
LLMs are Greedy Agents: Effects of RL Fine-tuning on Decision-Making Abilities

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The success of LLMs has sparked interest in various agentic applications. A
2 key hypothesis is that LLMs, leveraging common sense and Chain-of-Thought
3 (CoT) reasoning, can effectively explore and efficiently solve complex domains.
4 However, LLM agents have been found to suffer from sub-optimal exploration and
5 the knowing-doing gap, the inability to effectively act on knowledge present in the
6 model. In this work, we systematically study *why* LLMs perform sub-optimally
7 in decision-making scenarios. In particular, we closely examine three prevalent
8 failure modes: greediness, frequency bias, and the knowing-doing gap. We propose
9 mitigation of these shortcomings by fine-tuning via Reinforcement Learning (RL)
10 on self-generated CoT rationales. Our experiments across multi-armed bandits,
11 contextual bandits, and Tic-tac-toe, demonstrate that RL fine-tuning enhances
12 the decision-making abilities of LLMs by increasing exploration and narrowing
13 the knowing-doing gap. Finally, we study both classic exploration mechanisms,
14 such as ϵ -greedy, and LLM-specific approaches, such as self-correction and self-
15 consistency, to enable more effective fine-tuning of LLMs for decision-making.

16 1 Introduction

17 Large Language Models (LLMs) pre-trained on massive internet-scale datasets have demonstrated
18 success across diverse domains, including text generation and language understanding [Radford et al.,
19 2019, Brown et al., 2020b, Team et al., 2023b, 2024a, Dubey et al., 2024]. Their broad pre-training
20 distribution, enables generalization to a wide range of scenarios including coding assistance [Li et al.,
21 2022], education [Team et al., 2024d], and medicine [Saab et al., 2024]. Therefore, their success has
22 sparked interest in using LLMs for decision-making problems [Chen et al., 2023, Krishnamurthy
23 et al., 2024, Nie et al., 2024] at the core of agentic AI systems [Durante et al., 2024].

24 One key hypothesis is that LLMs can generate informed action predictions without extensive en-
25 vironment interaction [Lu et al., 2024] due to “world knowledge” present in the model. Moreover,
26 Chain-of-Thought (CoT) [Wei et al., 2022] equips models with the ability to reason about the observed
27 history and their actions, which facilitates environment interaction. However, these advantages do not
28 seem to materialize into strong performance when LLMs are faced with decision-making scenarios.
29 Notably, Krishnamurthy et al. [2024] and Nie et al. [2024] found that LLMs do not robustly engage
30 in *exploration* resulting in sub-optimal behavior. Similar shortcomings of LLMs have been observed
31 by Paglieri et al. [2024] and Ruoss et al. [2024] on stateful environments commonly used in RL (e.g.,
32 grid-worlds, Atari). Both works attribute the shortcomings to the *knowing-doing gap*, which states
33 that models can possess knowledge about a task or can describe the consequences of their behavior
34 (i.e., they know what to do), but cannot materialize this knowledge when acting (i.e., incapable of
35 doing). Consequently, sub-optimal exploration and the knowing-doing gap are considerable obstacles
36 towards more powerful agentic LLMs.

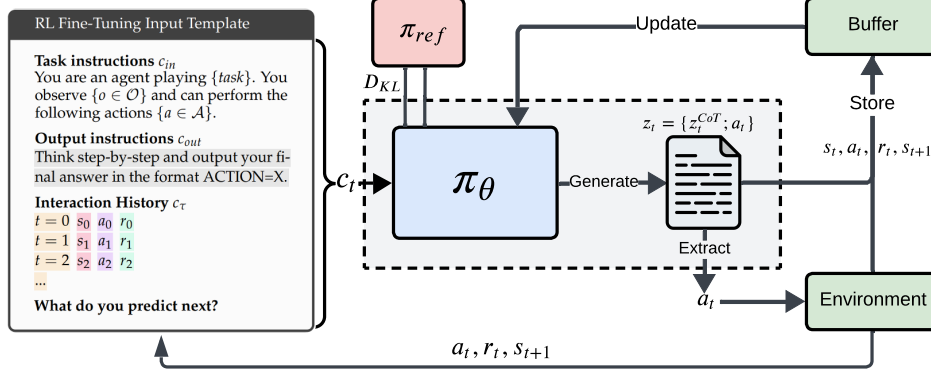


Figure 1: Illustration of our **Reinforcement Learning Fine Tuning (RLFT)** pipeline. We fine-tune a pre-trained LLM π_θ via self-generated Chain-of-Thought (CoT) rationales on environment rewards.

In this work, we aim to understand *why* LLMs often perform sub-optimally in simple decision-making scenarios. In particular, we systematically study three prevalent failure modes in small-to-medium-scale LLMs: greediness, frequency bias, and the knowing-doing gap (see Section 4.2). Our analysis shows that final performance often remains sub-optimal, because LLMs prematurely commit to greedy action selection strategies leading to stagnating action coverage that leave a large part of the action space unexplored (up to 55%). Moreover, we observe that small-scale LLMs (2B) tend to copy the most frequent actions in the context regardless of their respective reward, which we refer to as frequency bias. In contrast, larger LLMs (27B) mostly diminish the frequency bias, yet they remain prone to greedy behavior at the cost of exploration. Similarly, we quantify the knowing-doing gap and find that LLMs often know how to solve a task (87% correct rationales) but fail at acting on this knowledge as they prioritize greedy actions (64% of actions when rationale is correct).

To overcome these shortcomings, we propose Reinforcement Learning Fine-Tuning (RLFT) on self-generated CoT rationales. RL is the pre-dominant learning paradigm in decision-making scenarios and has been successful in game-playing [Silver et al., 2016, Vinyals et al., 2019], robotics [Tirumala et al., 2025], plasma-control [Degraeve et al., 2022], or navigating stratospheric balloons [Bellemare et al., 2020]. We study the effects of RLFT on pre-trained Gemma2 models [Team et al., 2024b,c] in three sizes (2B, 9B, and 27B) in multi-arm bandit (MAB) and contextual bandit (CB) settings proposed by Nie et al. [2024], and the textual Tic-tac-toe environment released by Ruoss et al. [2024]. Across environments, we find that RLFT enhances the decision-making abilities of LLMs by increasing exploration and narrowing the knowing-doing gap. While RLFT positively affects exploration of LLM agents, their exploration strategies remain sub-optimal. Therefore, we empirically evaluate both “classic” exploration mechanisms commonly employed in RL, such as ϵ -greedy, and LLM-specific approaches, such as self-correction and self-consistency, to enable more effective fine-tuning for decision-making scenarios. Finally, in our ablations we investigate the importance of CoT reasoning for decision-making, highlight the effectiveness of leveraging expert data, and show the benefits of giving the agent more reasoning tokens to solve the decision-making problem.

In summary, we make the following **contributions**:

- We systematically examine three failure modes of small-to-medium scale LLMs in decision-making scenarios: greediness, frequency bias, and the knowing-doing gap.
- We study how RL fine-tuning on self-generated CoT rationales affects these shortcomings, highlighting positive effects of RLFT on exploration and decision-making abilities.
- We evaluate a variety of exploration mechanisms (e.g., ϵ -greedy) and LLM-specific approaches (e.g., self-consistency), to enable more effective RLFT for LLMs.

2 Related Work

We discuss related works in Appendix A.

72 3 Methodology

73 3.1 Background

74 **Reinforcement Learning.** We assume the standard RL formulation via a Markov Decision Process
 75 (MDP) represented by a tuple of $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R})$, where \mathcal{S} and \mathcal{A} denote state and action spaces,
 76 respectively. At every timestep t the agent observes state $s_t \in \mathcal{S}$, predicts action $a_t \in \mathcal{A}$, and receives
 77 a reward r_t given by the reward function $\mathcal{R}(s_t, a_t)$. $\mathcal{P}(s_{t+1} | s_t, a_t)$ defines the transition dynamics
 78 constituting a probability distribution over next states s_{t+1} . The goal of RL is to learn a policy
 79 $\pi_\theta(a_t | s_t)$ with parameters θ that predicts an action a_t in state s_t that maximizes cumulative reward.

80 **Reinforcement Learning from Human Feedback.** RLHF aims to fine-tune pre-trained models
 81 towards human preferences [Christiano et al., 2017]. Preferences are typically encoded via a reward
 82 model r_ϕ with parameters ϕ learned from a human annotated dataset \mathcal{D} consisting of query-response
 83 pairs x and y , respectively. RLHF optimizes a constrained REINFORCE estimator [Williams, 1992]:

$$\max_{\theta} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(\cdot | x)} [(r_\phi(x, y) - b) \nabla_{\theta} \log \pi_\theta(y | x) - \beta D_{KL}(\pi_\theta(\cdot | x) || \pi_{ref}(\cdot | x))] \quad (1)$$

84 Here π_{ref} is a reference policy, which is typically the frozen pre-trained model, and β is as weighting
 85 term. The baseline b represents a baseline to reduce variance and is commonly instantiated by a value
 86 function [Schulman et al., 2017, Ouyang et al., 2022] or a Monte-Carlo (MC) estimate of the returns
 87 [Ahmadian et al., 2024, Ramesh et al., 2024, Shao et al., 2024].

88 3.2 Reinforcement Learning Fine-Tuning (RLFT)

89 Our RLFT approach relies on fine-tuning on self-generated CoT rationales on rewards obtained from
 90 environment interaction. During RLFT the model learns to iteratively refine its reasoning process,
 91 favoring CoT patterns and actions that lead to higher rewards (see Figure 1). Our approach is similar
 92 to Guo et al. [2025] and Zhai et al. [2025], but specialized for decision-making scenarios.

93 **Context Representation.** The input tokens to our model at step t consists of input instructions
 94 c_t^{in} , output instructions c_t^{out} , and the most recent interaction history $c_t^{\tau_{t-C:t}}$ (see Figure 1). The
 95 history representation contains the trajectory $\tau_{t-C:t} = (s_{t-C}, a_{t-C}, r_{t-C}, \dots, s_t, a_t, r_t)$ of the C
 96 most recent states, actions, and rewards. We opt for task-specific instructions for c_t^{in} rather than a
 97 generic instruction template, providing the agent with information about the observations, the possible
 98 actions, and its objective. Consequently, c_t is represented by the concatenation of the instruction and
 99 history tokens $c_t = [c_t^{in}; c_t^{out}; c_t^{\tau_{t-C:t}}]$.

100 **Factorization of Action Tokens.** At every interaction step t , the agent generates action tokens
 101 $z_t = [z_t^{CoT}; a_t]$ containing both the CoT reasoning tokens z_t^{CoT} and the action to be executed in the
 102 environment a_t . To extract a_t from z_t , we make use of an extraction function $a_t = g(z_t)$. In practice,
 103 g consists of regular expressions to match the output pattern given by c_t^{out} . If no valid action is found
 104 a random action is executed. To allow for flexibility in refining the reasoning process, we opt for a
 105 permissive output template (i.e., ACTION=X), rather than enforcing a structured output template (e.g.,
 106 <thought> and <action> blocks). We employ a token generation budget of G tokens ($G = 256$ by
 107 default), therefore $|z_t| \leq G$.

108 **Reward Shaping for Valid Actions.** In addition to the environment reward r_t^{env} , we employ a reward
 109 shaping term r_t^{valid} to encourage the model to adhere to the output template, $r_t = r_t^{env} + r_t^{valid}$.
 110 More specifically, we make use of a reward penalty of -5 if g cannot extract a valid action, $r_t^{valid} =$
 111 $-5 \cdot 1(g(a_t^{act}) \notin \mathcal{A})$. To ensure that the reward penalty does not overly bias optimization, we employ
 112 reward normalization to the environment rewards.

113 **Fine-tuning objective.** We fine-tune using the clipping objective introduced by Schulman et al.
 114 [2017] with an additional KL constraint to the reference policy π_{ref} :

$$\max_{\theta} \mathbb{E}_{(c, z) \sim \mathcal{D}} \left[\min \left(\frac{\pi_\theta(z|c)}{\pi_{\theta_{old}}(z|c)} A_{adv}, \text{clip}_\epsilon \left(\frac{\pi_\theta(z|c)}{\pi_{\theta_{old}}(z|c)} \right) A_{adv} \right) - \beta D_{KL}(\pi_\theta(\cdot | c) || \pi_{ref}(\cdot | c)) \right] \quad (2)$$

115 Here $\pi_{\theta_{old}}$ refers to the rollout generating policy, D is the rollout buffer, and ϵ is a hyperparameter. To
 116 allow for memory efficient fine-tuning in environments with fixed episode lengths (bandits), we make
 117 use of a Monte Carlo baseline to estimate A_{adv} . Instead of exploiting multiple rollouts, as used by

118 [Ahmadian et al. \[2024\]](#) and [Ramesh et al. \[2024\]](#), we compute rewards-to-go. For environments with
 119 variable episode lengths (Tic-tac-toe), we learn a separate state-value head on top of the last layer
 120 LLM representations and make use of generalized advantage estimation [[Schulman et al., 2015](#)].
 121 We provide additional implementation and training details in Appendix C.

122 4 Experiments

123 We study the effect of fine-tuning Gemma2 [[Team et al., 2024b,c](#)] models in MAB and CB settings
 124 proposed by [Nie et al. \[2024\]](#), and on a text-based version of Tic-tac-toe released by [Paglieri et al.](#)
 125 [[2024](#)]. We describe our environments and baselines in Section 4.1. For our experiments, we compare
 126 Gemma2 [[Team et al., 2024c](#)] models at three model scales: 2B, 9B and 27B. In Section 4.2, we first
 127 analyze three common failure modes of LLM agents in MAB scenarios: (1) greediness, (2) frequency
 128 bias, and (3) the knowing-doing gap. Then we investigate the effects of fine-tuning on self-generated
 129 CoT rationales or expert rationales in MABs and CBs (see Section 4.3), and in Tic-tac-toe (see
 130 Section 4.5). In Section 4.4, we study the effects of a exploration mechanisms on the fine-tuning
 131 performance. Finally, in Section 4.5 we empirically examine important components of our approach.

