

# SCOPE: Language Models as One-Time Teacher for Hierarchical Planning in Text Environments

Haoye Lu, Pavan Seshadri, Kaheer Suleman  
{haoye,pavan.seshadri,kaheer}@skyfall.ai

## Abstract

Long-term planning in complex, text-based environments presents significant challenges due to open-ended action spaces, ambiguous observations, and sparse feedback. Recent research suggests that large language models (LLMs) encode rich semantic knowledge about the world, which can be valuable for guiding agents in high-level reasoning and planning across both embodied and purely textual settings. However, existing approaches often depend heavily on querying LLMs during training and inference, making them computationally expensive and difficult to deploy efficiently. In addition, these methods typically employ a pretrained, unaltered LLM whose parameters remain fixed throughout training, providing no opportunity for adaptation to the target task. To address these limitations, we introduce SCOPE (Subgoal-COnditioned Pretraining for Efficient planning), a one-shot hierarchical planner that leverages LLM-generated subgoals only at initialization to pretrain a lightweight student model. Unlike prior approaches that distill LLM knowledge by repeatedly prompting the model to adaptively generate subgoals during training, our method derives subgoals directly from example trajectories. This design removes the need for repeated LLM queries, significantly improving efficiency, though at the cost of reduced explainability and potentially suboptimal subgoals. Despite their suboptimality, our results on the TextCraft environment show that LLM-generated subgoals can still serve as a strong starting point for hierarchical goal decomposition in text-based planning tasks. Compared to the LLM-based hierarchical agent ADaPT (Prasad et al., 2024), which achieves a 0.52 success rate, our method reaches 0.56 and reduces inference time from 164.4 seconds to just 3.0 seconds.

## 1 Introduction

Developing agents capable of planning and reasoning over long horizons remains a central challenge in reinforcement learning (RL). Classical RL and planning methods have achieved notable success in environments with compact state representations and finite action spaces (Bercher et al., 2019; Pateria et al., 2021). However, such favourable conditions typically rely on substantial domain expertise to design effective state abstractions, carefully shaped rewards, and constrained action sets (Ibrahim et al., 2024). Extending these successes to more general and unstructured settings is considerably more difficult, as many real-world problems lack such engineered structure, leading to persistent challenges in exploration, credit assignment, and generalization.

Motivated by the empirical success of large language models (LLMs) and their strong capability for general-purpose text processing, recent RL research has begun leveraging their semantic knowledge to support high-level reasoning and planning in embodied agents. This approach has shown strong promise in both robotic control (Ahn et al., 2022) and long-term planning in open-world environments such as Minecraft (Wang et al., 2024). However, because these methods rely heavily on querying LLMs during both training and inference, they incur high computational costs and are difficult to deploy efficiently. Furthermore, although some approaches introduce additional trainable components to give the LLM-based planner limited adaptability, the LLM itself is typically kept frozen, preventing its parameters from being updated to better fit the target task.

To address the high computational cost and limited flexibility of LLM-based planners, Zhou et al. (2024); Li et al. (2026) propose an alternative approach in which the LLM serves as a teacher to supervise a small student network for high-level decision making. Early in training, the student closely follows the LLM’s guidance, but its reliance gradually diminishes, eventually allowing it to operate independently at inference time. Although this approach removes the need for the LLM during deployment, it still requires iterative LLM queries throughout training to provide ongoing guidance to the student.

Adopting a similar idea, in this work, we propose SCOPE (Subgoal-COnditioned Pretraining for Efficient planning), a more efficient approach that uses LLM-generated subgoals only once at initialization, assuming access to a set of suboptimal example trajectories. Unlike prior methods that repeatedly query the LLM to generate subgoals during training adaptively, SCOPE directly extracts subgoal sequences from demonstrations and uses them to pretrain a student planner. Although these subgoals are suboptimal as the LLM does not interact with the environment, our preliminary results on TextCraft (a simplified, text-only version of Minecraft) show that they still provide a strong starting point for the planner in text-based planning tasks. After fine-tuning the planner by maximizing expected return through interaction with a world model, SCOPE achieves strong performance. In particular, compared to the LLM-based hierarchical agent ADaPT (Prasad et al., 2024), which obtains a 0.52 success rate on TextCraft, SCOPE reaches 0.56. It is also significantly more efficient at inference time: SCOPE completes the game in an average of 3.0 seconds on a single NVIDIA A10 GPU, whereas ADaPT requires 164.4 seconds using a GPT-3.5 backend accessed via the OpenAI API under ideal network conditions.

## 2 Related Work

The hierarchical paradigm for planning and executing long-term tasks has been extensively explored over the past several decades (Bercher et al., 2019; Pateria et al., 2021). The central idea is to decompose complex, long-horizon decision-making problems into a hierarchy of simpler subtasks, where a high-level policy operates over abstract actions, selecting subtasks rather than primitive controls, to achieve the overall objective. Each subtask can itself be formulated as a reinforcement learning problem, solved by a lower-level policy that learns to accomplish it (Hengst, 2010). The interaction between these hierarchical policies collectively governs the agent’s overall behaviour.

By representing long-horizon problems in terms of temporally extended subtasks, hierarchical reinforcement learning and planning effectively shortens the decision horizon, a principle known as temporal abstraction (Barto & Mahadevan, 2003; Dietterich, 2000; Sutton et al., 1999). Temporal abstraction enables more efficient credit assignment across extended timescales (Vezhnevets et al., 2017), while decomposing tasks into subtasks typically simplifies learning and promotes more structured exploration during training (Nachum et al., 2019). These advantages make HRL a powerful framework for scaling reinforcement learning to long-horizon domains (Dayan & Hinton, 1992; Vezhnevets et al., 2017; Nachum et al., 2018). Empirical evidence shows that HRL consistently outperforms flat RL methods across diverse settings, including continuous control (Florensa et al., 2017; Levy et al., 2019), strategic games (Vezhnevets et al., 2017; Nachum et al., 2018), and robotic manipulation (Fox et al., 2017; Gupta et al., 2019). Notably, several studies attribute these improvements primarily to enhanced exploration facilitated by subgoal-based hierarchical structures (Jong et al., 2008; Nachum et al., 2019).

Given that large language models (LLMs) encode rich semantic knowledge about the world, such knowledge can be highly valuable for guiding embodied agents to perform high-level reasoning and planning. Motivated by this insight, Ahn et al. (2022) introduced a hierarchical planning framework for robot control, in which a low-level agent executes concrete motor actions, such as opening and closing the gripper, while a high-level, LLM-based agent issues abstract commands like “put the Coke on the counter.” Building on this idea, subsequent studies have extensively explored the use of LLMs as task planners. These approaches can be broadly categorized into two types: comprehensive (Tang et al., 2023; Dalal et al., 2024), where all subgoals are planned in advance, and incremental, where subgoals are generated adaptively, allowing the planner to make dynamic adjustments based on feedback or environmental changes (Ichter et al., 2023; Zhang et al., 2023; Wang et al., 2024).

Due to the vast number of trainable parameters in LLM-based agents, directly optimizing them through fine-tuning is often impractical and computationally demanding (Hu et al., 2024; Schoepp et al., 2025). Consequently, most existing approaches improve the behaviour of LLM agents indirectly, by refining their prompts based on environmental feedback and/or by allowing the LLM to generate multiple high-level action candidates, from which a separate value function selects the most appropriate one (Ahn et al., 2022; Wang et al., 2023). Beyond this common paradigm, Zhou et al. (2024); Li et al. (2026) propose an alternative approach that leverages the LLM as a teacher to guide a smaller student network in producing high-level decisions. Initially, the agent is trained to follow the LLM’s guidance closely; however, as training proceeds, its reliance on the LLM gradually diminishes, enabling it to make independent decisions. This design removes the need for the LLM during inference, allowing the student network to be further fine-tuned in the final training stage by directly maximizing the expected return.

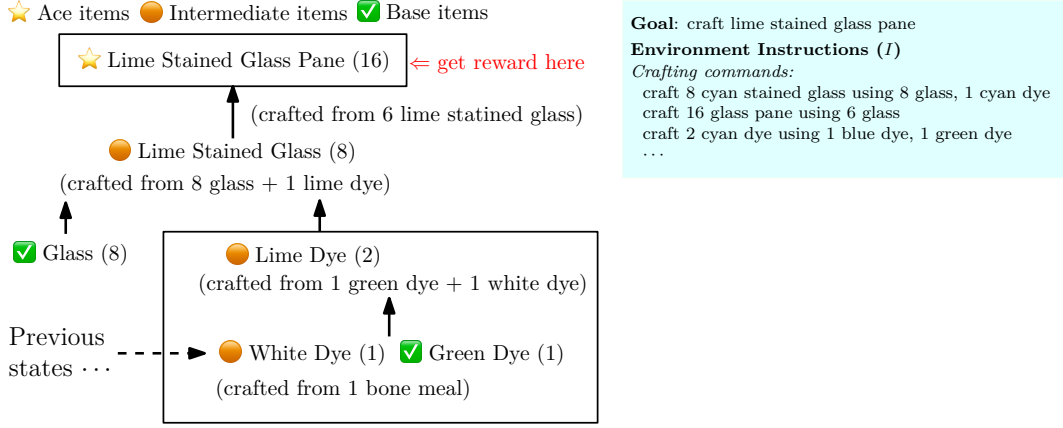


Figure 1: TextCraft Environment. An example crafting dependency chain for producing the ace item (lime stained glass pane). The agent must gather base items, synthesize intermediate items through the provided crafting commands, and execute the sequence in the correct order to obtain the final reward.

### 3 Backgrounds

In this work, we aim to develop an agent capable of solving long-horizon, goal-reaching tasks in a text-based environment, with only one-time guidance from an LLM prior to any agent training. Specifically, we consider a setting in which the system receives a user-provided instruction that includes a goal  $g$  that the agent should accomplish, along with additional descriptions of the environment configuration  $I$ . The agent receives a reward only upon successful completion of the task. We also assume access to a set of successful example goal-instruction-trajectory tuples,  $\mathcal{S} = \{(g_k, I_k, T_k) \mid k = 1, 2, \dots, N\}$ , provided by humans to help the agent learn how to interact with the environment. Each trajectory is defined as  $T_k = (s_0, a_0, s_1, \dots, a_{M_k-1}, s_{M_k})$ , where each  $a_i$  is the action executed when the agent is in state  $s_i$ , thereby forming a sequence of state-action transitions. Although all example trajectories are assumed to eventually reach the final goal, we do not assume they are optimal, as human demonstrations may include suboptimal or inefficient actions.

**Language Models as One-Time Teacher.** To establish an initial foundation for goal decomposition, an LLM-based planner analyzes the demonstration trajectories and heuristically partitions them into subtrajectories, each ending at a state where a subgoal is deemed achieved. We then train a low-level employee agent to accomplish these short-horizon subgoals, while a high-level manager agent is responsible for proposing them.

To ensure consistent goal decomposition across all trajectories, rather than directly asking the LLM planner to output subtrajectories, we provide it with 50 randomly sampled instruction-trajectory pairs and ask it to propose a Python function  $f_{dc}(T)$  that systematically decomposes a trajectory  $T = (s_0, a_0, s_1, \dots, a_{M-1}, s_M)$  into subtrajectories with associated subgoals. This function is paired with a corresponding subgoal completion function  $f_{sg}(s, \tilde{g})$ , which returns 1 if the subgoal  $\tilde{g}$  is achieved in state  $s$ , and 0 otherwise. Although the resulting subgoals may be suboptimal and lack interpretability, our empirical study demonstrates that they nevertheless serve as a strong initialization for subgoal decomposition and provide effective guidance for improving the efficiency of long-term planning training. (The prompts used and the generated Python programs are included in Sec D.) In particular, by applying  $f_{dc}$ , we decompose tuple  $(g, I, T)$  into

$$(\tilde{T}_k, \tilde{g}_k, I) \quad \text{with} \quad \tilde{T}_k = (s_{i_{k-1}}, a_{i_{k-1}}, \dots, s_{i_k}), \quad k = 1, 2, \dots, N_T, \quad (1)$$

where  $i_k$  denotes the index of the state at which subgoal  $\tilde{g}_k$  is achieved (i.e.,  $f_{sg}(s_{i_k}, \tilde{g}_k) = 1$ ), with  $i_0 = 0$ . Here,  $N_T$  denotes the total number of subgoals, and completing the final subgoal  $\tilde{g}_{N_T}$  corresponds to successfully achieving the ultimate goal  $g$ .

