

When and How Much to Imagine: Adaptive Test-Time Scaling with World Models for Visual Spatial Reasoning

Anonymous CVPR submission

Paper ID 51

Abstract

001 *Despite rapid progress in Multimodal Large Language*
002 *Models (MLLMs), visual spatial reasoning remains unreli-*
003 *able when correct answers depend on how a scene would*
004 *appear under unseen or alternative viewpoints. Recent*
005 *work addresses this by augmenting reasoning with world*
006 *models (WM) for visual imagination, but questions such as*
007 *when imagination is actually necessary, how much of it is*
008 *beneficial, and when it becomes harmful, remain poorly un-*
009 *derstood. In practice, indiscriminate imagination can in-*
010 *crease computation and even degrade performance by in-*
011 *roducing misleading evidence. In this work, we present*
012 *an in-depth analysis of test-time visual imagination as a*
013 *controllable resource for spatial reasoning. We first study*
014 *when static visual evidence is sufficient, when imagina-*
015 *tion improves reasoning, and how excessive or unneces-*
016 *sary imagination affects accuracy and efficiency. To support*
017 *this analysis, we then introduce AVIC, an adaptive test-*
018 *time framework with world models that explicitly reasons*
019 *about the sufficiency of current visual evidence before se-*
020 *lectively invoking and scaling visual imagination. Finally,*
021 *to further learn this gating and planning behavior without*
022 *any annotation of when and how much to imagine, we in-*
023 *troduce AVIC-R, which trains the policy end-to-end via*
024 *GRPO from QA-correctness rewards and penalties by imag-*
025 *ination cost. Across spatial reasoning benchmarks (SAT,*
026 *MMSI) and an embodied navigation benchmark (R2R), our*
027 *results reveal clear scenarios where imagination is critical,*
028 *marginal, or detrimental, and show that selective control*
029 *can match or outperform fixed imagination strategies with*
030 *substantially fewer world-model calls and language tokens.*
031 *Our AVIC-R surpasses strong proprietary baselines in-*
032 *cluding GPT-4o and GPT-4.1 while invoking the WM less*
033 *often. Overall, our findings highlight the importance of an-*
034 *alyzing and controlling test-time imagination for efficient*
035 *and reliable spatial reasoning*¹.

¹paper under peer review

1. Introduction

Recent advances in multimodal large language models (MLLMs) [19, 22] have led to impressive progress in visual understanding and reasoning across various tasks. These models can follow natural language instructions, perceive visual scenes, and reason over multimodal input to support decision making. Despite the progress, *visual spatial reasoning* remains a persistent challenge [7, 30, 33, 41], particularly for questions whose answer depends on unseen regions, viewpoint changes, or transformations that cannot be reliably inferred from a single static observation.

A natural way to address this challenge, mirroring how humans operate, is through *visual imagination* [18]: when the observed visual evidence is insufficient, people mentally simulate how a scene would appear from alternative viewpoints or after potential movements, leveraging strong world priors learned from years of physical interaction and visual experience. Inspired by this intuition, recent work [5, 28, 43] has begun to integrate MLLMs with visual world models that can generate controlled novel views conditioned on hypothetical action at inference time. However, existing approaches often invoke visual imagination using fixed and exhaustive strategies (see Figure 1), without first reasoning about whether additional imagination is necessary and helpful. This lack of deliberation can lead to problematic imagination, producing misleading (Figure 1 (b)) or redundant (Figure 1 (c)) views that not only incur substantial computational overhead but can also distract downstream reasoning and result in worse performance than relying on the original observation alone. Through a systematic analysis of always-on imagination (Section 3), we show that such strategies are both inefficient and unreliable, motivating the need for more adaptive use of world models.

Based on these observations, we aim to answer two fundamental questions for visual spatial reasoning with world model imagination: *when* should a model invoke visual imagination, and *how much* imagined visual evidence is necessary if imagination is required. Rather than treating visual imagination as an always-on operation, we seek to



Figure 1. Different cases in always-on visual imagination. Imagined views are generated independently for different beam-searched actions (shown by multiple arrows). Case 1 (Helpful): Visual imagination reveals previously unseen viewpoints, enabling helpful spatial reasoning. Case 2 (Misleading): Imagination fails to preserve task-relevant objects (e.g., the white table in the red box), resulting in incorrect spatial inference and wrong answers. Case 3 (Unnecessary): The required information is already clearly observable in the original view (e.g., the bathtub in the blue box), making additional imagined views redundant.

075 make it a controllable, self-adaptive component during inference time. In this paper, we introduce Adaptive Visual
 076 **Imagination Control (AVIC)**, a framework that gates and plans world-model usage with an explicit policy model.
 077 Given an observation and a question, the policy first reasons about the sufficiency of the available visual evidence and
 078 conditionally decides whether to invoke the world model. If it decides not to, it answers directly from the observed view;
 079 otherwise, it generates a dynamic-length action plan that specifies how the imagination should move or reorient to ac-
 080 quire informative viewpoints, which are rendered by the visual world model and consumed by a downstream reasoner.
 081 A trajectory-level verifier then selects the most informative imagined trajectory among multiple policy samples, in con-
 082 trast to prior beam-search approaches that score isolated keyframes. It enables instance-dependent test-time scaling,
 083 allowing us to move beyond fixed imagination strategies and to study visual spatial reasoning systematically.

093 While AVIC can be instantiated in a training-free manner via prompting and self-consistency sampling from a
 094

strong MLLM, training a stronger policy is non-trivial: no ground-truth supervision exists for what the optimal imag-
 ination trajectory should be, ruling out standard supervised fine-tuning or behavior cloning. To overcome this chal-
 lenge and learn this behavior end-to-end, we further propose **AVIC-R**, which trains the gating policy on top of
 small open-source model (Qwen2.5-VL [4]) via reinforcement learning, using a composite reward built around QA
 correctness that requires no human supervision on *when* imagination is necessary. Specifically, we adopt Group-
 Relative Policy Optimization (GRPO) [32] to compute group-normalized advantages over rollouts that share the
 same prompt; the reward augments QA correctness with an action-count cost that discourages over-imagination,
 a wrong-skip penalty that prevents the policy from collapsing to always-skip, and a format penalty for malformed
 outputs. Together, these designs keep training from collapsing and free the policy to discover optimal gating and
 planning behaviors directly from the reward signal.

We evaluate the proposed framework on challenging spa-

115 tial reasoning benchmarks (SAT [30], MMSI [42]), and the
116 navigation benchmark R2R [2]. Across these settings, adap-
117 tive test-time scaling achieves SoTA or competitive perfor-
118 mance while requiring substantially fewer extra language
119 tokens and world-model calls compared to fixed imagina-
120 tion strategies. Notably, with only a small training set and
121 LoRA updates, AVIC-R boosts a 7B open-source policy
122 enough that the resulting pipeline outperforms variants us-
123 ing GPT-4o or GPT-4.1 as the policy model. Overall, be-
124 yond improved performance, our results reveal that the ben-
125 efits of visual imagination are highly instance-dependent
126 and structured by the nature of the spatial reasoning query.
127 In particular, we find that world models are most beneficial
128 for action-conditioned spatial reasoning, where answers de-
129 pend on how a scene would evolve under specific move-
130 ments or viewpoint changes, while offering limited gains
131 for queries that can be resolved from existing observations.
132 Our analysis shows that effective visual spatial reasoning
133 typically requires only targeted imagination, and excessive
134 or indiscriminate simulation can introduce noise and de-
135 grade performance. Together, these findings indicate that
136 visual imagination is a selective, query-dependent test-time
137 resource, requiring adaptive, uncertainty-aware allocation
138 of world-model computation.

139 2. Related Work

140 **Visual Spatial Reasoning with MLLMs.** The rapid evolu-
141 tion of Multimodal Large Language Models (MLLMs) has
142 made significant progress in various downstream tasks [8,
143 10, 12, 21, 25, 29, 46–49, 54]. In particular, spatial
144 reasoning has attracted considerable attention due to its
145 critical role in bridging visual perception with embodied
146 tasks [36, 53, 56, 57]. However, recent comprehensive eval-
147 uations indicate that current MLLMs still struggle with ro-
148 bust spatial reasoning [7, 30, 33, 41]. While recent efforts
149 aim to enhance spatial capabilities through scaling training
150 data [6, 14, 23] or chain-of-thought prompting [17, 52], they
151 fundamentally process visual information as static 2D snap-
152 shots. In contrast, robust spatial reasoning requires a dy-
153 namic process where the agent can selectively acquire new
154 visual evidence, similar to human mental simulation.

155 **World Models and Visual Imagination.** Recent advances
156 in video generation have demonstrated the potential of serv-
157 ing as world models, enabling agents to imagine future
158 frames or outcomes for improved decision-making [9, 15,
159 20, 28, 60]. This capability is further boosted by the emer-
160 gence of controllable video generation, which allows for
161 action-conditioned simulation [3, 11, 37, 50]. Notably, re-
162 cent works such as MindJourney [43] have pioneered the
163 use of world models to enhance visual spatial reasoning
164 by synthesizing novel viewpoints. However, their model
165 blindly generates a set number of views regardless of the
166 question’s difficulty or necessity. In contrast, we show that

the utility of visual imagination is highly query-dependent,
motivating selective use of world models at test time.

Test-Time Scaling. Test-time scaling (TTS) improves per-
formance by allocating additional inference computation
without retraining. Prior work has explored various scaling
strategies in language/visual-language models, including
self-consistency [34], tree-based search [38, 44], verifier-
guided method [24, 51], and (multimodal) CoT [35, 39, 40].
In the visual spatial reasoning, recent works [5, 43] real-
ized through generating novel views and ensembling, but
typically apply uniform computation across instances. Our
method introduces adaptive visual test-time scaling, en-
abling targeted imagination only when necessary and im-
proving computational efficiency.