132 4.1 Environments & Baselines

133 **Multi-armed and Contextual Bandits.** MABs [[Slivkins et al., 2019](#), [Lattimore and Szepesvári,](#)
 134 [2020](#)] are a classic problem setting in RL that isolates the *exploration-exploitation* trade-
 135 off. For our MAB experiments, we leverage the text-based bandit scenarios released by [Nie](#)
 136 [et al. \[2024\]](#). We focus on the *continuous* and *button* variants, as illustrated in Figure 2.
 137 We report results for MAB with $k \in \{5, 10, 20\}$
 138 arms ($|\mathcal{A}| = k$) and payoffs of the arms being
 139 either Gaussian or Bernoulli distributed. In ad-
 140 dition, we consider three levels of stochasticity
 141 (low/medium/high) that determine the standard
 142 deviation or delta gap in Gaussian or Bernoulli
 143 bandits, respectively. For all MAB settings, we
 144 limit the horizon T to 50 interaction steps. We
 145 compare against two commonly used baselines
 146 for MABs: Upper-confidence Bound (UCB)
 147 [[Auer, 2002](#)] and a random agent that selects
 148 actions uniformly at random. UCB is consid-
 149 ered optimal and represents the upper-bound for
 150 agent performance, whereas the random base-
 151 line represents the lower bound. We provide
 152 more details on our MAB and CB setups in Ap-
 153 pendices B.1 and B.2, respectively.

154 **Tic-tac-toe.** In addition, we use the text-based
 155 Tic-tac-toe environment released by [Ruoss et al.](#)
 156 [[2024](#)], which exhibits proper state transitions.
 157 [Ruoss et al. \[2024\]](#) demonstrated that frontier
 158 models struggle to achieve strong performance
 159 in this environment and barely beat a random
 160 opponent. Consequently, it is a good target to in-
 161 vestigate the efficacy of RLFT. In Appendix B.3,
 162 we provide addition details on our environment
 163 and training setup.

164 4.2 Why do LLMs perform suboptimally in decision-making?

165 Prior works found that LLM agents perform suboptimally and fail to explore sufficiently in interactive
 166 settings [[Paglieri et al., 2024](#), [Ruoss et al., 2024](#)]. Therefore, we first examine *why* models perform
 167 suboptimally and identify three prevalent failure modes: (1) greediness, (2) frequency bias, and (3)
 168 the knowing-doing gap. In this section, we present analyses of Gemma2 models when given input

Button Multi-armed Bandit (Gaussian)

You are a bandit algorithm in a room with
 5 buttons labeled red, green, blue, yellow,
 orange. [...]. Your goal is to maximize the total
 reward. [More instructions]

Think step-by-step and output your final
 answer in the format ACTION=X where X is
 one of the arms listed above. IMPORTANT:
 Provide your (SHORT!) thinking process and
 your answer ACTION=X

So far you have tried/seen:

Step=0 Action=green Reward=0.3
 Step=1 Action=blue Reward=0.1
 Step=2 Action=orange Reward=-0.5
 Step=3 Action=red Reward=0.5
 Step=4 Action=green Reward=0.24
 ...
 What do you predict next?

Figure 2: Illustration of a Gaussian MAB for the *button* scenario from [[Nie et al., 2024](#)] using our context representation and instructions.

169 contexts that elucidate the failure modes. We conduct our analyses on the *button* instance of our
 170 MAB experiments at three model scales, and find that the failure modes persist across model scales
 171 (see Appendix D.1 for *continuous* instance).

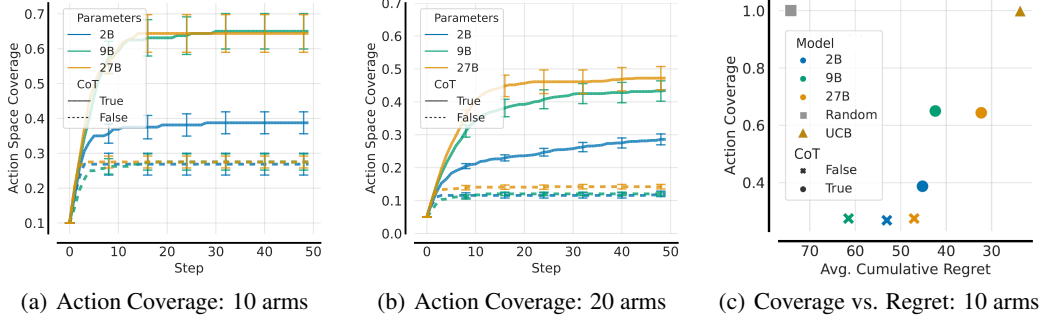


Figure 3: Illustration of Greediness. We show action coverage for Gemma2 2B/9B/27B w/ and w/o CoT for (a) 10 and (b) 20 arms over 50 interaction steps. Agents favor the best performing action among the set of selected actions, leading to stagnating action coverage, despite benefits of larger models and CoT. In (c), we plot cumulative regret against action coverage. The agents exhibit suboptimal regret, because of greedy action selection strategies.

172 **Greediness.** The first and most pervasive failure mode is *greediness*, which is characterized by the
 173 LLM overly favoring the best performing action among a small set of actions seen so far. To illustrate
 174 this failure mode, we show the average action coverage achieved by Gemma2 2B/9B/27B with and
 175 without CoT across 64 MABs with 10 and 20 arms over 50 interaction steps (see Figure 3 a and b).
 176 We define action coverage C_t at step t as the fraction of available actions that have been selected at
 177 least once, $C_t = \frac{|\{a \in \mathcal{A}: N_t(a) > 0\}|}{|\mathcal{A}|}$ with $N_t(a)$ representing the number of times action $a \in \mathcal{A}$ has been
 178 selected until t . For 10 arms and averaged over 64 parallel environments, we find that Gemma2 2B
 179 covers 40% of all actions, while 9B/27B cover 65% (i.e., 6.5 actions), leaving a significant part of
 180 the action space unexplored. Note that without CoT all models explore merely 25% of all actions in
 181 the 10 arms setting. The suboptimal coverage is caused by the model overly favoring high-reward
 182 actions (see Figure 12 in Appendix D.1.1). Consequently, the model prematurely commits to a greedy
 183 strategy leading to a stagnating action coverage beyond 10 steps. Increasing the number of arms
 184 makes the greediness even more apparent, with the largest models only covering 45% of all actions.
 185 Due to this, the regret remains high compared to UCB, even though the models improve significantly
 186 over a random agent (see Figure 3c).

187 **Frequency Bias.** The next prevalent failure mode is *frequency bias*, which is characterized by
 188 repeatedly selecting the most frequently occurring action in the context, even when that action gives
 189 low reward. To understand how the model’s behavior is influenced by the frequency of actions, we
 190 construct prefix histories using a random policy, vary the number of repetitions of the last action in the
 191 context history (0 to 100) and record the entropy over all actions (see Figure 4a and c). We provide
 192 details on the context generation in Appendix D.1.2. To quantify frequency bias, we categorize an
 193 action as *frequent* action $a_f = \arg \max_{a \in \mathcal{A}} N_T(a)$, *greedy* $a_g = \arg \max_{a \in \{a \in \mathcal{A}: N_T(a) > 0\}} R_T(a)$,
 194 or *other* if they are neither frequent nor greedy. Note that action is optimal with 10% probability.
 195 Subsequently, we compute the frequent F_f , greedy F_g and other F_o fractions as reported in Figure 4
 196 (see Appendix 4 for definitions).

197 Gemma2 2B heavily suffers from repeated actions, exhibiting a decreasing entropy with increasing
 198 repetitions (96% F_f , see Figure 4a). In contrast, 27B escapes the frequency bias (14%, see Figure
 199 4c) and interestingly becomes less certain of its action prediction with increasing repetitions. To
 200 examine this further, we show the bucketized fractions with 0-10, 45-55 and 90-100 repetitions for
 201 2B and 27B in Figure 4b. Indeed, for 2B F_f keeps increasing with increasing repetitions. While 27B
 202 escapes the frequency bias it suffers heavily from greediness. Similar biases have been identified in
 203 Behavior Cloning (BC) settings and termed *copycat* bias [Wen et al., 2020, Schmied et al., 2024b].
 204 This suggests that frequency bias is an artifact of supervised pre-training, and motivates the use of RL
 205 as a counter-measurement.

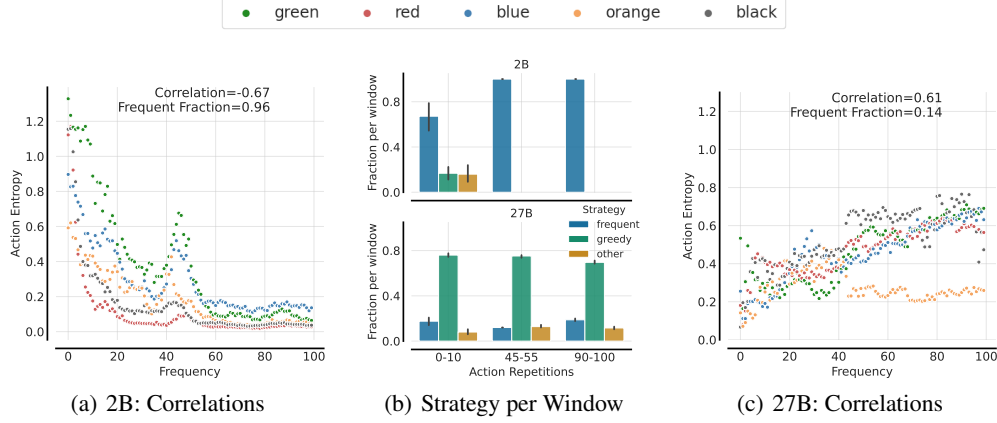


Figure 4: Illustration of Frequency Bias. We plot the frequency of the repeated action in the context against the action entropy across all actions for 10 armed MABs. (a) Gemma2 2B heavily suffers from frequency bias, becoming more certain of the most frequent action, the more often it occurs in the context. (c) Gemma2 27B overcomes the frequency bias, but instead behaves greedily. In (b) we show the action strategies for three repetition windows.

Knowing-Doing Gap. The *knowing-doing gap* has been observed by Paglieri et al. [2024] and Ruoss et al. [2024]. To illustrate the gap in our setting, we first task Gemma2 27B to produce the UCB algorithm, to compute the relevant quantities accordingly ("knowing"), and finally to act according to the computed quantities ("doing", see Figure 21 for the instructions and an agent response). We let Gemma2 27B interact with the environment (64 instances) for 50 timesteps with $G = 2048$ per step, and extract the UCB quantities from the rationales. To quantify "knowing", we compare the UCB values computed by the model against the real UCB values, and consider the rationale z_{CoT} as correct if the arm with the highest UCB values match (see Appendix D.1.3 for details). To quantify "doing", we categorize the generated actions as *optimal* action if the model selects the action with the highest UCB value, as *greedy* if it selects the action with the highest UCB value among the set of actions tried so far, and as *other* if the action is neither optimal nor greedy. Subsequently, we compute the percentages of greedy/optimal/other actions. The agent clearly knows how to solve the task, with 87% of all rationales being correct (see Figure 5). However, even for correctly computed rationales, the model often selects the greedy action (58%) over the optimal action (21%). This discrepancy highlights the shortcomings of the LLM when it comes to "acting" even when "knowing" the algorithm.

		"Doing": Action Strategy		
		greedy	optimal	other
"Knowing": Correctness	True	58.1 _{±2}	21.6 _{±2}	10.6 _{±1}
	False	7.2 _{±1}	0.6 _{±0}	2.0 _{±0}

Figure 5: Confusion matrix for the Knowing-Doing Gap of Gemma2 27B. The agent "knows" how to solve the task (87% correct rationales, sum of top row), but fails at "doing" (58% greedy actions among correct rationales). See Figure 21, for instructions and an agent response.

4.3 Effectiveness of RL Fine-Tuning

Next, we study the effects of RLFT on cumulative regret (w.r.t. optimal policy) and whether it alleviates the highlighted failure modes. We fine-tune Gemma2 2B and 9B on self-generated CoT rationales for 30K updates with an (accumulated) batch size of 128. To avoid memorization of reward distributions, we maintain a pool of 512 MABs and randomly select a subset of 16 MABs per rollout. We refer to Appendix C for training details and hyperparameters.

RLFT lowers regret. In Figure 6, we report the cumulative regrets across model sizes and arms for a medium noise $\sigma = 1.0$ scenario (see Appendix D.2 for low/high noise). Across environments, the LLMs clearly outperform the random baseline and RLFT lowers regret for both 2B and 9B. For 2B, RLFT narrows the gap to its larger counterparts and UCB. Similarly, RLFT lowers regret for

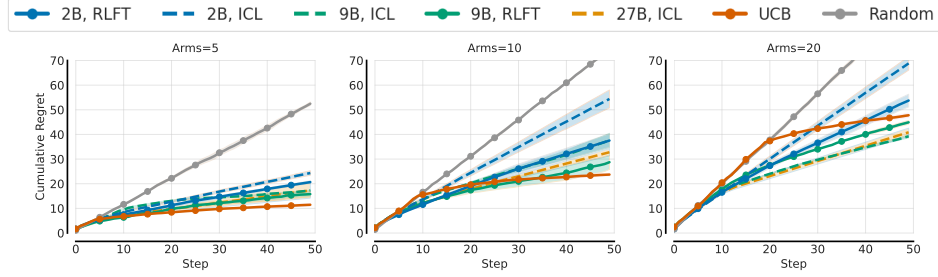


Figure 6: Main Comparison on **Gaussian MABs** button scenario in the medium noise ($\sigma = 1$) setting. We compare cumulative regrets (lower is better) of classic baselines against ICL and RLFT performances for 5, 10, and 20 arms. See Figure 16 for $\sigma = 0.1$ and $\sigma = 3$.

240 Gemma2 9B. Note, that the lower cumulative regret of Gemma2 9/27B compared to UCB after
 241 50 environment steps in the 20 arms scenario is an artifact of the limited interaction steps, but the
 242 trends remain clear. We repeat RLFT for CBs, and observe similar performance improvements for
 243 Gemma2 2B (see Appendix D.3). Consequently, reinforcing self-generated CoT rationales towards
 244 environment rewards improves performance on simple decision-making scenarios.

245 **RLFT mitigates greediness.** In Figure 14a, we report the action coverage for 2B after RLFT at
 246 different numbers of gradient steps (10K, 20K, 30K). Indeed, we observe that RLFT results in
 247 increased action coverage (+12%) after 30K updates. Interestingly, we first observe a decrease (at
 248 10K) followed by an increase in action coverage (20K, 30K). We observe similar effects for the 20
 249 arms scenario (see Figure 14b). Via RLFT the agent learns to explore and mitigates greediness.

250 **RLFT counteracts frequency bias.** We find that RLFT counteracts frequency bias (Figure 15).
 251 In particular, for 0-10 repetitions we observe a strong decrease in the fraction of frequent actions
 252 ($70\% \rightarrow 35\%$) and increase in "other" actions ($8\% \rightarrow 35\%$). However, F_f remains elevated for high
 253 repetitions. Consequently, RLFT counteracts frequency bias, but does not fully alleviate it.

