

---

# LLM-PriorCB: Textual Contextual Bandits with LLM-Induced Priors

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Many real-world decision-making problems, including recommendation, tool se-  
2 lection, and model routing, require choosing among actions described in natural  
3 language while observing feedback only for the selected action. Contextual bandits  
4 provide a principled framework for exploration under such partial feedback, but  
5 standard implementations typically start from uninformative reward-model priors  
6 and can adapt slowly in cold-start or evolving-action regimes. Large language  
7 models (LLMs) offer useful semantic knowledge from reward-free text, yet using  
8 an LLM directly as the online decision maker does not by itself provide calibrated  
9 uncertainty or principled exploration. We study textual contextual bandits, where  
10 both contexts and actions are represented by text, action descriptions are available  
11 without reward labels, and the action set may evolve over time. We propose LLM-  
12 PriorCB, a two-stage framework that separates semantic prior construction from  
13 online learning. In the offline stage, an LLM estimates rewards for reward-free  
14 context-action text pairs, and these estimates are distilled into action-specific prior  
15 parameters. In the online stage, LLM-PriorCB applies a disjoint linear UCB rule  
16 initialized by these action-specific priors, and updates only the selected action’s  
17 parameter using observed rewards, without further LLM queries. We derive prior-  
18 dependent regret guarantees for the resulting disjoint linear bandit, showing how  
19 prior misspecification affects the confidence radius while preserving no-regret  
20 learning. Experiments on MovieLens and GAIA show that LLM-PriorCB reduces  
21 cumulative regret relative to LinUCB, CBLLI, and prompted LLM policy baselines.

## 22 1 Introduction

23 Many real-world applications make decisions from natural-language inputs by selecting among  
24 alternatives also described in natural language, including recommender systems [21, 14, 26], LLM  
25 routers [18, 7], and tool-use or agentic systems [22, 4, 5]. These settings naturally involve bandit  
26 feedback, since only the outcome of the selected action is observed, and often feature evolving action  
27 sets as new items, models, APIs, or tools are introduced [15, 28]. Thus, effective decision making  
28 requires both exploration under partial feedback and rapid generalization to new actions.

29 LLMs are attractive in this setting because they encode broad linguistic, factual, and commonsense  
30 knowledge from large-scale pretraining [9, 20, 6], and have been used for recommendation, routing,  
31 and tool selection [21, 22, 18, 7]. However, directly using an LLM as the decision maker does not  
32 solve online learning [19]. Standard LLM pipelines lack decision-relevant posterior uncertainty,  
33 principled exploration from partial feedback, and regret guarantees. Moreover, empirical studies show  
34 that native in-context exploration by LLM agents is brittle and benefits from explicit algorithmic  
35 support [12, 17]. Semantic prior knowledge alone is therefore insufficient for reliable decision making  
36 under partial feedback.

37 Contextual bandits provide principled exploration strategies and theoretical guarantees for partial-  
 38 feedback learning [15, 8, 1, 2, 13], with recent extensions to neural function approximation and large  
 39 action spaces [27, 28]. While modern text encoders can provide useful semantic representations for  
 40 contexts and actions, existing bandit methods still typically initialize reward models from uninfor-  
 41 mative priors and adapt them only through sparse online feedback. As a result, even with pretrained  
 42 embeddings, the learner may require substantial interaction to identify reward-relevant structure,  
 43 particularly in cold-start settings and environments with evolving action spaces.

44 Prior work shows that warm-start signals can reduce early regret, but must be robust to mismatch  
 45 with online rewards [25]. Recent evidence also suggests that LLM-generated information can provide  
 46 useful prior knowledge for contextual bandits [3]. This raises our key question: can we use the  
 47 semantic prior knowledge of LLMs while preserving the uncertainty-aware online learning behavior  
 48 of contextual bandits?

49 We formalize *textual contextual bandits*, where both contexts and actions are natural-language  
 50 descriptions and the available action set may evolve over time. At each round, the learner observes a  
 51 textual context and a set of text-described actions, selects one action, and observes only its reward.  
 52 The goal is to use textual semantics to improve early decisions while allowing online feedback to  
 53 correct wrong prior beliefs.

54 We propose LLM-PriorCB, a framework that integrates LLM-induced priors with online contextual  
 55 bandit algorithms. The LLM constructs prior beliefs over context-action pairs from their textual  
 56 descriptions, while posterior updates rely on observed rewards. Crucially, the LLM does not directly  
 57 choose actions. Instead, it provides an informative initialization, while the bandit algorithm governs  
 58 online action selection and exploration. When the prior is informative, this initialization improves  
 59 sample efficiency and reduces early-stage regret. When the prior is misspecified, online reward  
 60 feedback can progressively correct the induced beliefs.