3. Analysis of Always-on World Model Calling

We consider a test-time setting where an MLLM is
equipped with a visual world model that generates imag-
ined observations from hypothetical viewpoints. Existing
methods commonly invoke the world model in an always-
on, which calling it on every instance and exhaustively ex-
ploring action branches, implicitly assuming that additional
imagination is consistently beneficial. In practice, however,
this incurs substantial computational cost and may yield
ambiguous or noisy observations, making imagined views
redundant when the answer is already evident from the ini-
tial observation and outright misleading when the world
model produces noise. To diagnose this strategy, we cate-
gorize each instance on SAT [30] into three cases (Fig. 1):

- **Case 1 (Imagination Helpful):** The model calls the
world model and produces a correct answer, indicating
that imagined views provide beneficial information.
- **Case 2 (Imagination Misleading):** The model calls the
world model, but produces a wrong answer because imag-
ined views introduce misleading or noisy information.
- **Case 3 (Imagination Unnecessary):** The model pro-
duces a correct answer without calling the world model,
suggesting that visual imagination is redundant.

We then examine three aspects of always-on behavior
in Figure 2: case distribution, performance vs. amount of
imagination, and computational cost. **(1) Case distribu-
tion.** As shown in Figure 2a, the majority of instances
(54%) fall into Case 3, where the model already answers
correctly without any world model invocation, while imag-
ination is genuinely helpful in only 14% (Case 1), indicat-
ing that always-on imagination is unnecessary for most in-
stances. **(2) Performance vs. amount of imagination.** Fig-
ure 2b shows that adding more imagined views does not
consistently improve accuracy and even degrades it, sug-
gesting that simply increasing imagination is not an ideal
strategy. **(3) Cost–accuracy trade-off.** While always-on
imagination yields only a 4.6% accuracy gain over the base-

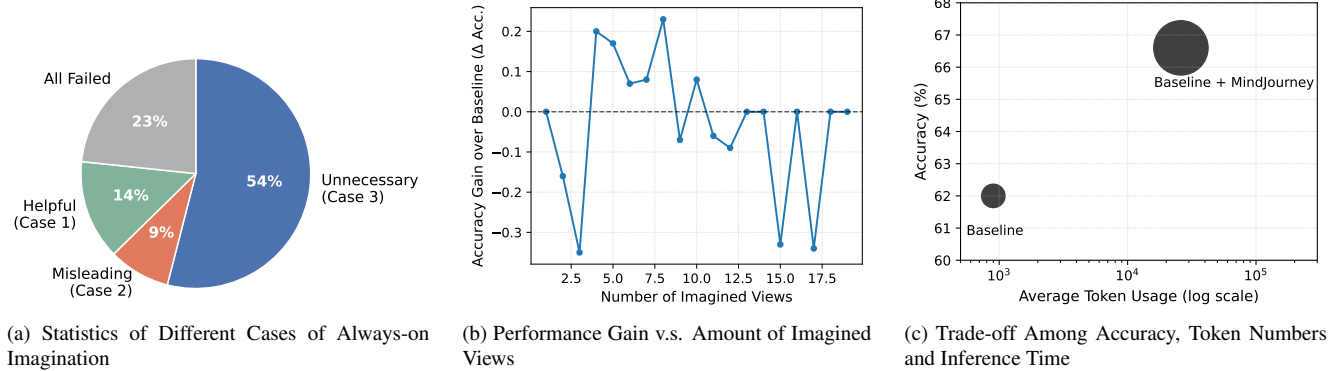


Figure 2. (a): In the majority of cases, visual imagination is unnecessary, while a smaller fraction is helpful or misleading, highlighting the need for selective invocation rather than uniform use. (b): Accuracy gain over the baseline over the number of imagined views. Performance improvements are non-monotonic, indicating that additional imagination does not consistently translate to better reasoning and may even degrade accuracy when there are too many generated views. (c): Accuracy v.s. average token usage. Bubble size indicates average running time. Fixed imagination strategies achieve higher accuracy at the cost of substantially increased computation, motivating adaptive test-time scaling that balances performance and efficiency.

line, it requires nearly two orders of magnitude more tokens and about $30\times$ higher inference time (Figure 2c), a steep computational price for limited return. (4) **Selective imagination upper bound.** We further quantify the potential of selective WM usage by assuming imagination is applied only when it leads to a correct prediction. The baseline reaches **62.0%** on SAT-Real and always-on imagination only marginally improves to **66.6%**, whereas this selective upper bound jumps to **75.3%**, demonstrating that selective imagination policies are strongly motivated. Overall, always-on WM calling is both inefficient and unreliable, motivating the need for selective, adaptive imagination.

4. Adaptive Visual Imagination Control

We now introduce **AVIC (ADAPTIVE VISUAL IMAGINATION CONTROL)**, an adaptive test-time framework that selectively invokes a world model only when additional visual evidence is likely to be useful (Fig. 3c), in contrast to the always-on baseline analyzed in Sec. 3. We first formalize the problem (Sec. 4.1), then describe the framework (Sec. 4.2), and finally describe an RL procedure that trains the policy from QA Model and WM rewards (Sec. 4.3).

4.1. Problem Formulation

We consider a visual spatial reasoning task defined by an input tuple $\langle I, q, \mathcal{A} \rangle$, where I is the current egocentric observation (one or multiple views), q is a multiple-choice question, and $\mathcal{A} = \{a_1, \dots, a_K\}$ is the answer set. The correct answer may depend on spatial relations that are ambiguous, occluded, or unobserved in I . The agent may optionally invoke a visual world model that renders novel imagined observations I_π from a sequence of egocentric actions π , with $I_\pi = \emptyset$ when no imagination is used. The predicted

answer is $\hat{a} = \arg \max_{a \in \mathcal{A}} P(a | I, I_\pi, q)$.

4.2. Adaptive Imagination Framework

AVIC couples a **policy model** that gates and plans imagination with a **trajectory-level verifier** that picks a single targeted imagined trajectory for downstream reasoning.

Policy gating with test-time scaling. A policy θ maps (I, q, \mathcal{A}) to a decision $d \in \{\text{skip}, \text{call_wm}\}$ together with a short discrete action plan π drawn from a fixed low-level egocentric action space \mathcal{U} :

$$(d, \pi) \sim \theta(d, \pi | I, q, \mathcal{A}),$$

$$\pi = \begin{cases} \emptyset, & d = \text{skip}, \\ (u_1, \dots, u_T), u_t \in \mathcal{U}, & d = \text{call_wm}. \end{cases} \quad (1)$$

To improve robustness, we sample the policy M times under independent decoding and aggregate d by majority voting, providing a simple form of self-consistency that reflects uncertainty in the necessity of imagination.

Action execution and trajectory selection. When $d = \text{call_wm}$, each sampled plan $\pi^{(m)}$ is executed by the world model W to render an imagined trajectory $\mathcal{I}_{\pi^{(m)}} = W(I, \pi^{(m)})$. Different policy samples produce trajectories of varying usefulness; unlike prior beam-search approaches [43] that score isolated keyframes, we evaluate the *entire* trajectory as a coherent unit via a verifier V , preserving temporal and geometric consistency across sequential actions. Such that:

$$s^{(m)} = V(I, q, \mathcal{I}_{\pi^{(m)}}), \quad \pi^* = \arg \max_{\pi^{(m)}} s^{(m)}, \quad (2)$$

Final prediction. A vision-language reasoner ϕ predicts the answer using the original observation and the selected

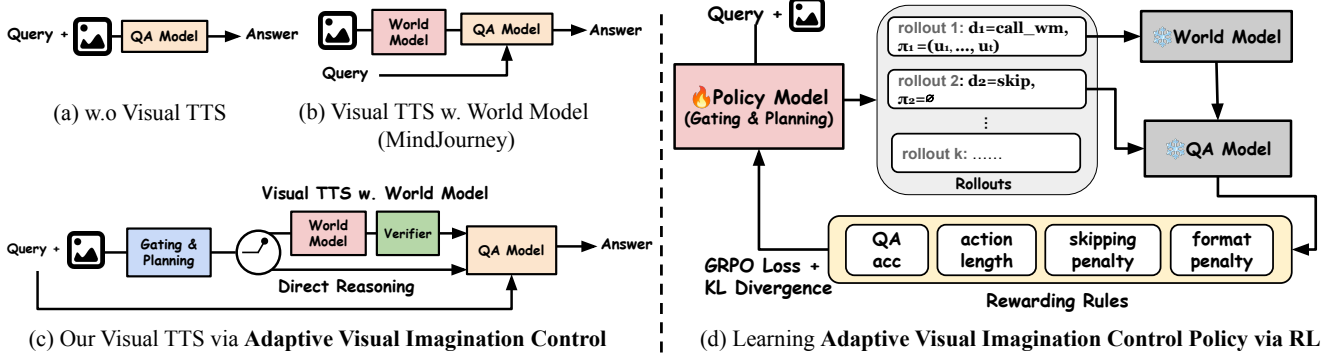


Figure 3. (a) Direct QA from the current observation. (b) Always-on world-model exploration. (c) Ours: a policy model decides *whether* and *how much* to imagine, selectively querying the world model only when warranted. (d) RL training loop for the gating policy: The policy θ samples K rollouts that either query the world model W (`call_wm`) or bypass it (`skip`). A frozen QA model ϕ answers them, and a four-component reward drives a GRPO update on θ .

275 imagined views (or I alone when the gate selected `skip`,
276 in which case $I_{\pi^*} = \emptyset$):

$$277 \quad \hat{a} = \arg \max_{a \in \mathcal{A}} P_{\phi}(a | I, I_{\pi^*}, q). \quad (3)$$

278 4.3. Learning When and How Much to Imagine via 279 RL

280 Learning when and how much to imagine is non-trivial,
281 as no ground-truth supervision exists for what the optimal
282 imagination trajectory should be, ruling out supervised fine-
283 tuning or behavior cloning. We therefore propose AVIC-
284 R, which trains θ end-to-end with reinforcement learning,
285 using QA correctness as the primary signal augmented by
286 lightweight reward shaping (Figure 3 (d)).

287 **Training loop.** For each question, the policy θ samples K
288 rollouts $\{(d^{(i)}, \pi^{(i)})\}$. Each rollout is executed by frozen
289 environment modules: `call_wm` rollouts query the world
290 model W to render imagined views, while `skip` rollouts
291 bypass it; both are answered by a frozen QA model ϕ . The
292 resulting reward signals are aggregated into a GRPO update
293 on the policy, while W and ϕ remain frozen during training.

