

# Supplementary Materials of ‘Exploring Beyond Curiosity Rewards: Language-Driven Exploration in RL’

**Editors:** Vu Nguyen and Hsuan-Tien Lin

## 1 Appendix A. Limitations and Discussions

2 A limitation of the present study that is worth acknowledging lies in the computational cost  
3 associated with LLMs. While our approach significantly enhances sample efficiency, it is  
4 important to recognize that querying massive LLMs may introduce an additional overhead  
5 in terms of time and computation. To mitigate this, a strategy could involve distilling the  
6 LLM into more compact actor networks, similar to [Parisotto and Salakhutdinov \(2021\)](#). In  
7 addition, strategies such as quantization, limiting the number of output tokens, or staged  
8 speculative decoding [Spector and Re \(2023\)](#) can be employed to accelerate their inference.  
9 It is important to note, however, that these computational constraints do not significantly  
10 impact the action generation task due to the minimal number of tokens it generates.

11 As pointed out in some recent papers [Nair et al. \(2018\)](#), exploring the task following  
12 sub-optimal demonstrations often results in poor performance or a suboptimal policy [Du  
13 et al. \(2023\)](#). Conversely, distilling a generative model through an exploration bonus has  
14 been shown to enable learning better-than-expert performance, even from a noisy expert [Yu  
15 et al. \(2020\)](#). Thus, iLLM relies on curiosity-driven exploration to enable better-than-llm  
16 performance. Moreover, this design also offers modularity, allowing the integration of any  
17 language tasks such as summarization or planning to guide the agent’s exploration.

18 In the action generation task, we select actions based on the highest log probability  
19 rather than generating actions directly. A significant challenge with direct action generation  
20 is the variability in the format of the generated actions, necessitating an additional step  
21 to map these actions to the environment’s action space. By selecting the action with the  
22 highest log probability, we bypass this step. Besides, the selected environments did not  
23 require complex action descriptions, alleviating the need to generate actions represented by  
24 sentences. We leave to future work to explore this direction further.

25 Another avenue for improvement is to leverage self-reflection on text observations and  
26 *translated* state-action pairs. Self-reflection mechanisms [Park et al. \(2023\)](#) that enable the  
27 agent to assess the quality and informativeness of textual observations could refine the  
28 summarization process. By employing self-assessment, the agent may learn to prioritize  
29 and extract key information, ultimately leading to more concise and informative textual  
30 representations.

31 We also aim to explore alternative formulations of intrinsic rewards. In our current  
32 implementation, these rewards are based on the similarity between the LLM’s predictions  
33 and the outputs from the action/history compression heads. In future work with richer  
34 environments, we anticipate using more sophisticated methods for estimating the similarity

---

**Algorithm 1** iLLM(hop) Algorithm

---

**Require:** Interaction budget  $N$ , Horizon  $T$ , Policy  $\pi$ , Large language model  $LLM$ , Action generation head  $f^a$ , History compression head  $f^{hc}$ , Environment  $env$

```
1: for  $i \leftarrow 1$  to  $N$  do
2:    $o_1 \leftarrow env.reset()$ 
3:   for  $t \leftarrow 1$  to  $T$  do
4:      $x_t \leftarrow \mathbf{E}^\top \text{softmax}(\beta \mathbf{E} \mathbf{P}(o_t \cdot a_{t-1}))$  {Translated representation}
5:      $Z_h \leftarrow LLM(c_{t-1}, x_t)$  {History representation}
6:      $Z_p \leftarrow \text{get\_action\_prompt}()$  {Action generation prompt}
7:      $Z_p'' \leftarrow \text{get\_history\_prompt}()$  {History compression prompt}
8:      $\bar{a}_t \leftarrow \max_{a^i \in \{1, \dots, K\}} LLM(a^i | Z_h, Z_p')$  {Get next action from the LLM}
9:      $\tilde{s}_t \leftarrow \sigma(LLM(Z_h, Z_p''))$  {Get mean-pooled representation of generated summary}
10:    Call policy  $\pi(CNN(o_t) \cdot Z_h)$  {Sample policy and get outputs of the two heads  $f^a$ ,  $f^{hc}$ }
11:    Compute intrinsic rewards  $r_t^a$  and  $r_t^{hc}$  based on head outputs, and  $\bar{a}_t$  and  $\tilde{s}_t$ , respectively {Intrinsic Rewards}
12:    Store  $r_t^* = r_t + b_t = r_t + r_t^a + r_t^{hc}$  along with standard experience in optimization batch  $B_i$  {Store agent's experience}
13:  end for
14:  Update( $\pi$ ) using  $B_i$  {Update the policy}
15: end for=0
```

---

35 between the agent's trajectory and those predictions, such as using a CLIP-based objective  
36 Yuan et al. (2023).

37 A promising research direction that emerges from our work is the refinement of the  
38 mapping process between state-action pairs and the LLM's internal representation. While  
39 iLLM successfully harnesses Hopfield networks to align them, there is room for further  
40 optimization. Advancements in multimodal techniques for semantically rich text encoding  
41 and decoding may play an important role in enhancing the agent's capacity to interpret and  
42 respond to *translated* token inputs Wang et al. (2023). We also plan to incorporate modern  
43 Hopfield networks to further mitigate the alignment gap Furst et al. (2022).

