



# Auto MC-Reward: Automated Dense Reward Design with Large Language Models for Minecraft

Hao Li<sup>1,2\*</sup>, Xue Yang<sup>2\*</sup>, Zhaokai Wang<sup>2,3\*</sup>, Xizhou Zhu<sup>4,5</sup>,  
Jie Zhou<sup>4</sup>, Yu Qiao<sup>2</sup>, Xiaogang Wang<sup>1,5</sup>, Hongsheng Li<sup>1</sup>, Lewei Lu<sup>5</sup>, Jifeng Dai<sup>4,2</sup>✉

<sup>1</sup>CUHK-SenseTime Joint Laboratory, The Chinese University of Hong Kong

<sup>2</sup>OpenGVLab, Shanghai AI Laboratory

<sup>3</sup>Shanghai Jiao Tong University <sup>4</sup>Tsinghua University <sup>5</sup>SenseTime Research

[https://yangxue0827.github.io/auto\\_mc-reward.html](https://yangxue0827.github.io/auto_mc-reward.html)

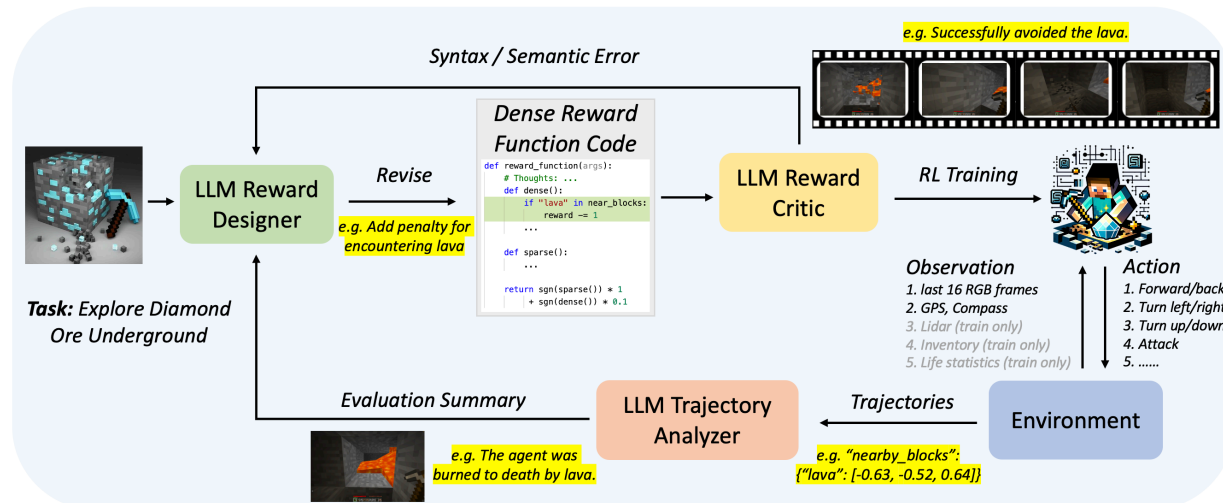


Figure 1. Overview of our Auto MC-Reward. Auto MC-Reward consists of three key LLM-based components: Reward Designer, Reward Critic, and Trajectory Analyzer. A suitable dense reward function is iterated through the continuous interaction between the agent and the environment for reinforcement learning training of specific tasks, so that the model can better complete the task. An example of exploring diamond ore is shown in the figure: i) Trajectory Analyzer finds that the agent dies from lava in the failed trajectory, and then gives suggestion for punishment when encountering lava; ii) Reward Designer adopts the suggestion and updates the reward function; iii) The revised reward function passes the review of Reward Critic, and finally the agent avoids the lava by turning left.

## Abstract

Many reinforcement learning environments (e.g., Minecraft) provide only sparse rewards that indicate task completion or failure with binary values. The challenge in exploration efficiency in such environments makes it difficult for reinforcement-learning-based agents to learn complex tasks. To address this, this paper introduces an advanced learning system, named Auto MC-Reward, that leverages Large Language Models (LLMs) to automatically design dense reward functions, thereby enhancing the learning efficiency. Auto MC-Reward consists of three important components: Reward Designer, Reward Critic, and Trajectory Analyzer. Given the environment information

and task descriptions, the Reward Designer first design the reward function by coding an executable Python function with predefined observation inputs. Then, our Reward Critic will be responsible for verifying the code, checking whether the code is self-consistent and free of syntax and semantic errors. Further, the Trajectory Analyzer summarizes possible failure causes and provides refinement suggestions according to collected trajectories. In the next round, Reward Designer will further refine and iterate the dense reward function based on feedback. Experiments demonstrate a significant improvement in the success rate and learning efficiency of our agents in complex tasks in Minecraft, such as obtaining diamond with the efficient ability to avoid lava, and efficiently explore trees and animals that are sparse in the plains biome.

\*Equal contribution. This work was completed by Hao Li and Zhaokai Wang during their internship at Shanghai Artificial Intelligence Laboratory. ✉Corresponding author: Jifeng Dai <daijifeng@tsinghua.edu.cn>.

## 1. Introduction

Minecraft, as the world’s best-selling game, offers a range of tasks from exploration, survival to creating. It has become an important environment for researching efficient Reinforcement Learning (RL) [16, 22]. In particular, its extreme sparsity of rewards and huge complexity of the decision space pose significant challenges for RL. Currently, the most effective learning strategy involves pre-training through behavior cloning [3], using learned behavioral priors to narrow the decision space. Nevertheless, it still requires billions of environmental interactions for effective learning due to the sparse reward nature of Minecraft.

On the other hand, previous researchers have proposed a variety of dense reward signals to enable efficient learning for specific sparse reward tasks [2, 27, 28, 35, 37]. However, their applicability on the complex and long-horizon tasks in Minecraft remains an open question. To deeply reveal the challenges in Minecraft, we examine on several representative challenging tasks, *e.g.* exploring underground for diamonds. We find that even after behavior cloning, most of these methods still fail to make significant progress on these tasks, further highlighting the difficulty of Minecraft and the limitations of existing dense reward methods.

