

Learning from Trials and Errors: Reflective Test-Time Planning for Embodied LLMs

Anonymous CVPR submission

Paper ID ****

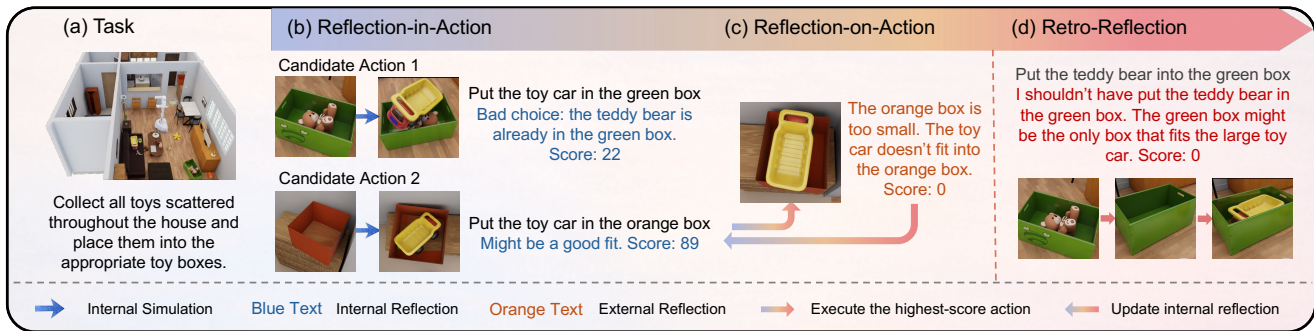


Figure 1. Conceptual overview of Reflective Test-Time Planning. The agent (a) receives a long-horizon task, (b) performs *reflection-in-action* by internally simulating and scoring candidate actions, (c) performs *reflection-on-action* by updating its beliefs and decision systems based on execution outcomes, and (d) conducts *retrospective reflection* to revise earlier decisions with hindsight.

Abstract

001 Embodied LLMs endow robots with high-level task reason-
 002 ing, but they cannot reflect on what went wrong or why,
 003 turning deployment into a sequence of independent trials
 004 where mistakes repeat rather than accumulate into experi-
 005 ence. Drawing upon human reflective practitioners, we in-
 006 troduce Reflective Test-Time Planning, which integrates two
 007 modes of reflection: *reflection-in-action*, where the agent
 008 uses test-time scaling to generate and score multiple can-
 009 didate actions using internal reflections before execution;
 010 and *reflection-on-action*, which uses test-time training to
 011 update both its internal reflection model and its action pol-
 012 icy based on external reflections after execution. We also
 013 include *retrospective reflection*, allowing the agent to re-
 014 evaluate earlier decisions and perform model updates with
 015 hindsight for proper long-horizon credit assignment. Ex-
 016 periments on our newly-designed Long-Horizon Household
 017 benchmark and MuJoCo Cupboard Fitting benchmark show
 018 significant gains over baseline models, with ablative studies
 019 validating the complementary roles of *reflection-in-action*
 020 and *reflection-on-action*. Qualitative analyses, including
 021 real-robot trials, highlight behavioral correction through
 022 reflection.

“Error isn’t simple darkness, it sheds a light of its
 own.” — Kathryn Schulz, *Being Wrong*

023
024

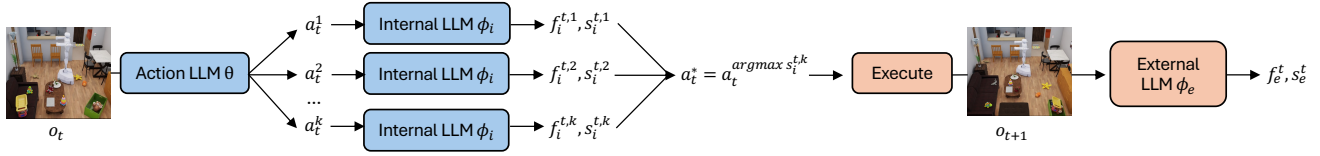
1. Introduction

025
 026 Embodied LLMs [10, 15, 41] equip agents with task plan-
 027 ning abilities, but they remain brittle static oracles that can-
 028 not learn from failures, turning deployment into indepen-
 029 dent trials of repeated mistakes rather than accumulated ex-
 030 perience. Humans, in contrast, are natural reflective prac-
 031 titioners. Drawing on Schön’s framework about reflective
 032 planning [28], humans fluidly alternate between two modes
 033 of reflection: through *reflection-in-action*, we engage in
 034 internal simulation, questioning whether our planned ap-
 035 proach will actually work given what we currently under-
 036 stand; through *reflection-on-action*, we use the actual out-
 037 comes to reshape both our beliefs about the environment
 038 and our strategies for acting within it. An illustrative exam-
 039 ple of these reflection modes is shown in Figure 1. This
 040 bidirectional flow allows us to learn not only from out-
 041 comes, but also from the very process of engaging with an
 042 uncertain world.

043 Current approaches, however, capture at best a superfi-
 044 cial version of one mode while neglecting the other. One
 045 line of work [21, 30] uses LLM-based verbal reflection,

046	generating natural-language critiques of past behavior to	099
047	condition future actions. While this enables reflection-on-	100
048	action at the level of reasoning traces, reflections are stored	101
049	only as contextual text and do not update the underlying de-	
050	cision process, making their effects transient and prone to	
051	repetition under distributional shift. A second line of work	
052	[5, 37] relies on internal world models to guide action selec-	
053	tion in embodied environments. These approaches support	
054	reflection-in-action through anticipated outcomes, but typi-	
055	cally assume fixed, pretrained dynamics models that may	
056	be wrong in ways only revealed during execution.	
057	To operationalize both reflection modes in embodied	
058	settings, we introduce Reflective Test-Time Planning, a	
059	framework that seamlessly unifies reflection-in-action and	
060	reflection-on-action for embodied agents during test-time	
061	deployment. Concretely, the framework employs three em-	
062	bodyed LLMs during deployment: an action generation	
063	model π_θ , an internal evaluator V_{ϕ_i} , and an external eval-	
064	uator V_{ϕ_e} . During reflection-in-action, the agent samples	
065	N candidate actions via high-temperature sampling, uses	
066	V_{ϕ_i} to generate internal reflections scoring each candidate,	
067	then executes the highest-scoring action. After execution,	
068	V_{ϕ_e} generates an external reflection, providing an immedi-	
069	ate, language-based evaluation of what happened and why.	
070	This immediate external reflection grounds the agent’s	
071	beliefs in reality, but remains inherently local—it only eval-	
072	uates consequences visible at the next timestep. Yet many	
073	embodied failures are non-local: an action that appears suc-	
074	cessful may later block progress, and a seemingly subopti-	
075	mal action may enable future success. To address this tem-	
076	poral credit assignment problem, we introduce retrospective	
077	reflection, where V_{ϕ_e} periodically re-evaluates earlier deci-	
078	sions with hindsight (e.g., at room transitions or after re-	
079	peated failures). These hindsight assessments provide self-	
080	supervised signals at deployment time, enabling two forms	
081	of test-time training: (1) policy gradient for π_θ to favor ac-	
082	tions that score well under hindsight, and (2) supervised	
083	learning for V_{ϕ_i} to anticipate what hindsight will reveal.	
084	Because these updates revise not only the action policy but	
085	also the predictive assumptions behind it, the process con-	
086	stitutes a form of double-loop learning [1], in which agents	
087	learn not merely from outcomes but from diagnosing and	
088	correcting the underlying causes of their errors.	
089	We evaluate our approach on two embodied benchmarks	
090	that we design to stress error-driven adaptation: (1) a Long-	
091	Horizon Household benchmark that requires failure recov-	
092	ery during multi-step planning across rooms, and (2) a con-	
093	trolled MuJoCo Cupboard Fitting benchmark that isolates	
094	geometric placement failures. Our framework achieves	
095	large gains over reflective language, RL and world-model	
096	baselines. Ablations indicate that improvement emerges	
097	only when both reflection-in-action and reflection-on-action	
098	take place, and when both action policy and internal re-	
	flexion model are updated during deployment. Qualitative	099
	analyses, including preliminary real-robot trials, show that	100
	reflection reduces repetitive failure modes in practice.	101
	2. Related Works	102
	Test-Time Adaptation (TTA) and Learning enables mod-	103
	els to adjust to distribution shifts during inference without	104
	source data [20, 31, 33]. Early methods minimize pre-	105
	diction entropy, with Tent [33] updating batch-norm pa-	106
	rameters online and later work adding calibrated objec-	107
	tives [26, 34]. Parameter-efficient tuning further boosts	108
	TTA: LoRA [11] enables low-rank weight updates with lit-	109
	tle memory [16], while bias-only tuning [3] offers alterna-	110
	tive efficiency–accuracy trade-offs. Recent extensions op-	111
	erate on hidden state representations [32], supporting long-	112
	context memory. For embodied settings, continual learning	113
	frameworks [17, 23] demonstrate viability in manipulation	114
	and navigation [8, 22]. We adapt models at test time via	115
	self-supervised reflective signals extracted from the agent’s	116
	own verbal assessments.	117
	Multimodal Embodied Large Language Models cou-	118
	ple visual perception with language understanding to en-	119
	able embodied planning. Recent foundation models lever-	120
	age large-scale robotic datasets for zero-shot generaliza-	121
	tion [2, 15, 41], with RT-2 transferring web knowledge	122
	and OpenVLA offering open-source support across hetero-	123
	geneous embodiments. 3D spatial understanding has be-	124
	come central, spanning point clouds [9], 3D patches [40],	125
	and lightweight point cloud injection [19]. Further exten-	126
	sions incorporate multisensory interaction [10], generative	127
	world models for manipulation [37], and long-term spa-	128
	tial–temporal embodied memory [12]. Architectural ad-	129
	vances explore interleaved multimodal instructions [4] and	130
	chain-of-thought reasoning [24, 36] for improved task de-	131
	composition. In contrast, we view deployment as a learning	132
	phase: instead of acting as a fixed policy, our embodied	133
	multimodal LLM reflects on its own actions and updates it-	134
	self via test-time training.	135
	Reflection and Self-Improvement in AI Agents. Re-	136
	flexion mechanisms enable agents to learn from failures	137
	through self-critique and refinement. Verbal self-reflection	138
	methods (e.g., Reflexion [30]) store natural-language cri-	139
	tiques to guide future actions, with extensions explor-	140
	ing self-refinement [21, 35], tool-assisted verification [6],	141
	curiosity-driven reflection [14], multi-agent systems [25],	142
	and robotics [13]. These approaches support reflection-on-	143
	action but store reflections only as text, influencing behavior	144
	indirectly without updating model parameters, and thus re-	145
	main brittle under distribution shift. A complementary line	146
	enables reflection-in-action via internal predictive models	147
	that anticipate outcomes [5, 7, 37, 38]. Such world mod-	148
	els are typically fixed despite evolving dynamics in em-	149
	bodyed settings. Our method unifies reflection-in-action	150

(a) Reflection-in-Action and External Reflection



(b) Reflection-on-Action

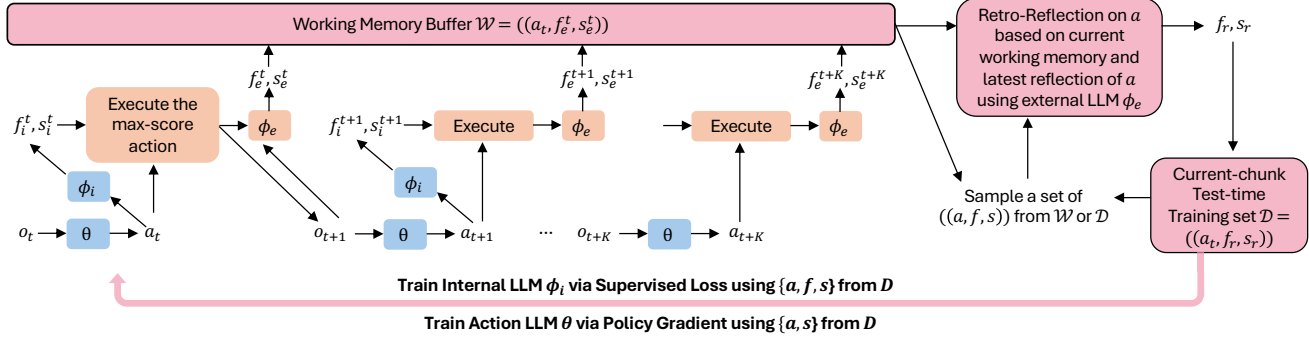


Figure 2. **Method overview.** (a) *Reflection-in-action*: multiple candidate actions are generated and scored by an internal reflection LLM prior to execution. (b) *Reflection-on-action*: iteratively invoked when working memory hits K or at key milestones. Executed actions are critiqued by an external reflection LLM and stored in a working memory buffer; at milestones, hindsight re-evaluation assigns long-horizon credit. The resulting verbal reflections form self-supervised training data to update both the internal reflection LLM (supervised loss) and the action LLM (policy gradient) via test-time training, enabling agents to learn from execution experience during deployment.

151 and reflection-on-action by converting reflections into self-
 152 supervised training signals for parameter updates during de-
 153 ployment.

154 3. Reflective Test-Time Planning

155 Consider an embodied agent with a multimodal large lan-
 156 guage model operating on task τ in a partially observ-
 157 able environment. At each timestep t , the model receives
 158 multimodal observation o_t , generates action a_t in nat-
 159 ural language, and receives execution feedback e_t indicat-
 160 ing whether the action executed successfully (e.g., object
 161 grasped, placement completed). Crucially, a positive e_t
 162 indicates successful execution but does not imply the action
 163 was strategically correct or contributes to full task comple-
 164 tion.