61 Our contributions are:

- 62 • We introduce *textual contextual bandits*, covering textual contexts, textual actions, partial  
 63 feedback, and evolving action sets.
- 64 • We propose LLM-PriorCB, which converts LLM-induced semantic information into priors  
 65 for online contextual bandit learning.
- 66 • We prove the asymptotic optimality of LLM-PriorCB, even under prior misspecification.
- 67 • We empirically show that LLM-PriorCB achieves stable and robust performance.

## 68 2 Preliminaries

69 This section introduces the contextual bandit notation used throughout the paper and recalls the  
 70 feature-based reward models on which our method builds.

71 **Contextual bandits.** We consider a contextual bandit over horizon  $T$ . At each round  $t \in [T]$ , the  
 72 learner observes a context  $c_t$  and a finite available action set  $\mathcal{A}_t$ , selects  $a_t \in \mathcal{A}_t$ , and observes only  
 73 the reward of the selected action. For each action  $a$ , let  $Y_t(a)$  denote its potential reward and define

$$\mu(c_t, a) := \mathbb{E}[Y_t(a) \mid c_t], \quad y_t = Y_t(a_t) = \mu(c_t, a_t) + \eta_t,$$

74 where  $\eta_t$  is conditionally mean-zero noise. The optimal action, optimal expected reward, instantaneous  
 75 regret, and cumulative regret are defined as

$$a_t^* \in \arg \max_{a \in \mathcal{A}_t} \mu(c_t, a), \quad \mu_t^* = \mu(c_t, a_t^*), \quad r_t = \mu_t^* - \mu(c_t, a_t), \quad R_T = \sum_{t=1}^T r_t.$$

76 A policy selects actions using the past interaction history and the current context-action set.

77 **Feature-based models.** Many contextual bandit algorithms represent each context-action pair by a  
 78 feature vector  $z_{t,a} \in \mathbb{R}^D$ . Linear methods assume

$$\mu(c_t, a) = \mu(z_{t,a}) = z_{t,a}^\top \theta^*,$$

79 where  $\theta^* \in \mathbb{R}^D$  is an unknown parameter vector. They then use confidence sets or posterior sampling  
 80 for exploration, as in LinUCB and Thompson sampling [15, 2]. For our purpose, the key limitation is  
 81 the lack of informative reward-aware priors over the bandit parameters. Even with strong pretrained  
 82 embeddings, contextual bandit algorithms typically initialize reward models from uninformative  
 83 priors and must adapt them using sparse online feedback.

### 84 3 Textual Contextual Bandits

85 We define a textual contextual bandit as a contextual bandit in which both contexts and actions are  
 86 described by natural language. At round  $t$ , the learner observes a textual context  $c_t \in \mathcal{C}$  and an  
 87 available action set  $\mathcal{A}_t$ , where each action  $a \in \mathcal{A}_t$  has a textual description  $s_a \in \mathcal{S}$ . Let

$$\psi : \mathcal{C} \cup \mathcal{S} \rightarrow \mathbb{R}^d$$

88 denote a text embedding model that maps textual contexts and action descriptions into vector  
 89 representations, and let

$$\Psi : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^D$$

90 denote a joint feature map that combines context and action representations. The textual information  
 91 is embedded and combined as

$$z_t = \psi(c_t), \quad z_a = \psi(s_a), \quad z_{t,a} = \Psi(z_t, z_a)$$

92 and the expected reward is written as  $\mu(z_{t,a}) = \mu(c_t, a)$ . Accordingly,

$$a_t^* \in \arg \max_{a \in \mathcal{A}_t} \mu(z_{t,a}), \quad \mu_t^* = \mu(z_{t,a_t^*}), \quad r_t = \mu_t^* - \mu(z_{t,a_t}), \quad R_T = \sum_{t=1}^T r_t.$$

93 For analysis, we consider either a linear model

$$\mu(z_{t,a}) = z_{t,a}^\top \theta^*$$

94 or a disjoint arm-specific model

$$\mu(z_{t,a}) = z_{t,a}^\top \theta_a^*.$$

95 In the disjoint setting,  $\theta_{0,a}$  denotes the prior parameter for arm  $a$ , and  $\hat{\theta}_{t,a}$  denotes its online estimate.  
 96 LLM-PriorCB uses the LLM as a source of prior information rather than as the final decision maker.  
 97 Specifically, the LLM uses textual descriptions to construct representations or prior parameters, while  
 98 the bandit algorithm selects actions and updates estimates from observed rewards. This separation  
 99 allows the learner to exploit semantic knowledge before sufficient feedback is available while  
 100 preserving uncertainty-aware exploration and regret analysis.

