

---

# Learning When to Plan: Efficiently Allocating Test-Time Compute for LLM Agents

---

**Davide Paglieri<sup>1\*</sup>, Bartłomiej Cupiał<sup>1,2,6\*</sup>, Jonathan Cook<sup>3</sup>,  
Ulyana Piterberg<sup>4</sup>, Jens Tuyls<sup>5</sup>, Edward Grefenstette<sup>1</sup>,  
Jakob Nicolaus Foerster<sup>3</sup>, Jack Parker-Holder<sup>1</sup>, Tim Rocktäschel<sup>1</sup>**  
<sup>1</sup>AI Centre, University College London, <sup>2</sup>IDEAS NCBR, <sup>3</sup>University of Oxford,  
<sup>4</sup>New York University, <sup>5</sup>Princeton University, <sup>6</sup>University of Warsaw,  
d.paglieri@cs.ucl.ac.uk

## Abstract

Training large language models (LLMs) to reason via reinforcement learning (RL) significantly improves their problem-solving capabilities. In agentic settings, existing methods like ReAct prompt LLMs to explicitly plan before every action; however, we demonstrate that always planning is computationally expensive and degrades performance on long-horizon tasks, while never planning further limits performance. To address this, we introduce a conceptual framework formalizing dynamic planning for LLM agents, enabling them to flexibly decide when to allocate test-time compute for planning. We propose a simple two-stage training pipeline: (1) supervised fine-tuning on diverse synthetic data to prime models for dynamic planning, and (2) RL to refine this capability in long-horizon environments. Experiments on the Crafter environment show that dynamic planning agents trained with this approach are more sample-efficient and consistently achieve more complex objectives. Additionally, we demonstrate that these agents can be effectively steered by human-written plans, surpassing their independent capabilities. To our knowledge, this work is the first to explore training LLM agents for dynamic test-time compute allocation in sequential decision-making tasks, paving the way for more efficient, adaptive, and controllable agentic systems.

## 1 Introduction

A key insight from recent work on LLM reasoning is the role of *test-time compute* — the ability to allocate additional computational resources to more difficult problems [Snell et al., 2025, Guo et al., 2025]. For humans, difficult tasks often require deliberate thinking. Similarly, LLMs benefit from dedicating extra processing to explicitly reason through steps via chain-of-thought [Wei et al., 2022]. In settings like math problem-solving and code generation, reasoning can enable models to explore possible answers before settling on a response in a manner akin to search [Xiang et al., 2025, Prystawski et al., 2023, Ruis et al., 2025]. Reasoning LLMs trained to effectively use additional test-time compute on single-step tasks have also been shown to make extremely effective zero-shot agents [Yao et al., 2023b]. However, a critical open question remains: *Can we further improve an LLM’s ability to effectively allocate test-time compute on sequential decision-making tasks?* On the challenging agentic benchmark BALROG [Paglieri et al., 2025a], reasoning models have thus far only shown marginal gains over models immediately producing the next action [Paglieri et al., 2025b].

In agentic tasks, planning naturally emerges as a multi-step analogue to single-step chain-of-thought reasoning. Rather than committing immediately to a single next action, an agent can invest computational resources to better understand the current state and anticipate the outcome of future actions

---

\*Equal contribution.

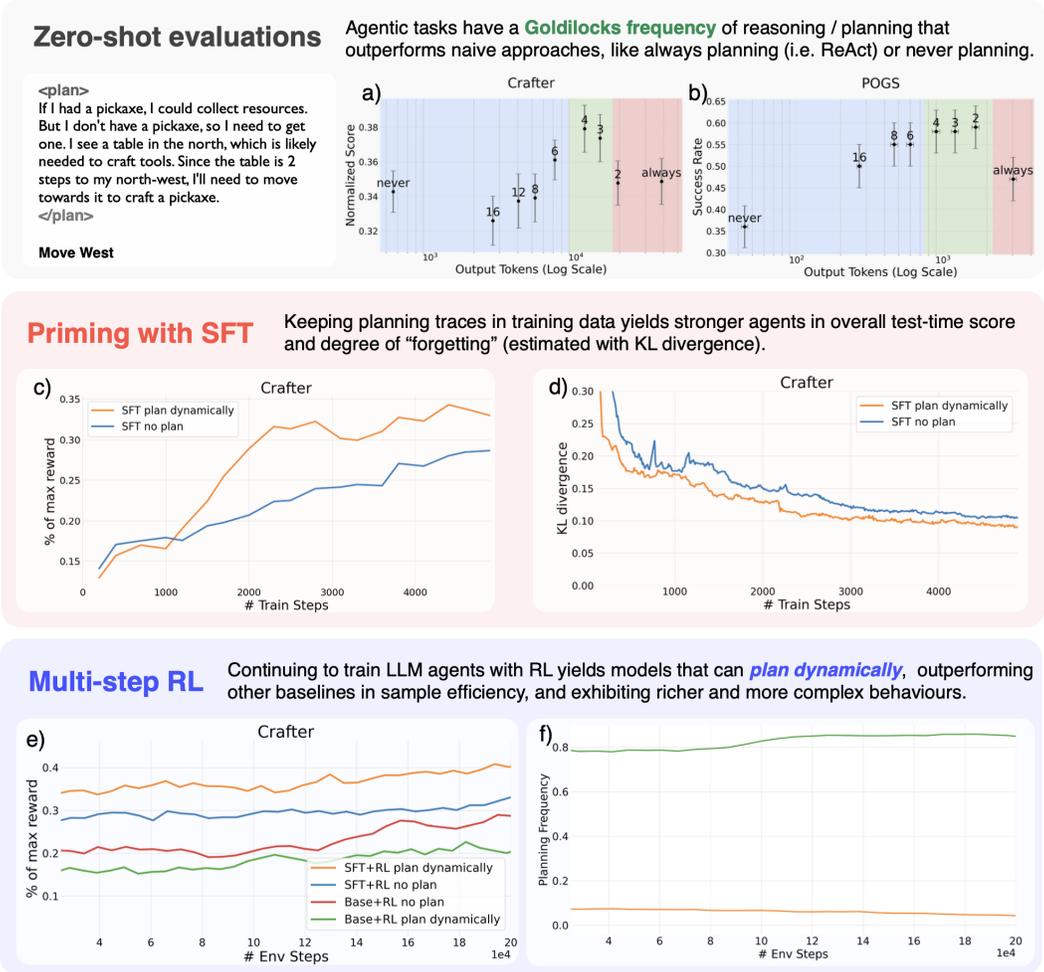


Figure 1: **Dynamic planning strategies across environments and training stages.** (a-b) Zero-shot results showing optimal "Goldilocks" planning frequency in Crafter and POGS (100 seeds, bars=standard-error). (c-d) SFT results demonstrating planning agents' improved performance with lower KL divergence from base model. (e-f) RL results where SFT-primed planning agents are more sample efficient than non-planning baselines and more consistently reach complex achievements.

sequences. Plans can serve as a guide for subsequent actions and improve the strategic coherence of future behaviour. However, introducing explicit planning presents its own critical challenge: deciding precisely when an agent should plan. This decision must carefully balance the performance improvements gained from more informed decision-making against the computational cost and additional variance in behaviour incurred by frequent replanning.

To formalise this problem, we develop a framework for modelling the cost-benefit trade-offs of planning in partially-observable environments. According to our framework, agents should allocate test-time compute for planning only when anticipated improvements in policy performance outweigh the associated computational costs and any instability or noise induced by excessive replanning.

We experimentally investigate these concepts in two distinct environments: Partially-Observable Graph Search (POGS), a synthetic environment that we design to systematically evaluate planning abilities, and Crafter, a Minecraft-inspired grid-world environment [Hafner, 2022]. Inspired by recent work showing that the presence of key inductive biases in training data are necessary for effective self-improvement [Gandhi et al., 2025], we develop a two-stage approach: first priming models with diverse planning behaviours through supervised fine-tuning (SFT), then applying RL. Using this approach, we successfully train agents that learn to plan strategically, execute their plans, and replan only when necessary, outperforming non-planning baselines trained via an equivalent two-stage pipeline. Furthermore, following the RL stage, agents that are trained to produce and follow their

own plans can be effectively steered by plans produced by humans to achieve performance that the agents cannot reach alone. In summary, our experiments yield four key insights:

1. Each task has a "Goldilocks" frequency for planning that clearly outperforms naive strategies of always planning or never planning.
2. SFT priming demonstrates that including explicit natural language plans in training data significantly improves imitation learning compared to using identical action sequences without plans.
3. RL fine-tuning after SFT priming yields planning agents that further outperform baselines in sample efficiency, and learn to plan, execute plans, and replan when necessary.
4. Planning agents can be collaboratively steered by humans that produce plans for them. This is only the case following RL and is not achieved by SFT priming alone.

Our work provides clear evidence that dynamic planning facilitates effective allocation of test-time compute in sequential decision making environments, showing that LLM agents can be trained to use additional computational resources intelligently. The ability to steer such planning agents—now capable enough to complete Crafter by collecting diamond under human guidance—marks a significant step towards safer and more collaborative LLM agents. Together, these findings suggest a promising path towards more capable, efficient, interpretable, and steerable agentic systems.

## 2 Related Work

**Classical Planning Methods** Historically, much progress in sequential decision making has involved systems that explicitly look ahead before acting. Monte Carlo Tree Search (MCTS)[Coulom, 2006], combined with deep neural networks, has driven landmark systems like AlphaGo and MuZero[Silver et al., 2017b,a, Schrittwieser et al., 2020]. Model Predictive Control (MPC) iteratively plans short horizons, adapting to new observations [Mayne et al., 2000]. Similarly, model-based RL methods, such as World Models, PlaNet, and Dreamer, use imagination rollouts in latent space for effective planning and learning [Ha and Schmidhuber, 2018, Hafner et al., 2019, 2020]. Collectively, these approaches underscore the strength of explicit planning, particularly when accurate internal or environmental models are available.

**LLM Reasoning and Planning** Large language models have demonstrated significant reasoning capabilities, particularly through techniques like chain-of-thought (CoT) prompting [Wei et al., 2022]. ReAct extends CoT into sequential settings, explicitly prompting models to reason before acting [Yao et al., 2023b]. Similar reasoning methods include self-reflective prompting [Shinn et al., 2023, Wang et al., 2023b, Hao et al., 2023, Besta et al., 2024, Yao et al., 2023a], automated prompt tuning [Fernando et al., 2024, Hu et al., 2025], and strategic planning demonstrated by CICERO in the Diplomacy game [, FAIR]. However, frequent replanning can cause behavioural instability, analogous to RL frame-skipping strategies that advocate less frequent action repetition for improved exploration and consistency [Sharma et al., 2017, Kalyanakrishnan et al., 2021]. Recent studies also show diminishing returns and increased brittleness from excessive reasoning [Stechly et al., 2024, Mizrahi et al., 2024, Sui et al., 2025], emphasizing the need for adaptive planning mechanisms. Behaviour cloning with LLMs on data that includes textual reasoning between actions has been shown to help imitation learning [Yang et al., 2022, Hu and Clune, 2023].