254 4.4 Effect of Exploration Mechanisms

255 For RLFT, we relied solely on the exploration properties for CoT reasoning. However, it is common
 256 practice to employ additional exploration strategies [Mnih et al., 2015, Schulman et al., 2017, Haarnoja
 257 et al., 2018]. Therefore, we study the effects of classic exploration mechanisms and LLM-specific
 258 strategies to encourage exploration (see Appendix C.4). Across model scales, we observe that the
 259 mechanisms result in varied effects on action coverage (see Figure 18). First, we find that the simple
 260 *try-all* strategy, which reduces the need for additional exploration by trying all actions, results in
 261 the biggest performance improvements. Second, a simple *exploration bonus* (+1 reward for untried
 262 actions during RLFT), significantly increases exploration ($50\% \rightarrow 70\%$) and lowers regret towards
 263 the expert compared to regular RLFT. This highlights the importance of reward shaping for fine-tuning
 264 LLMs to elucidate a desired behavior.

265 4.5 Ablations

266 We conduct additional ablations on the effect of RLFT in Tic-tac-toe, the importance of CoT for
 267 RLFT, leveraging pre-collected datasets via SFT, and the effect of "thinking" tokens in Appendix E.

268 5 Conclusion

269 We study *why* LLMs perform sub-optimally in decision-making scenarios and examine three prevalent
 270 failure modes: greediness, frequency bias, and the knowing-doing gap. We show that RLFT on CoT
 271 rationales mitigates greediness, counteracts frequency bias, and improves final performance. While
 272 RLFT improves the exploration, it remains sub-optimal compared to bandit algorithms. Therefore,
 273 we investigate a variety of mechanisms to improve exploration. Models act near-optimally if provided
 274 with sufficient information underscoring their shortcomings in exploration. Finally, we highlight the
 275 importance of reward shaping for RLFT. In Appendix F, we discuss limitations and future work.

References

- R. Agarwal, M. Schwarzer, P. S. Castro, A. C. Courville, and M. Bellemare. Deep reinforcement learning at the edge of the statistical precipice. *Advances in neural information processing systems*, 34:29304–29320, 2021.
- A. Ahmadian, C. Cremer, M. Gallé, M. Fadaee, J. Kreutzer, A. Üstün, and S. Hooker. Back to basics: Revisiting reinforce style optimization for learning from human feedback in llms. *arXiv preprint arXiv:2402.14740*, 2024.
- J. A. Arjona-Medina, M. Gillhofer, M. Widrich, T. Unterthiner, J. Brandstetter, and S. Hochreiter. Rudder: Return decomposition for delayed rewards. *Advances in Neural Information Processing Systems*, 32, 2019.
- P. Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, 3(Nov):397–422, 2002.
- M. Beck, K. Pöppel, M. Spanring, A. Auer, O. Prudnikova, M. Kopp, G. Klambauer, J. Brandstetter, and S. Hochreiter. xlstm: Extended long short-term memory. *Advances in Neural Information Processing Systems*, 37:107547–107603, 2025.
- M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279, 2013.
- M. G. Bellemare, S. Candido, P. S. Castro, J. Gong, M. C. Machado, S. Moitra, S. S. Ponda, and Z. Wang. Autonomous navigation of stratospheric balloons using reinforcement learning. *Nature*, 588(7836):77–82, 2020.
- A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, J. Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.
- A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choromanski, T. Ding, D. Driess, A. Dubey, C. Finn, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023.
- T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020a.
- T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020b.
- Y. Burda, H. Edwards, A. Storkey, and O. Klimov. Exploration by random network distillation. *arXiv preprint arXiv:1810.12894*, 2018.
- S. Chan, A. Santoro, A. K. Lampinen, J. Wang, A. Singh, P. H. Richemond, J. L. McClelland, and F. Hill. Data distributional properties drive emergent in-context learning in transformers. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022.
- L. Chen, L. Wang, H. Dong, Y. Du, J. Yan, F. Yang, S. Li, P. Zhao, S. Qin, S. Rajmohan, et al. Introspective tips: Large language model for in-context decision making. *arXiv preprint arXiv:2305.11598*, 2023.
- P. F. Christiano, J. Leike, T. Brown, M. Martic, S. Legg, and D. Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- W. Chu, L. Li, L. Reyzin, and R. Schapire. Contextual bandits with linear payoff functions. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 208–214. JMLR Workshop and Conference Proceedings, 2011.

324 K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plappert, J. Tworek, J. Hilton,
325 R. Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*,
326 2021.

327 R. Coulom. Efficient selectivity and backup operators in monte-carlo tree search. In *International*
328 *conference on computers and games*, pages 72–83. Springer, 2006.

329 C. Cundy and S. Ermon. Sequencematch: Imitation learning for autoregressive sequence modelling
330 with backtracking. *arXiv preprint arXiv:2306.05426*, 2023.

331 S. De, S. L. Smith, A. Fernando, A. Botev, G. Cristian-Muraru, A. Gu, R. Haroun, L. Berrada,
332 Y. Chen, S. Srinivasan, et al. Griffin: Mixing gated linear recurrences with local attention for
333 efficient language models. *arXiv preprint arXiv:2402.19427*, 2024.

334 J. Degraeve, F. Felici, J. Buchli, M. Neunert, B. Tracey, F. Carpanese, T. Ewalds, R. Hafner, A. Abdol-
335 maleki, D. de Las Casas, et al. Magnetic control of tokamak plasmas through deep reinforcement
336 learning. *Nature*, 602(7897):414–419, 2022.

337 Y. Duan, J. Schulman, X. Chen, P. L. Bartlett, I. Sutskever, and P. Abbeel. R12: Fast reinforcement
338 learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*, 2016.

339 A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang,
340 A. Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

341 Z. Durante, Q. Huang, N. Wake, R. Gong, J. S. Park, B. Sarkar, R. Taori, Y. Noda, D. Terzopoulos,
342 Y. Choi, et al. Agent ai: Surveying the horizons of multimodal interaction. *arXiv preprint*
343 *arXiv:2401.03568*, 2024.

344 A. Ecoffet, J. Huizinga, J. Lehman, K. O. Stanley, and J. Clune. Go-explore: a new approach for
345 hard-exploration problems. *arXiv preprint arXiv:1901.10995*, 2019.

346 C. Finn, P. Abbeel, and S. Levine. Model-agnostic meta-learning for fast adaptation of deep networks.
347 In *International conference on machine learning*, pages 1126–1135. PMLR, 2017.

348 S. Flennerhag, A. A. Rusu, R. Pascanu, F. Visin, H. Yin, and R. Hadsell. Meta-learning with warped
349 gradient descent. *arXiv preprint arXiv:1909.00025*, 2019.

350 O. Groth, M. Wulfmeier, G. Vezzani, V. Dasagi, T. Hertweck, R. Hafner, N. Heess, and M. Riedmiller.
351 Is curiosity all you need? on the utility of emergent behaviours from curious exploration. *arXiv*
352 *e-prints*, pages arXiv–2109, 2021.

353 D. Guo, D. Yang, H. Zhang, J. Song, R. Zhang, R. Xu, Q. Zhu, S. Ma, P. Wang, X. Bi, et al.
354 Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint*
355 *arXiv:2501.12948*, 2025.

356 T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta,
357 P. Abbeel, et al. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*,
358 2018.

359 D. Hafner. Benchmarking the spectrum of agent capabilities. *arXiv preprint arXiv:2109.06780*, 2021.

360 F. M. Harper and J. A. Konstan. The movielens datasets: History and context. *Acm transactions on*
361 *interactive intelligent systems (tiis)*, 5(4):1–19, 2015.

362 H. He, W. Yao, K. Ma, W. Yu, Y. Dai, H. Zhang, Z. Lan, and D. Yu. Webvoyager: Building an
363 end-to-end web agent with large multimodal models. *arXiv preprint arXiv:2401.13919*, 2024.

364 D. Hendrycks, C. Burns, S. Kadavath, A. Arora, S. Basart, E. Tang, D. Song, and J. Steinhardt.
365 Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*,
366 2021.

367 M. Hessel, J. Modayil, H. Van Hasselt, T. Schaul, G. Ostrovski, W. Dabney, D. Horgan, B. Piot,
368 M. Azar, and D. Silver. Rainbow: Combining improvements in deep reinforcement learning. In
369 *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.

370 S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780,
371 1997.

372 S. Hochreiter, A. S. Younger, and P. R. Conwell. Learning to learn using gradient descent. In *Artificial*
373 *Neural Networks—ICANN 2001: International Conference Vienna, Austria, August 21–25, 2001*
374 *Proceedings 11*, pages 87–94. Springer, 2001.

375 E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, W. Chen, et al. Lora: Low-rank
376 adaptation of large language models. *ICLR*, 1(2):3, 2022.

377 S. Hu and J. Clune. Thought cloning: Learning to think while acting by imitating human thinking.
378 *Advances in Neural Information Processing Systems*, 36:44451–44469, 2023.

379 L. Kirsch, S. van Steenkiste, and J. Schmidhuber. Improving generalization in meta reinforcement
380 learning using learned objectives. *arXiv preprint arXiv:1910.04098*, 2019.

381 L. Kirsch, J. Harrison, J. Sohl-Dickstein, and L. Metz. General-purpose in-context learning by
382 meta-learning transformers. *arXiv preprint arXiv:2212.04458*, 2022.

383 L. Kirsch, J. Harrison, C. Freeman, J. Sohl-Dickstein, and J. Schmidhuber. Towards general-purpose
384 in-context learning agents. In *NeurIPS 2023 Workshop on Generalization in Planning*, 2023.

385 M. Klissarov, P. D’Oro, S. Sodhani, R. Raileanu, P.-L. Bacon, P. Vincent, A. Zhang, and M. Henaff.
386 Motif: Intrinsic motivation from artificial intelligence feedback. *arXiv preprint arXiv:2310.00166*,
387 2023.

388 M. Klissarov, M. Henaff, R. Raileanu, S. Sodhani, P. Vincent, A. Zhang, P.-L. Bacon, D. Precup,
389 M. C. Machado, and P. D’Oro. Maestromotif: Skill design from artificial intelligence feedback.
390 *arXiv preprint arXiv:2412.08542*, 2024.

391 A. Krishnamurthy, K. Harris, D. J. Foster, C. Zhang, and A. Slivkins. Can large language models
392 explore in-context? *arXiv preprint arXiv:2403.15371*, 2024.

393 A. Kumar, V. Zhuang, R. Agarwal, Y. Su, J. D. Co-Reyes, A. Singh, K. Baumli, S. Iqbal, C. Bishop,
394 R. Roelofs, et al. Training language models to self-correct via reinforcement learning. *arXiv*
395 *preprint arXiv:2409.12917*, 2024.

396 M. Laskin, L. Wang, J. Oh, E. Parisotto, S. Spencer, R. Steigerwald, D. Strouse, S. Hansen, A. Filos,
397 E. Brooks, et al. In-context reinforcement learning with algorithm distillation. *arXiv preprint*
398 *arXiv:2210.14215*, 2022.

399 T. Lattimore and C. Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.

400 K.-H. Lee, O. Nachum, M. Yang, L. Lee, D. Freeman, W. Xu, S. Guadarrama, I. Fischer, E. Jang,
401 H. Michalewski, et al. Multi-game decision transformers. *arXiv preprint arXiv:2205.15241*, 2022.

402 Y. Li, D. Choi, J. Chung, N. Kushman, J. Schrittwieser, R. Leblond, T. Eccles, J. Keeling, F. Gimeno,
403 A. Dal Lago, et al. Competition-level code generation with alphacode. *Science*, 378(6624):
404 1092–1097, 2022.

405 C. Lu, S. Hu, and J. Clune. Intelligent go-explore: Standing on the shoulders of giant foundation
406 models. *arXiv preprint arXiv:2405.15143*, 2024.

407 G. Mialon, C. Fourrier, T. Wolf, Y. LeCun, and T. Scialom. Gaia: a benchmark for general ai
408 assistants. In *The Twelfth International Conference on Learning Representations*, 2023.

409 S. Mirchandani, F. Xia, P. Florence, B. Ichter, D. Driess, M. G. Arenas, K. Rao, D. Sadigh, and
410 A. Zeng. Large language models as general pattern machines. *arXiv preprint arXiv:2307.04721*,
411 2023.

412 N. Mishra, M. Rohaninejad, X. Chen, and P. Abbeel. A simple neural attentive meta-learner. In
413 *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada,*
414 *April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL <https://openreview.net/forum?id=B1DmUzWAW>.

416 V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Ried-
417 miller, A. K. Fidjeland, G. Ostrovski, et al. Human-level control through deep reinforcement
418 learning. *nature*, 518(7540):529–533, 2015.

419 G. Monea, A. Bosselut, K. Brantley, and Y. Artzi. Llms are in-context reinforcement learners. *arXiv*
420 *preprint arXiv:2410.05362*, 2024.

421 N. Muennighoff, Z. Yang, W. Shi, X. L. Li, L. Fei-Fei, H. Hajishirzi, L. Zettlemoyer, P. Liang,
422 E. Candès, and T. Hashimoto. s1: Simple test-time scaling. *arXiv preprint arXiv:2501.19393*,
423 2025.

424 A. Nie, Y. Su, B. Chang, J. N. Lee, E. H. Chi, Q. V. Le, and M. Chen. Evolve: Evaluating and
425 optimizing llms for exploration. *arXiv preprint arXiv:2410.06238*, 2024.

426 P.-Y. Oudeyer, F. Kaplan, and V. V. Hafner. Intrinsic motivation systems for autonomous mental
427 development. *IEEE transactions on evolutionary computation*, 11(2):265–286, 2007.

428 L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama,
429 A. Ray, et al. Training language models to follow instructions with human feedback. *Advances in*
430 *neural information processing systems*, 35:27730–27744, 2022.

431 D. Paglieri, B. Cupiał, S. Coward, U. Piterbarg, M. Wolczyk, A. Khan, E. Pignatelli, Ł. Kuciński,
432 L. Pinto, R. Fergus, et al. Balrog: Benchmarking agentic llm and vlm reasoning on games. *arXiv*
433 *preprint arXiv:2411.13543*, 2024.

434 D. Pathak, P. Agrawal, A. A. Efros, and T. Darrell. Curiosity-driven exploration by self-supervised
435 prediction. In *International conference on machine learning*, pages 2778–2787. PMLR, 2017.

