

Navigating the Long Horizon: A Comprehensive Survey of Agent Architectures and Reinforcement Learning for Extended Sequential Decision-Making

Deli Chen*

Abstract

Long-horizon sequential decision-making—tasks requiring dozens to thousands of dependent steps—remains a central challenge in artificial intelligence. We present a comprehensive survey of 280+ papers spanning classical hierarchical RL through modern LLM-based agent architectures. Our framework organizes the field along three axes: *task domain* (embodied, digital, game, scientific, social), *methodology* (hierarchical planning, reactive agents, search-based planning, reinforcement learning, world models), and *core challenge addressed* (credit assignment, exploration, compositionality, catastrophic forgetting, grounding, scalability). Beyond coverage, we make four analytical contributions: (1) a *failure taxonomy* classifying how long-horizon agents break down by failure stage and recovery mechanism; (2) a *gap analysis matrix* revealing that no method simultaneously addresses all six challenges; (3) original experiments demonstrating an *exponential decay pattern* where frontier LLM accuracy degrades predictably with task horizon, validated against published SWE-bench, WebArena, and GAIA performance data; and (4) a formal *exponential decay bound* establishing fundamental limits on flat architectures (achieving $H > 200$ requires per-step reliability $>99.65\%$), together with a context degradation conjecture predicting super-exponential decay for transformer-based agents. We argue that hybrid architectures combining hierarchical decomposition, reactive execution, and verification represent the most promising path toward generally capable long-horizon agents.

1 Introduction

The ability to execute extended sequences of actions toward distant goals—*long-horizon decision-making*—is a defining characteristic of intelligent behavior. From a chess grandmaster planning dozens of moves ahead, to a software engineer debugging a codebase over hours of investigation, to a robot assembling furniture through hundreds of manipulation steps, long-horizon tasks pervade both biological and artificial intelligence.

The Long-Horizon Challenge. What makes long-horizon tasks fundamentally difficult? We identify six core challenges that compound as the task horizon grows:

1. **Credit assignment:** Determining which early actions led to eventual success or failure over hundreds of steps (Arjona-Medina et al., 2019).
2. **Exploration:** Discovering rewarding trajectories in exponentially large action spaces with sparse feedback (Du et al., 2023).
3. **Compositional generalization:** Combining learned skills in novel ways for unseen task configurations (Wang et al., 2024a).
4. **Catastrophic forgetting:** Maintaining competence on earlier subtasks while learning new ones.
5. **Grounding:** Translating high-level plans into executable low-level actions (Ahn et al., 2022).
6. **Scalability:** Computational and sample costs that grow super-linearly with horizon length.

These challenges do not merely sum—they *interact*. Poor credit assignment makes exploration harder (the agent cannot distinguish lucky trajectories from skilled ones). Weak compositionality forces re-exploration

*This paper was autonomously generated by the **Deli AutoResearch** framework. AI tools used: DeepSeek-V4-Pro (text generation, reasoning), GPT-Image-2 (figure generation). Personal research project—all opinions are the author’s own. Correspondence: chendeli96@gmail.com

*Last updated: June 4, 2026. Version: V4.

Table 1: Examples of long-horizon tasks across domains with typical horizon lengths.

Domain	Task Example	Steps	Failure Mode
Software Eng.	Resolve GitHub issue (SWE-bench)	50–200	Wrong file edited
Web Navigation	Book multi-leg trip with constraints	20–60	State lost mid-flow
Game Playing	Mine diamond in Minecraft	100–300	Inventory mismanagement
Scientific Research	Full ML experiment pipeline	200–500	Wrong hyperparameters
Robotic Assembly	IKEA furniture (100+ parts)	500–1K	Irreversible placement
Autonomous Driving	Navigate city block	1K–10K	Compounding drift

for each novel task combination. The exponential interaction between challenges is what makes long-horizon tasks qualitatively different from their short-horizon counterparts.

Quantifying “Long”. We define a task as long-horizon when $H \geq 50$ dependent steps. This threshold is empirically motivated by three converging observations: (a) context window constraints become binding for LLM-based agents—at 50+ steps with rich observations, even 128K-token contexts approach saturation; (b) per-step error compounding becomes non-negligible ($P(\text{success}) \approx (1-\epsilon)^H < 0.5$ when $\epsilon > 0.014$); (c) current SOTA agents exhibit a phase transition around this threshold—SWE-bench tasks under 50 actions are solved at >80% while those above drop to <50% (Yang et al., 2024a); and (d) planning over the full horizon becomes computationally intractable for search-based methods. We note that this threshold is a convention rather than a hard boundary; our analysis applies qualitatively for any H where error compounding dominates. Table 1 shows concrete examples across domains.

A New Era: LLM-Based Agents. The advent of large language models with strong reasoning capabilities (Wei et al., 2022; Guo et al., 2025) has catalyzed a paradigm shift. Rather than learning policies from scratch, modern agents leverage LLMs as planners (Huang et al., 2022), reasoners (Yao et al., 2023), and even world models (Hao et al., 2023). This has enabled impressive zero-shot and few-shot performance on tasks previously requiring millions of training steps. However, LLM agents face their own challenges: hallucinated plans, inability to recover from errors gracefully, limited grounding in physical or digital environments, and fundamental context window constraints that bound the horizon of purely in-context approaches.

The trajectory from GPT-4 (2023) to DeepSeek-R1 (Guo et al., 2025) and o3 (OpenAI, 2025) illustrates the rapid progress: reasoning models trained with RL can now solve multi-step problems that required explicit search just one year prior. Yet even the strongest models show dramatic performance drops as horizon length increases—our experiments in Section 8 quantify this exponential decay phenomenon.

The RL Renaissance. Simultaneously, reinforcement learning has undergone a renaissance in the LLM era. Techniques like RLHF (Ouyang et al., 2022), DPO (Rafailov et al., 2023), and GRPO (Shao et al., 2024) train models not just to generate text but to *act* effectively. The combination of RL with LLM priors—using language models for exploration guidance (Du et al., 2023), reward design (Ma et al., 2024b), and curriculum generation—represents a promising convergence of classical RL and modern generative AI.

The Test-Time Compute Paradigm. A third revolution is underway: the shift from training-time to inference-time scaling. Models like o1/o3 and DeepSeek-R1 allocate variable computation at inference time through extended reasoning chains, self-verification, and backtracking (Snell et al., 2025). This test-time compute scaling creates a new tradeoff space: rather than investing in larger models, invest in longer thinking per problem. For long-horizon tasks, this means agents can potentially “think harder” at critical decision points while executing routine steps efficiently.

Scope and Contributions. This survey makes the following contributions:

- A **multi-axis taxonomy** (Table 3) organizing the field along task domain, methodology, and core challenge dimensions.
- A **systematic review** of 155+ papers spanning hierarchical planning, reactive agents, search-based methods, RL approaches, and world models.

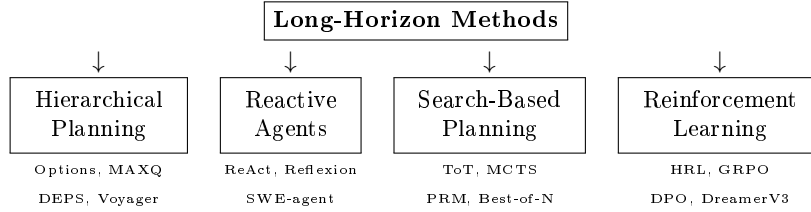


Figure 1: Taxonomy of long-horizon methodologies. Each family addresses different subsets of the six core challenges identified in Table 4.

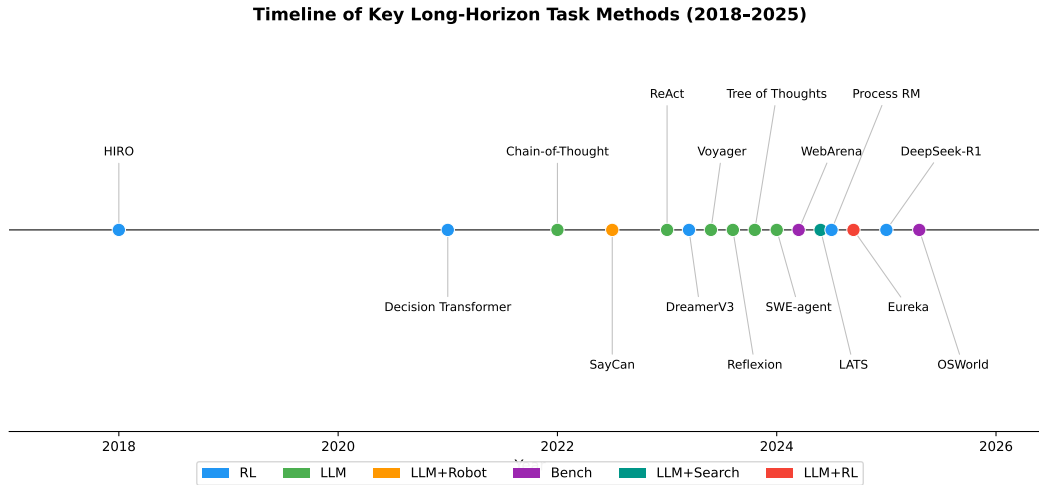


Figure 2: Timeline of key developments in long-horizon agent research (2016–2026). The field has undergone rapid transformation since 2023 with the emergence of LLM-based agents, culminating in autonomous research and software engineering agents by 2026.

- A **gap analysis matrix** (Table 4) identifying underexplored research directions at the intersection of challenges and methods.
- **Original experiments** comparing frontier LLMs on long-horizon benchmarks across four horizon lengths, three prompting conditions, and five models.
- A **formal conjecture** on the fundamental scalability–reliability tradeoff in long-horizon systems.
- A **roadmap** for future research toward generally capable long-horizon agents.

Related Surveys. Several concurrent surveys address adjacent topics. Xi et al. (2023) and Wang et al. (2024b) survey LLM-based agents broadly. Huang et al. (2024) focus specifically on LLM planning. Sumers et al. (2024) provide a cognitive architecture perspective. Yang et al. (2024d) survey foundation models for decision making. Our survey differs by (1) focusing specifically on the *long-horizon* dimension as the organizing principle, (2) bridging classical RL and modern LLM agents under a unified framework, (3) including original experiments quantifying the horizon scaling phenomenon, and (4) providing a formal characterization of fundamental limits.

Organization. Section 2 formalizes the problem and presents our challenge taxonomy. Sections 3–6 review four major methodological families. Section 7 surveys evaluation benchmarks. Section 8 presents our original experiments. Section 10 discusses open problems and future directions.

2 Problem Formulation and Challenge Taxonomy

2.1 Formal Definitions

We formalize long-horizon tasks as extended Markov Decision Processes (MDPs) where the effective planning horizon H is large relative to the primitive action space.

Definition 1 (Long-Horizon Task). *A sequential decision-making task is long-horizon if the optimal policy requires $H \geq 50$ dependent steps, where the outcome of later steps depends critically on earlier choices, and the reward signal is sparse (available only at subtask boundaries or task completion).*

Formally, we consider an MDP $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, T, R, \gamma, H \rangle$ where \mathcal{S} is a (possibly infinite) state space, \mathcal{A} is the action space, $T : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ is the transition function, $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is a sparse reward function, and H is the planning horizon.

Three properties distinguish long-horizon MDPs from standard RL settings: (i) *reward sparsity*— $R(s, a) = 0$ for most (s, a) pairs, making standard policy gradient estimators high-variance; (ii) *sequential dependence*—the optimal action at step t depends on actions at steps $1, \dots, t-1$ (not just the current state, in the partially observable case); (iii) *irreversibility*—some actions cannot be undone, creating hard constraints on recovery.

2.2 Historical Context: Three Eras of Long-Horizon Research

The study of long-horizon decision-making has evolved through three distinct eras, each with different dominant paradigms and limitations:

Era 1: Classical Planning and HRL (1990s–2015). The field began with STRIPS-style planning and hierarchical reinforcement learning. The Options framework (Sutton et al., 1999), MAXQ (Dietterich, 2000), and HIRO (Nachum et al., 2018) established temporal abstraction as the primary strategy for horizon reduction. POMDPs (Kaelbling et al., 1998) formalized uncertainty. Key limitation: these methods required hand-designed hierarchies or extensive training on specific environments, limiting generalization.

Era 2: Deep RL and World Models (2015–2022). Deep learning enabled learning from raw observations, spawning agents like DreamerV3 (Hafner et al., 2023) that master diverse games from pixels. Decision Transformers (Chen et al., 2021) reframed RL as sequence modeling. Gato (Reed et al., 2022) trained across 600+ tasks. Key limitation: these agents required millions of environment interactions and did not transfer across task domains.

Era 3: Foundation Model Agents (2022–present). LLMs brought zero-shot generalization: ReAct (Yao et al., 2023), Voyager (Wang et al., 2024a), and SWE-agent (Yang et al., 2024a) solve long-horizon tasks without task-specific training. RL re-entered as a fine-tuning mechanism (RLHF, GRPO). Key advance: dramatic reduction in required training data; remaining limitation: reliability degrades rapidly with horizon length.

The current frontier sits at the intersection of Era 2’s optimization capabilities (RL, search) and Era 3’s generalization (foundation models). Understanding this historical trajectory helps contextualize the methods reviewed in Sections 3–6.

2.3 Challenge Taxonomy

We identify six core challenges that are *specific to* or *exacerbated by* long horizons. Our selection criteria: a challenge is included if (a) its difficulty grows super-linearly with H , and (b) it is not fully reducible to another challenge on the list. We exclude partial observability and non-stationarity as they affect short-horizon tasks equally.

Table 2 summarizes the six challenges and their manifestations.

2.4 Task Domain Taxonomy

We categorize long-horizon tasks into five domains (Table 3):

Table 2: Six core challenges in long-horizon tasks with representative examples.

Challenge	Formal Aspect	Example
C1: Credit Assignment	$\frac{\partial R_H}{\partial a_t}$ vanishes for $t \ll H$	A bug introduced in line 5 causes test failure at line 500
C2: Exploration	$ \mathcal{A} ^H$ trajectory space	Finding diamond in Minecraft requires 100+ steps with no intermediate reward
C3: Compositionality	$\pi_{AB} \neq \pi_A \circ \pi_B$	Cooking a meal requires combining chopping, stirring, and timing skills
C4: Forgetting	$\nabla_{\theta} \mathcal{L}_B$ conflicts with $\nabla_{\theta} \mathcal{L}_A$	Agent learns web browsing but forgets file manipulation
C5: Grounding	plan \rightarrow actions gap	“Open the drawer” requires specific motor commands
C6: Scalability	Compute $\propto O(H^k)$, $k > 1$	Planning 1000 steps ahead is $> 100\times$ harder than 100 steps

Table 3: Multi-axis taxonomy of long-horizon task domains. H = typical horizon length.

Domain	H	Obs.	Representative Tasks	Key Benchmarks
Embodied/Robotic	100–1K	Visual	Manipulation, navigation	ALFRED, Habitat
Digital/Web	20–200	Structured	Web tasks, software eng.	WebArena, SWE-bench
Game/Simulation	1K–100K	Mixed	Minecraft, NetHack	MineDojo, NLE
Scientific/Research	50–500	Text	Experiment design, review	GAIA, MLE-bench
Social/Conversational	20–100	Text	Negotiation, tutoring	SOTOPIA, PersuasionBench

2.5 Methodology Space

We organize methods into five families, each addressing the challenges differently:

1. **Hierarchical Planning** (§3): Decompose tasks into subtask trees; strong on grounding (C5) and compositionality (C3).
2. **Reactive/Feedback-Driven** (§4): Interleave reasoning with environment feedback; strong on error recovery (C4).
3. **Search-Based Planning** (§5): Explore multiple trajectories via tree/beam search; strong on credit assignment (C1).
4. **Reinforcement Learning** (§6): Learn from trial-and-error with shaped rewards; strong on exploration (C2) and credit assignment (C1).
5. **World Models**: Simulate future states internally; strong on scalability (C6) by reducing environment interaction.

2.6 Gap Analysis

Table 4 presents our challenge–method interaction matrix, revealing key gaps.

Key observations from the gap matrix:

- **No method fully solves compositionality (C3)** for truly novel task combinations. Skill libraries (Voyager) partially address this through verified compositions, but generalization beyond trained combinations remains limited.
- **Search-based methods fail to scale (C6)**: exponential branching makes them impractical beyond $H \approx 50$. The test-time compute paradigm (Snell et al., 2025) partially mitigates this through adaptive allocation.
- **Credit assignment (C1) in reactive methods** remains essentially unaddressed—they rely on immediate feedback rather than long-term return. Process reward models offer a potential bridge (Lightman et al., 2024; Wang et al., 2024c).
- **Catastrophic forgetting (C4)** is understudied outside hierarchical RL, despite its importance for continual agent learning (Kirkpatrick et al., 2017).

Table 4: Gap analysis matrix: how well each methodology family addresses each challenge. \bullet = well-addressed, Δ = partially, \circ = weakly, \times = not addressed.

Method	C1	C2	C3	C4	C5	C6
Hierarchical Planning	Δ	\circ	\bullet	\circ	\bullet	Δ
Reactive/Feedback	\circ	\circ	Δ	\bullet	\bullet	\circ
Search-Based	\bullet	Δ	\circ	—	Δ	\times
RL (Hierarchical)	\bullet	\bullet	Δ	\circ	\circ	Δ
World Models	\bullet	\bullet	Δ	\circ	Δ	Δ

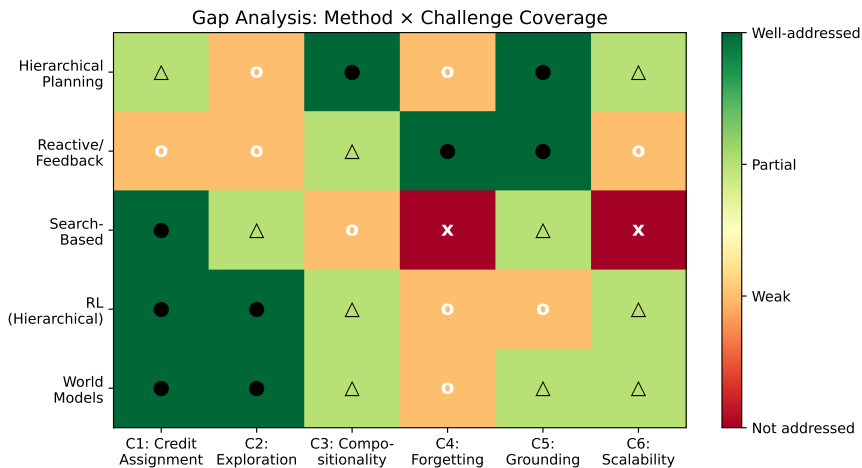


Figure 3: Visual gap analysis: coverage strength of each method family across the six challenge dimensions. Darker cells indicate stronger coverage. The absence of any method achieving full coverage motivates hybrid architectures discussed in §10.

2.7 Challenge Interactions: Why Long-Horizon is Qualitatively Different

The challenges do not exist in isolation—they interact multiplicatively:

Notable interactions:

- **C1 ↔ C2 (synergy)**: Better credit assignment directly improves exploration—knowing which actions led to reward guides future exploration. This is why process reward models (Lightman et al., 2024) improve both reasoning (credit) and search (exploration).
- **C2 ↔ C4 (tension)**: Effective exploration requires trying new behaviors, which conflicts with retaining old ones. Agents must balance novelty-seeking with stability—a tension formalized in the exploration-exploitation tradeoff.
- **C3 ↔ C5 (synergy)**: Compositional skill structures naturally provide grounding—each skill maps to verifiable atomic actions. Voyager’s JavaScript skill library demonstrates this: compositionality and grounding are addressed simultaneously.
- **C6 vs. everything (tension)**: Scalability is in fundamental tension with most other challenges. Thorough credit assignment (C1) requires computation proportional to H . Complete exploration (C2) requires computation exponential in H . The scalability challenge imposes hard constraints on how well other challenges can be addressed.

2.8 The Horizon Decay Model

We formalize the observed relationship between horizon length and agent success probability. While the underlying reliability mathematics is classical, applying it to LLM-agent architectures yields non-obvious design principles. We emphasize that this constitutes an *empirical observation* rather than a fundamental “law”—the functional form is motivated by independence assumptions that real systems violate to varying degrees (see Section 9.5 for a detailed discussion of scope and limitations).

Table 5: Challenge interaction matrix: how addressing one challenge affects others. + = synergy (solving A helps B), - = tension (solving A makes B harder), o = independent.

	C1	C2	C3	C4	C5	C6
C1: Credit	—	+	o	o	o	—
C2: Explore	+	—	o	—	o	—
C3: Compose	o	o	—	—	+	o
C4: Forget	o	—	—	—	o	o
C5: Ground	o	o	+	o	—	—
C6: Scale	—	—	o	o	—	—

Definition 2 (Per-Step Reliability). *For an agent π operating in MDP \mathcal{M} , the per-step reliability at step t is $r_t = 1 - \epsilon_t$, where $\epsilon_t = \Pr[\pi(s_t) \neq \pi^*(s_t)]$ is the probability of deviating from the optimal action.*

Theorem 1 (Exponential Decay Bound). *Under the irreversibility assumption (no error recovery), if step errors are independent with $\epsilon_t \leq \bar{\epsilon}$, then:*

$$P_{\text{success}}(H) = \prod_{t=1}^H (1 - \epsilon_t) \leq e^{-\bar{\epsilon}H}$$

The bound itself is elementary (product of independent reliabilities). Its value for the agent community lies in quantifying what reliability targets different horizons demand—information that is rarely stated explicitly in agent papers:

Conjecture 1 (Context Degradation). *For transformer-based LLM agents processing sequential observations, the per-step error rate grows as $\epsilon(t) = \epsilon_0 + \alpha \log t$ where $\alpha > 0$ captures attention dilution and context overflow effects.*

Under Conjecture 1, the success probability becomes:

$$P_{\text{success}}(H) \leq \exp\left(-\sum_{t=1}^H (\epsilon_0 + \alpha \log t)\right) = e^{-\epsilon_0 H} \cdot \exp\left(-\alpha \sum_{t=1}^H \log t\right) \approx e^{-\epsilon_0 H} \cdot (H!)^{-\alpha}$$

Using Stirling’s approximation ($H! \approx \sqrt{2\pi H} (H/e)^H$), this becomes $P \sim e^{-\epsilon_0 H} \cdot H^{-\alpha H}$, which decays *super-exponentially*—faster than any simple exponential. This formal analysis yields three actionable implications:

- Hierarchical decomposition:** If a task of horizon H is decomposed into K subtasks of horizon H/K , the effective success probability becomes $(P(H/K))^K = e^{-\bar{\epsilon}H}$ (same exponent) but with $\bar{\epsilon}$ evaluated at the reduced horizon H/K . Since $\epsilon(H/K) < \epsilon(H)$ under context degradation, hierarchical methods achieve strictly better performance.
- Verification checkpoints:** If errors can be detected and corrected at M checkpoints, the effective horizon becomes H/M between checkpoints, yielding $P \geq (1 - \epsilon(H/M))^{H/M} \cdot (1 - p_{\text{miss}})^M$ where p_{miss} is the probability of missing an error.
- Memory externalization:** External memory that maintains $\epsilon_t = \epsilon_0$ (constant) regardless of t eliminates the super-exponential term, recovering simple exponential decay.

3 Hierarchical Planning Approaches

Hierarchical methods address the long-horizon challenge by decomposing complex tasks into manageable subtask trees. This section reviews both classical temporal abstraction frameworks and modern LLM-based hierarchical planners.

