

Memo: An Open-World Agent with Multi-Scale Exploration and Memory-Driven Navigation in Minecraft

Anonymous CVPR submission

Paper ID 25998

Abstract

Autonomous completion of sparse, long-horizon tasks in open-world environments remains a key challenge for embodied intelligence. Minecraft, with its vast maps, dynamic entities, and complex task chains, provides an ideal benchmark. Existing hierarchical agents often suffer from inefficient exploration and limited semantic memory, causing repeated failures. To address this, we propose **Memo** (Multi-Scale Exploration and Memory-Driven Navigation), an embodied agent integrating multi-scale exploration with hierarchical memory-driven planning. **Memo**'s Multi-Scale Progressive Exploration efficiently covers global, regional, and local areas, while its Dynamic Planner with Long-Short-Term Memory constructs a semantic knowledge graph to enable context-aware memory retrieval and adaptive task prioritization. Experiments show **Memo** outperforms prior methods, increasing exploration coverage by 14%, improving sequential task success by 30% with a 36% reduction in completion time, and boosting non-sequential task success by 44% with a 42% shorter execution time. These results demonstrate **Memo**'s ability to combine efficient exploration, memory-driven reasoning, and adaptive planning for robust open-world task execution. We plan to release our code to support further research.

1. Introduction

With the rapid advancement of large-scale foundation models, embodied intelligence is shifting from single-task capabilities toward general-purpose agents that operate robustly across diverse, unstructured environments [4, 14, 15, 19]. A key requirement for such generality is the ability to autonomously solve sparse, multi-step instruction tasks in open-world settings. These tasks demand not only locating distant or rare targets across large areas but also adapting navigation and decision-making policies in response to evolving environmental conditions—posing a fundamental challenge for general embodied agents.

Minecraft, as a representative open-world sandbox environment, naturally exhibits these complexities. Its vast and personalized maps, dynamic entities, unknown resources, and long task chains—from resource gathering and crafting to construction and combat—closely mirror real-world challenges. Consequently, Minecraft has emerged as a standard testbed for evaluating the capability of embodied agents to perform long-horizon, sequential tasks [2]. Within such environments, the controller, responsible for exploration, navigation, and grounding high-level instructions, is critical to overall agent performance.

Most existing agent architectures adopt a hierarchical framework, where high-level planners decompose instructions into subtask sequences [16, 24], while controllers execute actions, explore the environment, and manage memory [17]. Despite their promise, current agents suffer from two major limitations. First they often exhibit limited semantic understanding and inflexible planning. Their memory system operate at a data level, lacking mechanisms to construct or utilize structures from historical experience. This prevents the effective identify and reuse of sparse but informative past observations, leading agents to repeatedly fail in similar situations. Moreover, their fixed decision pipelines struggle to adapt when target objects vary or when unexpected events arise, often resulting in suboptimal or failed trajectories. Second, existing exploration strategies tend to be locally greedy and globally inefficient, leaving large regions unvisited while generating redundant trajectories in already explored areas. Such inefficiency becomes especially problematic in large-scale open worlds, where balanced and comprehensive exploration is vital for task completion.

To address these bottlenecks, we introduce **Memo**, a new open-world agent with Multi-Scale Exploration and Memory-Driven Navigation. **Memo** is designed to robustly handle instruction-following tasks involving sparse, distant, and dynamically evolving targets in large-scale Minecraft environments. Its architecture comprises three core components: explorer, memory, and planner. The Multi-Scale Progressive Exploration framework operates across global,

076 regional, and local scales, enabling comprehensive environment coverage while optimizing exploration paths. Dynamic Planning with Hierarchical Memory Update combines long- and short-term memory to support semantic cognition. Unlike traditional data-level memory, this system constructs a structured semantic knowledge graph from real-time interactions, facilitating context-aware information retrieval. Retrieved information feeds directly into the planner, which dynamically prioritizes tasks based on environmental observations and memory cues, allowing flexible and efficient task execution that adapts to dynamic changes. Experimental results demonstrate Memo’s effectiveness: compared to the MrSteve(state of the art)[17] in large-scale exploration, coverage increased by 14% with notable efficiency gains; for sparse sequential tasks, success rates improved by 30% with a 35.8% reduction in completion time; and in sparse non-sequential tasks, success rates rose by 44% with a 42.5% shorter completion time.

094 Our main contributions are summarized as follows:

- 1 We present Memo, a new open-world embodied agent equipped with multi-scale exploration and memory-driven navigation. Memo effectively handles instruction-following tasks with sparse and distant targets in large-scale Minecraft environments.
- 2 We propose a multi-scale progressive exploration strategy that gathers complementary information at global, regional, and local levels, enabling efficient exploration and large-area coverage in open-world settings.
- 3 We develop a dynamic planning mechanism that continuously updates geometric and semantic memory to support adaptive navigation, particularly in tasks where target objects are diverse or can change during execution.

108 2. Related work

109 2.1. Low-Level Control in Minecraft

110 Early rule-based systems, such as CraftAssist [8], relied on predefined scripts for game operations and dialogue, 111 lacking learning ability and adaptability. MineRL [9] introduced 112 reinforcement learning with large-scale human 113 demonstrations, enabling skill learning but struggling with 114 long-horizon tasks. The data-driven shift began with VPT 115 [1], which mapped vision to actions through large-scale 116 video pretraining, enabling autonomous complex behaviors. 117 Building on this, Steve-1 [13] employed a UNCLIP-based 118 text-vision-action mapping, supporting open-text 119 instruction following but without long-term memory for 120 multi-scenario tasks. Memory-augmented approaches like 121 MrSteve [17] incorporated a Place-Event Memory 122 mechanism to link spatial, temporal, and event information, 123 enhancing long-sequence execution. However, semantic 124 understanding and memory management remain shallow, and 125 exploration strategies still lack global planning and flexibil-

126 ity.