436 D. A. Pomerleau. Alvin: An autonomous land vehicle in a neural network. *Advances in neural*
437 *information processing systems*, 1, 1988.

438 A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, et al. Language models are
439 unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.

440 R. Raileanu and T. Rocktäschel. Ride: Rewarding impact-driven exploration for procedurally-
441 generated environments. *arXiv preprint arXiv:2002.12292*, 2020.

442 S. S. Ramesh, Y. Hu, I. Chaimalas, V. Mehta, P. G. Sessa, H. B. Ammar, and I. Bogunovic. Group
443 robust preference optimization in reward-free rlhf. *arXiv preprint arXiv:2405.20304*, 2024.

444 D. Rao, F. Sadeghi, L. Hasenclever, M. Wulfmeier, M. Zambelli, G. Vezzani, D. Tirumala, Y. Aytar,
445 J. Merel, N. Heess, et al. Learning transferable motor skills with hierarchical latent mixture policies.
446 In *International Conference on Learning Representations*, 2021.

447 S. C. Raparthy, E. Hambro, R. Kirk, M. Henaff, and R. Raileanu. Generalization to new sequential
448 decision making tasks with in-context learning, 2023.

449 S. Reed, K. Zolna, E. Parisotto, S. G. Colmenarejo, A. Novikov, G. Barth-Maron, M. Gimenez,
450 Y. Sulsky, J. Kay, J. T. Springenberg, et al. A generalist agent. *arXiv preprint arXiv:2205.06175*,
451 2022.

452 A. Ruoss, F. Pardo, H. Chan, B. Li, V. Mnih, and T. Genewein. Lmact: A benchmark for in-context
453 imitation learning with long multimodal demonstrations. *arXiv preprint arXiv:2412.01441*, 2024.

454 K. Saab, T. Tu, W.-H. Weng, R. Tanno, D. Stutz, E. Wulczyn, F. Zhang, T. Strother, C. Park, E. Vedadi,
455 et al. Capabilities of gemini models in medicine. *arXiv preprint arXiv:2404.18416*, 2024.

456 A. Santoro, S. Bartunov, M. Botvinick, D. Wierstra, and T. Lillicrap. Meta-learning with memory-
457 augmented neural networks. In *International conference on machine learning*, pages 1842–1850.
458 PMLR, 2016.

459 J. Schmidhuber. Evolutionary principles in self-referential learning. on learning now to learn: The
460 meta-meta-meta...-hook. Diploma thesis, Technische Universität München, Germany, 14 May
461 1987.

462 J. Schmidhuber. Curious model-building control systems. In *Proc. international joint conference on*
463 *neural networks*, pages 1458–1463, 1991a.

464 J. Schmidhuber. A possibility for implementing curiosity and boredom in model-building neural
465 controllers. In *Proc. of the international conference on simulation of adaptive behavior: From*
466 *animals to animats*, pages 222–227, 1991b.

467 T. Schmied, M. Hofmarcher, F. Paischer, R. Pascanu, and S. Hochreiter. Learning to modulate
468 pre-trained models in rl. *Advances in Neural Information Processing Systems*, 36:38231–38265,
469 2023.

470 T. Schmied, T. Adler, V. Patil, M. Beck, K. Pöppel, J. Brandstetter, G. Klambauer, R. Pascanu, and
471 S. Hochreiter. A large recurrent action model: xlstm enables fast inference for robotics tasks. *arXiv*
472 *preprint arXiv:2410.22391*, 2024a.

473 T. Schmied, F. Paischer, V. Patil, M. Hofmarcher, R. Pascanu, and S. Hochreiter. Retrieval-augmented
474 decision transformer: External memory for in-context rl. *arXiv preprint arXiv:2410.07071*, 2024b.

475 J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel. High-dimensional continuous control
476 using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.

477 J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization
478 algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

479 Z. Shao, P. Wang, Q. Zhu, R. Xu, J. Song, X. Bi, H. Zhang, M. Zhang, Y. Li, Y. Wu, et al.
480 Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv*
481 *preprint arXiv:2402.03300*, 2024.

482 N. Shazeer and M. Stern. Adafactor: Adaptive learning rates with sublinear memory cost. In
483 *International Conference on Machine Learning*, pages 4596–4604. PMLR, 2018.

484 N. Shinn, F. Cassano, B. Labash, A. Gopinath, K. Narasimhan, and S. Yao. Reflexion: Language
485 agents with verbal reinforcement learning.(2023). *arXiv preprint cs.AI/2303.11366*, 2023.

486 D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser,
487 I. Antonoglou, V. Panneershelvam, M. Lanctot, et al. Mastering the game of go with deep neural
488 networks and tree search. *nature*, 529(7587):484–489, 2016.

489 A. Slivkins et al. Introduction to multi-armed bandits. *Foundations and Trends® in Machine Learning*,
490 12(1-2):1–286, 2019.

491 S. Still and D. Precup. An information-theoretic approach to curiosity-driven reinforcement learning.
492 *Theory in Biosciences*, 131(3):139–148, 2012.

493 Y. Tassa, Y. Doron, A. Muldal, T. Erez, Y. Li, D. d. L. Casas, D. Budden, A. Abdolmaleki, J. Merel,
494 A. Lefrancq, et al. Deepmind control suite. *arXiv preprint arXiv:1801.00690*, 2018.

495 A. A. Team, J. Bauer, K. Baumli, S. Baveja, F. M. P. Behbahani, A. Bhoopchand, N. Bradley-Schmieg,
496 M. Chang, N. Clay, A. Collister, V. Dasagi, L. Gonzalez, K. Gregor, E. Hughes, S. Kashem,
497 M. Loks-Thompson, H. Openshaw, J. Parker-Holder, S. Pathak, N. P. Nieves, N. Rakicevic,
498 T. Rocktäschel, Y. Schroecker, J. Sygnowski, K. Tuyls, S. York, A. Zacherl, and L. M. Zhang.
499 Human-timescale adaptation in an open-ended task space. In *International Conference on Machine*
500 *Learning*, 2023a.

501 G. Team, R. Anil, S. Borgeaud, J.-B. Alayrac, J. Yu, R. Soricut, J. Schalkwyk, A. M. Dai, A. Hauth,
502 K. Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint*
503 *arXiv:2312.11805*, 2023b.

504 G. Team, P. Georgiev, V. I. Lei, R. Burnell, L. Bai, A. Gulati, G. Tanzer, D. Vincent, Z. Pan, S. Wang,
505 et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv*
506 *preprint arXiv:2403.05530*, 2024a.

507 G. Team, T. Mesnard, C. Hardin, R. Dadashi, S. Bhupatiraju, S. Pathak, L. Sifre, M. Rivière, M. S.
508 Kale, J. Love, et al. Gemma: Open models based on gemini research and technology. *arXiv*
509 *preprint arXiv:2403.08295*, 2024b.

510 G. Team, M. Riviere, S. Pathak, P. G. Sessa, C. Hardin, S. Bhupatiraju, L. Hussenot, T. Mesnard,
511 B. Shahriari, A. Ramé, et al. Gemma 2: Improving open language models at a practical size. *arXiv*
512 *preprint arXiv:2408.00118*, 2024c.

513 L. Team, A. Modi, A. S. Veerubhotla, A. Rysbek, A. Huber, B. Wiltshire, B. Veprek, D. Gillick,
514 D. Kasenberg, D. Ahmed, et al. Learnlm: Improving gemini for learning. *arXiv preprint*
515 *arXiv:2412.16429*, 2024d.

516 D. Tirumala, M. Wulfmeier, B. Moran, S. Huang, J. Humplik, G. Lever, T. Haarnoja, L. Hasenclever,
517 A. Byravan, N. Batchelor, N. sreendra, K. Patel, M. Gwira, F. Nori, M. Riedmiller, and N. Heess.
518 Learning robot soccer from egocentric vision with deep reinforcement learning. In P. Agrawal,
519 O. Kroemer, and W. Burgard, editors, *Proceedings of The 8th Conference on Robot Learning*,
520 volume 270 of *Proceedings of Machine Learning Research*, pages 165–184. PMLR, 06–09 Nov
521 2025. URL <https://proceedings.mlr.press/v270/tirumala25a.html>.

522 O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi, R. Powell,
523 T. Ewalds, P. Georgiev, et al. Grandmaster level in starcraft ii using multi-agent reinforcement
524 learning. *nature*, 575(7782):350–354, 2019.

525 J. X. Wang, Z. Kurth-Nelson, D. Tirumala, H. Soyer, J. Z. Leibo, R. Munos, C. Blundell, D. Kumaran,
526 and M. Botvinick. Learning to reinforcement learn. *arXiv preprint arXiv:1611.05763*, 2016.

527 X. Wang, J. Wei, D. Schuurmans, Q. Le, E. Chi, S. Narang, A. Chowdhery, and D. Zhou.
528 Self-consistency improves chain of thought reasoning in language models. *arXiv preprint*
529 *arXiv:2203.11171*, 2022.

530 Y. Wang, X. Yue, and W. Chen. Critique fine-tuning: Learning to critique is more effective than
531 learning to imitate. *arXiv preprint arXiv:2501.17703*, 2025.

532 J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, et al. Chain-of-thought
533 prompting elicits reasoning in large language models. *Advances in neural information processing*
534 *systems*, 35:24824–24837, 2022.

535 S. Welleck, X. Lu, P. West, F. Brahman, T. Shen, D. Khashabi, and Y. Choi. Generating sequences by
536 learning to self-correct. *arXiv preprint arXiv:2211.00053*, 2022.

537 C. Wen, J. Lin, T. Darrell, D. Jayaraman, and Y. Gao. Fighting copycat agents in behavioral cloning
538 from observation histories. *Advances in Neural Information Processing Systems*, 33:2564–2575,
539 2020.

540 R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement
541 learning. *Machine learning*, 8:229–256, 1992.

542 M. Wulfmeier, M. Bloesch, N. Vieillard, A. Ahuja, J. Bornschein, S. Huang, A. Sokolov, M. Barnes,
543 G. Desjardins, A. Bewley, S. M. E. Bechtle, J. T. Springenberg, N. Momchev, O. Bachem,
544 M. Geist, and M. Riedmiller. Imitating language via scalable inverse reinforcement learning. In
545 A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors,
546 *Advances in Neural Information Processing Systems*, volume 37, pages 90714–90735. Curran
547 Associates, Inc., 2024. URL [https://proceedings.neurips.cc/paper_files/paper/2024/](https://proceedings.neurips.cc/paper_files/paper/2024/file/a5036c166e44b731f214f41813364d01-Paper-Conference.pdf)
548 [file/a5036c166e44b731f214f41813364d01-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/a5036c166e44b731f214f41813364d01-Paper-Conference.pdf).

549 S. Yao, J. Zhao, D. Yu, N. Du, I. Shafran, K. Narasimhan, and Y. Cao. React: Synergizing reasoning
550 and acting in language models. *arXiv preprint arXiv:2210.03629*, 2022.

551 E. Zelikman, Y. Wu, J. Mu, and N. Goodman. Star: Bootstrapping reasoning with reasoning.
552 *Advances in Neural Information Processing Systems*, 35:15476–15488, 2022.

553 E. Zelikman, G. Harik, Y. Shao, V. Jayasiri, N. Haber, and N. D. Goodman. Quiet-star: Language
554 models can teach themselves to think before speaking. *arXiv preprint arXiv:2403.09629*, 2024.

555 S. Zhai, H. Bai, Z. Lin, J. Pan, P. Tong, Y. Zhou, A. Suhr, S. Xie, Y. LeCun, Y. Ma, et al. Fine-tuning
556 large vision-language models as decision-making agents via reinforcement learning. *Advances in*
557 *Neural Information Processing Systems*, 37:110935–110971, 2025.

- 558 T. Zhang, A. Madaan, L. Gao, S. Zheng, S. Mishra, Y. Yang, N. Tandon, and U. Alon. In-context
559 principle learning from mistakes. *arXiv preprint arXiv:2402.05403*, 2024.
- 560 S. Zhou, F. F. Xu, H. Zhu, X. Zhou, R. Lo, A. Sridhar, X. Cheng, T. Ou, Y. Bisk, D. Fried,
561 et al. Webarena: A realistic web environment for building autonomous agents. *arXiv preprint*
562 *arXiv:2307.13854*, 2023.

563	Contents	
564	A Related Work	15
565	B Environments & Datasets	16
566	B.1 Multi-arm Bandits: BanditBench	16
567	B.1.1 Baselines.	16
568	B.1.2 SFT Datasets.	16
569	B.2 Contextual Bandits	17
570	B.3 Tic-tac-toe	18
571	C Experimental & Implementation Details	18
572	C.1 Training & Evaluation	18
573	C.2 RLFT	19
574	C.3 SFT	20
575	C.4 Exploration Mechanisms	20
576	D Additional Results	22
577	D.1 Failure Modes	22
578	D.1.1 Greediness	22
579	D.1.2 Frequency Bias	23
580	D.1.3 Knowing-Doing Gap	24
581	D.2 Multi-armed Bandits	25
582	D.3 Contextual Bandits	25
583	D.4 Effect of Exploration Mechanisms	25
584	E Ablations	25
585	E.1 RLFT in Tic-tac-toe.	26
586	E.2 Tic-tac-toe: Effect of Legal Actions in State	26
587	E.3 Importance of CoT for RLFT	26
588	E.4 Expert Behavior Cloning vs. Thought Cloning	27
589	E.5 “Thinking” Time	27
590	F Limitations & Future Work	28

591 A Related Work

592 **Exploration in RL and LLMs.** The trade-off between *exploration* and *exploitation* is a long-standing
593 challenge in the field of RL [Schmidhuber, 1991a,b, Still and Precup, 2012, Oudeyer et al., 2007].
594 Widely used RL agents have often relied on random schemes [Mnih et al., 2015], heuristics such
595 as state-visitation counts [Ecoffet et al., 2019, Raileanu and Rocktäschel, 2020], intrinsic curiosity
596 [Pathak et al., 2017, Burda et al., 2018, Groth et al., 2021], behavior priors [Rao et al., 2021], or
597 maximum entropy regularization [Schulman et al., 2017, Haarnoja et al., 2018]. Naturally, a number
598 of works looked into leveraging LLMs for improving exploration of RL agents either as a source of
599 rewards [Klissarov et al., 2023, Lu et al., 2024] or to orchestrate exploration strategies [Klissarov
600 et al., 2024]. Krishnamurthy et al. [2024] investigate the in-context exploration abilities of LLMs
601 when acting directly as a policy. Similarly, Nie et al. [2024] study the exploration abilities of LLMs
602 when fine-tuned on expert trajectories. In contrast, our work investigates the effects of RLFT on the
603 exploration abilities of LLMs and focuses on *why* models fail.