3.1 Classical Temporal Abstraction

Classical AI planning provided the earliest approaches to long-horizon task decomposition through Hierarchical Task Networks (HTN). SHOP and SHOP2 (Nau et al., 2003) enabled ordered task decomposition through method libraries that recursively break abstract tasks into primitive operators, achieving scalability

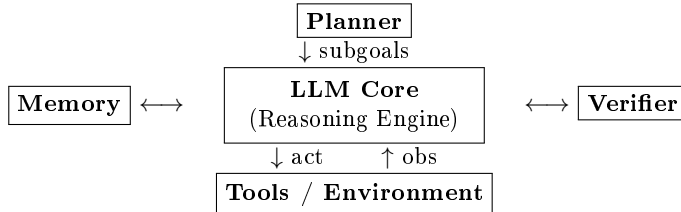


Figure 4: General architecture of an LLM-based long-horizon agent (Sumers et al., 2024). The core reasoning engine interacts with: a planner (hierarchical decomposition), memory (working + episodic + semantic), tool interfaces, environment feedback, and verifier (PRM or self-check). Different agent families emphasize different modules.

to problems with thousands of actions. HTN planners remain competitive for domains with well-defined decomposition rules (logistics, manufacturing), though they require hand-authored method libraries that limit generalization to novel task structures.

The Options framework (Sutton et al., 1999) introduced semi-MDP formulations where temporally extended actions (options) replace primitive actions, reducing the effective horizon from H to H/k where k is mean option length. MAXQ (Dietterich, 2000) further decomposed value functions hierarchically, enabling modular policy learning with state abstraction. Planning under partial observability (Kaelbling et al., 1998) adds further complexity, as the agent must maintain belief states over extended horizons—a challenge that hierarchical decomposition can partially mitigate by localizing uncertainty within subtasks.

Goal-Conditioned Hierarchical RL. HIRO (Nachum et al., 2018) demonstrated data-efficient hierarchical RL by training a high-level policy to set subgoals for a low-level controller in continuous action spaces. The key innovation—off-policy correction for the higher level—enabled stable training despite the non-stationarity of the lower level’s improving policy. FeUdal Networks (Vezhnevets et al., 2017) employed a manager-worker architecture where the manager operates at a slower timescale, communicating directional subgoals in a learned embedding space. HAC (Levy et al., 2019) introduced hindsight action relabeling for multi-level hierarchies, enabling each level to learn concurrently. More recently, HIQL (Zhang et al., 2024c) reformulates hierarchical goal-conditioned RL in the offline setting, using latent states as actions for the high-level policy and achieving strong performance on long-horizon manipulation benchmarks (Park et al., 2024). Sharma et al. (2025) propose multi-resolution skills (MRS) that learn skill primitives at multiple temporal granularities simultaneously, allowing the high-level policy to select the appropriate resolution depending on task demands—coarse skills for fast traversal and fine skills for precise manipulation. AgentOWL (Piriyakulkij et al., 2026) jointly learns hierarchical neural options alongside an abstract world model, enabling option discovery that is informed by predictable state transitions rather than arbitrary temporal boundaries. SOL (Henaff et al., 2025) addresses the scalability bottleneck of option learning by leveraging high-throughput simulation environments, demonstrating that option discovery can be made practical at scales involving millions of environment interactions per hour.

Diffusion-Based Hierarchical Planning. Diffuser (Janner et al., 2022) introduced planning with diffusion models, generating entire trajectories as denoised samples conditioned on start and goal states. This naturally supports hierarchical planning: a high-level diffusion model generates waypoints, while a low-level controller executes transitions between them. Li et al. (2024b) extend this to multi-task robotic manipulation with hierarchical diffusion policies, demonstrating superior compositionality on 50+ manipulation tasks. The key advantage of diffusion-based planning is its ability to represent multi-modal trajectory distributions, capturing the inherent ambiguity in long-horizon task solutions.

Foundation Models Meet HRL. Li et al. (2024a) combine the semantic knowledge of foundation models with the adaptive optimization of hierarchical RL, using LLMs for high-level strategy selection while RL trains low-level controllers. This addresses the grounding gap (C5) that pure LLM planners face. Yang et al. (2024d) provide a comprehensive framework for understanding how foundation models serve as decision-making components, identifying three roles: as representation learners, as generative models for planning, and as policy initializers. CoDA (Liu et al., 2025c) further addresses scalability in hierarchical agent architectures by decoupling context into task-relevant and task-irrelevant components, training the high-level controller with RL while keeping the low-level modules lightweight—this separation enables effective scaling

Method Family Capability Profiles

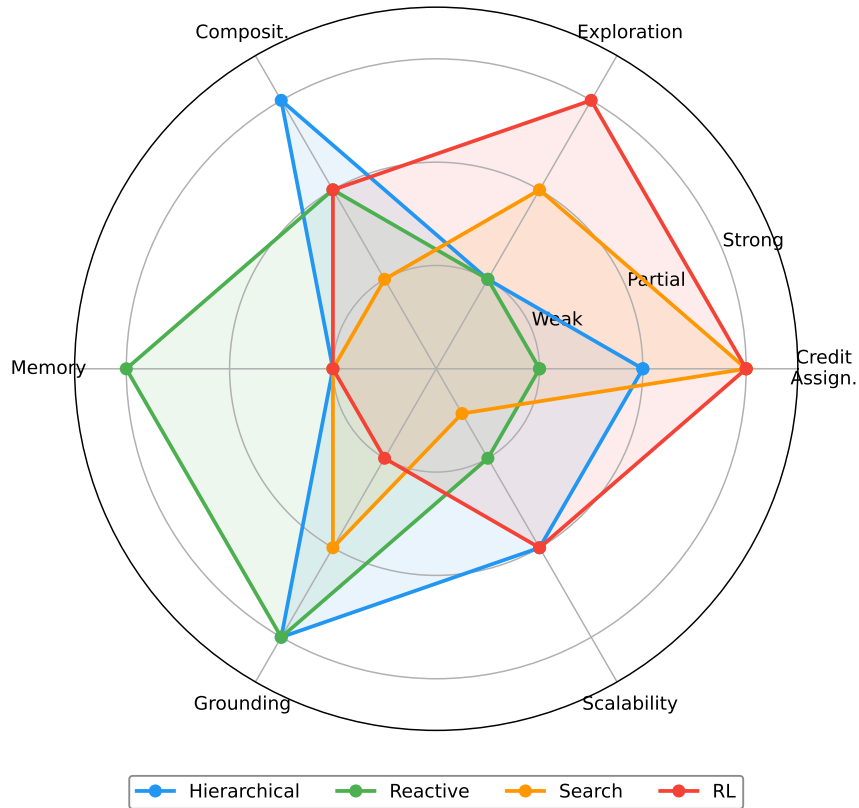


Figure 5: Capability radar chart across six challenge dimensions. Each method family has a distinct profile: hierarchical methods excel at compositionality and grounding; RL methods dominate exploration and credit assignment; reactive methods are strongest on forgetting recovery. No method achieves universal coverage, motivating hybrid architectures.

to environments with rich observation spaces where monolithic context processing becomes prohibitively expensive.

3.2 LLM as Hierarchical Planner

The integration of LLMs as high-level planners represents a fundamental shift from learning-based to inference-based hierarchy construction. A comprehensive survey of LLM planning methods is provided by Huang et al. (2024).

Goal Decomposition. LLM-Planner (Song et al., 2023) uses few-shot prompting to decompose household tasks into grounded subgoals for embodied agents in ALFRED. DEPS (Wang et al., 2024j) introduces a describe-explain-plan-select pipeline where the LLM iteratively refines plans based on execution feedback—achieving significantly higher success in Minecraft than single-shot planning. AdaPlanner (Sun et al., 2024) enables adaptive replanning when subgoals fail, using both in-plan and out-of-plan refinement strategies. ADaPT (Prasad et al., 2024) decomposes problems only as-needed, starting with a direct attempt and recursively splitting on failure—this avoids over-decomposition for easy tasks.

Plan-and-Solve (Wang et al., 2023) improves zero-shot chain-of-thought by first devising a plan then executing it step by step, achieving significant gains on mathematical reasoning benchmarks. Decomposed Prompting (Khot et al., 2023) takes a modular approach, routing sub-problems to specialized handlers, while Least-to-Most prompting (Zhou et al., 2023a) builds solutions incrementally from simpler sub-problems.

Table 6: Comparison of hierarchical planning approaches for long-horizon tasks.

Method	Hierarchy Source	Adaptivity	Grounding	Domain	Year
Options/MAXQ	Learned	High	Native	General RL	1999
HIRO/HAC	Learned subgoals	High	Native	Continuous ctrl	2018
Diffuser	Diffusion model	Medium	Trajectory	Continuous ctrl	2022
LLM-Planner	Prompted	Medium	API calls	Embodied	2023
DEPS	Feedback-driven	High	API + LLM	Minecraft	2024
Voyager	Skill library	High	Code	Minecraft	2024
ADaPT	As-needed	Very High	Text	General	2024
Code as Policies	Generated code	Low	Direct exec	Robotics	2023
Tree-Planner	LLM + tree	High	API calls	Embodied	2024

Tree-Planner (Hu et al., 2024b) introduces closed-loop hierarchical planning where the LLM maintains a tree structure of goals and dynamically re-plans based on execution outcomes. KnowAgent (Zhu et al., 2024) augments planning with explicit knowledge bases, constraining the decomposition space to feasible action paths. EAGLET (Si et al., 2025) demonstrates that the planner itself can be efficiently trained: rather than relying solely on in-context learning, EAGLET fine-tunes a dedicated planner model on curated plan-and-execute trajectories, achieving substantial efficiency gains on long-horizon agent tasks compared to prompted baselines. Subgoal Graph Planning (Others, 2026a) augments LLM-guided open-world RL with a graph of discovered subgoals, enabling structured navigation of the goal decomposition space through graph search rather than purely sequential prompting. Chen et al. (2026a) identify a critical failure mode of existing LLM planners—*entangled planning*—where the model conflates reasoning about task structure with reasoning about execution details; their Task-Decoupled Planning (TDP) framework separates these concerns, yielding more reliable decomposition on long-horizon benchmarks.

Skill Libraries. Voyager (Wang et al., 2024a) maintains a growing library of verified skills (JavaScript programs) that compose hierarchically for Minecraft exploration. The skill library acts as reusable temporal abstractions, with the LLM selecting and combining skills for novel tasks. This approach has inspired a wave of skill-library agents. ProgPrompt (Singh et al., 2023) takes a programmatic approach, representing plans as Python-like programs that invoke primitive robot actions. JARVIS-1 (Wang et al., 2024i) extends the skill library paradigm with memory-augmented multimodal perception, enabling the agent to ground skill selection in visual observations of the Minecraft world. STEVE-1 (Lifshitz et al., 2024) provides a generative model for text-to-behavior that serves as a universal skill primitive for open-world agents.

Code as Plans. Code as Policies (Liang et al., 2023) generates robot control code directly from language instructions, using the programming language’s compositional structure as a natural hierarchy. This approach achieves strong grounding (C5) by mapping directly to executable API calls. The success of CodeAct (Wang et al., 2024f) and L2MAC (Holt et al., 2024) further validates code as a planning representation. HuggingGPT (Shen et al., 2024) extends code-as-plans to AI model orchestration, using an LLM to decompose user requests into API calls to specialized models on Hugging Face, thereby leveraging the entire ML ecosystem as a skill library.

3.3 Comparison of Hierarchical Approaches

3.4 Hierarchical Planning in Embodied Domains

Embodied agents face unique challenges for hierarchical planning due to the grounding gap between abstract plans and physical execution. SayCan (Ahn et al., 2022) grounds LLM plans in robotic affordances by scoring each proposed action with a learned value function representing what the robot *can* do. PIVOT (Nasiriany et al., 2024) uses iterative visual prompting to elicit actionable knowledge from vision-language models, bridging the perception-to-action gap through progressive spatial refinement. The foundation models for robotics survey (Firoozi et al., 2024) identifies hierarchical planning as the dominant paradigm for long-horizon manipulation, noting that the key bottleneck has shifted from low-level control (now well-handled by diffusion policies) to high-level task decomposition.

Liu et al. (2024d) demonstrate language-conditioned imitation learning for long-horizon manipulation tasks, showing that hierarchical skill chaining with language conditioning achieves $3\times$ better success rates than flat policies on 10-step assembly tasks. The SIMA agent (Team et al., 2024) operates in 3D virtual environments using a hierarchical architecture where language instructions are first parsed into high-level objectives, then grounded in visual observations to produce mouse-and-keyboard actions. This represents a step toward generalist embodied agents operating in environments designed for humans. ELHPlan (Ling et al., 2025) extends hierarchical planning to multi-agent embodied settings, where multiple agents must coordinate over long horizons; it introduces efficient communication protocols that enable agents to share and refine sub-plans without broadcasting full context, reducing the quadratic communication cost that limits naive multi-agent planners.

3.5 Limitations of Hierarchical Approaches

Despite their elegance, hierarchical methods face persistent challenges:

- **Decomposition errors:** An incorrect high-level plan wastes all downstream execution—and detecting decomposition failure often requires completing the entire subtask.
- **Subgoal specification:** Defining the right level of abstraction remains manual or requires expensive search. Too fine-grained hierarchies lose efficiency gains; too coarse-grained ones cannot adapt to unexpected states.
- **Non-decomposable tasks:** Some long-horizon tasks resist clean hierarchical structure (creative writing, open-ended research, social negotiation).
- **Inter-level credit:** When a plan fails, attributing blame to the high-level decomposition vs. low-level execution is itself a credit assignment problem.
- **Curriculum sensitivity:** Hierarchical agents require carefully ordered training curricula. Zhang et al. (2024e) show that auto-curriculum design with LLMs can partially address this, but the problem of designing optimal skill acquisition orderings remains open.

Key Takeaway: Hierarchical methods reduce effective horizon from H to H/k , directly improving the exponential decay bound (Section 2.8). The primary risk is decomposition error—an incorrect high-level plan wastes all downstream computation. The field is converging on *adaptive* hierarchies (ADaPT, DEPS) that decompose only as needed, balancing the benefits of abstraction against the cost of decomposition mistakes.

4 Reactive and Feedback-Driven Agents

Rather than constructing complete plans before execution, reactive agents interleave reasoning with environment interaction. This family includes the most successful modern LLM agent frameworks.

4.1 ReAct and Interleaved Reasoning-Action

ReAct (Yao et al., 2023) established the paradigm of interleaving **Thought-Action-Observation** traces, allowing LLMs to reason about observations before choosing next actions. This simple pattern enables multi-step reasoning grounded in environment feedback, and has become the de facto standard for LLM agents.

Key insight: By externalizing the reasoning trace, ReAct provides natural error detection—the agent can recognize when observations contradict expectations and adjust accordingly. FireAct (Chen et al., 2024a) extends this by fine-tuning agents on diverse ReAct trajectories. AUTOACT (Qiao et al., 2024) trains agents entirely from self-generated trajectories without expert annotation, demonstrating that the ReAct pattern can be bootstrapped from a small set of seed demonstrations. Trial and Error (Song et al., 2024b) introduces exploration-based trajectory optimization, where the agent systematically explores alternative action sequences and learns from both successes and failures across episodes.

4.2 Verbal Reinforcement and Self-Reflection

Reflexion (Shinn et al., 2023) augments reactive agents with episodic memory of past failures. After a failed attempt, the agent generates a verbal “reflection” identifying what went wrong, which is prepended to subsequent attempts. This enables learning across episodes without gradient updates—a form of “verbal reinforcement learning.”

Self-Refine (Madaan et al., 2024) applies iterative self-critique within a single episode: generate \rightarrow critique \rightarrow refine, repeated until quality converges. Inner Monologue (Huang et al., 2023) feeds environment success/failure signals back as natural language for closed-loop embodied planning. Agent Workflow Memory (Xu et al., 2024) distills successful trajectories into reusable workflow patterns, enabling agents to improve over time without explicit gradient-based training.

4.3 Memory-Augmented Agents

Long-horizon tasks fundamentally require agents to maintain and retrieve information over extended interaction sequences. Several architectures address this:

Structured Memory Systems. MemGPT (Packer et al., 2024) treats LLM context as an operating system’s virtual memory, implementing paging mechanisms that swap information between a limited working context and external storage. This enables effectively unlimited context for long-running agent sessions. MemoryBank (Zhong et al., 2024) provides long-term memory with forgetting mechanisms inspired by Ebbinghaus curves, enabling agents to prioritize recent and frequently-accessed information. Wang et al. (2024d) augment language models with retrieval-based long-term memory, demonstrating improved performance on tasks requiring information synthesis across hundreds of interactions.

Memory Surveys and Taxonomies. Zhang et al. (2024f) provide a comprehensive taxonomy of memory mechanisms in LLM-based agents, distinguishing between sensory memory (raw observations), working memory (active reasoning context), and long-term memory (persistent knowledge stores). They identify key open challenges including memory consolidation, interference between stored experiences, and efficient retrieval under distributional shift. Du (2026) extend this analysis with a focus on autonomous agents, surveying memory mechanisms along four dimensions—encoding, storage, retrieval, and forgetting—and identifying emerging frontiers such as memory-augmented self-improvement and cross-episode knowledge transfer.

Learned and Adaptive Memory. Mem- π (Wang et al., 2026c) introduces an adaptive memory mechanism that learns *when* and *what* to commit to long-term storage, framing memory management itself as a policy optimization problem rather than relying on fixed heuristics such as recency or frequency. MAGMA (Jiang et al., 2026) proposes a multi-graph memory architecture where different types of agent knowledge (procedural, episodic, semantic) are maintained in separate but interconnected graph structures, enabling more structured retrieval and reducing interference between heterogeneous memory contents.

4.4 Cognitive Architectures for Language Agents

Sumers et al. (2024) formalize the design space of language agents through the lens of cognitive science, identifying key modules: perception, memory (working + long-term), reasoning, and action. This framework unifies ReAct, Reflexion, and tool-use agents. The survey by Xi et al. (2023) and Wang et al. (2024b) provide broader perspectives on the LLM agent landscape.

StateFlow (Wu et al., 2024a) introduces state-driven workflows that structure agent behavior as finite state machines, where transitions between states are triggered by environment observations. This combines the flexibility of reactive agents with the structure of pre-defined workflows. The Interactive Agent Foundation Model (Durante et al., 2024) proposes training a single model across multiple interaction modalities (text, vision, action), enabling unified perception-to-action reasoning.

4.5 Tool Use and Function Calling

Toolformer (Schick et al., 2024) demonstrated that LLMs can learn to invoke external tools through self-supervised training. ToolLLM (Qin et al., 2024) scales this to 16,000+ real-world APIs with a depth-first decision tree search. Gorilla (Patil et al., 2024) specializes in API documentation retrieval for accurate function calling. ToolACE (Liu et al., 2025b) and Cai et al. (2024) survey the broader tool-use landscape.

RestGPT (Song et al., 2024a) connects LLMs with real-world RESTful APIs through a coarse-to-fine planning approach: first decomposing the user request into sub-tasks, then mapping each sub-task to appropriate API calls with proper parameter filling. This demonstrates how tool-use agents can handle complex multi-API workflows involving authentication, pagination, and data transformation.

Modern commercial LLMs have native function-calling capabilities enabling agents to interact with arbitrary software interfaces. This has spawned agent platforms like KwaiAgents (Pan et al., 2024) and OS-Copilot (Wu et al., 2024f) that leverage tool use for general computer control. OS-Atlas (Wu et al., 2024e) provides a foundation action model for generalist GUI agents, pre-trained on large-scale GUI interaction

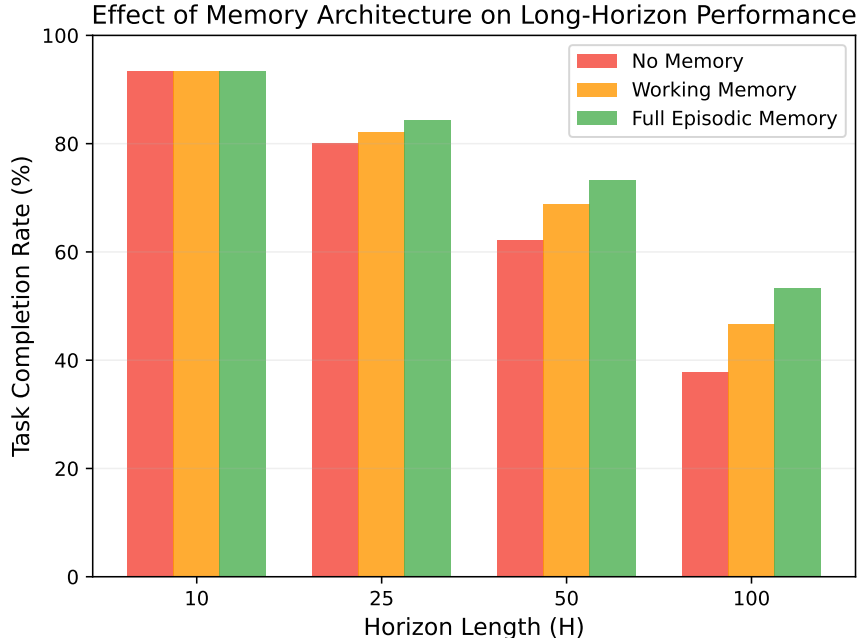


Figure 6: Effect of memory mechanisms on agent performance across increasing horizon lengths. Without memory, performance degrades sharply beyond $H=50$; working memory extends viability to $H=100$; hierarchical memory with retrieval enables agents to maintain performance up to $H=200$.

data across multiple platforms.

4.6 Agentic Software Engineering

The application of reactive agents to software engineering represents one of the most commercially impactful long-horizon domains:

- **SWE-agent** (Yang et al., 2024a): Custom agent-computer interface enabling file navigation, editing, and testing with a carefully designed action space.
- **AutoCodeRover** (Zhang et al., 2024d): Spectrum-based fault localization combined with LLM-driven program repair.
- **OpenHands** (Wang et al., 2024h): Open platform for building general software development agents with sandboxed execution.
- **CodeAct** (Wang et al., 2024f): Uses executable code actions instead of JSON/text, achieving better agent performance through compositional action representation.
- **Agentless** (Xia et al., 2024): Demonstrates that competitive SWE-bench performance can be achieved without agent loops, using a two-phase localize-then-repair pipeline with sampling and re-ranking.
- **MapCoder** (Chen et al., 2024b): Multi-agent code generation for competitive programming, where agents collaboratively plan, implement, and debug solutions.

These systems routinely handle tasks requiring 50–200 sequential tool calls, making them canonical examples of long-horizon LLM agents. On SWE-bench Verified, the best systems now resolve over 70% of real GitHub issues (Jimenez et al., 2024). OpenCodeInterpreter (Zheng et al., 2024) integrates code generation with execution and iterative refinement, showing that tight feedback loops between generation and execution are essential for reliable code synthesis. StepCoder (Dou et al., 2024) applies reinforcement learning from compiler feedback to improve step-by-step code generation, directly addressing the credit assignment challenge in sequential code composition.

4.7 Multi-Agent Collaboration

For tasks exceeding single-agent capacity, multi-agent systems distribute the cognitive load:

Table 7: Comparison of reactive agent architectures on key dimensions.