For later reference, each subtrajectory in eq. (1) is further decomposed into state-action-subgoal-instruction tuples of the form  $(s_{i_{k-1}}, a_{i_{k-1}}, \tilde{g}_k, I)$ , which are compiled into the dataset  $\Phi = \{(s, a, \tilde{g}, I)\}$ . In addition, for tuple  $(g, I, T) \in \mathcal{S}$ , let

$$\Phi_0^{g, I, T} = \{(s_{i_{k-1}}, \tilde{g}_k, g, I) : k = 1, 2, \dots, N_T\}, \quad (2)$$

where each element of  $\Phi_0^{g,I,T}$  corresponds to the initial state of a subtrajectory induced from  $T$ , paired with its associated subgoal  $\tilde{g}_k$ , ultimate goal  $g$  and instruction  $I$ . We further define

$$\Phi_0 = \bigcup_{(g,I,T) \in \mathcal{S}} \Phi_0^{g,I,T}. \quad (3)$$

**TextCraft Environment.** In this preliminary work, we use the TextCraft environment (Prasad et al., 2024) as the benchmark to evaluate our hierarchical text-based RL framework. TextCraft, inspired by Minecraft, tests an agent’s ability to perform compositional reasoning and long-term planning in crafting tasks. As shown in Fig 1, at the start of each episode, the agent receives a goal  $g$  to craft a specific ice item as well as additional instructions  $I$  containing a set of available crafting commands. The agent is expected to infer which base items are non-craftable and obtainable from the environment, then act through textual commands – **craft** using `<ingredients>` to produce intermediate or final items, and **get** to collect the required base materials. The agent tracks its inventory through the textual description of the current state. An episode terminates once the target item is successfully crafted, at which point the agent receives a reward; otherwise, no reward is provided. Although the environment appears simple, it is intentionally designed so that crafting the final item requires executing a precise sequence of actions in the correct order, including the creation of multiple intermediate items. This makes TextCraft a suitable testbed for evaluating the effectiveness of hierarchical planning algorithms in an open, purely text-based setting.

## 4 SCOPE

**Overview.** Similar to prior long-term planning work in open-world environments (Ahn et al., 2022; Wang et al., 2024), our approach, SCOPE, adopts a hierarchical agent design consisting of an employee agent for low-level execution and a manager agent for high-level planning. The manager agent proposes a high-level plan based on the current environment state and the ultimate goal, and then delegates control to the employee agent. The employee agent interacts directly with the environment to complete the assigned subgoals and returns control to the manager either upon completion or once a preset step limit is reached. This iterative process continues as the manager proposes subsequent subgoals, and terminates once the ultimate goal has been achieved.

The remainder of this section describes how both agents are trained. Each agent is first pretrained on the suboptimal trajectories, followed by a reinforcement learning-based stage to further improve the goal-achievement rate. The RL training relies on interactions with world models. For the employee agent, the elementary world model is obtained by training directly on the provided suboptimal trajectories. The world model used for training the manager agent is then constructed by combining the trained employee agent with the elementary world model; its implementation details will be presented in the manager training section.

### 4.1 Employee Agent

The employee agent is responsible for accomplishing individual subgoals and is trained in two stages: an imitation-based pretraining phase to mimic suboptimal trajectories, followed by RL fine-tuning to maximize subgoal completion rate.

**Pretraining.** For pretraining, the agent  $\pi_\theta^e$  is trained to imitate suboptimal trajectories collected from human players by minimizing

$$J^e(\theta) = - \sum_{(s,a,\tilde{g},I) \in \Phi} \log \pi_\theta^e(a|s,\tilde{g},I). \quad (4)$$

This objective encourages the policy  $\pi_\theta^e$  to reproduce the demonstrated behaviour, allowing it to learn an initial mapping from states and subgoals to actions under the given instruction.

**Employee World Model.** The RL training builds on a world model trained from suboptimal trajectories to predict the next state and check action validity. Given the current state  $s$ , action  $a$ , and instruction  $I$ , the elementary world model  $\text{EWM}(s,a;I)$  first verifies whether  $a$  is valid. In TextCraft, an action is invalid if it cannot be parsed (e.g., semantically incorrect) or executed, for example, when a crafting command is missing from  $I$  or when the inventory lacks the required materials. If valid, the model predicts the next state; otherwise, it returns the current state,

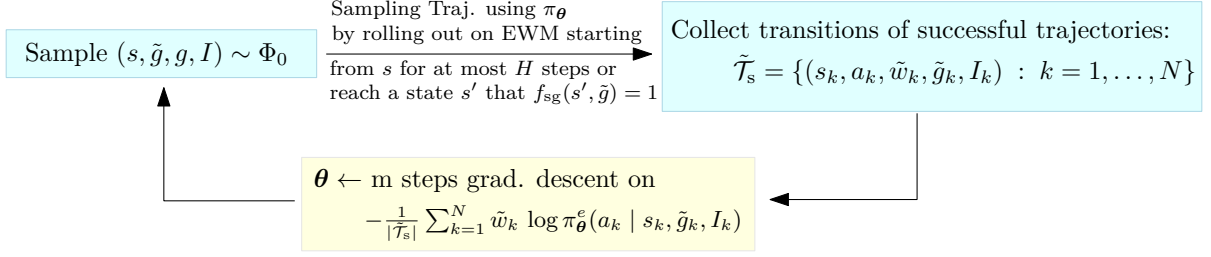


Figure 2: The RL training pipeline of the employee agent.

treating the action as rejected. It outputs both the next state and the validity flag. Since the state in TextCraft is represented as an inventory dictionary, we employ a neural network that takes  $(s, a; I)$  in text form to predict only the action validity, while the next-state update is implemented directly through set operations.

**RL Finetuning.** In this stage, the employee agent is further trained to maximize the expected subgoal achievement rate by interacting with the world model  $\text{EWM}(s, a; I)$ . We use  $f_{sg}$  function provided by the LLM-planner to check if a subgoal is completed. In our implementation, we adopt the Cross-Entropy Method (CEM, Mannor et al. 2003; de Boer et al. 2005), a widely used derivative-free model-based RL (MBRL) algorithm that has demonstrated strong performance across challenging benchmarks (Wang & Ba, 2020; Yang et al., 2020). As illustrated in Fig 2, at each iteration the employee agent samples a random tuple  $(s, \tilde{g}, g, I) \in \Phi_0$  and performs a rollout in the world model starting from  $s$ , which results in a new trajectory by executing a sequence of actions sampled from its policy network, conditioned on the current state, subgoal  $\tilde{g}$  and instruction  $I$ . A total of  $\mathcal{N}_e$  trajectories are sampled per iteration.

If a trajectory reaches a state  $s'$  such that  $f_{sg}(s', \tilde{g}) = 1$  within a preset number of steps, then each state-action pair  $(s, a)$  in that trajectory is assigned the normalized weight  $\tilde{w} = \frac{1}{|\text{traj}|}$ , where  $|\text{traj}|$  denotes the number of transitions; otherwise, the weight is set to zero. This scheme penalizes unnecessarily long trajectories and encourages achieving the subgoal with fewer steps. Collecting all weighted state-action pairs from successful trajectories yields

$$\tilde{\mathcal{T}}_s = \{(s_k, a_k, \tilde{w}_k, \tilde{g}_k, I_k) : k = 1, 2, \dots, \}$$

where each tuple contains a state-action pair together with its assigned weight  $\tilde{w}_k$ , the associated subgoal  $\tilde{g}_k$ , and instruction  $I_k$ . If  $\tilde{\mathcal{T}}_s$  is non-empty, the policy parameters  $\theta$  are then updated by minimizing the weighted negative log-likelihood over  $\tilde{\mathcal{T}}_s$ :

$$-\frac{1}{|\tilde{\mathcal{T}}_s|} \sum_k \tilde{w}_k \log \pi_{\theta}^e(a_k | s_k, \tilde{g}_k, I_k), \quad (5)$$

for a preset number of gradient descent steps, where  $|\tilde{\mathcal{T}}_s|$  denotes the total number of tuples in  $\tilde{\mathcal{T}}_s$ ; otherwise, the update is skipped. The process is repeated by resampling a new set of  $\mathcal{N}_e$  trajectories.

## 4.2 Manager Agent

Unlike the employee agent, which selects actions that directly interact with the environment, the manager agent proposes the next subgoal based on the current state and the ultimate goal, guiding the employee toward achieving that final goal.

**Pretraining.** For pretraining, the manager agent is trained to autoregressively reproduce the subgoal sequences generated by  $f_{dc}(I, T)$  for  $(g, I, T) \in \mathcal{S}$ . This can be achieved by minimizing

$$\mathcal{J}^m(\phi) = - \sum_{(g, I, T) \in \mathcal{S}} \sum_{k=1}^{N_T} \log \pi_{\phi}^m(\tilde{g}_k | s_{i_{k-1}}, g, I). \quad (6)$$

In this way, the manager agent acquires an initial ability to propose subgoals for each subtrajectory, conditioned on its initial state  $s_{i_{k-1}}$ , the ultimate goal  $g$ , and the accompanying instruction  $I$ . This pretrained policy serves as a reasonable starting point, but is not expected to be optimal.

After pretraining, the manager agent is further optimized using RL to maximize the ultimate goal-achievement rate. This stage involves rolling out trajectories within a world model that operates at the subgoal level. We show that this manager-level world model can be constructed by combining the employee world model EWM with the trained employee agent  $\pi_{\theta}^e$ , which the manager needs to coordinate with during planning.

**Manager World Model.** In contrast to the employee world model, the manager world model treats subgoals as actions. Therefore, correspondingly, the manager world model  $\text{MWM}(s, \tilde{g})$  takes the current state  $s$  and a proposed subgoal  $\tilde{g}$  as input and first checks whether  $\tilde{g}$  is achievable within a preset number of steps. If so, it returns the resulting state upon subgoal completion; otherwise, it returns  $s$ . Here, “achievable” means that the subgoal is both plausible and executable by the deployed employee agent. To implement this, we roll out the employee policy  $\pi_{\theta}^e$  on the employee world model EWM starting from  $s$ . If  $\tilde{g}$  is reached within the step limit, the final state is returned; otherwise, the initial state is returned.

To determine whether the ultimate goal is achieved, we can use an ultimate-goal completion function  $f_{\text{ug}}(s, g)$ , which returns 1 if the ultimate goal  $g$  is satisfied in state  $s$  and 0 otherwise. In general,  $f_{\text{ug}}$  can be implemented using a neural network trained via supervised learning on the provided suboptimal trajectories. In TextCraft, however, it can be computed directly by checking whether the inventory contains the ace item.

**RL Finetuning.** In the RL stage, the manager agent is trained in a manner similar to the employee agent, but operates at the subgoal-proposing level by interacting with MWM. In particular, at each iteration, the manager agent samples a random tuple  $(s, \tilde{g}, g, I) \in \Phi_0$  and performs a rollout in MWM starting from  $s$ .<sup>1</sup> The rollout ends when either a preset maximum number of steps is reached or the ultimate goal  $g$  is achieved. This produces a new state-subgoal trajectory  $(s_0, \tilde{g}_0, s_1, \tilde{g}_1, \dots, s_{l-1}, \tilde{g}_{l-1}, s_l)$ , where  $s_0 = s$  and  $l$  is the total number of proposed subgoals. If  $f_{\text{ug}}(s_l, g) = 1$ , the trajectory is considered successful and each state-subgoal pair  $(s_k, a_k)$  for  $k = 0, 1, \dots, l-1$  is assigned a normalized weight  $w = \frac{1}{l}$ . Otherwise, the trajectory is treated as a failure, and all state-subgoal pairs receive weight zero. Similar to the RL finetuning of the employee agent, this scheme penalizes unnecessarily long trajectories and encourages achieving the ultimate goal with a smaller number of subgoals. At each iteration,  $\mathcal{N}_m$  trajectories are generated. Gathering the weighted state-subgoal pairs from all successful trajectories yields

$$\mathcal{T}_s = \{(s_k, \tilde{g}_k, w_k, g_k, I_k) : k = 1, 2, \dots\},$$

where each tuple contains a state-subgoal pair together with its assigned weight  $w_k$ , the associated ultimate goal  $g_k$ , and instruction  $I_k$ . If  $\mathcal{T}_s$  is non-empty, the policy parameters  $\phi$  are then updated by minimizing the weighted negative log-likelihood over  $\mathcal{T}_s$ :

$$-\frac{1}{|\mathcal{T}_s|} \sum_k w_k \log \pi_{\phi}^m(\tilde{g}_k \mid s_k, g_k, I_k), \quad (7)$$

for a specified number of gradient updates, with  $|\mathcal{T}_s|$  indicating the number of tuples in the set. If  $\mathcal{T}_s$  is empty, no update is performed. The procedure then continues by sampling a fresh batch of  $\mathcal{N}_m$  trajectories for the next iteration.