165 Traditional multimodal embodied LLM systems keep
 166 model parameters fixed at inference, limiting adaptation to
 167 novel scenarios or recovery from failures. We depart from
 168 this static inference setting by building an adaptive test-time
 169 framework that employs three interacting models: an ac-
 170 tion generation LLM π_θ that produces actions given ob-
 171 servations, an internal reflection LLM V_{ϕ_i} that gener-
 172 ates pre-action evaluations, and an external reflection LLM
 173 V_{ϕ_e} that generates post-execution assessments. These mul-
 174 timodal LLMs are first initialized with basic capabilities for
 175 reasoning in specific embodied environments through mini-
 176 mum supervised fine-tuning on a small set of tasks, en-
 177 abling them to understand action formats, generate reflec-
 178 tions, and process 3D observations before they can effec-

tively learn from test-time experience. At deployment time,
 we introduce three reflection types: internal reflection f_i for
 pre-action scoring, external reflection f_e for post-execution
 assessment, and retrospective reflection f_r for hindsight re-
 evaluation. We combine test-time scaling (generating and
 scoring multiple candidate actions) for reflection-in-action,
 with test-time training (tuning π_θ and V_{ϕ_i}) for reflection-
 on-action. Figure 2 shows an overview of our method. We
 also provide a detailed method breakdown in Algorithm 1.

We choose verbal reflection as the representation for all
 reflection types inspired by “double-loop learning” [1]: by
 articulating what went wrong and why, the agent abstracts
 generalizable lessons transferable to future decisions rather
 than merely reporting that an action failed. These articu-
 lated lessons serve as supervisory signals during deploy-
 ment, providing interpretable feedback to be reused later.
 Thus, instead of only training the action model based on
 outcomes during test time (single loop), we also train the
 internal reflection LLM to align its pre-action internal ref-
 lections with post-execution external reflections, updating
 the underlying reasoning process behind the action itself.

3.1. Reflection-in-Action

Humans naturally deliberate under uncertainty by *mentally*
 simulating and reflecting on actions. We transfer this abil-
 ity to embodied agents via *reflection-in-action*: rather than
 greedily selecting the first plausible action, the agent sam-
 ples several candidates and reflects on each one before
 committing. We implement this through test-time scaling,

207 where we generate N diverse candidate actions and use the
208 internal reflection LLM to produce reflective evaluations for
209 each candidate, which guide action selection.

Algorithm 1 Reflective Test-Time Planning

Require: Task τ , initial observation o_1 ; action LLM π_θ , internal reflection LLM V_{ϕ_i} , external reflection LLM V_{ϕ_e} , window size K

- 1: Initialize $a_0, f_e^{(0)} = \text{None}$; working memory buffer, test-time training set, retro-buffer $\mathcal{W}, \mathcal{D}_{\text{train}}, \mathcal{D}_{\text{retro}} \leftarrow \emptyset$; Temperature T ; Max step \mathcal{T} ;
- 2: **for** $t = 1, 2, \dots, \mathcal{T}$ **do**
- 3: **// Reflection-in-Action**
- 4: Construct action prompt x_{action} from $\tau, o_t, a_{t-1}, f_e^{t-1}$
- 5: Sample N actions: $\{a_t^k\}_{k=1}^N \sim \pi_\theta(\cdot | x_{\text{action}}; T)$
- 6: **for** $k = 1, \dots, N$ **do**
- 7: Construct internal reflection prompt x_{internal}^k
- 8: Generate score: $f_i^{t,k}, s_i^{t,k} \leftarrow V_{\phi_i}(x_{\text{internal}}^k)$
- 9: **end for**
- 10: Select $a_t^* = a_t^{(\arg \max_k s_i^{t,k})}$
- 11: Execute a_t^* in environment \rightarrow observe (o_{t+1}, e_t)
- 12: Break if task complete
- 13: **// Reflection-on-Action**
- 14: Construct external reflection prompt x_{external}
- 15: Generate reflection feedback: $f_e^t, s_e^t \leftarrow V_{\phi_e}(x_{\text{external}})$
- 16: $\mathcal{W} \leftarrow \mathcal{W} \cup \{(o_t, a_t^*, f_e^t, s_e^t)\}$
- 17: **if** $|\mathcal{W}| = K$ **or** hit key milestone **then**
- 18: **// Retrospective Reflection: Re-evaluate with Hindsight**
- 19: **for** each $(o_j, a_j, f_j^j, s_j^j) \in \mathcal{W} \cup \mathcal{D}_{\text{retro}}$ **do**
- 20: Construct retro prompt x_{retro}^j from $\tau, a_j, \mathcal{W}, o_{t+1}, f_j^j, s_j^j$
- 21: Generate revised reflection: $f_r^j, s_r^j \leftarrow V_{\phi_e}(x_{\text{retro}}^j)$
- 22: $\mathcal{D}_{\text{retro}}[a_j] \leftarrow (a_j, f_r^j, s_r^j)$ // reflection overwrite for a_j
- 23: **end for**
- 24: $\mathcal{D}_{\text{train}} \leftarrow \mathcal{D}_{\text{retro}}$
- 25: **// Regularization: Prevent Catastrophic Forgetting**
- 26: **for** sampled unexplored action a_l **do**
- 27: Construct pre-action internal x_{internal}^l
- 28: $f_i^l, s_i^l \leftarrow V_{\phi_i}(x_{\text{internal}}^l)$ // Original output
- 29: $\mathcal{D}_{\text{train}} \leftarrow \mathcal{D}_{\text{train}} \cup \{(a_l, f_i^l, s_i^l)\}$
- 30: **end for**
- 31: **// Test-Time Training: Supervised Learning for Internal LLM**
- 32: **for** epoch $e = 1, \dots, E$ **do**
- 33: **for** $(a, f, s) \in \mathcal{D}_{\text{train}}$ **do**
- 34: Construct x_{internal} from a
- 35: Compute loss: $\ell_\phi = -\log p_{\phi_i}(f | x_{\text{internal}})$
- 36: Update: $\phi_i \leftarrow \phi_i - \eta_\phi \nabla_{\phi_i} \ell_\phi$
- 37: **end for**
- 38: **end for**
- 39: **// Test-Time Training: REINFORCE for Action LLM**
- 40: **for** RL step $= 1, \dots, \text{RLSteps}$ **do**
- 41: **for** $(a, f, s_r) \in \mathcal{D}_{\text{train}}$ **do**
- 42: Construct x_{action} from a . Reward $r = 2s_r/100 - 1$
- 43: Compute log probability: $\log p_\theta(a | x_{\text{action}})$
- 44: Compute REINFORCE loss: $\ell_\theta = -r \cdot \log p_\theta(a | x_{\text{action}})$
- 45: **end for**
- 46: Accumulate gradients and update: $\theta \leftarrow \theta - \eta_\theta \nabla_\theta \ell_\theta$
- 47: **end for**
- 48: Clear working memory: $\mathcal{W} \leftarrow \emptyset$
- 49: **end if**
- 50: **end for**

210 **Vanilla Action Generation.** Standard autoregressive
211 generation from the action LLM π_θ produces a single action
212 via greedy or low-temperature sampling:

$$213 \quad a_t = \arg \max_{a \in \mathcal{A}} p_\theta(a | o_t) \quad (1)$$

214 where o_t is the current observation, with some architectures
215 implicitly encoding historical context through mem-

ory structures (e.g., [12]). Greedy generation commits to
actions early without reflecting on potential consequences.

Candidates Generation. Different from the above, at
each decision step t , we construct an action generation
prompt x_{action} containing task description τ , current obser-
vation o_t , previous action a_{t-1} , and previous external reflection
 f_e^{t-1} (both initialized as None; external reflection
will be introduced later). We sample N diverse candidate
actions:

$$a_t^k \sim p_\theta(\cdot | x_{\text{action}}; T) \quad k = 1, \dots, N \quad (2)$$

where high temperature T encourages diversity in the gener-
ated candidates.

Internal Reflection Scoring. For each candidate, we
construct an internal reflection prompt x_{internal}^k , which is
identical to x_{action} except that it adds the candidate action
 a_t^k to be evaluated. The internal reflection LLM generates:

$$f_i^{t,k}, s_i^{t,k} = V_{\phi_i}(x_{\text{internal}}^k) \quad (3)$$

where $s_i \in [0, 100]$ is a numerical score and $f_i^{t,k}$ is natural
language reflection. Because this reflection occurs prior to
action execution, we refer to it as internal reflection.

Best Action Selection. We select the highest-scoring
candidate:

$$a_t^* = a_t^{(\arg \max_{k \in [N]} s_i^{t,k})} \quad (4)$$

Rather than executing the first feasible action, the agent
“mentally tries out” multiple options and chooses the one
it internally judges as most promising.

3.2. Reflection-on-Action

Reflection-in-action has a limitation: internal reflection op-
erates in imagination, not reality. The internal reflection
LLM may score an action highly based on plausible reason-
ing, yet the action fails due to unforeseen physical con-
straints or environmental dynamics. Reflection-on-action,
learning from experience after actions are executed, ad-
dresses this by grounding learning in actual execution out-
comes.

3.2.1. Multi-Scale External Reflection

External Reflection Generation. After executing a_t^* and
observing (o_{t+1}, e_t) , we construct x_{external} by extending
 x_{action} with a_t^*, e_t , and (o_t, o_{t+1}) . The external reflection
LLM generates:

$$f_e^t, s_e^t = V_{\phi_e}(x_{\text{external}}) \quad (5)$$

where s_e^t is a score and f_e^t is language feedback assess-
ing the immediate outcome and its cause. This provides
real-time assessment based on directly observable conse-
quences.

Working Memory Buffer. We maintain a buffer $\mathcal{W}_t =$
 $(o_j, a_j, e_j, f_e^j) \mid j = t - K + 1, \dots, t$ with K steps. This

263 buffer accumulates recent experience until reaching a key
264 milestone (e.g., exiting a room, detecting repeated failures
265 that need replanning) or when $|\mathcal{W}_t| = K$, at which point we
266 trigger memory consolidation and test-time training.

267 **Retrospective Reflection with Hindsight.** A critical
268 limitation of external reflection is that it evaluates actions
269 based only on immediate outcomes. An action may appear
270 successful initially but later prove problematic (e.g., placing
271 an object in an accessible compartment that blocks the only
272 space for a larger object). To address this credit assignment
273 problem, we introduce retrospective reflection.

274 Once we hit key milestone or reach the working memory
275 limit, for each action a_j that has been reflected before (ei-
276 ther by immediate external reflection or by previous retro-
277 reflection), the external reflection LLM re-evaluates the ac-
278 tion with full hindsight:

$$279 \quad f_r^j, s_r^j = V_{\phi_e}(x_{\text{retro}}^j) \quad (6)$$

280 where the retrospective prompt x_{retro}^j includes: (1) the com-
281 plete working memory window \mathcal{W}_t as context; (2) a histor-
282 ical action a_j to be retro-reflected based on the current out-
283 comes; (3) its most recent reflection f_{recent}^j (either f_e^j from
284 the current working memory \mathcal{W}_t if this is the first retrospec-
285 tive evaluation, or f_r^j from the previous retrospective round
286 and stored in a retro-buffer $\mathcal{D}_{\text{retro}}$; and (4) the current ob-
287 servation o_{t+1} . After retro-revision, we update $\mathcal{D}_{\text{retro}}$ and
288 store only the most recent retro-reflection for each action.
289 As $\mathcal{W} \cup \mathcal{D}_{\text{retro}}$ grows larger with actions, we may subsample
290 historical actions if necessary to keep it tractable.

291 3.2.2. Test-Time Training Dataset Construction.

292 We construct training dataset $\mathcal{D}_{\text{train}}$ with two types of data:

293 **Retro-supervised pairs:** For any action a_j that has been
294 retrospectively evaluated, we create training pairs using the
295 retrospective reflection f_r^j and s_r^j :

$$296 \quad \mathcal{D}_{\text{retro}} = (a_j, f_r^j, s_r^j) \quad (7)$$

297 These pairs use hindsight-corrected reflections and scores for
298 training both the internal LLM V_{ϕ_i} and action LLM π_{θ} .

299 **Regularization pairs:** To prevent catastrophic forget-
300 ting, we randomly sample unexplored actions a_l , construct
301 x_{internal}^l and use the internal LLM’s current predictions:

$$302 \quad \mathcal{D}_{\text{reg}} = (a_l, f_i^l, s_i^l) \quad (8)$$

303 where $f_i^l, s_i^l = V_{\phi_i}(x_{\text{internal}}^l)$ represents the model’s current
304 output for randomly sampled action a_l . This anchors the
305 model to its existing knowledge for actions not updated by
306 recent experience, preventing distribution shift caused by
307 training exclusively on retrospectively evaluated actions.