604 **In-context Learning for Decision-Making.** ICL is a form of Meta-learning, also referred to as
605 learning-to-learn [Schmidhuber, 1987]. While meta-learning is targeted via a meta-training phase
606 [Santoro et al., 2016, Mishra et al., 2018, Finn et al., 2017, Wang et al., 2016, Duan et al., 2016,

Kirsch et al., 2019, Flennerhag et al., 2019, Team et al., 2023a], ICL emerges as a result of the pre-training data distribution [Chan et al., 2022, Kirsch et al., 2022]. ICL has been rediscovered in LLMs [Brown et al., 2020a] after initial observations by Hochreiter et al. [2001] in LSTMs [Hochreiter and Schmidhuber, 1997]. Mirchandani et al. [2023] leverage the ICL abilities of LLMs to operate as general pattern machines. A number of works leverage the CoT abilities [Wei et al., 2022] of LLMs in simple text-based scenarios [Shinn et al., 2023, Yao et al., 2022]. Similar in-context abilities have been observed in decision-making with models trained from scratch, albeit in restricted environments [Laskin et al., 2022, Lee et al., 2022, Kirsch et al., 2023, Raparthy et al., 2023, Schmied et al., 2024b,a].

Self-Correction in LLMs. A critical component for LLM agents is the ability to self-correct over previously explored attempts. Existing works focus primarily on math benchmarks [Cobbe et al., 2021, Hendrycks et al., 2021, Welleck et al., 2022]. Zelikman et al. [2022] leverage hints to iteratively generate correct answers and fine-tune on the respective CoT rationales. Kumar et al. [2024] employ RLFT over multiple trials to induce self-correction. Similarly, Zelikman et al. [2024] make use of RL fine-tuning, but instead generate rationales at every token position. Instead of imitation, Wang et al. [2025] rely on critique fine-tuning to induce self-correction. Wulfmeier et al. [2024] make use of inverse RL to avoid compounding errors. Other works rely on ICL abilities to learn from previous mistakes [Zhang et al., 2024, Monea et al., 2024]. While conceptual corrections are possible, exact token-level correction is usually difficult for autoregressive generation [Cundy and Ermon, 2023].

B Environments & Datasets

We conduct experiments on three sets of environments: multi-armed bandits, contextual bandits and tic-tac-toe. For the SFT experiments reported in Section 4.5, we generate our own expert datasets. In this section, we provide additional details on our environments and datasets.

B.1 Multi-arm Bandits: BanditBench

MABs [Slivkins et al., 2019, Lattimore and Szepesvári, 2020] are a classic problem setting in RL that isolates the *exploration-exploitation* trade-off. In contrast, commonly used RL environments [Bellemare et al., 2013, Tassa et al., 2018] often conflate exploration with other RL-specific aspects, such as delayed rewards [Arjona-Medina et al., 2019]. We rely on the MAB scenarios released in BanditBench [Nie et al., 2024] and also used by [Krishnamurthy et al., 2024]. MABs come with a number of variable dimensions including the scenario type (textual description of the task), the type of reward distribution (Gaussian, Bernoulli) and its corresponding noise level (low/medium/high), the number of arms (i.e., actions), and the number of interaction steps per episode. Consequently, MABs are a good testbed for LLM agents.

We focus on the *continuous* and *button* variants released by Nie et al. [2024]. We report results for MAB with $k \in \{5, 10, 20\}$ arms ($|\mathcal{A}| = k$) for three levels of stochasticity (low/medium/high). In our experiments, for every arm the corresponding reward is sampled from a Gaussian distribution $r \sim \mathcal{N}(\mu, \sigma)$ where $\mu \sim \mathcal{U}(0, 1)$ and is a fixed scalar $\sigma \in \{0.1, 1, 3\}$ for the three levels of stochasticity, respectively. For all MAB settings, we limit the horizon T to 50 interaction steps. Limiting the horizon is necessary to handle the increasing lengths and consequently RAM requirements for fine-tuning. While we consider 50 interaction steps sufficient for 5 and 10 arms, it is insufficient for the 20 arms scenario. However, note that the general trends are well observable for the 20 arms scenario. In Figure 7, we show the continuous and button Gaussian MABs with CoT instructions for the agent. Similarly, in Figure 8 we show the same instances without CoT instructions.

B.1.1 Baselines.

We compare against two commonly used baselines for MABs: Upper-confidence Bound (UCB) [Auer, 2002] and a random agent that selects actions uniformly at random (see Appendix C for details). UCB is considered optimal and represents the upper-bound for agent performance, whereas the random baseline represents the lower bound. We provide implementations details for all baselines in Appendix C.

B.1.2 SFT Datasets.

In our main experiments, we focused on self-generated CoT rationales and action predictions produced by our fine-tuned agents, which do not require a pre-collected dataset. In contrast, for our SFT experiments reported in Section 4.5, we generated UCB expert datasets. In particular, we construct two dataset instances: a *behavior cloning* dataset that only contains expert actions and a

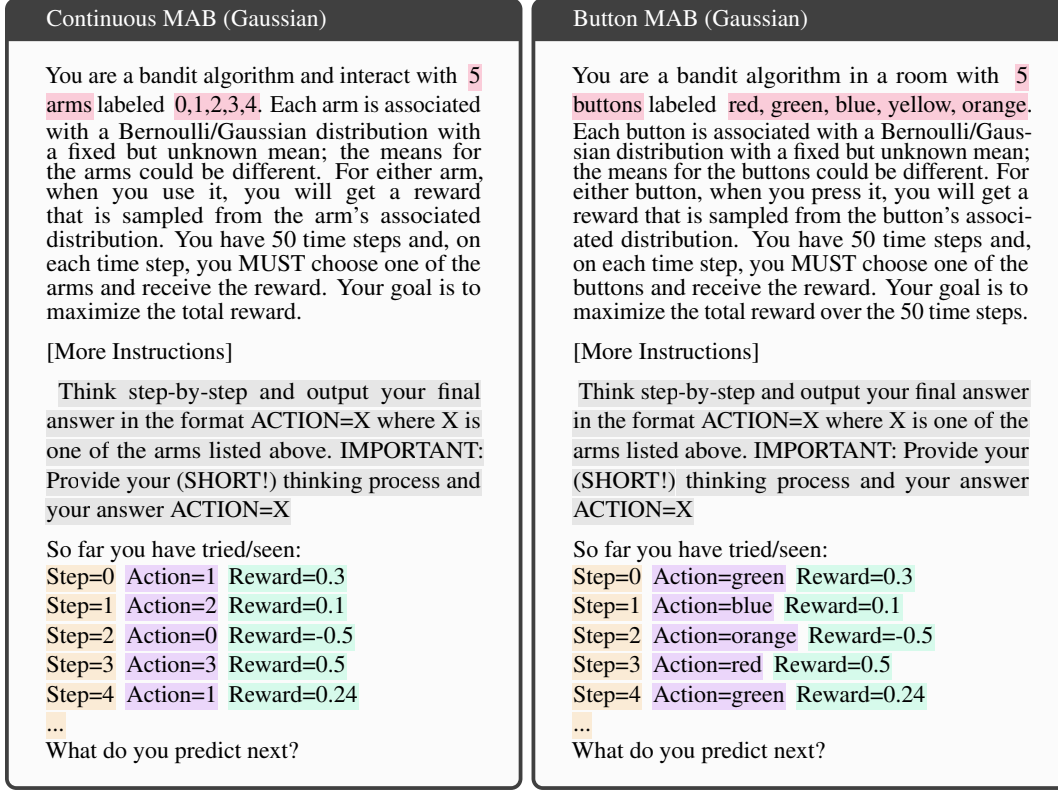


Figure 7: Illustration of *continuous* and *button* Gaussian multi-armed bandits scenarios from Bandit-Bench [Nie et al., 2024] using our context representation and **with CoT** instructions.

661 *thought cloning* (TC) dataset that incorporates expert actions alongside a thought process for coming
662 up with the respective action (i.e., w/ CoT). For every arm and stochasticity level combination,
663 we construct 32K rollouts from different MABs, which amounts to 1.6M transitions (state, action,
664 reward) per dataset.

665 To provide the thought process in the TC datasets, we reconstruct the computations of the UCB values
666 conducted by the UCB expert in textual format. The thought process ends with a final conclusion
667 why a particular action was selected (i.e., highest UCB value or exploratory action). Consequently,
668 the action z_t at step t contains the thought process z_{CoT} and the action to execute a_t (see Section
669 3.2). We illustrate the actions contained in the dataset for a trajectory at steps 4 and 11 in Figure 9.
670 The BC datasets do not contain the thought process. Instead, they only contain the final predictions
671 made by the model, for example ACTION=yellow as shown in Figure 9.

672 B.2 Contextual Bandits

673 MABs do not emit states. In contrast, CBs emit state representations at every interaction step, making
674 them contextual. Consequently, CBs are interesting to test abilities of LLMs to make use of the
675 given context when predicting the next action. For our CB experiments, we leverage the MovieLens
676 environment released by Nie et al. [2024], a semisynthetic bandit task based on the MovieLens
677 dataset [Harper and Konstan, 2015]. In this setting, the agent operates as a movie recommendation
678 engine given a contextual description of a user (10K users in total) and a list of K possible movies.
679 The context representation provides a textual description of the user to recommend the movie to.
680 This description includes the user's gender, age, profession, location, and a numeric description
681 of the user's preferences for each of the possible movies. As for MABs, we report results for
682 $K \in \{5, 10, 20\}$, limit the horizon to 50 interaction steps. In Figure 10, we provide an example for a
683 MovieLens CB with 5 actions with our context representation and CoT instructions.

684 **Baselines.** Similar to MABs, we compare against LinUCB [Chu et al., 2011] and an agent selecting
685 actions uniformly at random. We provide implementation details on our baselines in Appendix C.

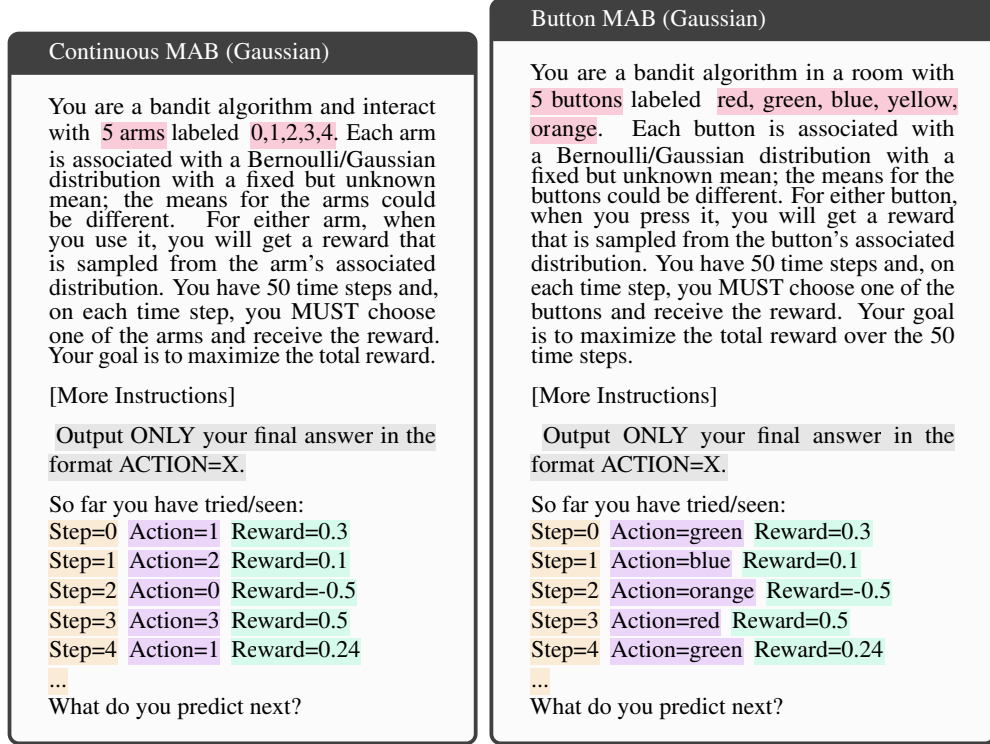


Figure 8: Illustration of *continuous* and *button* Gaussian multi-armed bandits scenarios from Bandit-Bench [Nie et al., 2024] using our context representation **without** CoT instructions.

686 B.3 Tic-tac-toe

687 Finally, we use the text-based Tic-tac-toe environment released by Ruoss et al. [2024] (see Figure
 688 11 for an example). Unlike MABs and CBs, Tic-tac-toe is a stateful environment with proper
 689 state transitions (i.e., action predicted at step t affects the state observed at step $t + 1$). The agent
 690 receives scalar rewards of 1, 0, and -1 for winning, drawing, and losing against its opponent,
 691 respectively. Episodes last until either of the players wins, draws or loses. To enable easy extraction
 692 of actions from the generated rationales, we represent the action space as a discrete set of 9 actions,
 693 corresponding to the grid positions on the 3×3 grid used in Tic-tac-toe ($|\mathcal{A}| = 9$). However, only at
 694 the start of an episode, all 9 actions are valid. Subsequently, only a subset is valid, because of taken
 695 board positions. We (optionally) provide the set of valid actions at a particular step in textual form
 696 in the context given to the agent. Ruoss et al. [2024] demonstrated that frontier models struggle to
 697 achieve strong performance in this environment and barely beat a random opponent. Consequently,
 698 we deem it a good target to investigate the efficacy of RLFT.

699 **Baselines.** Following Ruoss et al. [2024], we compare against a random agent by default. In addition,
 700 we also compare against (MCTS) [Coulom, 2006], and a noisy variant of MCTS that selects an action
 701 randomly with 50% chance and according to MCTS otherwise.

702 C Experimental & Implementation Details

703 C.1 Training & Evaluation

704 In our experiments, we fine-tune Gemma2 models in three model sizes (2B/9B/27B). For all exper-
 705 iments, we use the instruction-tuned versions of Gemma2 and leverage the respective instruction
 706 pre-and-postfixes. For bandits, we fine-tune all models for a total of 30K updates and evaluate after
 707 every 10K steps. with an accumulated batch size of 128. Similarly, we fine-tune for 12K updates and
 708 evaluate every 4K updates on Tic-tac-toe. We report the mean and 95% confidence intervals over
 709 three seeds, as suggested by Agarwal et al. [2021].