Upon completion of training, the manager and employee agents are combined in a hierarchical manner, as described earlier in this section, and deployed for evaluation in the environment. The effectiveness of this approach is demonstrated in Sec 5.

## 5 Empirical Results

In this section, we conduct a preliminary empirical study of SCOPE using the TextCraft environment.

**Dataset and Network Architecture.** To simulate suboptimal trajectories resembling those of human players, we generate 500K rollouts, reserving 1K each for validation and testing. The objective is to obtain diverse trajectories that generally progress toward the target goal in the TextCraft environment while including occasional suboptimal actions to mimic human-like exploration. Starting from a given state space (defined by available crafting commands) and a target goal item, we construct the recipe tree, derive the optimal crafting sequence, and inject random actions at a 10% rate to introduce variability. The resulting trajectories blend optimal and noisy transitions, producing behaviour that closely resembles human demonstrations.

<sup>1</sup>During RL fine-tuning, subgoals are drawn from  $\pi_{\phi}^m(\cdot \mid s, g, I)$ , and the  $\tilde{g}$  in  $\Phi_0$  is not used.

Table 3: Success rates for crafting ace items in TextCraft. **Top:** ADaPT (Prasad et al., 2024) is a hierarchical agent that uses a GPT-3.5-based planner (Brown et al., 2020). **Bottom:** Ablation results. SCOPE (hand-engineered-subgoal) replaces LLM-generated subgoals with manually constructed, interpretable ones. SCOPE (without manager RL-finetuning) removes the manager-level RL stage.

Setting	Success Rate	# Parameters
ADaPT (Prasad et al., 2024)	0.52	175B
SCOPE ( <b>ours</b> )	0.56	11.04M
SCOPE (hand-engineered-subgoal)	0.58	11.04M
SCOPE (without manager RL-finetuning)	0.24	11.04M

Table 4: Success rates for crafting ace items in TextCraft for ADaPT (Prasad et al., 2024) with different backends. The original implementation uses GPT-3.5 (Brown et al., 2020) as the backend.

Backend	Success Rate	# Parameters	Open Weight?
GPT-4o (OpenAI et al., 2024a)	0.58	1.8T*	No
Mistral Small 3 (Mistral AI, 2025)	0.58	24B	Yes
SCOPE ( <b>ours</b> )	0.56	11.04M	-
GPT-3.5 (Brown et al., 2020)	0.52	175B	No
GPT-4o mini (OpenAI et al., 2024a)	0.43	8B*	No
DeepSeek-R1-Distill-Qwen-32B (DeepSeek-AI et al., 2025)	0.13	32B	Yes
Claude-3 Haiku (Anthropic AI, 2024)	0.00	20B*	No

\* Rough estimates based on publicly discussed or widely believed parameter counts.

Although the employee and manager agents differ in their input-output formats, all signals are converted to text for unified processing, enabling both to share a variational sequence-to-sequence architecture (see Sec A for details). The hyperparameter settings are provided in Sec B.

## 5.1 Performance Comparison

As shown in Table 3, ADaPT (Prasad et al., 2024) is a hierarchical framework that relies on an LLM-based planner with a GPT-3.5 backend (Brown et al., 2020), and it achieves a success rate of 0.52. In comparison, SCOPE attains an ultimate-goal success rate of 0.56. This shows that a one-time LLM initialization, followed by RL fine-tuning, can outperform methods that depend on LLM inference throughout execution. The reduced model size also leads to substantially lower inference cost and latency. In particular, SCOPE agent completes the game in an average of 3.0 seconds on an NVIDIA A10 GPU, whereas ADaPT requires 164.4 seconds when using a GPT-3.5 backend accessed via the OpenAI API under ideal network conditions.

In addition to the original ADaPT implementation based on GPT-3.5 (Brown et al., 2020), we also report results using a broader selection of backends in Table 4. The table compares ADaPT with a range of LLM backends that vary in size and openness. GPT-4o is a large proprietary model developed by OpenAI (2024a), and GPT-4o mini is a smaller, more cost-efficient variant intended to reduce latency and compute while retaining much of GPT-4o’s capability. Claude-3 Haiku (Anthropic AI, 2024) is Anthropic’s lightweight model, designed to answer many requests quickly and reliably under tight latency or budget constraints. On the open-weight side, Mistral Small 3 (Mistral AI, 2025) is a 24B-parameter model in the “small” LLM category (below 70B) that delivers performance comparable to substantially larger models. DeepSeek-R1-Distill-Qwen-32B (DeepSeek-AI et al., 2025) is a mid-sized reasoning model distilled from DeepSeek-R1 based on Qwen2.5 (Qwen et al., 2024). It achieves competitive or superior results to OpenAI o1-mini (OpenAI et al., 2024b) on a range of benchmarks, representing a strong open-weight option in the dense-model regime. While using these more advanced backends improves ADaPT’s performance, SCOPE remains highly competitive: with only 11.04M parameters, it attains a success rate of 0.56, coming very close to the performance of GPT-4o and Mistral Small 3 (both 0.58 with tens of billions to trillions of parameters) and outperforming GPT-4o mini, DeepSeek-R1-Distill-Qwen-32B, and Claude-3 Haiku. This underscores how parameter-efficient SCOPE is compared to both closed and open-weight alternatives.

Provided Demonstration Trajectory	
Goal: craft birch trapdoor.	
Trajectory (action, state):	
craft 6 birch slab using 3 birch planks, {}	
get 2 birch logs, {'birch logs': 2}	
craft 4 birch planks using 1 birch logs, {'birch planks': 4, 'birch logs': 1}	
craft 4 birch planks using 1 birch logs, {'birch planks': 8}	
craft 2 birch trapdoor using 6 birch planks, {'birch trapdoor': 2, 'birch planks': 2}	

LLM-Generated Subgoals	Hand-Engineered Subgoals
{'birch logs': 1}	{'birch planks': 4, 'birch logs': 1}
{'birch planks': 6}	{'birch planks': 8}
{'birch trapdoor': 1}	{'birch trapdoor': 2, 'birch planks': 2}

Figure 5: Comparison of LLM-generated and hand-engineered subgoal decompositions for the demonstration trajectory shown above. Additional samples are provided in Sec D.3.

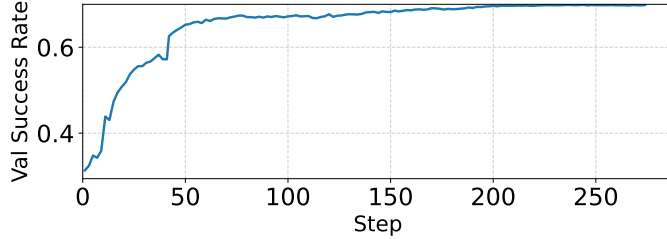


Figure 6: Validation trajectory success rate for the manager agent during RL fine-tuning (Hand-engineered-subgoal). The manager progressively adapts its subgoal proposals to compensate for employee imperfections, yielding steadily improving performance.

## 5.2 Impact of Suboptimal Subgoals on Performance

To better understand how much performance degradation results from suboptimal subgoals, in Fig 5, we compare the subgoals produced by  $f_{dc}$  with those generated through a hand-engineered, more interpretable procedure. In the hand-engineered version, each subgoal is defined as the inventory state immediately after an intermediate item is crafted in a demonstration trajectory. The subgoal is considered achieved once the agent’s inventory contains at least the same items as those target state. By sequentially completing these preset subgoals, the agent eventually crafts the ace item, as the final subgoal contains the ace item to be crafted. Note that while such hand-crafted and interpretable subgoal decomposition is feasible for this simple TextCraft environment, in more complex settings, LLM-generated subgoal decomposition may be the only practical option.

In contrast, the subgoals proposed by  $f_{dc}$  are far less interpretable, and it is unclear how the employee agent should leverage them as effective guidance for crafting the ace item. Interestingly, despite their limited explainability and potential suboptimality, these LLM-generated subgoals still provide sufficient structure for the hierarchical agent to learn and achieve strong performance. In particular, as shown in Table 3, using fully explainable, hand-engineered subgoals improves the success rate by only 2% compared to LLM-generated ones. This suggests that, despite being less interpretable, LLM-generated subgoals still provide sufficiently strong guidance to support effective hierarchical learning.

## 5.3 Agent Imperfections and Mechanisms to Overcome Them

Trajectory failures arise when the agent does not complete the task before exceeding the allotted step limit, and these failures may occur at either the manager or employee level. At the manager level, failures occur when the



chosen subgoals are infeasible under the available command set or are misaligned with the final objective. Moreover, even when a subgoal is valid and achievable, execution may still fail if the employee is unable to complete it due to its own imperfections.

At the employee level, failures typically stem from the agent’s inability to distinguish between elementary items that must be gathered directly from the environment and intermediate items that must be crafted from them, from attempting to craft items that are infeasible given the current inventory, or from issuing invalid or unsupported crafting commands. We further observe that these errors most often arise when the employee attempts to complete a subgoal from an unfamiliar or abnormal inventory state, substantially different from those encountered during training. Such inventory states are themselves a consequence of imperfect execution of earlier subgoals, leading to situations where the agent collects irrelevant items or accumulates an excessive number of resources beyond what is needed. Notably, RL-based finetuning enables the manager to adapt to and compensate for these imperfections effectively.

**Manager Agent Accommodates Employee Imperfections.** To illustrate this, Table 3 (last row) reports a variant of our framework in which the employee is guided directly by a fixed sequence of achievable subgoals, without any manager adaptation. These subgoals are extracted from successful suboptimal trajectories in the validation set, using the same construction as the hand-engineered version (chosen in this study for its high interpretability). In principle, a perfect employee, trained on hand-engineered subgoals, could simply execute these subgoals sequentially to craft the ace item. However, without a manager to adjust the plan, this variant cannot recover when the employee fails to complete a subgoal in the sequence, resulting in a much lower success rate of 0.24. In contrast, the RL-finetuned manager adapts to such failures: when a proposed subgoal is not achieved (due to the employee’s imperfect training), it receives no positive feedback and learns to propose alternative subgoals that the employee can successfully complete. Over time, the manager discovers easier, achievable subgoals that compensate for the employee’s limitations, as also reflected in the steadily increasing success rates shown in Fig 6.

#### Stronger Employee Leads to Higher Success Rates.

While the manager can compensate for employee imperfections through RL finetuning, a stronger employee still substantially improves the final goal-achievement rate. Fig 7 shows the relationship between subgoal success rate and ultimate goal success rate in SCOPE on the validation dataset, with weaker employee variants generated by evaluating partially trained checkpoints. As the figure illustrates, increases in subgoal success rate reliably translate into higher ultimate success, and this improvement accelerates when starting from a stronger employee. That is, higher subgoal success leads to disproportionately larger gains in overall goal completion. The effect is similar to compounding probabilities: the subgoal success rate applies repeatedly across steps, like a value raised to a power. When several subgoals must be achieved in sequence, even a small increase in subgoal reliability quickly compounds, resulting in much larger gains in final success.

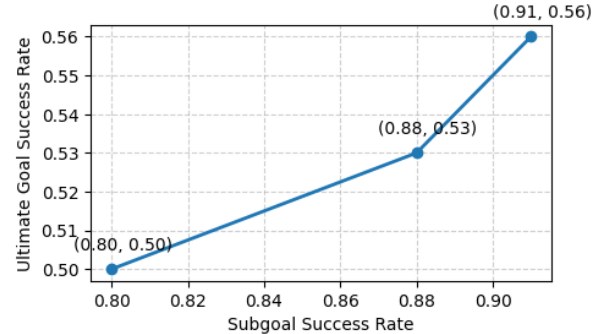


Figure 7: Subgoal vs. ultimate goal success rate in SCOPE.

## 5.4 How Subgoals Help the Agent Complete the Game

While the empirical study demonstrates that one-time LLM-generated subgoals enable SCOPE to surpass the LLM-driven hierarchical agent ADaPT, we now present additional ablations to analyze the source of this improvement more deeply.