308 3.2.3. Test-Time Training

309 **Internal Reflection LLM Training via Supervised**
310 **Learning.** We train the internal reflection LLM V_{ϕ_i} to

311 predict retrospective reflections using standard supervised
312 learning. The objective minimizes negative log-likelihood
313 over the combined dataset $\mathcal{D}_{\text{train}} = \mathcal{D}_{\text{retro}} \cup \mathcal{D}_{\text{reg}}$:

$$\mathcal{L}_{\text{internal}}(\phi_i) = \mathbb{E}(x_{\text{internal}}, f, s) \sim \mathcal{D}_{\text{train}} [-\log p_{\phi_i}(f|x)] \quad (9) \quad 314$$

315 Where we construct x_{internal} from each action a . We perform
316 E epochs of test-time training:

$$\phi_i^{(e+1)} = \phi_i^{(e)} - \eta_{\phi} \nabla_{\phi_i} \mathcal{L}_{\text{internal}} \quad (10) \quad 317$$

318 **Action LLM Training via RL.** The action LLM π_{θ} is
319 updated using policy gradient with retrospective scores as
320 rewards. We convert the retrospective score $s_r \in [0, 100]$
321 to a reward signal $r = 2 * (s_r/100) - 1$ mapping scores to
322 $[-1, 1]$. For each training example, we compute the log-
323 probability of the executed action sequence:

$$\log p_{\theta}(a|x_{\text{action}}) = \sum_{i=1}^{|a|} \log p_{\theta}(a_i|a_{<i}, x_{\text{action}}) \quad (11) \quad 324$$

325 where the sum is over action tokens. We construct x_{action}
326 from each action a . The REINFORCE loss is:

$$\ell_{\theta} = -r \cdot \log p_{\theta}(a|x_{\text{action}}) \quad (12) \quad 327$$

328 This gradient increases the probability of actions with posi-
329 tive rewards and decreases the probability of actions with
330 negative rewards. We accumulate gradients over all exam-
331 ples in $\mathcal{D}_{\text{train}}$ over several RL_steps:

$$\theta^{(s+1)} = \theta^{(s)} - \eta_{\theta} \nabla_{\theta} \sum_{(x_{\text{action}}, f, s_r) \in \mathcal{D}_{\text{train}}} \ell_{\theta}(x_{\text{action}}, f, s_r) \quad (13) \quad 332$$

333 4. Experiments on Long-Horizon Household

334 Tasks

335 4.1. Long-Horizon Household Task Construction

336 To evaluate our framework on tasks requiring multi-step
337 reasoning and failure recovery, we construct Long-Horizon
338 Household Tasks based on the BEHAVIOR-1K [18] envi-
339 ronments. Drawing inspiration from household scenarios,
340 we define four core task categories: (1) **Fitting**, where ob-
341 jects must be packed or placed into constrained contain-
342 ers or surfaces, stressing geometry, capacity and occlu-
343 sion failures; (2) **Selection**, where the agent compares and
344 retrieves the most suitable item (e.g., in terms of prefer-
345 ences or sizes). Failures occur when choices prove subopti-
346 mal upon discovering better alternatives; (3) **Preparation**,
347 where tasks require sequential constraints and dependencies
348 (e.g., assembling or nested placement), stressing sequen-
349 tial dependency and non-local failures; and (4) **Hybrid**, where
350 multiple modes appear within a single episode, stressing
351 mixed spatial, relational, and occlusion failures.

	Baselines						Ablations of Reflective Test-Time Planning (Ours)					
	Reflection	Self-Refine	ReflectVLM	PPO	DreamerV3	3DLLM-Mem	w/o RIA	w/o ROA	w/o RIA	w/o ROA	w/o Act. Loss	w/o Int. Loss
Fitting	8.51%	10.6%	2.12%	0%	4.26%	10.6%	0%	33.5%	6.38%	25.5%	12.8%	44.7%
Selection	8.82%	11.8%	5.88%	2.94%	11.8%	14.7%	17.6%	5.88%	11.8%	26.5%	8.82%	32.4%
Preparation	15.9%	12.7%	14.3%	7.94%	11.1%	9.52%	11.1%	3.17%	19.0%	20.6%	17.5%	31.7%
Hybrid	6.45%	9.68%	6.45%	3.23%	12.9%	9.68%	12.9%	3.23%	12.9%	16.1%	9.68%	25.8%
Average	9.92%	11.20%	7.19%	3.53%	10.02%	11.13%	10.40%	11.45%	12.52%	22.18%	12.20%	33.65%

Table 1. Baseline comparisons and ablations of Reflective Test-Time Planning on Long-Horizon Household Tasks. “RIA” is short for Reflection-in-action, “ROA” for reflection-on-action. Act. Loss means action model loss; Int. Loss means internal reflection model loss.

352 Task Generation & Execution. We employ GPT-5
353 to generate task specifications through a carefully struc-
354 tured prompting procedure, adapting existing task templates
355 that emphasize long-horizon reasoning and failure recovery
356 from the original BEHAVIOR-1K benchmark. Each gener-
357 ated task instance includes: (1) a task description, (2)
358 3-7 relevant rooms, (3) new objects with placement spec-
359 ifications, and (4) a complete trajectory with interleaved
360 thoughts, actions, and reflections scores. Since we include
361 scene graphs in the prompts to GPT-5, which provides ob-
362 ject properties such as bounding boxes, GPT-5 can deduce
363 potential failures (e.g., size mismatches) during generation.
364 We further execute the tasks generated by GPT-5 in BE-
365 HAVIOR simulators to ensure consistency with the actual
366 scene dynamics and provide ground-truth data for finetun-
367 ing embodied LLMs. We initialize the environment by
368 loading the BEHAVIOR-1K scenes, placing new objects at
369 designated locations, and positioning the robot at a default
370 starting pose. At each step, the agent executes the given
371 actions. After every interaction, the system captures RGB-
372 D observations, converted to point clouds which serve as
373 the inputs to our Embodied LLMs. The simulator then per-
374 forms physics-based verification and provides execution re-
375 sults that could be used for prompting external reflections.
376 We also perform task verification checks by comparing the
377 task execution results with GPT-5’s generated task speci-
378 fications (e.g., assumed task failures). Please refer to Ap-
379 pendix E for more details about data generation, models and
380 experiments.

381 Finetuning & Evaluation Sets. Each GPT-generated
382 and executed trajectory yields (observation, action, reflec-
383 tion, score) tuples extracted at corresponding timesteps,
384 which form supervised fine-tuning pairs to initialize our
385 three embodied LLMs with basic task competency before
386 test-time deployment. During evaluation, the agent is given
387 only the scene configuration and task description, with
388 all other components initialized as None, and must au-
389 tonomously generate and execute its trajectory step by step.
390 To ensure rigorous evaluation, finetuning and evaluation
391 sets have no overlap in task descriptions, scene configura-
392 tions, or object placements. We evaluate using task success
393 rate: the percentage of tasks where the agent completes the
394 objective within the action budget. A task succeeds if all
395 required objects reach target locations and constraints are
396 satisfied.

4.2. Model Architecture & Implementation Details

397 We build a 3D-LLM based on LLaVA-3D [40]. During
398 supervised fine-tuning, we train a single unified model on
399 all three modes (action generation, internal reflection, ex-
400 ternal reflection) for cross-mode learning. We fine-tune
401 LLaVA-3D-7B with learning rate 2×10^{-5} . At test-time,
402 we instantiate three copies of this model as π_θ , V_{ϕ_i} , and
403 V_{ϕ_e} , where π_θ and V_{ϕ_i} are updated via test-time training.
404 Following Hu et al. [12], we maintain fused observations
405 from previous-step point clouds. During test-time train-
406 ing, we apply LoRA (rank 4, alpha 8, dropout 0.15) tar-
407 geting `q_proj` and `v_proj` layers. The internal LLM
408 uses supervised learning (learning rate 5×10^{-3} , $E = 5$
409 epochs); the action LLM uses REINFORCE (learning rate
410 1×10^{-3} , 3 RL.Steps). We use gradient clipping (0.3 for su-
411 pervised, 0.5 for RL) to prevent catastrophic forgetting. For
412 reflection-in-action, we set $N = 4$ candidates and temper-
413 ature $T = 2.0$. We adopt four categories of baselines: (1)
414 verbal-reflection [21, 30]; (2) world-model based reflection
415 [5]; (3) RL-based baselines: Vanilla PPO [29] and world-
416 model based RL (DreamerV3 [7]) (4) 3D-LLM with inter-
417 action context [12].
418

4.3. Experimental Results & Analysis

419 Table 1 shows that our model achieves substantial improve-
420 ments over both prior methods and ablated versions across
421 all task categories. The improvement is most prominent on
422 Fitting tasks, where our approach reaches 44.7% success
423 rate compared to 10.6% for the strongest baseline (3DLLM-
424 Mem), 2.1% for ReflectVLM and 0% for PPO. Fitting tasks
425 benefit most from our method because they impose tight
426 spatial constraints that require iterative refinement through
427 trial and error—agents must reason about 3D geometry
428 from multiple attempts and continuously adjust their un-
429 derstanding based on execution feedback. The ablation
430 studies reveal that reflection-in-action (RIA) and reflection-
431 on-action (ROA) are mutually dependent. Removing ei-
432 ther causes performance degradation. Sometimes removing
433 just one component performs worse than removing both.
434 For instance, without RIA, Preparation falls to 3.17% and
435 Hybrid to 3.23%, while removing both components yields
436 11.1% and 12.9% respectively on these tasks. This is be-
437 cause RIA without ROA produces overconfident yet inaccur-
438 ate action scores with no hindsight correction mechanism,
439 while ROA without RIA wastes learning on poorly cho-
440

Baselines: [1] Reflexion [2] Self-refine [3] ReflectVLM [4] PPO [5] DreamerV3 [6] Qwen w/ Memory Context
 Ablations-w/o Ext. Reflection: [7] Ours w/o RIA w/o ROA w/o External Reflection [8] Ours w/o RIA w/o External Reflection
 Ablations-w/o RIA and/or ROA: [9] Ours w/o RIA w/o ROA [10] Ours w/o ROA [11] Ours w/o RIA (ROA: Base-Weight) [12] Ours w/o RIA (ROA: LoRA)
 Ablations-w/o Action loss: [13] Ours w/o Action loss (ROA: Base-Weight) [14] Ours w/o Action loss (ROA: LoRA)
 Full Model Variants: [15] Ours (ROA: Base-Weight) [16] Ours (ROA: LoRA)

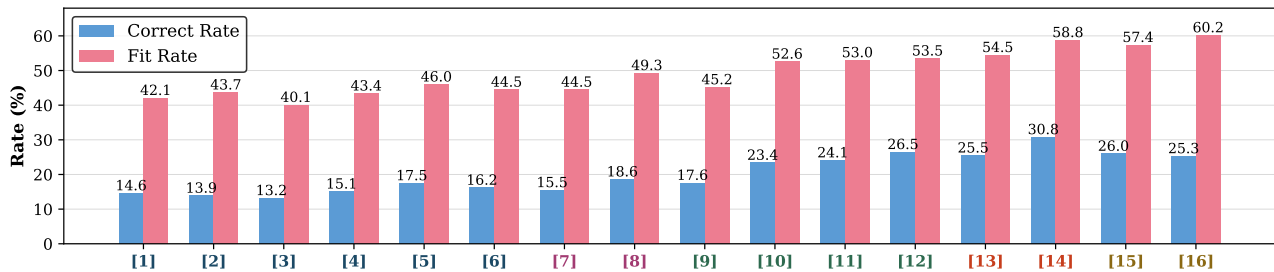


Figure 3. **Cupboard Fitting results.** Blue bars show correct-placement rate, pink bars show fit rate. RIA means Reflection-in-action; ROA means Reflection-on-action. “W/o external reflection” means that we don’t use external reflection as the input to the action generation LLM. We implement two test-time training variants for ROA: one is test-time training on all base weights; and the other is test-time training on LoRA parameters only. Reflective Test-Time Planning significantly improves both success metrics.

441 sen actions that fail to reveal true scene affordances. When
 442 combined, they form a virtuous cycle: best-of-N selection
 443 yields higher-quality trajectories for learning, and test-time
 444 training refines the internal model, leading to better future
 445 action selection. The loss ablations similarly demon-
 446 strate mutual dependency between the action loss and inter-
 447 nal reflection loss, which together constitute ROA. Remov-
 448 ing either loss causes performance degradation, and some-
 449 times removing just one loss performs worse than removing
 450 both (w/o ROA). Without the internal loss, Hybrid drops to
 451 9.68% compared to 12.9% for w/o ROA, and Selection falls
 452 to 8.82% compared to 11.8% for w/o ROA. This demon-
 453 strates that policy-gradient updates and supervised reflec-
 454 tion both contribute essential signals for effective adapta-
 455 tion. The qualitative example in Figure 5(a) demonstrates
 456 that our model supports continual learning and active per-
 457 ception driven by accumulated scene experience. The test-
 458 time cost analysis and compute-matched experiment in Ap-
 459 pendix A shows that even with substantially more test-time
 460 compute, the baselines continue to repeat failures and their
 461 performance even degrades. Appendix B also shows that
 462 the reflection mechanisms enable the agents to generalize
 463 to novel embodied environments like HM3D [27].

464 5. Experiments on the Cupboard Fitting Task

465 **Cupboard Fitting Task Design.** The Cupboard Fitting task
 466 serves as a complementary benchmark to Long-Horizon
 467 Household Tasks. While BEHAVIOR provides realistic
 468 multi-room environments with complex semantic reasoning
 469 challenges, it also introduces environmental uncertainties
 470 that make it difficult to isolate learning mechanisms. We
 471 design Cupboard Fitting as a controlled MuJoCo environ-
 472 ment where agents learn from placement failures, enabling
 473 precise measurement of reflective test-time training me-
 474 chanisms. The task environment consists of three key components.



Figure 4. Examples of the Cupboard Fitting Task.

