

# Self-Guided Plan Extraction for Instruction-Following Tasks with Goal-Conditional Reinforcement Learning

Anonymous ACL submission

## Abstract

We introduce SuperIgor, a framework for instruction-following tasks. Unlike prior methods that rely on predefined subtasks, SuperIgor enables a language model to generate and refine high-level plans through a self-learning mechanism, reducing the need for manual dataset annotation. Our approach involves iterative co-training: an RL agent is trained to follow the generated plans, while the language model adapts and modifies these plans based on RL feedback and preferences. This creates a feedback loop where both the agent and the planner improve jointly. We validate our framework in environments with rich dynamics and stochasticity. Results show that SuperIgor agents adhere to instructions more strictly than baseline methods, while also demonstrating strong generalization to previously unseen instructions.

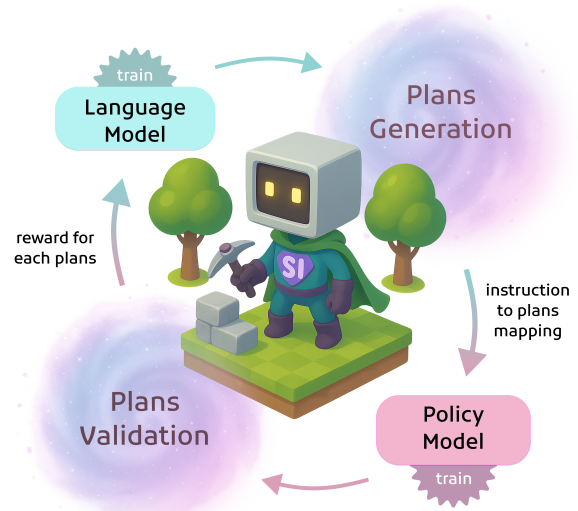


Figure 1: Conceptual diagram of the **SuperIgor** framework designed for Instruction Following

## 1 Introduction

The instruction-following task (Shridhar et al., 2020; Chevalier-Boisvert et al., 2018; Zhong et al., 2021) studies how an AI agent can achieve a goal specified through a natural-language instruction. Prior work in this area is commonly divided into two paradigms: learning from demonstrations and learning from experience. Demonstration-based approaches rely on expert trajectories (e.g., VPT (Baker et al., 2022), STEVE-1 (Lifshitz et al., 2023), Jarvis-VLA (Li et al., 2025)) and exhibit strong generalization to unseen tasks and environments (Fan et al., 2022; Zhou et al., 2024; Gray et al., 2019). However, their scalability is fundamentally constrained by the high cost of collecting large quantities of high-quality demonstrations, as highlighted in SIMA 2 (Bolton et al., 2025).

In contrast, experience-based methods, typically framed as reinforcement learning, learn directly from the agent’s own interactions with the environment without expert supervision (Hill et al., 2020; Mathur and Ahmed, 2025). While concep-

tually more scalable, these methods face substantially greater challenges, most notably the difficulty of grounding natural-language instructions in long-horizon behavior under sparse and delayed rewards (Chevalier-Boisvert et al., 2018; Wang and Narasimhan, 2021; Zhong et al., 2019). To make progress despite these challenges, most experience-based instruction-following methods have been studied in controlled and simplified environments, such as grid-world or cell-based domains, where observations are often symbolic rather than pixel-based, instructions are procedurally generated, and linguistic diversity is limited (Volovikova et al., 2025).

A natural way to strengthen experience-based instruction-following agents is to incorporate high-level planning through language models, which helps interpret complex and underspecified tasks. By leveraging world knowledge encoded in large language models, such planning can reduce the exploration space and accelerate learning (Jansen, 2020; Zhao et al., 2024; Zhou et al.). Existing planning-based approaches typically adopt a com-

mon design choice: planning is performed over a predefined and finite set of subgoals whose completion can be explicitly verified by the environment. This assumption enables generated subgoals to be grounded in an existing skill library via heuristics or similarity-based matching (Logeswaran et al., 2022), provides dense intermediate rewards, and allows automatic skill switching during execution, thereby bypassing the need to learn low-level policies from scratch (Ahn et al., 2022). Moreover, many planning-based systems rely on very large language models (100B+ parameters), which further limits their practicality and scalability.

These limitations raise an important question: how can we learn a low-level policy for instruction following in environments without predefined low-level skills? We address this challenge by proposing SuperIgor, a reinforcement-learning-based framework that integrates high-level planning through large language models. SuperIgor adopts a self-learning strategy that allows the agent to iteratively refine its plans based on its own experience, and to learn effectively from sparse and delayed instruction-level rewards provided only upon successful task completion. We further demonstrate that SuperIgor operates effectively in Crafter, a dynamic, partially observable, and open-ended environment, highlighting its applicability beyond simplified instruction-following benchmarks.

To conclude, our contributions are as follows:

- We propose a new self-supervised training paradigm for the instruction-following task, where high-level plans are generated and refined through interaction between a language model and a reinforcement learning agent—without requiring any manually annotated datasets.
- We introduce a special curriculum to train an RL agent to accurately follow the plan despite sparse reward conditions.
- We implement our approach in the Crafter benchmark and achieve state-of-the-art performance on out-of-distribution tasks, demonstrating the robustness and flexibility of our framework in dynamic and partially observable environments. The dataset and code for SuperIgor are publicly available<sup>1</sup>.

<sup>1</sup><https://anonymous.4open.science/r/SuperIgor-7A4F>

## 2 Related Work

**Instruction Following with RL.** One of the most common strategies for training agents with RL on instruction-following tasks is to jointly encode the instruction and the agent’s observations, enabling alignment between linguistic and perceptual modalities. A prominent line of work relies on shared representation models such as CLIP (Yao et al., 2022), or feature modulation techniques like FiLM layers (Perez et al., 2018; Chevalier-Boisvert et al., 2018; Zhong et al., 2019), to project language information into visual or state representations (Paischer et al., 2023; Wang and Narasimhan, 2021; Zhong et al., 2021). Alternatively, transformer-based architectures process multimodal inputs jointly to improve instruction understanding and execution. This includes embodied language models such as EmBERT (Suglia et al., 2021) as well as Vision-and-Language Navigation frameworks (Savva et al., 2019). Another direction explores model-based reinforcement learning, where agents learn structured policies conditioned on textual goals; Dynalang (Lin et al., 2023) is a representative example of this approach, emphasizing learning world dynamics alongside goal-conditioned behavior.

**Instruction Following and Planning.** Recent work has demonstrated that large language models, when fine-tuned on suitable datasets, are capable of generating detailed, coherent action plans for agents based on textual instructions (Jansen, 2020; Zhao et al., 2024; Zhou et al.). Building on this ability, subsequent approaches have shown that planning performance can be further improved by incorporating feedback from the environment (Wang et al., 2023; Huang et al., 2022; Volovikova et al., 2024). For example SayCan (Ahn et al., 2022) augments planning with affordance-based evaluation via offline reinforcement learning, while (Tan et al., 2024) leverages policy optimization and probability normalization to enhance learning through interaction. Beyond improving general plan quality, environment feedback can also enable personalized planning; for example, (Han et al., 2024) introduced Reinforced Self-Training to iteratively align agents’ behavior with user preferences in object rearrangement tasks. Alternatively, Logeswaran et al. (2022) proposed a different strategy by avoiding language model fine-tuning altogether, instead generating multiple candidate plans with a frozen model and ranking them using mutual information and a learned feasibility model.

### 3 Problem Statement

The environment is formalized as a goal-based Partially Observable Markov Decision Process (POMDP), defined by the tuple  $(\mathcal{S}, \mathcal{A}, \mathcal{O}, \mathcal{T}, \mathcal{R}, \mathcal{G}, \gamma)$ . The agent receives a natural language instruction  $I$  and must achieve the corresponding latent goal  $g \in \mathcal{G}$ . Each observation  $o \in \mathcal{O}$  contains partial information about both the environment and the instruction  $I$ . The agent learns a grounding function  $f_g(I)$  to infer the latent goal  $g = f_g(I)$ .

The policy  $\pi(a \mid o)$  selects actions based on observations to maximize the expected cumulative reward:

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^T \gamma^t R(s_t, a_t, g) \mid o_0 \right].$$

The environment involves stochastic transitions  $\mathcal{T}(s' \mid s, a)$  and partial observability, requiring the agent to infer goals and act effectively under uncertainty.

We extend this setup by **introducing plans**. In the planning-augmented formulation, the agent does not receive the instruction  $I$  directly. Instead, it is provided with a plan  $p = (p_1, p_2, \dots, p_n)$  derived from  $I$ , where each step  $p_i$  corresponds to an intermediate subgoal  $g_i = f_g(p_i)$ . At each timestep, the agent observes the environment together with the current plan step  $p_{\phi(t)}$ . The optimization objective becomes:

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^T \gamma^t R(s_t, a_t, g_{\phi(t)}) \mid o_0 \right],$$

where  $g_{\phi(t)}$  is the subgoal associated with the active plan step.

In contrast to settings with predefined subtasks and explicit intermediate rewards, our formulation introduces two key challenges:

- 1. Subtask alignment under sparse rewards.** The agent must discover how its behavior aligns with intermediate subgoals despite only receiving sparse, delayed feedback upon completing the full instruction. This exacerbates the credit assignment problem.
- 2. Extended action space.** The agent must also decide when to terminate the current subtask. This requires augmenting the action space with control operations (e.g., a *DONE* action),

which increases both exploration complexity and the difficulty of learning effective switching strategies.

### 4 Super Igor

Super Igor framework proposes a method for jointly training a large language model and a reinforcement learning agent to solve instruction-following tasks. The LLM is responsible for transforming natural language instructions into structured plans, i.e. sequences of subtasks. The RL agent learns to execute these plans in the environment by interacting with it and maximizing delayed rewards.