**Impact of Subgoal Vagueness and Explainability.** As discussed in Sec 5.2, LLM-generated subgoals are less explainable than hand-engineered ones, which likely contributes to the 2% performance gap observed between SCOPE and the hand-engineered-subgoal variant. To further study how subgoal vagueness and explainability influence SCOPE’s performance, we include two additional variants, as shown in Fig 8a. In SCOPE (no quantity), we remove the item-quantity information from LLM-generated subgoals, and a subgoal is considered satisfied once all listed item types have been collected at least once. Because this variant lacks quantity specifications, its subgoals are more ambiguous than those in standard SCOPE. We additionally include a limiting-case variant, Non-Hierarchical, which uses the same architecture and training setup as the employee agent (pretraining + RL fine-tuning) but is

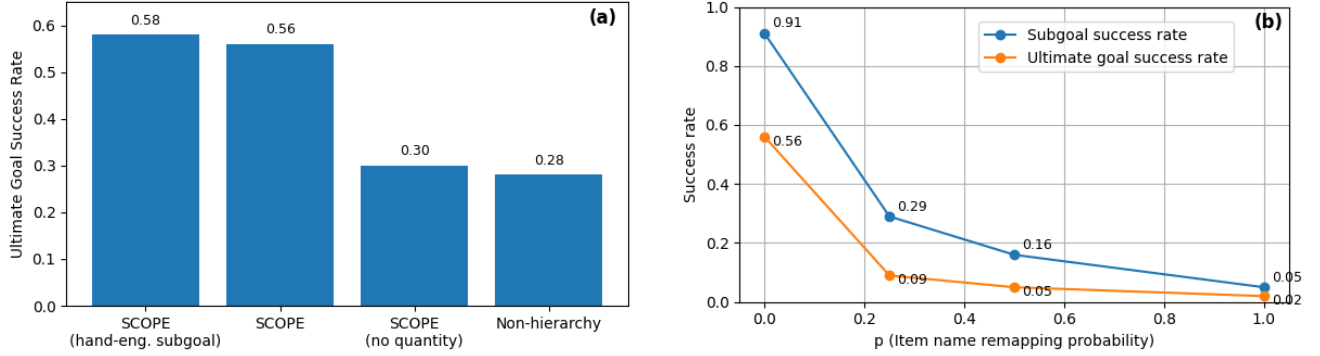


Figure 8: Impact of subgoal quality on agent performance. (a) Performance vs. subgoal explainability (decreasing from left to right). (b) Effect of subgoals when their alignment with the environment is disrupted. Alignment is progressively broken by randomly remapping a proportion  $p$  of item names in the LLM-generated subgoals, with  $p \in \{0.0, 0.25, 0.5, 1.0\}$ .

trained to pursue the final goal directly, without subgoals. In other words, the agent always conditions its decisions on the ultimate goal, rather than on any intermediate objectives. This represents the most vague scenario, as no subgoals are provided at all. As suggested by Fig 8a, less interpretable and more vague subgoals make it harder for SCOPE to extract useful guidance for long-term planning, leading to a clear decline in ultimate goal success as explainability decreases.

**Effect of Decoupling Subgoals from Environment Outcomes.** To test whether subgoals contribute because they genuinely correspond to the environment outcomes, we design a mechanism that gradually breaks this connection by randomly remapping a ratio  $p$  of item names in the LLM-generated subgoals. Specifically, when  $p = 0$ , subgoals are unchanged, which corresponds to the regular SCOPE; as  $p$  increases, more item names are replaced by different ones, making the subgoals increasingly misleading; and when  $p = 1.0$ , all item names are remapped, meaning the subgoals no longer match the items that are actually needed in the environment. The remapping is sampled once and then fixed for both training and testing. We modify only the output of the subgoal-decomposition function, while leaving the subgoal-completion process unchanged, allowing us to degrade subgoal quality without altering the true task objective. This setup helps us understand how performance changes as the alignment between LLM-generated subgoals and environment outcomes is systematically removed, showing the extent to which correct subgoals are causally necessary for SCOPE’s success.

Fig 8b shows how performance changes under different remapping probabilities  $p$ . When  $p = 0$ , subgoals remain unchanged, and both subgoal and ultimate success rates are those of the regular SCOPE. As we increase  $p$  and disrupt the correspondence between subgoals and the environment, performance drops rapidly. At  $p = 0.25$ , success falls to 0.29 (subgoal) and 0.09 (ultimate), and continues to decrease as remapping increases, reaching only 0.05 and 0.02 at  $p = 1.0$ . This sharp decline confirms that SCOPE relies critically on subgoals being causally aligned with the true environmental objectives. When this alignment is lost, both subgoal success and final goal completion collapse. Notably, even with only 25% of item names remapped, the ultimate success rate drops to 0.09 – substantially lower than the non-hierarchical agent’s 0.28, which operates without subgoals at all. This suggests that when subgoals lose alignment with the environment, they do not merely become unhelpful—they can actively mislead the agent. In contrast, earlier results show that vague but aligned subgoals (e.g., using LLM-generated subgoals instead of hand-engineered ones) remain helpful and degrade performance only mildly. It is the loss of alignment—not the loss of specificity—that is particularly harmful, underscoring that subgoals must remain causally grounded to benefit the agent.

## 6 Conclusion

In this work, we introduced SCOPE, an efficient hierarchical planning framework that uses LLM-generated subgoals only once at initialization, derived from suboptimal demonstration trajectories. Instead of repeatedly querying an LLM to adaptively produce subgoals during training, SCOPE directly extracts these imperfect subgoal sequences

and uses them to pretrain a student planner, followed by RL-based refinement using a world model. Although the subgoals are not optimal due to the lack of interaction between the LLM and the environment, our experiments on TextCraft show that they still provide a decent starting point for hierarchical goal decomposition. Empirically, our method outperforms the LLM-based ADaPT system (Prasad et al., 2024), increasing the success rate from 0.52 to 0.56. It also significantly improves efficiency: to complete the game, our agent runs in 3.0 seconds on a single NVIDIA A10 GPU, whereas ADaPT requires 164.4 seconds using a GPT-3.5 backend accessed through the OpenAI API with ideal network conditions. Overall, these results demonstrate that even suboptimal one-time LLM guidance can be highly effective when combined with RL fine-tuning, offering a practical and computationally efficient alternative to LLM-dependent hierarchical planning methods.

## References

- Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, et al. Do as i can, not as i say: Grounding language in robotic affordances. *Science Robotics*, 7(66), 2022. doi: 10.1126/scirobotics.abo0860. URL <https://www.science.org/doi/10.1126/scirobotics.abo0860>.
- Anthropic AI. The claude 3 model family: Opus, sonnet, haiku. Technical report, Anthropic, 2024. URL <https://www.anthropic.com/claude-3-model-card>. Model card / technical report.
- Andrew G. Barto and Sridhar Mahadevan. Recent advances in hierarchical reinforcement learning. *Discrete Event Dynamic Systems*, 13(4):341–379, 2003. doi: 10.1023/A:1025696116075. URL <https://doi.org/10.1023/A:1025696116075>.
- Pascal Bercher, Ron Alford, and Daniel Höller. A survey on hierarchical planning – one abstract idea, many concrete realizations. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence (IJCAI 2019)*, pp. 6267–6275. International Joint Conferences on Artificial Intelligence Organization, 7 2019. doi: 10.24963/ijcai.2019/875. URL <https://doi.org/10.24963/ijcai.2019/875>.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D. Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 1877–1901. Curran Associates, Inc., 2020.
- Murtaza Dalal, Tarun Chiruvolu, Devendra Chaplot, and Ruslan Salakhutdinov. Plan-seq-learn: Language model guided rl for solving long horizon robotics tasks. In *International Conference on Learning Representations*, 2024.
- Peter Dayan and Geoffrey E. Hinton. Feudal reinforcement learning. In *Advances in Neural Information Processing Systems 5 (NIPS 1992)*, pp. 271–278, San Francisco, CA, USA, 1992. Morgan Kaufmann Publishers Inc. ISBN 1558602747.
- Pieter-Tjerk de Boer, Dirk P. Kroese, Shie Mannor, and Reuven Y. Rubinstein. A tutorial on the cross-entropy method. *Annals of Operations Research*, 134(1):19–67, 2005. ISSN 1572-9338. doi: 10.1007/s10479-005-5724-z. URL <https://doi.org/10.1007/s10479-005-5724-z>.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z F Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J L Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R J Chen, R L Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S S Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T Wang, Wangding Zeng, Wanbiao Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W L Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X Q Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y K Li, Y Q Wang, Y X Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y X Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z Z Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu,

- Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. DeepSeek-R1: Incentivizing reasoning capability in LLMs via reinforcement learning. *arXiv preprint arXiv:2501.12948*, January 2025.
- Thomas G. Dietterich. Hierarchical reinforcement learning with the maxq value function decomposition. *Journal of Artificial Intelligence Research*, 13(1):227–303, 2000. ISSN 1076-9757.
- Carlos Florensa, Yan Duan, and Pieter Abbeel. Stochastic neural networks for hierarchical reinforcement learning. In *International Conference on Learning Representations*, 2017. URL <https://openreview.net/forum?id=B1oK8aoxe>.
- Roy Fox, Sanjay Krishnan, Ion Stoica, and Ken Goldberg. Multi-level discovery of deep options, 2017.
- Abhishek Gupta, Vikash Kumar, Corey Lynch, Sergey Levine, and Karol Hausman. Relay policy learning: Solving long horizon tasks via imitation and reinforcement learning. *Conference on Robot Learning (CoRL)*, 2019.
- Bernhard Hengst. Hierarchical reinforcement learning. In Claude Sammut and Geoffrey I. Webb (eds.), *Encyclopedia of Machine Learning*, pp. 495–502. Springer US, Boston, MA, 2010. ISBN 978-0-387-30164-8. doi: 10.1007/978-0-387-30164-8\_363. URL [https://doi.org/10.1007/978-0-387-30164-8\\_363](https://doi.org/10.1007/978-0-387-30164-8_363).
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- Sihao Hu, Tiansheng Huang, Gaowen Liu, Ramana Rao Kompella, Fatih Ilhan, Selim Furkan Tekin, Yichang Xu, Zachary Yahn, and Ling Liu. A survey on large language model-based game agents, 2024.
- Sinan Ibrahim, Mostafa Mostafa, Ali Jnadi, Hadi Salloum, and Pavel Osinenko. Comprehensive overview of reward engineering and shaping in advancing reinforcement learning applications. *IEEE Access*, 12:175473–175500, 2024. doi: 10.1109/ACCESS.2024.3504735.
- Brian Ichter, Anthony Brohan, Yevgen Chebotar, Chelsea Finn, Karol Hausman, Alexander Herzog, Daniel Ho, Julian Ibarz, Alex Irpan, Eric Jang, Ryan Julian, Dmitry Kalashnikov, Sergey Levine, Yao Lu, Carolina Parada, Kanishka Rao, Pierre Sermanet, Alexander T Toshev, Vincent Vanhoucke, Fei Xia, Ted Xiao, Peng Xu, Mengyuan Yan, Noah Brown, Michael Ahn, Omar Cortes, Nicolas Sievers, Clayton Tan, Sichun Xu, Diego Reyes, Jarek Rettinghouse, Jornell Quiambao, Peter Pastor, Linda Luu, Kuang-Huei Lee, Yuheng Kuang, Sally Jesmonth, Nikhil J. Joshi, Kyle Jeffrey, Rosario Jauregui Ruano, Jasmine Hsu, Keerthana Gopalakrishnan, Byron David, Andy Zeng, and Chuyuan Kelly Fu. Do as i can, not as i say: Grounding language in robotic affordances. In Karen Liu, Dana Kulic, and Jeff Ichnowski (eds.), *Proceedings of The 6th Conference on Robot Learning*, volume 205 of *Proceedings of Machine Learning Research*, pp. 287–318. PMLR, 14–18 Dec 2023. URL <https://proceedings.mlr.press/v205/ichter23a.html>.
- Nicholas K. Jong, Todd Hester, and Peter Stone. The utility of temporal abstraction in reinforcement learning. In *Proceedings of the 7th International Joint Conference on Autonomous Agents and Multiagent Systems (AAMAS '08), Volume 1, AAMAS '08*, pp. 299–306, Richland, SC, USA, 2008. International Foundation for Autonomous Agents and Multiagent Systems. ISBN 9780981738109. doi: 10.5555/1402383.1402429.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015. URL <https://arxiv.org/abs/1412.6980>.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *International Conference on Learning Representations*, 2014. URL <https://openreview.net/forum?id=33X9fd2-9FyZd>.
- Andrew Levy, Robert Platt, and Kate Saenko. Hierarchical reinforcement learning with hindsight. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=ryzECoAcY7>.
- Qianxi Li, Bao Pang, Yong Song, Hongze Fu, Qingyang Xu, Xianfeng Yuan, Xiaolong Xu, and Chengjin Zhang. Large language model-assisted hierarchical reinforcement learning training. *Information Sciences*, 723:122688, 2026. ISSN 0020-0255. doi: 10.1016/j.ins.2025.122688. URL <https://www.sciencedirect.com/science/article/pii/S0020025525008217>.
- Shie Mannor, Reuven Rubinstein, and Yoichi Gat. The cross-entropy method for fast policy search. In *Proceedings of the Twentieth International Conference on Machine Learning (ICML 2003)*, International Conference on Machine Learning, pp. 512–519, Washington, DC, USA, 2003. Association for the Advancement of Artificial Intelligence. ISBN 1-57735-189-4.