**Test-Time Compute Scaling** More recently, test-time scaling has shown great promise, spearheaded by the results of OpenAI o1 [Jaech et al., 2024] and DeepSeek R1 [Guo et al., 2025]. These gains arise when LLMs improve their own reasoning traces through RL training on tasks with verifiable rewards [Lambert et al., 2024]. Methods such as STaR [Zelikman et al., 2022], Quiet-STaR [Zelikman et al., 2024], ScoRE [Kumar et al., 2025], and Meta-CoT [Xiang et al., 2025] showcase iterative self-improvement. Simple prompting strategies like s1 [Muennighoff et al., 2025] and critical insights from [Gandhi et al., 2025], demonstrating the necessity of supervised fine-tuning (SFT) priming with reasoning examples, further support this direction. Moreover, emergent planning capabilities have been observed from RL-trained base models as comments in code tasks [Zhao et al., 2025]. While fixed always-planning hierarchical strategies exist [Erdogan et al., 2025], their rigidity motivates research toward adaptive, dynamic approaches.

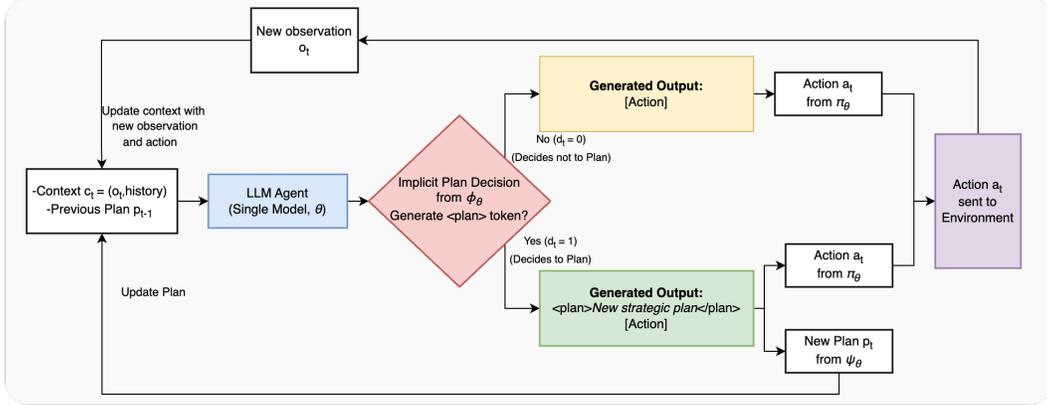


Figure 2: **Dynamic Planning Agent Architecture.** Our agent is a single, monolithic LLM whose conceptual policies are realized through its unified output format. The decision to plan ( $\phi_\theta$ ) is made implicitly by the model’s choice to begin its generation with a `<plan>` token. This single output string is then parsed to extract the action ( $a_t$ ) and, if present, the new plan ( $p_t$ ), thereby executing the acting ( $\pi_\theta$ ) and planning ( $\psi_\theta$ ) policies.

**Connection to Hierarchical Reinforcement Learning** Previous HRL work involving LLMs focused on discovering low-level skills [Wang et al., 2023a, Klissarov et al.] and making high-level decisions about which skills to use [Ahn et al., 2022, Klissarov et al., 2025]. Traditionally, HRL frameworks require the high-level controller to be invoked at every fixed interval or even every time step to decide which option or sub-policy to execute next; this repeated querying is common in modular HRL architectures and is a computational bottleneck [Sutton et al., 1999, Bacon et al., 2017]. This work diverges from the conventional HRL paradigm by focusing not on what sub-policy to choose, but on when to engage the costly planning process. By training an agent to dynamically decide when to plan, our approach directly addresses the computational constraints of using LLMs for strategic reasoning, optimizing the allocation of test-time compute.

**Steering LLM Agents** Recent studies have explored methods to steer LLM agents, such as influencing exploration through modulated representational uncertainty [Rahn et al., 2024], adaptively selecting reasoning modes based on task demands [Chen et al., 2024], and improving collaborative decision-making via step-wise RL evaluations [Zhou et al., 2025]. Our work demonstrates that LLM agents can be effectively steered through adaptive planning, enabling integration of external human-generated plans post-RL training.

### 3 A Conceptual Framework for Dynamic Planning with LLM Agents

Deciding when to allocate test-time compute for planning is a central challenge for LLM agents. To address this in a principled way, we first establish a conceptual framework that formalizes the underlying cost-benefit trade-offs. This framework provides the theoretical motivation for our practical training methodology, which uses reinforcement learning to teach an agent to implicitly master this dynamic planning skill.

Consider a sequential decision-making environment modelled as a Partially-Observable Markov Decision Process  $\langle S, A, T, R, O, \gamma \rangle$  (states, actions, stochastic transitions, rewards, observations, discount factor). An LLM agent with parameters  $\theta$  acts within this framework by generating tokens.

Specifically, at each timestep  $t$ , the agent receives an observation  $o_t$ , described in natural language, and maintains an internal context  $c_t = (o_t, history)$  which includes the current observation, a history of previous observations and actions, and any existing plan  $p_{t-1}$ . Formally, the agent’s behaviour is decomposed into a decision policy  $\phi_\theta$ , a planning policy  $\psi_\theta$ , and an acting policy  $\pi_\theta$ :

$$\phi_\theta(d_t | c_t, p_{t-1}), \quad \psi_\theta(p_t | c_t, p_{t-1}), \quad \pi_\theta(a_t | c_t, p_t)$$

Importantly, these three policies are not separate architectural components but rather a conceptual decomposition of the unified output from a single, monolithic LLM (Figure 2). The decision policy  $\phi_\theta$

corresponds to the decision  $d_t \in \{0, 1\}$ , where  $d_t = 1$  signifies that a new plan  $p_t$  will be generated by the planning policy  $\psi_\theta$ . If  $d_t = 0$ , the agent continues with the existing plan  $p_{t-1}$ . Thus the plan selection mechanism is:

$$p_t = d_t \cdot \psi_\theta(p_t \mid c_t, p_{t-1}) + (1 - d_t) \cdot p_{t-1}$$

Finally, the acting policy  $\pi_\theta$  generates action  $a_t$  based on  $c_t$  and  $p_t$ .

### 3.1 When Should an Agent Plan?

Intuitively, an agent should only plan when the expected benefit outweighs its cost. We quantify this state-dependent trade-off using a simple cost-benefit analysis:

The expected benefit of planning, or the *Planning Advantage*, measures how much the agent’s expected future rewards improve by adopting a new plan generated by  $\psi_\theta$  (i.e., if the decision  $d_t = 1$  is made), compared to continuing with the existing plan  $p_{t-1}$ . Conceptually, the value of generating a new plan is rooted in its potential to reduce the agent’s uncertainty about optimal future actions and to augment the context. By making strategic reasoning explicit, a new plan provides actionable insights that go beyond what is implicitly encoded in the agent’s internal representation (weights and activations). We formally define the planning advantage as the expected improvement in task-specific value, conditioned on the decision to generate a new plan:

$$A_{plan}(c_t) = \mathbb{E}_{p_t \sim \psi_\theta(\cdot \mid c_t, d_t=1)} [V^{\pi_\theta}(c_t, p_t) - V^{\pi_\theta}(c_t, p_{t-1})]$$

Here,  $V^{\pi_\theta}(c_t, p_t)$  represents the expected future rewards under the new plan  $p_t$ , and similarly for the existing plan  $p_{t-1}$ . While the agent does not explicitly compute  $A_{plan}(c_t)$  at each step, its decision policy  $\phi_\theta$  is trained to generate outputs that approximate this benefit, as detailed in Section 3.3.

The overall cost of planning,  $C_{plan}$ , arises from several sources:

$$C_{plan} = C_{tokens} + C_{latency} + C_{noise}$$

These components include:

**Computational Cost:** The direct cost of generating a plan, proportional to its token length:  $C_{tokens} = k_{tokens} \cdot |p_t|$ . This is a direct and measurable cost that we can explicitly penalize during training.

**Latency Cost:** The cost associated with the real-world time  $\Delta T_{plan}$  taken to plan. This is included for theoretical completeness, as it is a critical factor in time-sensitive applications like robotics, where its impact would be implicitly absorbed by the task reward. However, in the turn-based environments used in our experiments (POGS and Crafter), the environment pauses for the agent’s turn, so this cost is effectively zero ( $C_{latency} \approx 0$ ).

**Instability Cost:** This is a conceptual cost representing the performance degradation that can arise from erratic or excessive replanning. Frequent replanning, especially with imperfect or inconsistent plans, can introduce behavioral instability (e.g., inefficient backtracking, subgoal oscillation) that ultimately hinders task success. We model this conceptually as  $C_{noise} = k_{noise} \cdot f_p \cdot (1 - \bar{Q}_p)$ , where the negative impact of high planning frequency ( $f_p$ ) is magnified by low-quality plans (a low average plan quality  $\bar{Q}_p$ ). This cost is not explicitly calculated during training; instead, its effects are implicitly penalized because they naturally lead to lower task rewards. Our backtracking analysis in POGS (Appendix B) serves as an empirical proxy for this instability.

### 3.2 Plan Drift

The usefulness of an existing plan is not static; it typically diminishes over time as the agent acts and the environment evolves. This decay in relevance, or *plan drift*, makes replanning increasingly advantageous. Several factors contribute to how quickly plan drift occurs:

**Plan Abstraction Level:** High-level, conceptual plans (e.g., control the centre in chess) offer robustness against minor environmental shifts and remain relevant longer, though they provide less explicit guidance. Conversely, low-level detailed plans (e.g., specific move sequences) provide clearer direction but become outdated quickly.

**Planner and Model Accuracy:** Plans from highly accurate models tend to be robust and endure longer. In contrast, plans from imperfect models, like LLM natural language reasoning, may contain inaccuracies that accelerate their decay.

**Environment Dynamics:** The environment’s volatility significantly influences plan lifespan. In stable environments, plans retain value longer, while in dynamic environments with unpredictable shifts (e.g., an opponent’s unexpected move), existing plans can become instantly obsolete.

Understanding plan drift helps explain why agents must periodically reassess when to allocate compute to planning rather than following fixed planning strategies.