294 **Reward design.** We assemble a composite reward with 4
295 components, as illustrated in Figure 3 (d):

$$296 \quad r = \underbrace{\mathbf{1}_{\text{correct}}}_{\text{(i) QA correctness}} - \underbrace{c|\pi|}_{\text{(ii) action length}} \\ - \underbrace{\beta_s \mathbf{1}_{\text{wrong-skip}}}_{\text{(iii) wrong-skip}} - \underbrace{\beta_p \mathbf{1}_{\text{parse-fail}}}_{\text{(iv) format}} \quad (4)$$

297 where $\mathbf{1}_{\text{correct}}$, $\mathbf{1}_{\text{wrong-skip}}$, and $\mathbf{1}_{\text{parse-fail}}$ indicate, respectively,
298 that $\hat{a} = a^*$, that $d=\text{skip} \wedge \hat{a} \neq a^*$, and that the output
299 schema cannot be parsed; $|\pi|$ counts atomic actions; and
300 $c, \beta_s, \beta_p > 0$. Each term targets a specific failure mode:
301 **(1) QA correctness** is the only positive signal; all other
302 terms are deductions. **(2) Action-length cost** discourages

303 over-imagination, a longer correct chain still scores below
304 a shorter correct one, pushing the policy toward concise
305 plans. **(3) Wrong-skip penalty** is essential. Without it, an
306 incorrect `skip` pays nothing while an incorrect `call_wm`
307 still pays $c|\pi|$, biasing the policy toward skipping. **(4) For-**
308 **mat penalty** handles unrecoverable schema errors. We del-
309 iberately keep $\beta_p \approx \beta_s$ rather than larger: a heavier format
310 penalty would incentivize the policy to collapse onto trivial
311 outputs (e.g., always `skip` with empty π) merely to avoid
312 parse failures. We set $c = 0.1$ and $\beta_s = \beta_p = 0.5$, yielding
313 a clean qualitative ordering: a correct `skip` scores $+1.0$,
314 a correct `call_wm` with n actions scores $1 - 0.1n$, and
315 a wrong `skip` (-0.5) is strictly worse than a short wrong
316 `call_wm` (-0.1 to -0.3). The goal of this asymmetry is
317 to prevent the policy from collapsing to *always-skip*: when
318 in doubt, calling WM with a short plan is the safer bet than
319 skipping. The action cost $c|\pi|$ then handles the opposite
320 failure mode, preventing collapse onto over imagination.

321 **GRPO objective.** Without an external value function, we
322 adopt Group-Relative Policy Optimization [32], which esti-
323 mates advantages from the K rollouts sharing a prompt:

$$A_i = \frac{r_i - \mu_q}{\sigma_q + \epsilon}, \quad \mu_q = \frac{1}{K} \sum_j r_j, \quad \sigma_q^2 = \frac{1}{K} \sum_j (r_j - \mu_q)^2. \quad (5)$$

324 Group normalization absorbs per-question difficulty, where
325 uniformly easy or uniformly hard groups contribute zero
326 gradient. We optimize a token-level PPO-clipped objec-
327 tive [31] regularized by KL to a frozen reference θ_{ref} :
328

$$\mathcal{L} = -\mathbb{E}_{(q,i,t)} [\min(\rho_{i,t} A_i, \text{clip}(\rho_{i,t}, 1-\epsilon_c, 1+\epsilon_c) A_i)] \\ + \beta_{\text{KL}} \widehat{\text{KL}}(\theta \| \theta_{\text{ref}}), \quad (6)$$

329 with per-token importance ratio $\rho_{i,t} = \theta(y_{i,t} |$
330 $q, y_{i,<t}) / \theta_{\text{old}}(y_{i,t} | q, y_{i,<t})$. Only LoRA adapters in θ are
331

| Method | Policy Model | EgoM | ObjM | EgoAct | Goal | Pers | Avg. | # Token (K) | Avg. WM |
|--------------------|---------------|--------------|-------------|-------------|-------------|-------------|-------------|-------------|---------|
| InternVL3-14B [61] | – | 56.5 | <u>69.5</u> | 54.0 | 73.5 | <u>45.4</u> | 59.3 | 0.2 | 0 |
| + MindJourney | – | 69.6 | 60.9 | 78.4 | <u>79.4</u> | 42.4 | 66.7 | 2.5 | 12.34 |
| + AVIC | InternVL3-14B | 95.6 | 73.9 | 62.1 | 76.4 | 42.4 | <u>68.0</u> | 2.0 | 0.64 |
| + AVIC | Qwen2.5VL-7B | 73.9 | 47.8 | 67.5 | 73.5 | 42.4 | 61.3 | 4.4 | 1.81 |
| + AVIC-R | Qwen2.5VL-7B | <u>82.6</u> | 52.1 | <u>70.2</u> | 85.2 | 54.5 | 69.3 | 4.8 | 3.03 |
| GPT-4o [26] | – | 56.5 | 85.0 | 50.0 | 64.0 | 45.0 | 60.3 | 0.9 | 0 |
| + MindJourney | – | 78.3 | 60.9 | <u>78.4</u> | 70.6 | <u>57.5</u> | 69.3 | 26.0 | 12.34 |
| + AVIC | GPT-4o | 86.9 | 60.9 | 64.8 | <u>82.3</u> | 48.4 | 69.3 | 9.5 | 0.72 |
| + AVIC | Qwen2.5VL-7B | 65.2 | 73.9 | 64.8 | 91.1 | 60.6 | <u>71.3</u> | 5.0 | 1.81 |
| + AVIC-R | Qwen2.5VL-7B | <u>82.6</u> | <u>82.6</u> | 81.0 | 91.1 | 51.2 | 77.3 | 5.4 | 3.03 |
| GPT-4.1 [27] | – | <u>95.7</u> | 73.9 | 78.3 | 88.2 | 39.4 | 74.0 | 0.7 | 0 |
| + MindJourney | – | 100.0 | <u>82.6</u> | 86.5 | 79.4 | 45.4 | 77.3 | 67.1 | 12.34 |
| + AVIC | GPT-4.1 | 100.0 | 78.2 | <u>83.7</u> | <u>85.2</u> | <u>54.5</u> | <u>79.3</u> | 7.6 | 0.73 |
| + AVIC | Qwen2.5VL-7B | 82.6 | 86.9 | 75.6 | 88.2 | 36.3 | 72.6 | 4.8 | 1.81 |
| + AVIC-R | Qwen2.5VL-7B | 91.3 | 86.9 | <u>83.7</u> | <u>85.2</u> | 57.5 | 80.0 | 5.2 | 3.03 |
| o1 [16] | – | 78.3 | <u>82.6</u> | 73.0 | 73.5 | 69.7 | 74.6 | 1.4 | 0 |
| + MindJourney | – | 100.0 | 65.2 | 78.4 | 82.4 | 63.7 | 77.3 | 39.4 | 12.34 |
| + AVIC | o1 | 100.0 | 86.9 | 86.4 | <u>91.1</u> | 66.6 | 85.3 | 14.6 | 1.28 |
| + AVIC | Qwen2.5VL-7B | <u>86.9</u> | 65.2 | 78.3 | 94.1 | <u>69.6</u> | 79.3 | 5.7 | 1.81 |
| + AVIC-R | Qwen2.5VL-7B | <u>86.9</u> | 65.2 | <u>81.0</u> | 94.1 | <u>69.6</u> | <u>80.0</u> | 6.1 | 3.03 |

Table 1. Comparison between TTS methods on SAT-Real. The best results are denoted by **bold**, and the second-best are underlined. **Avg. WM**: average world model calling times over the dataset.

| Method | GPT-4o [26] | | GPT-4.1 [27] | |
|----------|-------------|-------------|--------------|-------------|
| | Base | + AVIC | Base | + AVIC |
| Accuracy | 30.3 | 32.3 | 30.9 | 33.8 |

Table 2. Results on MMSI.

| Methods | LLMs | NE↓ | OSR↑ | SR↑ | SPL↑ |
|---------------|---------|-------------|-------------|-------------|-------------|
| NavGPT [58] | GPT-3.5 | 8.02 | 26.4 | 16.7 | 13.0 |
| MapGPT [55] | GPT-4 | 5.80 | 61.6 | 41.2 | 25.4 |
| MapGPT | GPT-4o | 6.04 | 41.6 | 36.0 | 30.8 |
| MapGPT + AVIC | GPT-4o | 5.97 | 45.3 | 37.5 | 31.9 |

Table 3. Results on R2R embodied navigation dataset.

updated; the KL anchor preserves the base VLM’s general capabilities while shaping only gating and planning behaviors. We provide further details on training, data curation and output parsing in later Sec. 5.1 and Appendix.

5. Experiments

5.1. Experiment Setup

Datasets and Benchmarks. We validate our proposed framework on both visual spatial reasoning benchmarks and embodied navigation tasks, covering a range of spatial ambiguities and interaction requirements. For visual spatial reasoning, we evaluate on SAT [30] and MMSI [42],

two benchmarks for visual spatial reasoning with single/multiple images. We also evaluate on the Room-to-Room (R2R) [2] for the embodied navigation task.

Implementation Details. Our framework is implemented on top of a vision-language model and a pretrained visual world model, stable virtual camera (SVC) [59]. The policy model, verifier, and final QA model are instantiated using the same base MLLM in AVIC, with different prompting. All decisions in AVIC are made at test time without additional fine-tuning. We scale action planning by 5 times as the default. We adapt LoRA [13] finetuning for AVIC-R, with 8 LoRA rank, 16 LoRA alpha (more in the Appendix).

5.2. Main Results

AVIC beats the always-on baseline in both efficiency and effectiveness across all backbones. Table 1 compares our framework against baselines on SAT-Real [30] across five categories. Across all open-source and proprietary backbones, AVIC consistently improves over the base MLLM and matches or surpasses the always-on baseline (MindJourney) while using $\sim 10\%$ of the tokens and far fewer world-model calls. With GPT-4.1, accuracy rises from 74.0% to 79.3%; with o1, it reaches **85.3%** (+10.7% over base). Gains are most pronounced on Egocentric Movement, Action Consequence, and Perspective tasks, categories requiring action-conditioned spatial reasoning where selective imagination is most beneficial. Table 2 confirms these gains transfer to MMSI-Bench [42].