- Mistral AI. Mistral small 3. <https://mistral.ai/news/mistral-small-3>, 2025. Technical blog post describing the Mistral Small 3 model.
- Ofir Nachum, Shixiang (Shane) Gu, Honglak Lee, and Sergey Levine. Data-efficient hierarchical reinforcement learning. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- Ofir Nachum, Haoran Tang, Xingyu Lu, Shixiang Gu, Honglak Lee, and Sergey Levine. Why does hierarchy (sometimes) work so well in reinforcement learning?, 2019. URL <https://arxiv.org/abs/1909.10618>.
- OpenAI, Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, A J Ostrom, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Mądry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoochian, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein, Andrew Cann, Andrew Codisputi, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia, Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben Leimberger, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake Samic, Bob McGrew, Bobby Spero, Bogo Gierltler, Bowen Cheng, Brad Lightcap, Brandon Walkin, Brendan Quinn, Brian Guarraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo Lugaresi, Carroll Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson, Chak Li, Chan Jun Shern, Channing Conger, Charlotte Barette, Chelsea Voss, Chen Ding, Cheng Lu, Chong Zhang, Chris Beaumont, Chris Hallacy, Chris Koch, Christian Gibson, Christina Kim, Christine Choi, Christine McLeavey, Christopher Hesse, Claudia Fischer, Clemens Winter, Coley Czarnecki, Colin Jarvis, Colin Wei, Constantin Koumouzelis, Dane Sherburn, Daniel Kappler, Daniel Levin, Daniel Levy, David Carr, David Farhi, David Mely, David Robinson, David Sasaki, Denny Jin, Dev Valladares, Dimitris Tsipras, Doug Li, Duc Phong Nguyen, Duncan Findlay, Edede Oiwoh, Edmund Wong, Ehsan Asdar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow, Eric Kramer, Eric Peterson, Eric Sigler, Eric Wallace, Eugene Brevdo, Evan Mays, Farzad Khorasani, Felipe Petroski Such, Filippo Raso, Francis Zhang, Fred von Lohmann, Freddie Sulit, Gabriel Goh, Gene Oden, Geoff Salmon, Giulio Starace, Greg Brockman, Hadi Salman, Haiming Bao, Haitang Hu, Hannah Wong, Haoyu Wang, Heather Schmidt, Heather Whitney, Heewoo Jun, Hendrik Kirchner, Henrique Ponde de Oliveira Pinto, Hongyu Ren, Huiwen Chang, Hyung Won Chung, Ian Kivlichen, Ian O’Connell, Ian O’Connell, Ian Osband, Ian Silber, Ian Sohl, Ibrahim Okuyucu, Ikai Lan, Ilya Kostrikov, Ilya Sutskever, Ingmar Kanitscheider, Ishaan Gulrajani, Jacob Coxon, Jacob Menick, Jakub Pachocki, James Aung, James Betker, James Crooks, James Lennon, Jamie Kiros, Jan Leike, Jane Park, Jason Kwon, Jason Phang, Jason Teplitz, Jason Wei, Jason Wolfe, Jay Chen, Jeff Harris, Jenia Varavva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joanne Jang, Joaquin Quinonero Candela, Joe Beutler, Joe Landers, Joel Parish, Johannes Heidecke, John Schulman, Jonathan Lachman, Jonathan McKay, Jonathan Uesato, Jonathan Ward, Jong Wook Kim, Joost Huizinga, Jordan Sitkin, Jos Kraaijeveld, Josh Gross, Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao, Joyce Lee, Juntang Zhuang, Justyn Harriman, Kai Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kavın Karthik, Kayla Wood, Kendra Rimbach, Kenny Hsu, Kenny Nguyen, Keren Gu-Lemberg, Kevin Button, Kevin Liu, Kiel Howe, Krithika Muthukumar, Kyle Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lauren Workman, Leher Pathak, Leo Chen, Li Jing, Lia Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka, Lilian Weng, Lindsay McCallum, Lindsey Held, Long Ouyang, Louis Feuvrier, Lu Zhang, Lukas Kondraciuk, Lukasz Kaiser, Luke Hewitt, Luke Metz, Lyric Doshi, Mada Aflak, Maddie Simens, Madelaine Boyd, Madeleine Thompson, Marat Dukhan, Mark Chen, Mark Gray, Mark Hudnall, Marvin Zhang, Marwan Aljube, Mateusz Litwin, Matthew Zeng, Max Johnson, Maya Shetty, Mayank Gupta, Meghan Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao Zhong, Mia Glaese, Mianna Chen, Michael Janner, Michael Lampe, Michael Petrov, Michael Wu, Michele Wang, Michelle Fradin, Michelle Pokrass, Miguel Castro, Miguel Oom Temudo de Castro, Mikhail Pavlov, Miles Brundage, Miles Wang, Minal Khan, Mira Murati, Mo Bavarian, Molly Lin, Murat Yesildal, Nacho Soto, Natalia Gimelshein, Natalie Cone, Natalie Staudacher, Natalie Summers, Natan LaFontaine, Neil Chowdhury, Nick Ryder, Nick Stathas, Nick Turley, Nik Tezak, Niko Felix, Nithanth Kudige, Nitish Keskar, Noah Deutsch, Noel Bundick, Nora Puckett, Ofir Nachum, Ola Okelola, Oleg Boiko, Oleg Murk, Oliver Jaffe, Olivia Watkins, Olivier Godement, Owen Campbell-Moore, Patrick Chao, Paul McMillan, Pavel Belov, Peng Su, Peter Bak, Peter Bakkum, Peter Deng, Peter Dolan, Peter Hoeschele, Peter Welinder, Phil Tillet, Philip Pronin, Philippe Tillet, Prafulla Dhariwal, Qiming Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Rajan Troll, Randall Lin, Rapha Gontijo Lopes, Raul Puri, Reah Miyara, Reimar Leike, Renaud Gaubert, Reza Zamani, Ricky Wang, Rob Donnelly, Rob Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchandani, Romain Huet, Rory Carmichael, Rowan Zellers, Roy Chen,

Ruby Chen, Ruslan Nigmatullin, Ryan Cheu, Saachi Jain, Sam Altman, Sam Schoenholz, Sam Toizer, Samuel Miserendino, Sandhini Agarwal, Sara Culver, Scott Ethersmith, Scott Gray, Sean Grove, Sean Metzger, Shamez Hermani, Shantanu Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shirong Wu, Shuaiqi, Xia, Sonia Phene, Spencer Papay, Srinivas Narayanan, Steve Coffey, Steve Lee, Stewart Hall, Suchir Balaji, Tal Broda, Tal Stramer, Tao Xu, Tarun Gogineni, Taya Christianson, Ted Sanders, Tejal Patwardhan, Thomas Cunningham, Thomas Degry, Thomas Dimson, Thomas Raoux, Thomas Shadwell, Tianhao Zheng, Todd Underwood, Todor Markov, Toki Sherbakov, Tom Rubin, Tom Stasi, Tomer Kaftan, Tristan Heywood, Troy Peterson, Tyce Walters, Tyna Eloundou, Valerie Qi, Veit Moeller, Vinnie Monaco, Vishal Kuo, Vlad Fomenko, Wayne Chang, Weiyi Zheng, Wenda Zhou, Wesam Manassra, Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian, Yongjik Kim, Youlong Cheng, Yu Zhang, Yuchen He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and Yury Malkov. GPT-4o system card. *arXiv preprint arXiv:2410.21276*, October 2024a.

OpenAI, Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko, Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, Ally Bennett, Ananya Kumar, Andre Saraiva, Andrea Vallone, Andrew Duberstein, Andrew Kondrich, Andrey Mishchenko, Andy Applebaum, Angela Jiang, Ashvin Nair, Barret Zoph, Behrooz Ghorbani, Ben Rossen, Benjamin Sokolowsky, Boaz Barak, Bob McGrew, Borys Minaiev, Botao Hao, Bowen Baker, Brandon Houghton, Brandon McKinzie, Brydon Eastman, Camillo Lugaresi, Cary Bassin, Cary Hudson, Chak Ming Li, Charles de Bourcy, Chelsea Voss, Chen Shen, Chong Zhang, Chris Koch, Chris Orsinger, Christopher Hesse, Claudia Fischer, Clive Chan, Dan Roberts, Daniel Kappler, Daniel Levy, Daniel Selsam, David Dohan, David Farhi, David Mely, David Robinson, Dimitris Tsipras, Doug Li, Dragos Oprica, Eben Freeman, Eddie Zhang, Edmund Wong, Elizabeth Proehl, Enoch Cheung, Eric Mitchell, Eric Wallace, Erik Ritter, Evan Mays, Fan Wang, Felipe Petroski Such, Filippo Raso, Florencia Leoni, Foivos Tsimpourlas, Francis Song, Fred von Lohmann, Freddie Sulit, Geoff Salmon, Giambattista Parascandolo, Gildas Chabot, Grace Zhao, Greg Brockman, Guillaume Leclerc, Hadi Salman, Haiming Bao, Hao Sheng, Hart Andrin, Hessam Bagherinezhad, Hongyu Ren, Hunter Lightman, Hyung Won Chung, Ian Kivlichen, Ian O’Connell, Ian Osband, Ignasi Clavera Gilaberte, Ilge Akkaya, Ilya Kostrikov, Ilya Sutskever, Irina Kofman, Jakub Pachocki, James Lennon, Jason Wei, Jean Harb, Jerry Twore, Jiacheng Feng, Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joaquin Quiñonero Candela, Joe Palermo, Joel Parish, Johannes Heidecke, John Hallman, John Rizzo, Jonathan Gordon, Jonathan Uesato, Jonathan Ward, Joost Huizinga, Julie Wang, Kai Chen, Kai Xiao, Karan Singhal, Karina Nguyen, Karl Cobbe, Katy Shi, Kayla Wood, Kendra Rimbach, Keren Gu-Lemberg, Kevin Liu, Kevin Lu, Kevin Stone, Kevin Yu, Lama Ahmad, Lauren Yang, Leo Liu, Leon Maksin, Leyton Ho, Liam Fedus, Lilian Weng, Linden Li, Lindsay McCallum, Lindsey Held, Lorenz Kuhn, Lukas Kondraciuk, Lukasz Kaiser, Luke Metz, Madelaine Boyd, Maja Trebacz, Manas Joglekar, Mark Chen, Marko Tintor, Mason Meyer, Matt Jones, Matt Kaufer, Max Schwarzer, Meghan Shah, Mehmet Yatbaz, Melody Y Guan, Mengyuan Xu, Mengyuan Yan, Mia Glaese, Mianna Chen, Michael Lampe, Michael Malek, Michele Wang, Michelle Fradin, Mike McClay, Mikhail Pavlov, Miles Wang, Mingxuan Wang, Mira Murati, Mo Bavarian, Mostafa Rohaninejad, Nat McAleese, Neil Chowdhury, Neil Chowdhury, Nick Ryder, Nikolas Tezak, Noam Brown, Ofir Nachum, Oleg Boiko, Oleg Murk, Olivia Watkins, Patrick Chao, Paul Ashbourne, Pavel Izmailov, Peter Zhokhov, Rachel Dias, Rahul Arora, Randall Lin, Rapha Gontijo Lopes, Raz Gaon, Reah Miyara, Reimar Leike, Renny Hwang, Rhythm Garg, Robin Brown, Roshan James, Rui Shu, Ryan Cheu, Ryan Greene, Saachi Jain, Sam Altman, Sam Toizer, Sam Toyer, Samuel Miserendino, Sandhini Agarwal, Santiago Hernandez, Sasha Baker, Scott McKinney, Scottie Yan, Shengjia Zhao, Shengli Hu, Shibani Santurkar, Shraman Ray Chaudhuri, Shuyuan Zhang, Siyuan Fu, Spencer Papay, Steph Lin, Suchir Balaji, Suvansh Sanjeev, Szymon Sidor, Tal Broda, Aidan Clark, Tao Wang, Taylor Gordon, Ted Sanders, Tejal Patwardhan, Thibault Sottiaux, Thomas Degry, Thomas Dimson, Tianhao Zheng, Timur Garipov, Tom Stasi, Trapit Bansal, Trevor Creech, Troy Peterson, Tyna Eloundou, Valerie Qi, Vineet Kosaraju, Vinnie Monaco, Vitchyr Pong, Vlad Fomenko, Weiyi Zheng, Wenda Zhou, Wes McCabe, Wojciech Zaremba, Yann Dubois, Yinghai Lu, Yining Chen, Young Cha, Yu Bai, Yuchen He, Yuchen Zhang, Yunyun Wang, Zheng Shao, and Zhuohan Li. OpenAI o1 system card. *arXiv preprint arXiv:2412.16720*, December 2024b.