The training process proceeds through the following stages:

- 1. Plan Generation (4.1):** The LLM extracts possible subtasks from instructions and generates multiple candidate plans in natural language during the initial cycle (Cycle 1). In subsequent cycles (Cycle 2–N), the candidate pool is iteratively refined by filtering and re-prioritization, based on how well the plans align with the RL agent’s performance.
- 2. Policy Learning (4.2):** The RL agent is trained to execute the selected plans in the environment.
- 3. Plan Validation (4.3):** The quality of candidate plans is evaluated according to the RL agent’s success rate and execution trajectories.
- 4. LLM Fine-Tuning (4.4):** The language model is fine-tuned with feedback derived from validation, aligning its scoring of plans with the agent’s actual performance.

#### 4.1 Plan Generation

In our approach, we first generate all possible plans for the training set in zero-shot mode during the initial cycle. In subsequent cycles, we progressively reduce the set of candidate plans by filtering out those that perform poorly for the agent. Concretely, the initial cycle produces the complete pool of plans, while later cycles re-prioritize them using the LLM’s negative log-likelihood (NLL) score. Importantly, we leverage the agent’s performance feedback as a preference signal to fine-tune the LLM with DPO, so that the model learns to align its scoring with the agent’s actual success in executing the plans.

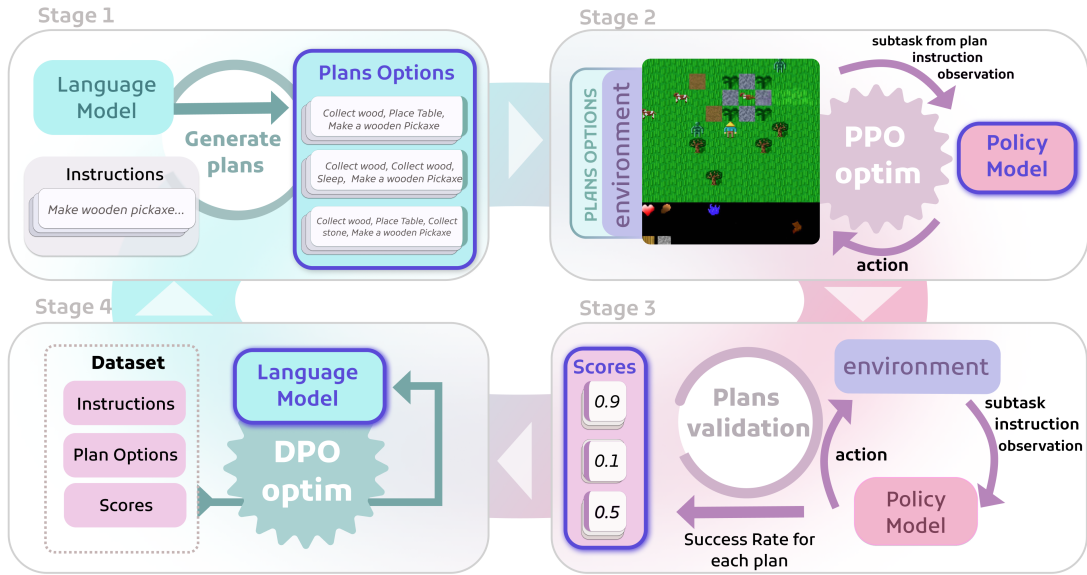


Figure 2: Super Igor Pipeline: The pipeline consists of four stages: (1) a language model generates multiple plan options for a given instruction; (2) a policy model is trained via PPO to execute each plan in the environment; (3) each plan is validated by measuring its execution success rate; (4) the language model is optimized using Direct Preference Optimization (DPO)(Rafailov et al., 2023) based on plan performance scores. This iterative loop refines both plan generation and execution.

### Training plans generation (Cycle 1).

Since the language model used for plan generation may not fully capture the exact dependencies and interaction rules of the target environment, we propose a structured procedure that separates the identification of goals from the reasoning about prerequisite constraints. The method unfolds in four steps.

First, we build a *subtask base* by extracting and canonicalizing possible subtasks from the instruction dataset, creating a unified vocabulary that reduces synonymy and ensures consistency. Each subtask is expressed in natural language, but in a strict normalized format that allows passing them one-by-one to the policy without ambiguity.

Second, the model generates a *goal-level plan*, producing for each instruction a single conceptual representation of its intended outcome, expressed in terms of the established subtask base. This step abstracts away from concrete execution details and captures only the high-level intent.

Third, we induce a *subtask ontology* that encodes the model’s hypotheses about prerequisite relations, i.e., which subtasks must be completed before others can be attempted. This provides a structured view of dependencies across the subtask base.

Finally, we perform *plan expansion*, where the single conceptual plan is unfolded into multiple detailed plans, with their number corresponding to

the hypotheses proposed by the model. The ontology ensures that these expanded variants remain consistent with prerequisite relations and avoid contradictions.

This approach provides two key benefits. First, it improves plan consistency by constructing plans from a shared set of subtasks and their relations, rather than from independent and potentially contradictory structures. Second, it supports partial normalization, since the model, when processing new instructions, tends to reuse previously identified subtasks, thereby reducing the proliferation of synonymous formulations. The details of the method and pseudocode are provided in Appendix A, and the prompts are presented in Appendix L.

**Plans re-prioritizing for RL-agent (Cycles 2-N).** After obtaining the initial feedback on agent performance for the generated plans and applying LLM fine-tuning (Subsection 4.4), subsequent cycles focus on re-prioritizing the candidate set. In each cycle, plans are rescored using the language model’s negative log-likelihood (NLL), which reflects how natural or plausible a plan is according to the model. Plans are then ranked by this score, and only the top-performing subset is retained for further training. As cycles progress, this iterative filtering process gradually narrows the candidate space, aligning the remaining plans both with the agent’s empirical success and with the model’s learned preferences.

## 4.2 Policy Learning

After the plans have been generated, we train a reinforcement learning agent using the stepwise plan observation setting (Subsection 3). At each timestep, the agent observes the environment and receives an embedding of the current plan step. It must learn to align actions with plan steps based on a delayed reward signal provided only upon successful completion of the entire plan. We use the PPO algorithm to train the policy.

---

### Algorithm 1 Skill Curriculum Learning

---

**Require:** Set of all plans  $\mathcal{P}$ , success-rate threshold  $\tau$

- 1: Initialize mastered skills  $\mathcal{M} \leftarrow \emptyset$
- 2: Initialize PPO agent  $\pi_\theta$
- 3: Initialize active plans

$$\mathcal{S} \leftarrow \{p \in \mathcal{P} \mid p \text{ contains exactly one skill}\}$$

- 4: **while** training not converged **do**
- 5:   Train  $\pi_\theta$  on active plans  $\mathcal{S}$  and collect rollouts
- 6:   For each skill  $s$ , compute success rate:

$$SR(s) = \frac{\# \text{ Successful episodes containing } s}{\# \text{ Total episodes containing } s}$$

- 7:   **if**  $SR(s) \geq \tau$  **then**
- 8:     Add to mastered skills  $\mathcal{M} \leftarrow \mathcal{M} \cup \{s\}$
- 9:   **end if**
- 10:   Update plans

$$\mathcal{S} \leftarrow \{p \in \mathcal{P} \mid p \text{ has at most one unmastered skill}\}$$

- 11: **end while**
  - 12: **return**  $\pi_\theta, \mathcal{M}$
- 

To address the sparse reward problem in training, we introduce **Skill Curriculum Learning**. The core principle is to create a dynamic curriculum that begins with the simplest single-subtask tasks, allowing the agent to learn foundational behaviors under a relatively dense reward signal.

As the agent trains, we monitor its Success Rate (SR) for each subtask. Once a subtask’s SR surpasses a predefined threshold  $\tau$ , it is marked as "mastered." This mastery triggers an update to the curriculum: the set of active training plans is expanded to include any plan composed of already mastered subtasks and, at most, one new, unmastered subtask. This incremental expansion, detailed in Algorithm 1, ensures a smooth learning gradient and prevents the agent from being overwhelmed.

## 4.3 Plan Validation

To evaluate each proposed plan, we run the RL agent multiple times using that plan as input. Because the environment is highly stochastic, a single

rollout is insufficient; instead, we aggregate metrics such as average success rate to obtain a reliable estimate of plan effectiveness.

## 4.4 LLM Fine-Tuning

In the first cycle, we warm-start the language model with supervised fine-tuning (SFT) to reproduce the plans generated in the zero-shot stage (Section 4.1), aligning it with the plan distribution of the target environment.

In subsequent cycles, we incorporate plan-level quality signals obtained from execution and validation. These signals are used to construct preference pairs of higher- and lower-scoring plans, which are then used for DPO fine-tuning. This allows the model to internalize the agent’s feedback and gradually improve plan generation.

In our framework, DPO acts as a lightweight plan-selection bias rather than a precise credit assignment mechanism. It increases the probability of plan structures that the RL agent can learn from early, implicitly forming an automated curriculum, while deprioritizing plans that yield little initial progress without explicitly labeling them as incorrect.

## 5 Experiments

In this section, we describe the experiments conducted to answer the following research questions (RQ):

**RQ1. (Effectiveness and Generalization of Auto-Generated Plans):** How well can the SuperIgor agent learn to follow instructions by leveraging LLM-generated plans, and how well does this learned behavior generalize to new instructions? We measure effectiveness as the agent’s final success rate on training tasks (Atomic and Combo splits). We measure generalization using final success rates on two test sets: Paraphrases (same goals, new wording) and New Objects (new goal combinations).

**RQ2. (Policy Training under Sparse Feedback):** How well can the SuperIgor policy model be trained to follow plans under sparse feedback? The primary metric for this is the final SR on the training tasks.

**RQ3. (Agent Effectiveness with Iterative SuperIgor Cycles):** How does the agent’s performance evolve over multiple iterations of the SuperIgor planning-training cycle?

## 5.1 Environment

We conduct our experiments on the CrafterText benchmark (Volovikova et al., 2025), which provides a unified testbed for evaluating instruction-following agents in a multimodal, dynamic, and partially observable open-ended environment. Leveraging the modular design of CrafterText, which allows for the procedural generation of custom tasks, we construct a specialized dataset named **FOCUS**.