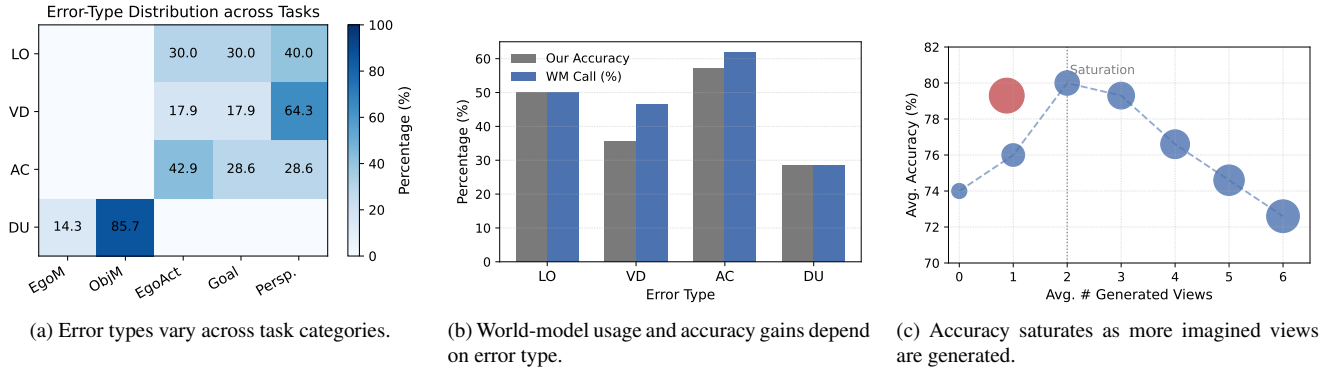


Figure 4. Analysis of when and how much to invoke world-model imagination.

370 **Lightweight RL turns a small policy into one that beats**
 371 **proprietary alternatives.** AVIC-R adds two further find-
 372 ings on top of AVIC. First, *a small RL-trained policy can*
 373 *drive a much larger backbone:* with only Qwen2.5VL-
 374 7B as the gating policy, AVIC-R outperforms AVIC vari-
 375 ants that use the proprietary backbone itself as the pol-
 376 icy on three of four backbones. On GPT-4o, AVIC-
 377 R reaches **77.3%**, beating AVIC with GPT-4o-as-policy
 378 by **8.0 points**; on GPT-4.1 it reaches **80.0%**, surpassing
 379 GPT-4.1-as-policy. Second, *RL is essential to making the*
 380 *small policy work:* prompting Qwen2.5VL-7B alone as the
 381 AVIC policy often underperforms even the always-on base-
 382 line (e.g., 61.3% vs. MindJourney’s 66.7% on InternVL3-
 383 14B; 72.6% vs. 77.3% on GPT-4.1), as the 7B model lacks
 384 the in-context reasoning to gate reliably. Our lightweight
 385 RL training fixes this, lifting the same Qwen2.5VL-7B by
 386 6–8 points across InternVL3-14B/GPT-4o/GPT-4.1 back-
 387 bones and turning a policy that previously *underperformed*
 388 the always-on baseline into one that drives the full pipeline
 389 to top results. More ablations about AVIC modules, rewards
 390 design and runtime are in the Appendix.

391 **AVIC’s selective imagination also transfers to embod-**
 392 **ied navigation.** Table 3 applies AVIC to embodied naviga-
 393 tion, integrated into MapGPT’s [55] step-wise framework
 394 on the 72-scene R2R [2] evaluation. At each step, our pol-
 395 icy model decides whether to invoke the world model on a
 396 subset of graph views; the imagined views are concatenated
 397 with original observations to inform the next-action predic-
 398 tion. Compared to MapGPT with GPT-4o, AVIC achieves
 399 higher OSR/SR/SPL and lower navigation error (NE), indi-
 400 cating more reliable goal reaching with shorter, less redun-
 401 dant trajectories, transferring the benefits of selective imag-
 402 ination from static spatial reasoning to embodied tasks.

403 5.3. When and How Much a World Model is Needed 404 for Visual Spatial Reasoning?

405 To diagnose *when* world-model imagination is necessary,
 406 we manually classify MLLM failures on SAT-Real into

four error types: **(1) Limited Observability (LO, 15.2%):**
 required information occluded or out of view; **(2) View-**
point Dependence (VD, 42.4%): answer depends on trans-
 forming between egocentric and object-centric frames; **(3)**
Action-Conditioned Reasoning (AC, 31.8%): answer de-
 pends on the scene state after a hypothetical action; and
Dynamics Understanding ((4) DU, 10.6%): temporal rea-
 soning about camera or object motion. Figure 4a shows
 that SAT task categories do not map one-to-one with error
 types but exhibit *compositional* patterns: EgoAct is domi-
 nated by AC errors (post-action viewpoints), Pers. by VD
 errors (reference-frame transformation), and ObjM by DU
 errors (temporal dynamics). LO errors appear across mul-
 tiple categories, indicating that occlusion and limited field-
 of-view are general failure sources. This decomposition lets
 us study WM utility through error structure rather than sur-
 face task labels.

RQ1: When to call WM?

WM should be used selectively, primarily when rea-
 soning requires predicting future states under hypo-
 theoretical actions rather than reinterpreting exist-
 ing visual evidence.

World models are most needed for action-conditioned
reasoning. Figure 4b shows that WM imagination yields

the largest gain (**+57.1%**) on AC errors, where the answer
 depends on the post-action scene state, e.g., counterfactual
 queries like “what if I turn left by 90°?”. By contrast, DU
 errors require only reference-frame transformation over the
 current view and benefit much less (+28.5%). LO and VD
 errors fall between these extremes: rendering can reveal oc-
 cluded content (LO) or visualize the scene from a differ-
 ent vantage (VD), but adds little when the transformation
 can be inferred symbolically from the original observation.
 These findings indicate that WM utility is highly instance-
 dependent. It is most useful when the question references a
 future or counterfactual scene, often unnecessary for static

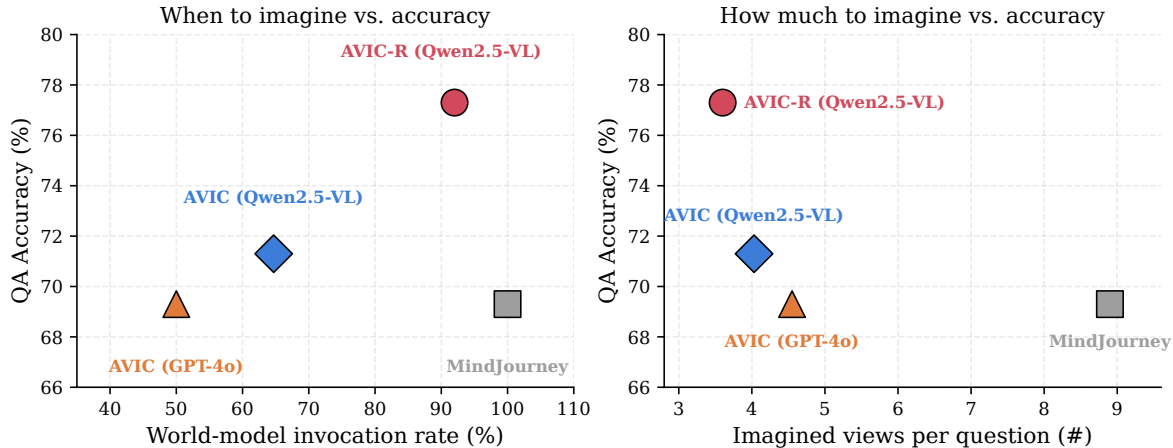


Figure 5. **Adaptive imagination achieves higher accuracy at lower cost on SAT-Real.** *Left:* world-model invocation rate (*when to imagine*) vs. QA accuracy. *Right:* imagined views per question (*how much to imagine*) vs. QA accuracy. AVIC-R achieves the highest accuracy (77.3%) using the fewest imagined views (3.60 vs. 8.90 for always-on).

439 reinterpretation of what is already observed.

RQ2: How much imagination is needed?

Visual spatial reasoning benefits from *targeted* rather than extensive WM imagination.

440

441 **Spatial reasoning requires limited imagination.** Figure
442 4c sweeps fixed-budget baselines with a predetermined
443 number of imagined views. Even a single targeted view
444 raises accuracy by roughly 4 points over the no-imagination
445 baseline (74.0% \rightarrow 76%), and a second view captures most
446 of the remaining headroom (76% \rightarrow 80%). Beyond two
447 views, additional rollouts bring no further gains and even-
448 tually degrade performance, as accumulated rendering ar-
449 tifacts and redundant content begin to confuse the down-
450 stream reasoner. Targeted, low-budget imagination is what
451 spatial reasoning actually needs, while exhaustive scaling is
452 wasteful at best and harmful at worst.

453 **AVIC-R learns better when and how much to imagine.**
454 The findings above point to a simple rule for visual test time
455 scaling with WM: we should call WM for visual spatial
456 reasoning mainly on action-conditioned questions, and use
457 only targeted views per call. Figure 5 (and Appendix Ta-
458 ble 8) shows that AVIC-R learns both parts, while alterna-
459 tive policies do not. On the *when* axis, AVIC-R calls WM
460 on 92% of questions on average, with strong differences
461 across categories: 100% on EgoAct (dominated by AC er-
462 rors) versus 78.8% on Pers.(dominated by VD errors). This
463 category-aware behavior emerges from QA-correctness and
464 WM-cost rewards alone, and no per-category labels are pro-
465 vided. Other alternatives are far less stable: AVIC with
466 GPT-4o calls 100% on EgoM but 0% on ObjM, and zero-

shot Qwen2.5VL over-skips at 64.7%, ending 6 points be-
hind AVIC-R (71.3% vs. 77.3%). Calling WM more often
is not the point, but calling it on the right questions is. On
the *how-much* axis, AVIC-R uses 3.60 views per question,
fewer than every selective baseline (zero-shot Qwen 4.03,
GPT-4o 4.55) and less than half of MindJourney’s 8.90,
while still reaching the best accuracy. This sits right at the
saturation point of Figure 4c: enough imagination to extract
the accuracy gain, none of the excess that begins to hurt
performance beyond 2–3 views. In short, both behaviors
emerge from our lightweight RL scheme: the policy learns
when to imagine and how much to imagine, just from QA
correctness and a WM cost penalty.