44 Finally, *action generation* and *history compression* tasks are less helpful in domains  
45 where human common sense is irrelevant or cannot be expressed in language (e.g., fine-  
46 grained manipulation), or action information is not naturally encoded as a language string.  
47 An avenue for future work is prompt engineering in order to craft more suitable language  
48 tasks.

## 49 Appendix B. Implementation Details

50 As stated previously, PPO was used as our policy learning method. We choose AdamW  
51 Loshchilov and Hutter (2019) as optimizer for all environments with default values for  
52 weight decay, and clip the norm of the gradients to 0.5. In all experiments, we utilized the  
53 4-bit quantized version of Transfo-XL 280M Dai et al. (2019). Default values were kept for

54 all parameters, except temperature, which was set to 0 to ensure deterministic responses.  
 55 This language model was employed to compute the embeddings of state-action descriptions  
 56  $Z_h$ . Similarly, we embedded the action generation and summarization prompts  $Z_p$  with the  
 57 same embedding model. To accelerate the LLM inference, we utilized Lamorel [lam \(2023\)](#)  
 58 and caching for the action generation and history compression tasks. Furthermore, querying  
 59 of the LLM for the two language tasks took place in-batch, following the completion of each  
 60 rollout. Note that due to image-based tasks, we could not report the results of Pangu and  
 61 PAE in some experiments.

62 In iLLM(obs), the policy takes as input the current observation encoded through a CNN.  
 63 In iLLM(hop), the current observation is encoded with a CNN and concatenated with  $Z_h$   
 64 to form the input for the policy. As mentioned above, the Hopfield module is specifically  
 65 designed to handle multimodal inputs, such as visual observations, textual descriptions, and  
 66 numerical data. By translating these diverse inputs into a unified token embedding space,  
 67 the Hopfield module allows the LLM to process and integrate information from multiple  
 68 modalities effectively.

69  $f^a$  and  $f^{hc}$  are parametrized by two fully connected layers with hidden dimensions  
 70 1024. To maintain a consistent scale of the intrinsic rewards, it is useful to normalize  
 71 them. This can be achieved by dividing the intrinsic rewards by a running estimate of the  
 72 standard deviations of the sum of discounted intrinsic rewards. We set  $\lambda = 0.8$  and  $\beta = 0.5$   
 73 to weight the intrinsic rewards. The simplified pseudo-code demonstrating the training  
 74 procedure of iLLM(hop) is depicted in Alg 1, where, for the sake of clarity, we purposely  
 75 show iLLM querying the LLM after each interaction with the environment. Note that action  
 76 generation and history compression queries (lines 8-9) can be performed in-batch following  
 77 the completion of each rollout to reduce the computational cost.

78 A graphical illustration of iLLM is provided in Figure 1 of the main manuscript. First,  
 79 the LLM first retrieves a representation of the recent (*action, observation*) pairs, which is  
 80 aligned with its internal representation (*input:  $(o_t \cdot a_{t-1})$ , output:  $Z_h$* ). If using text-based  
 81 observations (iLLM(obs)), only the second stage is needed. In this stage, a question prompt  
 82  $Z_p$  along with the aligned representation  $Z_h$  are passed through the LLM in order to obtain  
 83 either the next action or a summary of the observation (*input:  $[Z_p, Z_h]$ , output: an action*  
 84  $\bar{a}_t$  or a summary of action-observation pairs). Finally, an intrinsic reward is derived from  
 85 these next action and summary predictions.

## 86 B.1. State-Action History Representation

87 We now describe the representations of state-action history  $Z_h$  that were used in iLLM(obs)  
 88 and iLLM(hop).

89 In **iLLM(obs)**, the agent’s recent history consists of the last three state-action pairs. The  
 90 description of the agent’s history,  $Z_h$ , inputted into the LLM was formatted as follows:

```
91
92 Observation 0: {obs_0}
93 Action 0: {act_0}
94 Observation 1: {obs_1}
95 Action 1: {act_1}
96 Observation 2: {obs_2}
97 Action 2: {act_2}
98
```

99 where **obs** and **act** are the (text-based) descriptions of the observations at time  $t$  and ac-  
100 tions performed by the agent.

101

102 If using the Hopfield module (**iLLM(hop)**),  $Z_h$  represents the hidden states from the  
103 last hidden layer of the LLM, corresponding to the state-action token. In detail, the input  
104 to the Hopfield module consists of flattened grayscale observations concatenated with the  
105 previous action taken (one-hot encoded). The output of the Hopfield module  $x_t$  is then  
106 passed through the LLM to obtain the state-action token  $Z_h$ .

## 107 B.2. Action Generation and History Compression Tasks

108 This section provides details about the action generation and history compression prompts.  
109 As described above,  $Z_p$  refers to the text embeddings obtained from tokenized text prompts.  
110 That is, we employed the base Transfo-XL 280M tokenizer to tokenize the prompts and  
111 leveraged the word embeddings of the Transfo-XL model to extract the corresponding em-  
112 beddings.

113

114 The action generation prompt format was set as:

115

```
116 You are an expert player playing {task}  
117 Valid actions: {possible actions separated by commas}  
118 You see: {agent history}  
119 Suggest the best action the player can take. Do not recommend  
120 actions that are not possible or not desirable, such as ‘‘Eat door’’  
121 . Prioritize actions which involve the object you are facing or  
122 which the agent has not achieved before. What do you do?  
123
```

124 where **{agent history}** was replaced by  $Z_h$ .