476 First, a multi-compartment cupboard structure with 6-8
 477 compartments of varying sizes and colors, each defined by
 478 precise 3D bounding boxes. Second, a set of 6-10 colored
 479 geometric objects that must be placed into compartments.
 480 Third, a Franka Panda robotic arm controlled through high-
 481 level natural language commands such as “pick up the red
 482 apple” or “put the blue object in the green compartment.”
 483 The agent’s objective is to place all objects into the cup-
 484 board such that each fits completely within compartment
 485 boundaries, while multiple objects may share compartments
 486 when space permits. Each task has only one valid solution
 487 where all objects fit within their designated compartments.
 488 Success requires reasoning about object-compartment com-
 489 patibility, multi-object spatial packing, and long-horizon
 490 dependencies where early placement decisions affect later
 491 possibilities. After each action, the environment provides
 492 execution status through forward simulation, and updated
 493 visual observations. We define two evaluation metrics: cor-
 494 rect rate measures the percentage of objects placed in their
 495 correct target compartments, while fit rate measures the per-
 496 centage of objects successfully placed in any compartment.

497 **Model Architecture & Implementation Details.** We
 498 build upon Qwen2.5-VL-3B with unified supervised fine-
 499 tuning (batch size 128, learning rate 1×10^{-5} , weight decay
 500 0.1) on action generation, internal reflection, and external
 501 reflection. For test-time training, we implement two vari-
 502 ants: (1) *Base-Weight*: updates all non-visual parameters
 503 with SGD (learning rate 1×10^{-3} for action model, 5×10^{-5}

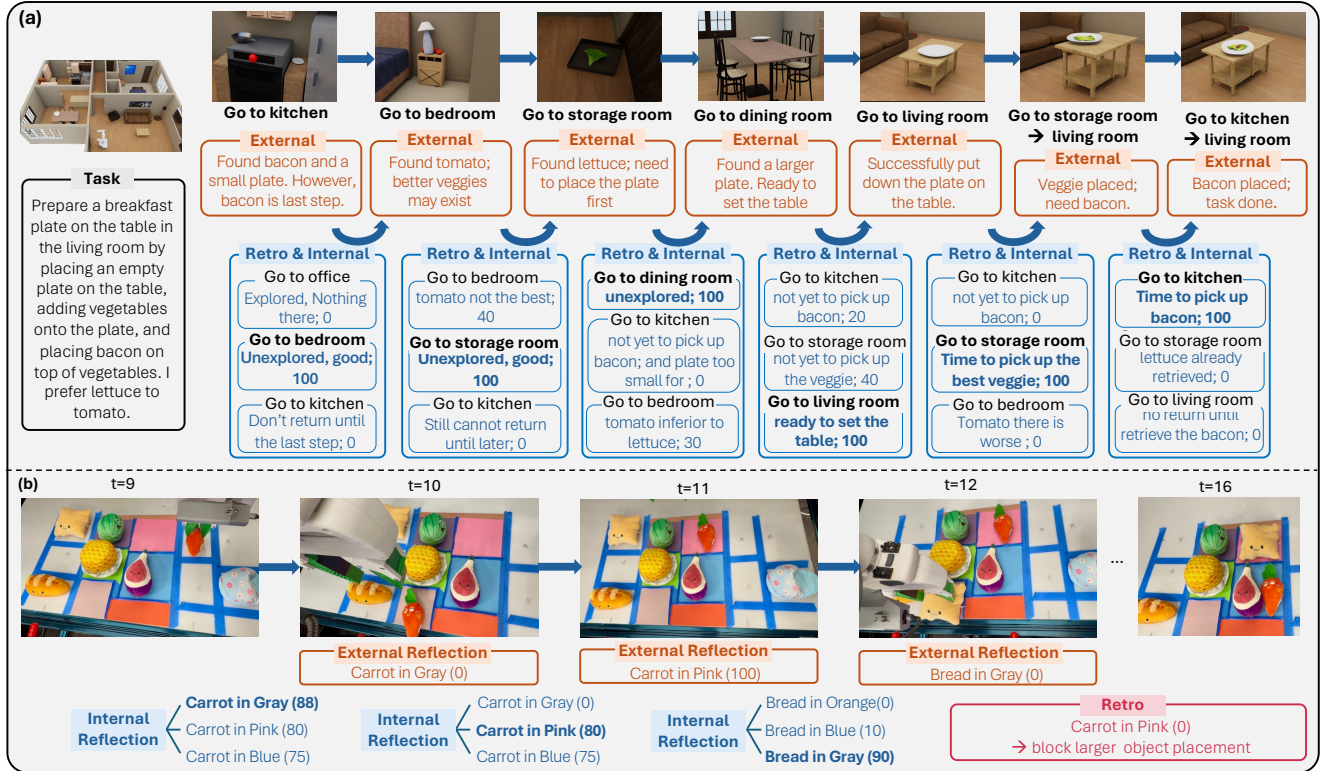


Figure 5. **Qualitative Examples.** Steps and reflections are simplified for better presentations. Blue text shows internal reflection used for candidate selection, orange text shows external reflection after execution, and red text suggests retrospective reflection. (a) Long-Horizon Household example. We use retro & internal because the generated retro reflection is also used to train the internal model. (b) Real-robot Cupboard Fitting example. We put reflection scores inside brackets, omit detailed reflections and only present the scores for simplicity.

504 for internal model); (2) *LoRA*: applies LoRA (rank 8, alpha
505 16, dropout 0.1) targeting all linear layers except `lm_head`,
506 `embed_tokens`, and visual encoder, using SGD with learning
507 rate 0.01 for action model and 0.2 for internal model.
508 Both train for 3 epochs. We limit episodes to 50 maximum
509 steps and update the model whenever an execution failure
510 occur. Hyperparameter selection ablations can be found in
511 Appendix C. We use the same baselines as the Long-Horizon
512 Household Task except that we replace 3D-LLM
513 with Qwen for the baseline model with context memory.

514 **Experimental Results & Analysis.** Figure 3 shows that
515 our full method with both reflection-in-action (RIA) and
516 reflection-on-action (ROA) using LoRA-based test-time
517 training achieves 60.2% fit rate and 25.3% correct rate, sub-
518 stantially outperforming LLM-based reflection baselines,
519 RL-based baselines and memory context baselines. Ablation
520 studies reveal that both reflection mechanisms are es-
521 sential: removing RIA drops performance to 53.5% fit rate,
522 while removing ROA reduces it to 45.2% fit rate. Removing
523 both mechanisms along with external reflection further
524 degrades performance to 44.5% fit rate. We also find that
525 LoRA-based test-time training (60.2% fit rate) performs
526 comparably to full base-weight training (57.4% fit rate)
527 while being more parameter-efficient. These results confirm

528 that our unified framework of reflection-in-action for candi-
529 date selection and reflection-on-action for test-time adapta-
530 tion enables effective learning from failures during deploy-
531 ment. In Figure 5(b), we show a qualitative result where we
532 generalize our model in the real-robot setting. In these tri-
533 als, reflective adaptation enables the robot to recover from
534 execution failures, avoid repeated placement errors, and
535 correct earlier decisions through retrospective reflection, il-
536 lustrating that the learned behaviors transfer to the physical
537 world with satisfying generalization abilities. Appendix D
538 also shows that our method enables long-horizon planning
539 by retrospective reflection that outperforms Receding Hor-
540 izon Planning while saving 5x test-time compute.

541 6. Conclusion

542 This paper introduces Reflective Test-Time Planning, which
543 couples reflection-in-action for pre-action evaluation with
544 reflection-on-action for post-execution assessments. Eval-
545 uations across two newly-designed embodied tasks demon-
546 strate strong gains and highlight the complementary roles of
547 the two reflection modes. Future work may extend reflective
548 adaptation to richer sensory modalities (*e.g.*, tactile).

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
594
595
596
597
598
599
600
601
602
603
604
605

References

- [1] Chris Argyris. Double loop learning in organizations. *Harvard Business Review*, 1977. 2, 3
- [2] Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar, Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Pete Florence. Palm-e: An embodied multimodal language model. In *arXiv preprint arXiv:2303.03378*, 2023. 2
- [3] Sri Harsha Dumpala, Chandramouli Sastry, and Sageev Oore. Test-time training for speech, 2023. 2
- [4] Cunxin Fan, Xiaosong Jia, Yihang Sun, Yixiao Wang, Jianglan Wei, Ziyang Gong, Xiangyu Zhao, Masayoshi Tomizuka, Xue Yang, Junchi Yan, and Mingyu Ding. Interleave-vla: Enhancing robot manipulation with interleaved image-text instructions. 2025. 2
- [5] Yunhai Feng, Jiaming Han, Zhuoran Yang, Xiangyu Yue, Sergey Levine, and Jianlan Luo. Reflective planning: Vision-language models for multi-stage long-horizon robotic manipulation, 2025. 2, 6, 9
- [6] Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujie Yang, Nan Duan, and Weizhu Chen. CRITIC: Large language models can self-correct with tool-interactive critiquing. *arXiv preprint arXiv:2305.11738*, 2024. 2
- [7] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models, 2024. 2, 6
- [8] Elvin Hajizada, Balachandran Swaminathan, and Yulia Sandamirskaya. Continual learning for autonomous robots: A prototype-based approach. 2024. 2
- [9] 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. In *Advances in Neural Information Processing Systems*, 2023. 2
- [10] Yining Hong, Zishuo Zheng, Peihao Chen, Yian Wang, Junyan Li, and Chuang Gan. MultiPLY: A multisensory object-centric embodied large language model in 3D world. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26406–26416, 2024. 1, 2
- [11] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 2
- [12] Wenbo Hu, Yining Hong, Yanjun Wang, Leison Gao, Zibu Wei, Xingcheng Yao, Nanyun Peng, Yonatan Bitton, Idan Szepes, and Kai-Wei Chang. 3dllm-mem: Long-term spatial-temporal memory for embodied 3d large language model, 2025. 2, 4, 6, 9
- [13] Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, Pierre Sermanet, Noah Brown, Tomas Jackson, Linda Luu, Sergey Levine, Karol Hausman, and Brian Ichter. Inner monologue: Embodied reasoning through planning with language models. In *arXiv preprint arXiv:2207.05608*, 2022. 2
- [14] Isaac Kauvar, Chris Doyle, Linqi Zhou, and Nick Haber. Curious replay for model-based adaptation. *International Conference on Machine Learning*, 2024. 2
- [15] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Perez Foster, Pannag Raj Sanketi, Quan Vuong, et al. OpenVLA: An open-source vision-language-action model. In *8th Annual Conference on Robot Learning*, 2024. 1, 2
- [16] Yuto Kojima, Jiarui Xu, Xueyan Zou, and Xiaolong Wang. Lora-ttt: Low-rank test-time training for vision-language models. 2025. 2
- [17] Timothée Lesort, Vincenzo Lomonaco, Andrei Stoian, Davide Maltoni, David Filliat, and Natalia Díaz-Rodríguez. Continual learning for robotics: Definition, framework, learning strategies, opportunities and challenges. *Information Fusion*, 58:52–68, 2020. 2
- [18] Chengshu Li, Ruohan Zhang, Josiah Wong, Cem Gokmen, Sanjana Srivastava, Roberto Martín-Martín, Chen Wang, Gabriel Levine, Wensi Ai, Benjamin Martinez, Hang Yin, Michael Lingelbach, Minjune Hwang, Ayano Hiranaka, Sujay Garlanka, Arman Aydin, Sharon Lee, Jiankai Sun, Mona Anvari, Manasi Sharma, Dhruva Bansal, Samuel Hunter, Kyu-Young Kim, Alan Lou, Caleb R Matthews, Ivan Villal-Renteria, Jerry Huayang Tang, Claire Tang, Fei Xia, Yunzhu Li, Silvio Savarese, Hyowon Gweon, C. Karen Liu, Jiajun Wu, and Li Fei-Fei. Behavior-1k: A human-centered, embodied ai benchmark with 1,000 everyday activities and realistic simulation. *arXiv preprint arXiv:2403.09227*, 2024. 5, 6
- [19] Chengmeng Li, Junjie Wen, Yan Peng, Yaxin Peng, Feifei Feng, and Yichen Zhu. Pointvla: Injecting the 3d world into vision-language-action models. 2025. 2
- [20] Jian Liang, Ran He, and Tieniu Tan. A comprehensive survey on test-time adaptation under distribution shifts. *International Journal of Computer Vision*, 133(1):31–64, 2025. 2
- [21] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36, 2023. 1, 2, 6, 9
- [22] Russell Mendonca, Emmanuel Panov, Bernadette Bucher, Jiguang Wang, and Deepak Pathak. Continuously improving mobile manipulation with autonomous real-world RL. In *Conference on Robot Learning*, 2024. 2
- [23] Yuan Meng, Zhenshan Bing, Xiangtong Yao, Kejia Chen, Kai Huang, Yang Gao, Fuchun Sun, and Alois Knoll. Preserving and combining knowledge in robotic lifelong reinforcement learning. *Nature Machine Intelligence*, pages 1–14, 2025. 2
- [24] Yao Mu, Qinglong Zhang, Mengkang Hu, Wenhai Wang, Mingyu Ding, Jun Jin, Bin Wang, Jifeng Dai, Yu Qiao, and Ping Luo. Embodiedgpt: Vision-language pre-training via embodied chain of thought, 2023. 2
- [25] Andrew Ng. Agentic design patterns part 2: Reflection. *DeepLearning.AI The Batch*, 2024. 2
- [26] Shuaicheng Niu, Jiayang Wu, Yifan Zhang, Yaofu Chen, Shuai Zheng, Peilin Zhao, and Mingkui Tan. Efficient