6. Conclusion

In this paper, we study visual spatial reasoning with world
models through the lens of adaptive test-time scaling, find-
ing that always-on imagination is often unnecessary and
even misleading. We introduce AVIC, which selectively
decides *when* and *how much* to imagine at inference time,
and AVIC-R, which trains this gating policy end-to-end via
lightweight RL from QA-correctness and WM-cost signals.
Across spatial reasoning and embodied navigation bench-
marks, our framework achieves competitive or state-of-the-
art results while substantially reducing world-model calls,
tokens, and inference time; notably, AVIC-R with a small
open-source policy outperforms pipelines that use propri-
etary backbones as the policy. Our analysis shows world-
model imagination is most beneficial for action-conditioned
reasoning and requires only limited, targeted views, high-
lighting the importance of instance-dependent TTS for effi-
cient and reliable reasoning with world models.

498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554**References**

- [1] Arslan Ali, Junjie Bai, Maciej Bala, Yogesh Balaji, Aaron Blakeman, Tiffany Cai, Jiaxin Cao, Tianshi Cao, Elizabeth Cha, Yu-Wei Chao, et al. World simulation with video foundation models for physical ai. *arXiv preprint arXiv:2511.00062*, 2025.
- [2] Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen Gould, and Anton Van Den Hengel. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3674–3683, 2018.
- [3] Sherwin Bahmani, Ivan Skorokhodov, Guocheng Qian, Aliaksandr Siarohin, Willi Menapace, Andrea Tagliasacchi, David B. Lindell, and Sergey Tulyakov. Ac3d: Analyzing and improving 3d camera control in video diffusion transformers. *2025 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 22875–22889, 2024.
- [4] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2.5-vl technical report, 2025.
- [5] Meng Cao, Xingyu Li, Xue Liu, Ian Reid, and Xiaodan Liang. Spatialdreamer: Incentivizing spatial reasoning via active mental imagery. *arXiv preprint arXiv:2512.07733*, 2025.
- [6] Boyuan Chen, Zhuo Xu, Sean Kirmani, Brian Ichter, Danny Driess, Pete Florence, Dorsa Sadigh, Leonidas J. Guibas, and Fei Xia. Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14455–14465, 2024.
- [7] 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 model. *ArXiv*, abs/2406.01584, 2024.
- [8] Andong Deng, Tongjia Chen, Shoubin Yu, Taojiannan Yang, Lincoln Spencer, Yapeng Tian, Ajmal Saeed Mian, Mohit Bansal, and Chen Chen. Motion-grounded video reasoning: Understanding and perceiving motion at pixel level. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 8625–8636, 2025.
- [9] Yilun Du, Mengjiao Yang, Bo Dai, Hanjun Dai, Ofir Nachum, Joshua B. Tenenbaum, Dale Schuurmans, and P. Abbeel. Learning universal policies via text-guided video generation. *ArXiv*, abs/2302.00111, 2023.
- [10] Xiao Guo, Xiufeng Song, Yue Zhang, Xiaohong Liu, and Xiaoming Liu. Rethinking vision-language model in face forensics: Multi-modal interpretable forged face detector. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 105–116, 2025.
- [11] Hao He, Yinghao Xu, Yuwei Guo, Gordon Wetzstein, Bo Dai, Hongsheng Li, and Ceyuan Yang. Cameractrl: Enabling camera control for text-to-video generation. *ArXiv*, abs/2404.02101, 2024.
- [12] Yining Hong, Haoyu Zhen, Peihao Chen, Shuhong Zheng, Yilun Du, Zhenfang Chen, and Chuang Gan. 3d-llm: Injecting the 3d world into large language models. *ArXiv*, abs/2307.12981, 2023.
- [13] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Liang Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *Iclr*, 1(2):3, 2022.
- [14] Jiangyong Huang, Silong Yong, Xiaojuan Ma, Xiongkun Linghu, Puhao Li, Yan Wang, Qing Li, Song-Chun Zhu, Baoxiong Jia, and Siyuan Huang. An embodied generalist agent in 3d world. *ArXiv*, abs/2311.12871, 2023.
- [15] Yidong Huang, Zun Wang, Han Lin, Dong-Ki Kim, Shayegan Omidshafiei, Jaehong Yoon, Yue Zhang, and Mohit Bansal. Planning with sketch-guided verification for physics-aware video generation. *arXiv preprint arXiv:2511.17450*, 2025.
- [16] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. OpenAI o1 system card, 2024.
- [17] Binbin Ji, Siddharth Agrawal, Qiance Tang, and Yvonne Wu. Enhancing spatial reasoning in vision-language models via chain-of-thought prompting and reinforcement learning. *ArXiv*, abs/2507.13362, 2025.
- [18] Stephen M Kosslyn, William L Thompson, and Giorgio Ganis. *The case for mental imagery*. Oxford University Press, 2006.
- [19] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Zhang, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. LLaVA-onevision: Easy visual task transfer, 2024.
- [20] Jialu Li and Mohit Bansal. Panogen: Text-conditioned panoramic environment generation for vision-and-language navigation. *ArXiv*, abs/2305.19195, 2023.
- [21] Junnan Li, Dongxu Li, Caiming Xiong, and Steven C. H. Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, 2022.
- [22] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023.
- [23] Jialu Li, Jaemin Cho, Yi-lin Sung, Jaehong Yoon, and Mohit Bansal. Selma: Learning and merging skill-specific text-to-image experts with auto-generated data. In *Neural Information Processing Systems*, 2024.
- [24] Shalev Lifshitz, Sheila A McIlraith, and Yilun Du. Multi-agent verification: Scaling test-time compute with multiple verifiers. *arXiv preprint arXiv:2502.20379*, 2025.
- [25] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *ArXiv*, abs/2304.08485, 2023.
- [26] OpenAI. Hello GPT-4o. OpenAI Blog, 2024.
- [27] OpenAI. GPT-4.1 technical overview, 2024.
- [28] Cheng Qian, Emre Can Acikgoz, Bingxuan Li, Xiusi Chen, Yuji Zhang, Bingxiang He, Qinyu Luo, Dilek Hakkani-Tür, Gokhan Tur, Yunzhu Li, et al. Current agents fail to leverage world model as tool for foresight. *arXiv preprint arXiv:2601.03905*, 2026.

- 613 [29] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya
614 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,
615 Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen
616 Krueger, and Ilya Sutskever. Learning transferable visual
617 models from natural language supervision. In *International
618 Conference on Machine Learning*, 2021.
- 619 [30] Arijit Ray, Jiafei Duan, Ellis Brown, Reuben Tan, Dina
620 Bashkirova, Rose Hendrix, Kiana Ehsani, Aniruddha Kem-
621 bhavi, Bryan A. Plummer, Ranjay Krishna, et al. SAT:
622 Dynamic spatial aptitude training for multimodal language
623 models, 2024.
- 624 [31] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Rad-
625 ford, and Oleg Klimov. Proximal policy optimization algo-
626 rithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 627 [32] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao
628 Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li,
629 Yang Wu, et al. Deepseekmath: Pushing the limits of math-
630 ematical reasoning in open language models. *arXiv preprint
631 arXiv:2402.03300*, 2024.
- 632 [33] Shengbang Tong, Ellis Brown, Penghao Wu, Sanghyun
633 Woo, Manoj Middepogu, Sai Charitha Akula, Jihan Yang,
634 Shusheng Yang, Adithya Iyer, Xichen Pan, Austin Wang,
635 Rob Fergus, Yann LeCun, and Saining Xie. Cambrian-1:
636 A fully open, vision-centric exploration of multimodal llms.
637 *ArXiv*, abs/2406.16860, 2024.
- 638 [34] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed
639 Chi, Sharan Narang, Aakanksha Chowdhery, and Denny
640 Zhou. Self-consistency improves chain of thought reason-
641 ing in language models. *arXiv preprint arXiv:2203.11171*,
642 2022.
- 643 [35] Yaoting Wang, Shengqiong Wu, Yuecheng Zhang,
644 Shuicheng Yan, Ziwei Liu, Jiebo Luo, and Hao Fei.
645 Multimodal chain-of-thought reasoning: A comprehensive
646 survey. *arXiv preprint arXiv:2503.12605*, 2025.
- 647 [36] Zun Wang, Jialu Li, Yicong Hong, Yi Wang, Qi Wu, Mo-
648 hit Bansal, Stephen Gould, Hao Tan, and Yu Qiao. Scal-
649 ing data generation in vision-and-language navigation. *2023
650 IEEE/CVF International Conference on Computer Vision
651 (ICCV)*, pages 11975–11986, 2023.
- 652 [37] Zun Wang, Jaemin Cho, Jialu Li, Han Lin, Jaehong Yoon,
653 Yue Zhang, and Mohit Bansal. Epic: Efficient video camera
654 control learning with precise anchor-video guidance. *ArXiv*,
655 abs/2505.21876, 2025.
- 656 [38] Ziyang Wang, Jaehong Yoon, Shoubin Yu, Mohaiminul Is-
657 lam, Gedas Bertasius, and Mohit Bansal. Video-rtts: Rethink-
658 ing reinforcement learning and test-time scaling for efficient
659 and enhanced video reasoning. In *Conference on Empirical
660 Methods in Natural Language Processing*, 2025.
- 661 [39] Yige Xu, Xu Guo, Zhiwei Zeng, and Chunyan Miao. Soft-
662 cot++: Test-time scaling with soft chain-of-thought reason-
663 ing. *arXiv preprint arXiv:2505.11484*, 2025.
- 664 [40] Ziang Yan, Xinhao Li, Yinan He, Zhengrong Yue, Xiangyu
665 Zeng, Yali Wang, Yu Qiao, Limin Wang, and Yi Wang.
666 Videochat-r1. 5: Visual test-time scaling to reinforce mul-
667 timodal reasoning by iterative perception. *arXiv preprint
668 arXiv:2509.21100*, 2025.
- 669 [41] Jihan Yang, Shusheng Yang, Anjali W. Gupta, Rilyn Han,
670 Fei-Fei Li, and Saining Xie. Thinking in space: How mul-
timodal large language models see, remember, and recall
spaces. *2025 IEEE/CVF Conference on Computer Vision and
Pattern Recognition (CVPR)*, pages 10632–10643, 2024.
- [42] Sihan Yang, Runsen Xu, Yiman Xie, Sizhe Yang, Mo Li,
Jingli Lin, Chenming Zhu, Xiaochen Chen, Haodong Duan,
Xiangyu Yue, et al. Mmsi-bench: A benchmark for multi-
image spatial intelligence. *arXiv preprint arXiv:2505.23764*,
2025.
- [43] Yuncong Yang, Jiageng Liu, Zheyuan Zhang, Siyuan Zhou,
Reuben Tan, Jianwei Yang, Yilun Du, and Chuang Gan.
Mindjourney: Test-time scaling with world models for spa-
tial reasoning, 2025.
- [44] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom
Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of
thoughts: Deliberate problem solving with large language
models. *Advances in Neural Information Processing Sys-
tems*, 36:11809–11822, 2023.
- [45] Baiqiao Yin, Qineng Wang, Pingyue Zhang, Jianshu Zhang,
Kangrui Wang, Zihan Wang, Jieyu Zhang, Keshigeyan Chan-
drasegaran, Han Liu, Ranjay Krishna, et al. Spatial mental
modeling from limited views, 2025. Structural Priors for Vi-
sion Workshop at ICCV 2025.
- [46] Jaehong Yoon, Shoubin Yu, and Mohit Bansal. Raccoon:
A versatile instructional video editing framework with auto-
generated narratives. *arXiv preprint arXiv:2405.18406*,
2024.
- [47] Shoubin Yu, Difan Liu, Ziqiao Ma, Yicong Hong, Yang
Zhou, Hao Tan, Joyce Chai, and Mohit Bansal. Veg-
gie: Instructional editing and reasoning video concepts with
grounded generation. In *Proceedings of the IEEE/CVF In-
ternational Conference on Computer Vision*, pages 15147–
15158, 2025.
- [48] Shoubin Yu, Jaehong Yoon, and Mohit Bansal. Crema: Gen-
eralizable and efficient video-language reasoning via mul-
timodal modular fusion. In *International Conference on
Learning Representations*, 2025.
- [49] Shoubin Yu, Yue Zhang, Ziyang Wang, Jaehong Yoon, and
Mohit Bansal. Mexa: Towards general multimodal reasoning
with dynamic multi-expert aggregation. *Findings of the 2025
Conference on Empirical Methods in Natural Language Pro-
cessing*, 2025.
- [50] Wangbo Yu, Jinbo Xing, Li Yuan, Wenbo Hu, Xiaoyu Li,
Zhipeng Huang, Xiangjun Gao, Tien-Tsin Wong, Ying Shan,
and Yonghong Tian. Viewcrafter: Taming video diffusion
models for high-fidelity novel view synthesis. *IEEE transac-
tions on pattern analysis and machine intelligence*, PP, 2024.
- [51] Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran
Kazemi, Aviral Kumar, and Rishabh Agarwal. Generative
verifiers: Reward modeling as next-token prediction. *arXiv
preprint arXiv:2408.15240*, 2024.
- [52] Ruohong Zhang, Bowen Zhang, Yanghao Li, Haotian Zhang,
Zhiqing Sun, Zhe Gan, Yinfei Yang, Ruoming Pang, and
Yiming Yang. Improve vision language model chain-of-
thought reasoning. In *Annual Meeting of the Association for
Computational Linguistics*, 2024.
- [53] Yue Zhang and Parisa Kordjamshidi. Vln-trans: Translator
for the vision and language navigation agent. In *The 61st*