125

126 The history compression (i.e., summarization) prompt format had the following structure:

127

```
128 You are an expert player playing {task}  
129 Recent player’s history:: {agent history}  
130 Summarize the main points of the player’s history into a short  
131 text:  
132
```

133 where **{agent history}** was replaced by  $Z_h$ . We restricted the number of tokens produced  
134 to  $L = 64$  in the history compression task.

## 135 Appendix C. Environments

### 136 C.1. BabyAI-Text

137 BabyAI-Text is a text-based environment that encapsulates BabyAI, providing a textual  
138 description of each observation [Chevalier-Boisvert et al. \(2018a\)](#). We evaluate iLLM on a  
139 set of nine tasks in the BabyAI-Text environment [Chevalier-Boisvert et al. \(2018a\)](#). The  
140 agent must navigate in procedurally generated rooms that include distractors — useless  
141 objects for completing the task.

142 We selected the following tasks:

- 143 • **Key corridor**, a task that requires to pick up an object which is behind a locked door.  
144 The key is hidden in another room, and the agent has to explore the environment to  
145 find it.
- 146 • **Obstructed maze**, a navigation task where a blue ball is hidden in one of the 4  
147 corners of a maze. Doors are locked, doors are obstructed by a ball and keys are  
148 hidden in boxes.
- 149 • **Go to**, a navigation task that requires reasoning abilities in order to reach the goal  
150 object.
- 151 • **Pick up**, a navigation task that combines navigation tasks and picking up the object.
- 152 • **Put object A next to object B (put next)**, a sequence of 3 tasks, including  
153 reaching object A then reaching object B and finally dropping object A next to object  
154 B.
- 155 • **Open door**, a task that requires inferring that a key is useful for unlocking a door,  
156 finding another key, and finally using the toggle action with the key in the door.

157 A textual description consists of a list of template descriptions with the following struc-  
158 ture:

```
159 "You see a <object> <location>" if the object is a key, a ball, a
160 box, or a wall.
161 "You see a(n) open/closed door <location>" , if the agent sees a
162 door.
163 "You carry a <object>", if the agent carries an object.
164
165
```

166 where the <object> is composed of an adjective (among six possible colors: blue, red,  
167 green, grey, purple, and yellow) and a noun (among four possible: door, key, ball, box).  
168 The <location> is given as the number of steps right, left, and or forward from the agent  
169 to the object.

## 170 C.2. MiniHack

171 The MiniHack environment based on NetHack (Kuttler et al., 2020) features a larger action  
172 space compared to BabyAI, with up to 75 distinct actions. Observations are composed  
173 of a  $21 \times 79$  matrix containing glyph identifiers and a 21-dimensional vector capturing  
174 agent statistics such as location and health. Additionally, real natural language messages  
175 received during gameplay are included in a 256-dimensional vector termed as a “message”,  
176 representing the on-screen display at the screen’s top. Each glyph corresponds to a unique  
177 entity, denoted by integers ranging from 0 to 5991. Our study focused on five MiniHack  
178 tasks, encompassing navigation challenges like River-Monster and Multiroom-N4-Monster,  
179 along with skill acquisition tasks such as LavaCross-Ring, LavaCross-Potion, and LavaCross-  
180 Full. The navigation tasks in MiniHack test the agent’s ability to navigate diverse obstacles,  
181 from maneuvering boulders to crossing rivers, while skill acquisition tasks exploit NetHack’s  
182 vast array of objects, monsters, and dungeon features, exploring their interactions and  
183 complexities.

Model	N=4	N=8	N=16
RND	0.59± 0.16	0.51± 0.22	0.09± 0.10
NGU	0.51± 0.20	0.42± 0.25	0.11± 0.13
ELLM	0.78± 0.03	0.66± 0.01	0.64± 0.05
APT	0.57± 0.19	0.48± 0.17	0.44± 0.20
ChibiT	0.56± 0.23	0.51± 0.15	0.48± 0.17
iLLM(obs)	<b>0.96± 0.01</b>	<b>0.94± 0.00</b>	0.91 ± 0.02
iLLM(hop)	<b>0.96± 0.00</b>	0.92± 0.01	<b>0.94 ± 0.01</b>

Table 1: Success rate for the agents on the *Go To* task for different number of distractors ( $N \in \{4, 8, 16\}$ ). The success rate is provided over 10 seeds with standard deviation after 100M training steps.

### 184 C.3. Crafter

185 Crafter is an open-ended environment in which exploration is required to discover long-  
186 term survival strategies [Hafner \(2021\)](#). It is a 2D variant inspired by Minecraft, featuring  
187 a procedurally generated and partially observable world world. Crafter enables collecting  
188 and creating a set of artifacts organized along an achievement tree, which lists all possible  
189 achievements and their respective prerequisites. Despite lacking a single main task, tracking  
190 the agent’s advancements along the achievement tree provides insights into its progress  
191 within Crafter.

## 192 Appendix D. Ablation Results

### 193 D.1. Impact of the Number of Distractors

194 This section aims to assess how much distractors impact the proposed method. These  
195 evaluations were performed in an environment with one room on the *Go To* task (BabyAI-  
196 Text). We report in Table 1 the average success rate for 4, 8, and 16 distractors. RND  
197 and NGU significantly degrade as the number of distractors increases, with a success rate  
198 decreasing by  $\approx 65\%$  from 4 to 16 distractors. We also observe a slight performance loss  
199 in iLLM agents when the number of distractors is larger than 8. A potential rationale  
200 behind this phenomenon is that the LLM effectively directs the iLLM agent’s attention  
201 to the relevant aspects of the environment, rapidly discarding distractors and noise in the  
202 observations. On the other hand, model-based curiosity is more impacted by distractors as  
203 it favors a full exploration of the state space, resulting in the exploration of a larger number  
204 of irrelevant behaviors. Overall, the present algorithm appears to be reasonably robust to  
205 distractors.