127 2.2. Exploration and Navigation Strategies

128 Efficient exploration is critical for open-world agents. 129 Frequency-based methods [3, 23, 27] are simple but suffer 130 from local redundancy and scalability issues. Video- 131 pretrained navigation (e.g., VPT-Nav) adapts to dynamic 132 environments but relies on local observations, leading 133 to suboptimal trajectories in unseen scenes. Active 134 information-gain strategies [4, 5, 29] discover novel 135 areas effectively but are computationally intensive, violating 136 Minecraft’s real-time (≤ 20 ms/step) constraints. Coverage- 137 based approaches ensure completeness but struggle to balance 138 coverage with task priorities and 3D spatial reasoning. 139

140 2.3. Memory Systems for Open-World Agents

141 Memory is essential for long-horizon reasoning and planning 142 [10, 20–22, 28]. Existing systems often use single- 143 modal or loosely fused multi-modal storage: Voyager [25] 144 stores skills as text; GITM [31] integrates textual memory 145 for efficient reasoning; MP5 [18] and JARVIS-1 [26] 146 maintain multi-modal context for situational retrieval; Optimus-1 147 [12] enhances storage with a multi-modal experience pool. 148 MrSteve [17] improves memory via a Place-Event 149 mechanism, yet retrieval accuracy, semantic abstraction, and 150 generalization remain limited. Overall, current memory 151 systems lack robust semantic adaptability and dynamic updating, 152 constraining long-term task performance in open-world 153 Minecraft settings.

154 3. Method

155 In this section, we present Memo, a novel embodied 156 agent designed to overcome the challenges of semantic 157 cognition and inefficient open-world exploration. As il- 158 lustrated in Fig. 1, Memo is built upon two core mod- 159 ules: 1) Multi-Scale Progressive Exploration, which plans 160 efficient exploration trajectories through a coordinated 161 global-regional-local structure, enabling both rapid cov- 162 erage and fine-grained target search. 2) Dynamic Planning 163 with Hierarchical Memory, which integrates fast, data-level 164 short-term memory and slow, knowledge-level long-term 165 memory. This allows the agent to rapidly retrieve recent 166 experiences, associate high-level semantic knowledge, and 167 dynamically reorder task priorities based on real-time ob- 168 servations. During operation, the exploration module pro- 169 poses cross-scale exploration paths, while the hierarchical 170 memory module continuously stores and retrieves environ- 171 mental information. The dynamic planning module then 172 fuses exploration objectives with memory feedback to adapt 173 the agent’s behavior, ensuring robust and efficient task ex- 174 ecution in complex open-world environments.

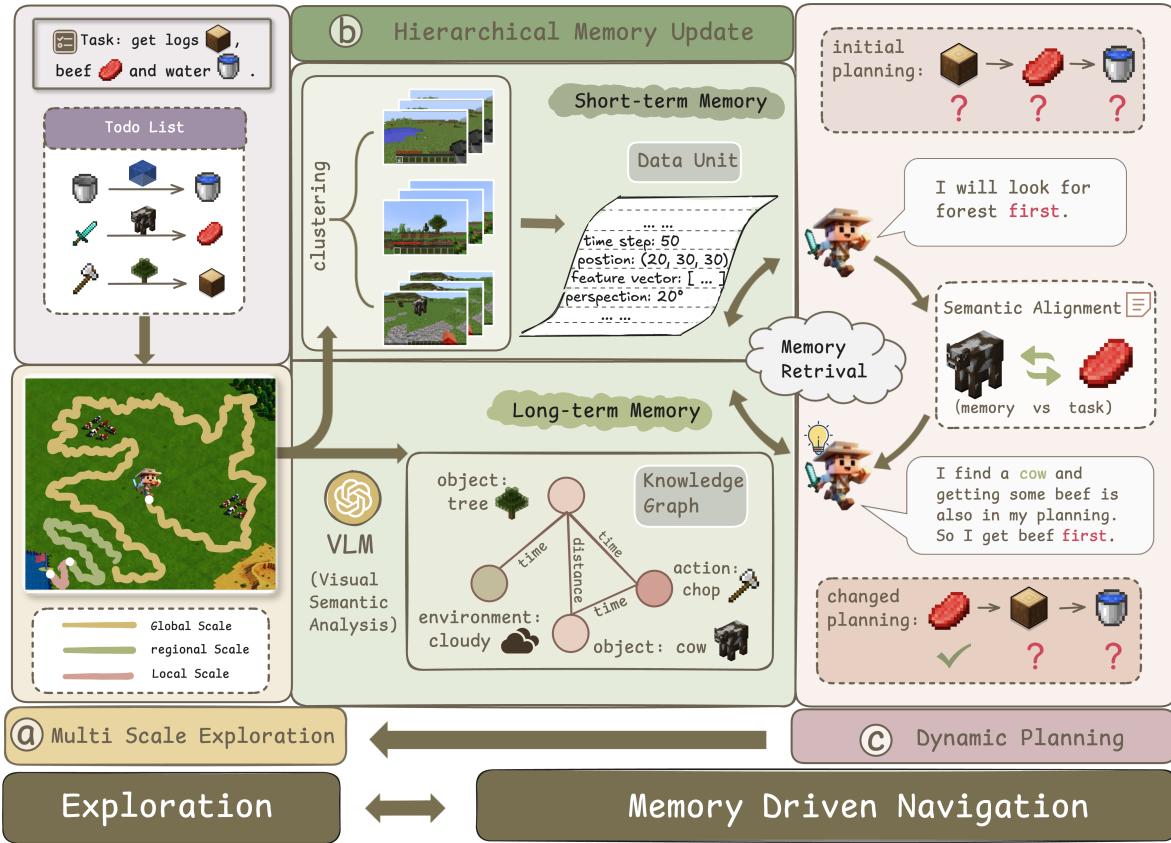


Figure 1. The overall architecture of Memo, which integrates exploration, memory, and planning for unified exploration–navigation. Tasks are decomposed into sub-tasks and executed through multi-scale progressive exploration: global rapid search, regional frontier refinement, and local execution. Newly observed information triggers hierarchical memory updates, storing short-term geometric cues for fast retrieval and long-term semantic knowledge in a graph. When the environment changes (e.g., encountering a cow while searching for a forest), Memo performs dynamic planning to reorder tasks and take advantage of new opportunities.