### 3.3 Training the Dynamic Planning Agent

To enable the agent to learn when to plan, thereby implicitly performing the cost-benefit analysis outlined in Section 3.1, we use RL fine-tuning. The agent’s policy parameters  $\theta$  (governing the decision policy  $\phi_\theta$ , planning policy  $\psi_\theta$  and acting policy  $\pi_\theta$ ) are optimized to maximize the expected discounted sum of task rewards, adjusted by a penalty for the computational cost of planning:

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{\tau \sim \theta} \left[ \sum_{t=0}^H \gamma^t (R_{task}(s_t, a_t) - d_t \cdot C_{tokens,t}) \right]$$

Other detrimental effects of poor planning strategies, such as those arising from excessive latency or instability (conceptualized as  $C_{latency}$  and  $C_{noise}$ ), are implicitly discouraged as they naturally lead to lower task rewards  $R_{task}(s_t, a_t)$ . Thus, by optimizing this objective, the decision policy  $\phi_\theta$ , planning policy  $\psi_\theta$ , and acting policy  $\pi_\theta$  should jointly learn to optimally decide *when* to plan ( $d_t = 1$ ), *how* to output plans ( $p_t$ ) that are beneficial, and how to output actions ( $a_t$ ) that effectively *follow* these plans, ensuring that the expected improvement in future task rewards (i.e., the empirical benefit corresponding to the conceptual  $A_{plan}(c_t)$ ) outweighs the explicit cost  $C_{tokens,t}$  as well as any implicit degradation of  $R_{task}$  due to poor planning.

## 4 Experimental Setup

Our experiments evaluate planning agents across diverse settings. In this section, we detail the environments used, the core evaluation protocol, and the specific setups for evaluation, SFT, and RL.

### 4.1 Environments

To evaluate dynamic planning across different conditions, we selected two complementary environments. First, **Partially Observable Graph Search (POGS)** is our custom synthetic environment designed to isolate planning under uncertainty. Agents navigate procedurally generated graph using only local observations, which require adaptive replanning upon discovering new nodes or dead ends. POGS allows measurement of exploration efficiency via backtracking statistics. Second, **Crafter** [Hafner, 2022] is a complex 2D grid-world, long-horizon benchmark inspired by Minecraft. It demands multi-scale planning for survival, resource management, and crafting, testing both short-term tactical decisions and long-term strategic choices. Interaction in both environments occurs via natural language. Figure 3 showcases the environments, and full technical details are provided in the Appendix B.

### 4.2 Evaluation Protocol

We utilize the BALROG benchmark [Paglieri et al., 2025a] for standardized agent evaluation and environment interaction. At each timestep  $t$ , the agent receives its history and current observation  $o_t$  within a chat template, guided by a system prompt outlining the task. The agent’s response must include a natural language action command  $a_t$ . Our dynamic planning agents are instructed to decide at each step whether to plan. If they choose not to plan, they output only the action command [Action]. If they choose to plan, they output the plan followed by the action, using the format `<plan> [natural language plan] </plan> [Action]`. BALROG parses this output, identifying a planning decision ( $d_t = 1$ ) if the `<plan>` block is present and using its content as the current plan  $p_t$  in subsequent context. The [Action] is always extracted and executed. Fallback mechanisms ensure robustness against invalid outputs. Appendix A provides detailed prompts.

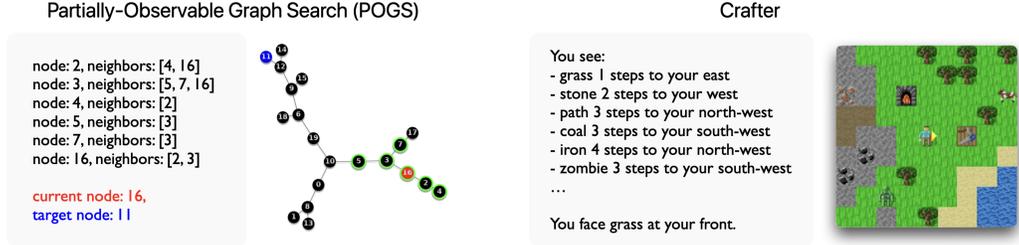


Figure 3: **POGS and Crafter environments.** POGS (left): Agent navigates a procedurally generated graph with partial visibility. Crafter (right): Agent receives natural language descriptions of terrain, resources, and creatures with their relative positions.

### 4.3 Zero-Shot Evaluation

To understand baseline capabilities and the raw effect of planning frequency, we perform zero-shot evaluations using Llama-3.3-70B-Instruct [Grattafiori et al., 2024] on POGS and Crafter (100 seeds each). We compare different prompting strategies without any fine-tuning. We test a **Naive Agent**, prompted to only output actions and thus never plan. We also test **Fixed-Frequency Planners**, which are prompted to plan every  $k$  steps for various  $k \in \{1, 2, 4, 8, \dots\}$ ; these agents are instructed to output a plan followed by an action every  $k$  steps, and only an action otherwise. These evaluations measure performance trade-offs associated purely with inference-time planning strategies.

### 4.4 Supervised Fine-Tuning (SFT)

To prepare Llama-3.1-8B-Instruct for RL, we first perform SFT priming as in [Gandhi et al., 2025].

**Data Generation:** We created a dataset of 1024 Crafter trajectories using Llama-3.3-70B-Instruct as a teacher. To ensure diversity, the teacher planned every  $K$  steps ( $K \sim U[2, 12]$  per trajectory) using 16 different planning prompts (Appendix A).

**SFT Priming:** The Llama-3.1-8B model is fine-tuned on this data, aligning the SFT process with the target RL configuration. For the **SFT+RL plan dynamically** agent, SFT targets use the dynamic format `<plan> . . . </plan> [Action]` if the teacher planned at that step, otherwise just `[Action]` – along with a dynamic prompt encouraging this choice. Conversely, for the **SFT+RL no plan** agent, SFT targets include only the actions (`[Action]`, with all plan blocks removed), paired with a naive prompt focused solely on action prediction. This prepares the model appropriately for the subsequent RL phase. More details on the prompt in the Appendix A.

### 4.5 Reinforcement Learning (RL)

We then use Proximal Policy Optimization (PPO) [Schulman et al., 2017] to fine-tune Llama-3.1-8B-Instruct agents in Crafter, optimizing task rewards possibly adjusted for planning costs (Sec. 3.3). We compare four key configurations:

- **Base+RL plan dynamically:** RL on the base model using a dynamic planning prompt.
- **Base+RL no plan:** RL on the base model using the naive (action-only) prompt.
- **SFT+RL plan dynamically:** RL on the SFT-primed dynamic planning model, using the dynamic planning prompt.
- **SFT+RL no plan:** RL on the naive SFT-primed naive model, using the naive prompt.

This isolates learned dynamic planning from fixed strategies, and the benefit of SFT priming with versus without explicit plan information. Further training details are in the Appendix C.

## 5 Results

We present findings analyzing the impact of planning frequency in zero-shot settings and the effectiveness of our SFT priming approach, and the RL results.



Figure 4: **Human-Agent collaboration in Crafter.** We show an example where a human guides the agent with high-level plans to clear a cave from a skeleton, and create a shelter to survive the night, a complex behaviour that was not observed in any of the training runs otherwise.

### 5.1 Zero-Shot Evaluation

We evaluated the zero-shot performance of Llama-3.3-70B-Instruct at various planning frequencies, as shown in Figure 1 (a) and (b) for POGS and Crafter. Average task progression is plotted against mean output tokens (log scale, proxying computational cost), comparing agents that plan at different intervals—from never to every step.

Results in both environments reveal a non-monotonic relationship: some planning improves over naive agents, but excessive planning reduces performance. Performance peaks at intermediate frequencies (e.g., every 4 steps in Crafter), then declines as planning becomes more frequent. The always-plan approach, similar to standard ReAct, consistently underperforms compared to less frequent planning, despite higher computational cost. This demonstrates an optimal, task-dependent "Goldilocks" zone for planning frequency. This supports our framework (Section 3), which introduces an instability cost ( $C_{noise}$ ) from frequent planning. While planning reduces uncertainty, replanning too often with noisy or inconsistent plans introduces instability, causing inefficient behaviours like oscillating between subgoals or abandoning viable actions.

Quantitative backtracking analysis in POGS (Appendix B) provides clear evidence: the always-plan agent showed the most backtracking, highlighting its tendency to oscillate or revisit states rather than explore efficiently. Agents planning at intermediate frequencies achieved better progression with significantly less backtracking, reflecting more systematic exploration. Crafter shows similar patterns: frequent replanning causes course changes that stall progress and induce inefficiencies.

Our initial attempts to elicit adaptive dynamic planning behaviour directly via complex zero-shot prompting proved challenging and unreliable. Models struggled to consistently interpret instructions like "plan when necessary," often defaulting to fixed patterns. This difficulty underscores the need for learning-based approaches, like SFT and RL, to effectively teach agents this meta-cognitive skill of deciding *when* to allocate test-time compute for planning.

### 5.2 SFT Priming

Having established the importance of planning frequency, we next evaluate our SFT priming stage (Section 4.4), designed to prepare the Llama-3.1-8B model for subsequent RL fine-tuning. We compare the performance of the 'Primed-Dynamic' and 'Primed-Naive' models during the SFT process itself, evaluating their Crafter progression periodically throughout training (Figure 1c).

The key comparison here is that both models are fine-tuned on datasets derived from the *exact same underlying action sequences*. The sole difference lies in the target outputs: 'Primed-Dynamic' is trained to predict the explicit natural language plans generated by the teacher model (when present) alongside the actions, using the dynamic planning format, while 'Primed-Naive' is trained to predict only the actions, with all plan information removed.

Despite the identical action supervision, the 'Primed-Dynamic' model consistently achieves higher task progression throughout the SFT phase. This result strongly suggests that access to the explicit natural language plans significantly aids the imitation learning process itself. We hypothesize several contributing factors for this observation:

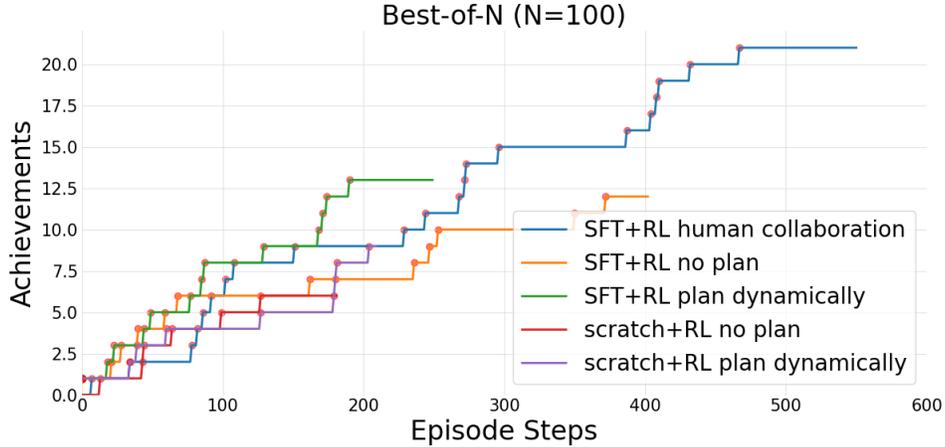


Figure 5: **Best-of-N (N=100) comparison on Crafter.** The plot shows achievements versus episode steps for each method, evaluated on the best of N trajectories. Human collaboration (N = 20) demonstrates the strongest progression, followed by SFT+RL plan dynamically and SFT+RL no plan, with base RL baselines lagging behind.