### 206 D.2. Language Model Tasks

207 As described in section 3, iLLM relies on action generation and history compression tasks.  
208 In this experiment, we evaluate the performance of the proposed method with two other

Table 2: Final mean performance ( $\pm$  std) of iLLM(hop) trained with different language tasks, including action generation ( $\clubsuit$ ), history compression ( $\heartsuit$ ), goal generation ( $\spadesuit$ ), and plan generation ( $\diamond$ ). For the sake of generality, we report results where the recent history  $Z_h$  consists of translated state-action pairs, iLLM(hop). Averages over 10 runs.

Method	BabyAI		Atari		MiniHack	
	GoToObj	PutNextLocal	MR	PrivateEye	LavaCross-Full	River-Monster
iLLM + $\clubsuit$	0.85 $\pm$ 0.01	0.43 $\pm$ 0.09	2,118 $\pm$ 329	3,981 $\pm$ 420	0.91 $\pm$ 0.01	0.35 $\pm$ 0.09
iLLM + $\heartsuit$	0.90 $\pm$ 0.01	0.41 $\pm$ 0.10	2,456 $\pm$ 198	3,565 $\pm$ 454	0.98 $\pm$ 0.03	0.35 $\pm$ 0.07
iLLM + $\spadesuit$	0.85 $\pm$ 0.02	0.42 $\pm$ 0.08	2,024 $\pm$ 176	3,429 $\pm$ 500	0.96 $\pm$ 0.02	0.32 $\pm$ 0.10
iLLM + $\diamond$	0.71 $\pm$ 0.07	0.12 $\pm$ 0.13	1,365 $\pm$ 202	2,560 $\pm$ 321	0.86 $\pm$ 0.09	0.26 $\pm$ 0.06
iLLM + $\clubsuit$ + $\heartsuit$	0.92 $\pm$ 0.01	0.49 $\pm$ 0.12	2,632 $\pm$ 277	4,422 $\pm$ 376	0.98 $\pm$ 0.03	0.38 $\pm$ 0.12
iLLM + $\spadesuit$ + $\diamond$	0.83 $\pm$ 0.05	0.42 $\pm$ 0.06	2,139 $\pm$ 251	2,945 $\pm$ 400	0.89 $\pm$ 0.07	0.30 $\pm$ 0.08
iLLM + $\clubsuit$ + $\diamond$	0.84 $\pm$ 0.03	0.45 $\pm$ 0.11	2,299 $\pm$ 244	3,647 $\pm$ 178	0.91 $\pm$ 0.01	0.31 $\pm$ 0.05
iLLM + $\clubsuit$ + $\spadesuit$	0.93 $\pm$ 0.02	0.46 $\pm$ 0.05	2,432 $\pm$ 199	4,295 $\pm$ 420	0.95 $\pm$ 0.01	0.34 $\pm$ 0.06
iLLM + $\heartsuit$ + $\spadesuit$	0.93 $\pm$ 0.01	0.48 $\pm$ 0.09	2,312 $\pm$ 271	3,999 $\pm$ 392	0.98 $\pm$ 0.02	0.36 $\pm$ 0.08
iLLM + $\heartsuit$ + $\diamond$	0.86 $\pm$ 0.07	0.39 $\pm$ 0.08	2,202 $\pm$ 302	2,876 $\pm$ 287	0.96 $\pm$ 0.01	0.31 $\pm$ 0.06
iLLM + $\clubsuit$ + $\heartsuit$ + $\spadesuit$	0.93 $\pm$ 0.03	0.51 $\pm$ 0.08	2,553 $\pm$ 298	4,500 $\pm$ 356	0.97 $\pm$ 0.02	0.39 $\pm$ 0.10
iLLM + $\clubsuit$ + $\heartsuit$ + $\diamond$	0.89 $\pm$ 0.02	0.47 $\pm$ 0.08	2,421 $\pm$ 230	4,053 $\pm$ 312	0.98 $\pm$ 0.03	0.36 $\pm$ 0.09
iLLM + $\spadesuit$ + $\heartsuit$ + $\diamond$	0.88 $\pm$ 0.04	0.45 $\pm$ 0.12	2,376 $\pm$ 255	3,971 $\pm$ 253	0.96 $\pm$ 0.01	0.34 $\pm$ 0.09
iLLM + $\clubsuit$ + $\spadesuit$ + $\diamond$	0.90 $\pm$ 0.03	0.50 $\pm$ 0.06	2,112 $\pm$ 303	4,421 $\pm$ 443	0.94 $\pm$ 0.02	0.34 $\pm$ 0.07
iLLM + $\clubsuit$ + $\spadesuit$ + $\heartsuit$ + $\diamond$	0.93 $\pm$ 0.02	0.50 $\pm$ 0.05	2,510 $\pm$ 300	4,499 $\pm$ 398	0.96 $\pm$ 0.02	0.37 $\pm$ 0.05

language tasks, including goal generation and plan generation. For **goal generation**, we queried the LLM with the following instruction:

```

209 You are an expert player playing {task}
210 Recent player's history: {agent history}
211 Suggest the next goal to reach based on the things you see and
212 previous actions. A goal should either be a single valid word or
213 a phrase. Only make suggestions that are reasonable given the
214 current scene (e.g. only 'Open door' if a door is visible).
215 Prioritize goals that involve the object you are facing or that
216 the agent has not achieved before.
217
218
219
220