- 728 *Annual Meeting Of The Association For Computational Lin-*
729 *guistics*, 2023.
- 730 [54] Yue Zhang, Ben Colman, Xiao Guo, Ali Shahriyari, and
731 Gaurav Bharaj. Common sense reasoning for deepfake de-
732 tection. In *European conference on computer vision*, pages
733 399–415. Springer, 2024.
- 734 [55] Yifan Zhang, Zhengting He, Jingxuan Li, Jianfeng Lin,
735 Qingfeng Guan, and Wenhao Yu. Mapgpt: an autonomous
736 framework for mapping by integrating large language model
737 and cartographic tools. *Cartography and Geographic Infor-*
738 *mation Science*, 51(6):717–743, 2024.
- 739 [56] Yue Zhang, Ziqiao Ma, Jialu Li, Yanyuan Qiao, Zun Wang,
740 Joyce Chai, Qi Wu, Mohit Bansal, and Parisa Kordjamshidi.
741 Vision-and-language navigation today and tomorrow: A sur-
742 vey in the era of foundation models. *Trans. Mach. Learn.*
743 *Res.*, 2024, 2024.
- 744 [57] Yue Zhang, Zhiyang Xu, Ying Shen, Parisa Kordjamshidi,
745 and Lifu Huang. Spartun3d: Situated spatial understand-
746 ing of 3d world in large language models. *arXiv preprint*
747 *arXiv:2410.03878*, 2024.
- 748 [58] Gengze Zhou, Yicong Hong, and Qi Wu. Navgpt: Explicit
749 reasoning in vision-and-language navigation with large lan-
750 guage models. In *Proceedings of the AAAI Conference on*
751 *Artificial Intelligence*, pages 7641–7649, 2024.
- 752 [59] Jensen Zhou, Hang Gao, Vikram Voleti, Aaryaman Vasishta,
753 Chun-Han Yao, Mark Boss, Philip Torr, Christian Rupprecht,
754 and Varun Jampani. Stable virtual camera: Generative view
755 synthesis with diffusion models, 2025.
- 756 [60] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe,
757 and Noah Snavely. Stereo magnification. *ACM Transactions*
758 *on Graphics (TOG)*, 37:1 – 12, 2018.
- 759 [61] Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shen-
760 glong Ye, Lixin Gu, Yuchen Duan, Hao Tian, Weijie Su, Jie
761 Shao, et al. Internv13: Exploring advanced training and test-
762 time recipes for open-source multimodal models, 2025.

763 A. Implementation Details

764 **AVIC-R Training Details.** We post-train Qwen2.5-VL-
765 7B-Instruct with LoRA [13] ($r=8$, $\alpha=16$, dropout 0.05, ap-
766 plied to the $q/k/v/o$ projections; $\sim 5M$ trainable param-
767 eters out of 8.3B) using online GRPO [32]. For each question
768 we sample $K=16$ rollouts from the policy with tempera-
769 ture 1.0, top- $p=0.95$, top- $k=50$, and a 512-token response
770 budget; advantages are computed by per-question (group)
771 reward normalization. The PPO surrogate uses a clip ra-
772 tio of $\epsilon=0.2$ and a K3 KL penalty to the frozen base pol-
773 icy with $\beta=0.1$, where the same backbone with the LoRA
774 adapter *disabled* serves as the reference policy (no separate
775 frozen copy). Each rollout’s reward combines task correct-
776 ness, judged by GPT-4o on the imagined views, with an
777 action cost of 0.1 per atomic step, a parse-failure penalty of
778 -0.5 , and an additional -0.5 *skip-wrong* penalty when the
779 policy bypasses the world model and answers incorrectly;
780 this last term is essential to prevent the policy from col-
781 lapsing to “always skip”. The action vocabulary consists of
782 9° rotations and 0.25 m forward steps, capped at 6 atomic
783 actions per plan. We optimize with AdamW, lr= 2×10^{-5} ,
784 weight decay 0, gradient clipping at 1.0, per-device batch
785 size 1 and no gradient accumulation, yielding an effective
786 batch of 8 questions per optimizer step on $8\times A100$ -80 GB
787 GPUs (DDP via `torchrun`). The world model is Stable
788 Virtual Camera [59] run in `img2trajvids-prob` mode
789 (CFG 4.0, 8 target views, trajectory prior, `interp` chunk-
790 ing, short side 576). Training data is a curated 30/70 mix of
791 GPT-4o-prescored *easy-skip* and *needs-WM* questions; we
792 train for up to 300 optimizer steps and select the checkpoint
793 that achieves the best held-out accuracy. We train AVIC-
794 R with signal from GPT-4o, and zero-shot transfer to other
795 backbone models test.

796 **Balanced data curation.** Random training sampling is
797 inefficient, most SAT questions are already answered cor-
798 rectly without imagination (Case 3), so every reasonable
799 rollout produces $r \approx +1$, the within-group variance van-
800 ishes, and Eq. (5) yields zero gradient. The questions that
801 drive learning are those the base VLM fails on but a tar-
802 geted imagined view would correct. We therefore pre-score
803 a candidate pool with a strong reference QA model under
804 the `skip` policy and split each question into *easy-skip* (base
805 correct) or *needs-imagination* (base wrong). The final train-
806 ing set (3000 examples) retains every needs-imagination in-
807 stance and sub-samples easy-skip ones at a 30% ratio, bal-
808 ancing learning signal against an anchor that prevents col-
809 lapse onto always-call-WM.