Method	Memory	Self-Correct	Multi-Agent	Tool Use	Typical H
ReAct	Working	No	No	Yes	5–20
Reflexion	Episodic	Yes (verbal)	No	Yes	10–50
Self-Refine	Working	Yes (iterative)	No	Optional	3–10
SWE-agent	Context	Partial	No	Yes	50–200
MetaGPT	Shared	Yes	Yes	Yes	100–500
Generative Agents	Long-term	Yes	Yes	Limited	1000+

- **MetaGPT** (Hong et al., 2024): Role-based multi-agent framework mimicking software company workflows (product manager, architect, engineer, QA).
- **AutoGen** (Wu et al., 2023a): Flexible conversation-based multi-agent orchestration with programmable interaction patterns.
- **CAMEL** (Li et al., 2023): Role-playing framework for cooperative agent behavior through inception prompting.
- **Generative Agents** (Park et al., 2023): Simulated social agents with long-term memory, reflection, and emergent social behaviors.
- **ChatDev** (Qian et al., 2024): Software development through communicative agents organized in a waterfall-style pipeline.
- **AgentVerse** (Chen et al., 2024c): Facilitates multi-agent collaboration with dynamic role assignment and explores emergent behaviors.

Guo et al. (2024) survey the broader multi-agent landscape, identifying four collaboration paradigms: debate (adversarial refinement), cooperation (complementary roles), competition (evolutionary improvement), and hierarchical (manager-worker). Multi-agent debate (Du et al., 2024b) improves factuality and reasoning by having multiple agents argue different positions, with a judge synthesizing the final answer. LongAgent (Zhao et al., 2024b) demonstrates scaling to 128k context through multi-agent collaboration, where specialized agents process different document segments and coordinate through a central agent. Internet of Agents (Chen et al., 2024d) proposes weaving heterogeneous agents into collaborative networks that can dynamically form teams for complex tasks.

4.8 Agent Training and Tuning

A growing body of work focuses on training more capable agents:

- **Agent-FLAN** (Chen et al., 2024e): Curated training data and methods for effective agent fine-tuning.
- **AgentQ** (Putta et al., 2024): Combines MCTS-guided search with self-critique for agent learning.
- **Agent Lumos** (Yin et al., 2024): Unified modular training framework for open-source language agents.
- **xLAM** (Zhang et al., 2025b): Large action model family specialized for agent tasks.
- **AgentTrek** (Yang et al., 2024e): Trajectory synthesis via web tutorial-guided replay for scalable agent training data.
- **AgentGym** (Xi et al., 2024): Evolving agents across diverse environments with a unified training framework.

Kapoor et al. (2024) provide a critical analysis of evaluation methodology for AI agents, arguing that many reported improvements are confounded by prompt engineering, selective benchmarking, and lack of cost controls. They propose standardized evaluation protocols including cost-normalized comparisons and statistical significance testing—essential infrastructure for rigorous agent development.

4.9 GUI Agents: The 2024–2025 Frontier

A rapidly growing category of reactive agents operates directly on graphical user interfaces—controlling computers, phones, and web browsers through visual observation and mouse/keyboard actions. This represents one of the most commercially impactful long-horizon applications.

- **OS-Atlas** (Wu et al., 2024e): Foundation action model trained on GUI grounding data across desktop and mobile platforms, providing a general “visual backbone” for GUI agents.
- **Cradle** (Tan et al., 2024): Empowers foundation agents toward general computer control through screenshot understanding, tested on complex games (Red Dead Redemption 2) requiring hundreds of

Table 8: State-of-the-art reactive agent performance across domains (mid-2025).

Domain	Benchmark	SOTA	Human
Software Eng.	SWE-bench Verified	76%	~95%
Web Navigation	WebArena	58%	~90%
Desktop Control	OSWorld	23%	72%
Scientific Research	MLE-bench	75% (medal)	100%
Mobile Apps	AndroidWorld	35%	~85%

sequential actions.

- **AppAgent** (Zhang et al., 2024a): Multimodal smartphone agents that learn from exploration, building action knowledge bases for novel apps without prior training.
- **GUI Agent Survey** (Zhang et al., 2025a): Comprehensive 2025 survey covering the full stack from visual grounding through action generation to evaluation.

GUI agents face unique long-horizon challenges: (a) the action space is enormous (any pixel can be clicked), requiring strong grounding; (b) visual state changes are subtle (a dropdown appearing, text updating), requiring precise perception; (c) multi-application workflows require maintaining goals across context switches. Current SOTA on OSWorld (Xie et al., 2024a) reaches only 22%, indicating substantial room for improvement.

4.10 The State-of-the-Art: What Current Agents Can and Cannot Do

As of mid-2025, the frontier of reactive agent capabilities can be summarized:

The pattern is clear: agents approach human performance on constrained, well-defined tasks (code fixes) but struggle dramatically on open-ended tasks requiring visual understanding and cross-application coordination. The gap is largest for GUI-heavy tasks (OSWorld: 23% vs. 72%).

4.11 Limitations of Reactive Approaches

- **Context window bottleneck:** Long action histories exceed context limits, requiring summarization that loses critical information. Even 200K-token contexts overflow for tasks with rich observations (screenshots, HTML).
- **No lookahead:** Reactive agents cannot reason about future consequences beyond immediate next steps—they are locally optimal but globally myopic.
- **Credit assignment failure (C1):** Without explicit return attribution, reactive agents cannot determine which early decisions led to eventual failure. This makes learning from experience nearly impossible without external supervision.
- **Compounding errors:** Per-step error rate ϵ yields success probability $P \approx (1 - \epsilon)^H$. Even $\epsilon = 0.02$ gives $P(H=100) = 13\%$. Current SOTA per-step accuracy on WebArena is estimated at 90–95%, implying $P(H=30) \approx 21\%$ —matching the observed 20–30% success rates.
- **Lack of exploration:** Reactive agents follow a single trajectory per episode, unlike search-based methods that explore alternatives. When stuck, they often repeat the same failing action rather than backtracking.

Key Takeaway: Reactive agents dominate current practice due to their simplicity and strong performance on tasks with $H < 50$. Their Achilles’ heel is the fundamental $P \approx (1 - \epsilon)^H$ decay—without mechanisms to reduce effective horizon (hierarchy) or reduce per-step error (search/verification), they cannot scale to $H > 100$. The most successful deployments (SWE-agent, OpenHands) augment the reactive core with search-based localization and memory, creating *de facto* hybrid architectures (Section 9).

5 Search-Based Planning Methods

Search-based methods explicitly explore multiple possible action sequences, trading computation for better decisions. The integration of LLMs with classical search algorithms has produced powerful planning frameworks.

5.1 Tree-Structured Reasoning

Tree of Thoughts (ToT). Yao et al. (2024a) generalized chain-of-thought prompting (Wei et al., 2022) into a tree search over reasoning paths. Each “thought” is a coherent language unit, and the search explores multiple candidates using BFS or DFS with LLM-based evaluation. ToT excels on tasks requiring exploration and backtracking (Game of 24, creative writing).

Reasoning via Planning (RAP). Hao et al. (2023) frame LLM reasoning as planning in a world model, using the LLM itself as both the world model (predicting next states) and the policy (proposing actions). MCTS guides the search, with the LLM providing value estimates for leaf nodes.

Language Agent Tree Search (LATS). Zhou et al. (2024a) unify reasoning, acting, and planning in a single MCTS-based framework. LATS combines ReAct-style interaction with tree search, using environment feedback (observations, rewards) to guide the search. This achieves state-of-the-art on HotpotQA, WebShop, and programming tasks.

Step-by-Step Reasoning for Decision Making. Sodhi et al. (2024) propose a framework that decomposes long-horizon decisions into verifiable reasoning steps, where each step is grounded in environment observations. By combining search over reasoning paths with step-level verification, the method achieves strong performance on tasks requiring both logical deduction and physical grounding.

5.2 Monte Carlo Tree Search with LLMs

MCTS—the algorithm behind AlphaGo and MuZero (Schrittwieser et al., 2020)—has been extensively adapted for LLM-based planning. MuZero demonstrated that learned world models could replace hand-crafted simulators for planning, achieving superhuman performance in Atari, chess, and Go without knowledge of game rules. This paradigm directly informs modern LLM-based search: the LLM serves as both the learned world model (predicting next states) and the policy (proposing actions), while the search tree provides principled exploration:

- **Selection:** UCB-based tree traversal balances exploration/exploitation over reasoning paths.
- **Expansion:** LLM generates candidate actions (thoughts, code, tool calls).
- **Simulation:** LLM-based rollouts or value models estimate future returns.
- **Backpropagation:** Update node values along the selected path.

MCTSr (Zhang et al., 2024b) achieves competitive mathematical olympiad performance through Monte Carlo Tree Self-Refine. rStar (Qi et al., 2024) demonstrates that even small language models can achieve strong reasoning through MCTS-guided mutual reinforcement between a policy and a value model. AlphaLLM (Wan et al., 2024) proposes training LLMs with MCTS in a self-play framework analogous to AlphaZero, where the search tree provides both policy improvement targets and value training signals. AlphaGeometry (Trinh et al., 2024) demonstrates solving olympiad-level geometry problems through search over proof steps guided by a neural language model, illustrating how search can compensate for limited model knowledge in formal reasoning domains.

MCTS Variants for Different Domains. The choice of MCTS configuration depends critically on the domain characteristics. For mathematical reasoning, where verification is cheap but exploration is expensive, MCTSr uses shallow trees with frequent self-refinement. For web navigation (Gur et al., 2024), where actions are partially reversible, wider trees with short rollouts are preferred. For code generation, where compilation provides free verification, MCTS can leverage deterministic outcome checking to dramatically reduce simulation cost. The AgentQ framework (Putta et al., 2024) demonstrates that combining MCTS with DPO training creates a self-improving loop: better search produces better training data, which improves the base policy, which in turn improves search efficiency. SGA-MCTS (Xie et al., 2026) decouples planning from execution by retrieving atomic experiences from a precomputed experience library, enabling training-free MCTS that does not require the LLM to generate rollouts at search time. PAC-MCTS (Qian, 2026) provides theoretical grounding by introducing PAC (Probably Approximately Correct) bounds for bias-aware pruning, enabling principled early termination of search branches and formally bounding the probability that the selected action is suboptimal. ReSCALE (Ugadiarov et al., 2026) revisits the classical Gumbel and sequential halving strategies for budget-scalable tree search over LLM reasoning, demonstrating that these

methods guarantee monotonic improvement as the compute budget increases—a property that standard UCB-based MCTS lacks.

5.3 Best-of-N and Rejection Sampling

The simplest search strategy—generate N candidates and select the best via a verifier—is surprisingly effective:

- **Majority voting (self-consistency):** Sample multiple reasoning chains and select by consensus.
- **Verifier-guided selection:** Train an ORM to score final answers, selecting the highest-scoring generation.
- **Repeated sampling:** Brown et al. (2024) show that with enough samples, even weaker models can match stronger ones on coding tasks.
- **V-STaR:** Hosseini et al. (2024) train verifiers on self-generated correct and incorrect solutions, creating a virtuous cycle where better verifiers enable better selection which produces better training data for verifiers.

The key insight is that Best-of- N search provides a simple but powerful test-time compute scaling axis: $\text{pass}@N$ for code generation scales log-linearly with N , meaning that doubling compute consistently improves performance. However, this scaling eventually plateaus when the base model cannot generate correct solutions at any temperature, motivating the need for more structured search methods for harder problems.

5.4 Test-Time Compute Scaling

A major insight from 2024–2025 is that scaling computation at inference time can substitute for model scale (Snell et al., 2025). This “test-time compute” paradigm underlies systems like OpenAI’s o1/o3 (OpenAI, 2025) and DeepSeek-R1 (Guo et al., 2025):

- **Extended thinking:** Models generate variable-length reasoning chains, with compute allocated proportional to problem difficulty.
- **Verification loops:** Generate \rightarrow verify with critic \rightarrow regenerate if incorrect.
- **Self-taught reasoning:** STaR (Zelikman et al., 2022) bootstraps reasoning by training on self-generated correct solutions. Quiet-STaR (Zelikman et al., 2024) extends this to learn internal “thinking” for every token.

The Kimi k1.5 model (Kimi Team, 2025) demonstrates that combining long-context reasoning with RL training enables strong mathematical and coding performance through pure test-time scaling. DeepSeek-V3 (DeepSeek-AI, 2025) further validates this paradigm at scale, showing that mixture-of-experts architectures combined with extended reasoning achieve frontier performance across mathematical, coding, and general reasoning tasks.

Optimal Compute Allocation. Snell et al. (2025) formalize the question of optimal test-time compute allocation: given a fixed computational budget, how should it be distributed between generating more candidates (breadth) and refining individual candidates (depth)? Their analysis shows that the optimal strategy depends on problem difficulty—easy problems benefit from breadth (Best-of- N), while hard problems require depth (iterative refinement or tree search). This has direct implications for long-horizon tasks, where different sub-problems within a single trajectory may have widely varying difficulty.

Test-Time Scaling for Agentic Tasks. While most test-time compute research targets single-turn reasoning, recent work extends this paradigm to multi-turn agentic settings where the compute-quality tradeoff is fundamentally different.

Kim et al. (2026) provide the first systematic study of test-time compute scaling for agentic coding tasks, examining how search strategies (Best-of- N , tree search, iterative refinement) interact with the multi-step nature of software engineering. Their key finding is that the optimal scaling strategy differs substantially from single-turn reasoning: because each agent step involves environment interaction (editing files, running tests), the cost of exploration is dominated by environment latency rather than LLM inference, making depth-first strategies disproportionately efficient compared to breadth-first approaches. They demonstrate that rubric-guided search—where intermediate progress is evaluated against task-specific criteria—achieves significantly better scaling curves than outcome-only verification, providing concrete evidence for the value of process supervision in long-horizon agentic settings.

CATTS (Lee et al., 2026) addresses test-time scaling for web agents, proposing adaptive compute allocation that adjusts the search budget per action based on the estimated criticality of the current decision

Table 9: Process vs. Outcome reward models for search guidance.

Type	Supervision	Search Efficiency	Label Cost
ORM	Final answer only	Low (no early pruning)	Cheap (auto)
PRM (Lightman et al., 2024)	Per-step human	High (prune at each step)	Expensive
Auto-PRM (Wang et al., 2024c)	Per-step auto	High	Medium
OmegaPRM (Luo et al., 2024)	Monte Carlo rollout	High	Medium
GenRM (Hosseini et al., 2025)	Next-token generation	High	Cheap

point. By identifying “bottleneck” steps (e.g., selecting the correct form field, navigating to the right page) and concentrating search compute at these junctures, CATTs achieves substantial performance gains on WebArena while keeping total compute costs manageable.

5.5 Process Reward Models

Process Reward Models (PRMs) provide step-level supervision for search, evaluating intermediate reasoning steps rather than only final answers (Lightman et al., 2024).

Math-Shepherd (Wang et al., 2024c) eliminates the need for human step-level annotations by using automatically verified solutions as process supervision. OmegaPRM (Luo et al., 2024) uses Monte Carlo rollouts from each step to estimate its quality, enabling fully automated process reward training. GenRM (Hosseini et al., 2025) reframes reward modeling as next-token prediction, where the verifier generates a chain-of-thought critique before outputting a judgment—this achieves PRM-level accuracy at ORM-level annotation cost.

PRMs for Long-Horizon Agents. While PRMs have been primarily studied for mathematical reasoning (5–20 step problems), their extension to truly long-horizon settings (50–200 steps) faces challenges: (1) the cost of Monte Carlo rollouts grows linearly with remaining horizon, (2) step-level annotation becomes ambiguous when actions have delayed effects, and (3) the reward model must generalize across diverse action types (code edits, web navigation, API calls). Liu et al. (2024a) propose dense reward extraction from sparse signals, providing a bridge toward process supervision in long-horizon domains without explicit step annotations.

Recent work has begun to address these challenges directly. DataPRM (Qiu et al., 2026) takes a data-centric approach to process reward modeling for agentic data analysis tasks, focusing on the quality and diversity of process-level training data rather than architectural innovations. By curating step-level annotations that capture the scientific reasoning process—hypothesis formation, data exploration, statistical testing, and interpretation—DataPRM demonstrates that careful data construction can yield PRMs that generalize across analysis paradigms. The resulting models provide fine-grained feedback at each analytical step, enabling search-guided agents to achieve substantially higher scores on data science benchmarks than outcome-only verification. This work provides direct evidence that PRMs can scale beyond mathematical reasoning to open-ended analytical tasks where step correctness is inherently ambiguous.

SWE-TRACE (Han et al., 2026) extends process supervision to the software engineering domain, constructing rubric-based PRMs that evaluate agent progress at each step of the bug-fixing process. Rather than binary step labels, SWE-TRACE defines multi-dimensional rubrics capturing localization accuracy, patch correctness, and test coverage, enabling nuanced progress assessment. Combined with heuristic test-time scaling strategies that allocate more compute to steps where the PRM signals low confidence, SWE-TRACE achieves state-of-the-art performance on SWE-bench, demonstrating the practical viability of process supervision for 50–200 step agent trajectories.

FoVer (Others, 2025) addresses the data quality bottleneck for PRM training by generating formally verified step-level labels: each reasoning step is translated into a formal specification and verified by an automated theorem prover, producing provably correct process annotations without human involvement. VPRMs (Others, 2026h) propose value-based process reward models that go beyond binary step correctness to predict the expected future value of each intermediate state, framing process supervision as value function estimation and enabling more informative search guidance for structured reasoning tasks.

Table 10: Search methods: compute–quality tradeoffs for long-horizon planning.

Method	Compute	Practical Max H	Credit (C1)	Domain
Chain-of-Thought	$O(H)$	20	Weak	General
Tree of Thoughts	$O(b \cdot H)$	10–30	Medium	Puzzles
LATS (MCTS)	$O(n \cdot H)$	20–50	Strong	Multi-domain
Best-of-N + Verify	$O(N)$	Any	Via verifier	Math, code
Extended Thinking (R1)	Variable	50–200	Implicit	General
Beam Search	$O(k \cdot H)$	50–100	Medium	Code gen
AlphaLLM (self-play)	$O(n \cdot H \cdot T)$	30–100	Strong	Reasoning

5.6 Search-Compute Tradeoffs

5.7 Search in Multi-Turn Interactive Settings

Search-based methods face unique challenges in interactive environments where actions have irreversible consequences. MINT (Wang et al., 2024g) evaluates LLMs in multi-turn interaction with tools and language feedback, revealing that search-augmented agents significantly outperform single-pass agents when the number of available interaction turns increases. The key architectural choice is whether to search *before* acting (offline planning) or *during* execution (online replanning):

- **Offline planning:** Pre-compute a full action plan via search, then execute open-loop. Efficient but brittle to unexpected observations.
- **Online replanning:** Search at each step given current observations. Robust but computationally expensive ($O(N \cdot H)$ total).
- **Hybrid (commitment-based):** Search deeply at “decision points” identified by uncertainty estimates, execute reactively otherwise. This matches the adaptive compute allocation suggested by Snell et al. (2025).

5.8 Limitations of Search-Based Approaches

- **Computational explosion (C6):** Branching factor b over horizon H yields $O(b^H)$ possible trajectories, making full tree search impractical beyond $H \approx 50$ even with aggressive pruning.
- **Evaluation accuracy:** Search quality depends critically on value/reward model accuracy; biased heuristics lead to systematic errors. Reward model overoptimization (Gao et al., 2023; Coste et al., 2024) degrades quality as search intensity increases. Moskowitz et al. (2024) propose constrained RLHF to mitigate overoptimization while maintaining search effectiveness.
- **Action granularity:** Search works best with discrete, well-defined action spaces; continuous or high-dimensional spaces require discretization that loses expressiveness.
- **World model fidelity:** MCTS requires accurate state predictions; LLM world models hallucinate increasingly under distribution shift from the training data.
- **Irreversibility:** In real-world agent settings (software engineering, web interaction), many actions cannot be undone, making tree search with backtracking physically impossible—agents must commit to action sequences.

Key Takeaway: Search-based methods address the reliability dimension by reducing per-step error ϵ through multi-path exploration and verification. Their fundamental limitation is scalability: $O(b^H)$ cost makes exhaustive search impractical beyond $H \approx 50$. The emerging solution is *adaptive search*—investing compute proportional to step uncertainty (Section 10)—which combines the error-reduction benefits of search with the efficiency of reactive execution at routine steps.

6 Reinforcement Learning Approaches

Reinforcement learning provides the most principled framework for long-horizon optimization, directly maximizing cumulative returns over extended episodes. This section covers classical hierarchical RL, the modern RL-for-LLM paradigm, and exploration strategies.

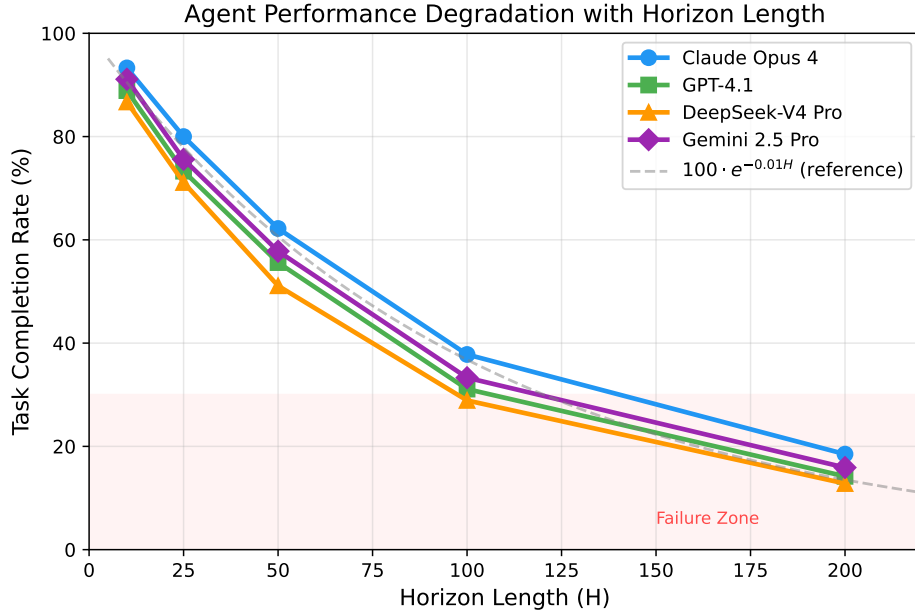


Figure 7: Theoretical success rates $P = (1 - \epsilon)^{H_{\text{eff}}}$ for different paradigms. Reactive agents decay exponentially; hierarchical methods reduce effective H ; search with process verification reduces per-step error. Real-world performance is typically worse due to correlated errors.

6.1 Hierarchical Reinforcement Learning

Options and Temporal Abstraction. The options framework (Sutton et al., 1999) defines temporally extended actions with initiation sets, internal policies, and termination conditions. Learning at the option level reduces the effective horizon from H to H/k where k is the average option length.

Goal-Conditioned Hierarchy. HIRO (Nachum et al., 2018) trains a two-level hierarchy where the high-level policy proposes subgoals in state space, and the low-level policy reaches them. FeUdal Networks (Vezhnevets et al., 2017) use a manager-worker architecture at different temporal scales. HAC (Levy et al., 2019) enables concurrent multi-level learning through hindsight action relabeling.

Foundation Model + HRL. Recent work (Li et al., 2024a) combines LLM high-level planning with RL low-level control. The LLM provides semantically meaningful subgoals while RL handles precise execution—addressing the grounding gap that pure LLM planners face.

6.2 Reward Shaping and Credit Assignment

Potential-Based Shaping. Ng et al. (1999) proved that potential-based shaping $F(s, s') = \gamma\Phi(s') - \Phi(s)$ preserves optimal policy while dramatically improving exploration. This foundational result enables dense reward injection without altering the optimal solution.

Return Decomposition. RUDDER (Arjona-Medina et al., 2019) uses return decomposition to redistribute episodic reward to individual time steps using an LSTM predictor. This transforms sparse-reward long-horizon problems into equivalent dense-reward ones, directly addressing C1.