It is noteworthy that for human players, Minecraft is a relatively simple casual game [12]. The advantage of human lies in their ability to summarize based on practice. For example, an accidental burning death from lava can teach human to avoid getting too close to it. Such summaries, based on life experience and practice, are key to human intelligence [40, 43]. Most existing RL methods overlook this ability. On the other hand, Large Language Models (LLMs) have recently demonstrated human-like common sense and reasoning capabilities [18, 36, 52]. We find that leveraging LLMs can help RL agents simulate the practice summarization abilities of human. Based on the historical action trajectories and success-failure signals of the agents, LLMs can automatically design and refine corresponding auxiliary rewards, effectively overcoming the sparse reward challenge in Minecraft.

According to above analysis, we propose an automated method named Auto MC-Reward, to design and improve auxiliary reward functions according to task descriptions and historical action trajectories. This method utilizes the task understanding and experience summarization abilities of LLMs, providing detailed and immediate rewards for learning guidance. Specifically We first use LLMs to design task-related dense reward functions based on basic descriptions of the environment and tasks, named as Reward Designer. These reward functions are used to train agents after self-verification, *i.e.* Reward Critic. To address potential biases or oversights in LLM’s understanding, we also propose a LLM-based Trajectory Analyzer to analyze and summarize collected trajectories from the trained agent, and

help Reward Designer to improve the reward functions.

We verify the effectiveness of Auto MC-Reward on a series of representative benchmarks, including horizontal exploration for diamonds underground and approaching trees and animals in the plains biome. Experiments show that Auto MC-Reward achieves significantly better results on these tasks compared to original sparse reward and existing dense reward methods, showing its advanced ability of empowering efficient learning on sparse reward tasks. By iteratively refining the design of rewards functions, Auto MC-Reward enables the agent to efficiently learn new behaviors that is beneficial to the corresponding tasks, *e.g.* avoiding lava, which greatly improves the success rate. Moreover, Auto MC-Reward achieves a high diamond obtaining success rate (36.5%) with only raw information, demonstrating its ability of solving long-horizon tasks.

## 2. Related Work

**Minecraft Agents** are intelligent agents designed to accomplish various tasks while playing the game Minecraft. Most of previous works adopt reinforcement learning for agent training. Due to the extremely sparse rewards and complex decision space of Minecraft tasks, early attempts have tried hierarchical RL [30, 34, 41, 42], curriculum learning [23], and imitation learning [1] to empower more efficient RL training. To narrow the decision space, recent work [3] build a foundation model by performing imitation on YouTube videos. DreamerV3 [17] instead learns a world model that explores the environment efficiently. As the LLMs demonstrate their general planning ability, a series of research [19, 46, 47, 51, 53] leverage LLMs as high-level planners that decompose long-term complex tasks as basic skills and implement the skills with RL agents or handcrafted scripts.

Auto MC-Reward aims to design dense rewards for Minecraft tasks automatically using LLMs, which is orthogonal to previous works on Minecraft agents that mainly focus on RL learning algorithms or high-level planning.

**Efficient Learning in Sparse Reward Tasks** is a long-standing challenge in RL due to the lack of effective learning signals [26]. A common solution is to handcraft dense reward functions that provide intermediate reward signals based on human expertise, which requires time-consuming trial-and-error for each environment and task. Another line of previous research focus on proposing general-purpose dense auxiliary reward functions, such as curiosity-driven exploration [6, 20, 28, 37], self-imitation learning [35], and goal-conditioned reinforcement learning [2, 24, 27, 45]. Despite the success on certain specific tasks, the applicability of these methods in the complex environment of Minecraft remains uncertain. Recent works [10, 11, 14, 31, 32] also propose to use pre-trained models to assign reward

to intermediate states of completing tasks. However, these approaches produce reward values in a black-box manner, which cannot be interpreted and improved based on the experience of the agents, and the generalizability of these models on new tasks is not guaranteed.

In contrast, Auto MC-Reward automatically produces explainable reward functions according to the task descriptions. Moreover, the reward functions can be improved to be more precise based on the experience of the agent.

**Automated Reward Function Design** aims to find an optimal reward function that drives efficient reinforcement learning for interested tasks. Previous works [9, 15] employ evolutionary algorithm for searching optimal reward functions for specific tasks. Most of these attempts have a highly constrained search space that only adjusts parameters of task-specific handcrafted reward templates. Recently, a series of research [7, 11, 25, 29] employs LLMs for integrating human preference into open-domain tasks without clear completion criteria by directly prompting LLMs with environment trajectories and natural language task descriptions. The reward values are generated on-the-fly by LLMs, which is black-box and has heavy computational cost due to the nature of LLMs. In contrast, Auto MC-Reward employs LLMs to generate white-box code-form reward functions.

Concurrent works [33, 49, 50] also propose to use LLMs as a coder to generate reward functions for robotics control tasks. Specifically, L2R [50] needs to prepare reward function templates and cannot cope with unexpected situations in open worlds. Text2Reward [49] and EUREKA [33] require complete environment code or description and rely on human feedback, which are not available in open worlds. Different from these methods, Auto MC-Reward considers more complex Minecraft environments that has diverse scenarios and high uncertainty, requiring more precise and thorough reward designing.

### 3. Method

Auto MC-Reward consists of three components: Reward Designer, Reward Critic, and Trajectory Analyzer. Given the environment information and task descriptions, the Reward Designer proposes the reward function by coding an executable Python function with pre-defined observation inputs. The Reward Critic verifies if the proposed reward function is self-consistent and if it meets the format requirements. The designed reward function which passes the Reward Critic is used to train agents in the environment. To improve the designed reward function according to empirical experience, the Trajectory Analyzer is proposed to summarize possible failure causes and provide refinement suggestions on the reward function based on the inference trajectories of the trained agents. Then the Reward Designer modifies the reward function based on the feedback from

Trajectory Analyzer. Figure 1 shows the overview of the Auto MC-Reward.

#### 3.1. Reward Designer

We utilize a Reward Designer to generate the reward function code to provide intermediate instructive learning signals to the agent. It takes as input task descriptions, game information, and reward function requirements, generating reward functions in executable code form. When updating reward function, we also provide analysis of the agent’s performance when interacting with the game environment. The input prompt is introduced in Section 4.2.

The generated reward function uses a pre-defined observation format as input. This includes the nearest distance of each block type within the visible range in the current and previous steps, changes in inventory between adjacent steps, health points, and the agent’s location in each past step. These parameters can provide information on the agent’s current and historical states, assisting the reward function in various situations.