The FOCUS dataset is built upon CrafterText’s Achievement scenarios but is specifically engineered to rigorously test instruction adherence under strict constraints. It contains over 900 instructions with a vocabulary of more than 1,500 unique words. The training set is composed of two instruction types: **Atomic**, which specify a single, indivisible goal (e.g., “craft a furnace”), and **Combo**, which combine multiple atomic goals into a sequence of actions (e.g., “craft a furnace and then collect wood”).

A critical challenge in this domain is that task composition often involves overlapping subtasks. For example, crafting a furnace first requires making a wooden pickaxe and collecting stone — steps that are also required for other goals. Consequently, agents may learn generic subroutines that maximize reward without truly grounding the linguistic instruction. To address this, FOCUS enforces a *strict evaluation protocol*: an instruction is considered successful only if all its specified goals are completed precisely as requested, with no extraneous achievements triggered. This prevents the agent from exploiting broad, non-specific strategies.

To evaluate generalization, FOCUS employs two distinct test sets. The **Paraphrases** set consists of Combo instructions from the training set reformulated with novel vocabulary and syntax. The **New Objects** set introduces new combinations of Atomic goals that appeared during training but never occurred together in a single instruction. This structure allows FOCUS to assess both robustness to linguistic variation and compositional generalization.

## 5.2 Experiments Setup

In our pipeline, we generate plans using Qwen2.5-14B-Instruct<sup>2</sup>, fine-tune it for one epoch with DPO ( $\beta = 0.5, lr = 1 \times 10^{-5}$ ) to stabilize lo-

<sup>2</sup><https://huggingface.co/Qwen/Qwen2.5-14B-Instruct>

cal updates, and then train policies with PPO-T ( $lr = 0.001, \epsilon = 0.02$ ) and Skill Curriculum Learning for 2.5B steps. We validate by executing 10 plans across 50 seeds to assess robustness. Two full cycles were conducted, with evaluations before and after LLM fine-tuning, and results compared against baselines at 2.5B and 5B steps (Figure 3). Additional hyperparameters are described in more detail in Appendix J.

## 5.3 Baselines

For our comparative analysis, we use several established baselines from the original CrafterText study (Volovikova et al., 2025). PPO-T (Text-Augmented PPO) augments PPO with textual grounding: instructions are encoded using a frozen DistilBERT [CLS] embedding, concatenated with CNN-based visual features, and processed by a GRU to maintain temporal context. PPO-T+ (Plan-Augmented PPO) extends this by first translating each instruction into a structured plan with GPT-4, and then providing the agent with a plan embedding instead of the raw instruction.

FiLM (Perez et al., 2018) offers an alternative integration of language and vision. Here, instruction embeddings generate parameters that modulate CNN outputs via Feature-wise Linear Modulation layers, allowing textual context to directly shape visual feature processing.

To ensure consistency, all baselines follow a strict protocol requiring the DONE action to signal task completion, with success only counted when both the instruction is satisfied and DONE invoked. We also evaluate an *Auto-DONE (Soft-)* variant, where episodes terminate automatically upon completion, and include an Oracle agent trained with PPO-T and Skill Curriculum Learning on human-written ground-truth plans.

## 5.4 Experiments Result

### RQ1. Effectiveness and Generalization of Auto-Generated Plans in the SuperIgor Pipeline

a) **Auto-generated plans train agents far more effectively than instruction-only baselines.** On Atomic tasks (Figure 3(a)), SuperIgor agents (SI-DPO / SI-SFT) reach 0.35–0.45, compared to only 0.10–0.19 for instruction-only RL baselines. Oracle remains higher at 0.56–0.65, but the SuperIgor  $\rightarrow$  Oracle gap ( $\approx 0.20$ ) is much smaller than the Baselines  $\rightarrow$  SuperIgor gap ( $\approx 0.25$ –0.30), clearly showing the value of plan supervision. On Combo

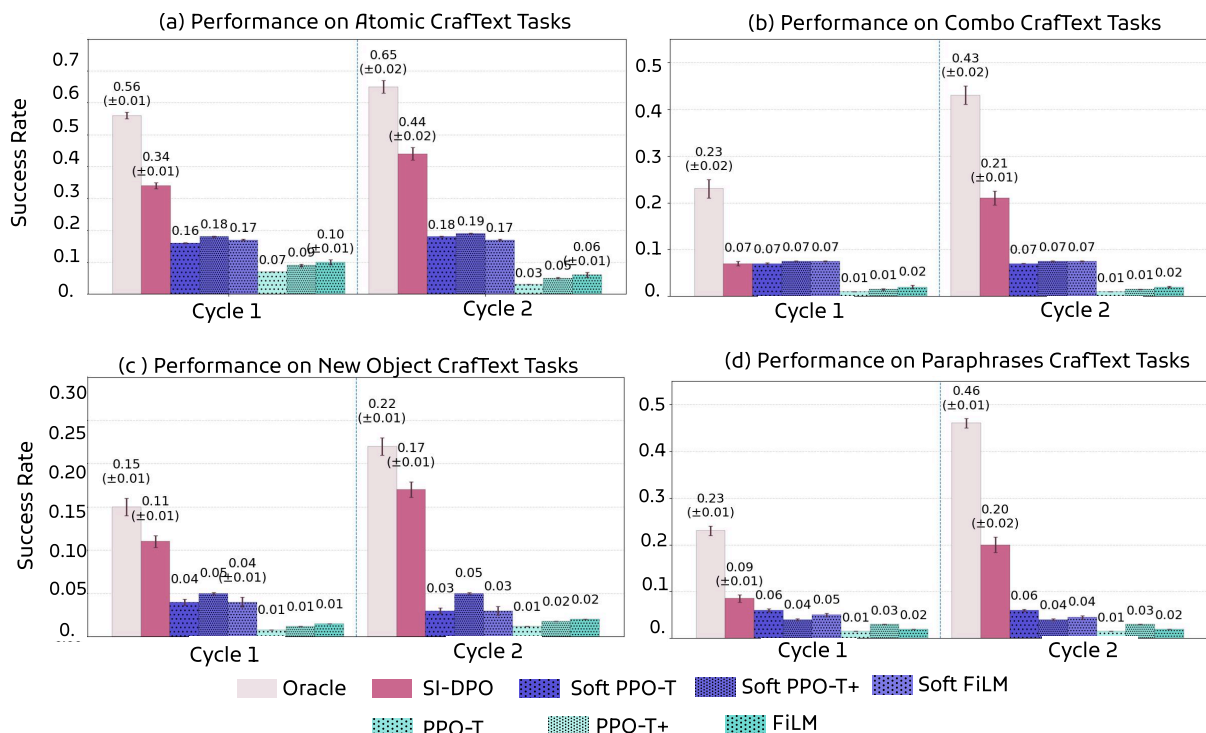


Figure 3: Comparison of SuperIgor and baseline performance on CraffText tasks (Atomic / Combo / New Objects / Paraphrases). SI-SFT denotes SuperIgor validated on plans generated after LLM supervised fine-tuning, while SI-DPO denotes SuperIgor validated after LLM DPO fine-tuning. All agents were evaluated at 2.5 billion steps (corresponding to the first cycle in the SuperIgor approach) and 5 billion steps (corresponding to the second cycle).

tasks (Figure 3(b)), SuperIgor achieves 0.21, outperforming baselines at 0.08, while Oracle reaches 0.46. The wider gap to Oracle here can be explained by the fact that SI agents must simultaneously learn up to 20 alternative plans, whereas Oracle is trained on a single expert-aligned plan, which simplifies optimization.

#### b) Agents trained with auto-generated plans generalize on unseen goals better than those trained with Oracle plans.

On Combo tasks, Oracle achieves 0.46, while SuperIgor reaches 0.21. But on New Object tasks (Figure 3(c)), Oracle drops sharply to  $\approx 0.22$ , while SI decreases more moderately to 0.12–0.17. Thus, although SI lags in absolute terms, its performance is more stable: the Oracle–SI gap shrinks from 0.25 on Combo to only 0.05–0.10 on New Object tasks. We attribute this stronger generalization precisely to the fact that SI agents learn from multiple alternative plans per instruction, which exposes them to richer variability during training.

#### c) Agents trained with auto-generated plans do not lose performance when instructions are paraphrased.

Paraphrases reuse (Figure 3 (d)) the same goals as in Combo tasks but are expressed in different

linguistic forms. In Cycle 1, SI-DPO performance increases from 0.07 on Combo to 0.09 on Paraphrases. In Cycle 2, SI-DPO remains stable, with 0.21 on Combo and 0.20 on Paraphrases. This shows that SuperIgor agents can successfully transfer their learned strategies to differently worded instructions, maintaining performance even when the language of the goal changes.

### RQ2. Policy Training under Sparse Feedback

#### a) Skill Curriculum Learning enhances agent to learn more subtasks compared to unstructured training

We evaluate the training process by the number of unique subtasks the agent masters over time. A subtask is considered "mastered" once its success rate surpasses a 70% threshold. This metric provides a clearer insight into the agent's growing capabilities and its ability to handle compositional tasks. We compare three configurations, with the results visualized in Figure 4.

The agent trained with **Skill Curriculum on Oracle Plans** sets a practical upper bound for performance. By the 10 billion step mark, it successfully masters **14 distinct subtasks**. It signifies that the agent has acquired almost the entire 'mining' technology tree: all the achievements from col-

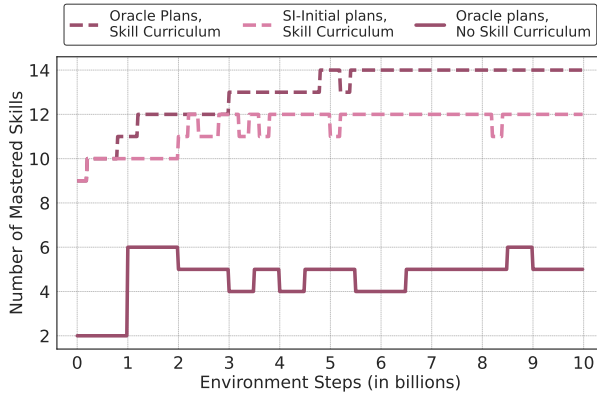


Figure 4: A comparative analysis of the number of mastered subtasks over 10 billion environment steps. The results highlight the critical role of the Skill Curriculum, as agents trained without it fail to learn, even with optimal Oracle Plans.

lecting wood to collecting iron. Furthermore, it demonstrates the ability to execute complex, combined instructions that require interleaving subtasks from different progression branches, such as eating, drinking, and collecting resources within a single, coherent plan.