710 **General.** We train all agents with an accumulated batch size of 128. We use a learning rate of
 711 $1e^{-4}$, 100 linear warm-up steps followed by a cosine decay to $1e^{-6}$. To allow for memory-efficient

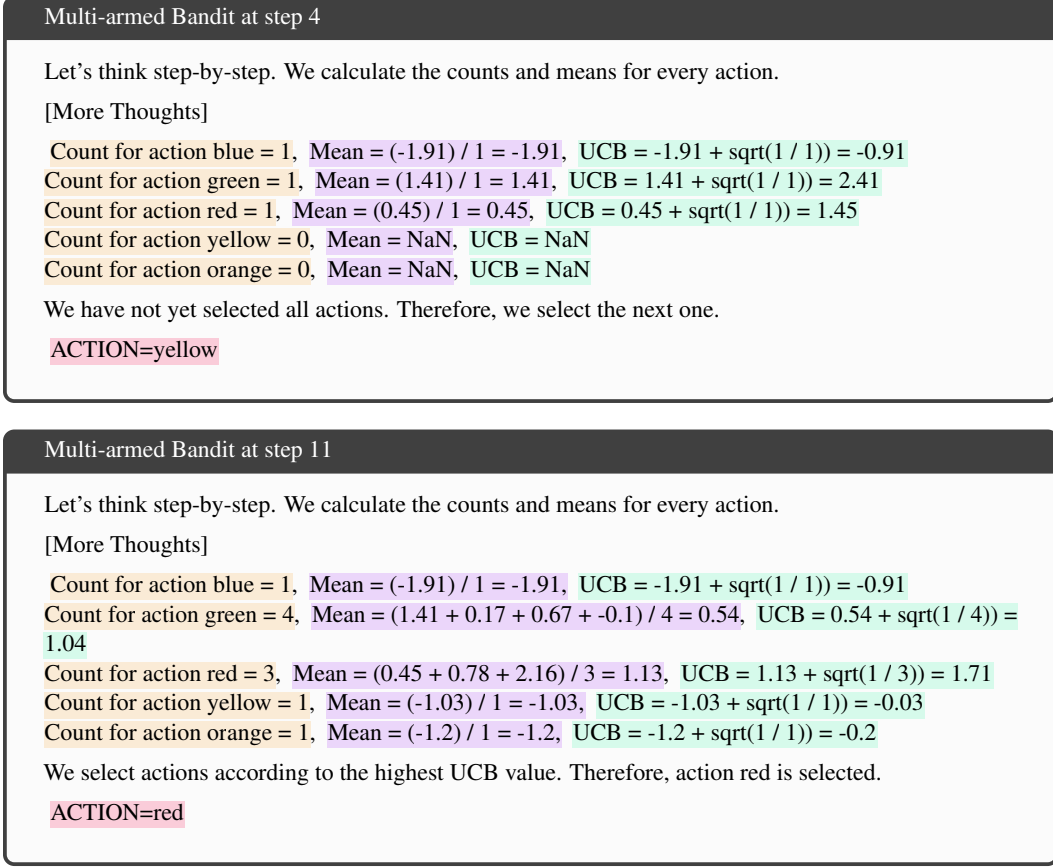


Figure 9: Illustration of UCB rationales contained in our **SFT expert datasets** at two timesteps (4 and 11) in the same trajectory. Both examples show the Thought Cloning dataset instance containing both the produced CoT rationale along with the predicted action. The Behavior Cloning instances contain only the final action prediction (in red).

fine-tuning of 2B and 9B models, we train using the AdaFactor optimizer [Shazeer and Stern, 2018]. We experiment with LoRA [Hu et al., 2022] for fine-tuning the 9B and 27B models but found it insufficient for improving the agent’s decision-making abilities in our setting. However, LoRA considerably reduces the amount of memory required for RLFT and has been shown to work well for supervised fine-tuning of decision-making agents [Schmied et al., 2023]. Therefore, we deem it a promising candidate for RLFT in decision-making scenarios. Furthermore, we employ gradient clipping of 1.0. We list all hyperparameters in Table 1.

Context Lengths & Generation Budget. For all model sizes and tasks, we use a context length of 1792 for the input context. By default, we set the generation budget to 256 tokens, except for the knowing-doing gap analyses reported in Section 4.2, which require a larger budget of 2048 tokens. Consequently, the effective sequence length for fine-tuning is 2048.

Hardware Setup. We train all models on a server equipped with $8 \times$ H100 GPUs.

C.2 RLFT

For our RLFT experiments on bandits, we employ the context representation, action factorization, reward shaping terms, and training objectives described in Section 3.2. To extract the target action a_t from z_t , we make use of a stack of regex expressions against the target pattern (i.e., ACTION=X) and consider the last match in the generated tokens as a_t . In addition to being fairly robust, we found that this approach allows for more flexibility during the RLFT process and led to better outcomes than a more structured approach. Furthermore, across model sizes, we found it essential to introduce a reward shaping term to penalize rationales that contain no valid actions. By default, we use a reward penalty of -5 for invalid actions. Empirically, we found that this reward shaping term is sufficient for the models to produce valid actions early on in the training.

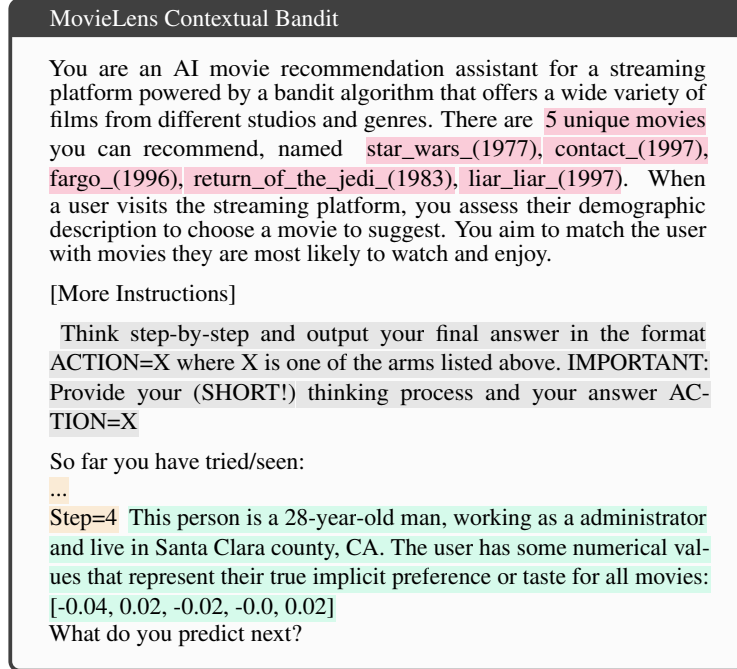


Figure 10: Illustration of **contextual MovieLens scenario** from BanditBench [Nie et al., 2024] using our context representation and instructions.

We fine-tune using the clipping objective introduced by Schulman et al. [2017] with an additional KL constraint to the reference policy π_{ref} . We set $\beta = 0.05$ and $\epsilon = 0.2$ for all experiments. We make use of the approximated (per-token) KL divergence instead of computing the full KL. While we found that computing the full KL slightly improves performance, it slows down training considerably. In contrast to Ahmadian et al. [2024] and Ramesh et al. [2024], we do not rely on producing multiple rollouts, because it is impractical for the multi-step nature of decision-making tasks. While generating multiple actions at a particular timestep is possible for simulated environments, it requires environment resets. Therefore, we rely on standard MC-baselines to estimate A_{adv} .

For bandit experiments, we maintain a pool of 512 stochastic MABs. For every rollout, we let the agent interact with a subset of 16 bandits for a single episode (50 timesteps). Consequently, every rollout contains 800 transitions. Similarly, for Tic-tac-toe, we maintain 16 parallel environments and collect 2048 rollout steps. We conduct 1 and 2 update epochs over the rollout buffer for bandits and Tic-tac-toe, respectively.

C.3 SFT

For our SFT experiments on MABs, we fine-tune on either on the expert action or expert rationales produced by UCB. We employ standard SFT training using a cross-entropy objective on the target tokens.

C.4 Exploration Mechanisms

In Section 18, we compare a variety of classic exploration mechanisms and LLM-specific approaches and study their effects on agent performance on Gaussian MABs with 10 arms. Here, we provide a description for each mechanism.

Try-all. The try-all strategy is inspired by UCB, which incorporates an initial phase for trying all untried actions. This is because the UCB values for all untried actions are ∞ . Therefore, we incorporate the same exploration phase when performing ICL and RLFT at the beginning of every episode. To enable fine-tuning on exploration actions, we provide an action rationale template to the model (e.g., Action X has not been tried yet, let’s explore it. ACTION=X). While simple, we find that this try-all strategy is effective for lowering regret across all model sizes (see Figure 18). This suggests that the model is able to select appropriate actions if given sufficient information, but struggles to explore.

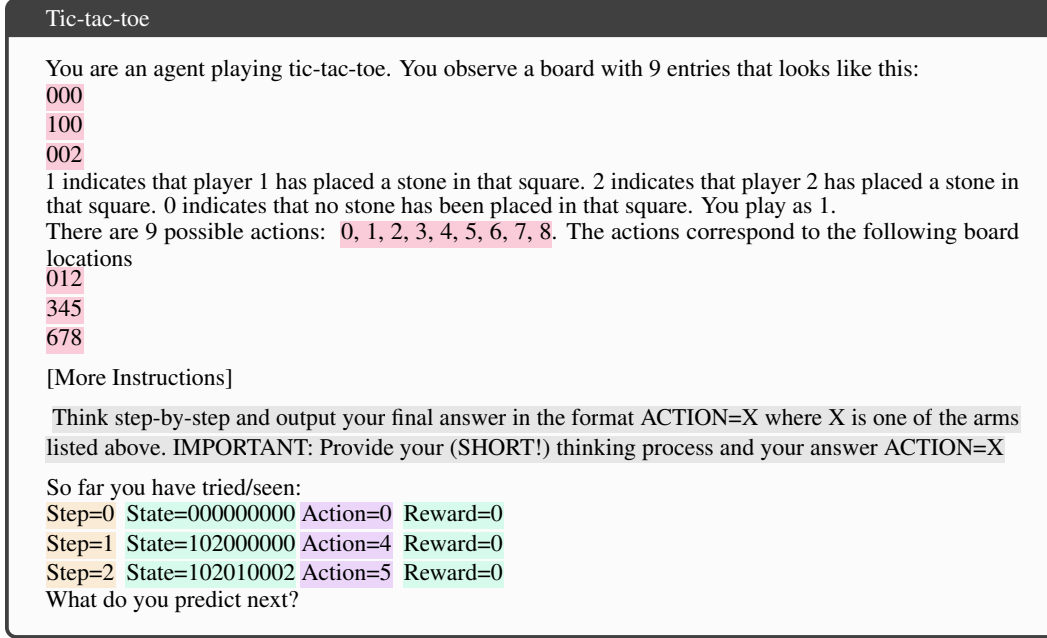


Figure 11: Illustration of the text-based **Tic-tac-toe** environment.

763 **ϵ -greedy.** ϵ -greedy is classic exploration mechanism and commonly used in RL algorithms [Mnih
764 et al., 2015, Hessel et al., 2018]. For our experiments, we use $\epsilon = 0.1$ both during training and
765 evaluation. We explored other values for ϵ but did not observe performance improvements. As for the
766 try-all strategy, we provide an action rationale template to enable fine-tuning on exploration actions.

767 **Context Randomization.** Context Randomization is an LLM-specific mechanism designed to
768 introduce randomness in the action predictions by modifying the context representation. At every
769 interaction step, we construct a mapping from the original action labels to a shuffled list of the same
770 action labels. Subsequently, we remap action in the context history according to the constructed
771 mapping. Finally, the predicted action is mapped back to the original action label space and executed
772 environment. Besides introducing randomness, context randomization acts as a control mechanism to
773 ensure that the observed biases do not only stem from biases towards particular action-tokens (e.g.,
774 blue occurs more often than magenta in the pre-training dataset).

775 **Context Summary.** Similar to Krishnamurthy et al. [2024] and Nie et al. [2024], we evaluate the
776 effects of providing a context summary to the agent. After the context history, we provide the model
777 with a summary of that history that contains the number of times every action has been selected so
778 far, along with their respective mean rewards.

779 **Self-Correction.** Inspired by Kumar et al. [2024] and Wang et al. [2025], we employ self-correction to
780 the model’s predicted actions. First, we let the model generate its initial rationale and corresponding
781 action prediction. Then we append the generated rationale along with a self-correction message
782 (similar to Kumar et al. [2024]) to the input context, and repeat the action generation. Finally, we
783 extract the action from the final response and execute it in the environment. For RLFT, we only
784 fine-tune on the final response, but retain the initial response along with the self-correction message
785 in the context.

786 **Self-Consistency.** Instead of generating a single answer, self-consistency [Wang et al., 2022] relies on
787 generating multiple responses. Subsequently, self-consistency employs a majority voting mechanism
788 to determine the final response. For our experiments in Figure 18, we report results for self-consistency
789 with 16 generated responses. Instead of majority voting, we experimented with sampling from the
790 respective response distribution or random mechanisms.

791 **Exploration Bonus.** Finally, we evaluate a reward shaping mechanism in the form of an exploration
792 bonus. In particular, we give an exploration bonus of +1 if the agents selects an action not yet tried
793 within the respective episode. While simple, we find that the exploration bonus effectively narrows
794 the gap to the UCB expert. This highlights the importance of reward shaping for fine-tuning LLMs in
795 decision-making scenarios.

Name	Value	Description
Training		
training_steps	30K or 12K	Number of training steps.
eval_freq	10K or 4K	Evaluation frequency (in updates).
batch_size	128	Accumulated batch size.
lr_scheduler	Linear + cosine	Learning rate scheduler
warmup_steps	100	Warmup steps.
lr	$1e^4$ to $1e^6$	Maximum learning rate.
optimizer	AdaFactor	Optimizer.
Sequence Length & Generation Budget		
context_length	1792	Input context length.
num_tokens	256	Generation budget.
RLFT		
rollout_steps	800 or 2048	Rollout steps in-between updates.
update_epochs	1 or 2	Update epochs over rollout-buffer.
reward_penalty	-5	Reward penalty for invalid actions.
loss	PPO clipping objective + KL constraint	Objective function.
baseline	MC-baseline or state-value head	Baseline.
envs	16	Number of parallel envs.
ϵ	0.2	Clipping value.
β	0.05	KL coefficient.
reward_norm	True	Whether reward normalization is used.
train_temp	1.0	Sampling temp during rollouts.
eval_temp	0.0	Sampling temp during evaluation.
top_p	1.0	Sampling top-p.
Hardware		
accelerator	$8 \times \text{H100}$	Hardware accelerator.