**Multi-Step Memory.** Long-term tasks require the transfer of information across multiple steps. Thus, we introduce a multi-step memory mechanism. It is provided to Reward Designer as a empty dictionary at the beginning, and the reward function can save necessary data into the memory to be used in future steps. In the actual reward function of the explore-tree task, we observed that the agent records the distance to a tree at each step, thereby encouraging getting closer with the tree than the previous step.

**Chain of Thought.** We require the LLM to first describe its design thoughts, such as considering potential failure reasons and the details of the reward function design. These thoughts are to be written as comments at the beginning of the code. This is a mechanism similar to Chain of Thought (CoT) [48], where the thought process precedes the coding implementation. In the specific code implementation, necessary comments will also be generated every few lines (e.g., “Check if lava is in the field of view in the previous step”). This approach not only allows Reward Designer to refer to the text-form thoughts during reward function initialization, but also assists subsequent Reward Critic in assessing the code’s rationality, and helps Reward Designer to understand the current reward function’s purpose when updating the reward function.

**Scale Constraints.** We impose a specific scale constraint for the reward function, where the LLM generates two sub-functions: dense and sparse. *Sparse* denotes rewards for achieving the final goal or heavy penalties (like death), while *dense* represents dense intermediate signals during the task completion process. We preset their numerical values and only allow the LLM to determine their positivity or negativity, limiting *sparse* to  $\{1, 0, -1\}$ , and *dense* to

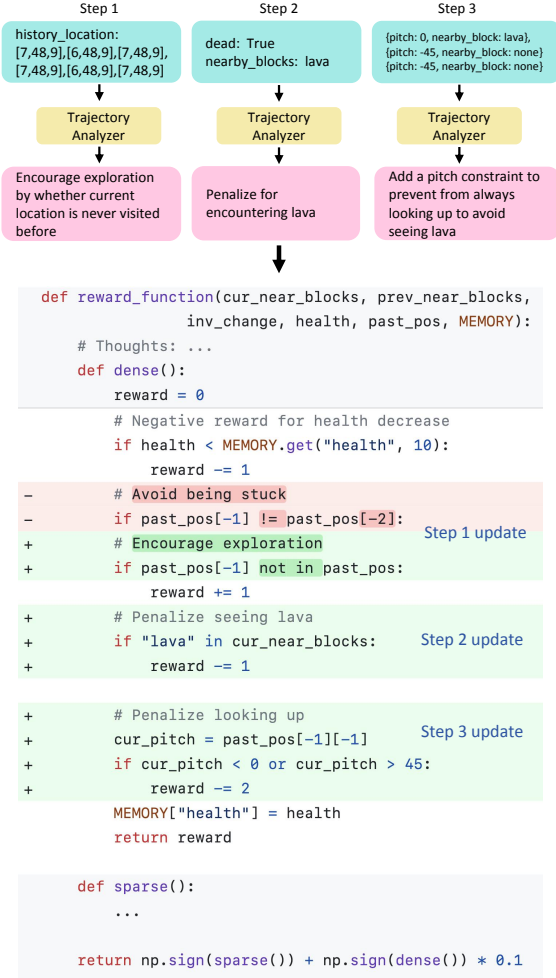


Figure 2. Example of updating the reward function. Trajectory Analyzer provides analysis for three scenarios at different steps, and then Reward Designer update the reward function based on the suggestions. We only display part of the trajectory data for brevity. **Step 1**: rewrite the code of encouraging exploration to avoid going back and forth. **Step 2**: add lava penalty to avoid falling into lava. **Step 3**: add pitch constraint to avoid constantly looking up to avoid lava.

$\{0.1, 0, -0.1\}$ . They are then added together for the final reward. Therefore, the final reward values can be one value of  $\{\pm 1.1, \pm 1.0, \pm 0.9, \pm 0.1, 0\}$ . The final reward is calculated as  $R = \text{sgn}(\text{sparse}) * 1 + \text{sgn}(\text{dense}) * 0.1$ , where  $\text{sgn}$  denotes the sign function. This is to keep the reward within a reasonable range, allowing the LLM to focus on various scenarios that need to be considered in the reward function, rather than trivial tasks like adjusting the reward value.

### 3.2. Reward Critic

In practice, it is difficult for LLM to generate a relatively complete reward function in the beginning. There may be errors in understanding parameter formats and data types (**syntax errors**), failure to consider game-specific informa-

tion, or misunderstanding of tasks (**semantic errors**), etc.

In order to eliminate above errors that are not easy to find, we design a LLM based Reward Critic to automatically review the designed reward function. In addition to checking for syntax errors, Reward Critic is also asked to check the quality of the reward function to further eliminate semantic errors. Specifically, we require Reward Critic to check whether the current code implementation matches its thoughts, whether it meets the reward function design requirements, and whether it takes game information into account. If the review fails, the Critic will provide a critique, and the Reward Designer will then modify the reward function based on the criterion and submit it for review again. The above process is repeated up to 3 times.

If an error occurs during the execution of the reward function in the process of interacting with the environment, the Python traceback of the error message will be fed back to Reward Designer for modification. These errors may include misunderstandings of input parameters, list index out of range, uninitialized keys in dictionaries, and other such issues. Some runtime errors only appear during the actual execution of the code.

### 3.3. Trajectory Analyzer

LLMs have the ability to understand environmental information and task instructions through in-context prompts to generate dense rewards. However, this zero-shot approach completely relies on LLMs understanding of the task and imagination of the problems it may face, and it is difficult to ensure the effectiveness of the designed reward. Take the the yellow highlighted part in Figure 1 as an example, in the initially designed reward function, Reward Designer does not consider the situation where the agent would encounter lava and be burned to death. Thus, in order to introduce empirical improvements to the designed dense reward, we propose to use LLMs, named as Trajectory Analyzer, to summarize the historical information of the interaction between the trained agent and the environment and use it to guide the revision of the reward function. The division of labor of Reward Designer and Trajectory Analyzer allows for independent operations of data analysis and reward function updates. Trajectory Analyzer does not need to know the details of the reward function, and Reward Designer does not need to process complex trajectory data.

Specifically, the current trained model is used to interact with the environment and obtain  $K$  trajectories. Then, we truncate these trajectories and use a LLM to summarize the observations of the last consecutive  $L$  frames of each failed trajectory to automatically infer its possible failure reasons. Based on the analysis of the reasons for the failure, the LLM Trajectory Analyzer is asked to propose key points that Reward Designer needs to consider in the next round of reward function revision. For instance, failure sce-





narios where punishment is not considered, misalignment of dense reward and sparse reward causes the agent’s behavior to deviate from the final goal, etc.