Agent trained on Oracle plans without the Skill Curriculum perform worse with only mastered 5 basic subtasks. Even with a flawless plan, the agent fails to learn without a structured progression that allows it to build foundational skills first. This finding confirms that Skill Curriculum helps to overcome sparse feedback problem and enables agent abilities to learn more subtasks.

#### b) SI-Initial plans are a good initial approximation of optimal plans

Skill Curriculum with SI-Initial plans graph follows this trajectory closely, mastering 12 subtasks within the same timeframe. This demonstrates the high quality of our SI-INITIAL plan generation, as it enables the agent to acquire most of the subtasks achievable even with perfect plans. The gap between these two curves represents the remaining challenge in our automated plan generation.

In conclusion, the curriculum is not just beneficial, it is *critical* for meaningful skill acquisition in this environment. The ablation clearly shows that our Skill Curriculum Learning framework is the key enabler of learning, while our SI-INITIAL procedure generates plans of sufficient quality to unlock a significant portion of the agent’s potential and a good baseline for futhermore plan generation improvement using SuperIgor framework.

### RQ3. Agent Effectiveness with Iterative SuperIgor Cycles

a) Plan-following quality improves across cycles. On Atomic tasks (training, Figure 3, (a)), SI-DPO increases from 0.34 in Cycle 1 to 0.43 in Cycle 2. On Combo tasks (training, Figure 3, (b)), SI-DPO grows from 0.06 in Cycle 1 to  $\approx 0.21$  in Cycle 2. On New Object tasks (testing, Figure 3, (c)), SI-DPO declines only slightly from  $\approx 0.21$  to 0.12–0.17, showing that performance improves with additional SuperIgor cycles on both training and testing setups and remains relatively stable when moving to unseen goals.

b) Plan reprioritization under DPO illustrates the process by which language models are incrementally grounded in the agent’s behavior and the underlying environment mechanics.. The re-ranking visualization (Appendix E, Figure 9) shows how plans shift across SFT, DPO-C1, and DPO-C2. Success Rates range from 0.68 to 0.86. A plan with  $SR = 0.86$  steadily climbs to the top across cycles, while weaker plans with  $SR \approx 0.68$  remain consistently at the bottom. These changes are gradual rather than abrupt, suggesting that DPO provides a soft grounding signal that progressively aligns plan priorities with the agent’s execution success. This interpretation is further supported by the ablation study (Appendix B.1), which shows that integrating DPO accelerates agent learning in the second cycle, indicating that gradual plan re-prioritization translates into more effective behavioral grounding.

## 6 Conclusion

In this work, we introduced SuperIgor, a framework for training agents to follow complex natural-language instructions in dynamic, multimodal environments. We showed that effective instruction-following can be learned by integrating language-model-based planning with reinforcement learning, even in the absence of predefined subtasks and without explicit verification of intermediate sub-goal completion. This is achieved through the integration of ontology-based plan generation, which provides the structural foundation for learning and enables both Skill Curriculum Learning to handle sparse and delayed rewards and cyclic DPO-based language model adaptation that prioritizes effective plans and accelerates agent training. Together, these components enable SuperIgor to outperform instruction-only baselines, generalize near Oracle on out-of-distribution tasks, and remain robust to instruction paraphrasing.

## 612 **Limitations**

613 While SuperIgor demonstrates a scalable approach  
614 to integrating language-model-based planning with  
615 reinforcement learning, it has several limitations.  
616 First, the method depends on the quality of the ini-  
617 tial plan structure induced by the language model.  
618 The ontology graph is generated fully automati-  
619 cally, without human supervision, which may lead  
620 to redundant, inconsistent, or cyclic dependencies  
621 between subtasks. In practice, we mitigate this is-  
622 sue empirically by evaluating different language  
623 models and selecting those that produce stable and  
624 coherent plan graphs (Appendix B.3); our experi-  
625 ments show that several mid-sized models already  
626 perform well in this setting. Nevertheless, robust-  
627 ness to poorly structured plan spaces remains a  
628 limitation of the current approach.

629 In addition, the interaction between plan qual-  
630 ity and policy learning is not always interpretable.  
631 While our validation protocol provides a clear rela-  
632 tive preference signal between candidate plans  
633 for the same instruction (via repeated evaluation  
634 across seeds), it can still be difficult to diagnose  
635 training failures at the system level. In particu-  
636 lar, when learning stalls, it is unclear whether the  
637 bottleneck is caused by limitations of the current  
638 policy optimization and exploration dynamics or  
639 by systematic issues in the generated plan space  
640 (e.g., missing prerequisites or overly difficult de-  
641 compositions).

## 642 **References**

643 Michael Ahn, Anthony Brohan, Noah Brown, Yevgen  
644 Chebotar, Omar Cortes, Byron David, Chelsea Finn,  
645 Chuyuan Fu, Keerthana Gopalakrishnan, Karol Haus-  
646 man, et al. 2022. Do as i can, not as i say: Ground-  
647 ing language in robotic affordances. *arXiv preprint*  
648 *arXiv:2204.01691*.

649 Bowen Baker, Ilge Akkaya, Peter Zhokov, Joost  
650 Huizinga, Jie Tang, Adrien Ecoffet, Brandon  
651 Houghton, Raul Sampedro, and Jeff Clune. 2022.  
652 Video pretraining (vpt): Learning to act by watch-  
653 ing unlabeled online videos. *Advances in Neural*  
654 *Information Processing Systems*, 35:24639–24654.

655 Adrian Bolton, Alexander Lerchner, Alexandra  
656 Cordell, Alexandre Moufarek, Andrew Bolt, Andrew  
657 Lampinen, Anna Mitenkova, Arne Olav Hallingstad,  
658 Bojan Vujatovic, Bonnie Li, et al. 2025. Sima 2: A  
659 generalist embodied agent for virtual worlds. *arXiv*  
660 *preprint arXiv:2512.04797*.

661 Maxime Chevalier-Boisvert, Dzmitry Bahdanau, Salem  
662 Lahlou, Lucas Willems, Chitwan Saharia, Thien Huu

Nguyen, and Yoshua Bengio. 2018. Babyai: A plat-  
form to study the sample efficiency of grounded lan-  
guage learning. *arXiv preprint arXiv:1810.08272*.

Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Man-  
dlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang,  
De-An Huang, Yuke Zhu, and Anima Anandkumar.  
2022. Minedojo: Building open-ended embodied  
agents with internet-scale knowledge. *Advances in*  
*Neural Information Processing Systems*, 35:18343–  
18362.

Jonathan Gray, Kavya Srinet, Yacine Jernite, Hao-  
nan Yu, Zhuoyuan Chen, Demi Guo, Siddharth  
Goyal, C. Lawrence Zitnick, and Arthur Szlam. 2019.  
[Craftassist: A framework for dialogue-enabled inter-  
active agents](#).

Dongge Han, Trevor McInroe, Adam Jelley, Stefano V  
Albrecht, Peter Bell, and Amos Storkey. 2024. Llm-  
personalize: Aligning llm planners with human pref-  
erences via reinforced self-training for housekeeping  
robots. *arXiv preprint arXiv:2404.14285*.

Felix Hill, Olivier Tieleman, Tamara Von Glehn,  
Nathaniel Wong, Hamza Merzic, and Stephen Clark.  
2020. Grounded language learning fast and slow.  
*arXiv preprint arXiv:2009.01719*.

Wenlong Huang, Pieter Abbeel, Deepak Pathak, and  
Igor Mordatch. 2022. Language models as zero-shot  
planners: Extracting actionable knowledge for em-  
bodied agents. In *International conference on ma-  
chine learning*, page 9118–9147. PMLR.

Peter A Jansen. 2020. Visually-grounded planning  
without vision: Language models infer detailed  
plans from high-level instructions. *arXiv preprint*  
*arXiv:2009.14259*.

Muyao Li, Zihao Wang, Kaichen He, Xiaojian Ma,  
and Yitao Liang. 2025. Jarvis-vla: Post-training  
large-scale vision language models to play visual  
games with keyboards and mouse. *arXiv preprint*  
*arXiv:2503.16365*.

Shalev Lifshitz, Keiran Paster, Harris Chan, Jimmy Ba,  
and Sheila McIlraith. 2023. Steve-1: A generative  
model for text-to-behavior in minecraft. *Advances in*  
*Neural Information Processing Systems*, 36:69900–  
69929.

Jessy Lin, Yuqing Du, Olivia Watkins, Danijar Hafner,  
P. Abbeel, Dan Klein, and Anca D. Dragan. 2023.  
[Learning to model the world with language](#). *ArXiv*,  
abs/2308.01399.

Lajanugen Logeswaran, Yao Fu, Moontae Lee, and  
Honglak Lee. 2022. Few-shot subgoal planning with  
language models. In *Proceedings of the 2022 Con-  
ference of the North American Chapter of the Asso-  
ciation for Computational Linguistics: Human Lan-  
guage Technologies*, pages 5493–5506.