Table 1: Default **hyperparameters** used in our experiments.

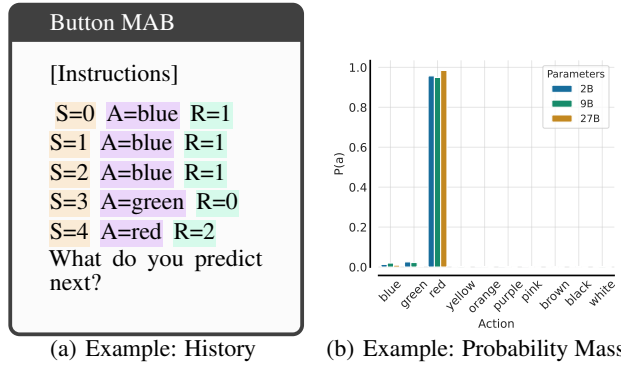


Figure 12: Illustration of **action probabilities** leading to greediness behavior. Models exhibit overly high action probabilities in the presence of rewards, potentially resulting in repeatedly selecting sub-optimal actions

D Additional Results

D.1 Failure Modes

D.1.1 Greediness

Greediness is characterized by the LLM overly favoring the best performing action among a small set of actions seen so far. We define action coverage C_t at step t as the fraction of available actions that have been selected at least once, $C_t = \frac{|\{a \in \mathcal{A}: N_t(a) > 0\}|}{|\mathcal{A}|}$ with $N_t(a)$ representing the number of times action a has been selected until t .

Action probabilities. The suboptimal action coverage reported in Section 4.2 is caused by the model overly favoring high-reward actions (i.e., overly high action probabilities). In Figure D.1.1, we provide an illustration of the action probabilities for a given input history. Across model sizes, Gemma2

806 exhibits overly high action probabilities in the presence of reward, which results in repeatedly
 807 selecting a potentially suboptimal action.

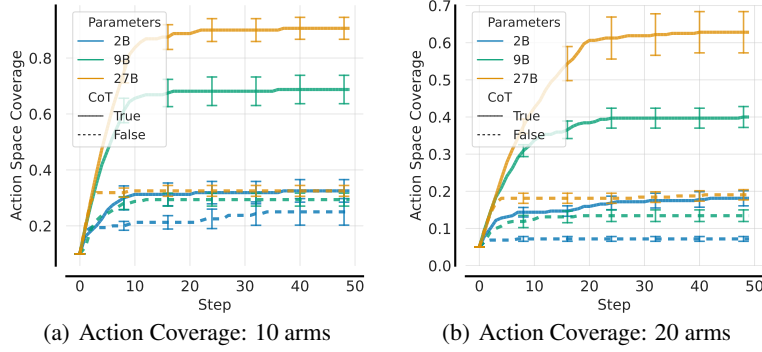


Figure 13: Illustration of greediness for the **numbers** scenario.

807
 808 **Greediness on Continuous MABs.** We repeat the analyses conducted in Section 4.2 using *numbers*
 809 instead of *buttons* as the possible actions. Indeed, we find that the same trends hold. Without CoT the
 810 performance remains low. For Gemma2 27B, we observe an increase in the action coverage to almost
 811 90% for the 10 arms scenario, and to 60% for the 20 arms scenario.

812 **Post RLFT.** In line with Figure 14a, we present the post RLFT action coverage on the 20 arms
 813 scenario in Figure 14b. Similar to the effects on the 10 arms scenario, we observe that RLFT improves
 the action coverage by 13%.

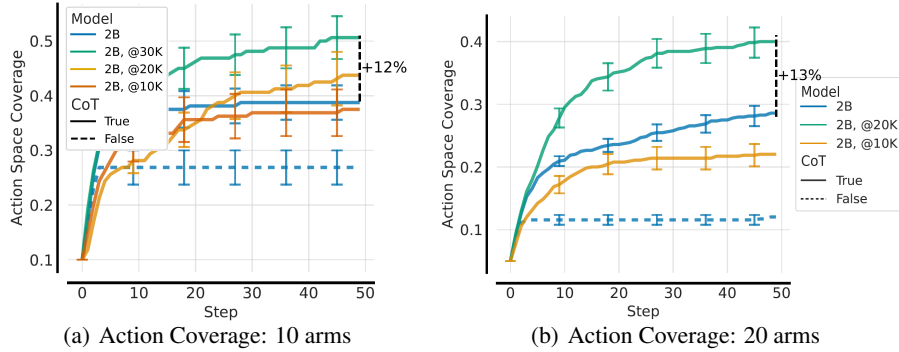


Figure 14: Effect of RLFT on **greediness** for 2B on 10 and 20 arms (medium noise).

814

815 D.1.2 Frequency Bias

816 Frequency bias is characterized by repeatedly selecting the most frequently occurring actions in
 817 the context, even when the dominant action gives low reward. To measure frequency bias, we first
 818 construct a variety of interaction histories (occurred during environment interaction) containing
 819 between 2 and 10 transitions. This interaction history is collected using a random policy. Given an
 820 initial interaction history, we repeat the last action in the history, which we also refer to as target
 821 action, between 0 and 100 times. Finally, we report the entropy all actions, $H(\theta) = -\sum_{a \in A} \pi_\theta(a | \tau) \log \pi_\theta(a | \tau)$. To achieve this, we conduct a separate forward pass for every possible action
 822 in the action space and report the respective log probabilities. We repeat the same procedure for
 823 different interaction histories and target actions (see Figure 4a and c). For the 10 arms scenario, every
 824 interaction history therefore results in 1000 (10 arms * 100 repetitions of the target action) forward
 825 passes. We repeat this procedure for the 5 target actions reported in Figure 4 using 5 interaction
 826 histories per action, accumulating to a total of 25K model forward passes (1000 * 5 * 5) per figure.
 827

828 To quantify frequency bias, we categorize the resulting actions as *frequent* action, *greedy*, or *other* if
 829 they are neither frequent nor greedy. Subsequently, we compute the frequent F_f , greedy F_g and other

830 F_o fractions as reported in Figure 4:

$$F_f = \frac{N_T(a_f)}{N}; \quad F_g = \frac{N_T(a_g)}{N}; \quad F_o = \frac{\sum_{a \in A \setminus \{a_f, a_g\}} N_T(a)}{N}, \quad \text{with } N = \sum_{a \in A} N_T(a). \quad (3)$$

831 Note that there can be an overlap between greedy and frequent actions. In these (rare) cases, the
 832 greedy action category is dominant, i.e., we categorize the action as greedy even if it would also
 833 be the frequent action. This implies that the actions classified as frequent in Figure 4, are always
 834 suboptimal/bad compared to the respective greedy action. Consequently, a high F_f indicates that the
 835 model prefers the most frequent action even when observing a better action in the context.

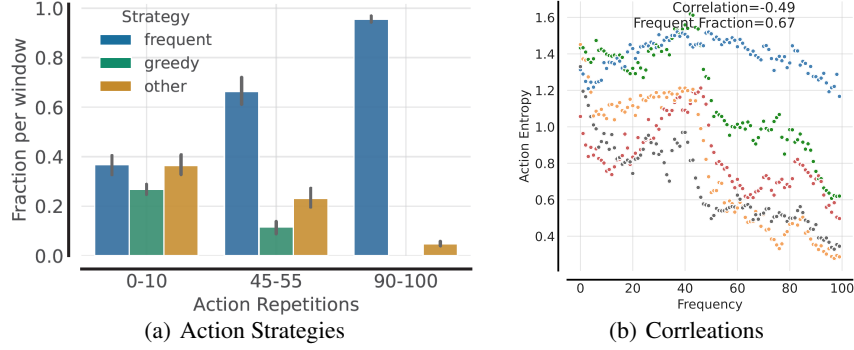


Figure 15: Effect of RLFT on frequency bias for 2B (10 arms, medium noise).

836 **Post RLFT.** In Section 4.3, we observed that RLFT counteracts frequency bias. In addition to
 837 frequency buckets reported in Figure 15a, we provide the plot frequency against action entropy post
 838 RLFT in Figure 15b. Compared to Figure 4a, we observe that after RLFT the models maintains a
 839 higher action entropy for longer. Only at high repetition frequencies the action entropy decreases
 840 severely. Consequently, RLFT counteracts frequency bias, but does not completely alleviate it.

841 D.1.3 Knowing-Doing Gap

842 The *knowing-doing gap* has been observed by Paglieri et al. [2024] and Ruoss et al. [2024]. It states
 843 that models possess knowledge about a task or consequences of their behavior (i.e., they *know* what
 844 to do), but cannot materialize this knowledge when acting (i.e., they are incapable of *doing*). We
 845 illustrate the knowing-doing gap empirically in Figure 5. To this end, we first task Gemma2 27B
 846 to produce the UCB algorithm and to compute the relevant quantities required to act according to
 847 UCB (“knowing”). This involves counting how often every action was selected, computing the mean
 848 rewards for every action, and computing the final UCB values. After producing the quantities, the
 849 model is tasked to act according to them (i.e., “doing”). In Figure 21, we present an example of the
 850 respective instructions given to the model along with a response produced by Gemma2 27B.

851 To evaluate performance empirically, we let Gemma2 27B interact with the environment (64 parallel
 852 instances) for 50 timesteps. We extend the token generation budget to 2048 tokens per step, to
 853 accommodate the additional required computations. Every produced action z contains both the CoT
 854 rationale z_{CoT} and the final selected action a . We first extract the computed UCB values from the
 855 produced rationale z_{CoT} . To achieve this, we task Gemma2 27B to enclose the computed values
 856 by `<ucb_values>` and `</ucb_values>` blocks. Then we extract the selected action a and execute it
 857 in the environment. For this experiment, we use Gemma2 27B, because we found that 2B and 9B
 858 struggled with computing the relevant UCB quantities and with enclosing them appropriately under
 859 the desired blocks.

860 **Quantifying “Knowing”.** To quantify “knowing”, we compare the UCB values computed by the
 861 model and extracted from z_{CoT} against the real UCB values. To this end, we recompute the real UCB
 862 values for every action at every time-step given the observed history. We consider the rationale as
 863 correct if the arm with the highest UCB values match. We opt for this choice rather than checking for
 864 exact equality, because we observed that the model struggles with exact calculations for complex
 865 operations. This is expected, because the necessary computations involve logarithm and square roots
 866 of floating point values. While tool use (e.g., calculator) could mitigate this issue, we observed that
 867 Gemma2 27B gets the quantities approximately right, resulting in valid rationales. Thus, the fraction
 868 of correct rationales is $F_c = \frac{1}{N} \sum_{i=1}^N g(z_{CoT}^i)$ given a classifier g .

869 **Quantifying “Doing”.** To quantify “doing”, we categorize the generated actions as *optimal* action
870 if the model selects the action with the highest UCB value, as *greedy* if it selects the action with
871 the highest UCB value among the set of actions tried so far, and as *other* if the action is neither
872 optimal nor greedy. It is possible that the greedy action is the optimal action. However, in this case
873 the action is considered optimal instead of greedy. Subsequently, we compute the percentages of
874 greedy/optimal/other actions (e.g., $F_g \times 100$). We find that the model clearly *knows* how to solve the
875 task, with 89% of all rationales being correct (see Figure 5).

876 D.2 Multi-armed Bandits

877 In Figure 6, we report the cumulative regrets across model sizes and arms for a medium noise
878 ($\sigma = 1.0$) scenario. In addition, we repeat the same experiment in the low-noise ($\sigma = 0.1$) and the
879 high-noise $\sigma = 3.0$ setting in Figure 16. For both noise levels, we observe similar trends as for the
880 medium noise setting. In particular, we observe that LLMs clearly outperform the random baseline
881 and RLFT lowers the cumulative regret for Gemma2 2B across all arm scenarios.

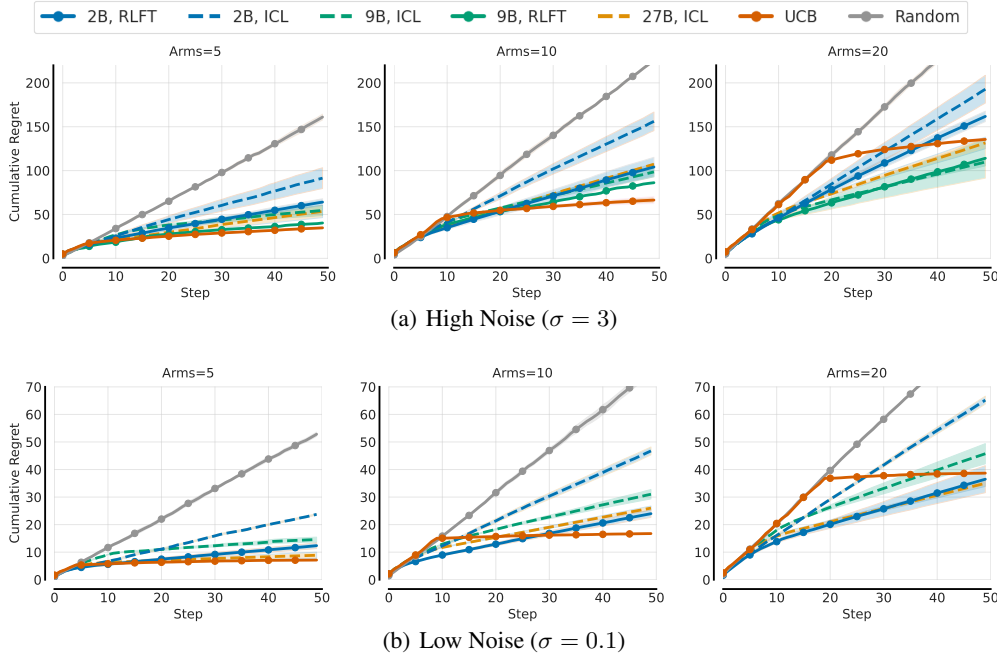


Figure 16: Main Comparison on **Gaussian MABs** button scenario in the (a) high $\sigma = 3$ and (a) low $\sigma = 0.1$ noise settings. We compare cumulative regrets (lower is better) of classic baselines against ICL and RLFT performances for Gemma2 2/9/27B for 5, 10, and 20 arms.