Shubham Pateria, Budhitama Subagdja, Ah-hwee Tan, and Chai Quek. Hierarchical reinforcement learning: A comprehensive survey. *ACM Comput. Surv.*, 54(5), June 2021. ISSN 0360-0300. doi: 10.1145/3453160. URL <https://doi.org/10.1145/3453160>.

Archiki Prasad, Alexander Koller, Mareike Hartmann, Peter Clark, Ashish Sabharwal, Mohit Bansal, and Tushar Khot. ADaPT: As-needed decomposition and planning with language models. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Findings of the Association for Computational Linguistics: NAACL*

- 2024, pp. 4226–4252, Mexico City, Mexico, June 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-naacl.264. URL <https://aclanthology.org/2024.findings-naacl.264/>.
- Qwen, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, December 2024.
- Sheila Schoepp, Masoud Jafaripour, Yingyue Cao, Tianpei Yang, Fatemeh Abdollahi, Shadan Golestan, Zahin Suifyan, Osmar R. Zaiane, and Matthew E. Taylor. The evolving landscape of llm- and vlm-integrated reinforcement learning, 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):181–211, 1999. ISSN 0004-3702. doi: 10.1016/S0004-3702(99)00052-1. URL <https://www.sciencedirect.com/science/article/pii/S0004370299000521>.
- Yujin Tang, Wenhao Yu, Jie Tan, Heiga Zen, Aleksandra Faust, and Tatsuya Harada. Saytap: Language to quadrupedal locomotion. In Jie Tan, Marc Toussaint, and Kourosh Darvish (eds.), *Proceedings of the 7th Conference on Robot Learning*, volume 229 of *Proceedings of Machine Learning Research*, pp. 3556–3570. PMLR, 06–09 Nov 2023. URL <https://proceedings.mlr.press/v229/tang23a.html>. <https://saytap.github.io>.
- Alexander Sasha Vezhnevets, Simon Osindero, Tom Schaul, Nicolas Heess, Max Jaderberg, David Silver, and Koray Kavukcuoglu. Feudal networks for hierarchical reinforcement learning. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning (ICML 2017)*, volume 70 of *Proceedings of Machine Learning Research*, pp. 3540–3549. PMLR, August 2017. URL <https://proceedings.mlr.press/v70/vezhnevets17a.html>.
- Tingwu Wang and Jimmy Ba. Exploring model-based planning with policy networks. In *International Conference on Learning Representations (ICLR)*, 2020. URL <https://openreview.net/forum?id=H1exf64KwH>.
- Zihao Wang, Shaofei Cai, Guanzhou Chen, Anji Liu, Xiaojian Ma, Yitao Liang, and Team CraftJarvis. Describe, explain, plan and select: interactive planning with large language models enables open-world multi-task agents. In *Proceedings of the 37th International Conference on Neural Information Processing Systems (NeurIPS 2023)*, 2023. URL <https://openreview.net/forum?id=ceI01w0PmT>.
- Zihao Wang, Shaofei Cai, Zhancun Mu, Haowei Lin, Ceyao Zhang, Xuejie Liu, Qing Li, Anji Liu, Xiaojian Ma, and Yitao Liang. OmniJARVIS: Unified vision-language-action tokenization enables open-world instruction following agents. In *Proceedings of the Thirty-eighth Annual Conference on Neural Information Processing Systems (NeurIPS 2024)*, 2024. URL <https://openreview.net/forum?id=ceI01w0PmT>.
- Yuxiang Yang, Ken Caluwaerts, Atil Iscen, Tingnan Zhang, Jie Tan, and Vikas Sindhwani. Data-efficient reinforcement learning for legged robots. In *Proceedings of the Conference on Robot Learning (CoRL 2020)*. PMLR, 2020.
- Xiang Zhang, Junjie Hu, and Qiang Liu. Llm-powered hierarchical language agent for real-time collaboration. In *OpenReview 2023*, 2023. URL <https://openreview.net/forum?id=HLA2023>.
- Zihao Zhou, Bin Hu, Chenyang Zhao, Pu Zhang, and Bin Liu. Large language model as a policy teacher for training reinforcement learning agents. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence (IJCAI '24)*, 2024. ISBN 978-1-956792-04-1. doi: 10.24963/ijcai.2024/627. URL <https://doi.org/10.24963/ijcai.2024/627>.



## A Model architectures

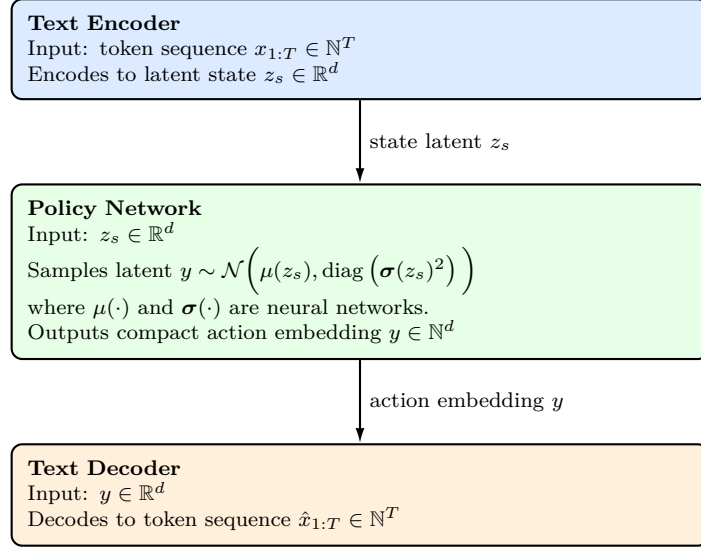


Figure 9: The Network Architecture of the employ and manager agents.

Although the employee and manager agents have different input–output formats, we convert all signals into text to enable unified processing under the full text-based input–output setting. This design allows both agents to share a universal variational sequence-to-sequence architecture. Examples of output are provided in Appx. C.

As illustrated in Fig 9, all networks are implemented using LSTMs (Hochreiter & Schmidhuber, 1997). The LSTM-based text encoder compresses the textual world-state description  $x_{1:T} \in \mathbb{N}^T$  into a latent representation  $z_s \in \mathbb{R}^d$ . Conditioned on  $z_s$ , the policy network produces an action embedding  $y \in \mathbb{R}^d$  via a stochastic latent variable

$$z_a \sim \mathcal{N}(\mu(z_s), \text{diag}(\sigma(z_s)^2)),$$

where  $\mu(\cdot)$  and  $\sigma(\cdot)$  are neural networks, and sampling is performed using the reparameterization trick (Kingma & Welling, 2014). The text decoder then conditions on  $y$  to generate the next action sequence  $\hat{x}_{1:T} \in \mathbb{N}^T$  in natural language form.

## B Hyperparameters

### B.1 Autoencoder

We train an LSTM autoencoder to obtain compact 64-dimensional latent representations of textual state-action sequences in TextCraft. Both the encoder and decoder operate over 256-dimensional token embeddings and employ a single-layer LSTM with a 512-dimensional hidden state. The model is optimized for a maximum of 2000 epochs using Adam (Kingma & Ba, 2015) with a learning rate of  $10^{-3}$  and a batch size of 128. During training, reconstructed sequences are truncated to a maximum length of 30 tokens.

### B.2 World Model

We employ a Transformer encoder operating over 64-dimensional latent states produced by an LSTM autoencoder. The model consists of a 2-layer Transformer with a hidden size of 256, eight attention heads, 1024-dimensional feedforward blocks, and a dropout rate of 0.1, with a [CLS] token used for sequence-level aggregation. Latent encodings are obtained via a single-layer LSTM autoencoder with 256-dimensional embeddings, a 512-dimensional hidden state, and a 64-dimensional bottleneck, trained using teacher forcing. Optimization is performed using Adam

(Kingma & Ba, 2015) with a learning rate of  $3 \times 10^{-4}$  and weight decay of  $1 \times 10^{-5}$ , together with a step-decay scheduler that halves the learning rate every 50 epochs. The world model is trained for 200 epochs with a batch size of 64, conditioning on two past states as temporal context.

### B.3 Policy Agent

Both policy agents take as input the 64-dimensional latent states derived from the pretrained LSTM autoencoder and use a unidirectional LSTM with a 512-dimensional hidden state, 0.1 dropout. The policy is initialized for RL training from a pre-trained checkpoint and trained for a maximum of 2000 epochs with a batch size of 256 for the employee agent, and 16 for the manager agent, rolling out sequences of length 8 and using EMA smoothing ( $\tau = 0.1$ ) together with 10% random exploration. Training incorporates a replay buffer of capacity 10,000 with reciprocal transition weighting, random replacement, and mixed sampling consisting of 90% replayed transitions and 10% fresh rollouts, with replay sampling enabled after 4096 collected transitions. We retain 15% of failed trajectories to improve robustness and perform rollouts every 10 update steps. Optimization uses Adam (Kingma & Ba, 2015) with a learning rate of  $1 \times 10^{-4}$ . Action execution permits one retry per failed step, and all sequences condition on two past states as temporal context.

## C Sample outputs

### Trajectory 1. Final goal: red wool

*Available commands:*

1. craft 1 red wool using 1 red dye, 1 white wool
2. craft 4 magenta dye using 1 blue dye, 2 red dye, 1 white dye
3. craft 1 loom using 2 planks, 2 string
4. craft 1 red dye using 1 poppy
5. craft 1 white wool using 4 string
6. craft 2 red dye using 1 rose bush
7. craft 8 red terracotta using 8 terracotta, 1 red dye
8. craft 1 bow using 3 stick, 3 string
9. craft 1 crossbow using 2 string, 3 stick, 1 iron ingot, 1 tripwire hook
10. craft 1 red dye using 1 red tulip
11. craft 1 green wool using 1 green dye, 1 white wool
12. craft 1 white bed using 3 white wool, 3 planks
13. craft 1 red dye using 1 beetroot
14. craft 1 black wool using 1 black dye, 1 white wool
15. craft 2 lead using 4 string, 1 slime ball
16. craft 1 lime wool using 1 lime dye, 1 white wool

#### Policy inference:

**Manager:** Goal: red dye (1)

Employee: get 2 string

Employee: get 1 beetroot

Employee: craft 1 red dye using 1 beetroot

**Manager:** Goal: red dye (1), white wool (1)

Employee : get 8 string

Employee : craft 1 white wool using 4 string

**Manager:** Goal: red wool (1)

Employee : craft 1 red wool using 1 red dye, 1 white wool

### Trajectory 2. Final goal: spruce slab

*Available commands:*

1. craft 1 spruce button using 1 spruce planks
2. craft 4 spruce planks using 1 spruce logs
3. craft 6 spruce slab using 3 spruce planks
4. craft 3 spruce sign using 6 spruce planks, 1 stick
5. craft 2 spruce trapdoor using 6 spruce planks
6. craft 1 spruce boat using 5 spruce planks
7. craft 1 spruce pressure plate using 2 spruce planks
8. craft 1 spruce fence gate using 4 stick, 2 spruce planks

#### Policy inference:

**Manager :** Goal: spruce logs (1)

Employee : get 1 spruce logs

**Manager :** Goal: spruce planks (3)

Employee : craft 4 spruce planks using 1 spruce logs

Manager : **Goal:** spruce slab (1)  
Employee : **craft** 6 spruce slab using 3 spruce planks

## D LLM Prompt and Output for Subgoal Proposal and Verification

### D.1 Subgoal Decomposition – $f_{dc}$

[User]

#### Definition

A subgoal is defined as a high-level action or decision that, when performed correctly,  
↪ significantly advances progress toward successfully completing the overall task. Subgoals may  
↪ consist of one or more intermediate steps that must be achieved in order to reach the ultimate  
↪ goal.