175

3.1. Multi-Scale Progressive Exploration

176

To overcome the inefficiency and locality issues of existing exploration strategies in large-scale environments, we propose a progressive exploration framework operating at the global, regional, and local scales. This framework balances exploration value and navigation cost, ensuring both wide-area coverage and fine-grained refinement.

177

178

179

180

181

182

3.1.1. Global-Scale Rapid Search

183

The global scale performs macro-level planning to identify promising exploration regions. Given Minecraft’s effectively infinite world, we adopt an adaptive quadtree spatial index, where the root node covers the full map and child nodes represent candidate exploration regions. The quadtree depth and subdivision adapt dynamically based on visitation frequency. To select the next region, we compute a priority score:

184

185

186

187

188

189

190

191

$$\text{Score}(N_i) = \alpha \cdot \text{Depth}(N_i) + \beta \cdot \text{Dist}(P, N_i) \quad (1)$$

Where $\text{Depth}(N_i)$ is the depth of the current node, quantifying the retrieval efficiency of the node; $\text{Dist}(P, N_i)$ is the Euclidean distance of the path length from the current position P to N_i , quantifying the navigation cost. This method reduces the time complexity of spatial query from $O(n)$ to $O(\log n)$, ensuring the efficiency of global planning.

3.1.2. Regional-Scale Frontier Refinement

At the regional scale, the agent focuses on frontier exploration, i.e., areas at the boundary between known and unknown space. Using a morphological-dilation-based frontier detector, we extract frontier points that maximize information gain. Each frontier point (f_i) is scored as:

$$\text{Score}(f_i) = \alpha \cdot \text{Cover}(f_i) + \beta \cdot \text{Nove}(f_i) + \gamma \cdot \text{Entro}(f_i) \quad (2)$$

Here, $\text{Cover}(f_i)$ is the area of the uncovered area within the visible range of the current viewpoint. $\text{Nove}(f_i)$ is the information entropy gain of the visible range of the current viewpoint. $\text{Entro}(f_i)$ is the distance between the current viewpoint and historical viewpoints.

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210 **3.1.3. Local-Scale Area Complement**

211 To achieve a precise coverage of unexplored areas at a
 212 local level, we utilize a Voronoi region decomposition
 213 algorithm[11]. This algorithm divides the map into n
 214 Voronoi cells $\{C_1, C_2, \dots, C_n\}$, based on the visited set
 215 $V = \{v_1, v_2, \dots, v_n\}$ (where $v_i = (x_i, y_i, t_i)$, x_i and y_i
 216 are 2D coordinates, t_i is the timestamp). Each cell C_i is de-
 217 fined as:

$$218 C_i = \{p \in \Omega \mid \|p - v_i\|_2 < \|p - v_j\|_2, \forall j \neq i\} \quad (3)$$

219 Here, Ω denotes the environmental space, and $\|\cdot\|_2$ repre-
 220 sent the Euclidean distance. This definition ensures that
 221 the distance from any point p in cell C_i to the visited
 222 point v_i is smaller than its distance to other visited points,
 223 thus satisfying the local optimality of spatial division. Un-
 224 covered regions U can be automatically identified, where
 225 $U = \Omega \setminus \bigcup_{i=1}^n C_i$ denotes the set difference operation, i.e.,
 226 regions in Ω that do not belong to any Voronoi cell C_i . The
 227 scoring formula for the i -th unexplored region U_i is given
 228 by:

$$229 Score(U_i) = \alpha \cdot Area(U_i) + \beta \cdot Dist(P, U_i) \quad (4)$$

230 Where $Area(U_i)$ is the uncovered area, and $Dist(P, U_i)$ is
 231 the navigation cost. This drives the agent to cover large,
 232 nearby gaps.

233 We also implement a local escape mechanism to prevent
 234 stagnation. If the agent moves less than 5 blocks in 30 sec-
 235 onds, it replans locally. Repeated failure to reach a target
 236 triggers a target reset from memory, and repeated revisiting
 237 of the same location initiates a rollback to a key node. These
 238 mechanisms prevent local traps and improve exploration re-
 239 silience.

240 **3.2. Dynamic Planning with Hierarchical Memory
241 Update**

242 To overcome the limitations of fixed planning strategies, we
 243 introduce a hierarchical memory–driven dynamic planning
 244 framework. By integrating multi-level memory with adap-
 245 tive decision-making, the agent can continuously refine its
 246 strategy based on both recent interactions and long-term se-
 247 mantic understanding, enabling flexible and efficient execu-
 248 tion of sparse, long-horizon tasks.

249 **3.2.1. Hierarchical Memory Update**

250 Traditional memory systems rely heavily on raw data stor-
 251 age without semantic structuring, limiting their ability to
 252 support long-sequence reasoning. To address this, we de-
 253 sign a hierarchical memory system: 1) Short-term memory
 254 provides rapid retrieval of recent experiences in terms of
 255 scene geometry. 2) Long-term memory constructs a struc-
 256 tured semantic knowledge graph. This combination offers
 257 both real-time responsiveness and long-range associative
 258 reasoning.

259 **Short-Term Geometric Memory for Rapid Experi-**

260 **ence Retrieval:** The short-term memory focuses on pro-
 261 cessing and organizing recent interaction data to provide
 262 quick access to comparable past scenes. Built upon the
 263 Place Event Memory mechanism, it first extracts visual fea-
 264 tures from gameplay frames using MineClip [7]. These
 265 features are then clustered via the DP-Means algorithm
 266 [6], with positional information incorporated to distinguish
 267 event types more effectively. To handle the redundancy and
 268 real-time nature of Minecraft interactions, a sliding-window
 269 strategy is applied to automatically discard outdated experi-
 270 ences, ensuring that the memory pool preserves only high-
 271 value recent information.