```

and for **plan generation**:

```

221 You are an expert player playing {task}
222 Recent player's history: {agent history}
223 Suggest the best sequence of actions the player can take. An
224 action should either be a single valid word or a phrase. Only
225 make suggestions that are reasonable given the current scene
226 (e.g., only 'Open door' if a door is visible). Prioritize
227 actions which involve the object you are facing or which the
228 agent has not achieved before. What do you do (include 2-7
229 actions)?
230
231
232

```

The associated heads were trained in the same fashion as the history compression head. As shown in Table 2, agents leveraging action generation and history compression tasks demonstrate reasonable scores on all tasks. Notably, *goal generation* achieves high performance on BabyAI-text, but the average return decreases on more complex games such as

Method	MR	PrivateEye	Gravitar	Pitfall	Seaquest
PPO	11±4	0.0±0.0	120±14	-7±2	1,245±199
iLLM(hop)	<b>2,632±277</b>	<b>4,422±376</b>	<b>4,044±559</b>	<b>125±24</b>	<b>18,851±2,930</b>
iLLM(hop)(no reward)	1,452±201	2,871±265	2,450±287	3±5	15,888±3,012

Table 3: Performance of iLLM that solely distills an LLM via the action generation and history compression heads on Atari tasks. Under this setting, iLLM(hop)(no reward) does not receive any intrinsic rewards. All methods are tested with 10 random seeds. Averages over 10 runs for 100 million training steps.

237 Atari. Besides, we find that the model performance saturates when the number of language  
238 tasks goes larger. Namely, leveraging more language tasks cannot monotonically promote  
239 performance when the number is larger than 2. Therefore, we empirically set the number  
240 of tasks to 2 by default. In addition, it is evident that *history compression* is the reward  
241 that accelerates the most exploration. In contrast, both *action generation* and *goal* gener-  
242 ation result in slightly lower improvements. Finally, *plan generation* proved to be the least  
243 effective reward, which may be attributed to the difficulty of generating meaningful plans  
244 with a small language model like Transfo-XL.

### 245 D.3. State Visitation

246 To test the good exploration coverage of our approach, we trained iLLM(hop) on a procedurally-  
247 generated environment: MiniGrid-MultiRoom-N6-v0 [Chevalier-Boisvert et al. \(2018b\)](#). Uti-  
248 lizing MiniGrid-MultiRoom-N6-v0, which involves sequential visits to multiple rooms, fa-  
249 cilitates a clearer measurement of the exploration progress. Fig. 1 demonstrates that in  
250 order to remain curious, the agent is pushed to explore distant regions of the state space,  
251 which entails that its coverage increases over time. Experimental results highlight that  
252 in procedurally-generated tasks, iLLM provides enough exploration incentive for learning  
253 useful behaviors. That is, in order to remain curious, agents are pushed to explore diverse  
254 behaviors, enabling the discovery of new rooms and interactions with the objects. Unlike  
255 RND and NGU, which spend time to exploring local behaviors that have a low interest,  
256 iLLM(hop) rapidly drives the agent towards the goal being pursued (red cell).

### 257 D.4. Diverse Exploration

258 In this set of experiments, we aim to evaluate the contribution of distilling an LLM to  
259 the efficacy of our proposed method. One of our assumptions is that appending  $f^a$  and  
260  $f^{hc}$  to the policy and training them to match the pretrained LLM’s outputs enables iLLM  
261 to leverage world knowledge through their respective gradients. In order to validate this  
262 hypothesis, we report the results of PPO, iLLM(hop), and iLLM(hop)(no reward) in the five  
263 Atari games. The latter method was trained without access to the intrinsic rewards. Thus,  
264 LLM’s prior knowledge is only distilled through gradients of  $f^a$  and  $f^{hc}$  — the agent does  
265 not receive an explicit incentive to explore. As indicated in Table 3, in the five tasks, we  
266 find that distilling an LLM, as done by *iLLM(hop)(no reward)*, leads to higher performance  
267 compared to plain PPO. This highlights that even in the absence of any intrinsic rewards, the