**Explanatory Power:** Plans provide semantic context and rationale for subsequent actions, potentially simplifying the behavioural cloning objective by making the action sequence more predictable.

**Learning Planning Structure:** Exposure to diverse planning data likely helps the model learn generalizable representations related to “look-ahead,” improving performance beyond simple imitation.

**Regularization and Grounding:** Natural language plans may act as a form of cognitive grounding, constraining the model’s behaviour during fine-tuning and preventing it from deviating too drastically from the base model’s capabilities in unproductive ways.

Supporting the grounding hypothesis, preliminary analyses indicate that the “Primed-Dynamic” model exhibits lower KL divergence from the base pre-trained model during SFT compared to the “Primed-Naive” model (see Figure 1d). This suggests that training with plans may support a less disruptive fine-tuning process. Overall, these findings highlight the value of incorporating explicit reasoning in the form of natural language plans directly into training data – not merely as a target behaviour, but as an effective mechanism to improve the learning process itself.

### 5.3 RL Fine-Tuning

Following SFT priming, we fine-tuned all agent variants through RL. Our results reveal that SFT-primed agent equipped with dynamic planning (SFT+RL plan dynamically) substantially outperforms the non-planning counterparts (SFT+RL no plan) during the initial phases of gameplay. Notably, qualitative analysis shows that the SFT+RL plan dynamically agent is able to generate and execute plans at multiple levels of abstraction, flexibly shifting between high-level strategic objectives and lower-level tactical actions depending on the demands of the current situation. Furthermore, the agent demonstrates the ability to plan only when necessary, autonomously executing plans and replanning in response to changing circumstances. Qualitative results in Appendix D.

However, as the environment becomes more challenging, the improvement rates for both SFT+RL agents diminish, resulting in eventual plateauing of their progression curves at comparable levels. This plateau suggests that while dynamic planning confers a clear advantage in early and mid-game scenarios, both agent types ultimately encounter similar bottlenecks in late-game survival and high-level achievement. Qualitative analysis attributes this stagnation to persistent deficiencies in survival skills such as acquiring food and water or building a shelter, indicating that further progress likely requires additional training beyond our current computational budget.

When comparing agents initialized from the base Llama-3.1-8B-Instruct model (Base+RL plan dynamically and Base+RL no plan), we observe inferior performance throughout training. Notably, the Base+RL plan dynamically agent exhibits the most limited progression, struggling to leverage planning in the absence of SFT priming. We hypothesize that the 8B model, without prior SFT, may possess insufficient scale or inherent capability to autonomously generate planning traces that

are consistently effective for solving the task. Consequently, for models of this scale, structured, expert-like planning examples provided through SFT appear crucial for enabling effective downstream RL optimization and the emergence of dynamic planning behaviour.

**Instruction following:** Finally, we explore human-agent collaboration by steering agents with human-written plans. The base model exhibits significant grounding issues, often attempting invalid actions that lead to unproductive behaviour and quick death when fighting a monster. The SFT-primed model, while much more efficient than the base model, still exhibits numerous errors, doesn't learn from feedback, and struggles to adhere to novel plans, often defaulting to previously learned strategies due to excessive memorization. For example, an SFT agent might ignore a human plan to mine stone, instead prioritizing combat with nearby monsters. In contrast, RL-trained agents both excel at the game and reliably follow external plans; one such agent was successfully guided by human plans to collect diamond and solve Crafter—an achievement unseen in autonomous agents. We also show a complex plan following example in Figure 4, while we refer the reader to the Appendix D for more qualitative results. To quantify the upper bound of performance, we conducted a Best-of-N analysis, comparing the best runs attained by each method across 100 independent runs, and 20 runs for the human collaboration. As shown in Figure 5, the SFT+RL when guided by human plans, successfully solves Crafter by mining a diamond, significantly outperforming all of the methods. This result demonstrates the impact of human steering and further validates the benefits of conditioning the models with plans.

In summary, our RL experiments demonstrate that SFT priming is essential for enabling dynamic planning to emerge in LLM agents. The resulting SFT+RL plan dynamically agent exhibits better sample efficiency and superior performance on advanced metrics. Nevertheless, persistent challenges in late-game survival and the observed performance plateau indicate that further advances will require additional training resources or other innovations. Addressing these limitations will be crucial for developing the next wave of reasoning agents for agentic tasks.

## 6 Discussion, Limitations & Conclusion

The ability for LLMs to leverage test-time compute has been transformative, yet efficiently allocating these resources in agentic settings stands as a critical, largely unexplored challenge. To the best of our knowledge, this work presents the first systematic investigation and learned solution enabling LLM agents to effectively allocate test-time compute in sequential decision-making tasks. Our findings reveal that prevailing always-plan (e.g. ReAct) and never-plan approaches are suboptimal. Instead, a "Goldilocks" effect emerges whereby intermediate frequencies outperform both extremes, likely avoiding instability ( $C_{noise}$  in our framework) induced by excessive replanning. This highlights the need for LLM agents to strategically allocate, rather than naively scale, deliberation resources. Our two-stage SFT+RL methodology demonstrates that agents can learn this meta-cognitive skill, moving beyond fixed heuristics towards the adaptive, efficient behaviour essential for sustained autonomy. Moreover, the resulting agents become sufficiently adept at planning and execution to be effectively steered by human-written plans towards remarkable feats, including the full completion of Crafter through diamond collection, significantly impacting human-AI collaboration and safety.

While our results establish the value of dynamic test-time compute allocation, several limitations suggest directions for future work. Our experiments focused on specific models at certain scales (Llama-3.1-8B-Instruct for fine-tuning, Llama-3.3-70B-Instruct for evaluation) due to computational constraints. Investigating how optimal compute allocation strategies scale with model parameters would provide valuable insights. Additionally, extending this work beyond our current environments (POGS and Crafter) to more diverse domains would further validate the generality of our approach. Future research could also explore more sophisticated compute allocation mechanisms, scale up our experiments, or investigate methods to more explicitly integrate our conceptual framework's insights into novel RL algorithms.

In summary, this paper establishes that in zero-shot evaluation, per-timestep planning or reasoning strategies akin to ReAct are outperformed by "Goldilocks" planning frequencies, likely due to instability that results from excessive planning or "overthinking." By using a two-stage training methodology combining SFT and RL, we successfully train agents to dynamically allocate planning resources at test-time. This approach yields more effective and efficient behaviour compared to fixed planning strategies, marking a step towards more autonomous and scalable agentic systems.

## References

- Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- Reza Yazdani Aminabadi, Samyam Rajbhandari, Ammar Ahmad Awan, Cheng Li, Du Li, Elton Zheng, Olatunji Ruwase, Shaden Smith, Minjia Zhang, Jeff Rasley, and Yuxiong He. Deepspeed-inference: Enabling efficient inference of transformer models at unprecedented scale. In *SC22: International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–15, 2022. doi: 10.1109/SC41404.2022.00051.
- Pierre-Luc Bacon, Jean Harb, and Doina Precup. The option-critic architecture. In *Proceedings of the AAAI conference on artificial intelligence*, volume 31, 2017.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17682–17690, 2024.
- Yongchao Chen, Harsh Jhamtani, Srinagesh Sharma, Chuchu Fan, and Chi Wang. Steering large language models between code execution and textual reasoning. *CoRR*, abs/2410.03524, 2024. URL <https://doi.org/10.48550/arXiv.2410.03524>.
- Rémi Coulom. Efficient selectivity and backup operators in monte-carlo tree search. In *International conference on computers and games*, pages 72–83. Springer, 2006.
- Lutfi Eren Erdogan, Nicholas Lee, Sehoon Kim, Suhong Moon, Hiroki Furuta, Gopala Anumanchipalli, Kurt Keutzer, and Amir Gholami. Plan-and-act: Improving planning of agents for long-horizon tasks. *arXiv preprint arXiv:2503.09572*, 2025.
- Meta Fundamental AI Research Diplomacy Team (FAIR)<sup>†</sup>, Anton Bakhtin, Noam Brown, Emily Dinan, Gabriele Farina, Colin Flaherty, Daniel Fried, Andrew Goff, Jonathan Gray, Hengyuan Hu, et al. Human-level play in the game of diplomacy by combining language models with strategic reasoning. *Science*, 378(6624):1067–1074, 2022.
- Chrisantha Fernando, Dylan Sunil Banarse, Henryk Michalewski, Simon Osindero, and Tim Rocktäschel. Promptbreeder: Self-referential self-improvement via prompt evolution, 2024. URL <https://openreview.net/forum?id=HKkiX32Zw1>.
- Kanishk Gandhi, Ayush Chakravarthy, Anikait Singh, Nathan Lile, and Noah D. Goodman. Cognitive behaviors that enable self-improving reasoners, or, four habits of highly effective stars. *CoRR*, abs/2503.01307, 2025. doi: 10.48550/ARXIV.2503.01307. URL <https://doi.org/10.48550/arXiv.2503.01307>.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Daya Guo, Dejian Yang, Haowei 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.
- David Ha and Jürgen Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018.
- Danijar Hafner. Benchmarking the spectrum of agent capabilities. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=1W0z96MFEoH>.
- Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James Davidson. Learning latent dynamics for planning from pixels. In *International conference on machine learning*, pages 2555–2565. PMLR, 2019.
- Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=S110TC4tDS>.

- Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. Reasoning with language model is planning with world model. In *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023. URL <https://openreview.net/forum?id=VTWwvYtF1R>.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.
- Shengran Hu and Jeff Clune. Thought cloning: Learning to think while acting by imitating human thinking. In *Advances in Neural Information Processing Systems 37: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2023, New Orleans, 2023*, 2023.
- Shengran Hu, Cong Lu, and Jeff Clune. Automated design of agentic systems. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=t9U3LW7JVX>.
- Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.
- Damjan Kalajdzievski. A rank stabilization scaling factor for fine-tuning with lora. *arXiv preprint arXiv:2312.03732*, 2023.
- Shivaram Kalyanakrishnan, Siddharth Aravindan, Vishwajeet Bagawat, Varun Bhatt, Harshith Goka, Archit Gupta, Kalpesh Krishna, and Vihari Piratla. An analysis of frame-skipping in reinforcement learning. *arXiv preprint arXiv:2102.03718*, 2021.
- Martin Klissarov, Mikael Henaff, Roberta Raileanu, Shagun Sodhani, Pascal Vincent, Amy Zhang, Pierre-Luc Bacon, Doina Precup, Marlos C Machado, and Pierluca D’Oro. Maestromotif: Skill design from artificial intelligence feedback, 2024. URL <https://arxiv.org/abs/2412.8542>.
- Martin Klissarov, Akhil Bagaria, Ziyang Luo, George Konidaris, Doina Precup, and Marlos C Machado. Discovering temporal structure: An overview of hierarchical reinforcement learning. *arXiv preprint arXiv:2506.14045*, 2025.
- Aviral Kumar, Vincent Zhuang, Rishabh Agarwal, Yi Su, John D Co-Reyes, Avi Singh, Kate Baumli, Shariq Iqbal, Colton Bishop, Rebecca Roelofs, Lei M Zhang, Kay McKinney, Disha Shrivastava, Cosmin Paduraru, George Tucker, Doina Precup, Feryal Behbahani, and Aleksandra Faust. Training language models to self-correct via reinforcement learning. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=CjwERcAU7w>.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. Tulu 3: Pushing frontiers in open language model post-training. *arXiv preprint arXiv:2411.15124*, 2024.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- David Q Mayne, James B Rawlings, Christopher V Rao, and Pierre OM Scokaert. Constrained model predictive control: Stability and optimality. *Automatica*, 36(6):789–814, 2000.
- Moran Mizrahi, Guy Kaplan, Dan Malkin, Rotem Dror, Dafna Shahaf, and Gabriel Stanovsky. State of what art? a call for multi-prompt llm evaluation. *Transactions of the Association for Computational Linguistics*, 12:933–949, 2024.

- Philipp Moritz, Robert Nishihara, Stephanie Wang, Alexey Tumanov, Richard Liaw, Eric Liang, Melih Elibol, Zongheng Yang, William Paul, Michael I Jordan, et al. Ray: A distributed framework for emerging {AI} applications. In *13th USENIX symposium on operating systems design and implementation (OSDI 18)*, pages 561–577, 2018.
- Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time scaling. *arXiv preprint arXiv:2501.19393*, 2025.
- Davide Paglieri, Bartłomiej Cupiał, Samuel Coward, Ulyana Piterbarg, Maciej Wolczyk, Akbir Khan, Eduardo Pignatelli, Łukasz Kuciński, Lerrel Pinto, Rob Fergus, Jakob Nicolaus Foerster, Jack Parker-Holder, and Tim Rocktäschel. BALROG: Benchmarking agentic LLM and VLM reasoning on games. In *The Thirteenth International Conference on Learning Representations*, 2025a. URL <https://openreview.net/forum?id=fp6t3F669F>.
- Davide Paglieri, Bartłomiej Cupiał, Samuel Coward, Ulyana Piterbarg, Maciej Wolczyk, Akbir Khan, Eduardo Pignatelli, Łukasz Kuciński, Lerrel Pinto, Rob Fergus, et al. Balrog leaderboard, 2025b. URL <https://balrogai.com>.
- Ben Prystawski, Michael Li, and Noah Goodman. Why think step by step? reasoning emerges from the locality of experience. *Advances in Neural Information Processing Systems*, 36:70926–70947, 2023.
- Nate Rahn, Pierluca D’Oro, and Marc G. Bellemare. Controlling large language model agents with entropic activation steering. *CoRR*, abs/2406.00244, 2024. URL <https://doi.org/10.48550/arXiv.2406.00244>.
- Laura Ruis, Maximilian Mozes, Juhan Bae, Siddhartha Rao Kamalakara, Dwaraknath Gnaneshwar, Acyr Locatelli, Robert Kirk, Tim Rocktäschel, Edward Grefenstette, and Max Bartolo. Procedural knowledge in pretraining drives reasoning in large language models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=1hQKHHUsMx>.
- Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering atari, go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Sahil Sharma, Aravind S. Lakshminarayanan, and Balaraman Ravindran. Learning to repeat: Fine grained action repetition for deep reinforcement learning. In *International Conference on Learning Representations*, 2017. URL <https://openreview.net/forum?id=B1G0WV5eg>.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36:8634–8652, 2023.
- David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dharshan Kumaran, Thore Graepel, et al. Mastering chess and shogi by self-play with a general reinforcement learning algorithm. *arXiv preprint arXiv:1712.01815*, 2017a.
- David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017b.
- Charlie Victor Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling LLM test-time compute optimally can be more effective than scaling parameters for reasoning. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=4FWAwZtd2n>.

- Kaya Stechly, Karthik Valmeekam, and Subbarao Kambhampati. Chain of thoughtlessness? an analysis of cot in planning. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang, editors, *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024. URL [http://papers.nips.cc/paper\\_files/paper/2024/hash/3365d974ce309623bd8151082d78206c-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2024/hash/3365d974ce309623bd8151082d78206c-Abstract-Conference.html).
- Yang Sui, Yu-Neng Chuang, Guanchu Wang, Jiamu Zhang, Tianyi Zhang, Jiayi Yuan, Hongyi Liu, Andrew Wen, Hanjie Chen, Xia Hu, et al. Stop overthinking: A survey on efficient reasoning for large language models. *arXiv preprint arXiv:2503.16419*, 2025.
- Richard S Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):181–211, 1999.
- Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlkar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*, 2023a.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*, 2023b. URL <https://openreview.net/forum?id=1PL1NIMMrw>.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL [https://openreview.net/forum?id=\\_VjQ1MeSB\\_J](https://openreview.net/forum?id=_VjQ1MeSB_J).
- Violet Xiang, Charlie Snell, Kanishk Gandhi, Alon Albalak, Anikait Singh, Chase Blagden, Duy Phung, Rafael Rafailov, Nathan Lile, Dakota Mahan, Louis Castricato, Jan-Philipp Fränken, Nick Haber, and Chelsea Finn. Towards system 2 reasoning in llms: Learning how to think with meta chain-of-thought. *CoRR*, abs/2501.04682, 2025. doi: 10.48550/ARXIV.2501.04682. URL <https://doi.org/10.48550/arXiv.2501.04682>.
- Mengjiao Yang, Dale Schuurmans, Pieter Abbeel, and Ofir Nachum. Chain of thought imitation with procedure cloning. *CoRR*, abs/2205.10816, 2022.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik R Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023a. URL <https://openreview.net/forum?id=5Xc1ecx01h>.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R. Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023b. URL [https://openreview.net/forum?id=WE\\_vluYUL-X](https://openreview.net/forum?id=WE_vluYUL-X).
- Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. Star: Bootstrapping reasoning with reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488, 2022.
- Eric Zelikman, Georges Harik, Yijia Shao, Varuna Jayasiri, Nick Haber, and Noah D Goodman. Quiet-star: Language models can teach themselves to think before speaking. *arXiv preprint arXiv:2403.09629*, 2024.
- Andrew Zhao, Yiran Wu, Yang Yue, Tong Wu, Quentin Xu, Yang Yue, Matthieu Lin, Shenzhi Wang, Qingyun Wu, Zilong Zheng, and Gao Huang. Absolute zero: Reinforced self-play reasoning with zero data. *arXiv preprint arXiv:2505.03335*, 2025.
- Yifei Zhou, Song Jiang, Yuandong Tian, Jason Weston, Sergey Levine, Sainbayar Sukhbaatar, and Xian Li. SWEET-RL: training multi-turn LLM agents on collaborative reasoning tasks. *CoRR*, abs/2503.15478, 2025. URL <https://doi.org/10.48550/arXiv.2503.15478>.

## A Appendix

### A.1 Planning Prompts

To ensure behavioural diversity in our SFT dataset, we designed 16 distinct planning prompts, each eliciting a different style of planning. These include prompts that ask the model to identify immediate subgoals, describe short-term action sequences, re-evaluate high-level strategies, or address gaps in information. The idea is to expose the model to a wide range of planning behaviours so that, during RL fine-tuning, it can learn when and how to use the most appropriate planning strategy depending on the situation. When planning is not required, a separate instruction prompt (Prompt 17) is used to guide the model to select the next action based solely on prior plans and observations.

#### 1. Immediate-Subgoal Planner

Identify the immediate next subgoal required to progress towards the overall task completion.

Outline your plan to achieve this specific subgoal, including any necessary reasoning.

Output your plan strictly in the following format:

```
<plan>YOUR_PLAN_FOR_SUBGOAL </plan>
```

Replace YOUR\_PLAN\_FOR\_SUBGOAL with your own plan.

Keep your plan relatively brief, only focusing on important information.

After your plan, choose exactly ONE action from the allowed actions list that initiates this plan.

Output no other text.

Prompt 1: Prompt encouraging subgoal identification and targeted execution.

#### 2. Milestone-Focused Planner

Consider the overall objective.

What is the most crucial intermediate milestone to achieve next?

Explain why reaching this milestone is important for the overall task, and outline the steps you'll take to get there.

Output your reasoning and plan strictly in the following format:

```
<plan>YOUR_REASONING_AND_SUBGOAL_PLAN</plan>
```

Replace YOUR\_REASONING\_AND\_SUBGOAL\_PLAN with your own plan.

Keep your plan relatively brief, only focusing on important information.

After your plan, choose exactly ONE action from the allowed actions list to start working towards this milestone.

Output no other text.

Prompt 2: Prompt focused on achieving the next key milestone and explaining its relevance.

### 3. Short-Term Sequence Planner

Detail the specific sequence of actions you intend to take over the next few steps.

Explain the purpose of this sequence in relation to the current situation.

Output your detailed short-term plan strictly in the following format:

`<plan>YOUR_DETAILED_SHORT_TERM_PLAN</plan>`

Replace YOUR\_DETAILED\_SHORT\_TERM\_PLAN with your own plan.

Keep your plan relatively brief, only focusing on important information.

After your plan, choose exactly ONE action from the allowed actions list, which should be the first step in your detailed plan.

Output no other text.

Prompt 3: Prompt directing the agent to outline a short, purposeful action sequence.

### 4. Step-by-Step Immediate Planner

Think step-by-step for the immediate future.

What actions are needed right now and why?

Describe the logic connecting these immediate actions to the next phase of the task.

Output your step-by-step thinking and plan strictly in the following format:

`<plan>YOUR_STEP_BY_STEP_LOGIC_AND_PLAN</plan>`

Replace YOUR\_STEP\_BY\_STEP\_LOGIC\_AND\_PLAN with your own plan.

Keep your plan relatively brief, only focusing on important information.

After your plan, choose exactly ONE action from the allowed actions list that represents the very next concrete step.

Output no other text.

Prompt 4: Prompt guiding step-by-step reasoning for immediate next actions.