882 D.3 Contextual Bandits

883 We repeat the same fine-tuning experiment for the contextual MovieLens bandits described in Section
884 B.2. In Figure 17, we report the cumulative regrets attained by Gemma2B across different model
885 sizes and for 5, 10 and 20 arms. Furthermore, we compare against a LinearUCB and a Random
886 baseline. Overall, we observe similar performance improvements for RLFT on CBs as on MABs.
887 While the ICL performances barely attain the same performance as a Random agent, RLFT fine-tuned
888 Gemma2 2B perform similar to UCB.

889 D.4 Effect of Exploration Mechanisms

890 E Ablations

891 Finally, we provide additional details on the ablations conducted in this work.

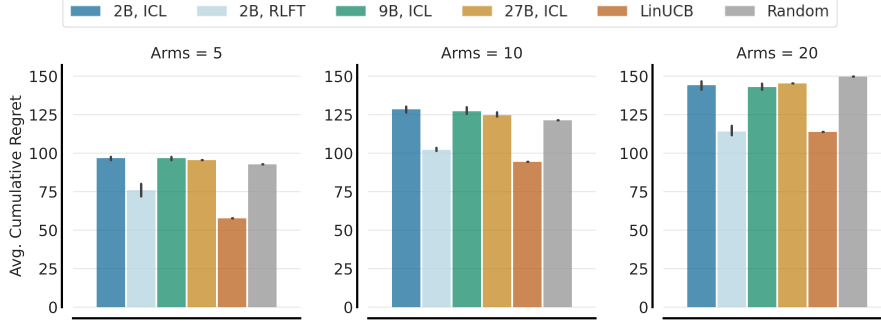


Figure 17: Main Comparison on Gaussian MovieLens CBs for (a) 5, (b) 10, and (c) 20 arms. We compare classic baselines against ICL and RLFT performances for Gemma2 2/9/27B.

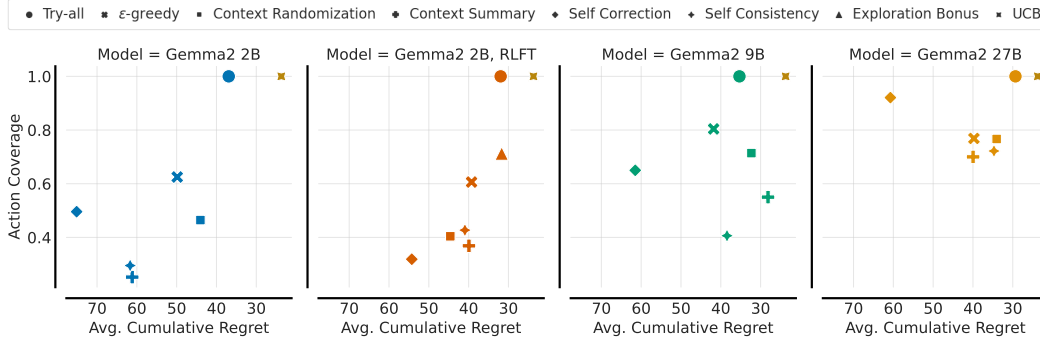


Figure 18: Effect of exploration mechanisms on action coverage and cumulative regret.

E.1 RLFT in Tic-tac-toe.

To investigate the efficacy of RLFT in stateful environments, we evaluate on Tic-tac-toe from Ruoss et al. [2024], in which frontier models struggle to achieve strong performance (see Appendix C for training details). We fine-tune against three opponents: a random agent, Monte Carlo Tree Search (MCTS) [Coulom, 2006], and noisy MCTS (50% of actions selected at random). We find that RLFT significantly enhances the win-rate of Gemma2 2B against all opponent compared to ICL (see Figure 19a). Against the random agent, RLFT elevates the average return from 0.15 (i.e., winning 15% of games) to 0.75. Notably, the agent even manages to draw against the optimal MCTS baseline ($-0.95 \rightarrow 0.0$), underscoring the effectiveness of RLFT for decision-making. However, for high performance it is essential to provide the legal actions in the context (see Figure 20).

E.2 Tic-tac-toe: Effect of Legal Actions in State

By default, we provided the legal actions available at the current turn within the input context to the agent. We found this design choice to be essential for effective fine-tuning compared to training without legal actions (see Figure 19b). Without legal actions in the context the average return drops from 0.75 (w/ legal actions) to 0.45. This suggests that the LLM fails at identifying the appropriate actions among the set of all possible actions when not given legal actions at the current state. In contrast, when provided with sufficient information (i.e., legal actions), the LLM is able to select actions appropriately (similar to Section 4.4). Providing the legal actions in the agent’s context alleviates the need to explore/identify invalid actions. Consequently, this shortcoming further highlights the need for principled exploration strategies for LLMs in decision-making scenarios.

E.3 Importance of CoT for RLFT

CoT reasoning is critical for ICL performance (see Figure 3), but the question remains how CoT influences RLFT. Therefore, we run RLFT on Gemma2 2B on the 10 arms Gaussian MAB both w/ and w/o CoT (see Figure 19b, RLFT). Indeed, without CoT, RLFT barely attains the performance of ICL w/ CoT. This highlights the function of CoT as a vital exploration and rationalization mechanism

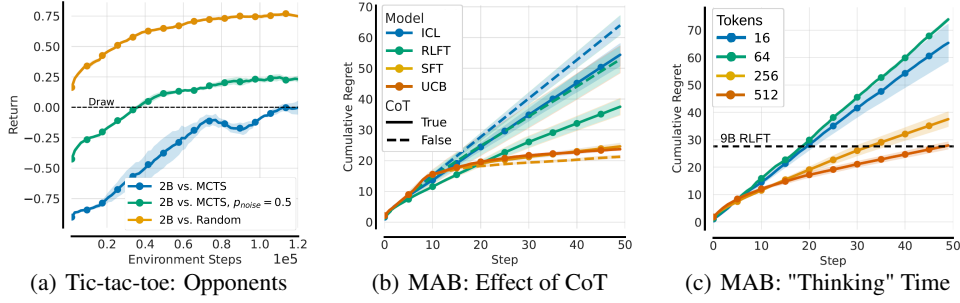


Figure 19: Ablations. (a) Effect of RLFT in Tic-tac-toe from Ruoss et al. [2024]. (b) Effect of CoT on ICL, RLFT and SFT (expert data) performance on MABs. (c) Effect of increasing the number of "thinking" tokens to generate during RLFT.

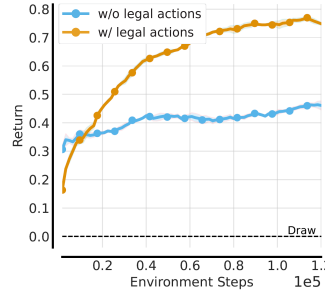


Figure 20: Effect of removing legal actions from the context in Tic-tac-toe.

for decision-making. For our results without CoT reported in Figure 19b, we remove the CoT instructions given to our agents. Instead, we instruct the agents to not perform any reasoning steps and to only produce the action to execute in the environment a . In addition, we limit the token generation budget G to 16 to avoid that the model ignores the instructions and makes use of the additional tokens. Furthermore, this considerably speeds up training due to faster rollout times and shorter context lengths.

E.4 Expert Behavior Cloning vs. Thought Cloning

A prevalent approach in sequence models for decision-making is behavior cloning (BC) [Pomerleau, 1988, Reed et al., 2022, Brohan et al., 2022, 2023], which relies on expert datasets. Consequently, we construct two UCB expert dataset comprising 32K rollouts either w/o CoT (*behavior cloning*) or w/ CoT (*thought cloning*), as described in Appendix B.1. Notably, both SFT variants successfully mimic the expert achieving comparable regret to the UCB expert (see Figure 19b, SFT). This result underscores the efficacy of training on expert data in decision-making scenarios when available, echoing recent findings in reasoning tasks [Muennighoff et al., 2025]. While BC and TC attain similar performance levels on the simplistic MABs, we anticipate that TC is advantageous in more complex decision-making scenarios as found by Hu and Clune [2023].

E.5 "Thinking" Time

We investigate the effect of giving the agent more time to "think" in Figure 19c. To achieve this, we vary the maximal number of tokens that the agent can generate per action $G \in \{16, 64, 256, 512\}$. By default, we set G to 256. Indeed, we observe that the performance improves consistently with more thinking tokens. Decreasing G to 16 or 64 results in poor performance, because the agent is unable to rationalize its decisions within the restricted generation budget. This is similar to the performance without CoT, but in contrast, the agent is instructed to produce the reasoning process. Over the course of RLFT, the agents learn to produce short rationales z_{CoT} , including the action a due to our reward shaping mechanism (see Section 3.2). However, the produced short rationales are unhelpful to improving agent performance.

In contrast, doubling G from 256 to 512 results in a considerable performance increase to the level of Gemma2 9B with RLFT (see Figure 6). We observe an increase in the average sequence length over the course of the RLFT process. This suggests that the agent learns to effectively leverage the additional “thinking time” and reflects recent observations in mathematical reasoning [Guo et al., 2025]. However, the increased performance comes with additional training cost due to the multistep nature of decision-making scenarios. In fact, we observed that rollout generation can make up the majority of the training time required by the RLFT process. This is because the agent has to produce more tokens at every environment interaction step. For example, for our default horizon of 50 timesteps and a generation budget of 500, the agent produces 25K tokens (at maximum).

F Limitations & Future Work

We focused our evaluation on the Gemma2 series and small-to-medium scale models. While we expect that our findings transfer to larger models, we deem research into frontier models important. Moreover, our MAB experiments were conducted with a limited horizon of 50 environment steps, which is sufficient for 5 and 10 arms, but insufficient for 20 arms. For future work, we believe that evaluating the exploration abilities of LLM agents is particularly interesting in environments that require targeted exploration towards an end-goal. First, this includes other stateful environments from Paglieri et al. [2024] and Ruoss et al. [2024], such as Crafter [Hafner, 2021]. Second, we deem a systematic investigation into exploration abilities of LLMs in existing agentic benchmarks [Mialon et al., 2023, He et al., 2024, Zhou et al., 2023] interesting. In our ablation studies, we found that LLMs benefit from additional “thinking” time and believe that allowing for a larger generation budget will become increasingly important for agentic scenarios, especially for scenarios with high-stakes decisions (e.g., economics or ethics). We deem investigations into such high-stakes scenarios fruitful for future work. While increasing “thinking” time improves performance, it comes with excessive computational cost at training time due to the rollout generation and the multi-step nature of decision-making. Therefore, modern recurrent architectures [De et al., 2024, Beck et al., 2025] that allow for faster inference may be promising alternatives for decision-making.

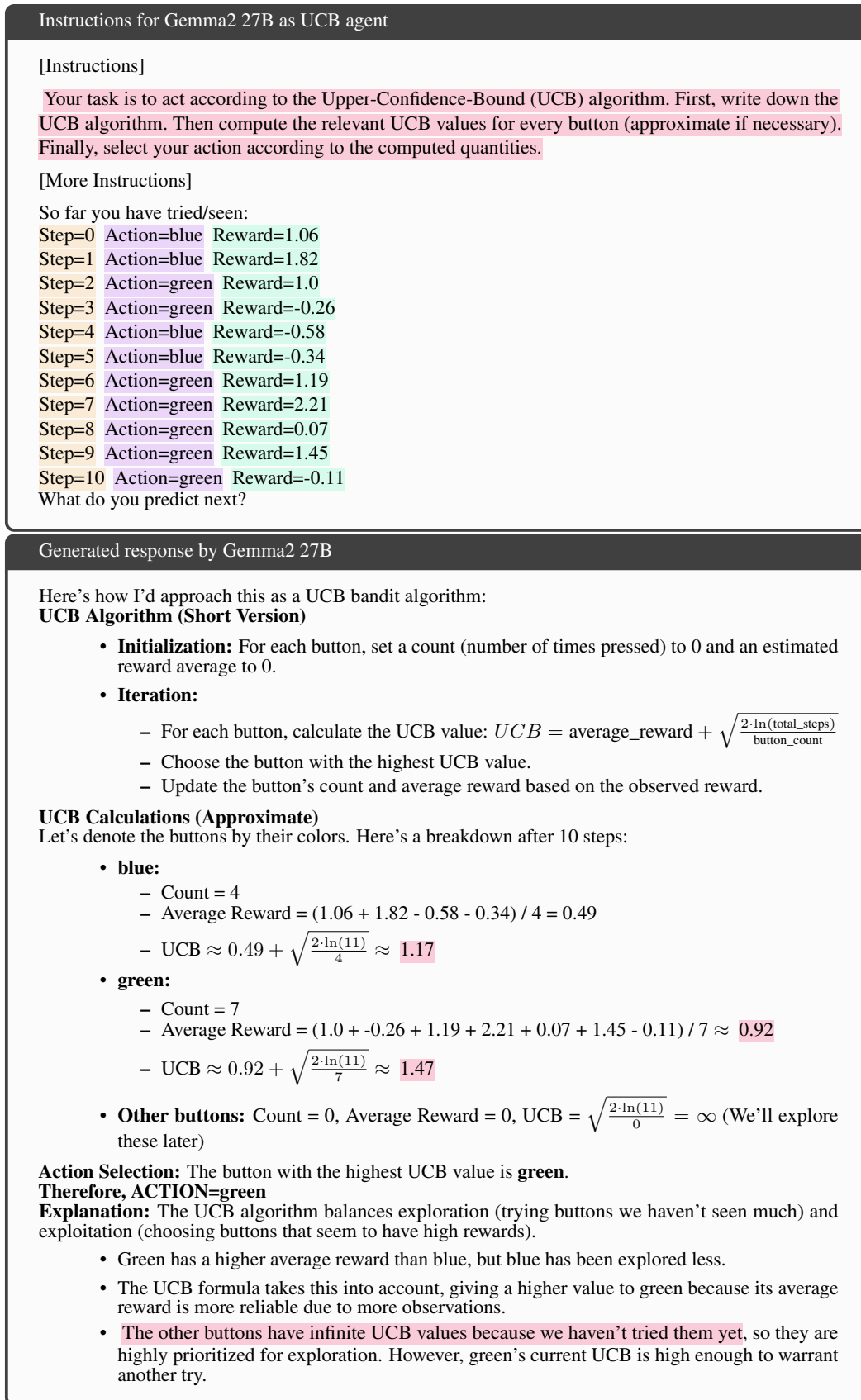


Figure 21: Illustration of the **knowing-doing gap**. (a) Instructions for the agent, which is prompted to act like a UCB algorithm. (b) The response generated by Gemma2 27B with greedy decoding (temperature=0). The LLM “knows” the UCB algorithm and computes the relevant quantities approximately correctly, but acts erroneously by selecting the next action greedily.