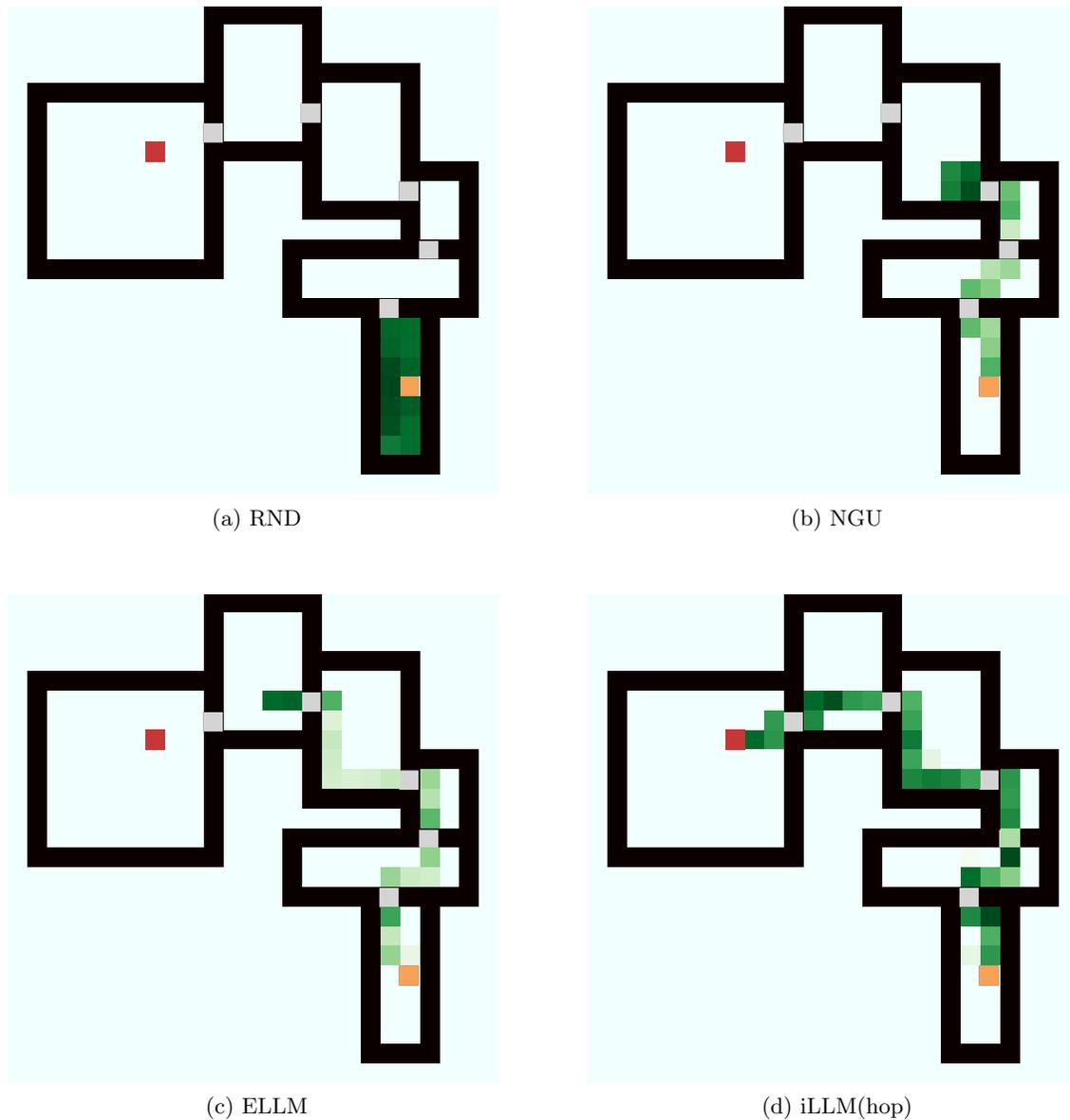


Figure 1: State visitation heatmaps over 10 runs for RND, NGU, ELLM, and iLLM(hop) on a random environment from the MiniGrid-MultiRoom-N6-v0 task. The initial agent's location is denoted with an orange cell, while the target is denoted with a red cell. We trained the models for 40m frames.

Method	MR	PrivateEye	Gravitar	Pitfall	Seaquest
iLLM(hop)	<b>2,632±277</b>	<b>4,422±376</b>	<b>4,044±559</b>	<b>125±24</b>	<b>18,851±2,930</b>
iLLM(hop) ( $\varrho = 0.50$ )	1,989±134	3,522±398	3,232±356	65±29	15,437±2,273
iLLM(hop) ( $\varrho = 1.0$ )	2,012±209	3,408±401	3,018±410	59±31	13,419±2,414

Table 4: LLM stability study. We report the average reward in randomized-versions of Atari games. All methods are tested with 10 random seeds. Averages over 10 runs for 100 million training steps.

two prediction heads attached to the policy provide sufficient domain-specific priors to speed up the training process. We also see here that iLLM(hop) is more competitive, suggesting that including intrinsic rewards is better suited for exploration than agents without access to intrinsic rewards.

### D.5. LLM Stability

One potential drawback of leveraging LLMs is their stability when facing various input conditions. As mentioned above, in our experiments, we set the temperature to 0 in order to enhance stability. In addition, we observed that limiting the number of tokens produced to  $L = 64$  in the history compression task makes predictions more stable. It is also important to note that the LLM in our approach is used to encourage exploration through a small reward incentive, rather than serving as an oracle providing plans or next actions. This makes iLLM inherently more stable than approaches that rely on the LLM for direct decision-making.

To further validate the stability of iLLM, we study the effect of adding perturbations to the observations. Namely, with a probability  $\varrho \in \{0.50, 1.0\}$ , a noise pattern ( $32 \times 32$ ) is displayed on the lower right of the observation - TV screen. The noise is sampled from  $[0, 255]$  independently for each pixel. As reported in Table 4, the performance iLLM deteriorates due to the stochasticity. Nevertheless, our approach is reasonably robust to randomized observations. As iLLM does not directly relies on the LLM’s output but is driven by a proxy reward obtained from the LLM, iLLM remains stable under various input conditions. Overall, LLM stability does not appear to be a concern for iLLM under our experimental design choices, but we leave it to future work to explore this direction further.