#### Task Description

1. Induce a consistent high-level plan or strategy—a sequence of subgoals—based on the provided set  
↪ of expert trajectories.
2. Summarize the key steps required to successfully complete the task.
3. The number of subgoals should be smaller than the total number of steps in the trajectory ( $0 < \text{num\_subgoals} < \text{num\_steps}$ ). In other words, not every step is a subgoal.
4. Use consistent methodology in determining these subgoals such that the logic is transferrable.

#### Output Requirements

- Produce a list of subgoals derived from the given trajectories that collectively lead to the  
↪ overall goal.
- Each subgoal must be expressed in the exact format:  
the required state needed to achieve the subgoal
- The ultimate goal must appear as the final subgoal in the list.

Please complete this task systematically, so that we can later design a program to automate the  
↪ process.

#### Input Format

The provided trajectories will have the following structure:

[goal]

[list of (action, outcome\_state) tuples]

The trajectories are below:

0.

Goal: craft birch trapdoor.

craft 6 birch slab using 3 birch planks

get 2 birch logs, {'birch logs': 2}

craft 4 birch planks using 1 birch logs, {'birch planks': 4, 'birch logs': 1}

craft 4 birch planks using 1 birch logs, {'birch planks': 8}

craft 2 birch trapdoor using 6 birch planks, {'birch trapdoor': 2, 'birch planks': 2}

1.

Goal: craft black terracotta.

get 8 terracotta, {'terracotta': 8}

get 1 wither rose, {'terracotta': 8, 'wither rose': 1}

craft 1 black dye using 1 wither rose, {'terracotta': 8, 'black dye': 1}

craft 8 black terracotta using 8 terracotta, 1 black dye, {'black terracotta': 8}

2.

Goal: craft birch button.

get 1 birch logs, {'birch logs': 1}

craft 4 birch planks using 1 birch logs, {'birch planks': 4}

craft 1 birch button using 1 birch planks, {'birch planks': 3, 'birch button': 1}

3.

Goal: craft pink banner.

(get 6 pink tulip ... craft 1 pink banner using 6 pink wool, 1 stick)

4.

Goal: craft crossbow.

```
(get 2 string ... craft 1 crossbow using 2 string, 3 stick, 1 iron ingot, 1 tripwire hook)
```

[GPT-5]

Here's a consistent, recipe-aware decomposition for each trajectory...

```
# 0) Goal: craft birch trapdoor
- {'birch logs': 1}
- {'birch planks': 6}
- {'birch trapdoor': 1}

# 1) Goal: craft black terracotta
- {'terracotta': 8}
- {'black dye': 1}
- {'black terracotta': 8}

# 2) Goal: craft birch button
- {'birch logs': 1}
- {'birch planks': 1}
- {'birch button': 1}

# 3) Goal: craft pink banner
- {'white wool': 6}
- {'pink dye': 6}
- {'pink wool': 6}
- {'stick': 1}
- {'pink banner': 1}

# 4) Goal: craft crossbow
- {'string': 2}
- {'stick': 4}
- {'iron ingot': 2}
- {'warped planks': 1}
- {'tripwire hook': 1}
- {'crossbow': 1}
```

[User]

Good, now create a Python function to programmatically perform this inference

[GPT-5]

```
import re
from collections import defaultdict, deque

CRAFT_RE = re.compile(r"^craft\s+(\d+)\s+([a-z ]+)\s+using\s+(.+)$")
GET_RE = re.compile(r"^get\s+(\d+)\s+([a-z ]+)$")

def _parse_ingredients(txt):
    """
    '8 terracotta, 1 black dye' -> list of (qty:int, item:str)
    """
    parts = [p.strip() for p in txt.split(",")]
    out = []
    for p in parts:
        m = re.match(r"^(\\d+)\s+([a-z ]+)$", p.strip())
        if not m:
            # be tolerant of stray punctuation (e.g., trailing periods)
```

```

        p = p.strip().rstrip(".,;:")
        m = re.match(r"^(\\d+)\\s+([a-z ]+?)$", p)
        if m:
            out.append((int(m.group(1)), m.group(2).strip()))
    return out

def _parse_action(line):
    """
    Returns a dict with keys:
        kind: 'craft' or 'get' or None
        out_qty, out_item, ing: for craft
        qty, item: for get
    """
    line = line.strip()
    # cut off trailing inventory dict if present: ", {...}"
    if ", {" in line:
        line = line[:line.index(", {" )].strip()

    m = CRAFT_RE.match(line)
    if m:
        out_qty = int(m.group(1))
        out_item = m.group(2).strip()
        ing_txt = m.group(3).strip()
        # tolerate a trailing inventory after a comma already removed; also strip trailing periods
        ing_txt = ing_txt.rstrip(".")
        ing = _parse_ingredients(ing_txt)
        return {"kind": "craft", "out_qty": out_qty, "out_item": out_item, "ing": ing}

    m = GET_RE.match(line)
    if m:
        return {"kind": "get", "qty": int(m.group(1)), "item": m.group(2).strip()}

    return {"kind": None}

def infer_subgoals_from_trajectory(goal_item, action_lines):
    """
    Parameters
    -----
    goal_item : str
        e.g., 'birch trapdoor', 'black terracotta'
    action_lines : list[str]
        Lines like:
        'get 8 terracotta, {...}'
        'craft 8 black terracotta using 8 terracotta, 1 black dye, {...}'

    Returns
    -----
    subgoals : list[dict]
        Each element is a required inventory state (e.g., {'terracotta': 8})
        Ordered from base → intermediates → final goal.
    """
    # 1) Parse actions into recipes and note which items can be crafted.
    recipes = {} # product_item → (out_qty, [(qty, item), ...])
    seen_get = set()
    all_actions = []
    for line in action_lines:
        a = _parse_action(line)
        all_actions.append(a)
        if a["kind"] == "craft":
            recipes[a["out_item"]] = (a["out_qty"], a["ing"])
        elif a["kind"] == "get":
            seen_get.add(a["item"])

```

```

# 2) Find the final craft that produces the goal_item (if present).
# If absent, we still assume we want at least 1 goal_item.
if goal_item not in recipes:
    # We can't expand dependencies; just return the final subgoal.
    return [{goal_item: 1}]

# 3) Recursive expansion to compute ingredient requirements for one goal item.
# We will compute a dependency graph and a "required counts" per node.
# For items with recipes, we compute how much of each ingredient is needed
# to yield the requested quantity (ceil div if needed).
import math

# store required quantities for intermediates (not bases)
required = defaultdict(int)
# frontier items to expand: (item, qty_needed)
frontier = deque([(goal_item, 1)])

# Graph for topo order: edges (ingredient -> product)
edges = defaultdict(set)
nodes = set()

# Track per-item total requirement (for all nodes, including base)
total_required = defaultdict(int)
total_required[goal_item] = 1

while frontier:
    item, need_qty = frontier.popleft()
    nodes.add(item)

    if item in recipes:
        out_qty, ing = recipes[item]
        # crafts needed to produce at least need_qty
        crafts_needed = math.ceil(need_qty / out_qty)
        # accumulate requirements for ingredients
        for q, ing_item in ing:
            req = q * crafts_needed
            total_required[ing_item] += req
            edges[ing_item].add(item)
            nodes.add(ing_item)

        # If ingredient itself is craftable, expand further; else it's base.
        if ing_item in recipes:
            # push what *this node* needs (not cumulative to avoid overcounting)
            frontier.append((ing_item, req))

# 4) Identify base vs intermediates (base = no recipe)
base_items = [n for n in nodes if n not in recipes and n != goal_item]
inter_items = [n for n in nodes if n in recipes and n != goal_item]

# 5) Topological order (Kahn) so that base -> intermediates -> goal.
indeg = {n: 0 for n in nodes}
for u in edges:
    for v in edges[u]:
        indeg[v] += 1

q = deque([n for n in nodes if indeg[n] == 0])
topo = []
while q:
    u = q.popleft()
    topo.append(u)
    for v in edges.get(u, []):
        indeg[v] -= 1
        if indeg[v] == 0:

```

```

        q.append(v)

# 6) Build subgoals:
#   - For each base item appearing before any of its consumers, include {item: required_count}
#   - For each intermediate item, include {item: required_count}
#   - End with {goal_item: 1}
subgoals = []
for n in topo:
    if n == goal_item:
        continue
    if n in total_required and total_required[n] > 0:
        subgoals.append({n: total_required[n]})

# Final goal state: have at least 1 of the goal item.
subgoals.append({goal_item: 1})
return subgoals

```

## D.2 Subgoal Completion – $f_{sg}$

[User]

### Definition

A subgoal is defined as a high-level action or decision that, when performed correctly,  
 ↪ significantly advances progress toward successfully completing the overall task. Subgoals may  
 ↪ consist of one or more intermediate steps that must be achieved in order to reach the ultimate  
 ↪ goal.

### Task Description

1. Induce a consistent high-level plan or strategy--a sequence of subgoals--based on the provided  
 ↪ set of expert trajectories.
2. Summarize the key steps required to successfully complete the task.
3. The number of subgoals should be smaller than the total number of steps in the trajectory ( $0 <$   
 ↪  $\text{num\_subgoals} < \text{num\_steps}$ ). In other words, not every step is a subgoal.
4. Use consistent methodology in determining these subgoals such that the logic is transferrable.

### Output Requirements

- Produce a list of subgoals derived from the given trajectories that collectively lead to the  
 ↪ overall goal.
- Each subgoal must be expressed in the exact format:  
 the required state needed to achieve the subgoal
- The ultimate goal must appear as the final subgoal in the list.
- Add a reward function such that, given a current state, returns whether a proposed subgoal has  
 ↪ been achieved or not (binary label)

Please complete this task systematically, so that we can later design a program to automate the  
 ↪ process.

### Input Format

The provided trajectories will have the following structure:

[goal]

[list of (action, outcome\_state) tuples]

The trajectories are below:

0.

Goal: craft birch trapdoor.

craft 6 birch slab using 3 birch planks

get 2 birch logs, {'birch logs': 2}

craft 4 birch planks using 1 birch logs, {'birch planks': 4, 'birch logs': 1}

craft 4 birch planks using 1 birch logs, {'birch planks': 8}

craft 2 birch trapdoor using 6 birch planks, {'birch trapdoor': 2, 'birch planks': 2}

1.