- 664 test-time model adaptation without forgetting. In *International Conference on Machine Learning*, pages 16888–
665 16905, 2022. 2
- 667 [27] Santhosh K. Ramakrishnan, Aaron Gokaslan, Erik Wijmans,
668 Oleksandr Maksymets, Alex Clegg, John Turner, Eric Un-
669 dersander, Wojciech Galuba, Andrew Westbury, Angel X.
670 Chang, Manolis Savva, Yili Zhao, and Dhruv Batra. Habitat-
671 matterport 3d dataset (hm3d): 1000 large-scale 3d environ-
672 ments for embodied ai, 2021. 7, 1
- 673 [28] Donald A. Schön. *The Reflective Practitioner: How Profes-*
674 *sionals Think in Action*. Basic Books, New York, NY, USA,
675 1992. 1
- 676 [29] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Rad-
677 ford, and Oleg Klimov. Proximal policy optimization algo-
678 rithms, 2017. 6, 9
- 679 [30] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik
680 Narasimhan, and Shunyu Yao. Reflexion: Language agents
681 with verbal reinforcement learning. *Advances in Neural In-*
682 *formation Processing Systems*, 36, 2023. 1, 2, 6, 9
- 683 [31] Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei A
684 Efros, and Moritz Hardt. Test-time training with self-
685 supervision for generalization under distribution shifts. In *In-*
686 *ternational Conference on Machine Learning*, pages 9229–
687 9248, 2020. 2
- 688 [32] Yu Sun, Xinhao Li, Karan Dalal, Jiarui Xu, Arjun Vikram,
689 Genghan Zhang, Yann Dubois, Xinlei Chen, Xiaolong Wang,
690 Sanmi Koyejo, Tatsunori Hashimoto, and Carlos Guestrin.
691 Learning to (learn at test time): RNNs with expressive hid-
692 den states. In *International Conference on Machine Learn-*
693 *ing*, 2024. 2
- 694 [33] Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Ol-
695 shausen, and Trevor Darrell. Tent: Fully test-time adapta-
696 tion by entropy minimization. In *International Conference on*
697 *Learning Representations*, 2021. 2
- 698 [34] Hao Yang, Min Wang, Jinshen Jiang, and Yun Zhou. To-
699 wards test time adaptation via calibrated entropy minimiza-
700 tion. In *Proceedings of the 30th ACM SIGKDD Con-*
701 *ference on Knowledge Discovery and Data Mining*, page
702 3736–3746, New York, NY, USA, 2024. Association for
703 Computing Machinery. 2
- 704 [35] Wenqi Zhang, Yongliang Shen, Linjuan Wu, Qiuying Peng,
705 Jun Wang, Yueting Zhuang, and Weiming Lu. Self-contrast:
706 Better reflection through inconsistent solving perspectives.
707 2024. 2
- 708 [36] Qingqing Zhao, Yao Lu, Moo Jin Kim, Zipeng Fu, Zhuoyang
709 Zhang, Yecheng Wu, Zhaoshuo Li, Qianli Ma, Song Han,
710 Chelsea Finn, et al. Cot-vla: Visual chain-of-thought rea-
711 soning for vision-language-action models. In *Proceedings*
712 *of the Computer Vision and Pattern Recognition Conference*,
713 pages 1702–1713, 2025. 2
- 714 [37] Haoyu Zhen, Xiaowen Qiu, Peihao Chen, Jincheng Yang,
715 Xin Yan, Yilun Du, Yining Hong, and Chuang Gan. 3D-
716 VLA: A 3D vision-language-action generative world model.
717 *arXiv preprint arXiv:2403.09631*, 2024. 2
- 718 [38] H. Zhen et al. Learning 4d embodied world models. In
719 *Proceedings of the IEEE/CVF International Conference on*
720 *Computer Vision (ICCV)*, 2025. 2
- [39] Siyuan Zhou, Yilun Du, Jiaben Chen, Yandong Li, Dit-Yan
Yeung, and Chuang Gan. Robodreamer: Learning composi-
tional world models for robot imagination, 2024. 9
- [40] Chenming Zhu, Tai Wang, Wenwei Zhang, Jiangmiao Pang,
and Xihui Liu. LLaVA-3D: A simple yet effective pathway
to empowering LMMs with 3D-awareness. In *International*
Conference on Computer Vision, 2025. 2, 6, 8
- [41] Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted
Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan Welker,
Ayzaan Wahid, et al. Rt-2: Vision-language-action models
transfer web knowledge to robotic control. In *Conference on*
Robot Learning, pages 2165–2183. PMLR, 2023. 1, 2

733 A. Test-Time Cost Analysis

734 A.1. Computational Comparison

735 Our full model introduces two additional sources of infer-
736 ence time cost: (i) candidate action sampling and internal
737 reflection scoring for reflection-in-action (RIA), and (ii) pe-
738 riodic reflection-on-action(ROA)-based LoRA test-time up-
739 dates.

740 On average across Long-Horizon Household Tasks and
741 Cupboard Fitting Tasks, we observe a $\sim 3\times$ increase in
742 per-step wall-clock time compared to the vanilla baseline:

$$743 \text{Time}_{\text{full}} \approx 3.0 \times \text{Time}_{\text{no-RIA/ROA}}.$$

744 Importantly, this latency is incurred at deployment and does
745 not require additional supervised data or environment roll-
746 outs beyond normal task execution.

747 A.2. Why the Overhead is Justified

748 Despite the additional test-time latency, the reflective over-
749 head is justified for three reasons:

750 **Deployment-Oriented Adaptation.** The extra time cost
751 occurs at deployment rather than pretraining. This aligns
752 with realistic embodied settings where robots adapt online
753 rather than relying on costly retraining cycles.

754 **Reduction of Execution Waste.** Vanilla agents fre-
755 quently repeat failures (e.g., placing incompatible objects,
756 revisiting rooms without intent). RIA reduces execution
757 waste by filtering poor actions before execution, while ROA
758 eliminates repeated failures through hindsight-driven up-
759 dates. A smaller number of higher-quality actions amortizes
760 the reflection cost.

761 **Conversion of Time into Learning.** Extra time is not
762 spent on naive rollouts, but on improving internal models.
763 RIA and ROA produce persistent behavioral improvements,
764 whereas baselines spend time without updating policies or
765 reasoning mechanisms. In short, reflective time is struc-
766 turally more valuable than mere rollout time.

767 A.3. Compute-Matched Experiment

768 A natural concern is whether the performance gap is merely
769 due to increased wall-clock time. To evaluate this, we con-
770 struct a time-matched variant where the vanilla baseline re-
771 ceives a $3\times$ steps budget, matching the approximate infer-
772 ence time of our full model:

$$773 \text{Steps}_{\text{baseline}}^{\text{expanded}} \approx 3 \times \text{Steps}_{\text{baseline}}.$$

774 From Table 2, we observe that even under a tripled time
775 budget, the baseline:

	w/o RIA w/o ROA		Ours (Full)
	Vanilla	Same Time Budget	
Fitting	0.00%	0.00%	44.7%
Selection	17.6%	14.7%	32.4%
Preparation	11.1%	12.7%	31.7%
Hybrid	6.45%	6.45%	25.8%
Average	8.79%	8.46%	33.65%

Table 2. Performance comparison between (1) vanilla ablation without RIA or ROA, (2) vanilla baseline with matched time budget ($3\times$ steps), and (3) our full reflective model on Long-Horizon Household Tasks.

- 776 • fails to correct early suboptimal decisions, 776
- 777 • frequently revisits states without strategic change, 777
- 778 • exhibits repeated placement/navigation failures, 778
- 779 • plateaus significantly below our full model on all bench- 779
- 780 marks and does not improve over the vanilla baseline. 780

781 This result supports that reflective time is not equivalent
782 to naive rollout time expansion: reflective updates change
783 the decision process itself, while rollout expansion merely
784 increases trajectory length without improving competence
785 or hindsight reasoning.

786 **Conclusion.** These findings indicate that although our
787 method incurs a $\sim 3\times$ test-time latency overhead, reflective
788 computation provides unique adaptation benefits that can-
789 not be recovered by proportional rollout-step scaling. The
790 overhead is therefore practical and justified in embodied de-
791 ployment settings.

792 B. Generalization to Habitat-Matterport 3D 793 Scenes

794 B.1. Experimental Setup

795 To evaluate the generalization capacity of our Reflec-
796 tive Test-Time Planning framework, we conduct addi-
797 tional experiments on the Habitat-Matterport 3D (HM3D)
798 dataset [27], which provides photorealistic 3D recon-
799 structions of real-world indoor environments. Unlike
800 BEHAVIOR-1K scenes which are primarily synthetic
801 household environments, HM3D offers diverse real-world
802 scenes with different spatial layouts, object distributions,
803 and visual appearances, providing a challenging domain
804 shift for testing our method’s transferability.

805 We focus specifically on **Preparation tasks**, which in-
806 volve sequential constraints and dependencies where ac-
807 tions must occur in specific orders. These tasks are particu-
808 larly challenging in HM3D environments due to: (1) more
809 complex spatial layouts with irregular room configurations,
810 (2) diverse object appearances and placements not seen dur-
811 ing training, (3) ambiguous spatial relationships that require

	Baselines						Ablations of Reflective Test-Time Planning (Ours)					Ours
	Reflexion	Self-Refine	ReflectVLM	PPO	DreamerV3	3DLLM-Mem	w/o RIA	w/o ROA	w/o ROA	w/o RIA	w/o Act. Loss	
Preparation	2.44%	4.88%	0%	0%	2.44%	7.32%	0%	7.32%	2.44%	14.6%	9.76%	19.5%

Table 3. Generalization results on Habitat-Matterport 3D (HM3D) environment for Preparation tasks (41 test cases). The significant performance drop compared to BEHAVIOR Long-Horizon Household Tasks (31.7% \rightarrow 14.6%) demonstrates the substantial domain gap. However, our reflective framework maintains relative advantages over baselines, with several methods (ReflectVLM, PPO, w/o RIA w/o ROA, w/o RIA) achieving 0% success rate. “RIA” is short for Reflection-in-action, “ROA” is short for reflection-on-action. Act. Loss means action model loss. Int. Loss means internal reflection model loss.

812 active exploration to resolve, and (4) longer navigation dis- 853
813 tances between task-relevant objects. 854

814 B.2. Task Construction and Adaptation 855

815 We adapt our task generation pipeline to HM3D scenes 856
816 while maintaining the same core task structure. We select 857
817 41 preparation task instances across diverse HM3D scenes, 858
818 ensuring coverage of various spatial configurations and ob- 859
819 ject arrangements. Each task requires the agent to: (1) navi- 860
820 gate through multiple rooms to locate task-relevant objects, 861
821 (2) retrieve objects in the correct sequential order, (3) per- 862
822 form placement or assembly actions with proper dependen- 863
823 cies, and (4) handle spatial constraints specific to real-world 864
824 scene layouts. 865

825 The key difference from Long-Horizon Household eval- 866
826 uation is the domain gap: our models are trained exclusively 867
827 on BEHAVIOR-1K synthetic scenes and must generalize to 868
828 HM3D’s photorealistic environments at test time. This tests 869
829 whether reflection-based adaptation can overcome distribu- 870
830 tion shift through deployment-time learning. 871

831 B.3. Implementation Details 872

832 We use the same model architecture (LLaVA-3D-7B) 873
833 and test-time training configuration as the Long-Horizon 874
834 Household experiments, with no additional fine-tuning on 875
835 HM3D scenes. Point cloud observations are extracted from 876
836 RGB-D sensors in the same manner, and the reflection gen- 877
837 eration prompts remain unchanged. This zero-shot transfer 878
838 setup provides a rigorous test of whether reflective mech- 879
839 anisms learned in synthetic environments transfer to real- 880
840 world scene understanding. 881

841 B.4. Results and Analysis 882

842 Table 3 presents the generalization results on HM3D Prepa- 883
843 ration tasks. Our method demonstrates robust generaliza- 884
844 tion to photorealistic real-world environments, achieving 885
845 19.5% success rate despite being trained exclusively on 886
846 synthetic BEHAVIOR-1K scenes. While this represents a 887
847 12.2 percentage point performance drop compared to Long- 888
848 Horizon Household Tasks (31.7% \rightarrow 19.5%), our reflect- 889
849 ive framework maintains over 60% of its original perfor- 890
850 mance and substantially outperforms all baselines, demon- 891
851 strating the effectiveness of reflection-based adaptation un- 892
852 der domain shift. The strongest baseline, 3DLLM-Mem, 893

853 achieves only 7.32%, while several methods completely 854
855 collapse to 0% (ReflectVLM, PPO, w/o RIA w/o ROA). 856
857 This stark contrast highlights that our reflection mecha- 858
859 nisms—particularly the combination of reflection-in-action 860
861 and reflection-on-action—provide critical robustness when 862
863 facing distribution shift from synthetic to photorealistic en- 864
865 vironments. The fact that our method maintains substantial 866
867 relative advantages over baselines suggests that the core re- 868
869 flective mechanisms—generating diverse candidates, scor- 870
871 ing them through internal simulation, learning from execu- 872
873 tion feedback, and updating both action and reflection mod- 874
875 els—transfer effectively across domain boundaries and pro- 876
877 vide a robust foundation for embodied agents generalizing 878
879 to novel environments. 880

867 C. Hyperparameter Analyses for Cupboard 868 869 Fitting Task 870

869 We conduct ablation studies on four key hyperparameters 870
871 that govern the reflection mechanism and show the results 872
873 in Figure 6. 874

872 C.1. Test-Time Scaling: Number of Candidate Ac- 873 874 tions 875

874 The reflection-in-action mechanism generates N candidate 875
876 actions via sampling, scores each using the internal reflec- 877
878 tion model V_{ϕ_i} , and executes the highest-scoring action. We 879
880 vary $N \in \{1, 2, 3, 4, 5, 6, 8, 9, 10\}$ to measure the impact of 881
882 candidate diversity on task performance (Figure 6, top left). 883