272 **Long-Term Semantic Memory for Knowledge Graph**

273 **Construction:** Given the computational cost of vi-
 274 sual–language models (VLMs), we first perform keyframe
 275 selection to reduce redundant inputs while preserving se-
 276 mantically meaningful information. To this end, we employ
 277 a hybrid keyframe extraction strategy that combines global
 278 detection based on image entropy with local detection based
 279 on feature clustering. The global detector identifies major
 280 environmental changes by computing the entropy of each t
 281 frame I_t :

$$282 H(I_t) = - \sum_{i=0}^{255} p_i \log_2 p_i \quad (5)$$

283 Here, p_i is the pixel intensity probability. If
 284 $|H(I_t) - H(I_{t-1})| > \theta_H$ (where θ_H is the thresh-
 285 old), the frame is considered a global keyframe. Local
 286 detection local , based on feature clustering, captures new
 287 events in more detail. we use MineClip to extract current
 288 frame features and compare them with the historical feature
 289 library using cosine similarity. If the minimum cosine
 290 similarity is less than θ_c (where θ_c is the threshold), it is
 291 determined as a local keyframe. This dual detection can
 292 avoid missing important information and effectively filter
 293 duplicate frames.

294 Using the selected keyframes, we construct a long-term
 295 semantic memory graph with GPT-4o [16], converting dis-
 296 crete observations into structured and interpretable knowl-
 297 edge. The resulting graph consists of three fundamental
 298 node types: entity nodes (e.g., cows, stones, zombies), ac-
 299 tion nodes (e.g., mining, attacking, crafting), and environ-
 300 ment nodes (e.g., darkness, rain, hunger states). Edges
 301 encode temporal relations (e.g., event sequences or trav-
 302 ersal time between locations), spatial relationships (e.g., rel-
 303 ative distances or co-occurrence), and causal dependen-
 304 cies. For instance, the entity nodes cow and beef are
 305 linked by a causal edge denoting “drops after killing.” To
 306 construct the graph, we first extract structured “environ-
 307 ment–object–action” triples from keyframes, and subse-
 308 quently apply a pruning stage to remove redundant or low-
 309 relevance nodes. This process yields a lightweight yet se-

310 mantically expressive knowledge graph that supports long-
311 horizon reasoning and interpretable memory retrieval.

312 **3.2.2. Scene Semantics Alignment for Memory Retrieval**

313 Memory retrieval in our framework proceeds through three
314 sequential stages: query understanding, memory matching,
315 and result ranking. To interpret user or task queries at a
316 semantic level, we employ a pre-trained Word2Vec model
317 to convert natural-language inputs into continuous semantic
318 vectors. The model is trained on a Minecraft-specific corpus
319 that comprehensively captures domain knowledge, includ-
320 ing creature drops, item crafting dependencies, and con-
321 tainer interactions. After tokenization and deduplication,
322 the corpus is used to construct resource–relation–resource
323 triples, resulting in a dense semantic space tailored to in-
324 game concepts. This specialization enables accurate sim-
325 ilarity estimation; for instance, the semantic closeness be-
326 tween “cow” and “beef” or between “water bucket” and
327 “water” exceeds 97.9%.

328 Building on this semantic embedding space, memory
329 retrieval integrates both short-term and long-term mem-
330 ory sources. Recent experiences are first queried from the
331 short-term geometric memory to capture the most imme-
332 diate contextual relevance. Subsequently, the query vector
333 is compared against nodes in the long-term semantic mem-
334 ory graph, enabling retrieval of historical knowledge with
335 deep semantic structure. Cross-referencing these two mem-
336 ory layers ensures that retrieval results balance immediate
337 perceptual context and long-horizon semantic associations.
338 Finally, all candidate results are jointly ranked according to
339 semantic similarity and temporal relevance, ensuring that
340 the retrieved memory is both contextually appropriate and
341 temporally aligned with current task demands.

342 **3.2.3. Dynamic Planning**

343 The dynamic planning module serves as the integrative core
344 that unifies perception, memory, and action, enabling the
345 agent to adaptively regulate its behavior in complex envi-
346 ronments. Inspired by human cognitive decision-making,
347 this module operates through three tightly coordinated
348 phases.

349 In the information fusion phase, the planner synthesizes
350 heterogeneous inputs—including real-time environmental
351 observations, exploration goals produced by the multi-scale
352 exploration hierarchy, and semantic cues retrieved from
353 memory—to form a holistic representation of the current
354 decision context.

355 In the task priority sorting phase, pending subtasks are
356 dynamically organized to avoid inefficiencies associated
357 with static or predetermined execution orders. Each sub-
358 task is evaluated along three dimensions: urgency, semantic
359 relevance (as informed by memory retrieval), and naviga-
360 tion cost. This multi-criteria prioritization ensures that the

361 agent allocates resources effectively while maintaining re-
362 sponsiveness to evolving task requirements.

363 In the optimal action generation phase, the planner pro-
364 duces an action sequence that best satisfies the prioritized
365 objectives under the current environmental conditions. By
366 jointly considering task ordering, real-time constraints, and
367 memory-derived semantic cues, the module closes the gap
368 between strategic planning and low-level execution. This
369 process enables coherent integration of exploration and task
370 execution, effectively overcoming the disconnect that often
371 arises in traditional controller-based systems.

372 **4. Experiments**

373 We evaluate the proposed Memo agent in sparse-object
374 tasks within open-world environments along three di-
375 mensions: (i) large-scale exploration (Section 4.1), (ii)
376 task-solving performance in sequential and non-sequential
377 sparse tasks (Sections 4.2–4.3), and (iii) ablation studies to
378 quantify the contributions of core modules (Section 4.4).

379 In open-world settings, critical resources are often
380 sparsely distributed. For example, crafting a water bucket
381 may require locating both water and iron across distant re-
382 gions. Solving such tasks requires memory of resource
383 locations and rational task sequencing. Accordingly, we
384 design two primary task categories: sequential and non-
385 sequential tasks.