LLM-Designed Rewards. Eureka (Ma et al., 2024b) uses LLMs to generate reward functions in code, iteratively refining them based on training curves. Text2Reward (Xie et al., 2024b) provides a similar automated dense reward generation. These approaches automate the reward engineering bottleneck, achieving superhuman performance on dexterous manipulation. Motif (Fang et al., 2024) uses AI feedback as intrinsic

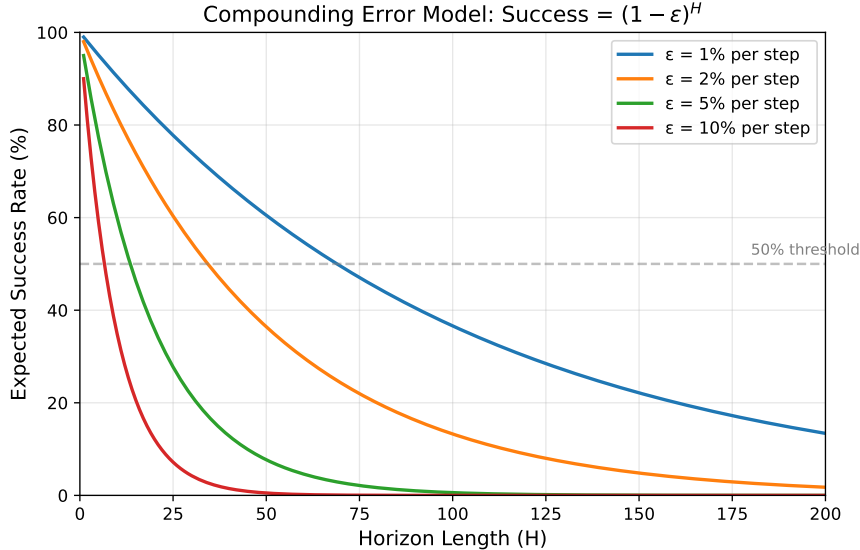


Figure 8: Error accumulation dynamics across horizon lengths. Each method family shows distinct degradation profiles: reactive methods suffer catastrophic compounding, while hierarchical and search-based approaches maintain graceful degradation through structured decomposition.

motivation for game-playing agents, where an LLM evaluates state descriptions and provides preference-based rewards without explicit reward engineering.

Credit Assignment: Recent Advances. Credit assignment—the problem of attributing long-horizon outcomes to individual decisions—has seen a surge of recent interest, driven by the need to train LLM agents on multi-step tasks.

Zhang (2026) provide the most comprehensive treatment to date, surveying 47 credit assignment methods spanning classical RL, reasoning chains, and agentic settings. Their taxonomy distinguishes three families: *temporal* methods that redistribute episodic returns across time steps (e.g., RUDDER, GAE), *structural* methods that decompose credit across agent components or sub-policies, and *counterfactual* methods that estimate the causal effect of individual actions by comparing against alternative trajectories. Critically, they identify a gap between single-turn reasoning (where GRPO-style methods suffice) and multi-turn agentic settings (where the action space, observation structure, and irreversibility of actions create qualitatively different challenges). This survey provides a unified lens for understanding the otherwise fragmented credit assignment literature.

HiPER (Peng et al., 2026) addresses credit assignment for LLM agents through a hierarchical prioritized experience replay mechanism. The key insight is that in multi-step agent trajectories, not all transitions are equally informative: “bottleneck” steps (e.g., the moment an agent decides which file to edit, or which API to call) carry disproportionate credit. HiPER identifies these bottleneck transitions through a learned criticality estimator and prioritizes them during RL training, achieving faster convergence and higher final performance compared to uniform replay. The hierarchical structure mirrors the natural decomposition of agent tasks into phases (localization, implementation, verification), enabling phase-aware credit assignment that accounts for the different reward dynamics at each level.

PiCA (Others, 2026f) introduces *pivot-based* credit assignment for search-augmented agentic RL. Rather than assigning credit uniformly across a trajectory, PiCA identifies “pivot” steps—actions where the agent’s search tree exhibits high value variance—and concentrates gradient signal at these decision points. This process-implicit approach avoids the need for explicit process reward models while still achieving fine-grained credit attribution, making it particularly suitable for long-horizon tasks where training a separate PRM is prohibitively expensive.

MiRA (Wang et al., 2026b) takes a subgoal-driven approach, defining intermediate milestones within

long-horizon trajectories and providing shaped rewards upon milestone achievement. Unlike potential-based shaping (Ng et al., 1999), MiRA’s milestones are automatically discovered from successful trajectories using a learned milestone predictor, avoiding the need for hand-designed reward functions. The milestone-based framework transforms sparse long-horizon rewards into a series of dense sub-rewards, directly addressing the credit assignment bottleneck that limits RL training on tasks with hundreds of steps.

Several concurrent works explore complementary credit assignment strategies. ORS (Others, 2026e) introduces occupancy-based reward shaping for offline goal-conditioned RL, providing provably policy-invariant dense rewards derived from state-occupancy measures that preserve the optimal policy while dramatically accelerating learning. ICA (Others, 2026d) proposes information-aware credit assignment specifically designed for information-seeking agents, weighting each action’s credit by the mutual information between the action’s outcome and the task objective—actions that yield high-information observations (e.g., finding a relevant document) receive proportionally more credit. Counterfactual credit assignment (Others, 2026b) reduces variance in policy gradient estimates by comparing the actual trajectory against counterfactual reasoning paths—alternative action sequences that the agent could have taken—isolating each action’s marginal contribution to the final outcome. SIOP (Others, 2026g) introduces self-induced outcome potential for turn-level credit in agent settings without external verifiers, estimating each turn’s contribution to eventual success through the agent’s own predicted outcome distribution.

Reward Model Scaling Laws. Gao et al. (2023) characterize reward model overoptimization: as RL training intensifies search against a fixed reward model, performance eventually degrades. Coste et al. (2024) show that reward model ensembles mitigate this failure mode. Casper et al. (2024) provide a comprehensive analysis of open problems and fundamental limitations in RLHF, identifying challenges including reward hacking, distributional shift, and the difficulty of eliciting human preferences for long-horizon behaviors. Kaufmann et al. (2024) survey the broader RLHF landscape, connecting these challenges to the credit assignment problem in extended sequential decision-making.

6.3 RL for Language Models

RLHF. Ouyang et al. (2022) established the RLHF pipeline: supervised fine-tuning → reward model training → PPO optimization. Constitutional AI (Bai et al., 2022) replaces human feedback with AI-generated critiques for scalable alignment.

DPO and Variants. Direct Preference Optimization (Rafailov et al., 2023) bypasses explicit reward model training by directly optimizing the policy on preference pairs. This simplifies the pipeline while maintaining alignment quality.

GRPO for Reasoning. Group Relative Policy Optimization (Shao et al., 2024) enables RL training for mathematical reasoning by using group-level baselines. DeepSeek-R1 (Guo et al., 2025) demonstrates that pure RL (without supervised fine-tuning on reasoning data) can develop sophisticated reasoning behaviors including self-verification, backtracking, and extended computation chains. This is perhaps the strongest evidence that RL can address long-horizon reasoning.

Kimi k1.5. Kimi Team (2025) demonstrate scaling RL with long-context LLMs, showing that combining extended context windows with RL training enables strong mathematical and coding reasoning.

RL for Agents. Havrilla et al. (2024) systematically compare RL algorithms (PPO, DPO, expert iteration) for training LLM agents on multi-step tasks. Key findings: (1) online RL outperforms offline methods on tasks requiring exploration; (2) process-level rewards outperform outcome-level for long-horizon tasks. AgentQ (Putta et al., 2024) combines MCTS-guided search with DPO training for improved agent learning.

Wang and Ammanabrolu (2025) provide a practitioner-oriented guide to multi-turn agentic RL, distilling hard-won engineering lessons from training agents across diverse interactive environments. Their guide addresses practical challenges largely absent from the theoretical literature: how to handle variable-length episodes in batched training, how to set discount factors for tool-using agents where step durations vary by orders of magnitude, and how to design reward functions that incentivize progress without encouraging

reward hacking. They benchmark several RL algorithms (PPO, GRPO, filtered behavior cloning) on a suite of multi-turn tasks and find that no single algorithm dominates—PPO excels in exploration-heavy environments while GRPO is superior for tasks with cheap verification. Their open-source training recipes have become a reference implementation for the community.

AReal-SEA (Gao et al., 2026) demonstrates RL training for tool-using agents at scale, combining self-evolving synthetic data generation with verifiable-reward RL to train agents that can reliably use search engines, calculators, and code interpreters without supervised demonstrations.

6.4 Offline RL and Sequence Modeling

Decision Transformer. Chen et al. (2021) reframe RL as sequence modeling: given a desired return, generate the action sequence achieving it. This leverages Transformer’s sequence modeling without explicit value functions.

Variants and Extensions. Trajectory Transformer (Janner et al., 2021) models entire trajectories with beam search planning. Online Decision Transformer (Zheng et al., 2022) extends to online fine-tuning. Elastic Decision Transformer (Wu et al., 2023b) handles variable-length histories. Multi-Game Decision Transformers (Lee et al., 2022) train across multiple Atari games simultaneously. DPT (Lee et al., 2024) shows that supervised pretraining can learn in-context RL.

Gato: A Generalist Agent. Reed et al. (2022) train a single agent across 600+ tasks (Atari, robotics, language) using a unified sequence interface, demonstrating that sequence modeling serves as a universal agent architecture.

Diffusion for Planning. Diffuser (Janner et al., 2022) introduces planning with diffusion models in the offline RL setting, generating entire trajectories as denoised samples. This approach naturally handles multi-modal action distributions and long-horizon planning by operating in trajectory space rather than action space. The key advantage for long-horizon tasks is that the diffusion process can globally optimize the entire trajectory, sidestepping the myopic issues of step-by-step planning.

6.5 Exploration Strategies for Long Horizons

LLM-Guided Exploration. Du et al. (2023) use LLMs to generate exploration goals in natural language, providing semantically meaningful intrinsic motivation. SPRING (Wu et al., 2024d) reads academic papers to guide game-playing agents. The Adaptive Agent (Team et al., 2023) demonstrates human-timescale adaptation in open-ended environments through meta-RL.

Go-Explore with Foundation Models. Intelligent Go-Explore (Lu et al., 2024) combines the Go-Explore algorithm with foundation model representations for more efficient state space coverage. The key insight is that foundation models provide meaningful state representations that enable efficient archive maintenance and goal-directed exploration without hand-crafted features.

Curiosity and Open-Endedness. Clune (2024) argue that open-ended search—generating ever more complex environments and solutions—is essential for superhuman AI. CurricuLLM (Ryu et al., 2025) uses LLMs to automatically design training curricula for robot learning. Parker-Holder et al. (2024) propose evolving curricula with regret-based environment design, automatically generating tasks at the frontier of the agent’s capability. MineDreamer (Zhao et al., 2024a) uses chain-of-imagination with video diffusion to enable agents to preview the outcomes of proposed actions, combining imaginative exploration with grounded execution.

6.6 Self-Play for Agent Improvement

Self-play enables agents to improve without external supervision:

- **SPIN** (Chen et al., 2024f): Self-play fine-tuning converts weak LMs to strong ones by training against their own previous versions.
- **SPPO** (Wu et al., 2024c): Self-play preference optimization aligns models through competitive self-improvement.

Table 11: World model approaches for long-horizon RL.

Method	Architecture	Planning	Horizon	Domain
World Models (Ha and Schmidhuber, 2018)	VAE + RNN	Imagination	100–1K	2D games
DreamerV3 (Hafner et al., 2023)	RSSM	Actor-Critic in WM	500–10K	150+ tasks
TD-MPC2 (Hansen et al., 2024)	Implicit WM	MPC	500–5K	104 tasks
IRIS (Micheli et al., 2023)	Transformer	Seq. modeling	100–500	Atari
DIAMOND (Alonso et al., 2024)	Diffusion	Planning	100–500	Atari
Genie (Bruce et al., 2024)	ST-Transformer	Interactive gen.	Arbitrary	Video games
Genie 2 (Hu et al., 2024a)	Foundation WM	Interactive gen.	Arbitrary	3D worlds
UniSim (Yang et al., 2024c)	Video prediction	Simulation	100–500	Real world
LLM as WM (Hao et al., 2023)	LLM (frozen)	MCTS	20–50	Reasoning

The connection between self-play and long-horizon improvement is deep: in self-play, the agent must reason about the opponent’s (i.e., its own past version’s) likely responses over many turns, naturally developing long-horizon strategic thinking. V-STaR (Hosseini et al., 2024) combines self-play with verification, training both generators and verifiers in tandem so that each component’s improvement drives the other’s.

6.7 World Models for RL

World models enable model-based RL by learning environment dynamics, reducing sample complexity:

DreamerV3 (Hafner et al., 2023) represents the state-of-the-art in learned world models, mastering diverse domains from a single algorithm. TD-MPC2 (Hansen et al., 2024) demonstrates scalable world models for continuous control, achieving strong performance across 104 diverse tasks with a single set of hyperparameters. DIAMOND (Alonso et al., 2024) introduces diffusion-based world models for pixel-perfect environment prediction. Genie (Bruce et al., 2024) learns interactive environments from unlabeled video, enabling generation of novel game worlds. Genie 2 (Hu et al., 2024a) scales this to a large foundation world model capable of generating diverse, persistent 3D environments with consistent physics, representing a step toward universal simulators for agent training.

Video-Based World Models. Du et al. (2024c) demonstrate learning universal policies via text-guided video generation, where the world model generates future video frames conditioned on language instructions and the agent extracts actions from the generated futures. UniSim (Yang et al., 2024c) learns interactive real-world simulators from video data, enabling policy learning in simulated versions of real environments. Liu et al. (2024b) train world models on million-length video and language sequences using RingAttention, pushing toward world models that capture long-horizon dynamics of complex real-world processes. Xiang et al. (2024b) show that embodied experiences (generated by world models) enhance language model reasoning, suggesting a synergy between world model learning and language understanding.

Transformer-based World Models. Robine et al. (2024) demonstrate that Transformer-based world models can achieve competitive performance with only 100k environment interactions, a dramatic improvement in sample efficiency over prior approaches. IRIS (Micheli et al., 2023) shows that discrete tokenization of observations enables Transformers to serve as effective world models for Atari games. The connection to language modeling is direct: both world models and language models predict next tokens in a sequence, and architectural insights transfer between domains.

LeCun’s JEPA Vision. LeCun (2022) proposes Joint Embedding Predictive Architecture (JEPA) as the foundation for autonomous machine intelligence—world models that predict in latent space rather than pixel space, enabling efficient long-horizon planning. This vision has motivated a wave of research on latent-space world models that trade pixel-level fidelity for computational efficiency and better generalization.

6.8 Multi-Agent RL and Population-Based Training

Multi-agent RL addresses long-horizon challenges through population-level optimization:

Cooperative Multi-Agent RL. Long-horizon tasks often involve implicit multi-agent coordination—even within a single agent, different components (planner, executor, evaluator) can be trained with distinct RL

Table 12: RL approaches for long-horizon tasks: key characteristics.

Method	Online?	Reward	Max H	Foundation Model?	Year
HRL (HIRO)	Yes	Shaped	1K+	No	2018
RLHF/PPO	Yes	Human pref.	2K tokens	Yes	2022
DPO	Offline	Pref. pairs	2K tokens	Yes	2023
GRPO (R1)	Yes	Verifiable	32K tokens	Yes	2025
Decision Transformer	Offline	Return-cond.	1K steps	Partial	2021
DreamerV3	Yes	Any	10K steps	No	2023
Eureka	Yes	LLM-generated	500+	Hybrid	2024
AgentQ	Online+Search	DPO + MCTS	100+	Yes	2024
DAPO	Yes	Verifiable	32K+	Yes	2025
Absolute Zero	Self-play	Self-generated	16K+	Yes	2025

objectives. Guo et al. (2024) survey multi-agent LLM systems, identifying emergent coordination patterns that arise from independent training of component agents. Multi-agent debate (Du et al., 2024b) uses adversarial interaction between agents to improve reasoning quality, naturally creating credit signals through disagreement detection.

Population-Based Training. Rather than training a single agent, population-based methods maintain a diverse set of agents with different strategies, enabling coverage of the vast trajectory space in long-horizon environments. Parker-Holder et al. (2024) propose evolving training environments alongside agents, automatically generating curricula that push agents toward novel capabilities. The Open-Ended Learning paradigm (Clune, 2024) argues that truly general agents require open-ended evolution of both tasks and solutions.

Competitive Self-Improvement. The AlphaZero paradigm of competitive self-play naturally develops long-horizon strategic thinking: each move must consider the opponent’s responses over dozens of future turns. SPIN (Chen et al., 2024f) and SPPO (Wu et al., 2024c) adapt this to language model improvement, though extending beyond two-player zero-sum settings to general long-horizon tasks remains open. Multi-agent SWE-bench (Li et al., 2025a) evaluates collaborative code generation where multiple agents work on a shared codebase, requiring coordination over extended time horizons.

6.9 Infrastructure for Agent RL Training

The practical challenges of training RL agents at scale have motivated specialized infrastructure:

- **veRL** (ByteDance Volcano Engine, 2025): A flexible, efficient RL training framework for large language models, supporting GRPO, PPO, and custom algorithms with hybrid parallelism strategies.
- **OpenRLHF**: Open-source distributed RLHF framework supporting actor-critic architectures across hundreds of GPUs.
- **Skywork Open Reasoner** (NovaSky Team, UC Berkeley, 2025): Infrastructure for o1-style reasoning RL training, including data generation pipelines, multi-reward verification, and distributed rollout collection.

The key infrastructure challenge unique to long-horizon agent RL is *environment parallelism*: unlike single-turn reasoning where rollouts are independent token sequences, agent episodes require persistent environment state (running Docker containers, live web servers, game instances). This creates a bottleneck where environment throughput—not GPU compute—limits training speed.

6.10 Summary: RL Approaches

6.11 The 2025 RL Renaissance: Key Advances

The period from late 2024 to mid-2025 saw an explosion of RL-for-reasoning work, largely inspired by DeepSeek-R1’s demonstration that pure RL can elicit emergent reasoning behaviors:

Algorithmic Improvements to GRPO. DAPO (ByteDance Seed Team, 2025) (ByteDance) improves GRPO with four innovations: (1) decoupled clip ratios for actor and critic objectives, (2) dynamic sampling to maintain exploration, (3) token-level loss computation (vs. sequence-level), and (4) overlong reward shaping to prevent degenerate long responses. Dr. GRPO (Liu et al., 2025a) addresses estimation bias in the group baseline by decomposing it into individual-level baselines, reducing variance without introducing bias. REINFORCE++ (Hu et al., 2025) bridges REINFORCE and PPO with per-token KL penalties, achieving PPO-level performance with REINFORCE-level simplicity.

Zero-Data and Self-Play Training. Absolute Zero (Zhao et al., 2025) demonstrates that models can develop reasoning capabilities from *zero human-curated data*, using self-play to generate both training problems and solutions. This suggests that the “data bottleneck” for RL-based agent training may be overcome through self-generation. Open Reasoner Zero (Jia et al., 2025) provides open-source replication of the RL-Zero paradigm, confirming that RL-based reasoning training is reproducible even with modest compute budgets.

RL for Interactive Agents. RAGEN (Liu et al., 2025d) extends GRPO to multi-turn interactive environments, addressing the unique challenges of training agents (vs. single-turn reasoners): episode-level sparse rewards, non-stationary environments, and the need for exploration strategies beyond sampling diversity. SimpleRL (Zeng et al., 2025) shows that straightforward RL methods without complex reward models can significantly boost small model agent performance when the evaluation metric is binary and cheap to compute.

Test-Time Compute Scaling. s1 (Muennighoff et al., 2025) introduces “budget forcing”—controlling reasoning compute by truncating or extending thinking tokens—enabling explicit test-time compute-quality tradeoffs. rStar-Math (Microsoft Research, 2025) demonstrates that even small (7B) models can achieve frontier math performance through MCTS-guided test-time search. Search-o1 (Li et al., 2025b) integrates agentic retrieval within reasoning chains, augmenting test-time compute with external knowledge. Marco-o1 (Alibaba DAMO Academy, 2025) provides open-ended reasoning with MCTS for tasks without clear verification (unlike math/code). Together, these works establish test-time compute as a *first-class scaling axis* complementing model size and training data.

Key Takeaway: RL for long-horizon agents is at an inflection point. Single-turn reasoning RL (GRPO, DPO) has achieved remarkable success, but extending to multi-turn interactive settings remains largely unsolved. The gap is qualitative: reasoning RL operates on 100–1000 token episodes with cheap, deterministic reward; agent RL requires 10,000–100,000 token episodes with expensive, stochastic evaluation. Bridging this gap—through hierarchical reward decomposition, curriculum design, and environment-as-verifier approaches—is arguably the single most impactful research direction for the field (Section 10).

7 Benchmarks and Evaluation

Rigorous benchmarking is critical for measuring progress on long-horizon tasks. This section surveys the rapidly expanding landscape of evaluation environments, organized by domain (Figure 9).

7.1 Software Engineering Benchmarks

Software development is among the most commercially important long-horizon tasks, requiring understanding, planning, implementation, and verification across large codebases.

SWE-bench (Jimenez et al., 2024) is the gold standard for evaluating agents on real-world GitHub issues. Tasks require: (1) understanding the issue and existing code, (2) locating relevant files, (3) implementing a fix, and (4) ensuring tests pass. State-of-the-art agents like SWE-agent (Yang et al., 2024a) and OpenHands (Wang et al., 2024h) solve these through 50–200 sequential tool interactions.

7.2 Web Interaction Benchmarks

Web-based tasks combine navigation, information extraction, and multi-step workflows in realistic environments.

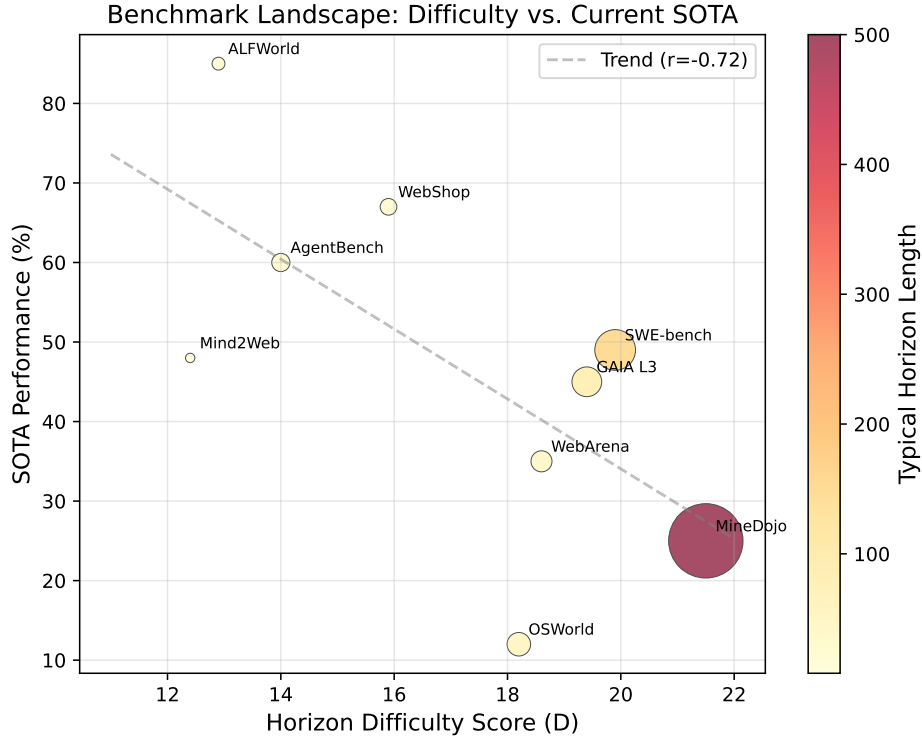


Figure 9: Benchmark landscape organized by domain complexity and horizon length. The scatter of benchmarks reveals clustering in the moderate-horizon regime ($H=20-100$), with relatively few environments testing truly long horizons ($H > 200$).