### D.6. Impact of the Dimension of the Action Space

This experiment aims to assess the sensitivity of the present reward to the size of the action space by implementing three action space settings on the *Go to* task:

- **Original:** the action space consists of the only 3 useful actions: turn left, turn right, go forward.
- **Augmented:** the action space consists of the 6 actions that can be performed in the environment. The agent can select one action among 3 useful and 3 useless actions that are pick up, drop and, toggle.
- **Irrelevant:** the action space consists of 9 actions (3 useful and 6 useless with pick up, drop, toggle, sleep, do nothing and think). The last three actions have been selected

Model	Original	Augmented	Irrelevant
RND	0.51± 0.22	0.42± 0.17	0.27± 0.08
NGU	0.42± 0.25	0.37± 0.21	0.25± 0.09
ELLM	0.66± 0.01	0.60± 0.05	0.55± 0.03
APT	0.48± 0.17	0.46± 0.18	0.42± 0.09
ChibiT	0.56± 0.23	0.51± 0.25	0.45± 0.12
iLLM(obs)	<b>0.94± 0.00</b>	0.88± 0.08	0.81± 0.04
iLLM(hop)	0.92± 0.01	<b>0.89± 0.11</b>	<b>0.83± 0.09</b>

Table 5: Effect of using different action space sizes on iLLM performance in the *Go to* task. We report results for three settings: original (3 useful actions), augmented (6 actions), and irrelevant (9 actions). Results are averaged over 10 trials.

299 such that they are irrelevant for solving the *Go To* task and therefore should not  
300 impact an agent that has knowledge about the world.

301 In Table 5, we present the evaluation conducted in an environment comprising 1 room and  
302 8 distractors. The result of this experiment is that using *augmented* or *irrelevant* actions  
303 does not significantly degrade the performance of the proposed approach, while the average  
304 reward of other agents decreases with the exception of APT. Upon looking at the videos,  
305 we observed that other agents tend to frequently select useless and irrelevant actions. On  
306 the other hand, from the onset of the training, our agents are rewarded for trying useful  
307 actions first, which allows us to achieve higher sample efficiency.

### 308 D.7. Use of LLM

309 At each time step the LLM is first employed to retrieve a representation of the (action,  
310 observation) pair, which is aligned with its internal representation (input:  $(o_t \cdot a_{t-1})$ , output:  
311  $Z_h$ ). Then, a prompt  $Z_p$  along with the aligned representation  $Z_h$  are passed through the  
312 LLM in order to obtain either the next action or a summary of the observation (input:  
313  $Z_p, Z_h$ , output: an action  $\bar{a}_t$  or a summary of the action-observation pair).

314 To illustrate the use of the LLM in our framework, Figure 2 provides an example specific  
315 to the action generation task. Although we focus on action generation here for brevity, it  
316 is important to note that task 1 (alignment) is performed only once per time step. The  
317 same aligned representation,  $Z_h$ , is then utilized for both the action generation and history  
318 compression tasks. For simplicity, we assume that  $Z_h$  is retrieved solely from the most  
319 recent state-action pair.

### 320 D.8. Choice of Foundation Model

321 We now seek to evaluate the performance of our methodology using various foundation  
322 models on Atari games. Specifically, we compare the results obtained by employing Transfo-  
323 XL 280M, GPT2 137M, GPT2 870M, Llama2-7b, and Llama-3 8b. The results, presented  
324 in Table 6, demonstrate that the performance of iLLM is generally robust across different  
325 foundation models. While LLama3 exhibits a significantly higher average return (t-test