879 **Results.** Performance improves monotonically from 880
881 53.0% ($N=1$, greedy decoding) to 60.0% ($N=6$), represent- 882
883 ing a 7 percentage point gain. This demonstrates that inter- 884
885 nal reflection scoring effectively identifies superior actions 886
887 from diverse candidate pools. Beyond $N=6$, performance 888
889 plateaus and slightly decreases to 58.8% at $N=10$, suggest- 890
891 ing diminishing returns from excessive candidates—likely 892
893 due to increased computational cost without proportional 894
895 quality improvements in the candidate pool. 896

888 **Analysis.** The optimal operating point at $N=6$ bal- 889
890 ances exploration breadth with computational efficiency. At 891
892 $N=1$, the model lacks opportunities to reconsider subopti- 893
894 mal greedy choices. At $N=2-5$, expanding the candidate 895
896 pool allows the internal reflection model to compare al- 897
898 ternatives and avoid locally optimal but globally poor ac- 899

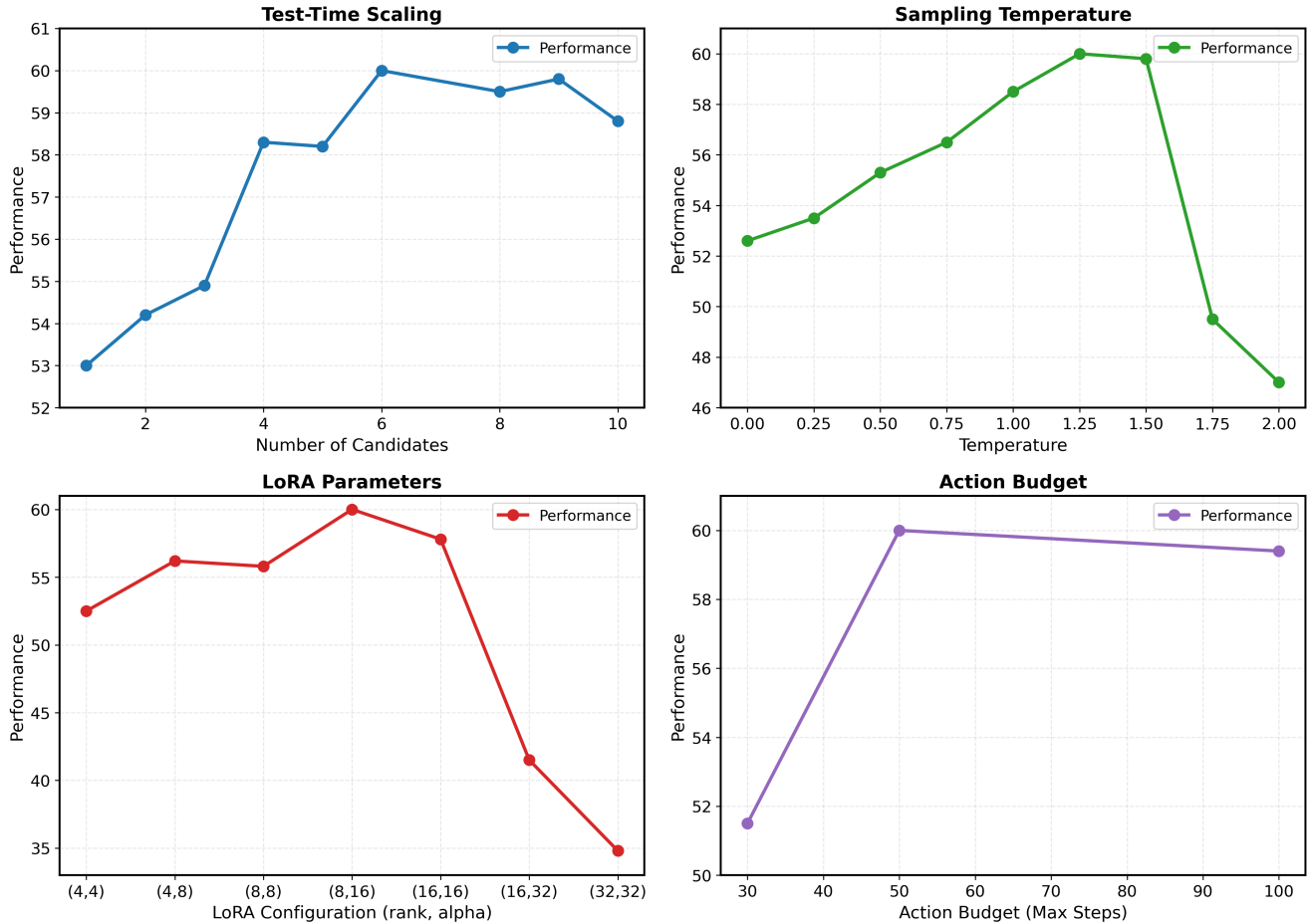


Figure 6. **Hyperparameter ablation studies on Cupboard Fitting.** *Top Left:* Performance vs. number of candidate actions. Peak performance (60.0%) occurs at $N=6$ candidates, demonstrating that internal reflection effectively identifies superior actions from diverse pools. Beyond $N=6$, performance plateaus as excessive candidates add computational cost without improving the best candidate quality. *Top Right:* Performance vs. sampling temperature. Optimal temperature range ($T=1.25-1.5$) balances candidate diversity with quality—temperatures below 0.5 produce overly similar candidates that limit reflection value, while temperatures above 1.75 generate incoherent actions that even accurate reflection cannot salvage. *Bottom Left:* Performance vs. LoRA configuration (rank, alpha). The optimal configuration ($r=8, \alpha=16$) achieves 60.0% performance, balancing adaptation capacity with training stability. Smaller configurations like (4,4) underfit with insufficient capacity (52.5%), while larger configurations cause mode collapse during test-time training—(16,32) drops to 41.5% and (32,32) collapses to 34.8% as the model begins predicting identical outputs for all inputs, losing the ability to distinguish between different spatial configurations and task contexts. *Bottom Right:* Performance vs. action budget (maximum steps). Performance improves dramatically from 30 steps (51.5%) to 50 steps (60.0%), but slightly degrades to 59.4% at 100 steps, suggesting that excessive action budgets allow suboptimal exploration strategies that accumulate errors over longer horizons.

894 tions (e.g., placing a small object in a large compartment
895 early). At $N \geq 6$, the candidate pool may include too many
896 low-quality options that add noise without improving the
897 best candidate’s quality, and the internal reflection model’s
898 scoring may become less reliable across excessively diverse
899 samples.

900 C.2. Sampling Temperature Analysis

901 Temperature T controls the sharpness of the probability
902 distribution during action generation: low temperatures

($T \rightarrow 0$) produce near-greedy sampling, while high temper- 903
atures ($T \rightarrow \infty$) approach uniform sampling. We evalu- 904
ate $T \in \{0.0, 0.25, 0.5, 0.75, 1.0, 1.25, 1.5, 1.75, 2.0\}$ (Fig- 905
ure 6, top right). 906

Results. Performance exhibits a clear inverted-U rela- 907
tionship with temperature. At $T=0.0$ (deterministic greedy 908
decoding), performance is 52.6%. Then it increases steadily 909
through moderate temperatures, reaching a peak of 60.0% 910
at $T=1.25$, and maintains near-peak performance of 59.8% 911
at $T=1.5$. However, performance drops sharply at higher 912

913 temperatures: 49.5% at $T=1.75$ and 47.0% at $T=2.0$.

914 **Analysis.** The optimal temperature range $T=1.25-1.5$
915 balances two competing factors: (1) *Diversity for reflec-*
916 *tion*: Sufficient randomness generates meaningfully differ-
917 ent candidates for the internal reflection model to com-
918 pare, enabling it to identify strategic failures (e.g., block-
919 ing future placements) that greedy or near-greedy decoding
920 would miss; (2) *Quality preservation*: Excessive random-
921 ness ($T \geq 1.5$) samples from the low-probability tail of the
922 distribution, producing incoherent or physically infeasible
923 actions that even accurate internal reflection cannot salvage.
924 The sharp drop at $T \geq 1.75$ indicates that overly stochastic
925 sampling overwhelms the reflection mechanism’s corrective
926 capacity—no amount of scoring can recover from funda-
927 mentally poor candidate pools.

928 C.3. LoRA Configuration Analysis

929 For parameter-efficient test-time training, we apply Low-
930 Rank Adaptation (LoRA) with varying rank r and scal-
931 ing factor α . We evaluate configurations $(r, \alpha) \in$
932 $\{(4, 4), (4, 8), (8, 8), (8, 16), (16, 16), (16, 32), (32, 32)\}$ to
933 study the trade-off between adaptation capacity and training
934 stability (Figure 6, bottom left).

935 **Results.** Performance exhibits a sharp peak at the in-
936 termediate configuration. Small configurations show lim-
937 ited performance: (4,4) achieves 52.5% and (4,8) reaches
938 56.2%, indicating insufficient adaptation capacity. Perfor-
939 mance peaks at ($r=8, \alpha=16$) with 60.0%, our optimal confi-
940 guration. However, larger configurations show dramatic per-
941 formance degradation: (8,8) maintains 55.8%, but (16,16)
942 drops to 57.8%, (16,32) plummets to 41.5%, and (32,32)
943 collapses catastrophically to 34.8%—below even the no-
944 adaptation baseline of 53.0%.

945 **Analysis.** The LoRA rank controls the expressiveness of
946 the adapter matrices, but excessive rank causes mode col-
947 lapse during aggressive test-time training. At ($r=4, \alpha=4$)
948 and ($r=4, \alpha=8$), the adapter capacity is too limited to cap-
949 ture nuanced spatial reasoning required for effective reflec-
950 tion—the model cannot adequately learn from retrospective
951 feedback about blocking placements or strategic failures.
952 The optimal configuration ($r=8, \alpha=16$) provides sufficient
953 capacity to update internal reflection scoring and action se-
954 lection policies based on task-specific feedback while main-
955 taining stable optimization under high learning rates (0.2 for
956 internal reflection, 0.01 for action model).

957 Larger configurations fail catastrophically due to mode
958 collapse. At ($r=16, \alpha=32$) and beyond, the increased pa-
959 rameter space combined with aggressive learning rates and
960 limited training data (10-15 retrospective examples per test-
961 time training iteration) causes the model to collapse to pre-
962 dicting identical outputs for all inputs. Rather than learn-
963 ing task-specific spatial reasoning, the overparameterized
964 adapters converge to degenerate solutions that ignore input

965 variations. At (32,32), the model effectively stops distin-
966 guishing between different object configurations, compart-
967 ment sizes, or task contexts—producing the same stereot-
968 yped action regardless of the actual scene state. This mode
969 collapse is particularly severe because test-time training
970 lacks the regularization and diverse data that prevent col-
971 lapse in offline training.

972 C.4. Action Budget Analysis

973 We evaluate the impact of maximum action budget on task
974 performance by varying the step limit from 30 to 100 steps
975 (Figure 6, bottom right).

976 **Results.** Performance shows a sharp initial gain fol-
977 lowed by slight degradation. At 30 steps, performance
978 is only 51.5%, indicating insufficient budget to complete
979 multi-object placement tasks. Performance peaks at 50
980 steps with 60.0%, our optimal setting. However, extending
981 the budget to 100 steps results in slight performance degra-
982 dation to 59.4%, a 0.6 percentage point drop.

983 **Analysis.** The low performance at 30 steps reflects
984 task incompleteness—agents frequently run out of steps be-
985 fore placing all objects, especially when early mistakes re-
986 quire exploration of alternative strategies. The 50-step bud-
987 get provides sufficient runway for the agent to explore the
988 environment, execute placements, recover from initial fail-
989 ures through reflection, and retry with corrected strategies
990 learned via test-time training.

991 The counterintuitive degradation at 100 steps reveals a
992 failure mode of excessive budgets: when given too many
993 steps, agents exhibit suboptimal exploration patterns that
994 accumulate errors over longer horizons. With a gener-
995 ous budget, the model may attempt more exploratory ac-
996 tions rather than committing to placements, leading to in-
997 efficient trajectories. Additionally, longer episodes pro-
998 vide more opportunities for compounding errors—a single
999 poor placement early in a 100-step episode has more down-
1000 stream consequences than in a tightly constrained 50-step
1001 episode where the agent must act decisively. This suggests
1002 that moderate action budgets not only improve computa-
1003 tional efficiency but also serve as a useful inductive bias
1004 that encourages focused, goal-directed behavior in embod-
1005 ied agents.

1006 D. Single-Step Action Generation vs. Receding 1007 Horizon Planning

1008 A key design decision in our framework is single-step action
1009 generation: the agent generates and executes one action at a
1010 time, rather than planning action sequences through reced-
1011 ing horizon control. We validate this choice through abla-
1012 tion experiments on the Cupboard Fitting benchmark.