716	Aryan Mathur and Asaduddin Ahmed. 2025. Adapting interleaved encoders with ppo for language-guided reinforcement learning in babyai. <i>arXiv preprint arXiv:2510.23148</i> .			
717				
718				
719				
720	Fabian Paischer, Thomas Adler, Markus Hofmarcher, and Sepp Hochreiter. 2023. Semantic helm: A human-readable memory for reinforcement learning. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .			
721				
722				
723				
724				
725	Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. 2018. Film: Visual reasoning with a general conditioning layer. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 32.			
726				
727				
728				
729				
730	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. <i>Advances in neural information processing systems</i> , 36:53728–53741.			
731				
732				
733				
734				
735	Manolis Savva, Jitendra Malik, Devi Parikh, Dhruv Batra, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, and Vladlen Koltun. 2019. <a href="#">Habitat: A platform for embodied AI research</a> . In <i>2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019</i> , pages 9338–9346. IEEE.			
736				
737				
738				
739				
740				
741				
742				
743	Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and Dieter Fox. 2020. Alfred: A benchmark for interpreting grounded instructions for everyday tasks. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 10740–10749.			
744				
745				
746				
747				
748				
749				
750	Alessandro Suglia, Qiaozhi Gao, Jesse Thomason, Govind Thattai, and Gaurav S. Sukhatme. 2021. <a href="#">Embodied bert: A transformer model for embodied, language-guided visual task completion</a> . <i>ArXiv</i> , abs/2108.04927.			
751				
752				
753				
754				
755	Weihaio Tan, Wentao Zhang, Shanqi Liu, Longtao Zheng, Xinrun Wang, and Bo An. 2024. True knowledge comes from practice: Aligning llms with embodied environments via reinforcement learning. <i>arXiv preprint arXiv:2401.14151</i> .			
756				
757				
758				
759				
760	Zoya Volovikova, Gregory Gorbov, Petr Kuderov, Aleksandr Panov, and Alexey Skrynnik. 2025. <a href="#">CrafText benchmark: Advancing instruction following in complex multimodal open-ended world</a> . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 26131–26151, Vienna, Austria. Association for Computational Linguistics.			
761				
762				
763				
764				
765				
766				
767				
768	Zoya Volovikova, Alexey Skrynnik, Petr Kuderov, and Aleksandr I Panov. 2024. Instruction following with goal-conditioned reinforcement learning in virtual environments. In <i>ECAI 2024</i> , pages 650–657. IOS Press.			
769				
770				
771				
772				
		H. J. Austin Wang and Karthik Narasimhan. 2021. <a href="#">Grounding language to entities and dynamics for generalization in reinforcement learning</a> . <i>ArXiv</i> , abs/2101.07393.		773 774 775 776
		Zihao Wang, Shaofei Cai, Guanzhou Chen, Anji Liu, Xiaojian Ma, and Yitao Liang. 2023. Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents. <i>arXiv preprint arXiv:2302.01560</i> .		777 778 779 780 781
		Lewei Yao, Jianhua Han, Youpeng Wen, Xiaodan Liang, Dan Xu, Wei Zhang, Zhenguo Li, Chunjing Xu, and Hang Xu. 2022. Detclip: Dictionary-enriched visual-concept paralleled pre-training for open-world detection. <i>Advances in Neural Information Processing Systems</i> , 35:9125–9138.		782 783 784 785 786 787
		Andrew Zhao, Daniel Huang, Quentin Xu, Matthieu Lin, Yong-Jin Liu, and Gao Huang. 2024. Expel: Llm agents are experiential learners. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 19632–19642.		788 789 790 791 792
		Victor Zhong, Austin W Hanjie, Sida I Wang, Karthik Narasimhan, and Luke Zettlemoyer. 2021. <a href="#">Silg: The multi-environment symbolic interactive language grounding benchmark</a> . <i>arXiv preprint arXiv:2110.10661</i> .		793 794 795 796 797
		Victor Zhong, Tim Rocktäschel, and Edward Grefenstette. 2019. Rtfm: Generalising to novel environment dynamics via reading. <i>arXiv preprint arXiv:1910.08210</i> .		798 799 800 801
		Enshen Zhou, Yiran Qin, Zhenfei Yin, Yuzhou Huang, Ruimao Zhang, Lu Sheng, Yu Qiao, and Jing Shao. 2024. Minedreamer: Learning to follow instructions via chain-of-imagination for simulated-world control. <i>arXiv preprint arXiv:2403.12037</i> .		802 803 804 805 806
		Zhiyuan Zhou, Pranav Atreya, Abraham Lee, Homer Walke, Oier Mees, and Sergey Levine. Autonomous improvement of instruction following skills via foundation models, 2024. <i>URL https://arxiv.org/abs/2407.20635</i> .		807 808 809 810 811

## A Plans Generation: candidates generation

In our plan extraction method the goal is to elicit the model’s hypotheses about the dependencies between subtasks in the environment. In the first step we construct a subtask bank  $B$ , i.e., the set of all candidate subtasks derived from the instruction set. For each instruction  $I \in \mathcal{D}$ , we prompt the language model  $f_{\text{LLM}}$  to generate a goals plan  $\mathcal{P}[I]$ , i.e., the set of goal subtasks directly required by the instruction. The model is provided with the current contents of the subtask bank  $B$ , which encourages reuse of already known subtasks and reduces the introduction of redundant synonyms. If the generated goals contain subtasks not yet present in  $B$ , they are added. At the initial iteration the bank is empty, so all subtasks generated by the model are included. The complete process is summarized in Algorithm 2.

Once a sufficiently rich subtask bank  $B$  has been established, ontological dependencies between subtasks are extracted. For each target subtask  $t \in B$ , the language model is queried multiple times to determine which elements from  $B$  are required for the completion of  $t$ . For every candidate dependency ( $r \rightarrow t$ ), its probability is estimated as

$$P(r \rightarrow t) = \frac{k_t}{N},$$

where  $k_t$  denotes the number of times subtask  $r$  was identified as necessary for  $t$  and  $N$  is the number of queries. To filter out spurious associations, the Wilson confidence interval is applied to the resulting probabilities. The procedure is carried out in two passes: first over the entire bank  $B$ , and then restricted to the subtasks previously identified as relevant, which refines the weighting of relations. The final output is an ontology graph  $G = (V, E)$  that encodes the model’s hypothesized structure of interrelations among subtasks. The full procedure is summarized in Algorithm 3.

After constructing the ontology  $G = (V, E)$ , each goal plan  $\mathcal{P}[I]$  is expanded with its dependencies. For every subtask  $s \in \mathcal{P}[I]$ , we recursively collect all prerequisites in  $G$ . The union of these subtasks with the original goals defines the plan’s vertices, which are then topologically sorted so that prerequisites precede dependents. The result is a linearized plan  $P$  containing the goals and all supporting subtasks (Algorithm 4).

---

### Algorithm 2 Subtask Bank Update

---

**Require:** Instruction stream  $\mathcal{D}$ , language model  $f_{\text{LLM}}$

**Ensure:** Subtask bank  $B$ , goals plans  $\mathcal{P}$

- 1: Initialize subtask bank  $B \leftarrow \emptyset$
- 2: Initialize goals plans  $\mathcal{P} \leftarrow \emptyset$
- 3: **for** each instruction  $I \in \mathcal{D}$  **do**
- 4:     Identify goal subtasks conditioned on  $B$ :

$$S \leftarrow f_{\text{LLM}}(I, B)$$

- 5:     **for** each subtask  $s \in S$  **do**
  - 6:         **if**  $s \notin B$  **then**
  - 7:              $B \leftarrow B \cup \{s\}$
  - 8:         **end if**
  - 9:     **end for**
  - 10:     Goals plan for  $I$ :  $\mathcal{P}[I] \leftarrow S$
  - 11: **end for**
  - 12: **return**  $B, \mathcal{P}$
- 

---

### Algorithm 3 Ontology Construction

---

**Require:** Subtask bank  $B$ , language model  $f_{\text{LLM}}$ , queries per pass  $N$ , threshold  $\tau$

**Ensure:** Ontology graph  $G = (V, E)$

- 1: Initialize counts  $count(r, t) \leftarrow 0$  for all  $r, t \in B, r \neq t$
- 2: **for** each target subtask  $t \in B$  **do**
- 3:     **for** two passes **do**
- 4:         Define candidate set  $C$ :

$$C \leftarrow \begin{cases} B \setminus \{t\}, & \text{pass 1} \\ \{r \in B : count(r, t) > 0\}, & \text{pass 2} \end{cases}$$

- 5:     **for**  $i = 1 \dots N$  **do**
- 6:         Query prerequisites:

$$R \leftarrow f_{\text{LLM}}(t, C)$$

- 7:     **for** each  $r \in R$  **do**
- 8:          $count(r, t) \leftarrow count(r, t) + 1$
- 9:     **end for**
- 10:     **end for**
- 11:     **end for**
- 12: **end for**
- 13: Initialize edge set  $E \leftarrow \emptyset$
- 14: **for** each pair  $(r, t)$  **do**
- 15:     Compute probability:

$$\hat{p}(r \rightarrow t) = \frac{count(r, t)}{N}$$

- 16:     Compute Wilson lower bound  $LB(\hat{p}, N)$
  - 17:     **if**  $LB \geq \tau$  **then**
  - 18:          $E \leftarrow E \cup \{(r \rightarrow t)\}$
  - 19:     **end if**
  - 20: **end for**
  - 21: **return**  $G = (V = B, E)$
-