### 5. Gap-Bridging Phase Planner

Analyze the current state and the final goal.  
Formulate a plan that bridges the gap, focusing on the most logical next phase of work.  
Explain how this phase contributes to the overall objective.

Output your analysis and plan strictly in the following format:  
<plan>YOUR\_BRIDGING\_PLAN\_AND\_REASONING</plan>  
Replace YOUR\_BRIDGING\_PLAN\_AND\_REASONING with your own plan.

Keep your plan relatively brief, only focusing on important information.

After your plan, choose exactly ONE action from the allowed actions list to begin executing this phase.  
Output no other text.

Prompt 5: Prompt for bridging the gap between the current state and final objective.

### 6. High-Level Strategic Planner

Re-evaluate the overall strategy.  
Outline your current high-level plan or strategic direction for completing the task from this point forward, focusing on the major phases ahead.

Output your strategic plan strictly in the following format:  
<plan>YOUR\_HIGH\_LEVEL\_STRATEGIC\_PLAN</plan>  
Replace YOUR\_HIGH\_LEVEL\_STRATEGIC\_PLAN with your own plan.

Keep your plan relatively brief, only focusing on important information.

After your plan, choose exactly ONE action from the allowed actions list that aligns with the first step of this strategy.  
Output no other text.

Prompt 6: Prompt eliciting a broad, high-level plan covering major task phases.

### 7. Justified Approach Planner

Propose a plan for the next stage of the task.  
Critically, justify *\*why\** this sequence of steps (or this approach) is the most sensible course of action right now.  
Output your plan and justification strictly in the following format:

Output your plan and justification strictly in the following format:  
<plan>YOUR\_PLAN\_WITH\_JUSTIFICATION</plan>  
Replace YOUR\_PLAN\_WITH\_JUSTIFICATION with your own plan.

Keep your plan relatively brief, only focusing on important information.

After your plan, choose exactly ONE action from the allowed actions list that initiates your justified plan.  
Output no other text.

Prompt 7: Prompt requiring justification for the chosen course of action.

### 8. Reasoning-Process Planner

Verbalize your thought process for deciding what to do next.  
Explain your reasoning, considering the current situation and the ultimate goal, and then state your resulting plan for the near term.

Output your reasoning process and plan strictly in the following format:  
<plan>YOUR\_REASONING\_PROCESS\_AND\_PLAN</plan>  
Replace YOUR\_REASONING\_PROCESS\_AND\_PLAN with your own plan.

Keep your plan relatively brief, only focusing on important information.

After your plan, choose exactly ONE action from the allowed actions list based on your reasoning.  
Output no other text.

Prompt 8: Prompt asking the agent to verbalize its reasoning before planning.

### 9. Approach-Comparison Planner

Briefly consider possible approaches for the next steps.  
State the approach you choose to take and why it seems preferable to alternatives right now.  
Outline the plan based on this chosen approach.

Output your chosen approach, rationale, and plan strictly in the following format:

`<plan>YOUR_CHOSEN_APPROACH_RATIONALE_AND_PLAN</plan>`

Replace YOUR\_CHOSEN\_APPROACH\_RATIONALE\_AND\_PLAN with your own plan.

Keep your plan relatively brief, only focusing on important information.

After your plan, choose exactly ONE action from the allowed actions list that corresponds to your chosen approach.  
Output no other text.

Prompt 9: Prompt comparing alternative approaches and selecting one.

### 10. Efficiency-Driven Planner

Devise a plan to make progress efficiently.  
What is the most direct path to achieving the next significant step or subgoal?  
Outline this efficient path.

Output your efficiency-focused plan strictly in the following format:

`<plan>YOUR_EFFICIENT_PLAN</plan>`

Replace YOUR\_EFFICIENT\_PLAN with your own plan.

Keep your plan relatively brief, only focusing on important information.

After your plan, choose exactly ONE action from the allowed actions list that represents the first step on this path.  
Output no other text.

Prompt 10: Prompt focused on generating the most direct, efficient plan forward.

## 11. Information-Gap Planner

Is there critical information missing?  
If so, formulate a plan focused on gathering the necessary information or resolving key uncertainties before proceeding with the main task execution.  
If not, state your plan for the next execution steps.

Output your information-gathering or execution plan strictly in the following format:  
`<plan>YOUR_INFORMATION_OR_EXECUTION_PLAN</plan>`  
Replace YOUR\_INFORMATION\_OR\_EXECUTION\_PLAN with your own plan.

Keep your plan relatively brief, only focusing on important information.

After your plan, choose exactly ONE action from the allowed actions list relevant to this plan.  
Output no other text.

Prompt 11: Prompt addressing missing information or uncertainty before acting.

## 12. Logical-Sequence Planner

Describe the logical sequence of operations you intend to perform next.  
Explain the dependency: why does step B follow step A? Focus on the immediate sequence.

Output your logical sequence and rationale strictly in the following format:  
`<plan>YOUR_LOGICAL_SEQUENCE_PLAN</plan>`  
Replace YOUR\_LOGICAL\_SEQUENCE\_PLAN with your own plan.

Keep your plan relatively brief, only focusing on important information.

After your plan, choose exactly ONE action from the allowed actions list representing the first operation in your sequence.  
Output no other text.

Prompt 12: Prompt outlining a logically dependent sequence of actions.

### 13. Short-term goal Planner

Define your immediate goal for the next few actions.  
Construct a plan specifically aimed at achieving this immediate goal.

Output your immediate goal and plan strictly in the following format:

`<plan>YOUR_IMMEDIATE_GOAL_AND_PLAN</plan>`

Replace YOUR\_IMMEDIATE\_GOAL\_AND\_PLAN with your own plan.

Keep your plan relatively brief, only focusing on important information.

After your plan, choose exactly ONE action from the allowed actions list that starts this plan.

Output no other text.

Prompt 13: Prompt focused on defining and achieving an immediate goal.

### 14. Practical-Progress Planner

Considering the available actions and the task objective, formulate a practical plan for the next steps.

What needs to be done now to make steady progress?

Output your practical plan strictly in the following format:

`<plan>YOUR_PRACTICAL_PROGRESS_PLAN</plan>`

Replace YOUR\_PRACTICAL\_PROGRESS\_PLAN with your own plan.

Keep your plan relatively brief, only focusing on important information.

After your plan, choose exactly ONE action from the allowed actions list to implement the first step.

Output no other text.

Prompt 14: Prompt aimed at formulating a grounded, actionable next step.

### 15. Intent-and-Approach Planner

State your intention for the next phase of action.

What do you aim to accomplish in the near future, and what's the general approach?

Output your statement of intent and approach strictly in the following format:

`<plan>YOUR_INTENTION_AND_APPROACH</plan>`

Replace YOUR\_INTENTION\_AND\_APPROACH with your own plan.

Keep your plan relatively brief, only focusing on important information.

After your plan, choose exactly ONE action from the allowed actions list that reflects this intention.

Output no other text.

Prompt 15: Prompt stating the agent's intent and general method for proceeding.

## 16. Next-Steps Planner

Outline your plan for what to do next.  
Keep it focused on the immediate steps required.

Output your plan strictly in the following format:  
<plan>YOUR\_NEXT\_STEPS\_PLAN</plan>  
Replace YOUR\_NEXT\_STEPS\_PLAN with your own plan.

Keep your plan relatively brief, only focusing on important information.

After your plan, choose exactly ONE action from the allowed actions list to start.  
Output no other text.

Prompt 16: Prompt asking for a concise plan of immediate next actions.

## Act instruction

Look at your previous plan and observations, then choose exactly ONE action from the allowed actions listed previously.  
Output no other text.

Prompt 17: Prompt instructing the model to directly choose the next action based on previous plan and observations.

## A.2 Eval Zero Shot Prompts

Prompt 18 is used for the zero-shot evaluation agent that never plans and is simply tasked with outputting the next action based on previous observations. Prompt 19 is used for the plan-every-k-steps agents, where the model is instructed to generate a plan at fixed intervals and to act directly at other times. In such cases, Prompt 17 is used at steps where planning is not required. These zero-shot evaluations help assess the impact of planning frequency in the absence of fine-tuning.

### Never Plan Prompt

```
Look at your previous observations, then choose exactly ONE action
from the allowed actions listed previously.
Output no other text.
```

Prompt 18: Prompt instructing the model to directly choose the next action based on previous observations.

### Plan Every K Prompt

```
Review your previous observations and plan, then make a high-level
plan for completing the task. Your plan can include reasoning about
how to solve the task.
After this planning phase, you will be asked to take actions one at a
time.
```

```
Output your plan strictly in the following format:
```

```
<plan>YOUR_PLAN</plan>
```

```
Replace YOUR_PLAN with your own thinking and plan.
```

```
After your plan, choose exactly ONE action from the allowed actions
listed previously.
Output no other text.
```

Prompt 19: Prompt used in fixed-frequency planning; elicits a plan at regular intervals.

## A.3 Dynamic Planning Prompts

Finally, Prompt 20 is used to replace all the aforementioned instruction prompts when creating the SFT dataset. This ensures the agent observes that the same prompt can elicit different planning strategies, with the goal that RL fine-tuning will enable it to learn to choose the appropriate type of plan. The adaptive nature of this prompt allows the model to autonomously decide when to create or update plans, facilitating flexible and context-aware decision-making.

### Dynamic Planning Prompt

Review your current plan and observations.

- If you do not have a plan yet, create one.
- If your plan is outdated or needs changes, create a new plan.

If you create a new plan, output it in the following format:

```
<plan>YOUR_NEW_PLAN</plan>
```

Replace YOUR\_NEW\_PLAN with your revised plan.

If your current plan is still valid, proceed without outputting it again.

After this evaluation (and any necessary replanning), output exactly ONE allowed action.

Output nothing else except an optional `<plan>...</plan>` block and that single action.

Prompt 20: Prompt allowing the model to decide when to plan based on task context.

## B Appendix

### B.1 Crafter

Crafter Hafner [2022] is an open-world survival game inspired by Minecraft. The environment is designed to benchmark reinforcement learning agents on tasks requiring generalization and long-term reasoning without extensive prior game knowledge. Crafter presents a procedurally generated 2D world where the agent must gather resources, craft tools, and defend against creatures to survive and unlock achievements. The environment provides observations as a combination of top-down pixel views of the agent’s local surroundings and non-visual data, including the agent’s health, inventory, and currently held item. The agent can perform 17 discrete actions, and its progress is measured by unlocking 22 predefined achievements, each yielding a sparse reward signal.