1013 **Experimental Setup.** We compare two variants of our
1014 full method: (1) *Single-Step (Ours)*: generates one action,
1015 executes it, observes outcome, performs test-time training,

1016 then generates the next action; (2) *Receding Horizon*: gen- 1058
 1017 erates a complete action sequence (5-10 actions), executes 1059
 1018 only the first action, observes outcome, performs test-time 1060
 1019 training, then replans a new sequence from the updated 1061
 1020 state. Both variants use identical model architecture, train- 1062
 1021 ing procedures, and test-time training mechanisms, differ- 1063
 1022 ing only in planning granularity. The receding horizon ap- 1064
 1023 proach requires approximately 5× more computation per 1065
 1024 step due to generating full sequences at each decision point. 1066

Table 4. Single-step action generation vs. receding horizon planning on Cupboard Fitting. Despite using 5× more computation to plan sequences at each step, receding horizon shows degraded performance, demonstrating that single-step action generation is more effective for our reflective learning framework. 1067

Method	Fit	Correct	Compute
Ours (Receding Horizon)	57.8%	25.8%	5.0×
Ours (Single-Step)	60.2%	25.3%	1.0×

1025 **Results.** Table 4 shows that receding horizon planning 1058
 1026 achieves only 57.8% fit rate compared to 60.0% for single- 1059
 1027 step action generation—a comparable performance but at 1060
 1028 the cost of 4x more computation per step. This demon- 1061
 1029 strates that planning full action sequences at each decision 1062
 1030 point not only increases computational cost but actually 1063
 1031 harms performance in our framework. 1064

1032 **Why Receding Horizon is Incompatible with Test-** 1058
 1033 **Time Training.** The performance gap reveals fundamental 1059
 1034 incompatibilities between sequence planning and reflective 1060
 1035 test-time training: 1061

1036 (1) *Wasted computation on unpredictable futures:* Re- 1058
 1037 ceding horizon generates 5-action sequences at each step 1059
 1038 but executes only the first action, discarding 80% of the 1060
 1039 computation. Critically, in our error-driven tasks, the out- 1061
 1040 comes of actions are inherently unpredictable before exe- 1062
 1041 cution—whether an object fits in a compartment, whether 1063
 1042 placement will be stable, or whether grasping succeeds de- 1064
 1043 pends on precise physical interactions that cannot be reli- 1065
 1044 ably simulated. The model imagines 4 future actions based 1066
 1045 on assumed execution outcomes, but these assumptions are 1067
 1046 frequently wrong. When the first action fails or succeeds 1068
 1047 differently than predicted, the entire planned sequence be- 1069
 1048 comes invalid. This wasted computation could instead gen- 1070
 1049 erate more candidates for reflection-in-action (increasing N 1071
 1050 from 3 to 6) or perform additional test-time training epochs, 1072
 1051 both of which operate on actual execution outcomes rather 1073
 1052 than unreliable predictions. 1074

1053 (2) *Learning from imagination conflicts with learning* 1058
 1054 *from reality:* Test-time training fundamentally relies on 1059
 1055 learning from actual execution feedback—the model up- 1060
 1056 dates its understanding of spatial constraints only after 1061
 1057 physically attempting placements and observing real out- 1062

1058 comes (does the object fit? does it block other spaces?). 1059
 1059 However, generating 5-action sequences forces the model 1060
 1060 to predict these outcomes before they occur: "If I place ob- 1061
 1061 ject A here, then I can place object B there, then object 1062
 1062 C..." These predictions are made with the model's current 1063
 1063 (pre-update) understanding, which is precisely what test- 1064
 1064 time training aims to improve. The model must simulta- 1065
 1065 neously optimize two conflicting objectives: (a) accurately 1066
 1066 predicting hypothetical future states for sequence genera- 1067
 1067 tion, and (b) updating its beliefs based on actual execution 1068
 1068 outcomes that contradict those predictions. This creates op- 1069
 1069 timization interference where gradients from test-time train- 1070
 1070 ing (learned from reality) fight against the inductive bias 1071
 1071 from sequence generation (learned from imagination). 1072

1072 Our single-action approach eliminates this conflict by 1073
 1073 committing only to decisions that can be made based on 1074
 1074 current observations, executing the action to obtain ground 1075
 1075 truth feedback, updating the model with real outcomes, then 1076
 1076 making the next decision with improved understanding. 1077
 1077 This aligns perfectly with the test-time training paradigm: 1078
 1078 learn from actual experience, not imagined futures. 1079

1079 **Retrospective Reflection Provides Implicit Long-** 1058
 1080 **Horizon Planning.** A potential concern is that single-step 1059
 1081 action generation lacks the foresight that explicit sequence 1060
 1082 planning provides. However, retrospective reflection ad- 1061
 1083 dresses this through a different mechanism. While reced- 1062
 1084 ing horizon achieves long-horizon reasoning by explicitly 1063
 1084 generating future action sequences, our approach achieves 1064
 1085 it through learned anticipation: retrospective reflection 1065
 1085 re-evaluates past actions with hindsight about their long-term 1066
 1086 consequences, creating training signals that teach the inter- 1067
 1087 nal reflection model to anticipate multi-step outcomes be- 1068
 1088 fore execution. 1069

1090 For example, when placing objects in a cupboard, an ac- 1058
 1091 tion that initially appears successful may be retrospectively 1059
 1092 downgraded when the agent discovers that this placement 1060
 1093 blocks the only space for a larger object. These retrospec- 1061
 1094 tive scores train the internal reflection model to predict such 1062
 1095 long-horizon consequences at decision time—effectively 1063
 1096 distilling multi-step lookahead into single-step action eval- 1064
 1097 uation. This learned implicit planning is more sample- 1065
 1098 efficient than explicit sequence generation: rather than ex- 1066
 1099 ploring all possible future sequences at every step (most of 1067
 1100 which will be discarded), the agent learns which single ac- 1068
 1101 tions lead to favorable long-term outcomes through accu- 1069
 1102 mulated experience. 1070

1101 **Computational Efficiency Analysis.** The 5× computa- 1058
 1102 tional cost of receding horizon stems from generating full 1059
 1103 5-action sequences at each decision step, only to execute 1060
 1104 the first action and discard the rest. This is particularly 1061
 1105 inefficient in our test-time training setting, where: (1) Se- 1062
 1106 quence generation requires forward passes through the lan- 1063
 1107 guage model for all actions in the sequence; (2) The ad- 1064
 1108 1109 1110

ditional compute does not improve learning quality—test-time training still operates on single executed actions, so planning discarded sequences provides no learning benefit; (3) The saved computation in single-step action generation can be reallocated to improvements that actually help: generating more candidates (N=6 vs. N=3), performing more test-time training epochs, or using larger working memory windows for retrospective reflection.

Our results demonstrate that effective long-horizon reasoning in embodied agents need not come from explicit sequence planning. Instead, combining single-step action generation with retrospective reflection achieves superior performance at 5× lower computational cost by learning to anticipate long-term consequences rather than exhaustively simulating future possibilities.

E. Experiments on Long-Horizon Household Tasks: More Details

In this section, we provide comprehensive implementation details for our Long-Horizon Household Tasks, including task generation, validation, data construction, model architecture, training procedures, and evaluation protocols.

E.1. Why Build Upon BEHAVIOR?

BEHAVIOR-1K [18] provides an excellent foundation for embodied AI research with several key strengths: (1) **Photorealistic environments**: BEHAVIOR-1K features high-fidelity household scenes with realistic object models, physics simulation, and diverse room layouts across 1,000+ scenes; (2) **Rich object diversity**: The benchmark includes hundreds of object categories with varied sizes, shapes, and physical properties, enabling complex manipulation tasks; (3) **Standardized infrastructure**: BEHAVIOR-1K provides well-maintained simulation infrastructure, observation APIs, and action spaces that facilitate reproducible research.

However, BEHAVIOR-1K’s original task design does not systematically stress two critical capabilities for reflective learning: (1) **Learning from failures**: Most BEHAVIOR-1K tasks are designed to be solvable with correct initial planning, without requiring agents to recover from or learn from execution failures. Tasks rarely include scenarios where early actions create downstream failures that only become apparent after multiple steps (e.g., placing a small object first that later blocks the only space for a larger object). (2) **Long-term dependencies**: The original benchmark emphasizes task completion but does not specifically design tasks around sequential dependencies where action order critically determines success, or where consequences of early actions remain hidden until much later in the episode.

To evaluate our reflective test-time training framework—which specifically learns from execution failures

through retrospective reflection—we adapt BEHAVIOR-1K environments to create tasks that systematically incorporate these failure modes. We retain BEHAVIOR-1K’s photorealistic scenes and simulation infrastructure while introducing task specifications that stress failure recovery, long-horizon credit assignment, and dependency reasoning. This allows us to leverage the strengths of BEHAVIOR-1K (realism, diversity, standardization) while evaluating capabilities (learning from failures, retrospective reasoning) that the original benchmark was not designed to measure.

E.2. Task Categories

We develop a systematic pipeline to generate Long-Horizon Household Tasks that stress error-driven adaptation in embodied agents. Our tasks are designed around four core failure modes common in real life: spatial reasoning errors, object selection mistakes, sequential dependency violations, and non-local planning failures.

Task Categories and Failure Modes. We define four task categories, each targeting specific failure patterns:

Fitting Tasks require agents to pack or place objects into constrained containers or surfaces. These tasks stress geometric reasoning and capacity constraints. Common failure modes include: (1) attempting to place oversized objects in small compartments, (2) blocking access to larger storage spaces with premature small-object placements, and (3) failing to recognize occlusion after placement. For example, placing a toy car in a box already containing a teddy bear may succeed physically but block future placements of larger items.

Selection Tasks require agents to compare and retrieve items based on preferences or constraints. Failure modes include: (1) selecting suboptimal items when better alternatives exist in unexplored rooms, (2) committing to choices before exploring all options, and (3) failing to revise decisions when new information becomes available. A typical scenario involves retrieving vegetables for a meal where lettuce (preferred) is in one room and tomato (less preferred) is in another—agents that explore insufficiently may select the inferior option.

Preparation Tasks involve sequential constraints and dependencies where actions must occur in specific orders. Common failures include: (1) attempting steps out of sequence (e.g., adding toppings before placing the base plate), (2) violating prerequisite conditions (e.g., trying to cook without retrieving ingredients first), and (3) missing intermediate setup steps. These tasks require agents to maintain action dependencies across multiple rooms and objects.

Hybrid Tasks combine multiple failure modes within a single episode, requiring agents to simultaneously reason about spatial constraints, object preferences, and sequential dependencies across long horizons.

The distribution of task categories can be found in Figure

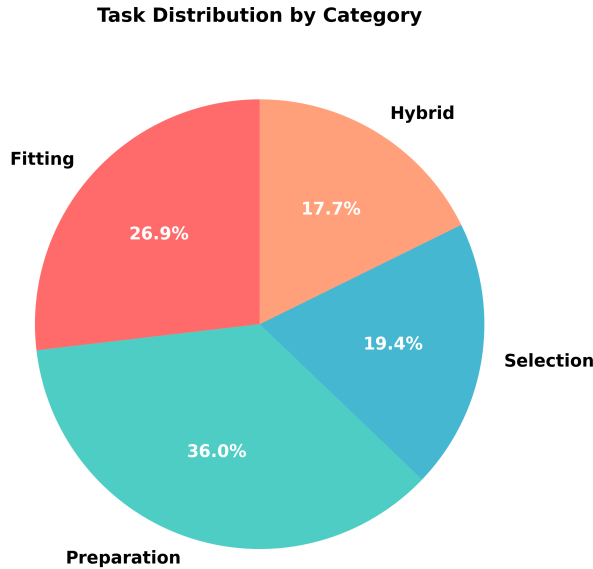


Figure 7. Distribution of task categories in the dataset.

1214 7.

1215 E.3. GPT-5 Task Generation Prompting Strategy

1216 We employ a carefully structured prompting approach to
1217 generate high-quality, physically plausible long-horizon
1218 tasks. Our prompting strategy consists of three key compo-
1219 nents: (1) comprehensive scene context provision, (2) ex-
1220 plicit failure mode specification, and (3) structured output
1221 formatting with reflection annotations.

1222 **Scene Context Provision.** Each task generation prompt
1223 includes the complete scene graph from BEHAVIOR-1K,
1224 containing: room layouts with spatial relationships, exist-
1225 ing objects in each room with 3D bounding boxes and af-
1226 fordance properties, navigable connections between rooms,
1227 and furniture placement. This rich context allows GPT-5
1228 to reason about physical feasibility when proposing object
1229 placements and action sequences. For example, knowing
1230 that a table has dimensions $1.2\text{m} \times 0.8\text{m} \times 0.75\text{m}$ allows the
1231 model to avoid proposing placements of oversized objects.

1232 **Failure Mode Specification.** We explicitly instruct
1233 GPT-5 to incorporate specific failure scenarios into the gen-
1234 erated trajectories. For Fitting tasks, we request scenarios
1235 where early placements block later optimal choices. For
1236 Selection tasks, we specify that preferred items should be
1237 in rooms requiring more exploration. For Preparation tasks,
1238 we request action sequences with complex dependencies
1239 where naive sequential execution fails. For each task cat-
1240 egory, we provide 2-3 concrete examples of desired failure
1241 patterns in few-shot demonstrations.