Figure 2 shows an example of multiple rounds of improving the reward function during the search for diamonds. In the first step, through analysis of the trajectory, Trajectory Analyzer finds that the agent would opportunistically find a shortcut to increase the reward, that is, move back and forth to deceive the reward function into thinking that the agent is moving actively. Therefore, the Reward Designer modifies the code snippet that encourages the agent to move, *i.e.* encourage the agent to appear in unvisited locations as much as possible. Although the initially designed reward function has taken into account the penalty for the loss of the agent’s health, the agent still cannot effectively learn to avoid lava. When modifying the reward function in the second round, Trajectory Analyzer discovers through the failed trajectory that the agent may die from lava, so it is suggested that Reward Designer increase the penalty for encountering lava, as shown in the step 2 update in Figure 2. According to the interactive experience, Reward Designer explicitly punishes the continuous appearance of lava in the field of view. However, the excessive punishment of lava caused the agent to choose to turn its perspective upward or downward to avoid the appearance of lava in the visible range, making it impossible for the agent to continue effective exploration, which deviates from the ultimate goal. To this end, Reward Designer further constrain the agent’s perspective in step 3, so that the lava disappeared from the agent’s perspective by turning left/right while continuing to search for diamonds, which is the desired strategy. Figure 3(a) shows the successful trajectory of avoiding lava: The agent sees the lava after breaking the stone ahead using iron pickaxe, and then turn left to avoid the lava through the mining tunnel.

## 4. Experiment

### 4.1. Environment Setup

We mainly use the harvest mode in the MineDojo [14] environment to verify the model’s ability to play Minecraft. The training pseudo code of Auto MC-Reward is shown in Algorithm 1. We set up the following challenging tasks for model performance comparison and ablation study:

- **Exploring diamond ore  on the 11-th floor underground:** Initially, the agent  is equipped with an iron pickaxe on the 11-th floor underground. When the diamond ore is within the visible range and the distance is less than 2 distance units, the task is deemed completed. The difficulty of the task lies in the fact that diamonds are very rare, lava frequently appears during exploration, leading to death, and the maximum number of steps is limited to 60,000. When steps exceed the limit, the tra-

---


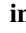


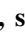

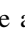

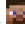
### Algorithm 1 Auto MC-Reward Training Pseudo Code

---

**Require:** Task ( $T$ ), Initial Agent ( $A_0$ ), Environment (Env), Max number of Critic reviews ( $N_{\text{Critic}}$ )  
**Ensure:** Final Agent ( $A_N$ ), Final Reward ( $R_N$ )  
 Summary = None  
 Critique = None  
 $R_0 = \text{None}$   
**for**  $i = 1, \dots, N$  **do**  
    $R_i = \text{RewardDesigner}(\text{Summary}, \text{Critique}, T, R_{i-1})$   
   **for**  $j = 1, \dots, N_{\text{Critic}}$  **do**  
     Critique, Done = RewardCritic( $R_i$ )  
     **if** Done **then**  
       break  
     **else**  
        $R_i = \text{RewardDesigner}(\text{Summary}, \text{Critique}, T, R_i)$   
     **end if**  
**end for**  
 $A_i = \text{TrainAgent}(A_0, R_i, \text{Env}, T)$   
 $\text{Traj}_i, \text{Stat}_i = \text{Eval}(\text{Env}, A_i)$   
 Summary = TrajectoryAnalyzer( $\text{Traj}_i, \text{Stat}_i$ )  
 Critique = None  
**end for**

---

jectory is considered failed. Long-term exploration can demonstrate the advantages of dense rewards.

- **Approaching tree  in plains biome **: The task is considered successful if the tree is within the agent’s  visible range and the distance is less than 1 distance unit. The difficulty of the task lies in the fact that the trees are very sparse on plains, which is extremely detrimental to sparse reward functions. The maximum number of steps is limited to 2,000 steps.
- **Approaching specific animal (e.g. cow , sheep ) in plains biome **: The task is considered successful if the animal is within the agent’s  visible range and the distance is less than 2 distance unit. The difficulty of the task is that the animals are constantly moving. The maximum number of steps is limited to 2,000 steps.
- **Obtaining diamond **: The agent  needs to complete the whole process of mining diamonds, including key behaviors such as finding and obtaining materials on the surface, crafting, digging down, going back to the ground, and mining stone/iron ore/diamond ore. The tech tree is shown in Figure 4.

### 4.2. Implementation Details

**LLM Prompt.** The components of the input prompts for Trajectory Analyzer include task description, game information, statistical metrics, and information on failed trajectories. Components of the input prompts for Reward Designer and Reward Critic includes task description, game information, input parameters, and reward function requirements and format. We use GPT-4 [36] for all the LLM components, and set temperature to 0.3. Since the LLMs are only used once for each whole agent training instead of each action, their computation overhead is negligible.



Figure 3. The trajectories of the new behaviors. (a) Avoid lava when exploring for diamond ore. (b) Attack cow in plains.

- **Instruction:** Instructions on initializing, updating and handling execution error of reward function for Reward Designer, reviewing function for Reward Critic, and analyzing trajectory for Trajectory Analyzer.
- **Task description:** The objective, initial conditions, success criteria, and task flow. For example, for the explore diamond task, the objective is “to find and approach a diamond, achieving a high success rate while avoiding death.” The initial condition is “agent at y level 11 with an iron pickaxe.” The success criteria is “being less than 1 meter from the nearest diamond block”, and the task flow is “horizontally explore to find a diamond, face it, and approach it”. In the task description, we do not provide prior game strategy information (task challenges, DFS exploration strategies, or avoiding lava, etc.) to ensure the method’s versatility.
- **Game information:** Game version, block names, field of view, action space, and units of measurement. Game information provides knowledge about the game’s simulation environment, not game strategy.
- **Statistical metrics and information on failed trajectories:** success rates, and actions sequences, reward sequences, final inventory and nearby blocks of  $K = 10$  failed trajectories. If a trajectory exceeds  $L = 32$  steps, it is truncated to the last 32 steps.
- **Input parameters:** The nearest distance of each block type within the visible range in the current and previous steps, changes in inventory between adjacent steps, health points, and the agent’s location in each past step. The memory is also provided as an input parameter for storing information to monitor changes across different steps. We provide explanation and examples of the parameters in the input prompt.
- **Reward function requirements and format:** We require the Designer to write a dense function and a sparse function, and consider only the sign of the two functions’ return values, not the magnitude. The detail of the scale constraints is in Section 3.1.