810 **Lenient schema parsing.** High-temperature sampling
811 occasionally yields JSON outputs that are syntactically bro-
812 ken but semantically recoverable: d set to an action verb

(turn-left) rather than a meta-decision, or π written 813
as free-form strings ("turn-right 9 degrees"). A 814
strict parser would reject these as parse failures, conflating 815
format errors with semantic confusion. We instead use a le- 816
nient parser with a salvage stage that infers $d = \text{call_wm}$ 817
whenever π is non-empty regardless of the literal d value, 818
and parses natural-language action strings into structured 819
records via regex. Truly unrecoverable outputs still incur 820
 r_{parse} . The finer-grained reward distinguishes “intent re- 821
covered, format mangled” (taught via the call-WM reward) 822
from “no plan at all” (parse failure), supplying GRPO with 823
a smoother gradient. 824

825 **Metrics.** We evaluate SAT/MMSI spatial reasoning 825
benchmark with multiple-choice QA accuracy. In the R2R 826
navigation setting, the agent must follow natural language 827
instructions in indoor environments. We integrate our adap- 828
tive visual test-time scaling framework into the navigation 829
pipeline and measure performance. We evaluate navigation 830
performance using four standard metrics: Navigation Er- 831
ror (NE), Oracle Success Rate (OSR), Success Rate (SR), 832
and Success weighted by Path Length (SPL). NE measures 833
the geodesic distance between the agent’s final position and 834
the target, while SR reports the fraction of episodes where 835
the final position is within a predefined success threshold. 836
OSR measures whether the agent ever reaches within the 837
success threshold at any point along its trajectory, reflect- 838
ing exploration ability independent of stopping. SPL jointly 839
evaluates success and efficiency by weighting successful 840
episodes by the ratio between shortest-path length and the 841
actual trajectory length. 842

843 **Prompts.** We provide extra technical details of our adap- 843
tive visual test time scaling framework. In Tab. 12, we pro- 844
vide verifier prompts that are used to score each generated 845
trajectory, and in Tab. 13, we provide prompts for world 846
model gating and action planning in our policy model. 847

848 B. Extra Experiments

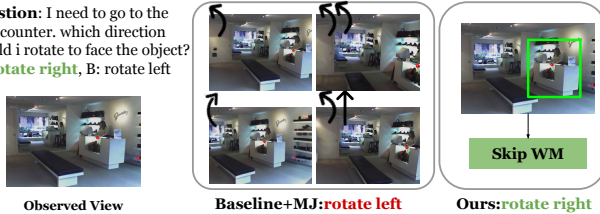
849 **Effect of selective gating and action-level scaling.** As 849
listed in Tab. 5, we analyze the contributions of world- 850
model (WM) imagination, gating, and action-level test-time 851
scaling. Vanilla baseline achieves 74.0% w.o any scaling. 852
Always-on invoking the WM with spatial beam search and 853
without gating or action scaling improves performance to 854
77.3%, but at the cost of excessive computation, requiring 855
an average of 12.34 WM calls. Introducing a gating mech- 856
anism via a policy model alone drastically reduces WM us- 857
age (0.51 calls) but also hurts accuracy (73.3%), indicating 858
that binary WM invocation without action-level control is 859
insufficient and can suppress necessary imagination. In con- 860
trast, our full method that combines gating with action scal- 861
ing achieves the best performance (79.3%) while keeping 862
WM usage low (0.73 calls). This demonstrates that *when* 863
to invoke the WM and *how* to use it are both critical: se- 864

| | ego-mov. | obj-mov. | goal-aim | action-conseq. | perspective | overall |
|-----------------------|--------------|--------------|--------------|----------------|--------------|--------------|
| AVIC-R w/o skip-wrong | 65.22 | 65.22 | 79.41 | 62.16 | 42.42 | 62.67 |
| AVIC-R (full) | 82.61 | 82.61 | 91.18 | 81.08 | 51.52 | 77.33 |
| Δ | +17.39 | +17.39 | +11.77 | +18.92 | +9.10 | +14.66 |

Table 4. **Reward ablation: the skip-wrong penalty.** Removing the -0.5 penalty for an incorrect `skip` causes the policy to collapse to never querying the world model, dropping overall accuracy by 14.66 points and degrading every question type. Values are test accuracy (%).

Visual Spatial Reasoning Example 1

Question: I need to go to the cash counter. which direction should i rotate to face the object?
A: rotate right, B: rotate left



Visual Spatial Reasoning Example 2

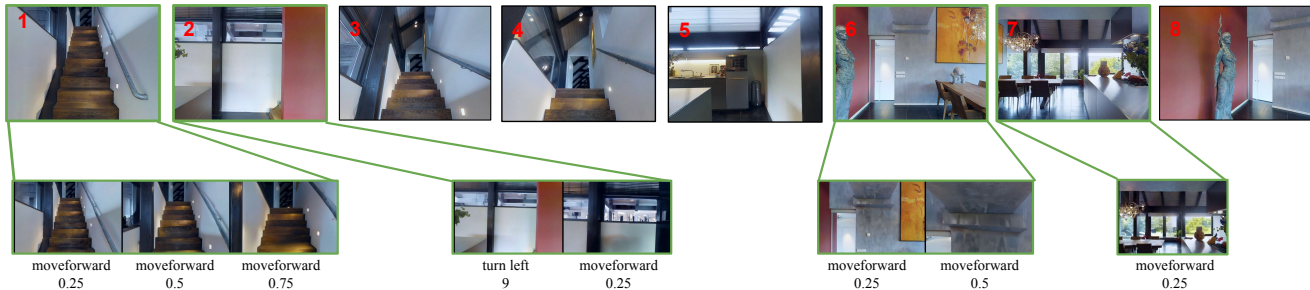
Question: if I stand in front of the trash bin, is there two person to my right?
A.No. B.Yes



Navigation Example

Global Instruction: Enter the house, and go into the kitchen. Stop next to the first counter on your left.

Observations + Adaptive Imagination w. World Model



Action Options

A: stop **B:** turn around to Place 1 which is corresponding to Image 1 **C:** turn around to Place 7 which is corresponding to Image 7
D: turn left to Place 5 which is corresponding to Image 5 **E:** go forward to Place 8 which is corresponding to Image 8

Figure 6. Qualitative examples on SAT of the always-on imagination method and our adaptive method, as well as the R2R navigation task. In the navigation example, the green option is selected by the model with adaptive imagination via our method, while the red one is without world model imagination.

| Action Scaling | Gating | WM | Avg. WM | Acc. (%) |
|----------------|--------|----|---------|----------|
| - | - | - | 0 | 74.0 |
| - | - | ✓ | 12.34 | 77.3 |
| - | ✓ | ✓ | 0.51 | 73.3 |
| ✓ | ✓ | ✓ | 0.73 | 79.3 |

Table 5. Ablation over action scaling, gating, and world model. Based on GPT-4.1.

| Policy | QA | Acc. (%) |
|--------|--------|----------|
| GPT-4o | GPT-4o | 68.0 |
| GPT-4o | o1 | 80.0 |
| o1 | o1 | 81.3 |
| o1 | GPT-4o | 68.6 |

Table 6. Effect of policy and QA model choice on SAT-Real.

865 lective gating must be paired with targeted action planning
 866 to ensure that limited imagination is informative rather than
 867 restrictive. Overall, it highlights that effective visual TTS
 868 requires control over both WM invocation and action plan-
 869 ning.

Effect of policy and QA model choice. As listed in Tab. 6,
 we compare different combinations of policy models and
 QA models on SAT-Real. We find upgrading the QA
 model yields substantial improvements regardless of the
 policy model used (68.0% \rightarrow 80.0%). These results indi-
 cate that SAT performance is primarily bottlenecked by

870
 871
 872
 873
 874
 875

876 the QA model’s spatial reasoning capability, but a stronger
877 policy model can also bring improvements. It also implies
878 that policy modelling mainly affects *efficiency and control*
879 of world-model invocation.

880 **The skip-wrong penalty is essential to RL training.** We
881 isolate the contribution of the *skip-wrong* reward term, a
882 -0.5 penalty applied whenever the policy chose *skip*
883 (bypassing the world model) and answered incorrectly.
884 Without this term, a wrong *skip* costs nothing while a
885 wrong *call_wm* costs at least the per-action cost (0.30 for
886 a typical 3-step plan with *action_cost=0.1*), so the opti-
887 mizer prefers skipping under any uncertainty and the pol-
888 icy quickly collapses to “always skip”. Reintroducing the
889 penalty makes a wrong skip strictly worse than the most ex-
890 pensive wrong WM call, restoring the incentive to query
891 the world model. Table 4 reports the result. The skip-
892 wrong term lifts overall test accuracy from 62.67% to
893 77.33% (+14.66 points), with the largest gains on *action-*
894 *consequence* (+18.92), *ego-movement* (+17.39), and *obj-*
895 *movement* (+17.39); without it, the policy never learns to
896 invoke the world model and its accuracy drifts toward the
897 no-WM baseline.

898 **Qualitative Analysis.** We also provide qualitative exam-
899 ples as illustrated at the top of Fig. 6. We compare our
900 adaptive visual TTS method with the always-on imagination
901 method, MindJourney (MJ). In the first example, the target
902 object (the *cash counter*) is already clearly visible in the ob-
903 served view. Our method correctly identifies that additional
904 visual imagination is unnecessary and directly skips world
905 model. In contrast, MJ indiscriminately invokes the world
906 model, generating multiple imagined views that introduce
907 misleading evidence and ultimately lead to an incorrect pre-
908 diction. In the second example, AVIC yields the correct an-
909 swer by selectively imagining the state where the agent is in
910 front of the *trash bin*. In contrast, MJ performs dense imagi-
911 nation and generates views that do not accurately reflect this
912 critical spatial condition, leading to an incorrect prediction.
913 Furthermore, we present a qualitative navigation example at
914 the bottom of Fig. 6. Our adaptive visual test-time scaling
915 selectively augments informative indoor observations (e.g.,
916 zooming in or turning to explore nearby views), enabling
917 the agent to better inspect the environment and align its ac-
918 tions with the global instruction (“*go to the kitchen*”). In
919 contrast, the baseline without visual imagination lacks suf-
920 ficient perceptual evidence and consequently chooses an in-
921 correct direction.

922 **Sensitivity to Errors from World Models.** Our work is
923 motivated by the observation that imperfect imagination can
924 introduce misleading or noisy evidence that may degrade
925 performance (Sec. 3, Figs. 1–2). To further evaluate robust-
926 ness, we conduct additional experiments using *Cosmos* [1]

in Tab. 10 as an alternative world model, which tends to
produce less visually stable camera trajectories and noisier
geometric structures compared to SVC.

Stage-wise Computation Cost. The decision module
adds a lightweight inference step whose cost is small rela-
tive to world-model (WM) generation. Rather than al-
ways invoking expensive imagination, it predicts *when* and
how to use the WM, reducing unnecessary calls. The
table shows that policy cost is minor compared to WM
cost and is offset by large savings. Although AVIC in-
troduces ~ 14.9 s of policy overhead, it reduces WM time
by ~ 153.7 s ($163.32 \rightarrow 9.59$), yielding a $\sim 6\times$ **reduction**
in total time ($177.84 \rightarrow 29.04$) and $\sim 20\times$ **fewer tokens**
($162.6 \rightarrow 7.6$ k), while improving accuracy by 2.0 points
over always-on. The decision module thus accounts for a
small fraction of total compute and is more than amortized
by the reduction in expensive WM calls. AVIC achieves
both higher accuracy and substantially lower overall com-
putation.