386 **4.1. Exploration Evaluation in Large-scale Envi- 387 ronments**

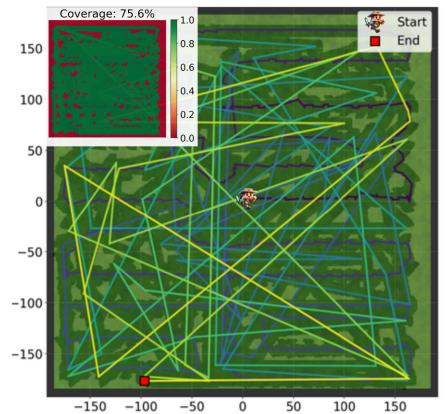
388 **Experimental Setup** Experiments are conducted in both
389 simulated and real Minecraft environments. The simulated
390 environment removes stochastic factors such as mobs and
391 terrain variations to ensure reproducibility and isolate al-
392 gorithmic performance. The real environment retains na-
393 tive interactions to validate practical applicability. Two map
394 scales are used: small (128×128 blocks) for localized eval-
395 uation and large (384×384 blocks) for assessing open-world
396 scalability.

397 We compare Memo against Steve-1 and MrSteve base-
398 lines, and a variant, Memo w/global, which only uses
399 quadtree-based global planning. Steve-1 follows “Go Ex-
400 plore” instructions [2, 30], MrSteve combines count-based
401 exploration with VPT-Nav [17], and Memo implements the
402 proposed global–regional–local three-layer progressive ex-
403 ploration with VPT-Nav. Core metrics include map cover-
404 age rate and coverage efficiency.

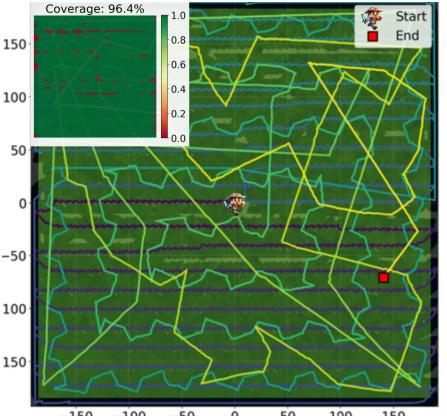
405 **Results and Discussion** Fig. 2 visually presents the ex-
406 ploration trajectories of MrSteve and Memo (ours) in a
407 simulated environment. MrSteve frequently converges to
408 local optima, leaving edges and key regions unexplored.
409 In contrast, Memo’s multi-layer exploration achieves more

410 complete coverage, reflecting effective coordination across
 411 global, regional, and local layers.

412 Tab. 1a and Tab. 1b reports coverage rates over time
 413 for MrSteve, Memo with global-only planning , and Memo
 414 across different map scales. Memo consistently achieves
 415 the highest coverage and efficiency. While global-only
 416 planning improves efficiency it can miss certain areas.
 417 Interestingly, as map size increases, MrSteve’s coverage
 418 slightly improves because the proportion of missed edge re-
 419 gions decreases relative to the total map.



(a) Exploration trajectory of Mrsteve(baseline).



(b) Exploration trajectory of Memo(ours).

Figure 2. The trajectory of exploration over 10,000 timesteps in a simulated environment on a 384×384 map. Under the same conditions, Memo has more than 20% of map coverage than Mrsteve’s and have smarter exploration strategy with more reasonable trajectory.

420 As shown in Tab. 1c , The results from real environ-
 421 ments corroborate those from the simulated experiments.
 422 On the small real map , Memo’s coverage rate reached 0.97
 423 after 10,000 steps, significantly outperforming Steve-1
 424 (0.19) and MrSteve (0.67), with respective improvements
 425 of 78% and 30%. Similarly, on the large real map, Memo
 426 achieved a coverage rate of 0.83 after 40,000 steps, while

Method	128×128			
	500	1000	2000	6000
Mrsteve	0.46	0.57	0.60	0.59
Memo*(ours)	0.65	0.87	0.92	0.93
Memo(ours)	0.66	0.85	0.93	0.99

(a) Map coverage on 128×128 map in simulated environment.

Method	384×384			
	1000	4000	10000	20000
Mrsteve	0.17	0.53	0.75	0.83
Memo*(ours)	0.19	0.72	0.83	0.83
Memo(ours)	0.19	0.73	0.96	0.98

(b) Map coverage on 384×384 map in simulated environment.

Method	128×128		384×384	
	5,000	10,000	20,000	40,000
Steve-1	0.11	0.19	0.06	0.13
Mrsteve	0.64	0.67	0.39	0.69
Memo*(ours)	0.72	0.75	0.44	0.74
Memo(ours)	0.71	0.97	0.44	0.83

(c) Map coverage of different exploration policies in real environment.

Memo* means just use global level search when exploring. Map coverage $\in [0, 1]$. Best results in each column are highlighted in **bold**.

Table 1. Map coverage of different exploration policies. We use different steps on small (128×128) and large (384×384) maps to measure map coverage of different methods. From the result, we can see Memo outperform other methods typically when the scale of map and timestep is large. And when timestep is small, Memo without whole scale exploration, just using global level methods, is more effective.

baseline models struggled to effectively cover the extensive map due to their limited exploration efficiency. These experimental findings confirm that Memo’s proposed exploration method is optimal, as Steve-1 exhibits low map coverage and efficiency. Although MrSteve’s map coverage improved compared to Steve-1, it remains prone to missing locations and falling into local optima.

Furthermore, an independent evaluation of the global coverage layer revealed that, compared to MrSteve’s count-based target method, it yields a higher coverage rate and faster coverage efficiency. However, it still exhibits sub-optimal exploration of local locations, underscoring the necessity of Memo’s multi-scale approach. These results collectively demonstrate that Memo’s progressive exploration framework effectively explores local positions and achieves more comprehensive coverage in complex, open-world environments.

