48.42	44.21	47.62
<i>Spatial VLMs</i>					
Spatial-MLLM (Wu et al., 2025a)	Passive, 2D+3D	32.06	38.39	20.92	32.82
Space-Qwen (Chen et al., 2024a)	Passive	33.28	38.02	33.71	26.32
MindCube <sub>Qwen2.5-VL-3B</sub> (Yin et al., 2025)	Passive, Cog. Map	70.7	48.0	79.2	68.4
3DThinker <sub>Qwen2.5-VL-3B</sub> (Chen et al., 2026)	Passive, 3D Rec.	<u>75.2</u>	<u>55.5</u>	<b>81.8</b>	<u>75.2</u>
Ours	Active, Cog. Map	<b>80.5</b> ↑5.3	<b>85.0</b> ↑29.5	<u>81.0</u> ↓0.8	<b>75.6</b> ↑0.4

- Active perception pipeline outperforms passive
- All reward components contribute to the final performance
- The cognitive map provides clear gains over plain context

Setting	Acc. ↑	Pre. ↑	Rec. ↑	F1 ↑
<i>Study of Perception Pipeline</i>				
Passive	27.5	26.7	27.3	26.6
Active (Ours)	<b>38.5</b> ↑11	<b>30.6</b> ↑3.9	<b>29.2</b> ↑1.9	<b>29.0</b> ↑2.4
<i>Reward Component Ablation</i>				
Ours Full	<b>80.4</b>	<b>79.2</b>	<b>79.3</b>	<b>79.2</b>
Ours w/o $R_{\text{retrieval}}$	72.5↓7.9	57.4↓21.8	57.0↓22.3	57.1↓22.1
Ours w/o $R_{\text{cogmap}}$	72.6↓7.8	57.3↓21.9	57.0↓22.3	57.1↓22.1
Ours w/o $R_{\text{SAC}}$	70.2↓10.2	55.5↓23.7	54.5↓24.8	54.9↓24.3
<i>Study of Memory Mechanism</i>				
Context	50.9	36.2	34.8	34.8
Cog. Map (Ours)	<b>54.2</b> ↑3.3	<b>40.6</b> ↑4.4	<b>39.4</b> ↑4.6	<b>39.5</b> ↑4.7

- Both of SFT and RFT post-training stages are essential
- Our adaptive retrieval strategy always surpasses the greedy one.





- Re-formulate spatial reasoning and propose an agentic pipeline
- Propose SAC for RLVR
- Achieve state-of-the-art performance on MindCube



**ICML**  
International Conference  
On Machine Learning

**Thank You!**