Robustness across runs. To assess the stability of AVIC-
R, we repeat our main evaluation three times with inde-
pendently sampled rollouts; all other hyperparameters are
held fixed. The overall accuracy gave a mean of 69.33 with
a sample standard deviation of 0.77 (standard error of the
mean 0.44). The narrow overall spread (< 1 point) confirms
that the gains reported in Tab. 1 are not driven by a single
lucky run.

While overall performance slightly degrades with *Cosmos*
due to increased noise and inconsistencies in the imag-
ined views, *the relative improvement from AVIC remains*
consistent (+2.7 over the GPT-o1 baseline), indicating that
our method is robust to moderate world-model errors. We
attribute this robustness to two design factors:

- **Gating.** The adaptive imagination mechanism selectively
invokes the world model only when informative, rather
than relying on all generated samples.
- **Action plan execution and trajectory selection.** The
reasoning process operates over multiple imagined per-
spectives, mitigating the impact of occasional erroneous
generations.

Evaluations on Additional Benchmark. We extend our
evaluation to *MindCube* [45] in Tab. 9. AVIC improves
performance on this fine-grained spatial reasoning bench-
mark, indicating that our method generalizes across tasks.

Framework Error Analysis. Our adaptive world-model
(WM) invocation policy does not call the WM uniformly
across tasks. It triggers WM imagination most frequently
for Egocentric Movement tasks (**EgoM**, 82.6%) and Action
Consequence tasks (**EgoAct**, 70.2%), while being much

| Method | Policy (s) | WM (s) | QA (s) | Total (s) | Tokens (k) | WM Calls | Acc. (%) |
|-------------------------|------------|--------|--------|-----------|------------|----------|----------|
| Baseline (no WM) | 0 | 0 | 3.9 | 3.9 | 0.7 | 0 | 74.0 |
| MindJourney (always-on) | 0 | 163.32 | 14.52 | 177.84 | 162.6 | 12.42 | 77.3 |
| AVIC | 14.92 | 9.59 | 4.53 | 29.04 | 7.6 | 0.73 | 79.3 |

Table 7. **Inference cost and accuracy on SAT-Real.** AVIC adds modest policy-inference overhead but substantially reduces world-model time and token usage, achieving the highest accuracy at $\sim 6\times$ lower total time and $\sim 20\times$ fewer tokens than always-on imagination.

| Method | EgoM | ObjM | EgoAct | Goal | Pers | Avg. |
|---|-------------|-------------|-------------|-------------|-------------|-------------|
| <i>when to call world model (call_wm %)</i> | | | | | | |
| MindJourney (always-on) | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 |
| AVIC (GPT-4o) | 100.0 | 0.0 | 62.2 | 32.4 | 54.5 | 50.0 |
| AVIC (Qwen2.5VL) | 47.8 | 60.9 | 78.4 | 73.5 | 54.5 | 64.7 |
| AVIC-R (Qwen2.5VL) | 91.3 | 91.3 | 100.0 | 97.1 | 78.8 | 92.0 |
| <i>how much to imagine (# img per question)</i> | | | | | | |
| MindJourney (always-on) | 11.13 | 4.83 | 8.37 | 11.03 | 8.55 | 8.90 |
| AVIC (GPT-4o) | 3.30 | 0.00 | 5.30 | 4.18 | 5.39 | 4.55 |
| AVIC (Qwen2.5VL) | 4.09 | 3.07 | 4.90 | 3.28 | 4.39 | 4.03 |
| AVIC-R (Qwen2.5VL) | 2.95 | 3.00 | 4.59 | 2.79 | 4.23 | 3.60 |
| <i>task accuracy (%)</i> | | | | | | |
| MindJourney (always-on) | 78.3 | 60.9 | 78.4 | 70.6 | 57.5 | 69.3 |
| AVIC (GPT-4o) | 87.0 | 69.6 | 64.9 | 82.4 | 48.5 | 69.3 |
| AVIC (Qwen2.5VL) | 65.2 | 73.9 | 64.9 | 91.2 | 60.6 | 71.3 |
| AVIC-R (Qwen2.5VL) | 82.6 | 82.6 | 81.1 | 91.2 | 51.5 | 77.3 |

Table 8. **Analysis of when and how much to imagine on SAT-Real.** We compare four imagination policies along three axes: fraction of questions on which the world model is invoked (*top*), average number of imagined views per question (*middle*), and resulting QA accuracy (*bottom*).

| Method | Acc. (%) |
|---------------|--------------------|
| GPT-4o | 36.5 |
| GPT-4o + AVIC | 38.7 (+2.2) |

Table 9. Results on MindCube-Tiny

| Method | Acc. (%) |
|--------------------|----------|
| o1 | 74.6 |
| o1 + AVIC (SVC) | 85.3 |
| o1 + AVIC (Cosmos) | 77.3 |

Table 10. Performance across world models. Based on o1.

| Metric | EgoM | ObjM | EgoAct | Goal | Pers. | Avg. |
|-----------|-------|------|--------|------|-------|------|
| Recall | 100.0 | 33.3 | 55.6 | 33.3 | 52.6 | 43.9 |
| Precision | 5.6 | 50.0 | 19.2 | 22.2 | 62.5 | 27.1 |

Table 11. **Gating recall and precision (%)** per category. *Recall*: fraction of WM-needed questions on which the policy calls WM. *Precision*: fraction of WM-called questions that actually benefit from WM.

C. Impact Statement

This work studies world-model-based visual imagination in visual spatial reasoning and highlights the limitations of existing always-on test-time imagination methods. Through systematic analysis, we show that indiscriminate visual

976 more conservative for goal-oriented tasks (**Goal**, 26.4%).
 977 While frequent WM usage on **EgoM** improves accuracy,
 978 it is misaligned with the dominant error sources identi-
 979 fied manually in Observation 1 and 2, where many fail-
 980 ures instead stem from action-conditioned and viewpoint-
 981 dependent reasoning. This mismatch results in low recall
 982 and precision for cases that truly require world-model im-
 983 agination, as we reported in Tab. 11. It indicates that the
 984 current policy design remains a significant chance for im-
 985 provement. Overall, these results reveal substantial room
 986 for improving adaptive WM calling strategies, motivating
 987 future work on error-aware and state-aware invocation poli-
 988 cies that better align WM usage with underlying reasoning
 989 demands.

990

991

992

993

994

995 imagination can be computationally inefficient and, in some
996 cases, harmful due to misleading or redundant imagined
997 views. Our findings emphasize the importance of adaptive
998 test-time computation, demonstrating that effective spatial
999 reasoning requires selectively invoking visual imagination
1000 only when necessary and scaling it appropriately. Beyond
1001 the specific benchmarks studied, our insights are broadly
1002 applicable to multimodal agents that rely on test-time sim-
1003 ulation, including embodied AI and interactive systems.

1004 **D. Limitations**

1005 Our work focuses on adaptive imagination control for vi-
1006 sual spatial reasoning and short-horizon embodied naviga-
1007 tion; extending the framework to longer-horizon decision-
1008 making, manipulation, and broader multi-modal tasks is a
1009 natural next direction. The framework also assumes access
1010 to a separate visual world model and a fixed discrete action
1011 space, leaving room for richer extensions such as contin-
1012 uous action spaces, multiple specialized world models, or
1013 joint training of the gating policy with the world model it-
1014 self. Incorporating world-model uncertainty into the reward
1015 signal is another promising direction for improving robust-
1016 ness under noisier imagined rollouts.

1017 **E. License**

1018 We will make our code and models publicly accessible. We
1019 use standard licenses from the community and provide the
1020 following links to the licenses for the datasets, codes, and
1021 models that we used in this paper. For further information,
1022 please refer to the specific link.

1023 **QWen2.5VL [4]:** Apache-2.0

1024 **SAT [30]:** MIT

1025 **MMSI [42]:** CC-BY-4.0

1026 **R2R [2]:** MIT

Role. You are an *independent evaluator* for visual spatial reasoning.

Input. A multiple-choice question, answer options, the current observation image(s), and *one* candidate action plan. The plan includes imagined views rendered by a world model.

Task. Score how *useful* the imagined views are for answering the question.

Score Range. Integer from 0 (not helpful, irrelevant, or low quality) to 9 (highly helpful and informative).

Scoring Guidelines.

- Assign higher scores if the imagined views reveal missing evidence needed to answer the question (e.g., resolving occlusion or viewpoint ambiguity).
- Assign higher scores if the imagined views are sharp, coherent, and visually consistent.
- Assign lower scores if the views are redundant, uninformative, distorted, or unrelated to the question.
- If the original observations are already sufficient, most plans should receive low scores.

Rules.

- Do **not** answer the question.
- Output **only** a single integer between 0 and 9.
- Do not output any additional text.

Output Example.

5

Table 12. Verifier prompts for scoring imagined view plans.

Role. You are a *policy model* for spatial reasoning in a 3D indoor environment. Your goal is to decide whether to invoke a world model (WM) and, if needed, plan actions that acquire the most informative imagined views.

Input. One or more images, a multiple-choice question, and answer options.

Tasks.

- Decide whether to *SKIP* or *CALL* the world model.
- If *CALL*, generate a short action plan (1–6 actions) to gather additional visual evidence.

Action Space (Discrete, Fixed).

- *move-forward* 0.25 meters
- *turn-left* 9 degrees
- *turn-right* 9 degrees

Action Composition Guidelines.

- Repeated turns approximate larger rotations (e.g., 2 turns $\approx 18^\circ$, 5 turns $\approx 45^\circ$, 10 turns $\approx 90^\circ$).
- When a question specifies a larger angle, approximate it using repeated 9° turns.

When to Call the World Model.

- The answer is not directly observable from the current view.
- The question depends on perspective, facing direction, rotation, or left/right relations.
- The question requires reasoning about motion or state changes over time.

Constraints.

- Do not generate cancelling or oscillating actions (e.g., left then right).
- If turning, choose a single direction and turn monotonically.

Output Format (JSON only).

```
{
  "decision": "skip" | "call_wm",
  "reason": "<one sentence>",
  "actions": [
    {"type": "move-forward" | "turn-left" | "turn-right",
     "value": <number>}
  ]
}
```

Table 13. Policy model prompts for world model gating and action planning.