427
428
429
430
431
432
433

434
435
436
437
438
439
440
441
442
443

444

4.2. Exploration and Navigation in Sequential Tasks

445 This experiment evaluates the agent’s ability to solve
 446 sparse-resource, long-horizon tasks that require returning
 447 to previously visited locations. We design an ABA-Sparse
 448 task in which the agent must complete an A–B–A se-
 449 quence—first locating resource A, then resource B, and
 450 finally returning to the initial A location. Successful execu-
 451 tion therefore depends on efficient exploration to find scarce
 452 targets and reliable memory to revisit earlier sites.

453 **Experimental Setup** We compare three
 454 baselines—Steve-1, MrSteve, and a memory-only variant
 455 Memo*—against our full Memo agent. Memo* retains
 456 the hierarchical memory module but removes multi-scale
 457 exploration and dynamic planning. Performance is mea-
 458 sured using task success rate and completion time. All
 459 experiments are conducted on a 128×128 sparse-resource
 460 map featuring ponds, forests, deserts, mountains, pits, and
 461 scattered ores that increase perception difficulty (Fig. 3).

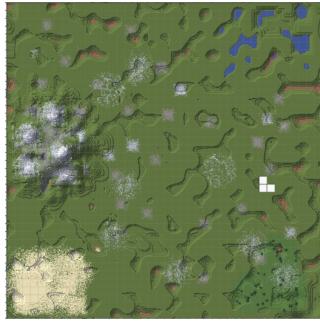


Figure 3. Real environment map in minecraft. Ponds are distributed in the upper right corner, forests in the lower right corner, high mountains on the left, and deserts in the lower left corner. Various minerals and deep pits are scattered across the map.

462 **Results and Discussion** As shown in Table 2, Steve-1
 463 achieves a 0% success rate, largely due to its inability to
 464 locate sparse targets or resume tasks from past experiences.
 465 MrSteve performs better but still suffers from limited ex-
 466 ploration and weak memory. Memo* further improves suc-
 467 cess rates, demonstrating the benefit of hierarchical memory
 468 alone.

469 The full Memo agent achieves the highest performance
 470 across all metrics. For example, in the scenario illus-
 471 trated in Fig. 4, MrSteve fails to recognize the pond un-
 472 til repeated exploration increases its field-of-view cov-
 473 erage, only then enabling water retrieval. In contrast,
 474 Memo’s VLM-enhanced semantic memory identifies and
 475 stores pond information immediately, enabling rapid return
 476 in the final “A” stage of the ABA sequence. While memory-
 477 only models help, their performance remains significantly
 478 below Memo’s. This indicates that efficient completion of

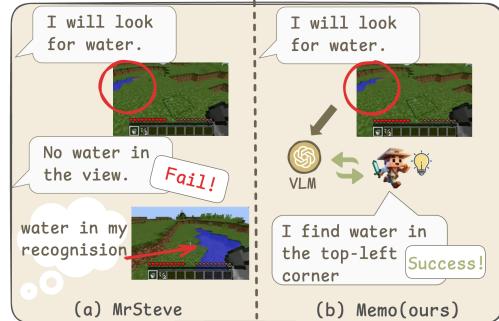


Figure 4. Exploration comparison in ABA sequence task. The left shows that MrSteve(baseline), based on PEM, fails to recognize the pond when view is limited. It only identifies the pond after its field of view is large enough. But Memo can recognize them both.

long-sequence tasks requires three components working in synergy: multi-scale exploration for rapid resource localization, hierarchical memory for accurate long-term recall, and dynamic planning for optimal execution order. Together, these modules address the core challenges of resource sparsity, unreliable recall, and weak task continuity, enabling robust performance in sequential and non-sequential sparse-resource tasks.

Method	Success Rate (%)	Average Timesteps
Steve-1	0	x
Mrsteve	32	8123
Memo*	47	8093
Memo(ours)	62	5981

(a) ABA sequential task result

Method	Success Rate (%)	Average Timesteps
Steve-1	0	x
Mrsteve	33	8040
Memo(ours)	77	4620

(b) ABC non-sequential task result

Memo* means without exploration.

Table 2. Result of success rates and average timesteps of different kinds of sparse task. Fig (a) represent the result of executing A task that has been performed historically. Memo is 1/4 faster than Mrsteve while achieving nearly twice success rate, indicating better utilization of historical information. Fig (b) also shows the efficiency of Memo compared to Mrsteve, because the sub-task sequency can be changed according to environment in our system.

4.3. Exploration and Navigation in Non-Sequential Tasks

To evaluate Memo’s autonomous decision-making in more complex sparse-resource scenarios, we designed the ABC non-sequential task. Here, task element C may appear while

479
480
481
482
483
484
485
486

487
488
489
490
491

492 the agent is executing task A, requiring the agent to dynamically
 493 adjust task priorities based on environmental observations and historical memory rather than following a fixed
 494 sequence. This setup tests the agent's capacity for adaptive
 495 decision-making and resource coordination.
 496

497 **Experimental Setup** We used the same baselines as in
 498 Section 4.2 and conducted experiments on a 128×128
 499 sparse-resource map, consistent with the sequential task
 500 setup.

501 **Results and Discussion** As shown in Table 2b, Steve-1
 502 fails completely (0% success), while MrSteve achieves a
 503 33% success rate with an average completion time of 8040
 504 steps, relying on limited memory and static exploration.

505 Memo substantially outperforms the baselines, achieving
 506 a 74% success rate and reducing the average completion
 507 time to 4620 steps. This improvement is enabled by
 508 the agent's dynamic planning module, which identifies that
 509 task C's resource lies along the current path during task A,
 510 dynamically reprioritizes tasks, and executes C first. By
 511 optimizing the sequence and leveraging multi-scale explora-
 512 tion for efficient resource localization, Memo minimizes
 513 redundant movements and task execution time. In con-
 514 trast, MrSteve's fixed-sequence strategy leads to repeated
 515 searches and reduced efficiency, highlighting the critical
 516 role of adaptive planning in complex, open-world tasks.