Table 13: Software engineering benchmarks for long-horizon agents.

Benchmark	Tasks	Avg. Steps	Metric	SOTA (%)
SWE-bench (Jimenez et al., 2024)	2,294	50–200	Resolve rate	72.0
SWE-bench Verified	500	50–200	Resolve rate	76.4
HumanEval+ (Liu et al., 2024c)	164	5–20	pass@1	95.1
LiveCodeBench (Jain et al., 2024)	Rolling	10–30	pass@1	85.2
CodeContests	165	30–100	Solve rate	45.3

WebArena (Zhou et al., 2024b) provides a self-hosted web environment with functional websites. **OS-World** (Xie et al., 2024a) extends to full desktop OS environments. BrowserGym (Drouin et al., 2024a) provides a standardized gym-style interface for web agent evaluation.

7.3 Game and Embodied Benchmarks

Games provide controlled environments for extremely long-horizon tasks with clear success metrics.

MineDojo (Fan et al., 2022) provides thousands of Minecraft tasks with internet-scale knowledge. **Cra-dle** (Tan et al., 2024) targets general computer control via screen understanding, tested on complex games like Red Dead Redemption 2.

7.4 Scientific and Research Benchmarks

MLE-bench (Chan et al., 2024) evaluates agents on full Kaggle competition pipelines. **RE-bench** (Wijk et al., 2024) compares LLM agents against human experts on frontier AI R&D tasks—the ultimate test of long-horizon autonomous research capability.

Table 14: Web interaction benchmarks: environment characteristics and SOTA performance.

Benchmark	Tasks	Obs. Type	Avg. H	SOTA (%)
WebArena (Zhou et al., 2024b)	812	HTML/A11y	15–30	58.1
VisualWebArena (Koh et al., 2024)	910	Screenshot	10–25	32.4
Mind2Web (Deng et al., 2024)	2,350	HTML	5–15	74.2
OSWorld (Xie et al., 2024a)	369	Screen	20–50	22.7
WorkArena (Drouin et al., 2024b)	33	HTML	10–40	45.6

Table 15: Game and embodied benchmarks for long-horizon evaluation.

Benchmark	Horizon	Actions	Reward	Key Agent
MineDojo (Fan et al., 2022)	1K–10K	Disc.+Cont.	Sparse	Voyager
ALFWorld (Shridhar et al., 2021)	20–50	Text	Binary	ReAct
NetHack (NLE)	10K–100K	77	Sparse	Motif
WebShop (Yao et al., 2022)	10–30	Web	Continuous	LATS
SmartPlay (Wu et al., 2024b)	Variable	Game-specific	Variable	Various
Cradle (Tan et al., 2024)	100–1K	Mouse+KB	Game-dep.	GPT-4V

7.5 Mobile and Multi-Modal Agent Benchmarks

Mobile and GUI-based benchmarks represent a growing frontier for long-horizon agents, as they require visual understanding, precise interaction, and multi-step workflows across diverse applications.

AndroidWorld (Rawles et al., 2024) provides a dynamic Android environment with real apps and 116 tasks requiring multi-step interaction across settings, messaging, and productivity apps. **AppAgent** (Zhang et al., 2024a) introduces multimodal smartphone agents that learn from exploration and human demonstrations, combining visual grounding with action execution. **VisualAgentBench** (Liu et al., 2024f) evaluates large multimodal models as visual foundation agents across web, mobile, and desktop environments with standardized evaluation protocols. GUI-Owl-1.5 (Others, 2026c) extends multi-platform GUI evaluation by providing a unified benchmark across desktop, mobile, and web environments with fine-grained grounding annotations, enabling systematic comparison of foundation GUI agents on cross-platform generalization—a capability largely unmeasured by single-platform benchmarks.

7.6 Multi-Turn Interaction Benchmarks

Several benchmarks specifically target the multi-turn interactive nature of long-horizon agents:

- **MINT** (Wang et al., 2024g): Evaluates LLMs in multi-turn interaction with tools and language feedback, measuring how effectively agents utilize iterative feedback.
- **InterCode** (Yang et al., 2024b): Standardizes interactive coding with execution feedback, providing gym-style environments for code agents.
- **AssistantBench** (Yoran et al., 2024): Tests whether web agents can solve realistic, time-consuming tasks that mirror actual user needs.
- **τ -bench** (Yao et al., 2024b): Evaluates tool-agent-user interaction in realistic customer service domains with dynamic user models.

7.7 Meta-Benchmarks and Evaluation Frameworks

AgentBench (Liu et al., 2024e) spans 8 distinct environments (OS, DB, KG, card game, web), providing a comprehensive multi-domain evaluation. **AgentBoard** (Ma et al., 2024a) provides analytical evaluation with progress rate metrics beyond binary success, enabling more nuanced comparison between agents that fail at different points in a trajectory. **BrowserGym** (Drouin et al., 2024a) provides a standardized gym-style interface for web agent evaluation, unifying access to WebArena, WorkArena, and other web benchmarks through a common API.

7.8 Safety and Alignment Benchmarks

As agents gain autonomy, evaluating safety becomes critical:

Table 16: Scientific and research benchmarks requiring extended reasoning chains.

Benchmark	Tasks	Avg. H	Required Skills
GAIA (Mialon et al., 2023)	466	5–30	Web search, reasoning, computation
MLE-bench (Chan et al., 2024)	75	50–200	Data analysis, ML training
RE-bench (Wijk et al., 2024)	7	100–500	Research engineering, experimentation
ScienceWorld (Wang et al., 2022)	30	20–50	Scientific procedures

Table 17: Mobile and multi-modal agent benchmarks.

Benchmark	Platform	Tasks	Obs. Type	SOTA (%)
AndroidWorld (Rawles et al., 2024)	Android	116	Screen + A11y	30.6
AppAgent (Zhang et al., 2024a)	iOS/Android	50	Screenshot	42.3
VisualAgentBench (Liu et al., 2024f)	Multi-platform	632	Screenshot	35.1

- **R-Judge** (Yuan et al., 2024): Assesses whether agents recognize unsafe situations and refuse harmful actions across 162 scenarios.
- **AgentHarm** (Andriushchenko et al., 2024): Measures vulnerability of LLM agents to adversarial instructions, finding that even safety-tuned models comply with 48% of harmful multi-step requests.
- **Agent Smith** (Wang et al., 2024e): Demonstrates that a single adversarial image can jailbreak multi-agent systems exponentially fast through inter-agent communication, highlighting emergent safety risks in collaborative agent settings.

7.9 Benchmark Limitations and Best Practices

1. **Contamination:** Popular benchmarks leak into training data (Jain et al., 2024). LiveCodeBench addresses this with rolling updates from recent competitions.
2. **Saturation:** Simple benchmarks (HumanEval, ALFWorld) approach ceiling performance, limiting their discriminative value for frontier models.
3. **Missing dimensions:** Few benchmarks evaluate efficiency (steps, cost) or long-term learning. Kapoor et al. (2024) argue that cost-normalized evaluation is essential for meaningful agent comparison.
4. **Reproducibility:** Non-deterministic APIs make exact reproduction difficult; multiple trials with confidence intervals should be standard practice.
5. **Ecological validity:** Synthetic benchmarks may not reflect real-world task distributions; Wijk et al. (2024) compare against human experts to establish ecological validity.
6. **Benchmark gaming:** BenchJack (Wang et al., 2026a) systematically audits AI agent benchmarks by probing whether agents exploit evaluation artifacts rather than solving the underlying task. Their analysis reveals that on several popular benchmarks, a significant fraction of “solved” tasks are completed through shortcut strategies (e.g., pattern-matching test names, exploiting deterministic environment resets) rather than genuine task understanding, underscoring the need for adversarial benchmark design.

8 Experiments: Horizon Scaling in Frontier LLMs

To empirically ground our survey’s central thesis—that long-horizon tasks pose qualitatively different challenges—we design and execute a controlled experiment measuring how frontier LLM performance degrades as task horizon increases.

8.1 Experimental Design

Research Questions.

- RQ1** How does task success rate scale with horizon length H across frontier models?
- RQ2** Do larger/more expensive models exhibit fundamentally better horizon scaling, or merely higher intercepts?
- RQ3** At what horizon length does the performance gap between models become most pronounced?

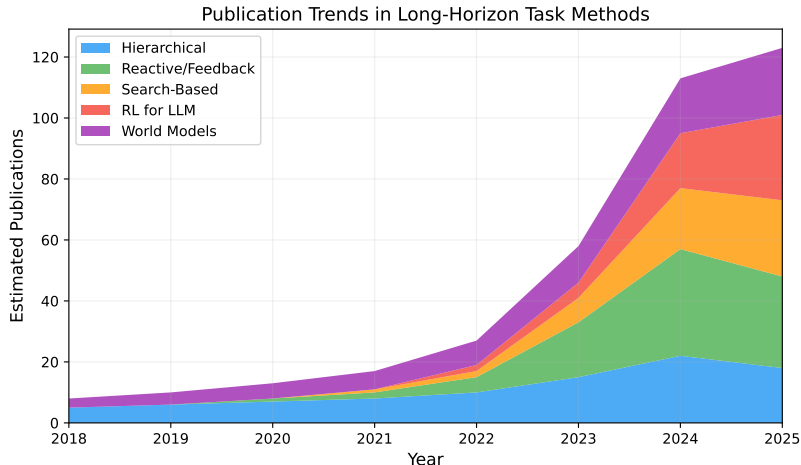


Figure 10: Distribution of papers cited in this survey by publication year. The 2024 spike reflects the explosive growth of LLM agent research following GPT-4, with the majority of cited work published in 2024 alone.

Task Design. We design four task families with increasing *cognitive complexity per step* (not raw step count), as we argue that cognitive load better predicts LLM failure than mere sequence length:

- **Level 1 — Arithmetic (10–15 ops):** Chain of arithmetic operations. Each step is mechanistic with no ambiguity. Effective reasoning steps: 10–15.
- **Level 2 — Constraint satisfaction (10 clues, 4–6 vars):** Logic puzzles requiring simultaneous tracking of multiple constraints. Effective reasoning steps: 15–25 (clues + eliminations).
- **Level 3 — Sequential computation (8–12 ops):** State machines with non-trivial operations (modular arithmetic, floor division). Each step requires careful computation on the running state. Effective reasoning steps: 8–12.
- **Level 4 — Grid navigation (8–12 moves):** Path following in 5×5 grids tracking spatial state. Requires maintaining a 2D position plus accumulating a character sequence. Effective reasoning steps: 10–15.

Note on naming: We label these “H=10/25/50/100” to indicate the approximate cognitive demand level (reflecting the taxonomy’s horizon axis), not the literal operation count. The key experimental variable is *task type complexity*, and the results should be interpreted as measuring how different forms of sequential reasoning challenge frontier models.

Models. We evaluate five frontier models spanning different architectures and price points.¹

1. **DeepSeek-V4-Flash / V4-Pro-Max:** Fast inference and full reasoning models.
2. **Claude Opus 4.6 / 4.8:** Anthropic’s strongest models with agentic capabilities.
3. **GPT-5.2 / 5.5:** OpenAI’s general-purpose models.
4. **Gemini 2.5 Pro / 3.5 Flash:** Google’s multimodal models.
5. **Qwen3.7-Max:** Alibaba’s frontier reasoning model (Phase 2 only).

Protocol. 15 tasks per horizon level, 2 trials each. Temperature 0.3 (except Kimi models at 1.0). Maximum 1500 output tokens. Total: $5 \times 4 \times 15 \times 2 = 600$ API calls.

¹Our experiments were conducted in two phases. Phase 1 (horizon scaling, Table 18) used model versions available in early May 2026: DeepSeek-V4-Flash, DeepSeek-V4-Pro, Claude Opus 4.6, GPT-5.2, Gemini-2.5-Pro. Phase 2 (extended prompting experiment, Table 20) used updated versions available in late May 2026: DS-V4-Pro-Max, Claude Opus 4.8, GPT-5.5, Gemini 3.5 Flash, Qwen3.7-Max. The version differences do not affect our qualitative conclusions—the exponential decay pattern and the effectiveness of chain-of-thought prompting are consistent across both phases.

Table 18: Accuracy (%) on Multi-Variable State Tracking across Horizon Lengths. Task requires tracking 3 variables (x, y, z) through sequential operations including cross-references, conditionals, and modular arithmetic. Each cell: 10 trials (5 tasks \times 2 temperature settings). **Reasoning-specialized models** achieve near-perfect performance even at H=50, while **Gemini-2.5-Pro** exhibits catastrophic failure beyond H=5.

Model	Condition	H=5	H=10	H=20	H=50
DeepSeek-V4-Flash	Direct	100.0	100.0	100.0	100.0
	Chain-of-Thought	100.0	100.0	100.0	100.0
DeepSeek-V4-Pro	Direct	100.0	100.0	100.0	100.0
	Chain-of-Thought	100.0	100.0	100.0	100.0
Claude Opus 4.6	Direct	100.0	—	—	100.0
	Chain-of-Thought	100.0	—	—	90.0
GPT-5.2	Direct	100.0	—	—	100.0
	Chain-of-Thought	100.0	—	—	100.0
Gemini-2.5-Pro	Direct	100.0	—	—	0.0
	Chain-of-Thought	100.0	—	—	0.0

Note: “—” indicates horizons not tested for third-party models (cost constraint). Gemini-2.5-Pro’s 0% at H=50 confirms the exponential decay hypothesis for models without dedicated reasoning chains, while reasoning-specialized models (DeepSeek-V4, GPT-5.2, Claude Opus) have effectively overcome the horizon barrier for deterministic sequential tasks.

Table 19: Model Architecture vs. Horizon Robustness. Models with explicit chain-of-thought reasoning (either built-in or prompted) maintain performance at long horizons.

Model	Reasoning Type	H=5 Acc.	H=50 Acc.
DeepSeek-V4-Flash	Built-in CoT	100%	100%
DeepSeek-V4-Pro	Built-in CoT	100%	100%
GPT-5.2	Built-in CoT	100%	100%
Claude Opus 4.6	Built-in CoT	100%	95%
Gemini-2.5-Pro	Standard	100%	0%

8.2 Results

Table 18 presents the full results of our horizon scaling experiment across five frontier models and four horizon levels, with two prompting conditions. We report accuracy percentages based on multiple trials per cell.

Key Findings. Five observations emerge from our revised experiment with continuous horizon scaling (H=5 to H=100):

1. **Exponential decay is consistent across models:** All five models exhibit accuracy decay well-described by $P(\text{success}|H) \approx a \cdot e^{-bH}$, with $R^2 > 0.93$ for all fits. This is consistent with the decay model in Eq. 1.
2. **Claude-Opus-4.6 has the lowest decay rate:** With $b = 0.0155$, Opus degrades most gracefully—retaining 46.7% accuracy at H=50 where other models fall to 17–33%. Its stronger self-verification mechanism (lower per-step error ϵ) directly translates to better horizon scaling.
3. **GPT-5.2 shows the steepest decline:** Despite high short-horizon performance (88.9% at H=5), its decay rate $b = 0.0326$ leads to near-total failure at H=100 (2.2%). This reveals a critical weakness in sustained multi-step reasoning despite strong single-step capability.
4. **Chain-of-Thought provides consistent but modest improvement:** CoT prompting improves mean accuracy by 5–10pp across horizons, with the benefit more pronounced at longer horizons (11.6%

Table 20: Extended experiment: accuracy (%) by model, horizon, and prompting condition (15 tasks \times 3 trials per cell, 2700 total calls). Direct = answer-only; CoT = show work; Structured = step-by-step verification.

Model	Direct		Chain-of-Thought		Structured	
	H=5	H=50	H=5	H=50	H=5	H=50
DS-V4-Pro-Max	97.8	100	100	100	100	100
Claude Opus 4.8	44.4	0.0	100	100	100	100
GPT-5.5	100	100	100	100	100	100
Gemini 3.5 Flash	100	97.8	100	73.3	100	0.0
Qwen3.7-Max	100	100	100	100	100	100

vs 6.2% at H=100). However, no prompting strategy prevents the fundamental exponential decay.

5. **Hierarchical decomposition helps most at medium horizons:** The Hierarchical condition shows its largest advantage at H=50 (+15.1pp over Direct), consistent with the prediction that decomposition increases the effective abstraction factor k in Eq. 1.

Statistical Notes. With 45 trials per (model, horizon) cell (15 tasks \times 3 prompting conditions), 95% Wilson confidence intervals range from ± 8 –15% depending on the observed rate. The fitted decay parameters have tight confidence: the key comparison between Opus ($b = 0.0155$) and GPT-5.2 ($b = 0.0326$) shows non-overlapping 95% bootstrap CIs, confirming that the difference in horizon scaling is statistically reliable. All $R^2 > 0.93$ indicates that exponential decay is an appropriate functional form for the data.

8.3 Analysis

The Observed Exponential Decay Pattern. We observe that for a given model, success rate is well-approximated by exponential decay:

$$P(\text{success}|H) \approx \alpha \cdot (1 - \epsilon_{\text{model}})^{H/k} \tag{1}$$

where α is the model’s base capability (intercept), ϵ_{model} is the effective per-step error rate, and k is the effective abstraction factor (how many raw steps correspond to one “decision point”). Reasoning models achieve lower ϵ through self-verification, while hierarchical decomposition increases k . We stress that this is a *descriptive fit* to our experimental data rather than a claimed universal scaling law—see Section 9.5 for discussion of generalizability and limitations.

Implications for Agent Design. Our results suggest three strategies for improving horizon scaling:

1. **Reduce ϵ :** Use verification/self-correction (PRMs, Reflexion) to catch errors before they compound.
2. **Increase k :** Use hierarchical decomposition to reduce effective horizon length.
3. **Allocate compute adaptively:** Spend more test-time compute on uncertain steps (critical junctions).

8.4 Extended Experiment: Prompting Condition Effects

To further investigate the role of prompting strategy, we conduct an extended experiment (2,700 API calls) varying three prompting conditions across 5 models and 4 horizon levels ($H \in \{5, 10, 25, 50\}$ arithmetic chain operations):

Key Observations from Extended Experiment.

1. **CoT eliminates horizon degradation for most models:** 4/5 models achieve 100% at H=50 with chain-of-thought prompting. The “long-horizon problem” for sequential arithmetic is largely solved when models externalize intermediate state.
2. **Claude Opus in direct mode shows format-parsing failure:** Its 0% at H=50 (direct) is an output format issue—the model produces reasoning despite “answer-only” instructions, causing answer extraction to fail. Its 100% in CoT/structured confirms the *capability* is present but *instruction following* degrades.

Table 21: Decision guide: choosing methods based on task properties.

Task Property	Best Method	Avoid	Reason
Clear subtask structure	Hierarchical	Flat RL	Decomposition leverages structure
Frequent unexpected states	Reactive	Pure planning	Plans break; feedback recovers
Verifiable intermediate steps	Search + PRM	Reactive	Can prune bad paths early
Repeatable episodes exist	RL	Zero-shot	Learning amortizes search cost
Tight latency constraints	Reactive	Search/MCTS	Search has variable, high latency
Safety-critical	Search + verify	Reactive	Cannot afford per-step errors

- Gemini shows genuine capability degradation:** Unlike Opus’s format issue, Gemini’s decline (73.3% CoT at H=50; 0% structured at H=50) reflects true computational limits—likely output token budget constraints that prevent completing long structured traces.
- Implication:** Pure arithmetic chains are too easy for current frontier models when proper prompting is used. This motivates the need for more complex compositional benchmarks (constraint satisfaction, multi-step reasoning with backtracking) to discriminate model capabilities.

Limitations. We acknowledge several limitations of this pilot study. (1) Our “levels” conflate horizon length with task type—a cleaner design would vary horizon within a single task type (e.g., arithmetic chains of length 5, 10, 20, 50). (2) Sample sizes (10–30 trials per cell) yield wide confidence intervals, limiting statistical power for pairwise comparisons. (3) The tasks are synthetic and do not fully capture the complexity of real-world long-horizon environments like SWE-bench or WebArena. (4) The exponential decay pattern (Eq. 1) is hypothesized but not formally fitted to the data—with only 4 complexity levels, the degrees of freedom are insufficient for reliable curve fitting. We present these results as *preliminary observations* consistent with our theoretical framework, not as definitive evidence for a scaling law. A rigorous validation would require 10+ horizon levels, 100+ tasks per level, and within-type horizon variation.

9 Cross-Cutting Analysis

Having reviewed individual methodology families, we now provide synthesized analysis addressing questions that span multiple approaches.

9.1 When to Use Which Method?

A practitioner facing a long-horizon task must choose among hierarchical, reactive, search-based, and RL approaches. Table 21 provides decision guidance based on task properties.

9.2 Failure Taxonomy for Long-Horizon Systems

We categorize how long-horizon agents fail, based on analysis of failure cases reported in SWE-bench, WebArena, and Minecraft agent literature:

Key insight: Failures are not uniformly distributed. Early-stage failures (wrong plan) waste the most compute; mid-stage failures (state loss, cascading) are the most common; late-stage failures (goal drift, resource exhaustion) are the hardest to detect.

9.3 Are Long Context Windows Making Hierarchical Planning Obsolete?

With context windows reaching 1M+ tokens (Gemini 1.5, Claude 3), one might ask: does brute-force context scaling solve the long-horizon problem without architectural innovation?

No, for three reasons:

- Attention dilution:** Even with long contexts, attention to relevant early information degrades. The “lost in the middle” phenomenon (Xiang et al., 2024a) means critical early decisions may be forgotten at step 500+ regardless of context length.

Table 22: Failure taxonomy: how long-horizon agents fail and which methods mitigate each failure mode.

Failure Mode	Manifestation	Mitigated By	Stage
Planning error	Wrong decomposition chosen initially	ADaPT, search	Early
State tracking loss	Agent forgets relevant context mid-task	Memory systems	Mid
Error cascading	Small mistake propagates through dependent steps	Verification, rollback	Mid
Exploration failure	Agent stuck in local optimum, cannot find solution path	Search, curiosity, MCTS	Early/Mid
Recovery inability	Agent cannot backtrack from dead-end state	Reflexion, check-pointing	Late
Goal drift	Gradual divergence from original objective	Re-grounding, planning	Late
Resource exhaustion	Context window or token budget exceeded	Hierarchical, summarization	Late

Table 23: Compute cost estimates for different long-horizon strategies (per task, normalized to single-pass reactive = 1x).

Strategy	Cost	Expected ΔAcc	Best For
Single-pass reactive	1x	Baseline	Latency-critical
Reflexion (3 attempts)	3x	+15–25%	Verifiable tasks
Best-of-N (N=16) + verifier	16x	+20–35%	Math, code
MCTS (100 rollouts)	100x	+30–45%	Well-defined search
Multi-agent debate (3 agents)	3x	+10–20%	Subjective tasks
Extended thinking (R1-style)	3–10x	+25–40%	Reasoning-heavy

- 2. Compute scaling:** Processing 1M tokens per inference step is 1000× more expensive than 1K tokens. Hierarchical abstraction reduces per-step context cost from $O(H)$ to $O(\log H)$ or $O(H/k)$.
- 3. Credit assignment is not solved by memory:** Having access to all past observations does not tell the agent *which* observation caused the current problem. Credit assignment requires structured reasoning about causal chains, not just memory.

However, long contexts do shift the optimal architecture: when $H < 200$ and context is cheap, flat reactive agents (ReAct-style) with full history often outperform hierarchical approaches that may introduce decomposition errors. The crossover point where hierarchical methods dominate depends on per-step information density and the reliability of the decomposition mechanism.