Our agent interacts with the Crafter environment through the BALROG Paglieri et al. [2025a] framework, which acts as a standardized wrapper. BALROG facilitates the communication between our agent, conceptualized as an underlying Large Language Model (LLM) combined with a specific prompting strategy, and the Crafter game environment. At each timestep, Crafter, via BALROG, relays the current observation to our agent. The agent, which internally maintains a history of past observations and actions, incorporates this new observation into its context. A dedicated prompt builder component within the agent updates this interaction history and formats it into a chat template.

This prompt is then passed to the LLM, which processes the contextual information and generates the subsequent action as a natural language string. BALROG then translates this string into a game-compatible command for Crafter. Further details regarding the specifics of the system prompt, the wrapper’s mechanics, and the handling of action validation can be found in BALROG Paglieri et al. [2025a].

### B.2 POGS

The Partially-Observable Graph Search (POGS) environment is a custom-designed, synthetic environment created specifically to isolate and evaluate planning under uncertainty. In POGS, agents navigate procedurally generated graphs, seeing only a limited area around their current node, and must find a path to a target node. This partial observability means agents often discover new sections of the graph or hit dead ends, forcing them to backtrack and adjust their plans. A key feature of POGS is its ability to quantify exploration efficiency by tracking backtracking, specifically defined as the number of times an agent visited the node it was on two steps prior; a lower backtracking

count indicates better efficiency in solving the environment. To integrate POGS into the BALROG framework, we implemented a dedicated system prompt (see Prompt 21) that clearly specifies the agent’s objectives, valid actions, and observation format, as demonstrated in the natural language observation example in Figure 3.

```


The system prompt provided to the agent in POGS



You are an AI agent designed to navigate the Partially Observable Graph Search (POGS) environment. Your primary objective is to find and reach a specific target node.



The following are the only valid actions you can take in the game:  
{list(range(env.num_nodes))}



In a moment I will present you with an observation containing:  
- Adjacency list showing the neighbors of all currently visible nodes  
- Your current node position  
- The target node you need to reach



The graph has {env.num_nodes} nodes and is partially observable, meaning you can only see connections within a k-nearest neighbor radius of your current position. In this episode k={env.k_nearest}.



Your action should be a single integer representing the label of the node you want to travel to. This node must be directly connected to your current node.



PLAY


```

Prompt 21: The system prompt for POGS. This prompt guides the AI agent by defining its objective (to find and reach a specific target node in the POGS environment), outlining the valid actions (moving to an adjacent node), and detailing the structure of observations (adjacency list of visible nodes, current position, and target node).

### B.3 Backtracking on POGS

Figure 6 visualizes additional metrics for POGS: success rate, length, backtrack count, and the number of output tokens. Increasing planning frequency leads not only to higher costs (more output tokens) but also to a higher backtrack count, with the 'always-plan' agent performing the most backtracking. Notably, this excessive backtracking correlates with a reduced success rate, as 'always-plan' agents exhibit lower performance compared to those planning less frequently. This suggests that planning too often causes the agent to continually change its mind, ultimately hindering its ability to efficiently navigate the graph and reach the target.

In contrast, agents that never plan also exhibit a high backtrack count, underscoring the value of having a plan. Planning appears to make the agent’s behaviour more consistent and less erratic, which in turn leads to higher success rates. Overall, these results highlight the importance of balanced planning: planning too frequently or not at all hinders performance, while moderate planning improves efficiency and success.

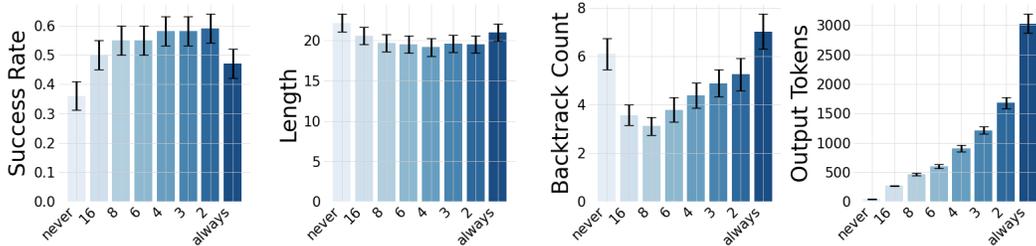


Figure 6: **Planning frequency affects exploration in POGS.** Intermediate planning frequency yields higher success rates and efficiency, while 'always-plan' agents show increased backtracking ( $C_{noise}$ ).

## C Appendix

### C.1 Supervised Fine-Tuning

For the Supervised Fine-Tuning (SFT) phase, we generated a dataset of 1,024 trajectories using Llama 3.3 70B instruct. To ensure data diversity, we employed 16 distinct planning prompts (detailed in Appendix A), which were sampled uniformly. The planning frequency for trajectories in the dataset was also sampled uniformly from the range [2, 12]. The Llama 3.1 8B instruct model served as the baseline for all SFT experiments, with no Low-Rank Adaptation (LoRA) adapters applied during this stage [Hu et al., 2022]. We used AdamW [Loshchilov and Hutter, 2019] as the optimizer throughout.

We considered two objective functions: direct action prediction and full world modeling. Our findings, presented in Figures 7 and 8, indicate that omitting world modeling yields superior results in terms of both performance and reduced catastrophic forgetting. While we initially tested world modeling as a potential regularization technique to mitigate overfitting, it ultimately proved to be a distractor for the agents. To additionally model the environment dynamics, the weights of the model had to be changed more significantly, which is reflected in the higher KL divergence, and consequently led to increased forgetting and reduced model steerability. It is worth noting that for the agent to plan dynamically during SFT, it must also model the plans themselves, exposing it to a greater number of tokens during this stage. Training was performed using the DeepSpeed ZeRO Stage 3 optimizer [Aminabadi et al., 2022] to effectively manage memory and scale operations across multiple GPUs.

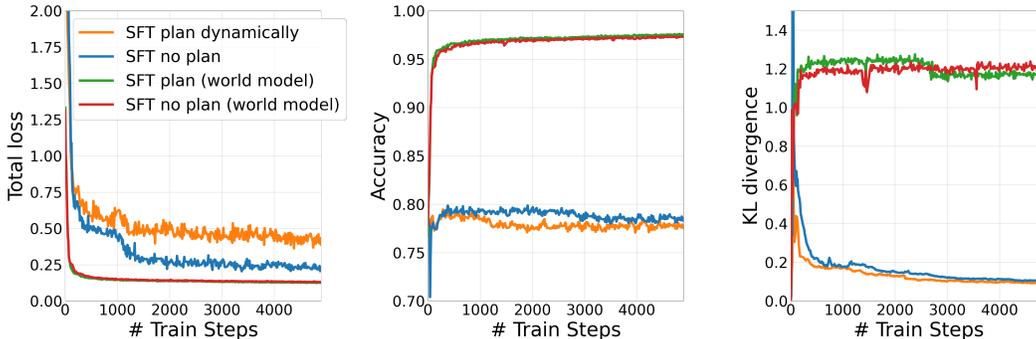


Figure 7: **SFT Training Metrics.** Comparison of (left) total loss, (center) accuracy, and (right) KL divergence across training steps for four SFT configurations. Models incorporating world modeling (green, red) exhibit lower training loss and higher accuracy but also show higher KL divergence. Among the configurations without world modeling, 'SFT plan dynamically' (orange) demonstrates highest total loss as modeling plans is much harder than modeling the world or just actions.

The SFT hyperparameters are shown in Table 1.

### C.2 RLFT

During the Reinforcement Learning Fine-Tuning (RLFT) phase, we applied RSLoRA adapters [Kalajdzievski, 2023] separately to actor and critic models. Data collection involved using vLLM [Kwon et al., 2023] to host model weights and BALROG [Paglieri et al., 2025a] integrated with Ray [Moritz et al., 2018] to process outputs and handle agent-environment interactions efficiently. A

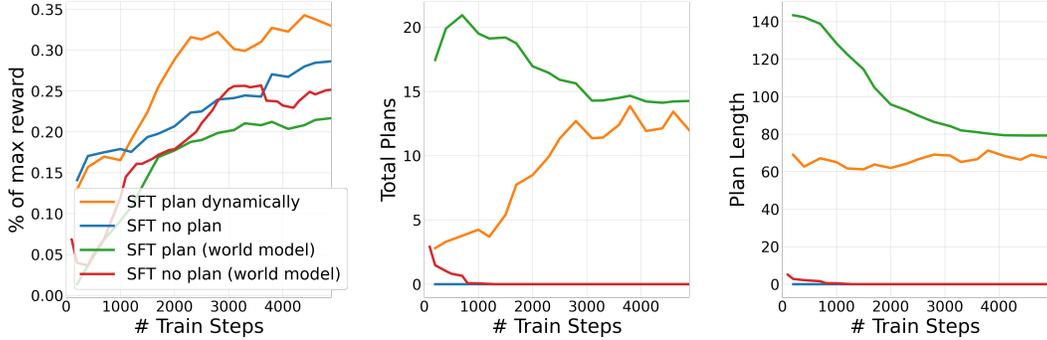


Figure 8: **SFT Evaluation Metrics.** Comparison of (left) normalized game score, (center) total number of plans generated, and (right) average plan length across training steps for the same four SFT configurations as in Figure 7. The configurations without world modeling (orange, blue) consistently achieve better task performance in terms of normalized score than configurations which additionally model the world (green, red).

Table 1: SFT Training Confgs

Data & Model	
Model Size	8B
Max Context Window	8192
LoRA	False
Dataset	1024 trajectories
Optimization	
Learning Rate	5e-6
Beta (KL Coefficient)	0.1
Number of rollouts in batch	384
Rollout length	16
History length	16

custom Proximal Policy Optimization (PPO) trainer [Schulman et al., 2017] was implemented, with model weights broadcasted back to vLLM after each training epoch. The actor model was based on the 8B parameter Llama 3.1 architecture, while the critic model utilized the 1B parameter Llama 3.3 architecture [Grattafiori et al., 2024]. Similar to the SFT phase, the DeepSpeed ZeRO Stage 3 optimizer was employed for memory-efficient and scalable training.

Experiments were conducted using a total of 8 GPUs, allocating 6 GPUs for training and 2 for data collection. Most experiments were executed on a node with 8xH100 GPUs and typically completed within 24–48 hours. Initially, we experimented with reward penalties targeting invalid actions, excessively long responses, and overly frequent planning. However, such explicit reward shaping often resulted in agents refraining from planning entirely. Recognizing optimal planning frequencies ("Goldilocks zones"), such as planning every four steps rather than at every opportunity, we removed explicit reward shaping, thus allowing agents to autonomously learn optimal planning frequencies. As illustrated in Figure 9, agents trained under this strategy gradually improved their efficiency. Specifically, agents learned to execute plans less frequently (center panel) but with increased effectiveness, leading to shorter, more concise plans (right panel), and improved overall performance (left panel).

Detailed hyperparameters for the RL experiments are provided in Table 2.