517 4.4. Ablation Study

518 To assess the individual contributions of Memo's core mod-
 519 ules, we conducted ablation experiments in the challenging
 520 non-sequential (ABC) task scenario on a 128×384 real map,
 521 using task success rate and average completion time as eval-
 522 uation metrics (Tab. 3).

Method	Success Rates(%)	Timesteps
Memo w/o Memory	31	8034
Memo w/o Exploration	45	7942
Memo w/o Planner	62	6981
Memo(ours)	77	4620

Table 3. Results of ablation experiment.

523 The results clearly demonstrate the critical role of
 524 each module. Removing the hierarchical memory module
 525 (Memo w/o Memory) severely impairs the agent's abil-
 526 ity to store and retrieve location information, reducing the
 527 success rate to 31% and nearly doubling task completion
 528 time. Excluding the multi-scale exploration module (Memo
 529 w/o Exploration) leads to inefficient local exploration and
 530 slower resource localization, yielding a 45% success rate

531 and longer completion time. Without the dynamic decision-
 532 making module (Memo w/o Planner), the agent executes
 533 tasks in a fixed sequence, increasing path redundancy and
 534 resulting in a 62% success rate.

535 The full Memo model, integrating all modules, achieves
 536 the highest success rate (77%) and shortest completion
 537 time (4620 steps). As illustrated in Fig. 5, Memo dy-
 538 namically reprioritizes tasks—for example, switching from
 539 "water-log-beef" to "beef-water-log" when encountering
 540 a cow—showcasing the synergistic operation of memory,
 541 exploration, and planning. These findings confirm that
 542 the hierarchical memory module provides accurate histori-
 543 cal knowledge, multi-scale exploration ensures efficient re-
 544 source localization, and dynamic planning enables optimal
 545 task scheduling, collectively forming Memo's core advan-
 546 tage in managing long-sequence tasks in open-world envi-
 547 ronments.

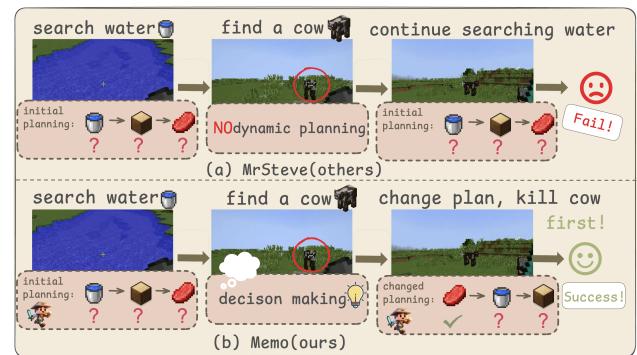


Figure 5. Planning comparison in ABC's non-sequence. The initial task is 'water-log-beef'. When the Memo(ours) first find a cow, it switches its plan to 'beef-water-log' and prioritizes 'killing a cow' to obtain beef. However, due to the lack of dynamic planning ability, other agents can only continue searching for water.

5. Conclusion

549 In this paper, we presented Memo, a novel open-world agent
 550 featuring a hierarchical memory architecture for robust ex-
 551 ploration, navigation, and task execution. Memo integrates
 552 three functional modules—explorer, memory, and plan-
 553 ner—to dynamically coordinate perception, memory, and
 554 action, enabling effective completion of sparse sequence
 555 tasks in large-scale environments. Experimental results
 556 demonstrate that this integration significantly enhances task
 557 success and efficiency, highlighting Memo's capability to
 558 handle both environment- and memory-dependent chal-
 559 lenges. Despite these advances, two limitations remain.
 560 First, reliance on GPT-4o incurs high computational costs,
 561 which could be mitigated with open-source alternatives.
 562 Second, our evaluation is currently limited to Minecraft,
 563 leaving cross-platform generalization as a promising direc-
 564 tion for future work.