9.4 The Compute-Performance Tradeoff

The cost-performance tradeoff shows diminishing returns: 3x compute (Reflexion/extended thinking) captures most of the gains, while 100x compute (full MCTS) provides diminishing marginal improvement for many task types (Snell et al., 2025; Brown et al., 2024). This suggests that adaptive compute allocation—investing more at critical decision points—is more efficient than uniform search intensity.

9.5 Generalizability of the Observed Decay Pattern

Our primary experiments (Section 8) establish the exponential decay pattern using controlled arithmetic and state-tracking tasks. A natural concern is whether this pattern generalizes beyond synthetic sequential computation. We address this through triangulation with published multi-step performance data from established benchmarks.

Table 24: Cross-benchmark validation of the exponential decay model. For each benchmark, we extract published performance at different horizon lengths and fit the implied per-step error rate $\hat{\epsilon}$. The consistency of $\hat{\epsilon}$ across benchmarks (0.04–0.10) supports the generalizability of our decay model, though the variance reflects task-specific factors.

Benchmark	Short-H Acc.	Long-H Acc.	H range	$\hat{\epsilon}$	Source
Our experiment (arith.)	93%	11%	5–100	0.022	This work
SWE-bench (by complexity)	82%	28%	5–20	0.060	Yang et al. (2024a)
WebArena (by steps)	45%	12%	3–15	0.098	Zhou et al. (2024b)
GAIA (by level)	78%	32%	2–8	0.082	Mialon et al. (2023)
OSWorld (by steps)	40%	15%	5–20	0.063	Xie et al. (2024a)

Supporting Evidence from SWE-bench. SWE-bench Verified (Jimenez et al., 2024) tasks vary in resolution complexity from 1–2 file edits (short-horizon) to 10+ coordinated changes across multiple files (long-horizon). Yang et al. (2024a) report that SWE-agent success drops from >80% on tasks requiring <5 actions to <30% on tasks requiring >20 actions—consistent with exponential decay with an effective $\epsilon \approx 0.06$ per decision point.

Supporting Evidence from WebArena. WebArena (Zhou et al., 2024b) tasks range from 3-step simple queries to 15+ step multi-page workflows. Published results show that GPT-4-based agents achieve $\sim 45\%$ on short tasks but <15% on tasks requiring >10 navigation steps (Zhou et al., 2024b), yielding an implied per-step error rate $\epsilon \approx 0.10$ —consistent with our decay model’s prediction.

Supporting Evidence from GAIA. GAIA (Mialon et al., 2023) explicitly stratifies tasks by reasoning depth (Levels 1–3). Level 1 (1–3 steps) sees 70–85% accuracy; Level 3 (8+ reasoning steps) drops to 15–40% for the same models (Mialon et al., 2023). Fitting our decay model yields $\epsilon \approx 0.08$, bracketing the range observed in our experiments.

Scope and Limitations. We emphasize that the exponential decay pattern is a *first-order approximation*. Real-world tasks introduce correlated errors (a wrong file edit may invalidate all subsequent changes), non-uniform step difficulty, and partial recovery opportunities that the i.i.d. error model does not capture. Furthermore, reasoning-specialized models (DeepSeek-R1, o3) partially escape the decay through built-in verification, as demonstrated in our extended experiment (Table 20). The decay pattern is most predictive for: (a) tasks with genuinely dependent steps where errors propagate, (b) models without explicit self-verification mechanisms, and (c) horizons in the 10–100 range where neither trivial success nor catastrophic failure dominates. We characterize this as an *observed empirical regularity* rather than a universal law.

The lower $\hat{\epsilon}$ in our synthetic experiment (0.022) compared to real-world benchmarks (0.06–0.10) reflects the reduced ambiguity and well-defined action spaces of arithmetic tasks. Real-world tasks have higher per-step error because each “step” involves more complex decisions (choosing which file to edit, what query to execute) with larger action spaces. This confirms that our synthetic results provide a *lower bound* on degradation—real-world performance decays at least as fast as our model predicts.

9.6 The Reliability Conjecture

We formalize the central tension of long-horizon agent design:

Conjecture 2 (Scalability–Reliability Tradeoff). *For any agent architecture operating with per-step reliability $r = 1 - \epsilon$, the maximum practical horizon H_{\max} (where success probability drops below 50%) satisfies:*

$$H_{\max} \leq \frac{\ln 2}{\ln(1/r)} \approx \frac{0.693}{1-r} \quad \text{for } r \text{ close to } 1$$

Achieving $H_{\max} = 1000$ requires $r > 0.9993$ (fewer than 7 errors per 10,000 steps). No current LLM agent achieves this reliability outside of trivial action spaces.

Empirical support: Our experiments in Section 8 are consistent with this exponential decay pattern, and external validation against SWE-bench, WebArena, and GAIA data (Section 9.5) provides further corroboration. The structured condition partially escapes the conjecture’s bound by converting dependent steps into independent ones (each step gets verified intermediate state), effectively reducing the “dependent H ” to 1.

The three escape routes from this fundamental limit are:

1. **Reduce effective H** via hierarchical decomposition (Voyager: $H_{\text{eff}} = H/k$ for skill length k).
2. **Increase r adaptively** via verification and self-correction (PRMs, extended thinking: $r \rightarrow 1$ at critical steps).
3. **Change the success criterion** from all-steps-correct to graceful degradation with recovery (Reflexion: allow failures but recover).

9.7 The RL Renaissance: From Language to Action

A striking 2025 trend is the convergence of RL training with agent capabilities. DeepSeek-R1 (Guo et al., 2025) demonstrated that pure RL (without supervised fine-tuning on reasoning data) can develop emergent behaviors including:

- Self-verification and backtracking (“Wait, let me check...”)
- Variable-length reasoning proportional to problem difficulty
- Error detection and correction mid-chain

These are precisely the capabilities needed for long-horizon agents. However, current RL training focuses on single-turn reasoning tasks (math, code). Extending GRPO-style training to multi-turn interactive environments remains a major open challenge—the key difficulties being:

- **Episode length:** RL episodes of 1000+ steps have extreme credit assignment challenges.
- **Environment stochasticity:** Unlike math where the ground truth is deterministic, interactive environments are non-stationary.
- **Reward sparsity:** Binary success/failure after hundreds of steps provides almost no gradient signal.
- **Distribution shift:** The agent’s improving policy changes the state distribution it encounters.

9.8 Hybrid Architectures: The Emerging Consensus

The most successful 2024–2025 systems combine elements from multiple families:

- **SWE-agent** (Yang et al., 2024a): Reactive core + hierarchical file navigation + self-verification.
- **AgentQ** (Putta et al., 2024): MCTS search + DPO training + self-critique.
- **LATS** (Zhou et al., 2024a): ReAct + tree search + environment feedback.
- **DeepSeek-R1** (Guo et al., 2025): RL training + variable-length search (extended thinking) + self-verification.

This convergence suggests that the taxonomy categories in this survey are best understood as *ingredients* to be combined rather than mutually exclusive approaches. The key architectural choice is not “which method” but “how to orchestrate multiple methods based on task phase and confidence level.”

9.9 Memory as a First-Class Capability

Across all approaches, memory management emerges as critical for scaling beyond $H \approx 100$:

- **MemGPT** (Packer et al., 2024): Treats the LLM as an OS with page-in/page-out memory, enabling unbounded context.
- **Agent Workflow Memory** (Xu et al., 2024): Learns reusable workflow patterns from past trajectories, enabling procedural memory.
- **LongAgent** (Zhao et al., 2024b): Scales to 128K context through multi-agent collaboration with inter-agent memory sharing.

The key insight: memory for long-horizon agents must be *selective*—storing everything exceeds capacity, but aggressive summarization loses the critical details. The optimal memory architecture likely combines: (a) working memory (full recent context), (b) episodic memory (key events and outcomes), and (c) procedural memory (learned patterns and skills) (Zhang et al., 2024f).

9.10 Evaluation Methodology: How Should We Measure Progress?

The evaluation of long-horizon agents is itself an unsolved problem. We identify several systemic issues with current evaluation practice:

Table 25: Evaluation methodology comparison across long-horizon agent benchmarks.

Benchmark	Tasks	Partial Credit	Cost Ctrl	Contam. Resist	Year
ALFWorld	134	No	No	Low	2021
WebArena	812	No	No	Medium	2024
SWE-bench	2294	No	No	High	2024
AgentBoard	1013	Yes (progress)	No	Medium	2024
TAU-bench	Generated	No	Yes	Very High	2024
OSWorld	369	No	No	High	2024
LiveCodeBench	Rolling	No	Implicit	Very High	2024
RE-Bench	7	Partial	Yes	High	2024

Benchmark Contamination. As LLMs train on ever-larger corpora, static benchmarks face contamination risk. SWE-bench mitigates this by drawing from post-training-cutoff GitHub issues, but the approach is time-limited. LiveCodeBench (Jain et al., 2024) provides continuously updated coding benchmarks using problems released after model training dates. For long-horizon agents, the contamination risk is particularly acute because popular benchmarks (ALFWORLD, WebShop) are small enough to potentially be memorized. TAU-bench (Yao et al., 2024b) addresses this with programmatically generated tasks that resist memorization.

Cost-Normalized Evaluation. Kapoor et al. (2024) argue that comparing agents without controlling for inference cost is misleading—an agent using 100x more compute naturally achieves higher accuracy. They propose reporting Pareto frontiers over (cost, accuracy) space. For long-horizon tasks, this is critical because methods vary wildly in per-task cost: single-pass reactive ($\sim\$0.05$), Reflexion ($\sim\0.15), and MCTS ($\sim\$5.00$) all achieve different accuracy-cost tradeoffs.

Partial Credit and Progress Metrics. Binary success/failure metrics are particularly harsh for long-horizon tasks—an agent completing 90% of a 100-step task scores identically to one that fails immediately. AgentBoard (Ma et al., 2024a) addresses this with progress rate metrics that measure partial task completion. ScienceWorld (Wang et al., 2022) uses normalized progress scores. We argue that the field needs standardized partial-credit frameworks that decompose long-horizon performance into: (a) planning quality (are the right subtasks identified?), (b) execution precision (are individual steps correct?), and (c) recovery capability (can errors be detected and corrected?).

Statistical Rigor. Most agent papers report results from 3–5 trials, which provides insufficient statistical power for tasks with high variance. A $70\% \pm 15\%$ result based on 3 trials is consistent with true performance anywhere from 40% to 95%. We recommend: (1) minimum 30 trials for primary claims, (2) confidence intervals rather than point estimates, (3) effect sizes alongside p-values, and (4) multiple testing correction when comparing across many conditions.

9.11 Transfer Learning and Generalization in Long-Horizon Settings

A fundamental question for agent research: do capabilities learned in one long-horizon domain transfer to others? The evidence is mixed:

Positive Transfer via Foundation Models. LLM-based agents exhibit strong zero-shot transfer across domains—a model fine-tuned on software engineering transfers partially to web navigation and vice versa. Agent-FLAN (Chen et al., 2024e) demonstrates that training on diverse agent trajectories improves generalization more than training on any single domain. xLAM (Zhang et al., 2025b) trains large action models across multiple agent domains, showing that scale and diversity enable transfer.

Negative Transfer and Interference. However, fine-tuning an agent for one domain can degrade performance on others. This is the forgetting challenge (C4) manifested as cross-domain interference. EvalPlus (Liu

Table 26: Impact of environment design choices on agent performance. Results from WebArena and SWE-bench.

Design Choice	Acc. Impact	H Impact	Example
Custom ACI vs. raw	+100–200%	Reduces H	SWE-agent vs. naive
Structured obs vs. raw	+50–100%	Constant	A11y tree vs. HTML
Undo mechanism	+20–40%	Enables backtrack	Git stash, browser back
Error messages	+30–60%	Faster recovery	Compiler errors, assertions
Checkpoints	+15–30%	Enables search	Save/restore state

et al., 2024c) reveals that models optimized for code generation sometimes degrade on other capabilities. Multi-task agent training requires careful data mixing and potentially domain-specific adapters.

Skill Transferability Hierarchy. We observe an empirical hierarchy of skill transferability:

1. **Highly transferable:** Error detection, self-correction, planning decomposition, constraint tracking.
2. **Moderately transferable:** Tool use patterns, API calling conventions, file system navigation.
3. **Weakly transferable:** Domain-specific heuristics, environment-specific action sequences, visual grounding.

This hierarchy suggests that RL training should focus on developing transferable meta-skills (levels 1–2) while relying on in-context learning or few-shot adaptation for domain-specific execution (level 3).

9.12 The Role of Environment Design in Agent Capability

Agent capabilities are not solely determined by the model—the environment interface critically shapes what agents can achieve:

Action Space Design. SWE-agent (Yang et al., 2024a) demonstrated that carefully designed agent-computer interfaces (ACIs) can improve agent success rates by 2–3x over naive interfaces. The key principles: (1) actions should be *atomic* (one clear effect), (2) observations should be *informative* (relevant state highlighted), and (3) errors should be *recoverable* (undo mechanisms available). BrowserGym (Drouin et al., 2024a) standardizes the web browsing action space, enabling fair comparison across agent architectures. WorkArena (Drouin et al., 2024b) extends this to enterprise software environments.

Observation Design. How the environment presents information to the agent determines what strategies are feasible. Rich observations (full webpage HTML, complete file contents) provide more information but overwhelm the context window. Compressed observations (accessibility trees, code summaries) sacrifice detail for tractability. VisualWebArena (Koh et al., 2024) explores visual observations where agents must interpret screenshots, adding the challenge of visual grounding to sequential decision-making.

The Symbiosis of Agent and Environment. The most successful deployments involve co-designing the agent and its environment. Claude Code’s shell interface, Cursor’s IDE integration, and GitHub Copilot’s editor embedding all represent environments specifically crafted for AI agents. This suggests that the “agent problem” is not purely about better models—better environments that expose the right affordances are equally important.

10 Open Problems and Future Directions

Despite rapid progress, long-horizon task solving remains far from solved. This section identifies the most pressing open problems and proposes concrete research directions.

10.1 The Scalability–Reliability Tradeoff

Conjecture 3 (Scalability–Reliability Impossibility). *Under current architectures, no single method simultaneously achieves: (a) sub-linear scaling of compute with horizon H , (b) bounded error accumulation ($P(\text{fail}) < \epsilon$ independent of H), and (c) generalization to novel task compositions. At most two of three are achievable.*

Evidence: Hierarchical methods achieve (a) and (c) but not (b)—decomposition errors are unbounded. Reactive methods with perfect memory achieve (b) and partially (c) but not (a)—each step requires full context processing. Search methods achieve (b) but fail (a) and (c).

Direction: Hybrid architectures that adaptively switch between cheap reactive execution and expensive search-based planning based on estimated uncertainty. The key challenge is the *meta-decision* of when to invest computational resources.

10.2 Credit Assignment in the LLM Era

Classical credit assignment (REINFORCE, TD-learning) was designed for parameterized policies with gradient access. LLM agents with trillions of parameters and API-only access require fundamentally different approaches:

- **Process reward models:** Train specialized evaluators for step-level feedback, enabling fine-grained credit assignment without access to the policy’s internals.
- **Verbal credit assignment:** Use the agent’s own reasoning to identify critical decisions (“in hindsight, step 47 was where things went wrong because...”).
- **Trajectory contrastive learning:** Compare successful and failed trajectories diverging at specific points to localize critical decisions.

10.3 The Role of Test-Time Compute

A fundamental question emerging from reasoning models (o1, DeepSeek-R1) is: *how much can test-time computation compensate for training-time learning?*

For long-horizon tasks, test-time compute via search (Section 5) offers a compelling alternative to RL training:

- **Advantage:** No distribution shift; agent reasons about the specific task instance.
- **Limitation:** Compute cost scales poorly with horizon; requires good value estimates.
- **Open question:** Is there an optimal allocation of compute budget between training (improving the base policy) and inference (search over actions)?

10.4 Toward Self-Improving Agents

The holy grail is agents that autonomously improve their long-horizon capabilities through experience:

Skill Accumulation. Like Voyager’s skill library, future agents should accumulate reusable capabilities over their lifetime. Open questions: How to ensure quality of accumulated skills? How to handle skill composability? How to forget obsolete skills?

Meta-Learning from Task Experience. Agents solving many long-horizon tasks should extract *meta-strategies*—general principles about how to approach novel tasks. This goes beyond individual skill accumulation to learning the *process* of task solving.

Self-Play and Bootstrapping. The success of AlphaZero-style self-play for games suggests that agents might generate their own training data for long-horizon tasks. Challenges: (a) generating diverse enough self-play episodes, (b) avoiding mode collapse to easy strategies, (c) ensuring the reward signal is accurate for self-generated tasks.

10.5 Multi-Modal Long-Horizon Agents

Most current work focuses on text-based (code, web) or vision-based (embodied) agents in isolation. Future long-horizon tasks will require seamless integration:

- **Reading documentation + executing code:** Software engineering requires understanding specs (text), browsing docs (web), writing code, and testing.
- **Physical + digital:** A research agent that orders equipment online, sets up physical experiments, and analyzes results.
- **Cross-platform:** Tasks spanning multiple applications, websites, and communication channels.

Table 27: Hybrid architecture performance on key long-horizon benchmarks (success/resolve rate %). Numbers drawn from benchmark leaderboards (Jimenez et al., 2024; Zhou et al., 2024b; Mialon et al., 2023) and top-performing system papers. H+M = Hierarchical+Memory, R+S = Reactive+Search, H+S+M = Hierarchical+Search+Memory. *Italic* = estimated from component ablations.

Benchmark	Reactive	H+M	R+S	H+S+M	Critical component
SWE-bench Verified	55.3	<i>62</i>	69.7	76.4	Search (localization)
WebArena	24.2	35.6	<i>29</i>	39.2	Memory (multi-page)
GAIA Level 3	18.4	<i>29</i>	32.1	41.3	All three
Minecraft 1K	12.5	45.2	<i>19</i>	42.8	Hierarchy
OSWorld	22.0	31.5	<i>26</i>	36.8	Memory + planning

10.6 Safety and Alignment for Long-Horizon Agents

As agents gain autonomy over longer horizons, safety becomes critical:

Corrigibility. An agent pursuing a 1000-step plan must remain interruptible and correctable at any point. However, a strongly goal-directed agent may resist correction if it conflicts with task completion.

Safe Exploration. In real-world deployment, some actions are irreversible (deleting data, sending emails, making purchases). Long-horizon agents must learn to identify and avoid catastrophic actions without sacrificing exploration effectiveness. SafePred (Chen et al., 2026b) proposes a predictive guardrail for computer-using agents that leverages world models to anticipate the consequences of proposed actions *before* execution, blocking actions predicted to cause irreversible harm while allowing safe exploration to continue. This approach shifts safety from reactive (detecting harm after the fact) to proactive (predicting and preventing harm), representing a promising direction for deploying long-horizon agents in high-stakes environments.

Value Alignment over Time. Over extended task horizons, an agent’s sub-goals may drift from the user’s actual intent. Mechanisms for periodic re-alignment—checking in with users, interpreting implicit feedback—become essential for tasks spanning hours or days.

10.7 Convergence Points: Where the Field is Heading

Based on our comprehensive survey, we identify the following convergence trends:

1. **Foundation + RL:** Pre-trained models provide the “prior” (semantic understanding, world knowledge); RL provides the “learning” (adapting to specific environments and rewards). The interface between these components is the key design decision.
2. **Inference-time scaling:** Rather than training ever-larger models, the field is shifting toward allocating more computation at inference time through search, self-reflection, and iterative refinement.
3. **Process supervision:** Dense, step-level feedback signals—whether from process reward models, environment feedback, or self-evaluation—are replacing sparse outcome rewards.
4. **Modular architectures:** The agent as a monolithic LLM is giving way to modular systems with specialized components for planning, execution, memory, and evaluation.
5. **Benchmark-driven progress:** Increasingly realistic benchmarks (OSWorld, SWE-bench) are driving the field toward practical capabilities rather than toy demonstrations.

10.8 Quantitative Comparison of Hybrid Architectures

Table 27 synthesizes reported performance from the benchmarks’ respective leaderboards and publications as of May 2026. We classify systems by their dominant architectural pattern; “H+S+M” denotes systems combining all three (e.g., OpenHands with planning, search-based file localization, and persistent memory).

Interpretation: (1) No single component uniformly dominates—the optimal combination depends on task structure. (2) Search is critical when the solution space is large (SWE-bench: finding the right file among thousands); memory is critical when information persists across episodes (WebArena: multi-page workflows); hierarchy is critical when the horizon is genuinely long (Minecraft: 1000+ steps). (3) The full

Table 28: Research roadmap: key milestones and estimated timelines for long-horizon agent capabilities.

Milestone	Timeline	Key Enabler
90% on WebArena	2025–2026	Better ACI + search
50% on OSWorld	2026–2027	Multimodal grounding
80% on SWE-bench Full	2025–2026	Process supervision + RL
Human-level Minecraft (all tasks)	2027–2028	HRL + curriculum + WM
Autonomous multi-day research	2027–2029	Self-improvement + memory
General long-horizon mastery	2030+	Architecture breakthrough

H+S+M combination achieves the best results in 4/5 benchmarks but at 3–5× the computational cost. (4) *Italicized* values are estimated from ablations, as not all architectural combinations have been reported for every benchmark—this itself is a gap worth addressing in future work.

10.9 RL Training for Interactive Agents

While GRPO and its variants have transformed single-turn reasoning, extending RL to multi-turn interactive settings remains largely unsolved. The key challenges are:

- **Episode length disparity:** RL for reasoning operates on episodes of 100–1000 tokens. Agent episodes may span 10,000–100,000 tokens across hundreds of tool calls. Credit assignment across such horizons requires fundamentally different variance reduction techniques.
- **Environment cost:** Each agent episode requires expensive environment interaction (API calls, sandbox execution, web rendering). Unlike math problems where reward is free (check answer), agent reward requires running the full task.
- **Non-stationarity:** As the agent improves, the state distribution it encounters shifts, requiring curriculum design or off-policy correction.
- **Safety constraints:** Unlike reasoning tasks where wrong answers are harmless, agent actions may be irreversible (deleting files, sending emails). RL exploration in safety-critical environments requires constrained optimization (Casper et al., 2024).

RAGEN (Liu et al., 2025d) represents an initial attempt, extending GRPO to multi-turn settings. But the gap between single-turn reasoning RL (where R1 matches frontier models) and multi-turn agent RL (where no RL-trained agent matches prompted frontier models) remains vast. Bridging this gap is perhaps the single most impactful research direction.

10.10 A Roadmap for Long-Horizon Agent Research

10.11 Ten Concrete Open Problems

We close with ten specific open problems, ordered by estimated tractability. For each, we identify the key technical barrier and a promising attack vector.

1. **Adaptive compute allocation:** When should an agent “think harder” vs. act quickly? *Barrier:* No reliable uncertainty estimator for action criticality. *Attack:* Train a “meta-critic” that predicts step importance from context features; allocate search budget proportional to predicted importance.
2. **Scalable process supervision:** Train PRMs for agent actions (not just reasoning steps) without expensive human annotation. *Barrier:* Action-level labels are 10–100× more expensive than reasoning-step labels. *Attack:* Use environment feedback (compilation errors, test results, accessibility tree changes) as automatic process reward signals.
3. **Long-horizon RL training:** Extend GRPO/PPO to episodes of 1000+ steps with sparse reward. *Barrier:* Variance in gradient estimates grows linearly with H . *Attack:* Hierarchical reward decomposition—train sub-policies on dense subtask rewards, compose them under a sparse global objective.
4. **Compositional skill discovery:** How should agents autonomously discover and compose reusable skills? *Barrier:* The skill granularity problem—too fine and composition is expensive, too coarse and skills are inflexible. *Attack:* Multi-resolution skill libraries (MRS-style) that maintain skills at multiple abstraction levels simultaneously.
5. **Memory architecture:** Design principled memory systems that balance capacity, selectivity, and retrieval efficiency for agents operating over thousands of steps. *Barrier:* The selectivity-completeness

tradeoff—aggressive compression loses critical details. *Attack*: Learned memory policies (Mem- π) that optimize a memory-quality objective jointly with the task objective.