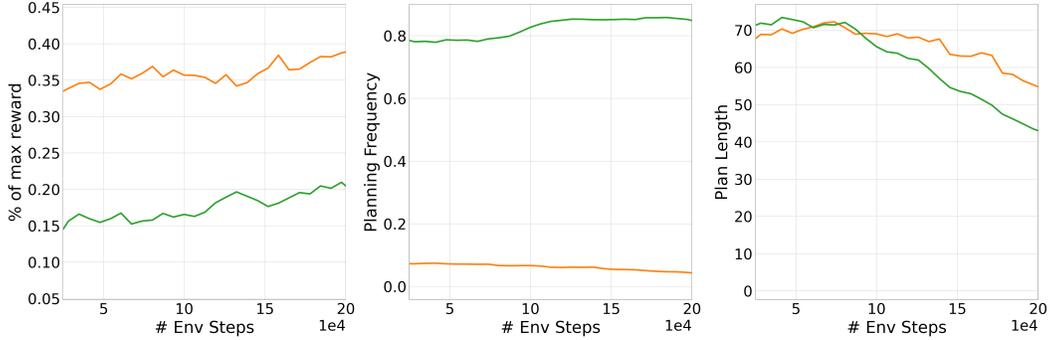


Figure 9: Comparison of (left) Normalized Score, (center) Planning Frequency, and (right) Plan Length across environment steps for two RL configurations. Both agents become more efficient with time; they learn to execute plans for longer, reflected in reduced planning frequency, and also generate more concise plans (reduced plan length).

Table 2: RL Training Configs

Data & Model	
Actor Size	8B
Critic Size	1B
Max Context Window	8192
Max output tokens	200
LoRA	True
LoRA	
R	128
Alpha	64
Dropout	0.0
Optimization	
Actor Learning Rate	5e-6
Critic Learning Rate	5e-6
Optimizer	AdamW
KL Coefficient	0.05
Number of rollouts in batch	192
Rollout length	16
History length	16
PPO Epochs	1
PPO policy clip	0.2
PPO value clip	1.0
Gamma	1.0
GAE	0.95
Value coefficient	0.1
Rollout Temperature	1.0

## D Appendix

### D.1 Qualitative results

In this section, we present further qualitative examples to illustrate the capabilities and limitations of our SFT+RL dynamically planning agent.

First, Figure 10 showcases successful human-agent collaboration, detailing how human-provided high-level plans guided the agent to a game-winning Crafter trajectory after approximately 20

attempts. Next, we demonstrate its autonomous dynamic planning. Figure 11 shows the agent interrupting an ongoing task and adaptively replanning to acquire critical food supplies when its health is low. Figure 12 further illustrates its dynamic planning capabilities, observing the agent employing multi-stage tactics: initially planning to craft a weapon, then devising a new plan to strategically position itself before engaging enemies. Finally, Figure 13 presents an execution failure, showing how an otherwise sound plan to craft an item falters because the agent fails to verify all necessary prerequisites—specifically, by not ensuring a required furnace is accessible. This is akin to the knowing-doing gap identified in BALROG [Paglieri et al., 2025a].

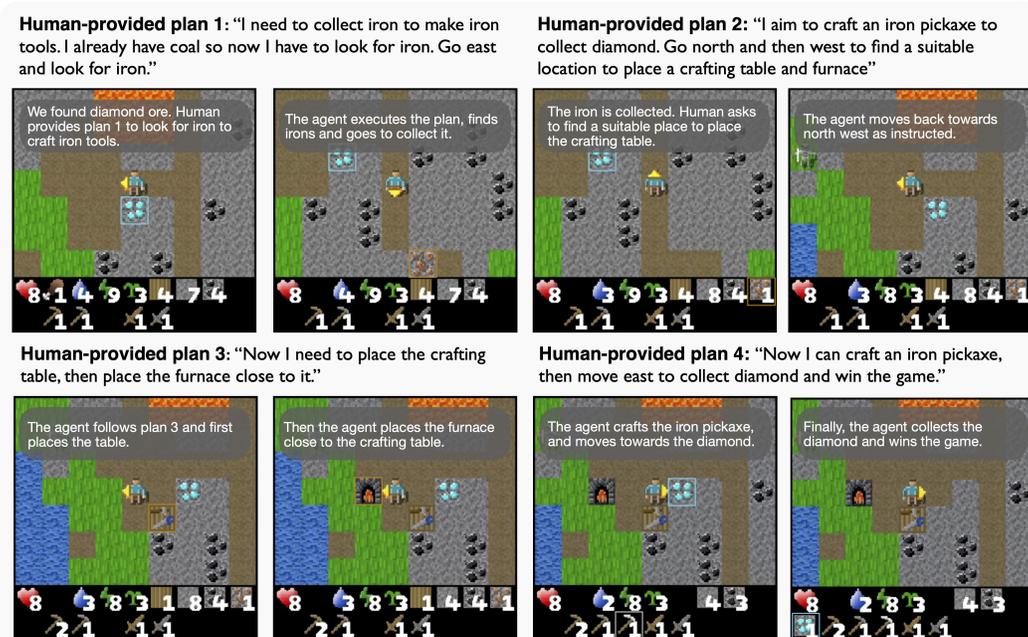


Figure 10: **Human-Agent Collaboration in Crafter.** This figure illustrates a successful human-agent collaboration, where human-provided plans guided the RL-trained planning agent to complete the game by mining diamond.

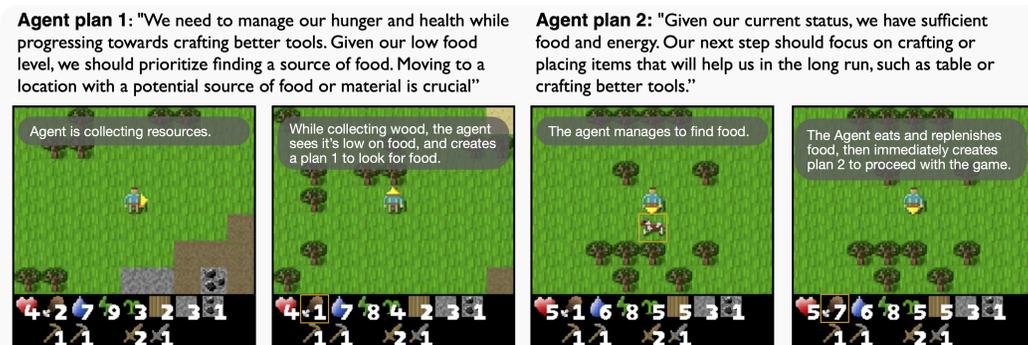


Figure 11: **Autonomous Replanning for Survival.** The SFT+RL plan dynamically agent demonstrates adaptive behavior by interrupting its current objective to address a critical need. It formulates a plan to acquire food when low on health, and upon replenishing, generates a new plan to resume game progression.

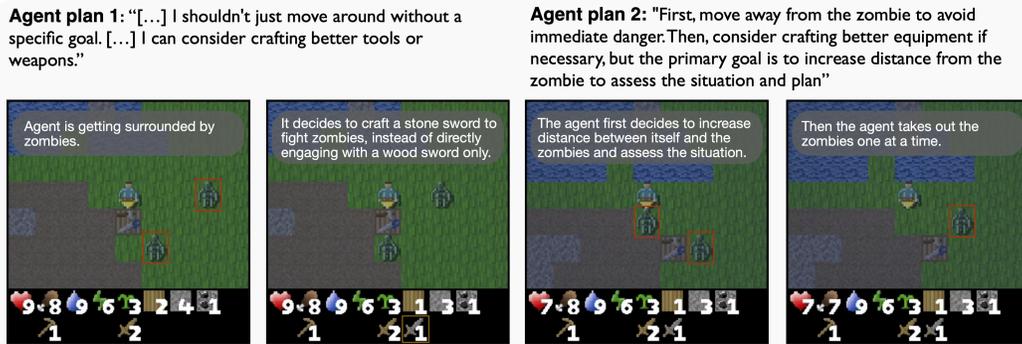


Figure 12: **Autonomous Multi-Stage Tactical Planning.** The agent showcases its ability to chain plans and adapt its tactics. After initially planning and crafting a stone sword to combat zombies, it reassesses and creates a subsequent plan to first gain distance, demonstrating a more strategic approach before engaging.



Figure 13: **Agent Failure Case.** The agent correctly plans to craft an iron pickaxe, identifying the necessary iron and furnace. However, after successfully collecting iron, it fails by attempting the craft action without ensuring the furnace is accessible, indicating a lapse in verifying all conditions of its plan.

## D.2 RL Cost Penalty Ablation

We experimented with different penalties on the cost of planning ( $C_{tokens}$ ) to analyze how agents adapt their behavior to computational constraints. As illustrated in Figure 14, the addition of these penalties led agents to reduce their planning frequency and plan length over time; however, we observed that the normalized score was largely unaffected regardless of penalty level. This is in contrast to our main findings, where explicit planning consistently enhances agent performance. This divergence suggests that, after sufficient training, agents may increasingly internalize planning behaviors; as they gain proficiency in the environment, much of the reasoning and strategy required for success can become implicit within the policy, reducing reliance on overt, explicit planning actions. This could help explain why further penalizing explicit plan generation produced only limited effects on final scores, even though explicit planning is crucial during earlier stages.

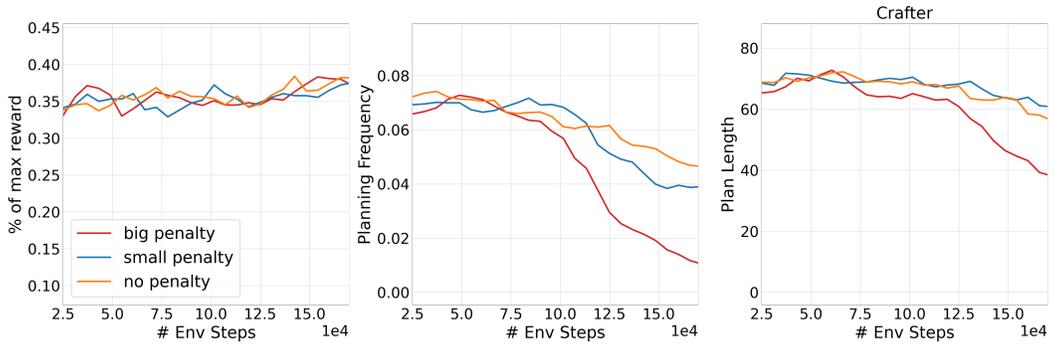


Figure 14: Comparison of training with different planning cost penalties. We compare the (left) Normalized Score, (center) Planning Frequency, and (right) Plan Length for agents trained with no penalty, a small penalty (-0.001 per token), and a big penalty (-0.005 per token).