565 **References**

- [1] Bowen Baker, Ilge Akkaya, Peter Zhokov, Joost Huizinga, Jie Tang, Adrien Ecoffet, Brandon Houghton, Raul Sampedro, and Jeff Clune. Video pretraining (vpt): Learning to act by watching unlabeled online videos. *Advances in Neural Information Processing Systems*, 35:24639–24654, 2022. [2](#)
- [2] Shaofei Cai, Bowei Zhang, Zihao Wang, Xiaojian Ma, Anji Liu, and Yitao Liang. Groot: Learning to follow instructions by watching gameplay videos. *arXiv preprint arXiv:2310.08235*, 2023. [1, 5](#)
- [3] Matthew Chang, Theophile Gervet, Mukul Khanna, Sriram Yenamandra, Dhruv Shah, So Yeon Min, Kavit Shah, Chris Paxton, Saurabh Gupta, Dhruv Batra, et al. Goat: Go to any thing. *arXiv preprint arXiv:2311.06430*, 2023. [2](#)
- [4] Devendra Singh Chaplot, Dhiraj Gandhi, Saurabh Gupta, Abhinav Gupta, and Ruslan Salakhutdinov. Learning to explore using active neural slam. *arXiv preprint arXiv:2004.05155*, 2020. [1, 2](#)
- [5] Devendra Singh Chaplot, Ruslan Salakhutdinov, Abhinav Gupta, and Saurabh Gupta. Neural topological slam for visual navigation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12875–12884, 2020. [2](#)
- [6] Or Dinari and Oren Freifeld. Revisiting dp-means: fast scalable algorithms via parallelism and delayed cluster creation. In *Uncertainty in Artificial Intelligence*, pages 579–588. PMLR, 2022. [4](#)
- [7] Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. Minedojo: Building open-ended embodied agents with internet-scale knowledge. *Advances in Neural Information Processing Systems*, 35:18343–18362, 2022. [4](#)
- [8] Jonathan Gray, Kavya Srinet, Yacine Jernite, Haonan Yu, Zhuoyuan Chen, Demi Guo, Siddharth Goyal, C Lawrence Zitnick, and Arthur Szlam. Craftassist: A framework for dialogue-enabled interactive agents. *arXiv e-prints*, pages arXiv–1907, 2019. [2](#)
- [9] William H Guss, Brandon Houghton, Nicholay Topin, Phillip Wang, Cayden Codel, Manuela Veloso, and Ruslan Salakhutdinov. Minerl: A large-scale dataset of minecraft demonstrations. *arXiv preprint arXiv:1907.13440*, 2019. [2](#)
- [10] Tomoyuki Kagaya, Thong Jing Yuan, Yuxuan Lou, Jayashree Karlekar, Sugiri Pranata, Akira Kinose, Koki Oguri, Felix Wick, and Yang You. Rap: Retrieval-augmented planning with contextual memory for multimodal llm agents. *arXiv preprint arXiv:2402.03610*, 2024. [2](#)
- [11] Jimmy Krozel and Dominick Andrisani. Navigation path planning for autonomous aircraft: Voronoi diagram approach. *Journal of Guidance, Control, and Dynamics*, 13(6):1152–1154, 1990. [4](#)
- [12] Zaijing Li, Yuquan Xie, Rui Shao, Gongwei Chen, Dongmei Jiang, and Liqiang Nie. Optimus-1: Hybrid multimodal memory empowered agents excel in long-horizon tasks. *Advances in neural information processing systems*, 37:49881–49913, 2024. [2](#)
- [13] Shalev Lifshitz, Keiran Paster, Harris Chan, Jimmy Ba, and Sheila McIlraith. Steve-1: A generative model for text-to-behavior in minecraft. *Advances in Neural Information Processing Systems*, 36:69900–69929, 2023. [2](#)
- [14] Zichuan Lin, Junyou Li, Jianing Shi, Deheng Ye, Qiang Fu, and Wei Yang. Juewu-mc: Playing minecraft with sample-efficient hierarchical reinforcement learning. *arXiv preprint arXiv:2112.04907*, 2021. [1](#)
- [15] Junhyuk Oh, Satinder Singh, Honglak Lee, and Pushmeet Kohli. Zero-shot task generalization with multi-task deep reinforcement learning. In *International Conference on Machine Learning*, pages 2661–2670. PMLR, 2017. [1](#)
- [16] R OpenAI. Gpt-4 technical report. *arxiv 2303.08774. View in Article*, 2(5):1, 2023. [1, 4](#)
- [17] Junyeong Park, Junmo Cho, and Sungjin Ahn. Mrsteve: Instruction-following agents in minecraft with what-where-when memory. In *The Thirteenth International Conference on Learning Representations*, 2025. [1, 2, 5](#)
- [18] Yiran Qin, Enshan Zhou, Qichang Liu, Zhenfei Yin, Lu Sheng, Ruimao Zhang, Yu Qiao, and Jing Shao. Mp5: A multi-modal open-ended embodied system in minecraft via active perception. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16307–16316. IEEE, 2024. [2](#)
- [19] Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov, Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, et al. A generalist agent. *arXiv preprint arXiv:2205.06175*, 2022. [1](#)
- [20] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36:8634–8652, 2023. [2](#)
- [21] Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M Sadler, Wei-Lun Chao, and Yu Su. Llm-planner: Few-shot grounded planning for embodied agents with large language models. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2998–3009, 2023. [2](#)
- [22] Zhiyuan Sun, Haochen Shi, Marc-Alexandre Côté, Glen Berseth, Xingdi Yuan, and Bang Liu. Enhancing agent learning through world dynamics modeling. *arXiv preprint arXiv:2407.17695*, 2024. [2](#)
- [23] Haoran Tang, Rein Houthooft, Davis Foote, Adam Stooke, OpenAI Xi Chen, Yan Duan, John Schulman, Filip DeTurck, and Pieter Abbeel. # exploration: A study of count-based exploration for deep reinforcement learning. *Advances in neural information processing systems*, 30, 2017. [2](#)
- [24] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. [1](#)
- [25] Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*, 2023. [2](#)
- [26] Zihao Wang, Shaofei Cai, Anji Liu, Yonggang Jin, Jiningbing Hou, Bowei Zhang, Haowei Lin, Zhaofeng He, Zilong

- 679 Zheng, Yaodong Yang, et al. Jarvis-1: Open-world multi-
680 task agents with memory-augmented multimodal language
681 models. *IEEE Transactions on Pattern Analysis and Ma-*
682 *chine Intelligence*, 2024. 2
- 683 [27] Brian Yamauchi. Frontier-based exploration using multiple
684 robots. In *Proceedings of the second international confer-*
685 *ence on Autonomous agents*, pages 47–53, 1998. 2
- 686 [28] Jesse Zhang, Jiahui Zhang, Karl Pertsch, Ziyi Liu, Xiang
687 Ren, Minsuk Chang, Shao-Hua Sun, and Joseph J Lim. Boot-
688 strap your own skills: Learning to solve new tasks with large
689 language model guidance. *arXiv preprint arXiv:2310.10021*,
690 2023. 2
- 691 [29] Tianjun Zhang, Huazhe Xu, Xiaolong Wang, Yi Wu, Kurt
692 Keutzer, Joseph E Gonzalez, and Yuandong Tian. Noveld: A
693 simple yet effective exploration criterion. *Advances in Neu-*
694 *ral Information Processing Systems*, 34:25217–25230, 2021.
695 2
- 696 [30] Enshen Zhou, Yiran Qin, Zhenfei Yin, Yuzhou Huang,
697 Ruimao Zhang, Lu Sheng, Yu Qiao, and Jing Shao. Mine-
698 dreamer: Learning to follow instructions via chain-of-
699 imagination for simulated-world control. *arXiv preprint*
700 *arXiv:2403.12037*, 2024. 5
- 701 [31] Xizhou Zhu, Yuntao Chen, Hao Tian, Chenxin Tao, Wei-
702 jie Su, Chenyu Yang, Gao Huang, Bin Li, Lewei Lu, Xiao-
703 gang Wang, et al. Ghost in the minecraft: Generally capable
704 agents for open-world environments via large language mod-
705 els with text-based knowledge and memory. *arXiv preprint*
706 *arXiv:2305.17144*, 2023. 2