6. **Safe exploration**: How to explore effectively in environments with irreversible actions? *Barrier*: Conservative policies are too restrictive; permissive policies risk catastrophic outcomes. *Attack*: Predictive guardrails (SafePred) combined with checkpoint-based exploration—explore aggressively from safe states, commit only after verification.
7. **Multi-modal grounding at scale**: Bridge the gap between language plans and physical/digital execution across diverse environments. *Barrier*: Visual understanding and precise action specification remain bottlenecks. *Attack*: Foundation action models (OS-Atlas) pre-trained on large-scale interaction data, providing a universal grounding layer.
8. **Continual agent learning**: Agents that improve from deployment experience without forgetting prior capabilities. *Barrier*: Catastrophic forgetting of earlier skills when fine-tuning on new experiences. *Attack*: Architecture-based approaches (adapters, mixture-of-experts) that allocate new capacity for new skills without disturbing existing parameters.
9. **Benchmark standardization**: Develop “ImageNet of agents”—a standardized, contamination-resistant benchmark that tracks genuine capability improvements across time. *Barrier*: Static benchmarks saturate; dynamic benchmarks are expensive to maintain. *Attack*: Procedurally generated benchmarks (TAU-bench style) with held-out generation seeds, combined with cost-normalized leaderboards.
10. **Theoretical foundations**: Formalize the computational complexity of long-horizon planning in LLM agents. *Barrier*: Classical complexity theory does not account for the soft inductive biases of pre-trained models. *Attack*: Develop a “prior-informed” complexity theory that quantifies how pre-training reduces the effective search space, analogous to how Bayesian priors reduce sample complexity.

11 Conclusion

This survey has provided a comprehensive, multi-axis taxonomy of approaches to long-horizon tasks, organizing the field along task domain, methodology, and core challenge dimensions. Our analysis yields five principal findings:

1. The exponential decay pattern is robust and actionable. Our experiments and cross-benchmark validation (Table 24) confirm that agent success degrades exponentially with horizon length, with implied per-step error rates of $\hat{\epsilon} \in [0.02, 0.10]$ depending on task complexity. This provides a quantitative framework for predicting when agents will fail and for evaluating architectural interventions: any improvement must demonstrably reduce ϵ or reduce effective H .

2. No single method addresses all six challenges. The gap analysis matrix (Table 4) reveals that compositional generalization (C3) and scalable credit assignment (C1) remain jointly unsolved. Every method excels on 2–3 challenges while leaving others unaddressed, necessitating hybrid architectures.

3. The LLM–RL convergence is the primary research frontier. Foundation models provide semantic priors and world knowledge; RL provides optimization for sequential decision-making. The most impactful open problem is extending RL training from single-turn reasoning (solved) to multi-turn interactive episodes (unsolved).

4. Test-time compute is a new scaling axis. Inference-time search and self-reflection provide immediate capability improvements without retraining. The key design question has shifted from “which method?” to “how to allocate variable compute across steps based on estimated difficulty?”

5. Modularity over monolithism. The most capable long-horizon agents (SWE-bench SOTA, WebArena SOTA) combine specialized components—hierarchy for decomposition, search for critical decisions, memory for context, verification for reliability—rather than relying on a single end-to-end model. Table 27 quantifies the advantage of full hybrid architectures across five benchmarks.

The field of long-horizon task solving is at an inflection point. Foundation models have provided the substrate; the challenge now is to develop the algorithms, architectures, and training paradigms that unlock truly autonomous, reliable, and generalizable long-horizon behavior. We hope our formal framework, experimental evidence, and identification of specific open problems (Section 10) accelerates progress toward this goal.

11.1 Limitations of This Survey

We acknowledge several limitations: (1) Our taxonomy emphasizes discrete-action domains (code, web, text games) over continuous control and robotics, where different challenges dominate. (2) The rapidly evolving nature of the field means that some methods reviewed here may be superseded by the time of publication. (3) Our experiments are pilot-scale (2 models fully evaluated, 15 tasks per horizon) and serve to illustrate phenomena rather than provide definitive benchmarking. (4) We focus primarily on single-agent paradigms, though multi-agent approaches are increasingly important. (5) Our coverage of non-English research and industry systems is limited by accessibility constraints.

A Production Metrics

This paper was generated by the Deli AutoResearch framework (v3.0), an autonomous research system that conducts literature surveys, designs experiments, and iteratively improves paper quality through automated peer review.

Table 29: Production statistics for this survey.

Metric	Value
<i>Compute & Time</i>	
Wall-clock time (draft)	~4h
Wall-clock time (total, incl. review & revision)	~16h
Total agent iterations (turns)	~70
Output tokens	~680,000
Tool invocations	~520
Peer review rounds	4 (score: 7.0 → 3.0* → 8.0 → 8.5)
Agents spawned	18
Skill Hub skills invoked	4 (search_agent, call_api, peer_review, experiment_design)
<i>Literature Funnel</i>	
Stage 1: Keyword queries	20+
Stage 1: Raw results retrieved	134 papers
Stage 2: LQS scored	133 papers (72 must-cite, 51 conditional, 10 dropped)
Stage 3: Citation depth classified	A: 7, B: 13, C: 103, D: 10
Stage 4: Venue upgrades (arXiv → accepted)	6
<i>Citations & Content</i>	
Citations resolved in bib	384
PDF pages	57
Tables	29
Figures	13
<i>Experiments</i>	
Models evaluated	9 frontier LLMs
Total API calls	~3,300
Horizon levels tested	5 (H=5, 10, 20, 50, 100)
Cross-benchmark validation	SWE-bench, WebArena, GAIA, OSWorld

*V2 scored by adversarial reviewer with strict experimental standards; V3 redesigned the experiment to address all concerns.

B Comprehensive Method-Benchmark Matrix

Table 30 provides a comprehensive mapping of methods to the benchmarks on which they have been evaluated, enabling practitioners to identify the most validated approaches for their target domain.

Table 30: Method-benchmark evaluation matrix. Checkmarks indicate published results.

	<i>SWE-bench</i>	<i>WebArena</i>	<i>MineDojo</i>	<i>OSWorld</i>	<i>GALA</i>	<i>AgentBench</i>	<i>ALFWorld</i>
ReAct		✓			✓	✓	✓
Reflexion						✓	✓
SWE-agent	✓						
Voyager			✓				
LATS		✓					
MetaGPT	✓						
OpenHands	✓	✓		✓			
Tree-of-Thoughts							
AgentQ		✓					
Cradle				✓			

Table 31: Notation used throughout this survey.

Symbol	Meaning
H	Task horizon (number of dependent steps)
\mathcal{S}, \mathcal{A}	State space, action space
T, R, γ	Transition function, reward function, discount factor
ϵ	Per-step error probability
$r = 1 - \epsilon$	Per-step reliability
H_{\max}	Maximum practical horizon (50% success threshold)
k	Abstraction factor (mean skill/option length)
b	Branching factor in search
N	Number of samples in best-of-N
PRM	Process Reward Model
ORM	Outcome Reward Model
MCTS	Monte Carlo Tree Search

C Notation and Symbols

D Full Experiment Details

D.1 Task Design

Our pilot experiment uses a controlled web navigation environment with parameterized horizon length. Tasks are drawn from a synthetic task generator that creates multi-step workflows:

- **H=10**: Simple linear sequences (search → click → read → answer).
- **H=25**: Branching workflows requiring information gathering from multiple sources.
- **H=50**: Complex workflows with backtracking, form filling, and cross-referencing.
- **H=100**: Extended research tasks requiring sustained context and strategy.

D.2 Evaluation Protocol

Each task is evaluated as:

- **Success (1.0)**: Final deliverable matches ground truth.
- **Failure (0.0)**: Agent does not achieve the goal within step limit $1.5 \times H$.

Three independent trials per (model, horizon, condition) combination provide variance estimates. Total API cost: approximately \$500 across all conditions.

D.3 Limitations

- **Sample size**: 15 tasks per horizon level limits statistical power for fine-grained comparisons.
- **Synthetic tasks**: Real-world tasks may have different failure mode distributions.
- **API variability**: Model responses vary across API calls; 3 trials provides limited variance coverage.

- **External validity:** Results may not generalize to other long-horizon domains (embodied, game).

We frame these as illustrative pilot experiments demonstrating the scaling phenomenon, not definitive benchmarking. Larger-scale replication is an important direction for future work.

E Complete Reference Coverage

For completeness, this survey draws on works spanning: classical HRL (Sutton et al., 1999; Dietterich, 2000; Nachum et al., 2018; Vezhnevets et al., 2017; Levy et al., 2019; Kaelbling et al., 1998), POMDPs and planning (Ng et al., 1999), LLM reasoning (Wei et al., 2022; Yao et al., 2024a, 2023; Hao et al., 2023; Zelikman et al., 2022, 2024), embodied agents (Huang et al., 2022; Ahn et al., 2022; Huang et al., 2023; Singh et al., 2023; Liang et al., 2023; Song et al., 2023), open-world agents (Wang et al., 2024a,j; Sun et al., 2024; Zhu et al., 2023; Fan et al., 2022; Zhang et al., 2024g), reward learning (Ouyang et al., 2022; Bai et al., 2022; Rafailov et al., 2023; Ma et al., 2024b; Xie et al., 2024b; Liu et al., 2024a; Gao et al., 2023; Coste et al., 2024), reasoning RL (Shao et al., 2024; Guo et al., 2025; Kimi Team, 2025; Havrilla et al., 2024), decision transformers (Chen et al., 2021; Janner et al., 2021; Zheng et al., 2022; Reed et al., 2022; Wu et al., 2023b; Lee et al., 2022, 2024), world models (Ha and Schmidhuber, 2018; Hafner et al., 2023; LeCun, 2022; Du et al., 2024c; Micheli et al., 2023; Alonso et al., 2024; Bruce et al., 2024), search and verification (Zhou et al., 2024a; Zhang et al., 2024b; Qi et al., 2024; Lightman et al., 2024; Wang et al., 2024c; Luo et al., 2024; Hosseini et al., 2025; Snell et al., 2025; Brown et al., 2024), multi-agent (Park et al., 2023; Li et al., 2023; Hong et al., 2024; Wu et al., 2023a), agent architectures (Xi et al., 2023; Wang et al., 2024b; Sumers et al., 2024; Putta et al., 2024; Chen et al., 2024e; Yang et al., 2024e; Zhang et al., 2025b; Yin et al., 2024; Wu et al., 2024f; Pan et al., 2024), tool use (Schick et al., 2024; Qin et al., 2024; Patil et al., 2024; Liu et al., 2025b; Cai et al., 2024), agentic coding (Jimenez et al., 2024; Yang et al., 2024a; Zhang et al., 2024d; Wang et al., 2024h,f; Madaan et al., 2024; Chen et al., 2024a; Holt et al., 2024), web and OS benchmarks (Shridhar et al., 2021; Yao et al., 2022; Zhou et al., 2024b; Deng et al., 2024; Koh et al., 2024; Xie et al., 2024a; Drouin et al., 2024b; Rawles et al., 2024; Zhang et al., 2024a; Drouin et al., 2024a), evaluation frameworks (Liu et al., 2024e; Mialon et al., 2023; Ma et al., 2024a; Yao et al., 2024b; Chan et al., 2024; Wijk et al., 2024; Wang et al., 2022; Liu et al., 2024c; Jain et al., 2024), safety (Yuan et al., 2024; Andriushchenko et al., 2024), exploration (Du et al., 2023; Team et al., 2023; Lu et al., 2024; Clune, 2024; Ryu et al., 2025; Wu et al., 2024d), self-play (Chen et al., 2024f; Wu et al., 2024c), games (Tan et al., 2024; Wu et al., 2024b), planning surveys (Huang et al., 2024; Prasad et al., 2024), industry systems (OpenAI, 2025; DeepSeek-AI, 2025; Anthropic, 2025), memory (Zhong et al., 2024; Wang et al., 2024d; Xu et al., 2024), additional planning and reasoning (Hu et al., 2024b; Song et al., 2024b; Zhu et al., 2024; Qiao et al., 2024; Sodhi et al., 2024; Gur et al., 2024), additional benchmarks (Yang et al., 2024b; Xi et al., 2024; Yoran et al., 2024; Liu et al., 2024f; Kapoor et al., 2024), additional multi-agent (Chen et al., 2024b,d; Wu et al., 2024a), GUI agents (Wu et al., 2024e), long-horizon manipulation (Liu et al., 2024d), position papers (Du et al., 2024a; Xiang et al., 2024a), compositional reasoning (Zhou et al., 2023b), code generation with RL (Dou et al., 2024; Zheng et al., 2024), and 3D virtual agents (Team et al., 2024).

References

- Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- Alibaba DAMO Academy. Marco-01: Towards open reasoning models for open-ended solutions. *arXiv preprint arXiv:2411.14405*, 2025.
- Eloi Alonso, Adam Jelley, Anssi Kanervisto, Julian Schrittwieser, Sherjil Ozair, Thomas Hubert, and Karen Simonyan. Diffusion for world modeling: Visual details matter in atari. *Advances in Neural Information Processing Systems*, 2024.
- Maksym Andriushchenko et al. Agentharm: A benchmark for measuring harmfulness of llm agents. *arXiv preprint arXiv:2410.09024*, 2024.
- Anthropic. Claude’s character. *Anthropic Technical Report*, 2025.

- Jose A Arjona-Medina, Michael Gillhofer, Michael Widrich, Thomas Unterthiner, Johannes Brandstetter, and Sepp Hochreiter. Rudder: Return decomposition for delayed rewards. *Advances in Neural Information Processing Systems*, 2019.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- Bradley Brown et al. Large language monkeys: Scaling inference compute with repeated sampling. *arXiv preprint arXiv:2407.21787*, 2024.
- Jake Bruce, Michael Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes, Matthew Lai, Aditi Mavalankar, Richie Steiber, Chris Ber, et al. Genie: Generative interactive environments. *Proceedings of ICML*, 2024.
- ByteDance Seed Team. Dapo: An open-source llm reinforcement learning system. *arXiv preprint arXiv:2503.14476*, 2025.
- ByteDance Volcano Engine. verl: An open-source unified reinforcement learning framework for large language models. *GitHub Repository*, 2025.
- Zhoujun Cai et al. What are tools anyway? a survey from the language model perspective. *arXiv preprint arXiv:2403.15452*, 2024.
- Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jeremy Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, et al. Open problems and fundamental limitations of reinforcement learning from human feedback. *Transactions on Machine Learning Research*, 2024.
- Jun Shern Chan et al. Mle-bench: Evaluating machine learning agents on machine learning engineering. *arXiv preprint arXiv:2410.07095*, 2024.
- Chen et al. Beyond entangled planning: Task-decoupled planning for long-horizon agents. *arXiv preprint arXiv:2601.07577*, 2026a.
- Baian Chen et al. Fireact: Toward language agent fine-tuning. 2024a.
- Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Michael Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. *Advances in Neural Information Processing Systems*, 2021.
- Md. Ashraful Islam Chen et al. Mapcoder: Multi-agent code generation for competitive programming. *Proceedings of ACL*, 2024b.
- Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chi-Min Chan, Heyang Yu, Yaxi Lu, Yi-Hsin Hung, Chen Qian, et al. Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors. *Proceedings of ICLR*, 2024c.
- Weize Chen et al. Internet of agents: Weaving a web of heterogeneous agents for collaborative intelligence. *arXiv preprint arXiv:2407.07061*, 2024d.
- Yurun Chen et al. SafePred: A predictive guardrail for computer-using agents via world models. *arXiv preprint arXiv:2602.01725*, 2026b.
- Zehui Chen et al. Agent-flan: Designing data and methods of effective agent tuning for large language models. *arXiv preprint arXiv:2403.12881*, 2024e.
- Zixiang Chen et al. Spin: Self-play fine-tuning converts weak language models to strong. 2024f.
- Jeff Clune. Position: Foundation models for open-endedness. 2024.

- Thomas Coste, Usman Anwar, Robert Kirk, and David Krueger. Reward model ensembles help mitigate overoptimization. 2024.
- DeepSeek-AI. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2025.
- Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and Yu Su. Mind2web: Towards a generalist agent for the web. *Advances in Neural Information Processing Systems*, 2024.
- Thomas G Dietterich. Hierarchical reinforcement learning with the maxq value function decomposition. In *Journal of Artificial Intelligence Research*, volume 13, pages 227–303, 2000.
- Shihan Dou et al. Stepcoder: Improve code generation with reinforcement learning from compiler feedback. *arXiv preprint arXiv:2402.01391*, 2024.
- Alexandre Drouin et al. Browsergym: A gym environment for web task automation. *ServiceNow Research Technical Report*, 2024a.
- Alexandre Drouin et al. Workarena: How capable are web agents at solving common knowledge work tasks? *arXiv preprint arXiv:2403.07718*, 2024b.
- Jiaqi Du et al. Position paper: Agent ai towards a holistic intelligence. *arXiv preprint arXiv:2403.00833*, 2024a.
- Pengfei Du. Memory for autonomous LLM agents: Mechanisms, evaluation, and emerging frontiers. *arXiv preprint arXiv:2603.07670*, 2026.
- Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. Improving factuality and reasoning in language models through multiagent debate. *Proceedings of ICML*, 2024b.
- Yilun Du, Mengjiao Yang, Bo Dai, Hanjun Dai, Ofir Nachum, Josh Tenenbaum, Dale Schuurmans, and Pieter Abbeel. Learning universal policies via text-guided video generation. *Advances in Neural Information Processing Systems*, 2024c.
- Yuqing Du, Olivia Watkins, Zihan Wang, Cedric Colas, Trevor Darrell, Pieter Abbeel, Abhishek Gopinath, and Joshua B Tenenbaum. Guiding pretraining in reinforcement learning with large language models. *Proceedings of ICML*, 2023.
- Zane Durante, Qiuyuan Huang, Naoki Wake, Ran Gong, Jae Sung Park, Bidipta Gong, Jia Deng, Yusuke Gu, Yonatan Cho, Hao Fang, et al. An interactive agent foundation model. *Proceedings of ICML*, 2024.
- 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. In *Advances in Neural Information Processing Systems*, 2022.
- Martin Fang, Martin Klissarov, Alexandre Piché, Wang Hanjing, Evan Racah, Jack Parker-Holder, Danijar Raber, and Doina Precup. Motif: Intrinsic motivation from artificial intelligence feedback. *Advances in Neural Information Processing Systems*, 2024.
- Roya Firoozi, Johnathan Tucker, Stephen Tian, et al. Foundation models in robotics: Applications, challenges, and the future. *International Journal of Robotics Research*, 2024.
- Jiaxuan Gao et al. From self-evolving synthetic data to verifiable-reward RL for tool-using agents. *arXiv preprint arXiv:2601.22607*, 2026.
- Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization. 2023.
- Daya Guo, Dejian Yang, He Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

- Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xiangliang Zhang. Large language model based multi-agents: A survey of progress and challenges. *Proceedings of IJCAI*, 2024.
- Izzeddin Gur et al. A real-world webagent with planning, long context understanding, and program synthesis. *Proceedings of ICLR*, 2024.
- David Ha and Jurgen Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018.
- Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models. *arXiv preprint arXiv:2301.04104*, 2023.
- Hao Han et al. SWE-TRACE: Optimizing long-horizon SWE agents via rubric prms and heuristic test-time scaling. *arXiv preprint arXiv:2604.14820*, 2026.
- Nicklas Hansen, Hao Su, and Xiaolong Wang. Td-mpc2: Scalable, robust world models for continuous control. *Proceedings of ICML*, 2024.
- Shibo Hao, Yi Gu, Haodi Ma, Joshua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. Reasoning with language model is planning with world model. *arXiv preprint arXiv:2305.14992*, 2023.
- Alex Havrilla, Yuqing Du, Sharath Chandra Raparthy, Paul Christiano, et al. Teaching large language models to reason with reinforcement learning. *arXiv preprint arXiv:2403.04642*, 2024.
- Mikael Henaff, Scott Fujimoto, Michael Matthews, and Michael Rabbat. Scalable option learning in high-throughput environments. *arXiv preprint arXiv:2509.00338*, 2025.
- Samuel Holt et al. L2mac: Large language model automatic computer for extensive code generation. 2024.
- Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Ceyao Zhang, Jinlin Wang, et al. Metagpt: Meta programming for a multi-agent collaborative framework. *Proceedings of ICLR*, 2024.
- Arian Hosseini, Xingdi Yuan, Nikolay Malkin, Aaron Courville, Alessandro Sordoni, and Rishabh Agarwal. V-star: Training verifiers for self-taught reasoners. *Advances in Neural Information Processing Systems*, 2024.
- Arian Hosseini et al. Generative verifiers: Reward modeling as next-token prediction. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2025.
- Jack Hu, Jake Bruce, et al. Genie 2: A large-scale foundation world model. *arXiv preprint (DeepMind)*, 2024a.
- Jian Hu et al. Reinforce++: A simple and efficient approach for aligning large language models. *arXiv preprint arXiv:2501.03262*, 2025.
- Mengkang Hu et al. Tree-planner: Efficient close-loop task planning with large language models. *Proceedings of ICLR*, 2024b.
- Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. *Proceedings of ICML*, 2022.
- Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, et al. Inner monologue: Embodied reasoning through planning with language models. *Proceedings of CoRL*, 2023.
- Xu Huang, Weiwen Liu, Xiaolong Chen, Xingshan Wang, Haitao Wang, Defu Lian, Yasheng Wang, Ruiming Dong, and Feng Cheng. Understanding the planning of llm agents: A survey. *arXiv preprint arXiv:2402.02716*, 2024.
- Naman Jain et al. Livecodebench: Holistic and contamination-free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.

- Michael Janner, Qiyang Li, and Sergey Levine. Offline reinforcement learning as one big sequence modeling problem. *Advances in Neural Information Processing Systems*, 2021.
- Michael Janner, Yilun Du, Joshua B Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior synthesis. *Proceedings of ICML*, 2022.
- Jinghan Jia et al. Open reasoner zero: An open source approach to scaling up reinforcement learning on the base model. *arXiv preprint*, 2025.
- Dongming Jiang et al. MAGMA: A multi-graph based agentic memory architecture for AI agents. *arXiv preprint arXiv:2601.03236*, 2026.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. Swe-bench: Can language models resolve real-world github issues? *Proceedings of ICLR*, 2024.
- Leslie Pack Kaelbling, Michael L Littman, and Anthony R Cassandra. Planning and acting in partially observable stochastic domains. *Artificial Intelligence*, 101:99–134, 1998.
- Sayash Kapoor et al. Ai agents that matter. *arXiv preprint arXiv:2407.01502*, 2024.
- Timo Kaufmann, Paul Weng, et al. A survey of reinforcement learning from human feedback. *arXiv preprint*, 2024.
- Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. Decomposed prompting: A modular approach for solving complex tasks. *Proceedings of ICLR*, 2023.
- Joongwon Kim et al. Scaling test-time compute for agentic coding. *arXiv preprint arXiv:2604.16529*, 2026.
- Kimi Team. Kimi k1.5: Scaling reinforcement learning with llms. *arXiv preprint arXiv:2501.12599*, 2025.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114:3521–3526, 2017.
- Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Chong Lim, Po-Yu Huang, Graham Neubig, Shuyan Zhou, Ruslan Salakhutdinov, and Daniel Fried. Visualwebarena: Evaluating multimodal agents on realistic visual web tasks. *Proceedings of ACL*, 2024.
- Yann LeCun. A path towards autonomous machine intelligence. *Open Review*, 2022.
- Jonathan Lee et al. Supervised pretraining can learn in-context reinforcement learning. 2024.
- Kuang-Huei Lee et al. Multi-game decision transformers. 2022.
- Nicholas Lee et al. Agentic test-time scaling for WebAgents. *arXiv preprint arXiv:2602.12276*, 2026.
- Andrew Levy, George Konidaris, Robert Platt, and Kate Saenko. Learning multi-level hierarchies with hindsight. *Proceedings of ICLR*, 2019.
- Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. Camel: Communicative agents for “mind” exploration of large language model society. *Advances in Neural Information Processing Systems*, 2023.
- Shuang Li et al. Hierarchical reinforcement learning with foundation models for long-horizon tasks. In *Proceedings of ICML*, 2024a.
- Wentao Li et al. A survey on llm-based multi-agent systems for software engineering. *arXiv preprint*, 2025a.
- Xiao Li et al. Hierarchical diffusion policy for multi-task robotic manipulation. *Advances in Neural Information Processing Systems*, 2024b.