Goal: craft black terracotta.

get 8 terracotta, {'terracotta': 8}

get 1 wither rose, {'terracotta': 8, 'wither rose': 1}



```

craft 1 black dye using 1 wither rose, {'terracotta': 8, 'black dye': 1}
craft 8 black terracotta using 8 terracotta, 1 black dye, {'black terracotta': 8}
2.
Goal: craft birch button.
get 1 birch logs, {'birch logs': 1}
craft 4 birch planks using 1 birch logs, {'birch planks': 4}
craft 1 birch button using 1 birch planks, {'birch planks': 3, 'birch button': 1}
3.
Goal: craft pink banner.
get 6 pink tulip, {'pink tulip': 6}
craft 1 pink dye using 1 pink tulip, {'pink tulip': 5, 'pink dye': 1}
get 3 bone meal, {'pink tulip': 5, 'pink dye': 1, 'bone meal': 3}
craft 4 birch planks using 1 birch logs, {'pink tulip': 5, 'pink dye': 1, 'bone meal': 3, 'birch
↳ planks': 4}
craft 1 pink dye using 1 pink tulip, {'pink tulip': 4, 'pink dye': 2, 'bone meal': 3, 'birch
↳ planks': 4}
craft 1 pink dye using 1 pink tulip, {'pink tulip': 3, 'pink dye': 3, 'bone meal': 3, 'birch
↳ planks': 4}
craft 1 pink dye using 1 pink tulip, {'pink tulip': 2, 'pink dye': 4, 'bone meal': 3, 'birch
↳ planks': 4}
craft 1 pink dye using 1 pink tulip, {'pink tulip': 1, 'pink dye': 5, 'bone meal': 3, 'birch
↳ planks': 4}
craft 1 pink dye using 1 pink tulip, {'pink tulip': 0, 'pink dye': 6, 'bone meal': 3, 'birch
↳ planks': 4}
get 24 string, {'pink dye': 6, 'bone meal': 3, 'birch planks': 4, 'string': 24}
craft 1 white wool using 4 string, {'pink dye': 6, 'bone meal': 3, 'birch planks': 4, 'string': 20,
↳ 'white wool': 1}
craft 1 white wool using 4 string, {'pink dye': 6, 'bone meal': 3, 'birch planks': 4, 'string': 16,
↳ 'white wool': 2}
craft 1 white wool using 4 string, {'pink dye': 6, 'bone meal': 3, 'birch planks': 4, 'string': 12,
↳ 'white wool': 3}
craft 1 white wool using 4 string, {'pink dye': 6, 'bone meal': 3, 'birch planks': 4, 'string': 8,
↳ 'white wool': 4}
craft 1 white wool using 4 string, {'pink dye': 6, 'bone meal': 3, 'birch planks': 4, 'string': 4,
↳ 'white wool': 5}
craft 1 white wool using 4 string, {'pink dye': 6, 'bone meal': 3, 'birch planks': 4, 'string': 0,
↳ 'white wool': 6}
craft 1 pink wool using 1 pink dye, 1 white wool, {'pink dye': 5, 'bone meal': 3, 'birch planks':
↳ 4, 'white wool': 5, 'pink wool': 1}
craft 1 pink wool using 1 pink dye, 1 white wool, {'pink dye': 4, 'bone meal': 3, 'birch planks':
↳ 4, 'white wool': 4, 'pink wool': 2}
craft 1 pink wool using 1 pink dye, 1 white wool, {'pink dye': 3, 'bone meal': 3, 'birch planks':
↳ 4, 'white wool': 3, 'pink wool': 3}
craft 1 pink wool using 1 pink dye, 1 white wool, {'pink dye': 2, 'bone meal': 3, 'birch planks':
↳ 4, 'white wool': 2, 'pink wool': 4}
craft 1 pink wool using 1 pink dye, 1 white wool, {'pink dye': 1, 'bone meal': 3, 'birch planks':
↳ 4, 'white wool': 1, 'pink wool': 5}
craft 1 pink wool using 1 pink dye, 1 white wool, {'pink dye': 0, 'bone meal': 3, 'birch planks':
↳ 4, 'white wool': 0, 'pink wool': 6}
get 2 bamboo, {'bone meal': 3, 'birch planks': 4, 'pink wool': 6, 'bamboo': 2}
craft 1 stick using 2 bamboo, {'bone meal': 3, 'birch planks': 4, 'pink wool': 6, 'bamboo': 0,
↳ 'stick': 1}
craft 1 pink banner using 6 pink wool, 1 stick, {'bone meal': 3, 'birch planks': 4, 'pink banner':
↳ 1}
4.
Goal: craft crossbow.
get 2 string, {'string': 2}
get 6 bamboo, {'string': 2, 'bamboo': 6}
craft 1 stick using 2 bamboo, {'string': 2, 'bamboo': 4, 'stick': 1}
craft 4 spruce planks using 1 spruce logs, {'string': 2, 'bamboo': 4, 'stick': 1, 'spruce planks':
↳ 4}
craft 1 stick using 2 bamboo, {'string': 2, 'bamboo': 2, 'stick': 2, 'spruce planks': 4}
get 4 spruce planks, {'string': 2, 'bamboo': 2, 'stick': 2, 'spruce planks': 8}

```

```

craft 1 stick using 2 bamboo, {'string': 2, 'bamboo': 0, 'stick': 3, 'spruce planks': 8}
get 1 iron ingot, {'string': 2, 'bamboo': 0, 'stick': 3, 'spruce planks': 8, 'iron ingot': 1}
get 2 bamboo, {'string': 2, 'bamboo': 2, 'stick': 3, 'spruce planks': 8, 'iron ingot': 1}
craft 1 stick using 2 bamboo, {'string': 2, 'bamboo': 0, 'stick': 4, 'spruce planks': 8, 'iron
↳ ingot': 1}
get 1 iron ingot, {'string': 2, 'bamboo': 0, 'stick': 4, 'spruce planks': 8, 'iron ingot': 2}
get 1 warped stems, {'string': 2, 'bamboo': 0, 'stick': 4, 'spruce planks': 8, 'iron ingot': 2,
↳ 'warped stems': 1}
craft 4 warped planks using 1 warped stems, {'string': 2, 'bamboo': 0, 'stick': 4, 'spruce planks':
↳ 8, 'iron ingot': 2, 'warped planks': 4}
craft 2 tripwire hook using 1 stick, 1 iron ingot, 1 warped planks, {'string': 2, 'bamboo': 0,
↳ 'stick': 3, 'spruce planks': 8, 'iron ingot': 1, 'warped planks': 3, 'tripwire hook': 2}
craft 1 crossbow using 2 string, 3 stick, 1 iron ingot, 1 tripwire hook, {'spruce planks': 8,
↳ 'warped planks': 3, 'tripwire hook': 1, 'crossbow': 1}

```

## [GPT-5]

Got it -- here's a clean, transferable method plus the concrete subgoals for each trajectory and a  
↳ tiny reward checker you can drop into your pipeline.

# Method (used consistently below)

- 1) Read the final recipe and list its direct ingredients.
- 2) For any ingredient that is itself crafted, compress its internal chain into a single milestone  
↳ state (e.g., "have 6 pink wool"), rather than exposing every micro-step.
- 3) Keep only necessary inventory thresholds that enable the next craft.
- 4) Ensure  $0 < \text{\#subgoals} < \text{\#steps}$  and end with the ultimate goal state ( $\geq 1$  of target item).

---

## 0) Goal: craft birch trapdoor

**\*\*Key steps (summary):\*\*** convert birch logs -> planks; once you have  $\geq 6$  birch planks, craft  
↳ trapdoor.

**\*\*Subgoals (states only, final one is the goal):\*\***

1. {'birch planks': 6}
2. {'birch trapdoor': 1}

---

## 1) Goal: craft black terracotta

**\*\*Key steps (summary):\*\*** obtain terracotta (8); obtain a black dye (e.g., from wither rose); craft  
↳ black terracotta.

**\*\*Subgoals:\*\***

1. {'terracotta': 8}
2. {'black dye': 1}
3. {'black terracotta': 1}

---

## 2) Goal: craft birch button

**\*\*Key steps (summary):\*\*** make birch planks; craft button.

**\*\*Subgoals:\*\***

1. {'birch planks': 1}
2. {'birch button': 1}

---

## 3) Goal: craft pink banner

**\*\*Key steps (summary):\*\*** produce 6 pink wool (via white wool + pink dye), get a stick, craft  
↳ banner.

```

**Subgoals:**
1. {'pink wool': 6}
2. {'stick': 1}
3. {'pink banner': 1}

*(Note: This compresses the internal chains "string -> white wool" and "tulip -> pink dye -> pink
↳ wool" into the single milestone `pink wool: 6`.)*

---

## 4) Goal: craft crossbow
**Key steps (summary):** craft a tripwire hook; ensure remaining ingredients (2 string, 3 sticks, 1
↳ iron ingot) are on hand; craft crossbow.

**Subgoals:**
1. {'tripwire hook': 1}
2. {'string': 2, 'stick': 3, 'iron ingot': 1}
3. {'crossbow': 1}

*(Note: The hook's own chain (stick + iron ingot + planks) is compressed into the milestone ``have
↳ 1 tripwire hook``.)*

---

# Binary Reward Function (subgoal achieved?)
from typing import Dict

def subgoal_achieved(current_state: Dict[str, int], subgoal_state: Dict[str, int]) -> int:
    """
    Returns 1 if current_state satisfies (>=) every item requirement in subgoal_state; else 0.
    Both dicts map item_name -> nonnegative int quantity.
    """
    for item, req_qty in subgoal_state.items():
        if current_state.get(item, 0) < req_qty:
            return 0
    return 1

### Example usage
# Example for trajectory 3 (pink banner)
current = {'pink wool': 6, 'stick': 0}
sg1 = {'pink wool': 6}
sg2 = {'stick': 1}
sg3 = {'pink banner': 1}

assert subgoal_achieved(current, sg1) == 1
assert subgoal_achieved(current, sg2) == 0
assert subgoal_achieved({'pink banner': 1}, sg3) == 1

If you want, I can package these into a JSON schema (trajectories -> subgoals array) or wire the
↳ checker into your rollout collector.

```

## D.3 Comparison of LLM-generated and Hand-engineered Subgoal Decompositions

### D.3.1 Trajectory 1: Crafting Black Terracotta

#### Provided Demonstration Trajectory

Goal: craft black terracotta.

Trajectory (action, state):

get 8 terracotta, {'terracotta': 8}

get 1 wither rose, {'terracotta': 8, 'wither rose': 1}

craft 1 black dye using 1 wither rose, {'terracotta': 8, 'black dye': 1}

craft 8 black terracotta using 8 terracotta, 1 black dye, {'black terracotta': 8}

#### LLM-Generated Subgoals

{'terracotta': 8}

{'black dye': 1}

{'black terracotta': 8}

#### Hand-Engineered Subgoals

{'terracotta': 8, 'black dye': 1}

{'black terracotta': 8}

### D.3.2 Trajectory 2: Crafting Birch Button

#### Provided Demonstration Trajectory

Goal: craft birch button.

Trajectory (action, state):

get 1 birch logs, {'birch logs': 1}

craft 4 birch planks using 1 birch logs, {'birch planks': 4}

craft 1 birch button using 1 birch planks, {'birch planks': 3, 'birch button': 1}

#### LLM-Generated Subgoals

{'birch logs': 1}

{'birch planks': 1}

{'birch button': 1}

#### Hand-Engineered Subgoals

{'birch planks': 4}

{'birch button': 1, 'birch planks': 3}

### D.3.3 Trajectory 3: Crafting Crossbow

#### Provided Demonstration Trajectory

Goal: craft crossbow.

Trajectory (action, state):

```
get 2 string, {'string': 2}
get 6 bamboo, {'string': 2, 'bamboo': 6}
craft 1 stick using 2 bamboo, {'string': 2, 'bamboo': 4, 'stick': 1}
craft 4 spruce planks using 1 spruce logs, {'string': 2, 'bamboo': 4, 'stick': 1}
craft 1 stick using 2 bamboo, {'string': 2, 'bamboo': 2, 'stick': 2}
get 4 spruce planks, {'string': 2, 'bamboo': 2, 'stick': 2}
craft 1 stick using 2 bamboo, {'string': 2, 'stick': 3}
get 1 iron ingot, {'string': 2, 'stick': 3, 'iron ingot': 1}
get 2 bamboo, {'string': 2, 'bamboo': 2, 'stick': 3, 'iron ingot': 1}
craft 1 stick using 2 bamboo, {'string': 2, 'stick': 4, 'iron ingot': 1}
get 1 iron ingot, {'string': 2, 'stick': 4, 'iron ingot': 2}
get 1 warped stems, {'string': 2, 'stick': 4, 'iron ingot': 2, 'warped stems': 1}
craft 4 warped planks using 1 warped stems, {'string': 2, 'stick': 4, 'iron ingot': 2,
'warped planks': 4}
craft 2 tripwire hook using 1 stick, 1 iron ingot, 1 warped planks, {'string': 2,
'stick': 3, 'iron ingot': 1, 'warped planks': 3, 'tripwire hook': 2}
craft 1 crossbow using 2 string, 3 stick, 1 iron ingot, 1 tripwire hook, {'warped
planks': 3, 'tripwire hook': 1, 'crossbow': 1}
```

#### LLM-Generated Subgoals

```
{'string': 2}
{'stick': 4}
{'iron ingot': 2}
{'warped planks': 1}
{'tripwire hook': 1}
{'crossbow': 1}
```

#### Hand-Engineered Subgoals

```
{'string': 2, 'bamboo': 4, 'stick': 1}
{'string': 2, 'bamboo': 2, 'stick': 2}
{'string': 2, 'stick': 3}
{'string': 2, 'iron ingot': 1,
'stick': 4}
{'string': 2, 'stick': 4, 'iron
ingot': 2, 'warped planks': 4}
{'string': 2, 'stick': 3, 'iron
ingot': 1, 'warped planks': 3, 'tripwire
hook': 2}
{'warped planks': 3, 'tripwire hook': 1,
'crossbow': 1}
```