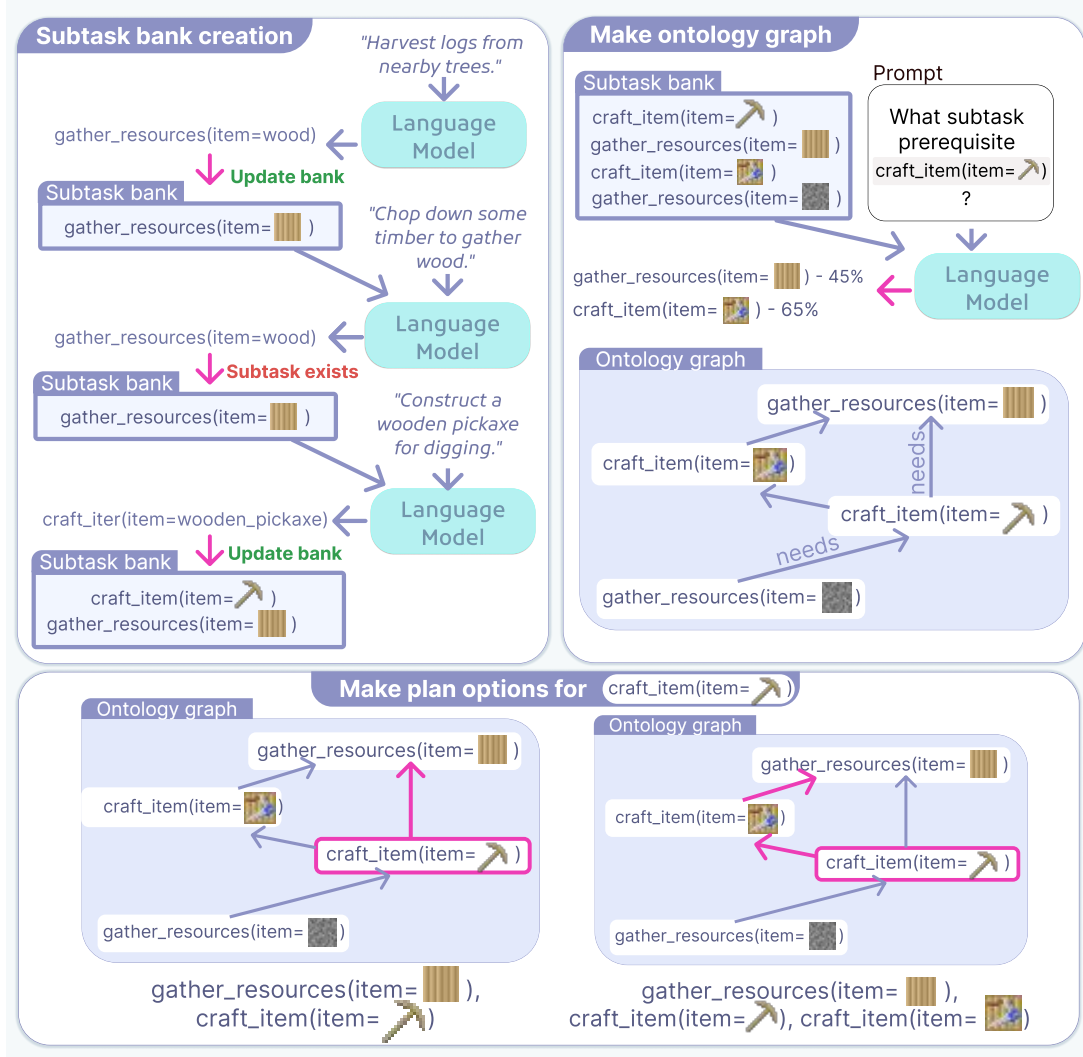


Figure 5: Training plans generation

**Algorithm 4** Final Plan Generation from Ontology

**Require:** Instruction  $I$ , goals mapping  $\mathcal{G}$ , goals plan  $\mathcal{P}$ , ontology  $G = (V, E)$

**Ensure:** Final plan  $P$

- 1: Retrieve goal subtasks:  $S \leftarrow \mathcal{G}[I]$
- 2: Initialize plan vertex set:  $U \leftarrow S$
- 3: **for** each  $s \in S$  **do**
- 4: Expand prerequisites via ontology:
 
$$D \leftarrow \text{PREREQCLOSURE}(s, G)$$
- 5:  $U \leftarrow U \cup D$
- 6: **end for**
- 7: Extract induced subgraph:  $G_U \leftarrow G[U]$
- 8: Topologically sort  $G_U$  to obtain ordered plan  $P$
- 9: **return**  $P$

**B Additional Experiments**

**B.1 Ablation Study: SuperIgor Framework**

To quantify the contribution of each module of the SuperIgor framework, we conduct an ablation study in which individual components are removed from the training pipeline. We evaluate the influence of four factors: (1) Ontology-Based Training Plan Generation, (2) Curriculum design in the RL stage, (3) LLM plan-model pretraining (SFT), and (4) DPO finetuning based on RL agent performance signals. Table 1 presents the results of this experiment, where we measure the SuperIgor agent’s SuccessRate on the Atom subset of the CrafterText instruction dataset. The analysis of the results yields two central findings.

(1) **Curriculum is effective only when paired with high-quality, ontology-structured plans.** Although a full-cycle evaluation may give the impres-

860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877

Table 1: Ablation study of the SuperIgor framework, measuring agent SuccessRate on the Atom subset of the CraFText dataset across two training cycles

Ontology	Curriculum	DPO	SFT	Cycle-1	Cycle-2
✗	✓	✓	✓	0.06	N/A
✓	✗	✓	✓	0.08	N/A
✓	✓	✗	✓	0.34	0.39
✓	✓	✓	✗	0.25	0.13
✓	✓	✓	✓	<b>0.35</b>	<b>0.45</b>

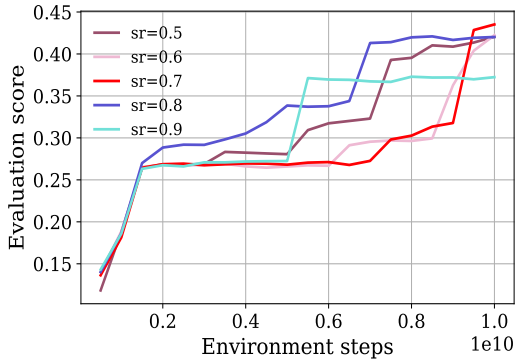


Figure 6: Ablation of the skill-mastery threshold  $\tau$ . The plot shows evaluation scores on the Atomic and Combo tasks during training for different  $\tau$  values.

878 sion that the primary gains come from curriculum  
 879 learning, the results of this ablation study show that  
 880 its effectiveness emerges only in combination with  
 881 ontology-guided plan generation. Without ontol-  
 882 ogy (i.e., without structured, hierarchical plans),  
 883 the curriculum has no meaningful ordering signal  
 884 and fails to provide improvement: Cycle-1 perfor-  
 885 mance drops to 0.06 when ontology is removed.

886 Ontology-based plans, however, naturally en-  
 887 code a hierarchy of instructions and goals, enabling  
 888 a principled progression from simpler to more com-  
 889 plex targets. This hierarchical structure is pre-  
 890 cisely what makes a curriculum implementable:  
 891 the RL agent can first master low-complexity goals  
 892 and then gradually advance to more difficult ones.  
 893 When ontology is present, this alignment between  
 894 plan structure and staged learning produces large  
 895 gains, improving Cycle-1 performance from 0.06  
 896 (no curriculum) to 0.35 (with curriculum).

897 **(2) DPO improves the RL agent by learning**  
 898 **to prioritize plans that lead to higher-quality be-**  
 899 **havior.** Unlike SFT, which is trained to reproduce  
 900 the ontology-induced distribution of plans, DPO  
 901 directly leverages RL performance as a preference  
 902 signal: it learns to rank plans higher when they  
 903 empirically yield better agent behavior. Removing

DPO results in weaker prioritization: the RL agent  
 reaches only 0.39 in Cycle-2 without DPO, com-  
 pared to 0.45 when DPO is included. Thus, DPO  
 systematically shifts the plan distribution toward  
 behaviorally effective plans, accelerating and am-  
 plifying the RL agent’s improvement across cycles.

## B.2 Ablation Study: Skill Mastery Threshold

We conducted an ablation study to analyze the sen-  
 sitivity of the Skill Curriculum Learning to the mas-  
 tery threshold parameter  $\tau$ . Figure 6 presents the  
 final performance of the Skill Curriculum Learn-  
 ing agent after 10 billion environment steps in  
 CraFText-Symbolic configuration for different val-  
 ues of  $\tau$  ranging from 0.5 to 0.9.

The results demonstrate that  $\tau = 0.7$  provides  
 an optimal balance for curriculum progression. We  
 hypothesize that lower thresholds ( $\tau = 0.5$ ) allow  
 the agent to progress too quickly to complex skills  
 before achieving reliable proficiency, while higher  
 thresholds ( $\tau = 0.9$ ) cause the agent to spend ex-  
 cessive time perfecting basic skills, slowing overall  
 learning. The  $\tau = 0.7$  value strikes an optimal  
 balance between progression speed and skill reli-  
 ability.

## B.3 Ablation Study: Choice of LLM for Ontology and Training-Plan Generation

To understand how the choice of language model  
 affects the quality of the ontology and the generated  
 training plans we conducted an ablation study com-  
 paring several families of LLMs. For each model,  
 we regenerated both the ontology and the full train-  
 ing dataset (plans), and then trained an RL agent  
 using our Skill Curriculum Learning procedure.

Table 2 reports the agent’s success rate on  
 the training split under different planner mod-  
 els. The experiment includes models from the  
 Qwen and Gemma families, as well as the larger  
 microsoft/NextCoder-32B model.

### (1) Larger models do not necessarily produce

Table 2: Ablation on the choice of LLM used for generating both ontology and training plans. We report success rate on the training set.

LLM	Qwen1.5-32B	NextCoder-32B	Qwen1.5-14B	Gemma-12B	Qwen-7B
SR (Train)	0.43	0.26	0.35	0.14	0.22

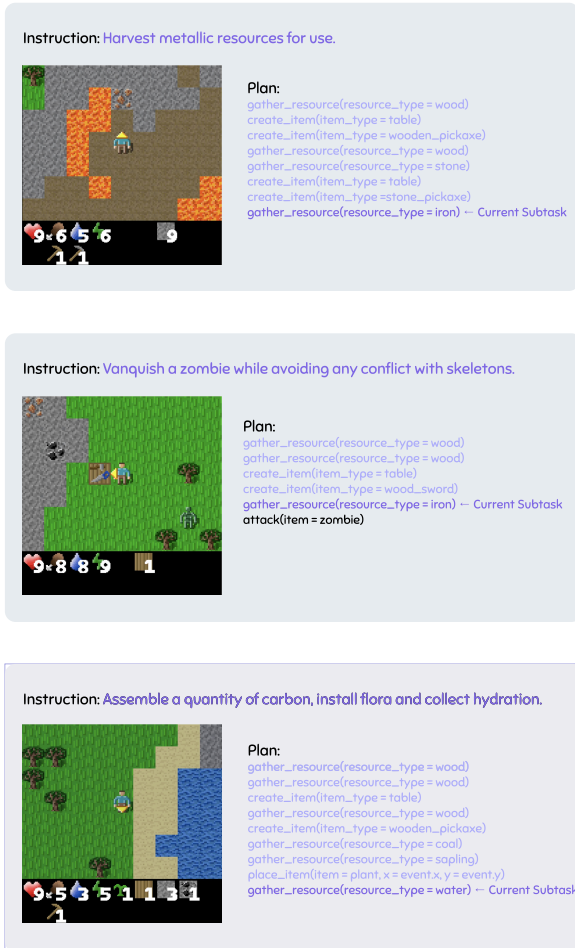


Figure 7: Example of instructions and corresponding plans

**better ontologies or plans.** Although one may expect the largest models to generate the most structured plans, but NextCoder-32B performance is surpassed by significantly smaller Qwen models. Qwen-32B yields the highest performance (0.43), and even Qwen-7B outperforms Gemma-12B, indicating that model family and training specialization matter more than raw parameter count.