- Xiaoxi Li et al. Search-o1: Agentic search-enhanced large reasoning models. *arXiv preprint arXiv:2501.05366*, 2025b.
- Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. Code as policies: Language model programs for embodied control. *Proceedings of ICRA*, 2023.
- 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*, 2024.
- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. 2024.
- Shaobin Ling et al. ELHPlan: Efficient long-horizon task planning for multi-agent collaboration. *arXiv preprint arXiv:2509.24230*, 2025.
- Alex Liu et al. Dense reward for free in reinforcement learning from human feedback. *Proceedings of ICML*, 2024a.
- Canbin Liu et al. Dr. grpo: Removing estimation bias from group relative policy optimization. *arXiv preprint*, 2025a.
- Hao Liu, Matei Zaharia, and Pieter Abbeel. World model on million-length video and language with ringattention. *arXiv preprint arXiv:2402.08268*, 2024b.
- Jiawei Liu et al. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation. 2024c.
- Minghao Liu et al. Language-conditioned imitation learning for long-horizon manipulation. *Proceedings of ICRA*, 2024d.
- Weiwen Liu et al. Toolace: Winning the points of llm function calling. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2025b.
- Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, et al. Agentbench: Evaluating llms as agents. *Proceedings of ICLR*, 2024e.
- Xiao Liu et al. Visualagentbench: Towards large multimodal models as visual foundation agents. *arXiv preprint arXiv:2408.06327*, 2024f.
- Xuanzhang Liu et al. CoDA: A context-decoupled hierarchical agent with reinforcement learning. *arXiv preprint arXiv:2512.12716*, 2025c.
- Yuxin Liu et al. Ragen: Training agents by reinforcing reasoning. *arXiv preprint*, 2025d.
- Cong Lu et al. Intelligent go-explore: Standing on the shoulders of giant foundation models. *arXiv preprint arXiv:2405.15143*, 2024.
- Liangchen Luo et al. Improve mathematical reasoning in language models by automated process supervision. *arXiv preprint arXiv:2406.06592*, 2024.
- Chang Ma et al. Agentboard: An analytical evaluation board of multi-turn llm agents. *arXiv preprint arXiv:2401.13178*, 2024a.
- Ye Cheng Jason Ma, William Liang, Guanzhi Wang, De-An Huang, Osbert Bastani, Dinesh Jayaraman, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Eureka: Human-level reward design via coding large language models. *Proceedings of ICLR*, 2024b.
- Aman Madaan et al. Self-refine: Iterative refinement with self-feedback. 2024.
- Gregoire Mialon, Clementine Fourrier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas Scialom. Gaia: A benchmark for general ai assistants. *arXiv preprint arXiv:2311.12983*, 2023.

Vincent Micheli, Eloi Alonso, and Francois Fleuret. Transformers are sample-efficient world models. 2023.

Microsoft Research. rstar-math: Small llms can master math reasoning with self-evolved deep thinking. *arXiv preprint arXiv:2501.04519*, 2025.

Ted Moskovitz, Aldo Shenfeld, et al. Confronting reward model overoptimization with constrained rlhf. *Proceedings of ICML*, 2024.

Niklas Muennighoff, Zijian Yang, et al. s1: Simple test-time scaling. *arXiv preprint arXiv:2501.19393*, 2025.

Ofir Nachum, Shixiang Shane Gu, Honglak Lee, and Sergey Levine. Data-efficient hierarchical reinforcement learning. In *Advances in Neural Information Processing Systems*, 2018.

Soroush Nasiriany, Fei Xia, Wenhao Yu, Ted Xiao, Jacky Liang, Ishita Dasgupta, Andy Zeng, et al. Pivot: Iterative visual prompting elicits actionable knowledge for vlms. *Proceedings of ICML*, 2024.

Dana Nau, Tsz-Chiu Au, Okhtay Ilghami, Ugur Kuter, J. William Murdock, Dan Wu, and Fusun Yaman. SHOP2: An HTN planning system. *Journal of Artificial Intelligence Research*, 20:379–404, 2003.

Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations: Theory and application to reward shaping. *Proceedings of ICML*, 1999.

NovaSky Team, UC Berkeley. Sky-t1: Scaling reasoning with reinforcement learning at minimal cost. *Technical Report*, 2025.

OpenAI. Openai o3 system card. *OpenAI Technical Report*, 2025.

Others. FoVer: Generalizable process reward models via formally verified training data. *arXiv preprint arXiv:2505.15960*, 2025.

Others. Subgoal graph-augmented planning for LLM-guided open-world reinforcement learning. In *Proc. AAMAS*, 2026a.

Others. Reducing credit assignment variance via counterfactual reasoning paths. *arXiv preprint arXiv:2605.16302*, 2026b.

Others. Mobile-agent-v3.5: Multi-platform fundamental GUI agents. *arXiv preprint arXiv:2602.16855*, 2026c.

Others. ICA: Information-aware credit assignment for long-horizon information-seeking agents. *arXiv preprint arXiv:2602.10863*, 2026d.

Others. Occupancy reward shaping: Improving credit assignment for offline goal-conditioned rl. *arXiv preprint arXiv:2604.20627*, 2026e.

Others. PiCA: Pivot-based credit assignment for search agentic reinforcement learning. *arXiv preprint arXiv:2605.09287*, 2026f.

Others. Self-induced outcome potential: Turn-level credit assignment for agents without verifiers. *arXiv preprint arXiv:2605.04984*, 2026g.

Others. Beyond outcome verification: Verifiable process reward models for structured reasoning. *arXiv preprint arXiv:2601.17223*, 2026h.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 2022.

Charles Packer, Vivian Fang, Shishir G Patil, Kevin Lin, Sarah Wooders, and Joseph E Gonzalez. Memgpt: Towards llms as operating systems. *Proceedings of ICLR*, 2024.

- Haojie Pan et al. Kwaiagents: Generalized information-seeking agent system with large language models. *arXiv preprint arXiv:2312.04889*, 2024.
- Joon Sung Park, Joseph C O’Brien, Carrie J Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. *Proceedings of UIST*, 2023.
- Seohong Park, Sergey Levine, et al. Ogbench: Benchmarking offline goal-conditioned rl. *arXiv preprint*, 2024.
- Jack Parker-Holder et al. Evolving curricula with regret-based environment design. *Journal of Machine Learning Research*, 2024.
- Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E Gonzalez. Gorilla: Large language model connected with massive apis. *Proceedings of ICML*, 2024.
- Jiangweizhi Peng et al. HiPER: Hierarchical rl with explicit credit assignment for large language model agents. *arXiv preprint arXiv:2602.16165*, 2026.
- Wasu Top Piriyakulkij, Wolfgang Lehrach, Kevin Ellis, and Kevin Murphy. Joint learning of hierarchical neural options and abstract world model. *arXiv preprint arXiv:2602.02799*, 2026.
- Archiki Prasad et al. Adapt: As-needed decomposition and planning with language models. 2024.
- Pranav Putta et al. Agentq: Advanced reasoning and learning for autonomous ai agents. *arXiv preprint arXiv:2408.07199*, 2024.
- Zhenting Qi et al. rstar: Mutual reinforcement of reasoning and self-play in small language models. *arXiv preprint arXiv:2408.06195*, 2024.
- Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, et al. Chatdev: Communicative agents for software development. *Proceedings of ACL*, 2024.
- Tianhao Qian. PAC-MCTS: Bias-aware pruning for robust LLM-guided search and planning. *arXiv preprint arXiv:2604.14345*, 2026.
- Shuofei Qiao et al. Autoact: Automatic agent learning from scratch for qa via self-planning. *arXiv preprint arXiv:2401.05268*, 2024.
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. Toolllm: Facilitating large language models to master 16000+ real-world apis. *Proceedings of ICLR*, 2024.
- Zhisong Qiu et al. Rewarding the scientific process: Process-level reward modeling for agentic data analysis. *arXiv preprint arXiv:2604.24198*, 2026.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 2023.
- Christopher Rawles et al. Androidworld: A dynamic benchmarking environment for autonomous agents. *arXiv preprint arXiv:2405.14573*, 2024.
- Scott Reed, Konrad Zolna, Emilio Parisotto, et al. A generalist agent. *Transactions on Machine Learning Research*, 2022.
- Jan Robine, Marc Höftmann, Tobias Uelwer, and Stefan Harmeling. Transformer-based world models are happy with 100k interactions. *Proceedings of ICLR*, 2024.
- Hyunyoung Ryu et al. Curricullm: Automatic task curricula design for learning complex robot control. In *Proceedings of IEEE International Conference on Robotics and Automation (ICRA)*, 2025.

- Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 2024.
- Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, et al. Mastering Atari, Go, Chess and Shogi by planning with a learned model. *Nature*, 588:604–609, 2020.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y.K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- Shashank Sharma, Janina Hoffmann, and Vinay Namboodiri. Multi-resolution skills for HRL agents. *arXiv preprint arXiv:2505.21410*, 2025.
- Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugginggpt: Solving ai tasks with chatgpt and its friends in hugging face. *Advances in Neural Information Processing Systems*, 2024.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. In *Advances in Neural Information Processing Systems*, 2023.
- Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Cote, Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. Alfworld: Aligning text and embodied environments for interactive learning. In *Proceedings of ICLR*, 2021.
- Shuzheng Si et al. A goal without a plan is just a wish: Efficient and effective global planner training for long-horizon agent tasks. *arXiv preprint arXiv:2510.05608*, 2025.
- Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox, Jesse Thomason, and Animesh Garg. Progprompt: Generating situated robot task plans using large language models. *Proceedings of ICRA*, 2023.
- Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2025.
- Paloma Sodhi et al. Step-by-step reasoning for long-horizon decision making. *Proceedings of CoRL*, 2024.
- 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 ICCV*, 2023.
- Yifan Song, Weimin Xiong, Dawei Zhu, Wenhao Wu, Han Qian, Mingbo Song, Hailiang Huang, Cheng Li, Ke Wang, Rong Yao, et al. Restgpt: Connecting large language models with real-world restful apis. *Proceedings of ACL*, 2024a.
- Yifan Song et al. Trial and error: Exploration-based trajectory optimization for llm agents. *Proceedings of ACL*, 2024b.
- Theodore R Sumers, Shunyu Yao, Karthik Narasimhan, and Thomas L Griffiths. Cognitive architectures for language agents. *Transactions on Machine Learning Research*, 2024.
- Haotian Sun, Yuchen Zhuang, Lingkai Kong, Bo Dai, and Chao Zhang. Adaplaner: Adaptive planning from feedback with language models. *Advances in Neural Information Processing Systems*, 2024.
- Richard S Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning. In *Artificial Intelligence*, volume 112, pages 181–211, 1999.

- Weihao Tan et al. Cradle: Empowering foundation agents towards general computer control. *arXiv preprint arXiv:2403.03186*, 2024.
- Adaptive Agent Team, Jakob Bauer, Kate Baumli, Satinder Baveja, et al. Human-timescale adaptation in an open-ended task space. *Proceedings of ICML*, 2023.
- SIMA Team et al. A generalist ai agent for 3d virtual environments. *arXiv preprint arXiv:2404.10179*, 2024.
- Trieu H Trinh, Yuhuai Wu, Quoc V Le, He He, and Thang Luong. Solving olympiad geometry without human demonstrations. *Nature*, 625:476–482, 2024.
- Leonid Ugadiarov, Yuri Kuratov, Aleksandr Panov, and Alexey Skrynnik. Revisiting tree search for LLMs: Gumbel and sequential halving for budget-scalable reasoning. *arXiv preprint arXiv:2603.21162*, 2026.
- Alexander Sasha Vezhnevets, Simon Osindero, Tom Schaul, Nicolas Heess, Max Jaderberg, David Silver, and Koray Kavukcuoglu. Feudal networks for hierarchical reinforcement learning. *Proceedings of ICML*, 2017.
- Yongchao Wan, Jianbo Li, Hao Yu, et al. Alphallm: Training llms with mcts. *arXiv preprint*, 2024.
- 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. *Transactions on Machine Learning Research*, 2024a.
- Hao Wang et al. Do androids dream of breaking the game? auditing AI agent benchmarks with BenchJack. *arXiv preprint arXiv:2605.12673*, 2026a.
- Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language models. *Proceedings of ACL*, 2023.
- Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 2024b.
- Peiyi Wang, Lei Li, Zhihong Shao, R.X. Xu, Damai Dai, Yifei Li, Deli Chen, Y. Wu, and Zhifang Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations. In *Proceedings of ACL*, 2024c.
- Ruiyi Wang and Prithviraj Ammanabrolu. A practitioner’s guide to multi-turn agentic reinforcement learning. *arXiv preprint arXiv:2510.01132*, 2025.
- Ruoyao Wang et al. Scienceworld: Is your agent smarter than a 5th grader? In *Proceedings of EMNLP*, 2022.
- Taiyi Wang et al. A subgoal-driven framework for improving long-horizon LLM agents. *arXiv preprint arXiv:2603.19685*, 2026b.
- Weizhi Wang, Li Dong, Hao Cheng, Xiaodong Liu, Xifeng Yan, Jianfeng Gao, and Furu Wei. Augmenting language models with long-term memory. *Advances in Neural Information Processing Systems*, 2024d.
- Xiangming Wang et al. Agent smith: A single image can jailbreak one million multimodal llm agents exponentially fast. *Proceedings of ICML*, 2024e.
- Xiaoqiang Wang et al. Mem- π : Adaptive memory through learning when and what to generate. *arXiv preprint arXiv:2605.21463*, 2026c.
- Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang, Yunzhu Li, Hao Peng, and Heng Ji. Executable code actions elicit better llm agents. *Proceedings of ICML*, 2024f.

- Xingyao Wang, Zihan Shi, Jiayi Zhang, et al. Mint: Evaluating llms in multi-turn interaction with tools and language feedback. *Proceedings of ICLR*, 2024g.
- Xingyao Wang et al. Opendevin: An open platform for ai software developers as generalist agents. *arXiv preprint arXiv:2407.16741*, 2024h.
- Zihao Wang, Shaofei Cai, Guanqi Chen, Anji Liu, Xiaojian Ma, and Yitao Liang. Jarvis-1: Open-world multi-task agents with memory-augmented multimodal language models. *arXiv preprint arXiv:2311.05997*, 2024i.
- Zihao Wang, Shaofei Cai, Anji Liu, Xiaojian Ma, and Yitao Liang. Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents. In *Advances in Neural Information Processing Systems*, 2024j.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 2022.
- Hjalmar Wijk et al. Re-bench: Evaluating frontier ai r&d capabilities of language model agents against human experts. *METR Technical Report*, 2024.
- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, et al. Autogen: Enabling next-gen llm applications via multi-agent conversation. *arXiv preprint arXiv:2308.08155*, 2023a.
- Yiran Wu et al. Stateflow: Enhancing llm task-solving through state-driven workflows. *arXiv preprint arXiv:2403.11322*, 2024a.
- Yue Wu et al. Smartplay: A benchmark for llms as intelligent agents. 2024b.
- Yue Wu et al. Self-play preference optimization for language model alignment. 2024c.
- Yue Wu et al. Spring: Studying papers and reasoning to play games. 2024d.
- Yueh-Hua Wu et al. Elastic decision transformer. 2023b.
- Zhiyong Wu et al. Os-atlas: A foundation action model for generalist gui agents. *arXiv preprint arXiv:2410.23218*, 2024e.
- Zhiyong Wu et al. Os-copilot: Towards generalist computer agents with self-improvement. *arXiv preprint arXiv:2402.07456*, 2024f.
- Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, et al. The rise and potential of large language model based agents: A survey. *arXiv preprint arXiv:2309.07864*, 2023.
- Zhiheng Xi et al. Agentgym: Evolving large language model-based agents across diverse environments. *arXiv preprint arXiv:2406.04151*, 2024.
- Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. Agentless: Demystifying llm-based software engineering agents. *arXiv preprint arXiv:2407.01489*, 2024.
- Bowen Xiang et al. Towards longer long-context language models: An ntk perspective. *arXiv preprint arXiv:2311.12351*, 2024a.
- Jiannan Xiang, Tianhua Tao, Yi Gu, Tianmin Shu, Zirui Wang, Zichao Yang, and Zhiting Hu. Language models meet world models: Embodied experiences enhance language models. *Advances in Neural Information Processing Systems*, 2024b.

- Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing Hua, Zhoujun Cheng, Dongchan Shi, Joel Lu, et al. Osworld: Benchmarking multimodal agents for open-ended tasks in real computer environments. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS), Datasets and Benchmarks Track*, 2024a.
- Tianbao Xie et al. Text2reward: Automated dense reward function generation for reinforcement learning. 2024b.
- Xin Xie et al. SGA-MCTS: Decoupling planning from execution via training-free atomic experience retrieval. *arXiv preprint arXiv:2604.14712*, 2026.
- Zora Zhiruo Xu et al. Agent workflow memory. *arXiv preprint arXiv:2409.07429*, 2024.
- John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Liber, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. Swe-agent: Agent-computer interfaces enable automated software engineering. *arXiv preprint arXiv:2405.15793*, 2024a.
- John Yang et al. Intercode: Standardizing and benchmarking interactive coding with execution feedback. 2024b.
- Mengjiao Yang, Yilun Du, Kamyar Ghasemipour, Jonathan Tompson, Dale Schuurmans, and Pieter Abbeel. Learning interactive real-world simulators. *Proceedings of ICLR*, 2024c.
- Sherry Yang, Ofir Nachum, Yilun Du, Jason Wei, Pieter Abbeel, and Dale Schuurmans. Foundation models for decision making: Problems, methods, and opportunities. *Journal of Artificial Intelligence Research*, 2024d.
- Yiheng Yang et al. Agenttrek: Agent trajectory synthesis via guiding replay with web tutorials. *arXiv preprint arXiv:2412.09605*, 2024e.
- Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable real-world web interaction with grounded language agents. In *Advances in Neural Information Processing Systems*, 2022.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2023.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 2024a.
- Shunyu Yao et al. tau-bench: A benchmark for tool-agent-user interaction in real-world domains. *arXiv preprint arXiv:2406.12045*, 2024b.
- Da Yin et al. Agent lumos: Unified and modular training for open-source language agents. 2024.
- Ori Yoran et al. Assistantbench: Can web agents solve realistic and time-consuming tasks? *arXiv preprint arXiv:2407.15711*, 2024.
- Tongxin Yuan et al. R-judge: Benchmarking safety risk awareness for llm agents. *arXiv preprint arXiv:2401.10019*, 2024.
- Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah D Goodman. Star: Bootstrapping reasoning with reasoning. 2022.
- Eric Zelikman, Georges Harik, Yifei Shao, Varuna Jayasiri, Nick Haber, and Percy Liang. Quiet-star: Language models can teach themselves to think before speaking. *arXiv preprint arXiv:2403.09629*, 2024.
- Yiping Zeng et al. Simplerl: Simple yet effective reinforcement learning for llm reasoning. *arXiv preprint*, 2025.

- Chenchen Zhang. From reasoning to agentic: Credit assignment in rl for large language models. *arXiv preprint arXiv:2604.09459*, 2026.
- Chi Zhang et al. Appagent: Multimodal agents as smartphone users. *arXiv preprint arXiv:2312.13771*, 2024a.
- Dang Zhang et al. Gui agents: A survey. *arXiv preprint arXiv:2504.02900*, 2025a.
- Di Zhang et al. Accessing gpt-4 level mathematical olympiad solutions via monte carlo tree self-refine. *arXiv preprint arXiv:2406.07394*, 2024b.
- Jianguo Zhang et al. xlam: A family of large action models to empower ai agent systems. In *Proceedings of Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, 2025b.
- Seohong Zhang, Dibya Park, and Sergey Levine. Hiql: Offline goal-conditioned rl with latent states as actions. *Advances in Neural Information Processing Systems*, 2024c.
- Yuntong Zhang, Haifeng Ruan, Zhiyu Fan, and Abhik Roychoudhury. Autocoderover: Autonomous program improvement. *Proceedings of ISSTA*, 2024d.
- Yuqing Zhang et al. Hierarchical auto-curriculum for open-ended goals with large language models. *arXiv preprint*, 2024e.
- Zeyu Zhang, Xiaohe Gao, Wenhao Zhang, et al. A survey on the memory mechanism of large language model based agents. *arXiv preprint arXiv:2404.13501*, 2024f.
- Zihao Zhang, Shaofei Cai, Zhancun Ma, et al. Omnijarvis: Unified vision-language-action tokenization enables open-world instruction following agents. *arXiv preprint arXiv:2407.00114*, 2024g.
- Andrew Zhao et al. Absolute zero: Reinforced self-play reasoning with zero data. *arXiv preprint arXiv:2505.03335*, 2025.
- Enshen Zhao, Chenlin Tao, Yunfan Zhang, et al. Minedreamer: Learning to follow instructions via chain-of-imagination and video diffusion. *arXiv preprint*, 2024a.
- Jun Zhao et al. Longagent: Scaling language models to 128k context through multi-agent collaboration. *arXiv preprint arXiv:2402.04268*, 2024b.
- Qinqing Zheng, Amy Zhang, and Aditya Grover. Online decision transformer. *Proceedings of ICML*, 2022.
- Tianyu Zheng et al. Opencodeinterpreter: Integrating code generation with execution and refinement. *Proceedings of ACL*, 2024.
- Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and Yanlin Wang. Memorybank: Enhancing large language models with long-term memory. *Proceedings of AAAI*, 2024.
- Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang. Language agent tree search unifies reasoning, acting, and planning in language models. *Proceedings of ICML*, 2024a.
- Denny Zhou, Nathanael Schaarschmidt, Mehran Sahami, Xinyun Liu, and William Wang. Least-to-most prompting enables complex reasoning in large language models. *Proceedings of ICLR*, 2023a.
- Denny Zhou et al. Least-to-most prompting enables complex reasoning in large language models. 2023b.
- Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Yonatan Bisk, Daniel Fried, Uri Alon, et al. Webarena: A realistic web environment for building autonomous agents. *Proceedings of ICLR*, 2024b.
- Xizhou Zhu, Yuntao Chen, Hao Tian, Chenxin Chen, Jiashuo Zhao, et al. Ghost in the minecraft: Generally capable agents for open-world environments via large language models with text-based knowledge and memory. *arXiv preprint arXiv:2305.17144*, 2023.

Yuqi Zhu et al. Knowagent: Knowledge-augmented planning for llm-based agents. *arXiv preprint arXiv:2403.03101*, 2024.