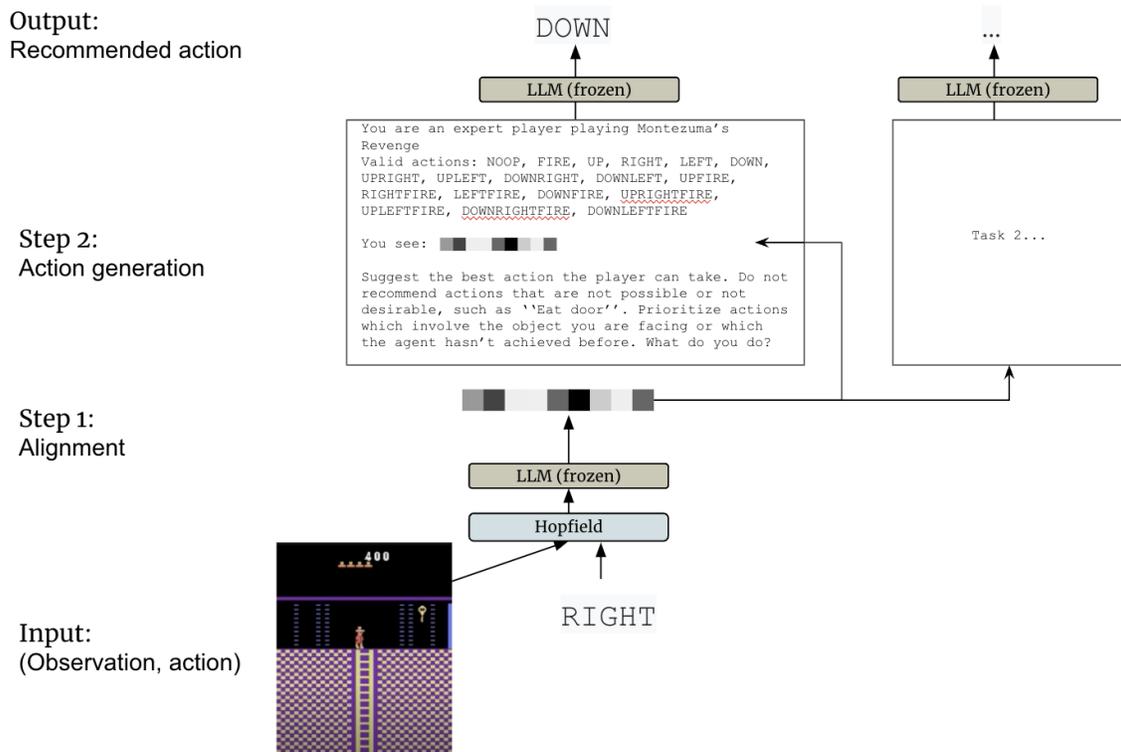


Figure 2: Example of LLM usage in iLLM for the action generation task. iLLM first aligns a (action, observation) pair with the internal representation of the LLM. Then, a prompt along with the aligned representation are passed through the LLM in order to obtain the next action.

Method	MR	PrivateEye	Gravitar	Pitfall	Seaquest
iLLM (hop · Transfo-XL)	2,632±277	4,422±376	4,044±559	125±24	18,851±2,930
iLLM (hop · GPT2 137)	2,299±256	4,134±350	3,756±498	98±31	16,648±2,615
iLLM (hop · GPT2 870)	2,524±312	4,288±321	3,877±522	110±22	17,461±2,646
iLLM (hop · Llama2)	2,953±269	4,954±341	4,644±488	176±31	20,129±2,424
iLLM (hop · Llama3)	<b>3,220±312</b>	<b>5,098±430</b>	<b>4,831±502</b>	<b>277±55</b>	<b>23,980±2,731</b>

Table 6: Performance of iLLM with different types of foundation LLMs on Atari tasks. All methods are tested with 10 random seeds. Averages over 10 runs for 100 million steps.

326  $p < 0.05$ ), Transfo-XL achieves similar performance but with a much lower inference time.  
327 Namely, Transfo-XL model has only 280M parameters compared to the 8B parameters of  
328 LLama3.

## 329 References

- 330 Lamorel. <https://github.com/flowersteam/lamorel>, 2023.
- 331 Maxime Chevalier-Boisvert, Dzmitry Bahdanau, Salem Lahlou, Lucas Willems, Chitwan  
332 Saharia, Thien Huu Nguyen, and Yoshua Bengio. Babyai: A platform to study the sample  
333 efficiency of grounded language learning. *arXiv preprint arXiv:1810.08272*, 2018a.
- 334 Maxime Chevalier-Boisvert, Lucas Willems, and Suman Pal. Minimalistic gridworld envi-  
335 ronment for openai gym. <https://github.com/maximecb/gym-minigrid>, 2018b.
- 336 Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V Le, and Ruslan Salakhut-  
337 dinov. Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv*  
338 *preprint arXiv:1901.02860*, 2019.
- 339 Yuqing Du, Olivia Watkins, Zihan Wang, Cédric Colas, Trevor Darrell, Pieter Abbeel,  
340 Abhishek Gupta, and Jacob Andreas. Guiding pretraining in reinforcement learning with  
341 large language models. *arXiv preprint arXiv:2302.06692*, 2023.
- 342 Andreas Furst, Elisabeth Rumetshofer, Johannes Lehner, Viet T Tran, Fei Tang, Hubert  
343 Ramsauer, David Kreil, Michael Kopp, Günter Klambauer, Angela Bitto, et al. Cloob:  
344 Modern hopfield networks with infoloob outperform clip. *Advances in neural information*  
345 *processing systems*, 35:20450–20468, 2022.
- 346 Danijar Hafner. Benchmarking the spectrum of agent capabilities. *arXiv preprint*  
347 *arXiv:2109.06780*, 2021.
- 348 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019.
- 349 Ashvin Nair, Bob McGrew, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel.  
350 Overcoming exploration in reinforcement learning with demonstrations. In *Proceedings of*  
351 *the IEEE International Conference on Robotics and Automation*, pages 6292–6299, 2018.
- 352 Emilio Parisotto and Ruslan Salakhutdinov. Efficient transformers in reinforcement learning  
353 using actor-learner distillation. *arXiv preprint arXiv:2104.01655*, 2021.
- 354 Joon Sung Park, Joseph C O’Brien, Carrie J Cai, Meredith Ringel Morris, Percy Liang,  
355 and Michael S Bernstein. Generative agents: Interactive simulacra of human behavior.  
356 *arXiv preprint arXiv:2304.03442*, 2023.
- 357 Benjamin Spector and Chris Re. Accelerating llm inference with staged speculative decod-  
358 ing. *arXiv preprint arXiv:2308.04623*, 2023.
- 359 Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo,  
360 Tong Lu, Jie Zhou, Yu Qiao, et al. Visionllm: Large language model is also an open-ended  
361 decoder for vision-centric tasks. *arXiv preprint arXiv:2305.11175*, 2023.

- 362 Xingrui Yu, Yueming Lyu, and Ivor Tsang. Intrinsic reward driven imitation learning via  
363 generative model. In *International conference on machine learning*, pages 10925–10935.  
364 PMLR, 2020.
- 365 Haoqi Yuan, Chi Zhang, Hongcheng Wang, Feiyang Xie, Penglin Cai, Hao Dong, and  
366 Zongqing Lu. Skill reinforcement learning and planning for open-world long-horizon  
367 tasks. *arXiv preprint arXiv:2303.16563*, 2023.