**(2) Qwen models produce more stable and semantically consistent plan structures.** Models from the Qwen family demonstrate higher robustness in generating hierarchical task decompositions that align with our ontology constraints. This leads to more reliable curriculum construction and more

effective RL training.

**(3) Some widely used LLMs fail to benefit from the alignment stage.** We also conducted experiments with several other well-known models, including *microsoft/phi-4*, *mistralai/Mistral-7B-Instruct-v0.2*, and *openai/gpt-oss-20b*, and found that the alignment stage does not provide any measurable benefit for them. Despite explicit prompt constraints on which subtasks should be used, these models tend to generate large numbers of synonymously similar subtasks. Consequently, the set of goals that the agent must recover becomes even larger than when instructions are provided directly, rendering it impractical to run the full pipeline with these models.

## C CraffText

To provide a concrete example of our method, Figure 7 visualizes the agent’s state at a single timestep. The CraffText environment, shown on the left, is a dynamic grid-world where the agent must gather resources, craft items, and navigate diverse terrains to survive and complete tasks.

The core of our approach lies in the hierarchical decomposition of complex goals. As shown on the right, a high-level instruction, which may be ambiguous or require long-term planning (e.g., "Craft an iron pickaxe."), is first translated into a deterministic, multi-step plan. Each step in this plan constitutes a distinct subtask.

Crucially, the agent’s policy is not conditioned on the entire plan. Instead, it focuses solely on the currently active subtask. This transforms a challenging long-horizon problem into a more tractable sequence of short-horizon tasks. The agent’s objective at any moment is to complete the highlighted subtask and then invoke the DONE action. For example, optimal agent can choose DONE action based on the inventory state (when completing subtasks such as collecting resources and crafting items), player status (for subtasks that are related to eating, drinking or sleeping) or map state (for subtasks such as placing blocks).

Upon successful completion, the framework provides the next subtask in the sequence, guiding the

1001 agent through the overall plan until the final goal is  
 1002 achieved.

1003 We utilize the **FOCUS** instruction dataset de-  
 1004 scribed in the main text. The structure of the dataset  
 1005 is as follows:

1006 Training Set:

- 1007 • Atomic: Single, indivisible goals (e.g., "Craft  
 1008 a furnace").
- 1009 • Combo: Sequences of multiple atomic goals  
 1010 (e.g., "Craft a furnace and then collect wood").
- 1011 • Crucially, each instruction in the training set  
 1012 also has a paraphrased version to encourage  
 1013 linguistic robustness from the start.

1014 Test Sets (Out-of-Distribution):

- 1015 • Paraphrases: Contains the same underlying  
 1016 goals as the Combo training set, but expressed  
 1017 with novel vocabulary and syntax. This tests  
 1018 robustness to linguistic variation.
  - 1019 – Training Combo: "Consume beef and  
 1020 create a stone pickaxe."
  - 1021 – Test Paraphrase: "Eat steak and forge a  
 1022 stone pickaxe." or "Devour cow meat and  
 1023 create a stone pickaxe."
- 1024 • New Objects: Introduces new combinations of  
 1025 atomic goals that appeared during training but  
 1026 never occurred together in a single instruction  
 1027 in the training set. This directly tests composi-  
 1028 tional generalization. These instructions also  
 1029 come with their own paraphrases.
- 1030 • Training contained: "Consume beef" and  
 1031 "Forge a stone pickaxe" and "Forge a stone  
 1032 blade" as separate atomic or part of other com-  
 1033 bos.
- 1034 • Test New Object: "Consume beef and forge  
 1035 a stone blade." or "Eat cow meat and create a  
 1036 sword from stone."

1037 This structure allows us to rigorously dissect  
 1038 the agent’s capabilities: learning from language  
 1039 (Atomic/Combo), generalizing to new phrasing  
 1040 (Paraphrases), and generalizing to new goal combi-  
 1041 nations (New Objects).

## 1042 D Complete SuperIgor Training Pipeline

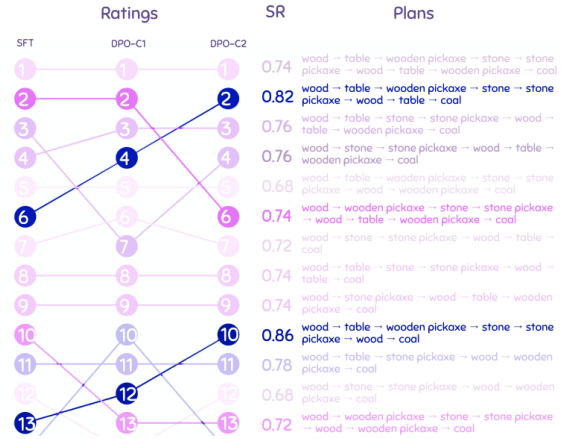
1043 The SuperIgor framework integrates multiple com-  
 1044 ponents that exchange specific inputs and outputs  
 1045 during training. Below we describe the key data  
 1046 flows between components:

### 1047 Component Interfaces:

- 1048 • **LLM Planner** ( $f_{LLM}$ )
  - 1049 – **Input:** Instruction  $I$
  - 1050 – **Output:** Candidate plan  $\mathcal{P}$  consisting  
 1051 of a sequence of subtasks from subtask  
 1052 bank  $\mathcal{B}$
- 1053 • **RL Policy** ( $\pi_\theta$ )
  - 1054 – **Input:** Environment observations  $o_t$ , cur-  
 1055 rent plan step DistilBERT [CLS] embed-  
 1056 ding  $p_{\phi(t)}$  of a plan  $\mathcal{P}$
  - 1057 – **Output:** Action  $a_t$  from extended action  
 1058 space containing default Crafttext actions  
 1059 and additional DONE action that gives  
 1060 the agent the next plan step embedding  
 1061  $p_{\phi(t+1)}$  of a plan  $\mathcal{P}$

1062 The complete training procedure integrating all  
 1063 components is summarized in Algorithm 5

## 1064 E DPO plans reprioritization



1065 Figure 8: Example of DPO plan reprioritization for the  
 1066 instruction: "Forge a stone pickaxe and mine coal"

## 1067 F Training Details: Policy Optimization

1068 Our low-level policy, which is responsible for exe-  
 1069 cuting individual subtasks, is trained using Proxi-  
 1070 mal Policy Optimization (PPO). The agent’s goal at  
 this stage is to learn an optimal strategy for complet-  
 ing a given subtask based on its visual observations.

---

**Algorithm 5** Complete SuperIgor Training Pipeline

---

**Require:**

- 1: Environment  $\mathcal{E}$
- 2: Instruction dataset  $\mathcal{D}_{\text{train}} = \{I_1, I_2, \dots, I_N\}$
- 3: Initial LLM planner  $f_{\text{LLM}}$  with parameters  $\theta_{\text{LLM}}$
- 4: Initial RL policy  $\pi_\theta$  with parameters  $\theta_{\text{RL}}$
- 5: Mastery threshold  $\tau$ , number of cycles  $C$

**Ensure:**

- 6: Optimized planner  $f_{\text{LLM}}^*$
  - 7: Trained policy  $\pi_\theta^*$
  - 8:
  - 9: Initialize subtask bank  $\mathcal{B} \leftarrow \emptyset$
  - 10: Initialize candidate plans  $\mathcal{P} \leftarrow \{\}$
  - 11: Initialize mastered subtasks  $\mathcal{M} \leftarrow \emptyset$
  - 12:
  - 13: **Initial Plan Generation (Cycle 1):**
  - 14: Extract subtasks:  $\mathcal{S} \leftarrow f_{\text{LLM}}(\mathcal{D}_{\text{train}})$
  - 15: Build ontology:  $\mathcal{O} \leftarrow \text{BuildOntology}(\mathcal{S}, f_{\text{LLM}})$
  - 16: Generate initial plans:  $\mathcal{P}_{\text{initial}} \leftarrow \text{ExpandPlans}(\mathcal{D}_{\text{train}}, \mathcal{O})$
  - 17:
  - 18: Fine-tune  $f_{\text{LLM}}$  on  $\mathcal{P}_{\text{initial}}$  using SFT
  - 19: Generate training plans:  $\mathcal{P} \leftarrow f_{\text{LLM}}(\mathcal{D}_{\text{train}})$
  - 20:
  - 21: **for** cycle  $c = 1$  **to**  $C$  **do**
  - 22:
  - 23:     **Policy Training with Skill Curriculum:**
  - 24:     Train  $\pi_\theta$  on  $\mathcal{P}$  using PPO with Skill Curriculum Learning
  - 25:     Update mastered subtasks  $\mathcal{M}$  based on success rates
  - 26:
  - 27:     **Plan Validation:**
  - 28:     Execute  $\pi_\theta$  with plans  $\mathcal{P}$  for multiple seeds
  - 29:     For every plan  $P$  in compute average success rate  $SR(p)$
  - 30:     Construct preference dataset  $\mathcal{D}_{\text{pref}}$
  - 31:
  - 32:     **LLM Fine-tuning:**
  - 33:     Fine-tune  $f_{\text{LLM}}$  on  $\mathcal{D}_{\text{pref}}$  using DPO
  - 34:
  - 35:     **Plan Generation:**
  - 36:     Select plans for new training epoch:
  - 37:      $\mathcal{P} \leftarrow \text{SelectPlans}(f_{\text{LLM}}, \mathcal{D}_{\text{train}}, \mathcal{P})$
  - 38: **end for**
  - 39:
  - 40: **return**  $f_{\text{LLM}}, \pi_\theta$
- 

The standard clipped surrogate objective for PPO is defined as:

$$\mathcal{L}^{\text{PPO}}(\theta) = \mathbb{E}_t \left[ \min \left( \rho_t(\theta) \hat{A}_t, \text{clip}(\rho_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right],$$

where  $\rho_t(\theta) = \frac{\pi_\theta(a_t|o_t)}{\pi_{\theta_{\text{old}}}(a_t|o_t)}$  is the probability ratio and  $\hat{A}_t$  is the estimated advantage at timestep  $t$ .