**Imitation Learning Details.** When large labeled datasets

do not exist, the canonical strategy for training capable agents is RL, which is inefficient and expensive to sample for hard-exploration problems [3, 4, 21], e.g. mining diamond in Minecraft. Therefore, in order to more efficiently explore the effectiveness of the LLM-based reward function design mechanism proposed in this paper, we pre-trained some foundation models through imitation learning as done by VPT [3]. Specifically, we use GITM [53] to continuously perform Diamond Mining task and record important observation data of each frame, such as RGB, action, inventory, GPS, compass, structured actions, etc. In the end, we collect about 11 million image data, totaling about 153 hours (the control frequency is 20 Hz) of game videos. Subsequently, we train these data through fully supervised learning by using Impala CNN [13] and Transformer [44] as backbone, and obtained several foundation models. The main differences between the foundation models are different biomes (forest and plains), temporal frames (16 and 128), and whether goal embedding is used. In subsequent experiments, these foundation models were used in two different purposes:




- Give the RL model preliminary basic Minecraft gameplay capabilities, e.g. forward/back, turn left/right, attack, etc. For some tasks that have not been learned (e.g. approaching cows in Figure 3(b) or not learned well (e.g. avoid lava in Figure 3(a), explore tree on plains) in imitation learning, RL algorithms can be studied more efficiently.
- In the Diamond Mining task, the diamond collection success rate, lava escape rate, death rate, etc. between the RL model and the imitation learning model are compared to demonstrate the superiority of the proposed method.

**RL Training Details.** We use proximal policy optimization (PPO) algorithm [39] with generalized advantage estimation (GAE) [38] to train our RL model. We use  $\gamma = 0.99$  and  $\lambda = 0.95$  for all of our experiments, and the total training frames is 256,000. To prevent catastrophically forgetting or overly aggressive policy update during RL training, we follow VPT [3] to apply an auxiliary Kullback-Leibler (KL) divergence loss between the RL model and the frozen pre-trained policy. We also normalize the reward based on the trajectory returns to constrain the gradient scales of different tasks. See Appendix for details.

### 4.3. Main Results

**Baselines.** We compare our Auto MC-Reward against the following methods:

- **Naive Handcraft:** The agent keeps moving (and mining for diamond exploring task) in one direction with a small probability of turning left/right.
- **Imitation Learning:** Our foundation model pre-trained with GITM-generated data, as introduced in Section 4.2.
- **RL with Sparse Reward:** Use only the reward from the original environment, i.e. only receives a reward when

Table 1. Comparison with other reward methods on three Minecraft tasks. Max steps for exploring tree  and cow  are set to 2000. <sup>†</sup>Sparse reward receives a low death rate because it is often stuck in the same place or move in a small area without encountering lava .



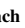

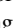
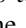
Method	Reward	Explore Diamond Ore  Underground				Approach Tree  on Plains		Approach Cow  on Plains	
		avg. dist. $\uparrow$	death (%) $\downarrow$	lava escape (%) $\uparrow$	succ. (%) $\uparrow$	avg. step $\downarrow$	succ. (%) $\uparrow$	avg. step $\downarrow$	succ. (%) $\uparrow$
Naive Handcraft	-	85.7	74.3	1.5	18.6	1993	2.1	1956	10.8
Imitation Learning	-	102.2	55.6	46.8	38.9	1988	2.5	1772	22.4
RL	Sparse	16.8	1.5 <sup>†</sup>	0	0.5	1936	4.3	1854	12.6
RL	Dense (Curiosity)	102.6	55.1	46.0	39.3	1672	45.8	1477	13.7
RL	Dense (Self-Imitation)	104.0	54.8	47.2	39.7	1532	42.5	1280	23.5
RL	Dense (MineCLIP)	105.9	54.0	47.8	40.5	1022	65.6	1206	44.9
Ours	Dense (LLM)	<b>142.8</b>	<b>45.2</b>	<b>70.0</b>	<b>45.2</b>	<b>972</b>	<b>73.4</b>	<b>1134</b>	<b>56.3</b>

Table 2. Comparison with previous methods on success rates of obtaining diamond . We list observations that are used in the inference phase. Auto MC-Reward achieves a remarkable success rate without exploiting unfair information (*i.e.* Lidar and Voxel) during inference.

Method	Controller	Observation	Diamond Succ. (%)
Human [3]	-	-	50.0
DreamerV3 [17]	RL	RGB, Status	0.01
DEPS [47]	IL	RGB, Status, Voxel	0.6
VPT [3]	IL + RL	RGB	20.0
GITM [53]	Handcraft	Lidar, Voxel, Status	55.0
Ours	IL (GITM-guided)	RGB, GPS	28.8
Ours	IL + RL	RGB, GPS	36.5

Table 3. Ablations on Reward Critic and Trajectory Analyzer for explore diamond ore  task. The first row corresponds to using the sparse reward from the original environment. <sup>†</sup>Sparse reward receives a low death rate because it is often stuck in the same place or move in a small area without encountering lava .

Designer	Critic	Analyzer	Avg. Dist. $\uparrow$	Death $\downarrow$	Lava Esc. $\uparrow$	Succ. $\uparrow$
			16.8	1.5 <sup>†</sup>	0	0.5
✓			75.8	58.2	30.4	35.1
✓	✓		95.2	49.3	40.7	40.5
✓		✓	130.6	47.8	64.8	43.1
✓	✓	✓	<b>142.8</b>	<b>45.2</b>	<b>70.0</b>	<b>45.2</b>

the success criteria is completed.

- **RL with Curiosity Dense Reward [37]:** Encourage the agent to discover and learn about parts of the environment that it has not encountered before.
- **RL with Self-Imitation Dense Reward [35]:** Encourage the agent to replicate its past actions that led to high rewards.
- **RL with MineCLIP [14] Dense Reward:** Use MineCLIP to calculate the dense reward based on the similarity between RGB frames and task objectives.








**Results on Diamond Ore  Exploring Task.** For the plain imitation learning model, fitting the training data makes it lack the awareness of avoiding lava, so it often dies in lava during the search for diamonds, and only has 38.9% success rate under the limit of 60,000 steps, as shown in Table 1. In contrast, our Auto MC-Reward makes the agent realize the importance of avoiding lava by continuously im-




Figure 4. The tech tree of obtaining diamond. The green squares are tasks to be optimized with Auto MC-Reward, *i.e.* obtaining log , cobblestone , iron ore  and diamond .

proving the dense reward function, and the final success rate has increased to 45.2% with 70% lava escape success rate. Figure 3(a) demonstrates good awareness of avoiding lava. Based on the same reinforcement learning algorithm, the disadvantages of sparse reward functions in long-horizon tasks are undoubtedly revealed. By watching the videos of collected trajectories, we find that using sparse functions often leads to irreversible behavior, such as being unable to break the surrounding ores to move due to maintaining a head-up posture. Although a low death rate of 1.5% is achieved, the actual average moving distance is only 16.8, and the success rate is only 0.5%. Due to the similar scenes underground, MineCLIP cannot give differentiated rewards, so its performance is close to the initial imitation learning model. Other baselines, like curiosity and self-imitation dense reward, also have mediocre performance and the success rate has not been significantly improved.

**Results on Tree  Approaching Task.** Since trees are extremely sparse on the plain, the imitation learning model and the RL model with sparse reward cannot perform well, with only 2.5% and 4.3% success rates respectively, and their average action steps are close to the maximum limit. MineCLIP dense reward receives a success rate of 65.6% since it can provide positive reward when tree is visible. Curiosity and self-imitation methods also achieve better results than imitation learning. For Auto MC-Reward, Reward Designer uses a strategy of giving positive rewards for getting closer and deducting rewards for going away, so that the agent learns to slowly approach the target, ultimately achieving 73.4% success rate with only 972 average steps.

**Results on Cow  Approaching Task.** The task of exploring for cows does not appear in the training data of imitation

learning, so the zero-shot ability on this task is not ideal, with about 22.4% success rate and average steps close to the maximum limit. By checking the videos, we find most of the successful cases are due to good luck without intention to actively approach the target. The same experimental conclusion is also obtained in the experiment of sparse reward function. Similar to the Tree Approaching Task, the superior dense reward function design mechanism makes our agent 43.7% (56.3% vs. 12.6%) higher than sparse reward, as listed in Table 1. Another dense reward MineCLIP also shows strong performance in this task, but due to the need to calculate the similarity of images and texts at all times during training, the efficiency is unacceptable.

**Results on Obtaining Diamond** . We verify the proposed method on a more difficult task, that is, the tech tree of collecting diamonds, as shown in Figure 4. As mentioned before, our foundation imitation learning model already has a certain ability from birth to diamond mining. We use the proposed method to optimize several key tasks in the process to increase the success rate of final diamond acquisition. The green parts in Figure 4 are the tasks that need to be optimized, *i.e.* obtaining log, cobblestone, iron ore and diamond. We conduct experiments in two biomes in Minecraft, and the cumulative success rate is shown in Figure 5. Specifically, the lower death rate allows our agent to have a higher success rate in mining iron ore and diamond, and ultimately achieves 36.5% success rate on forest biome, which is 7.7% higher than the imitation learning model. As for plains, the difficulty of obtaining log makes the imitation learning model unable to complete any tasks. Auto MC-Reward overcomes the difficulty of obtaining log, thus achieving a 28.1% success rate in obtaining diamonds. Table 2 provides a rough comparison of several different methods on the task of mining diamonds. We achieve a high success rate without exploiting unfair information (*i.e.* Lidar and Voxel) during the inference phase.

#### 4.4. Ablation studies

**Effectiveness of Reward Designer.** The first row of Table 3 is an RL experiment with a sparse reward function. As mentioned before, it cannot explore diamonds normally. After adding Reward Designer, it regained the ability to explore under a dense reward function.

**Effectiveness of Reward Critic.** As listed in Table 3, the success rate of exploring diamonds has increased from 35.1% to 40.5% by adding Reward Critic, because it can reduce the syntax and semantic errors in the code, making the training process more effective and sufficient. For example, the Trajectory Analyzer concludes that the agent died in lava and asks the Reward Designer to add relevant penalties. However, without being checked by Critic for semantic errors, it is possible that the added code snippet uses the word “magma” instead of the correct one “lava”. This will result

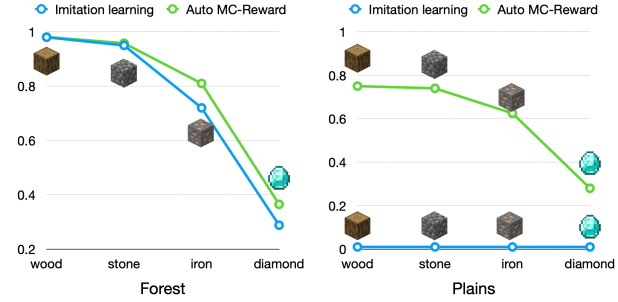




Figure 5. Cumulative success rates for 4 key items of obtaining diamond on forest  and plains . In terms of diamond, the performance comparison between imitation learning and Auto MC-Reward in two biomes are: 28.8% vs. 36.5%, and 0% vs. 28.1%.

in insufficient learning of lava avoidance, which is reflected in a 2.1% (43.1% vs. 45.2%) success rate difference.

**Effectiveness of Trajectory Analyzer.** As observed in Table 3, Trajectory Analyzer is the key to improve the success rate of completing tasks. It summarizes the reasons for failure to be fed into Reward Designer, allowing it to iteratively modify an appropriate dense reward function to guide the agent to overcome difficulties. In terms of Diamond Exploring Task, Trajectory Analyzer provides timely feedback on the potential risks of lava, which greatly improves the survival rate and moving distance, ultimately improving the success rate from 40.5% to 45.2%.

## 5. Conclusion

We proposed Auto MC-Reward, an automated dense reward design framework for addressing challenges caused by sparse reward and complex environment of Minecraft. It addresses the issue of sparse rewards by leveraging LLMs to automatically generate dense reward functions, enhancing learning efficiency. The system consists of three key components: Reward Designer, Reward Critic, and Trajectory Analyzer, which are used for the design, verification and analysis of the reward function respectively. Its capabilities are validated through experiments, demonstrating a remarkable improvement in complex tasks in Minecraft. Future work may deal with the limited trajectory length for analysis (last 32 frames) due to the context length of LLMs, which hinders the analysis of long-term failures (e.g., not exploring new areas, circling around lava). Auto MC-Reward humbly contributes to more effective learning in complex tasks through its automated dense reward function design. We hope it can pave the way for further research in reinforcement learning and its real-world applications.

## Acknowledgement

This work is supported by the National Key R&D Program of China (NO. 2022ZD0161300, NO. 2022ZD0160100), by the National Natural Science Foundation of China (62376134).



## References

- [1] Artemij Amiranashvili, Nicolai Dorka, Wolfram Burgard, Vladlen Koltun, and Thomas Brox. Scaling imitation learning in minecraft. *arXiv preprint arXiv:2007.02701*, 2020. [2](#)
- [2] Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay. *Advances in neural information processing systems*, 30, 2017. [2](#)
- [3] Bowen Baker, Ilge Akkaya, Peter Zhokov, Joost Huizinga, Jie Tang, Adrien Ecoffet, Brandon Houghton, Raul Sampe-dro, and Jeff Clune. Video pretraining (vpt): Learning to act by watching unlabeled online videos. *Advances in Neural Information Processing Systems*, 35:24639–24654, 2022. [2](#), [6](#), [7](#), [12](#), [15](#)
- [4] Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019. [6](#), [12](#)
- [5] Albert Bou, Matteo Bettini, Sebastian Dittert, Vikash Kumar, Shagun Sodhani, Xiaomeng Yang, Gianni De Fabritiis, and Vincent Moens. Torchrl: A data-driven decision-making library for pytorch, 2023. [15](#)
- [6] Yuri Burda, Harri Edwards, Deepak Pathak, Amos Storkey, Trevor Darrell, and Alexei A Efros. Large-scale study of curiosity-driven learning. *arXiv preprint arXiv:1808.04355*, 2018. [2](#)
- [7] Thomas Carta, Pierre-Yves Oudeyer, Olivier Sigaud, and Sylvain Lamprier. Eager: Asking and answering questions for automatic reward shaping in language-guided rl. *Advances in Neural Information Processing Systems*, 35: 12478–12490, 2022. [3](#)
- [8] Maxime Chevalier-Boisvert, Bolun Dai, Mark Towers, Rodrigo de Lazcano, Lucas Willems, Salem Lahlou, Suman Pal, Pablo Samuel Castro, and Jordan Terry. Minigrid & mini-world: Modular & customizable reinforcement learning environments for goal-oriented tasks. *CoRR*, abs/2306.13831, 2023. [11](#)
- [9] Hao-Tien Lewis Chiang, Aleksandra Faust, Marek Fiser, and Anthony Francis. Learning navigation behaviors end-to-end with autorl. *IEEE Robotics and Automation Letters*, 4(2): 2007–2014, 2019. [3](#)
- [10] Yuqing Du, Ksenia Konyushkova, Misha Denil, Akhil Raju, Jessica Landon, Felix Hill, Nando de Freitas, and Serkan Cabi. Vision-language models as success detectors. *arXiv preprint arXiv:2303.07280*, 2023. [2](#)
- [11] Yuqing Du, Olivia Watkins, Zihan Wang, Cédric Colas, Trevor Darrell, Pieter Abbeel, Abhishek Gupta, and Jacob Andreas. Guiding pretraining in reinforcement learning with large language models. *arXiv preprint arXiv:2302.06692*, 2023. [2](#), [3](#)
- [12] Sean C Duncan. Minecraft, beyond construction and survival. *Well Played: A Journal on Video Games, Value and Meaning*, 1(1), 2011. [2](#)
- [13] Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Vlad Mnih, Tom Ward, Yotam Doron, Vlad Firoiu, Tim Harley, Iain Dunning, et al. Impala: Scalable distributed deep-rl with importance weighted actor-learner architectures. In *International conference on machine learning*, pages 1407–1416. PMLR, 2018. [6](#), [13](#)
- [14] Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. Minedojo: Building open-ended embodied agents with internet-scale knowledge. *Advances in Neural Information Processing Systems*, 35: 18343–18362, 2022. [2](#), [5](#), [7](#)
- [15] Aleksandra Faust, Anthony Francis, and Dar Mehta. Evolving rewards to automate reinforcement learning. *arXiv preprint arXiv:1905.07628*, 2019. [3](#)
- [16] William H Guss, Brandon Houghton, Nicholay Topin, Phillip Wang, Cayden Codell, Manuela Veloso, and Ruslan Salakhutdinov. Minerl: A large-scale dataset of minecraft demonstrations. *arXiv preprint arXiv:1907.13440*, 2019. [2](#)
- [17] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models. *arXiv preprint arXiv:2301.04104*, 2023. [2](#), [7](#)
- [18] Jie Huang and Kevin Chen-Chuan Chang. Towards reasoning in large language models: A survey. *arXiv preprint arXiv:2212.10403*, 2022. [2](#)
- [19] Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, et al. Inner monologue: Embodied reasoning through planning with language models. *arXiv preprint arXiv:2207.05608*, 2022. [2](#)
- [20] Max Jaderberg, Volodymyr Mnih, Wojciech Marian Czarnecki, Tom Schaul, Joel Z Leibo, David Silver, and Koray Kavukcuoglu. Reinforcement learning with unsupervised auxiliary tasks. *arXiv preprint arXiv:1611.05397*, 2016. [2](#)
- [21] Max Jaderberg, Wojciech M Czarnecki, Iain Dunning, Luke Marris, Guy Lever, Antonio Garcia Castaneda, Charles Beattie, Neil C Rabinowitz, Ari S Morcos, Avraham Ruderman, et al. Human-level performance in 3d multiplayer games with population-based reinforcement learning. *Science*, 364 (6443):859–865, 2019. [6](#), [12](#)
- [22] Matthew Johnson, Katja Hofmann, Tim Hutton, and David Bignell. The malmo platform for artificial intelligence experimentation. In *Ijcai*, pages 4246–4247, 2016. [2](#)
- [23] Ingmar Kanitscheider, Joost Huizinga, David Farhi, William Hebgen Guss, Brandon Houghton, Raul Sampe-dro, Peter Zhokhov, Bowen Baker, Adrien Ecoffet, Jie Tang, et al. Multi-task curriculum learning in a complex, visual, hard-exploration domain: Minecraft. *arXiv preprint arXiv:2106.14876*, 2021. [2](#)
- [24] Tejas D Kulkarni, Karthik Narasimhan, Ardavan Saeedi, and Josh Tenenbaum. Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation. *Advances in neural information processing systems*, 29, 2016. [2](#)
- [25] Minae Kwon, Sang Michael Xie, Kalesha Bullard, and Dorsa Sadigh. Reward design with language models. *arXiv preprint arXiv:2303.00001*, 2023. [3](#)
- [26] Pawel Ladosz, Lilian Weng, Minwoo Kim, and Hyondong Oh. Exploration in deep reinforcement learning: A survey. *Information Fusion*, 85:1–22, 2022. [2](#)

- [27] Andrew Levy, George Konidaris, Robert Platt, and Kate Saenko. Learning multi-level hierarchies with hindsight. *arXiv preprint arXiv:1712.00948*, 2017. 2
- [28] Jing Li, Xinxin Shi, Jiehao Li, Xin Zhang, and Junzheng Wang. Random curiosity-driven exploration in deep reinforcement learning. *Neurocomputing*, 418:139–147, 2020. 2
- [29] Jessy Lin, Daniel Fried, Dan Klein, and Anca Dragan. Inferring rewards from language in context. *arXiv preprint arXiv:2204.02515*, 2022. 3
- [30] Zichuan Lin, Junyou Li, Jianing Shi, Deheng Ye, Qiang Fu, and Wei Yang. Juewu-mc: Playing minecraft with sample-efficient hierarchical reinforcement learning. In *International Joint Conference on Artificial Intelligence (IJCAI)*, 2022. 2
- [31] Yecheng Jason Ma, Shagun Sodhani, Dinesh Jayaraman, Osbert Bastani, Vikash Kumar, and Amy Zhang. Vip: Towards universal visual reward and representation via value-implicit pre-training. *arXiv preprint arXiv:2210.00030*, 2022. 2
- [32] Yecheng Jason Ma, William Liang, Vaidehi Som, Vikash Kumar, Amy Zhang, Osbert Bastani, and Dinesh Jayaraman. Liv: Language-image representations and rewards for robotic control. *arXiv preprint arXiv:2306.00958*, 2023. 2
- [33] Yecheng Jason Ma, William Liang, Guanzhi Wang, De-An Huang, Osbert Bastani, Dinesh Jayaraman, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Eureka: Human-level reward design via coding large language models. *arXiv preprint arXiv:2310.12931*, 2023. 3
- [34] Hangyu Mao, Chao Wang, Xiaotian Hao, Yihuan Mao, Yiming Lu, Chengjie Wu, Jianye Hao, Dong Li, and Pingzhong Tang. Seihai: A sample-efficient hierarchical ai for the minerl competition. In *Distributed Artificial Intelligence: Third International Conference, DAI 2021, Shanghai, China, December 17–18, 2021, Proceedings 3*, pages 38–51. Springer, 2022. 2
- [35] Junhyuk Oh, Yijie Guo, Satinder Singh, and Honglak Lee. Self-imitation learning. In *International Conference on Machine Learning*, pages 3878–3887. PMLR, 2018. 2, 7
- [36] OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 2, 5
- [37] Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-driven exploration by self-supervised prediction. In *International conference on machine learning*, pages 2778–2787. PMLR, 2017. 2, 7
- [38] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015. 6, 15
- [39] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017. 6, 15
- [40] Kyoungwon Seo, Joice Tang, Ido Roll, Sidney Fels, and Dongwook Yoon. The impact of artificial intelligence on learner-instructor interaction in online learning. *International journal of educational technology in higher education*, 18(1):1–23, 2021. 2
- [41] Alexey Skrynnik, Aleksey Staroverov, Ermek Aitygulov, Kirill Aksenov, Vasilii Davydov, and Aleksandr I Panov. Hierarchical deep q-network from imperfect demonstrations in minecraft. *Cognitive Systems Research*, 65:74–78, 2021. 2
- [42] Chen Tessler, Shahar Givony, Tom Zahavy, Daniel Mankowitz, and Shie Mannor. A deep hierarchical approach to lifelong learning in minecraft. In *Proceedings of the AAAI conference on artificial intelligence*, 2017. 2
- [43] Adrienne L Tierney and Charles A Nelson III. Brain development and the role of experience in the early years. *Zero to three*, 30(2):9, 2009. 2
- [44] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 6, 13
- [45] Alexander Sasha Vezhnevets, Simon Osindero, Tom Schaul, Nicolas Heess, Max Jaderberg, David Silver, and Koray Kavukcuoglu. Feudal networks for hierarchical reinforcement learning. In *International Conference on Machine Learning*, pages 3540–3549. PMLR, 2017. 2
- [46] Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*, 2023. 2
- [47] Zihao Wang, Shaofei Cai, Guanzhou Chen, Anji Liu, Xiaojian Ma, and Yitao Liang. Describe, explain, plan and select: interactive planning with llms enables open-world multi-task agents. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. 2, 7
- [48] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837, 2022. 3
- [49] Tianbao Xie, Siheng Zhao, Chen Henry Wu, Yitao Liu, Qian Luo, Victor Zhong, Yanchao Yang, and Tao Yu. Text2reward: Automated dense reward function generation for reinforcement learning. *arXiv preprint arXiv:2309.11489*, 2023. 3
- [50] Wenhao Yu, Nimrod Gileadi, Chuyuan Fu, Sean Kirmani, Kuang-Huei Lee, Montse Gonzalez Arenas, Hao-Tien Lewis Chiang, Tom Erez, Leonard Hasenclever, Jan Humplik, et al. Language to rewards for robotic skill synthesis. *arXiv preprint arXiv:2306.08647*, 2023. 3
- [51] Haoqi Yuan, Chi Zhang, Hongcheng Wang, Feiyang Xie, Penglin Cai, Hao Dong, and Zongqing Lu. Plan4mc: Skill reinforcement learning and planning for open-world minecraft tasks. *arXiv preprint arXiv:2303.16563*, 2023. 2
- [52] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023. 2
- [53] Xizhou Zhu, Yuntao Chen, Hao Tian, Chenxin Tao, Weijie Su, Chenyu Yang, Gao Huang, Bin Li, Lewei Lu, Xiaogang Wang, et al. Ghost in the minecraft: Generally capable agents for open-world environments via large language models with text-based knowledge and memory. *arXiv preprint arXiv:2305.17144*, 2023. 2, 6, 7, 12