1242 **Structured Output with Reflection Annotations.**
1243 The prompt requires GPT-5 to generate not just action

sequences, but complete trajectories with: (1) Each action 1244
wrapped in angle brackets (e.g., `<GO TO kitchen>`); 1245
(2) Internal reflections before each action with format 1246
`INTERNAL REFLECTION: [reasoning]` 1247
`| SCORE: [0-100]`; (3) Execution result annota- 1248
tions (execution: success or execution: 1249
fail); (4) External reflections after each executed ac- 1250
tion with format `PREVIOUS ACTION TO REFLECT` 1251
`ON: <action> | EXTERNAL REFLECTION:` 1252
`[assessment] | SCORE: [0-100]`; (5) Retro- 1253
spective reflections at room exits with format `PREVIOUS` 1254
`ACTION TO REFLECT ON (retro): <action>` 1255
`| EXTERNAL REFLECTION: [hindsight` 1256
`assessment] | SCORE: [0-100].` 1257

Actions marked with `INTERNAL REFLECTION:` 1258
`... | SCORE: [low score]` followed by the note 1259
“score low. don’t execute” represent strategically poor 1260
choices that should be avoided through reflection-in-action. 1261
These actions may be physically feasible but lead to subop- 1262
timal task outcomes (e.g., picking up an inferior item when 1263
better alternatives exist). Yet they provide valuable data to 1264
train the internal reflection LLM (these actions are not used 1265
for training action LLM though). 1266

Few-Shot Demonstrations. We provide 2-3 complete 1267
task examples for each task category, demonstrating the ex- 1268
pected output format, reflection structure, and failure pat- 1269
terns. These examples show diverse scenarios: a Fitting task 1270
where placing small objects first blocks large object storage, 1271
a Selection task where exploring only nearby rooms leads to 1272
selecting inferior items, and a Preparation task where vio- 1273
lating sequential dependencies causes failure. 1274

Prompt Iteration and Refinement. We iteratively re- 1275
fined the prompt through multiple rounds of generation and 1276
validation. Early versions produced tasks with: (1) Phys- 1277
ically implausible object placements (e.g., large furniture 1278
items on small shelves); (2) Insufficient failure diversity 1279
(most tasks had similar error patterns); (3) Inconsistent re- 1280
flection scores (high scores for objectively poor actions). 1281
We addressed these through: (1) Adding explicit bounding 1282
box information and size reasoning requirements; (2) Pro- 1283
viding diverse few-shot examples spanning different failure 1284
modes; (3) Including calibration guidelines for score as- 1285
signment (e.g., “score 0-30 for actions that fail or lead to 1286
dead ends; 70-100 for optimal strategic choices”). 1287

Reflection Score Calibration. To ensure consistent 1288
score semantics across generated tasks, we provide GPT- 1289
5 with explicit calibration guidelines: Scores 0-20 indicate 1290
actions that fail physically or lead to immediate dead ends; 1291
scores 21-40 indicate poor strategic choices that succeed 1292
physically but create future problems; scores 41-60 indi- 1293
cate suboptimal but acceptable actions; scores 61-80 indi- 1294
cate good strategic choices; scores 81-100 indicate optimal 1295
or near-optimal actions. This calibration ensures that scores 1296

1297 are meaningful training signals for the reflection models
1298 rather than arbitrary numbers.

1299 E.4. Physical Validation in BEHAVIOR.

1300 Raw GPT-5 outputs may contain physically implausible
1301 scenarios or inconsistent spatial reasoning. We validate
1302 each generated task through execution in BEHAVIOR Om-
1303 niGibson physics simulation:

1304 *Object Placement Validation:* We verify all new ob-
1305 jects can be physically placed at specified locations using
1306 sampling-based kinematics (`sample_kinematics`). Ob-
1307 jects that fail placement (due to size mismatches, collision
1308 constraints, or stability issues) trigger task rejection. We
1309 use uniform scaling for all objects to avoid non-orthogonal
1310 transform errors.

1311 *Trajectory Execution Verification:* We execute the
1312 ground-truth trajectory action-by-action, verifying that: (1)
1313 Navigation actions (`GO TO`) successfully place the robot
1314 in target rooms; (2) Manipulation actions (`PICK UP`, `PUT`
1315 `DOWN`) complete as annotated; (3) Expected failures actu-
1316 ally fail (confirming physical constraints match annota-
1317 tions); (4) Action sequences respect object affordances and
1318 spatial constraints. Tasks where any expected-success ac-
1319 tion fails, or any expected-failure action succeeds, are re-
1320 jected as inconsistent. This ensures our evaluation measures
1321 genuine agent learning rather than dataset annotation errors.

1322 *Observation Generation:* During validation, we generate
1323 3D point cloud observations for each room after every in-
1324 teraction step (navigation, pickup, placement). We capture
1325 observations from three external camera viewpoints posi-
1326 tioned around the robot at fixed relative poses. For each
1327 viewpoint, we: (1) Capture RGB-D images at 560×560 res-
1328 olution; (2) Convert depth to point clouds using camera
1329 intrinsics; (3) Transform point clouds from camera frame
1330 to robot base frame; (4) Store point clouds (stacked across
1331 viewpoints) as `.numpy` files; (5) Save corresponding RGB im-
1332 ages as `.png` files. Each room observation is stored in a
1333 unique directory named `{room_name}_{step_idx}`, al-
1334 lowing models to access the most recent observation per
1335 room at any point during inference.

1336 E.5. Training Data Construction

1337 From validated trajectories, we extract training data for
1338 three model components. All training examples include
1339 point cloud observations stored as file paths, which are
1340 loaded and processed by the 3D vision encoder during train-
1341 ing.

1342 **Action Training Data.** For each step t in a validated

trajectory, we create training examples of the form: 1343

Input: {Task, All Rooms, Explored Rooms, 1344

Current Room, Observations, 1345

Previous Action, Previous External Reflection} 1346

Output: $Action_t$ 1347

Observations consist of point cloud paths for each explored 1348
room (most recent per room). 1349

Internal Reflection Training Data. For each action (in- 1350
cluding non-executed low-score actions), we create exam- 1351
ples: 1352

Input: {Task, All Rooms, Explored Rooms, 1353

Current Room, Observations, 1354

Previous Action, Previous External Reflection} 1355

Potential $Action_t$ } 1356

Output: {Internal Reflection $_t$, Score $_t$ } 1357

This includes actions marked “don’t execute”—the model 1358
must learn to score these actions low during internal reflec- 1359
tion to prevent execution. 1360

External Reflection Training Data. After each exe- 1361
cuted action, we create examples: 1362

Input: {Task, All Rooms, Explored Rooms, 1363

Current Room, Observations (before and after), 1364

Previous Action, Previous External Reflection} 1365

Executed Action $_t$, Execution Result $_t$ } 1366

Output: {External Reflection $_t$, Score $_t$ } 1367

Execution results indicate `success` or `fail`, and obser- 1368
vations include both pre-action and post-action point clouds 1369
to enable change detection. 1370

Retrospective Reflection Training Data. For retro- 1371
spective reflection, we identify room transitions (marked by 1372
`EXIT` actions) and collect all actions taken in that room. 1373
For each historical action a_j , we create: 1374

Input: {Task, Context, Current Observations, 1375

Action $_j$, Room Action History, 1376

Last Reflection $_j$ } 1377

Output: {Retro Reflection $_j$, Updated Score $_j$ } 1378

The prompt includes all actions and their external reflec- 1379
tions from the current room window, allowing the model 1380
to re-evaluate a_j with hindsight about downstream conse- 1381
quences. 1382

1383 E.6. Model Architecture and Training

Base Architecture. We build our 3D vision-language- 1384
action model on LLaVA-3D [40], which processes point 1385

1386 clouds through a 3D encoder and fuses them with language
1387 instructions via a multimodal projector. The architecture
1388 consists of: (1) A 3D point cloud encoder that extracts spa-
1389 tial features; (2) A vision-language projector that aligns 3D
1390 features with language embeddings; (3) A Llama-based lan-
1391 guage model backbone (7B parameters) for reasoning and
1392 generation.

1393 E.7. Evaluation Protocol and Baselines

1394 **Train-Test Split.** We ensure zero overlap between fine-
1395 tuning and evaluation: (1) No shared task descriptions—
1396 evaluation tasks have completely different natural lan-
1397 guage instructions; (2) No shared scenes—different scene
1398 instances from BEHAVIOR-1K; (3) No shared object
1399 placements—all object configurations are unique; (4) No
1400 shared trajectories—action sequences and reflection pat-
1401 terns differ. This tests generalization to novel tasks rather
1402 than memorization.

1403 **Success Criteria.** A task succeeds if and only if: (1) All
1404 required objects reach target locations; (2) All spatial con-
1405 straints are satisfied (e.g., inside, on top of); (3) All prefer-
1406 ence constraints are met (e.g., selecting preferred items); (4)
1407 All sequential dependencies are respected; (5) Task comple-
1408 tion occurs within the action budget (30 steps).

1409 **Deployment Procedure.** For each test task, we: (1) Ini-
1410 tialize agent in the first room with task description; (2) Exe-
1411 cute action generation + internal reflection + external reflec-
1412 tion loop; (3) Trigger retrospective reflection at room exits
1413 or after $K = 5$ steps; (4) Perform test-time training when
1414 retrospective reflections accumulate; (5) Continue until task
1415 completion or action budget exhaustion. The agent receives
1416 no human feedback during deployment—all learning sig-
1417 nals come from self-generated reflections.

1418 **Baseline Implementations.** We implement the follow-
1419 ing baselines: *Reflexion* [30] maintains a text buffer of past
1420 reflections and includes them in prompts for future actions,
1421 generating verbal critiques after each step without param-
1422 eter updates. *Self-Refine* [21] iteratively improves actions
1423 through self-critique and revision cycles, allowing up to 3
1424 refinement iterations per action before execution. We revise
1425 the above two baselines to incorporate multimodal inputs.
1426 *ReflectVLM* adapts the reflection mechanism from [5] using
1427 learned value functions for action scoring, with a separate
1428 value head trained on our training data. *3DLLM-Mem* [12]
1429 maintains fused point cloud observations from all previous
1430 steps and previous-step execution results as context, provid-
1431 ing spatial memory without explicit reflection or test-time
1432 training. *PPO* [29] and *DreamerV3* [39] are trained as re-
1433 inforcement learning baselines: PPO uses on-policy policy
1434 gradient with clipped surrogate objective and GAE for ad-
1435 vantage estimation, while DreamerV3 learns a world model
1436 from observations and trains a policy in the learned latent
1437 space. Both RL baselines are trained for the same total

number of environment interactions as our supervised fine-
tuning phase to ensure fair compute comparison. All base-
lines use the same LLaVA-3D-7B backbone with identical
finetuning procedures where applicable for fair comparison.

F. Cupboard Fitting: More Details

F.1. Base Model and Supervised Fine-tuning.

We build upon Qwen2.5-VL-3B as our base vision-
language model for the Cupboard Fitting benchmark. We
employ unified multi-task supervised fine-tuning on action
generation, internal reflection, and external reflection tasks
using: global batch size 128, learning rate 1×10^{-5} , weight
decay 0.1, trained for 3 epochs using AdamW optimizer.
Task-specific prompts distinguish between the three reflec-
tion modes, enabling cross-task knowledge transfer.

F.2. Test-Time Training Variants.

We implement two test-time training variants to study the
trade-off between adaptation capacity and computational
efficiency:

Base-Weight Test-Time Training for Reflection-on-Action: This variant updates all non-visual parameters dur-
ing deployment. We freeze visual encoder parameters
(identified by `visual` or `vision` in parameter names)
but update all language model parameters. For the action
model trained via REINFORCE, we use SGD optimizer
with learning rate 1×10^{-3} , weight decay 1×10^{-4} , zero
momentum, training for 3 epochs. For the internal reflec-
tion model trained via supervised learning, we use SGD op-
timizer with learning rate 5×10^{-5} , weight decay 1×10^{-4} ,
zero momentum, training for 3 epochs. This approach pro-
vides maximum adaptation capacity but requires updating
millions of parameters.

LoRA Test-Time Training for Reflection-on-Action: For
memory-efficient adaptation, we apply Low-Rank Adapta-
tion with rank $r = 8$, alpha $\alpha = 16$, dropout rate 0.1, tar-
geting all linear layers except `lm_head`, `embed_tokens`,
and visual encoder components. The LoRA adapters are ap-
plied to both internal reflection and action models using the
PEFT library. For the action model, we use SGD with learn-
ing rate 1×10^{-2} , weight decay 1×10^{-4} , zero momentum,
trained for 3 epochs. For the internal reflection model, we
use SGD with learning rate 0.2, weight decay 1×10^{-4} , zero
momentum, trained for 3 epochs. The higher learning rates
compensate for the reduced parameter count—LoRA up-
dates only the low-rank adapter parameters, reducing train-
able parameters by over 95% while maintaining comparable
performance.

G. More Qualitative Examples

In Figure 8, we show more qualitative examples. We can
see that the model does get improved over reflection mech-



Figure 8. Steps and reflections are simplified for better presentations. Blue text shows internal reflection used for candidate selection, orange text shows external reflection after execution. We put reflection scores inside brackets. Red text shows retrospective reflection and model updates.

1487 anisms.

1488 **H. Physical Experiment Setup**

1489 We validate our framework on a physical Franka Panda
1490 robotic arm in a cupboard fitting scenario. The workspace
1491 consists of a multi-compartment cupboard with compart-
1492 ments of varying sizes, along with different shape ob-
1493 jects that must be placed into designated compartments.
1494 A camera mounted above the workspace captures RGB
1495 images from a top-down offset viewpoint, matching the
1496 ‘top_down_offset’ camera configuration used in simulation
1497 training. Before each task, we record the initial positions of
1498 objects and compartment locations in the robot base frame
1499 through calibration. Given high-level action commands
1500 from the model (e.g., “pick up the red cube and put it in the
1501 blue compartment”), the robot executes a two-phase manip-
1502 ulation: (1) moves to the object’s registered initial position
1503 and grasps it using the parallel-jaw gripper, (2) places the
1504 object at the target compartment location computed from
1505 the registered compartment position. After each action, the
1506 system captures a post-execution image from the same cam-
1507 era viewpoint and provides binary success/failure feedback
1508 based on whether the object remains within the target com-
1509 partment bounds, which is used to generate external reflec-
1510 tions following the same protocol as simulation.