Our agent’s policy and value functions are parameterized by a single neural network with a shared multimodal feature extractor and separate actor and critic heads. The visual stream processes the  $63 \times 63$  pixel image with 3 channels observations using a three-layer Convolutional Neural Network (CNN). Each convolutional layer utilizes 32

filters with a  $5 \times 5$  kernel, followed by a ReLU activation and max-pooling. For the language stream, textual instructions are encoded using a pre-trained BERT model (`bert-base-uncased`), and we use the embedding of the [CLS] token as the final text representation.

The flattened output of the CNN and the text embedding are then concatenated to form a unified multimodal representation. This combined feature vector is fed into two separate feed-forward networks: the **actor head**, which outputs the logits for the categorical action distribution, and the **critic head**, which outputs a scalar estimate of the state-value function.

## G Training Details: LLM Fine-Tuning

To improve the high-level planner (the LLM), we employ a reinforcement learning-based feedback loop. The planner generates a sequence of subtasks (a plan), which is then executed by the PPO agent. The final outcome of the agent’s execution (e.g., task success or failure, efficiency) serves as a signal to update the planner.

**Direct Preference Optimization (DPO).** This method aligns the model toward preferred completions using pairwise preference data. The DPO loss is:

$$\mathcal{L}^{\text{DPO}} = -\log \sigma \left( \beta \left( \log \pi(x^+ | q) - \log \pi(x^- | q) \right) \right),$$

where  $x^+$  and  $x^-$  are preferred and less preferred plans for instruction  $q$ , and  $\beta$  is a temperature parameter.

## H Agent’s plan following

Figure 9 presents the example of how the agent follows the plan and chooses actions. For each subtask, there are two frames: the first shows the observation up to the moment when the agent takes the action that completes the subtask, along with the action distribution at that time; the second shows a few timesteps later, when the agent decides to skip the subtask in order to solve a new one.

## I Compute Resources

All experiments were conducted on a high-performance computing cluster equipped with nodes containing 1 NVIDIA A100 GPU with 80 GB of VRAM. Each node was powered by an 12 CPU Cores CPU with 96 GB of system RAM.

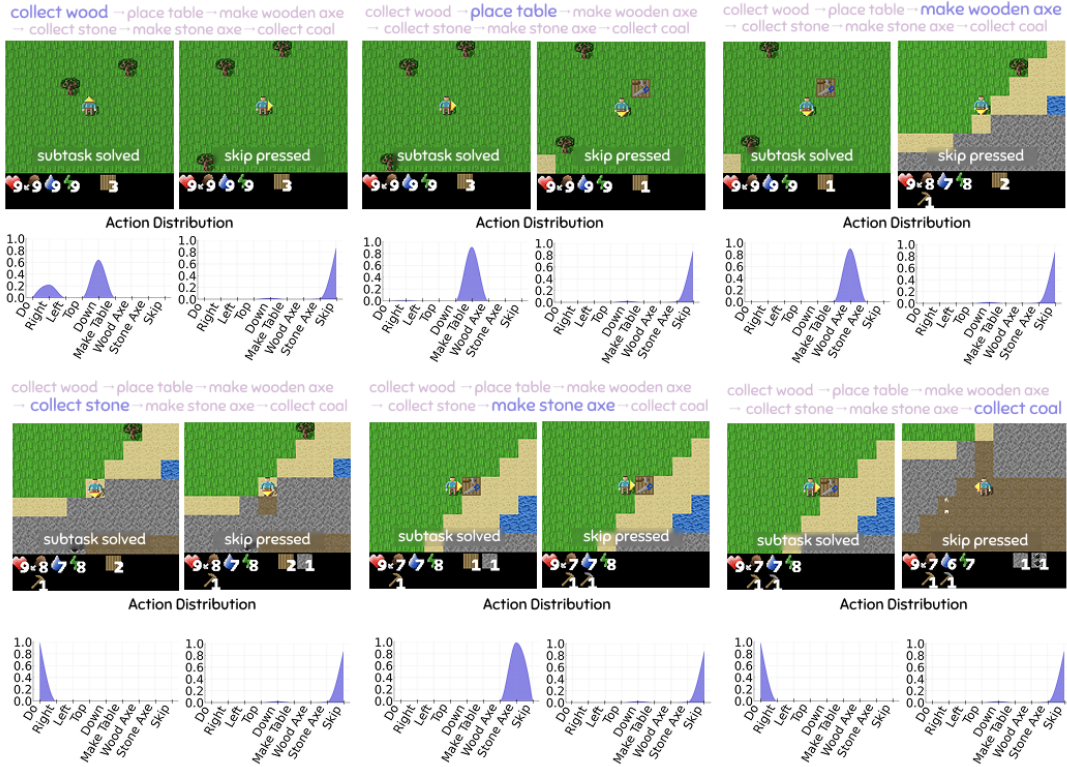


Figure 9: Example of how the agent follows the plan and chooses actions.

The total computational budget can be broken down into two primary stages:

**Policy Training and Evaluation.** The primary computational cost stems from training the PPO agent. Each full training run for a single configuration up to 10 billion environment steps took approximately 120-150 GPU-hours. Reproducing all presented experiments, including the baseline comparisons and ablation studies, required a total of 10 such training runs.

**LLM Training and Generation.** The initial generation of plans using the Qwen2.5-14B-Instruct model for the entire dataset required approximately X GPU-hours on a single NVIDIA A100 GPU. Epoch of finetuning LLM with DPO on evaluated plans takes approximately 15 GPU-hours.

In total, we estimate the full computational cost to reproduce all results presented in this paper to be approximately 2000-2500 GPU-hours.

## J Training Details: Hyperparameters

The hyperparameters for our experiments are detailed in Table 3. For the PPO agent training, we adopt the configuration from the original Crafter baseline study (Volovikova et al., 2025). For the LLM planner fine-tuning, we use a Q-LoRA approach with a comprehensive set of parameters

optimized for efficient large model training.

1155

Table 3: Hyperparameters used for training the low-level agent and fine-tuning the high-level planner.

Hyperparameter	Value
<i>PPO Agent Training</i>	
Learning rate	0.0002
Discount factor ( $\gamma$ )	0.99
GAE lambda ( $\lambda$ )	0.95
Clipping epsilon ( $\epsilon$ )	0.2
PPO epochs	4
Number of minibatches	8
Entropy coefficient	0.01
Value function coef.	0.5
Activation function	Tanh
Hidden layer size	512
<i>LLM Planner Fine-Tuning</i>	
Base model	Qwen2.5-14B-Instruct
Training epochs	1
Learning rate (SFT)	2e-4
Learning rate (DPO)	1e-5
Beta (DPO)	0.5
Optimizer	Paged AdamW (32-bit)
LR scheduler	Cosine
Warmup ratio	0.03
Batch size (per device)	16
Gradient accumulation	1
Gradient clipping norm	0.3
Weight decay	0.001
Mixed precision	bf16
<i>LoRA Configuration</i>	
LoRA rank ( $r$ )	64
LoRA alpha ( $\alpha$ )	16
LoRA dropout	0.1
<i>Quantization (4-bit)</i>	
Quantization type	nf4
Compute dtype	float16

## K AI Assistants Usage

AI assistants were used solely for language editing and text refinement. All scientific contributions and conclusions are the sole responsibility of the authors.

## L Plans Generation: Prompt for plan generation

You control an agent in a 2D game with simplified Minecraft environment.  
You will need to provide a detailed step-by-step plan for following the user's instructions.  
You must include all the preliminary steps that it needs to complete.

You are controlling an agent in a 2D game set within a simplified Minecraft-like environment.  
The agent starts from scratch with an empty inventory and no gathered resources.  
Your task is to generate a step-by-step plan that enables the agent to follow a given user instruction.

What you must do:

- Break down the instruction into atomic actions the agent needs to perform.
- Include all necessary preliminary steps, such as gathering or crafting resources.
- Assume the agent has nothing at the beginning -- you must plan from the ground up.
- Output your answer as a Python list of strings.
- Each string must represent one atomic skill invocation, written on a separate line.

Format for each step:

```
"skill_name(arg1 = value1, arg2 = value2, ...)"  
- skill_name: the name of the primitive skill or action the agent will execute.  
- Inside the parentheses, list all required arguments with their names and corresponding values.
```

Example:

```
gather_resource(resource_type = wood)
```

Each of the step agents will be implemented without knowledge of what it did before, so it can only rely on observation and the current step. Therefore, each step must be self-sufficient and not require knowledge of past steps.

"If the instruction doesn't specify what the agent needs to do and is more general--like 'Explore the world' or 'Go out and examine the world around you'--send explore(object=world). In this case, the plan should consist of only one step: "explore(object=world)"."

Send your answer as a python list.

Instruction: Make a pickaxe from wood

Answer:

```
["gather_resource(resource_type = wood)",  
"gather_resource(resource_type = wood)",  
"create_item(item_type = table)",  
"gather_resource(resource_type = wood)",  
"gather_resource(resource_type = wood)",  
"create_item(item_type = wooden_pickaxe)"]
```

Send your answer as a python list.

Instruction: \$INSTRUCTION\$

Answer: