

---

# LLM Augmented Hierarchical Agents

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Solving long-horizon, temporally-extended tasks using Reinforcement Learning  
2 (RL) is challenging, compounded by the common practice of learning without  
3 prior knowledge (or *tabula rasa* learning). Humans can generate and execute  
4 plans with temporally-extended actions and quickly learn to perform new tasks  
5 because we almost never solve problems from scratch. We want autonomous  
6 agents to have this same ability. Recently, LLMs have been shown to encode a  
7 tremendous amount of knowledge about the world and to perform impressive in-  
8 context learning and reasoning. However, using LLMs to solve real world problems  
9 is hard because they are not grounded in the current task. In this paper we exploit  
10 the planning capabilities of LLMs while using RL to provide learning from the  
11 environment, resulting in a hierarchical agent that uses LLMs to solve long-horizon  
12 tasks. Instead of completely relying on LLMs, they guide a high-level policy,  
13 making learning significantly more sample efficient. This approach is evaluated  
14 in simulation environments such as MiniGrid, SkillHack, and Crafter, and on a  
15 real robot arm in block manipulation tasks. We show that agents trained using our  
16 approach outperform other baselines methods and, once trained, don't need access  
17 to LLMs during deployment.

## 18 1 Introduction

19 Humans can generate and execute plans with temporally extended actions to perform complex tasks  
20 in a dynamic and uncertain world. We would like autonomous agents to have the same capabilities.  
21 Massive engineering efforts can lead to agents that are remarkably robust, such as the rovers in space,  
22 and surgical and industrial robots. In the absence of such resources, techniques such as Reinforcement  
23 Learning (RL) can be used to extract robust control policies from experience. However, RL has  
24 many challenges, such as exploration under sparse rewards, generalization, safety, etc. This makes  
25 it difficult to learn good policies in a sample efficient way. Popular ways to tackle these problems  
26 include using expert feedback [6, 24] and leveraging the hierarchical structure of complex tasks.  
27 There is significant prior work on learning hierarchical policies to break down tasks into smaller  
28 sub-tasks [22, 9, 2].

29 Hierarchical Reinforcement Learning (HRL) does indeed mitigate some of the problems mentioned  
30 above. However, as the number of options or skills increases, we face some of the same problems  
31 again. Using some form of supervision, such as providing details about the sub-tasks or intermediate  
32 rewards or high-level human guidance, is one approach [18, 14, 16].

33 One of the reasons that humans are so good at dealing with unfamiliar situations is that we almost  
34 never solve problems from scratch. Presented with a new task and a library of skills, we are able to  
35 choose a subset of skills that seem most relevant and explore from there. We might perform some  
36 trial and error exploration (as in RL), but we quickly learn the right subset of skills as well as the  
37 correct sequence in which they need to be executed. For example, the door handles on newer cars  
38 lie flat against the door, unlike most other car door handles in existence. That presents a problem

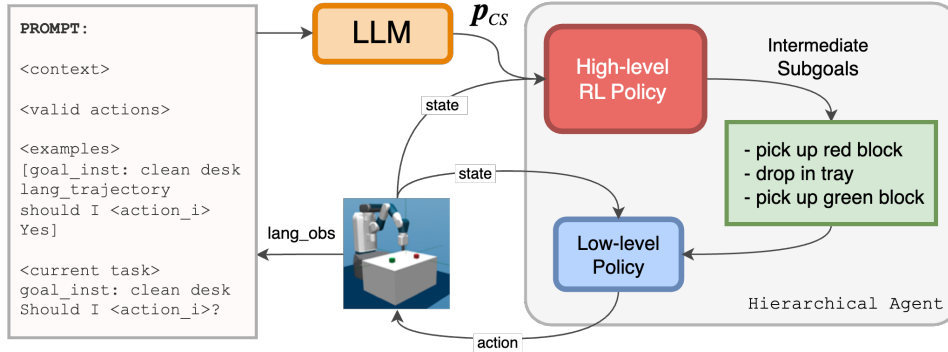


Figure 1: We use the LLM to guide the high-level policy and accelerate learning. The LLM is prompted with the context, some examples, and the current task and observation. We then use the LLM output to bias the high-level action selection

39 the first time you try to open one. Humans immediately narrow down to a few exploratory actions,  
 40 like trying to get a finger under the handle or pushing on it in different places. We don't, unlike most  
 41 RL algorithms might, tap the window or pull the side mirrors, as we believe that such options are  
 42 causally irrelevant based on deep world knowledge.

43 Large language models (LLMs) have been shown to encode a tremendous amount of knowledge  
 44 about the world by virtue of being trained on massive amounts of text. We hypothesize that this  
 45 knowledge can be leveraged to focus the training of hierarchical policies, making them significantly  
 46 more sample efficient. In particular, we explore how large pretrained language models can be used to  
 47 inject common sense priors into hierarchical agents.

48 In this approach we assume access to several low level skills. These can be, for example, engineered  
 49 planners or policies learned using RL and sub-task rewards. Based on a high-level task description  
 50 and current state, the LLMs guide the agent by suggesting the most likely courses of action. Instead  
 51 of random exploration, we use these suggestions to intelligently explore the various options. Because  
 52 LLMs are not grounded in the domain, they are only used to bias action selection and their influence  
 53 is reduced as training progresses. This results in a policy that can be deployed without depending  
 54 on the LLM at run time. We evaluate this approach on several simulated environments (MiniGrid  
 55 [5], SkillHack [17], and Crafter [11]), showing that it can learn to solve complex, long-horizon tasks  
 56 much faster than baseline methods. Experiments with a real robot arm in block manipulation tasks  
 57 using a tabular Q-learning version of the same algorithm show that it can learn policies much faster  
 58 with less experience in the domain. Our contributions are summarized as follows:

- 59 • an approach for using LLMs to guide exploration by extracting common sense priors;
- 60 • a hierarchical agent that uses these priors to solve long-horizon tasks;
- 61 • an evaluate of the framework in simulation as well as a simple real-world environment, show  
 62 that it performs significantly better than baselines;
- 63 • a discussion of (1) the advantages of our method compared prior work and (2) potential  
 64 future work.

## 65 2 Related Work

66 **Language and HRL** There is significant prior work on hierarchical RL where the standard MDP is  
 67 converted into a semi-Markov decision process (SMDP). The most common approach is to incorporate  
 68 temporally extended actions, also known as options or skills [2]. Typically, a low-level policy achieves  
 69 sub-tasks by executing primitive actions and a high-level policy plans over temporally extended  
 70 options or skills. Natural language is a popular way to specify sub-tasks and achieve generalization  
 71 due to its inherent compositionality and hierarchical structure [14, 18, 25, 12]. Most of these methods  
 72 specify or generate a high-level plan in natural language, which is then executed sequentially by a  
 73 separate low-level policy. These approaches face challenges when operating in high-dimensional

74 observation spaces. They also rely on manual data collection to train the high-level policies and  
75 therefore difficult to generalize to new tasks.

76 **RL and Foundation Models** Recently, large language models such as GPT-3 have been used to  
77 build agents capable of acting in the real world based on language instructions [4, 3]. The in-context  
78 learning and intelligent prompting strategies supported by these models have been used to design  
79 language-guided hierarchical agents. [13] use LLMs as zero-shot planners to enable embodied agents  
80 to act in real world scenarios. Similarly, [1] use LLMs along with affordance functions to generate  
81 feasible plans that guide a robot to achieve goals specified in natural language instructions. Our work  
82 is closely related to [7], where they improve exploration by using LLMs to provide intermediate  
83 rewards and encourage the agent to seek novel states.

## 84 3 Methods

### 85 3.1 Problem Statement

86 We consider a system that receives instructions in the form of natural language describing a task,  
87 similar to [1]. The instructions can be long, may contain warnings and constraints, and may not  
88 include all of the necessary individual steps. We also assume that the agent has access to a finite set of  
89 skills or sub-policies that can be executed in sequence to solve long-horizon tasks. These skills can be  
90 hand-coded, or trained using reinforcement learning or imitation learning with manual reward design.  
91 They must be accompanied by a simple description in natural language, such as "pick up red block"  
92 or "open blue door". They must also be able to detect sub-task completion to switch control back to  
93 the high-level policy. Given a finite set of options or skills, our objective is to obtain a high-level  
94 decision policy that selects among them.

### 95 3.2 Using LLMs to Guide High-level Policies

96 This section introduces our method for using LLMs to improve exploration in the high-level policy  
97 of an HRL system. The semantic knowledge and planning capabilities of LLMs improve high-level  
98 action selection given a task description and current state in the form of language. The core idea is to  
99 use LLMs to obtain a value that approximates the probability that a given skill or sub-task is relevant  
100 to achieve the larger goal. As mentioned earlier, each skill is accompanied by a language description  
101  $l_{skill}$  and the current trajectory is translated into language,  $l_{traj}$ . There is also a high level instruction,  
102  $l_{goal\_inst}$ , describing the larger goal along with optional constraints.

103 The LLM is used to evaluate the function  $f_{LLM}(l_{skill_i}, l_{goal\_inst}, l_{traj})$  for each skill at every high-  
104 level decision step. Essentially, the LLM answers the following question: given the task,  $l_{goal\_inst}$ ,  
105 and trajectory so far,  $l_{traj}$ , should we choose skill  $l_{skill}$ ? The output of the LLM, ‘yes’ or ‘no’, can  
106 easily be converted to an int (“0” or “1”). This kind of closed form question-answering prompt has  
107 been shown to work better than open ended prompts [7]. After evaluating this for each of the  $k$  skills,  
108 we get  $F_{LLM} = [f_{LLM_1}, f_{LLM_2}, f_{LLM_3}, \dots, f_{LLM_k}]$ . For example,  $F_{LLM} = [0, 1, 0, \dots, 0, 1, 0, 0]$ .  
109 A LOG SOFTMAX function is applied to these logits to get the common sense priors from the LLM  
110 denoted by  $p_{CS} = \log\_softmax(F_{LLM})$ . Relying entirely on  $p_{CS}$  is not enough to solve complex  
111 tasks. At the same time, using RL and exploring without any common sense intuition is inefficient.  
112 Therefore we still use RL and sparse rewards to obtain high-level policies but also use the common  
113 sense priors,  $p_{CS}$ , from the LLMs to guide exploration. More details about the RL algorithms used  
114 are in the Experiments section. The action selection in the exploration policy samples actions from a  
115 categorical distribution where the logits are obtained by the policy head processing the state. These  
116 logits are biased with the common sense priors  $p_{CS}$  and a weight factor  $\lambda$ . So the action selection  
117 looks like this:  $a = \text{Categorical}[\pi(s_t) + \lambda.p_{CS}(s_t)]$ . Here, the action  $a$  is the temporally extended  
118 macro action or the skill. The weight factor starts from  $\lambda = 1$  and is annealed gradually until it  
119 reaches zero by the end of training. This means that our trained agent does not continue to reply on  
120 the LLM during deployment. The process is summarized in Algorithm 1 and Figure 1.

121 **LLM Queries and Prompt Design.** We use the *gpt-3.5-turbo* GPT provided by OpenAI APIs. To re-  
122 duce the number of API calls, the LLM responses for all possible combinations of  $l_{goal\_inst}$  and  $l_{traj}$   
123 are cached. A simplified version of  $l_{traj}$  is used to denote the current trajectory history using the past  
124 two actions. The main prompt used in our experiments has the following structure `Goal: $l_{goal\_inst}$ ,  
125 So far I have:  $l_{traj}$ , Should I  $l_{skill_i}$ ? .` The LLM is shown a few examples of responses

---

**Algorithm 1**

---

```
high_inst  $\leftarrow$  high level goal in language
 $\pi_\theta \leftarrow$  high level policy
 $f_{LLM} \leftarrow$  common sense priors from LLM
procedure LLMxHRL(high_inst)
  init  $\pi_\theta$ 
  while  $\pi_\theta$  not converged do
    init  $\tau \leftarrow \{\}$ 
    for  $t \leftarrow 0$  to  $T$  do
       $p_{CS} \leftarrow f_{LLM}(high\_inst, \tau)$ 
       $a_t \leftarrow cat\_dist[\pi_\theta(\tau) + \lambda.p_{CS}(\tau)].sample()$ 
       $s_t, r_t \leftarrow act(a_t)$ 
       $\tau \leftarrow append(s_t, r_t, a_t)$ 
    end for
    update  $\pi_\theta$ 
  end while
  return  $\pi_\theta$ 
end procedure
```

---

126 to such queries and the prompt specifies that a one word Yes/No answer is required. Example prompts  
127 are in the Appendix.

## 128 4 Experiments

129 This section describes the experimental setup and results of testing the framework in three simulation  
130 environments and one real world robotic arm block manipulation task. The framework relies on  
131 communicating with the LLM using text. As mentioned earlier, each skill corresponds to a text  
132 description  $l_{skill_i}$  and the high level goal,  $l_{goal\_inst}$ . We assume access to a captioner which maps  
133 the current observation history to  $l_{traj}$ . This could be automated by using modern vision to language  
134 models such as [19], but that is left for future work. Instead we use a CLIP-based model along with  
135 an LLM in the experiments with a real robot to convert visual input to a discrete low dimensional  
136 state. More details about obtaining  $l_{traj}$  is in the Appendix. In each environment, our method is  
137 compared with baseline hierarchical agents without any guidance from LLMs, and an oracle and a  
138 SayCan-like agent without affordances.

Method	Description
LLM x HRL (ours)	Use LLMs to bias high-level action selection as explained in Section 3. Only receives reward in the end at task completion.
Vanilla HRL	A baseline hierarchical agent which has no guidance from the LLMs.
Shaped HRL	Same as the Vanilla HRL with no LLM guidance. But here agents receive shaped rewards for successful sub-task completion. Requires hand engineered reward functions.
Oracle	This is the upper bound. The high-level policy is an oracle state-machine which provides the right sub-tasks in the correct sequence. [10]
SayCan w/o Aff	A SayCan [1] like architecture but without an affordance function and blindly trusting the LLM. This method will depend on LLM access during deployment

Table 1: This table lists the all methods we compared

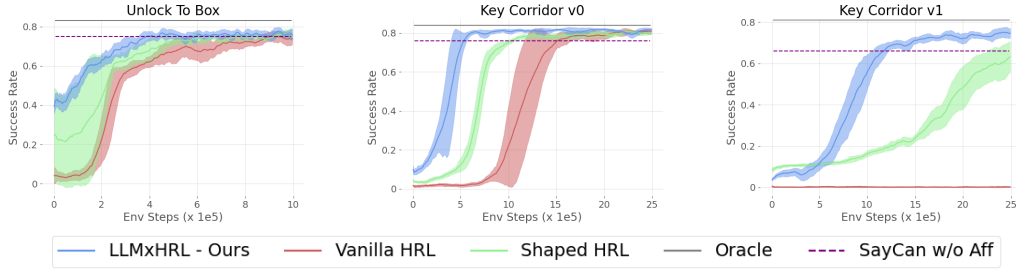


Figure 2: The plots show the success rate of different methods on the three tasks in the MiniGrid Environment.

#### 139 4.1 MiniGrid Experiments

140 **Setup** The experiments described in this section were performed on the MiniGrid environment by  
 141 [5], which is a simple grid world. The environment can be designed with multiple rooms with doors,  
 142 walls, and goal objects. These objects can have different colors and the agent and goal objects are  
 143 spawned at random locations. The action space is discrete which allows movement in the 4 compass  
 144 directions, opening and closing doors, and picking up and dropping objects. We designed multiple  
 145 tasks in this setup which can be broken down into smaller sub-tasks.

- 146 • *UnlockReach* task consists of a random object in a room which is behind a locked door.  
 147 The agent has to first find the right key based on the door color, unlock the door, and then  
 148 navigate to the goal object.
- 149 • *KeyCorridor v0* task consists of a corridor with multiple rooms on either side. A goal object  
 150 is inside a locked room whose key is in another room. The agent has to first find the key and  
 151 then unlock the door to ultimately reach the goal
- 152 • *KeyCorridor v1* is similar to v0, but some of the rooms have defective keys. The goal  
 153 instruction comes with the rooms to avoid. This task is much more difficult for standard  
 154 HRL methods.

155 Each task has a single reward that is only provided on successful task completion. The agents have  
 156 access to several temporally extended skills: *GoToObject*, *PickupObject*, *UnlockDoor*, *OpenBox*.  
 157 These are conditioned on the type and color of the objects. For example the *Object* may refer to  
 158 key, ball, or box, and the color can be *red*, *green*, *blue*, *yellow*, etc. These low-level sub-tasks were  
 159 pretrained and frozen using PPO [21] and manual reward specification. The high-level policies are  
 160 also trained using PPO where the . We compared against Vanilla HRL, Shaped HRL, an Oracle, and  
 161 a SayCan-like method as described in Table 1. The results are summarized in Figure 2. It’s clear that  
 162 our method outperform both the baseline HRL methods with and without shaped rewards. It is also  
 163 able to converge to the optimal policy much sooner than the other methods. The Oracle and SayCan  
 164 are not trained using RL and so we show their performance using the horizontal lines. Although they  
 165 are comparable to our method, one benefit of our method is that it does not rely on the LLM during  
 166 deployment.

#### 167 4.2 SkillHack

168 The NetHack Learning Environment [15] is an RL environment based on the classic game of NetHack.  
 169 It is notoriously difficult because of the large number on entities, actions, procedural generation, and  
 170 stochastic nature of the game. MiniHack [20] and SkillHack [17] are extensions of NetHack that  
 171 enable creation of custom levels and tasks. They are simpler than the full game while retaining most  
 172 of the interesting complexities. The SkillHack suite contains 16 skills such as *PickUp*, *Navigate*,  
 173 *Fight*, *Wear*, *Weild*, *Zap*, *Apply*, etc. More details are in the Appendix. These skills can be executed  
 174 sequentially to achieve larger tasks. We consider two such tasks - *Battle*, *FrozenLavaCross*.

- 175 • In the *Battle* task, the needs to *PickUp* a randomly placed Sword, *Wield* the Sword and  
 176 finally *Fight* and kill a Monster.

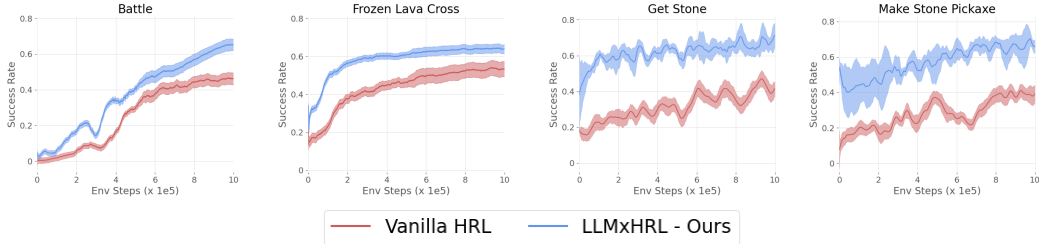


Figure 3: The 2 plots on the left show the success rate of different methods on the SkillHack - Battle and Frozen Lava Cross. The 2 plots on the right show the success rate of different methods on the Crafter - Get Stone and Make Stone Pickaxe

177 • In the *FrozenLavaCross* task, the needs to *PickUp* either a *WandOfCold* or a *FrostHorn*  
 178 based on what is available, then create a bridge across the lava with either *ZapWandOfCold*  
 179 or by *ApplyFrostHorn*. Finally, *NavigateLava* across your newly made bridge to reach the  
 180 staircase on the other side.

181 In this environment we compare against Vanilla HRL and an Oracle high-level policy. The low level  
 182 skills are trained using IMPALA [8] with the code provided by [17]. The high-level policy is also  
 183 trained using IMPALA where the policy skills are macro actions. As seen in the first two plots in  
 184 Figure 3, in both the tasks, *Battle* and *FrozenLavaCross*, our method clearly outperforms the HRL  
 185 agent without LLM guidance.

### 186 4.3 Crafter

187 Crafter [11] is a 2D version Minecraft which has the same complex dynamics but with a simpler  
 188 observation space and faster simulation speeds. Similar to Minecraft, it involves collecting and  
 189 building artifacts along an achievement tree. We modified the game slightly to make it easier by  
 190 slowing down health degradation and having fewer dangers to fight. We evaluated on two tasks that  
 191 have a natural hierarchical structure - *MakeWoodPickaxe* and *MakeStonePickaxe*. More details are in  
 192 the Appendix. Similar to our other experiments we pretrain policies for multiple skills using PPO.  
 193 The high level policy is then trained to select among these skills. The last two plots in 3 shows how  
 194 our method performs better than the baseline HRL method.

### 195 4.4 uArm Real Robot Experiments

196 We also test on a real robot arm on a simpler tabular Q-learning version of  
 197 our method. uArm Swift Pro [23] is an open-source desktop robotic arm.

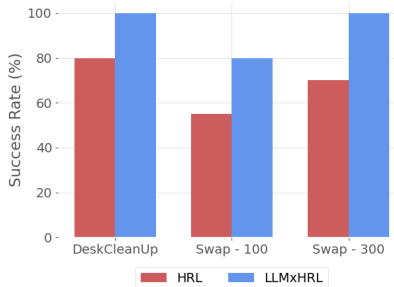


Figure 4: Robot Arm Results

210 the Appendix for more details.

212 Figure 4 show the results of our experiments. In the *DeskCleanUp* task, we have a 3 locations  
 213 where we have a tray and 2 blocks (red and green). The episode is initialized with blocks in random  
 214 locations. The goal is to pick up the blocks and place them in the tray, essentially cleaning the desk.  
 215 This task was trained for 100 episodes. In the *SwapBlocks* task again have 3 locations (or zones) and  
 216 2 blocks in 2 random locations. The goal is to swap the position of blocks. In the Figure 4, *Swap*

217 - 100 denotes performance after 100 episodes and *Swap* - 300, after 300 episodes. We can see that  
218 using LLMs to guide agent exploration give us better performance in fewer trials.

## 219 5 Discussion

220 In this work we present a framework for using LLMs to guide exploration in hierarchical agents.  
221 Instead of learning from random exploration without any prior knowledge, we use the LLMs to  
222 suggest high-level actions based on the task and current state. We evaluate our method on long  
223 horizon tasks which in simulation environments as well as a real robot. We show that our method  
224 perform better than baselines and does not require manual reward shaping. Moreover, once the agent  
225 is trained, we no longer depend on the LLM during deployment unlike some prior methods.

226 This work can be extended in several ways to make to more end to end. We Currently assume access  
227 to a function which provides us to language descriptions of the current trajectory and state. This can  
228 be automated using recent advancements in vision language models (VLM). It will also be interesting  
229 to extend this framework for more and one level or hierarchy to tackle longer tasks.

## 230 References

- 231 [1] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, C. Fu, K. Gopalakrishnan,  
232 K. Hausman, et al. Do as i can, not as i say: Grounding language in robotic affordances.  
233 *arXiv preprint arXiv:2204.01691*, 2022.
- 234 [2] P.-L. Bacon, J. Harb, and D. Precup. The option-critic architecture. In *Proceedings of the AAAI*  
235 *Conference on Artificial Intelligence*, volume 31, 2017.
- 236 [3] R. Bommasani, D. A. Hudson, E. Adeli, R. Altman, S. Arora, S. von Arx, M. S. Bernstein,  
237 J. Bohg, A. Bosselut, E. Brunskill, et al. On the opportunities and risks of foundation models.  
238 *arXiv preprint arXiv:2108.07258*, 2021.
- 239 [4] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam,  
240 G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in neural*  
241 *information processing systems*, 33:1877–1901, 2020.
- 242 [5] M. Chevalier-Boisvert, L. Willems, and S. Pal. Minimalistic gridworld environment for openai  
243 gym. <https://github.com/maximecb/gym-minigrid>, 2018.
- 244 [6] P. F. Christiano, J. Leike, T. Brown, M. Martic, S. Legg, and D. Amodei. Deep reinforcement  
245 learning from human preferences. In *Advances in Neural Information Processing Systems*,  
246 pages 4299–4307, 2017.
- 247 [7] Y. Du, O. Watkins, Z. Wang, C. Colas, T. Darrell, P. Abbeel, A. Gupta, and J. Andreas.  
248 Guiding pretraining in reinforcement learning with large language models. *arXiv preprint*  
249 *arXiv:2302.06692*, 2023.
- 250 [8] L. Espeholt, H. Soyer, R. Munos, K. Simonyan, V. Mnih, T. Ward, Y. Doron, V. Firoiu, T. Harley,  
251 I. Dunning, et al. Impala: Scalable distributed deep-rl with importance weighted actor-learner  
252 architectures. In *International conference on machine learning*, pages 1407–1416. PMLR,  
253 2018.
- 254 [9] R. Fruit and A. Lazaric. Exploration-exploitation in mdps with options. In *Artificial Intelligence*  
255 *and Statistics*, pages 576–584. PMLR, 2017.
- 256 [10] V. G. Goecks, N. Waytowich, D. Watkins-Valls, and B. Prakash. Combining learning from  
257 human feedback and knowledge engineering to solve hierarchical tasks in minecraft. *arXiv*  
258 *preprint arXiv:2112.03482*, 2021.
- 259 [11] D. Hafner. Benchmarking the spectrum of agent capabilities. *arXiv preprint arXiv:2109.06780*,  
260 2021.
- 261 [12] H. Hu, D. Yarats, Q. Gong, Y. Tian, and M. Lewis. Hierarchical decision making by generating  
262 and following natural language instructions. *Advances in neural information processing systems*,  
263 32, 2019.
- 264 [13] W. Huang, P. Abbeel, D. Pathak, and I. Mordatch. Language models as zero-shot planners:  
265 Extracting actionable knowledge for embodied agents. In *International Conference on Machine*  
266 *Learning*, pages 9118–9147. PMLR, 2022.

- 267 [14] Y. Jiang, S. S. Gu, K. P. Murphy, and C. Finn. Language as an abstraction for hierarchical deep  
268 reinforcement learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- 269 [15] H. Küttler, N. Nardelli, A. Miller, R. Raileanu, M. Selvatici, E. Grefenstette, and T. Rocktäschel.  
270 The nethack learning environment. *Advances in Neural Information Processing Systems*, 33:  
271 7671–7684, 2020.
- 272 [16] H. Le, N. Jiang, A. Agarwal, M. Dudik, Y. Yue, and H. Daumé III. Hierarchical imitation and  
273 reinforcement learning. In *International conference on machine learning*, pages 2917–2926.  
274 PMLR, 2018.
- 275 [17] M. Matthews, M. Samvelyan, J. Parker-Holder, E. Grefenstette, and T. Rocktäschel. Hierarchical  
276 kickstarting for skill transfer in reinforcement learning, 2022. URL [https://arxiv.org/  
277 abs/2207.11584](https://arxiv.org/abs/2207.11584).
- 278 [18] B. Prakash, N. Waytowich, T. Oates, and T. Mohsenin. Interactive hierarchical guidance using  
279 language. *arXiv preprint arXiv:2110.04649*, 2021.
- 280 [19] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,  
281 P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision.  
282 In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- 283 [20] M. Samvelyan, R. Kirk, V. Kurin, J. Parker-Holder, M. Jiang, E. Hambro, F. Petroni, H. Küttler,  
284 E. Grefenstette, and T. Rocktäschel. Minihack the planet: A sandbox for open-ended reinforcement  
285 learning research. In *Thirty-fifth Conference on Neural Information Processing Systems  
286 Datasets and Benchmarks Track (Round 1)*, 2021. URL [https://openreview.net/forum?  
287 id=skFwlyefkWJ](https://openreview.net/forum?id=skFwlyefkWJ).
- 288 [21] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization  
289 algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 290 [22] R. S. Sutton, D. Precup, and S. Singh. Between mdps and semi-mdps: A framework for temporal  
291 abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):181–211, 1999.
- 292 [23] uArm Developer. uArm-Python-SDK, 2018. URL [https://github.com/uArm-Developer/  
293 uArm-Python-SDK](https://github.com/uArm-Developer/uArm-Python-SDK).
- 294 [24] G. Warnell, N. Waytowich, V. Lawhern, and P. Stone. Deep tamer: Interactive agent shaping in  
295 high-dimensional state spaces. *AAAI Conference on Artificial Intelligence*, pages 1545–1553,  
296 2018. URL <https://aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/16200>.
- 297 [25] N. Waytowich, S. L. Barton, V. Lawhern, and G. Warnell. A narration-based reward shaping  
298 approach using grounded natural language commands, 2019.



## 299 A Appendix

### 300 A.1 Environment Details

#### 301 A.1.1 MiniGrid

302 Minigrid is an open source gridworld environment [5]. We use three tasks *UnlockReach*, *KeyCorridor v0* and *KeyCorridor v1*. Figure 5 shows the grid layouts for the three tasks. The observation is an  
303 encoded version of the grid which capture the each cell type, color and an optional door/box state.  
304 We consider the fully observable version of these tasks, which means the observations consists of the  
305 full grid - 13x13 in our case.  
306

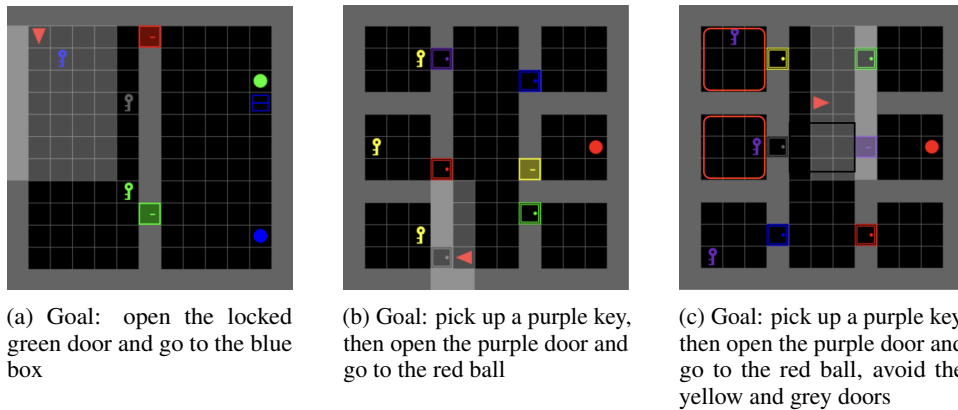


Figure 5: The agent is represented using the red triangle. Left: The *UnlockReach* task where the agent needs to get the right key and open a door and then go to the object in the right room. Middle: The *KeyCorridor-v0* task where the agent needs to reach the red ball in one of the locked rooms on the right. It first needs to get a the key to from one of the rooms on the left. Right: Similar to *v0* but the some of the rooms have defective keys shown in red. The agent does not see this, it only recieves this information in the text goal

#### 307 A.1.2 SkillHack

308 *SkillHack* [17] is an extensions on top of [15] which where you can design custom levels and get  
309 visual/spacial states along with text descriptions. Figure 6 shows the tasks we test where each of them  
310 require solving multiple sub-tasks. The state consists of a 2D map along with inventory information  
311 and text describing the effect of each action. This is very convenient for our method as we need to  
312 interact with the LLM using language.

#### 313 A.1.3 Crafter

314 *Crafter* [11] is a 2D version Minecraft, Figure 6. it has a very simple state representation encoding  
315 the items on the map, an inventory and the health of the agent. It is fairly easy to translate this into  
316 text using a hang-coded function. The high-level skills  $l_{skill_i}$  we consider, such as *chop tree*, *create*  
317 *table*, *make wood pickaxe* etc can also be naturally described using language phrases.

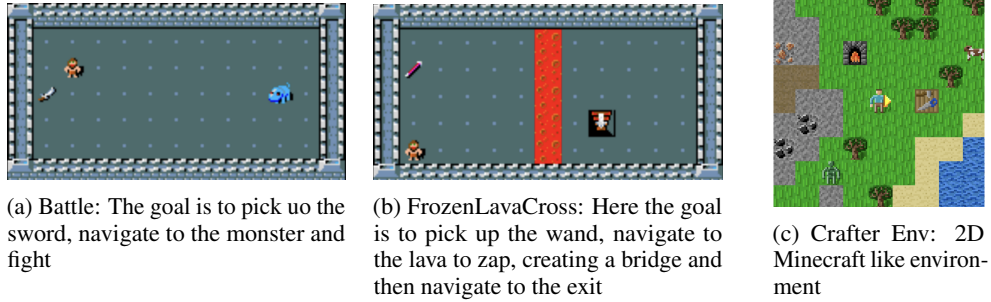


Figure 6: Left and Middle: The SkillHack environment suite based on the text based NetHack game. Right: Crafter: Here, the agent explores the open world environment to collect resources like wood, stone and create new objects and tools

#### 318 A.1.4 uArm Robot Arm

319 uArm Swift Pro is an open source research robotic arm. We use this to test the tabular Q-learning  
 320 version of our method using two tasks [23]. The environment is setup where the table has 3 zones as  
 321 shown in Figure 7 (b). The robot is able to pick and place objects from any of these zones. They are  
 322 considered the low level skills and are hard-coded. However with more resources, these can be trained  
 323 using RL making them more robust. In the *DeskCleanup* task, there is a tray in one of the zones as seen  
 324 in Figure 7 (a) and (b). The goal is to pick up the blocks and place them in the tray. In the *BlockSwap*  
 325 task the two blocks are placed in random zones Figure 7 (c). The goal is to swap the positions. This  
 326 can only be done by utilizing the empty zone. We use an Intel Realsense camera to convert the visual  
 327 information to a simplified discrete state for the tabular Q-learning. The state consists of 4 discrete  
 328 values,  $[arm\_location, holding, red\_location, green\_location]$ . We get the *arm\_location* from  
 329 the robot apis, which is one of the 3 zones. The *red\_location* and *green\_location* denote the  
 330 location of the red and green blocks respectively (one of 3 zones). *holding* denotes the block the  
 331 arm is currently holding if any. We extract *holding*, *red\_location* and *green\_location* from the  
 332 camera image. Instead of training or finetuning a separate object detection model we use CLIP and  
 333 text phrases describing the objects on image patches. We then get the co-ordinates and map them to  
 334 the right values in the state space.

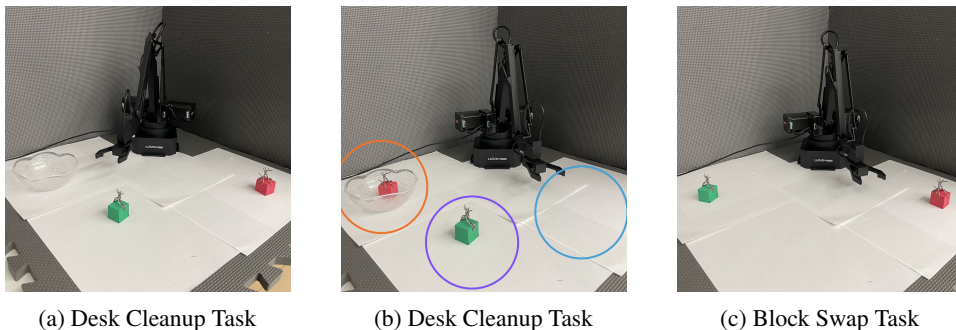


Figure 7: uArm: Robot arm environment

#### 335 A.2 Prompt Design

336 We use the following structure for all our prompts.

337

338 Goal:  $l_{goal\_inst}$

339 So far:  $l_{traj}$

340 Should I  $l_{skill_i}$ ?

341

342  $l_{traj}$  represents the trajectory history which captures what the agent has done so far. We  
 343 find that the 2 actions are sufficient to capture this as that the tasks we test on not extremely

344 complicated. As mentioned in previous sections,  $l_{goal\_inst}$  describes the goal along with details,  
345 warnings and constraints. Each skill is associated with a language description  $l_{skill_i}$ .

### 346 A.2.1 MiniGrid Prompts

```
347 You are a 2D maze-solving agent with access to a variety of low-level
348 skills such as "pick:red:ball", "pick:red:key", "pick:green:ball", "
349 pick:green:key", "pick:blue:ball"...
350
351 Goal: open the locked green door and go to the blue box
352 So far:
353 Question: Should I pick:red:ball?
354 Answer: No
355
356 Goal: open the locked green door and go to the blue box
357 So far:
358 Question: Should I goto:blue:box?
359 Answer: No
360
361 Goal: open the locked green door and go to the blue box
362 So far: pick:green:key
363 Question: Should I unlock:green:door?
364 Answer: Yes
365
366 Goal: open the locked green door and go to the blue box
367 So far: pick:green:key, unlock:green:door
368 Question: Should I goto:red:box?
369 Answer: No
370
371 [..few more examples..]
372
373 Goal: open the locked green door and go to the blue box
374 So far: pick:green:key, unlock:green:door
375 Question: Should I goto:blue:box?
376 Answer:
377
378
```

### 379 A.2.2 SkillHack Prompts

380 The NetHack environment is originally a text based game., So luckily we get the language description  
381 of our action and observations from the game engine.

```
382 You are a NetHack Agent equipped with the following skills:
383
384 ApplyFrostHorn: Use a frost horn to freeze some lava.
385 Eat: Eat an apple.
386 Fight: Hit a monster.
387 NavigateLava: Reach the staircase past random lava patches.
388 PickUp: Pick up a random item.
389 PutOn: Put on an amulet or towel.
390 TakeOff: Take off clothes.
391 Unlock: Use a key to unlock a locked door.
392 Wear: Wear a robe.
393 Wield: Wield a sword.
394 ZapWandOfCold: Use a wand of cold to freeze lava.
395
396 <list common item names>
397
398 Battle: PickUp a randomly placed Sword, Wield the Sword and finally
399 Fight and kill a Monster.
400
401 Goal: Battle the Monster
402 So far: I see a silver saber
403 Question: Should I TakeOff?
404
```

```
405 Answer: No
406
407 Goal: Battle the Monster
408 So far: picked up a silver saber
409 Question: Should I Wield?
410 Answer: Yes
411
412 Goal: Battle the Monster
413 So far: picked up a silver saber, weild a silver saber
414 Question: Should I ZapWandOfCold?
415 Answer: No
416
417 [..few more examples..]
418
419 Goal: Battle the Monster
420 So far: picked up sword, weild sword
421 Question: Should I Fight?
422 Answer:
423
```

```
424
425 You are a NetHack Agent equipped with the following skills:
426
427 [..same as above..]
428
429 Frozen Lava Cross: Either a Wand of Cold or a Frost Horn will spawn on
430 the near side of a river of lava. PickUp the item and then create a
431 bridge across the lava with either ZapWandOfCold or by ApplyFrostHorn
432 . Finally, NavigateLava across your newly made bridge to reach the
433 staircase on the other side.
434
435 <list common item names>
436
437 Goal: Reach the staircase past the lava
438 So far: I see a wand
439 Question: Should I PickUp?
440 Answer: Yes
441
442 [..few more examples..]
443
444 Goal: Reach the staircase past the lava
445 So far: picked up a wand, apply zap, the lava cools and solidifies
446 Question: Should I NavigateLava?
447 Answer:
448
```

### 449 A.2.3 Crafter Prompts

450 For this environment, the prompts are designed similar to [7].