101 Before online interaction, the learner may access a reward-free textual dataset

$$\mathcal{D}_0 = (\mathcal{C}_0, \mathcal{A}_0, \{s_a\}_{a \in \mathcal{A}_0})$$

102 which contains textual contexts and action descriptions but no logged reward tuples  $(c, a, y)$ . In our  
 103 framework,  $\mathcal{D}_0$  is used only to extract LLM-induced prior information before deployment.

### 104 4 LLM-PriorCB

105 We propose LLM-PriorCB, a two-stage framework for textual contextual bandits. In the offline stage,  
 106 an LLM is used to construct action-specific prior parameters from the reward-free textual dataset  $\mathcal{D}_0$   
 107 defined in Section 3. In the online stage, a prior-centered UCB algorithm selects actions and updates  
 108 its estimates using observed rewards. After online interaction begins, the LLM is no longer queried.  
 109 Thus, the LLM provides only prior information, while exploration and adaptation are handled by the  
 110 bandit algorithm.

111 **4.1 Offline Construction of LLM-Induced Priors**

112 We focus on the disjoint linear model from Section 3, where each action  $a$  has an unknown parameter  
 113  $\theta_a^*$  and reward model  $\mu(z_{t,a}) = z_{t,a}^\top \theta_a^*$ . The goal of the offline stage is to construct an informative  
 114 prior parameter  $\theta_{0,a} \in \mathbb{R}^D$  for each action  $a$ .

115 For each offline context  $c_i \in \mathcal{C}_0$  and action  $a \in \mathcal{A}_0$ , we query an LLM  $\mathcal{M}$  with the textual pair  
 116  $(c_i, s_a)$  to obtain an estimated expected reward. When using  $M$  samples, we average the LLM outputs  
 117 as

$$\hat{\mu}_{i,a} = \frac{1}{M} \sum_{m=1}^M \hat{\mu}_{i,a}^{(m)}.$$

118 Here,  $\hat{\mu}_{i,a}^{(m)}$  denotes the  $m$ -th LLM-generated estimate for  $(c_i, s_a)$ .

119 We introduce a prior encoder  $f_\phi : \mathbb{R}^d \rightarrow \mathbb{R}^D$  that maps the action embedding  $z_a = \psi(s_a)$  to the  
 120 prior parameter  $\theta_{0,a} = f_\phi(z_a)$ . The encoder is trained to make the induced linear predictor match the  
 121 LLM-generated reward estimates:

$$\min_{\phi} \frac{1}{|\mathcal{C}_0| |\mathcal{A}_0|} \sum_{c_i \in \mathcal{C}_0} \sum_{a \in \mathcal{A}_0} (\hat{\mu}_{i,a} - z_{i,a}^\top f_\phi(z_a))^2. \quad (1)$$

122 After training, we set  $\theta_{0,a} = f_\phi(z_a)$ . For a newly introduced online action  $a \notin \mathcal{A}_0$ , we initialize its  
 123 prior in the same way from its textual description  $s_a$ .

124 **4.2 Online Learning with Prior-Centered UCB**

125 Given the LLM-induced priors  $\{\theta_{0,a}\}$ , online learning proceeds with an action-wise prior-centered  
 126 UCB rule. For each action  $a$ , initialize

$$A_{0,a} = \lambda I_D, \quad b_{0,a} = \lambda \theta_{0,a}.$$

127 At round  $t$ , define the current estimate

$$\hat{\theta}_{t-1,a} = A_{t-1,a}^{-1} b_{t-1,a}.$$

128 For each available action  $a \in \mathcal{A}_t$ , compute

$$U_t(a) = z_{t,a}^\top \hat{\theta}_{t-1,a} + \beta_{t-1}^{\text{dis}}(\delta) \sqrt{z_{t,a}^\top A_{t-1,a}^{-1} z_{t,a}},$$

129 where  $\beta_{t-1}^{\text{dis}}(\delta)$  is an exploration coefficient controlling the uncertainty bonus. The learner selects

$$a_t \in \arg \max_{a \in \mathcal{A}_t} U_t(a).$$

130 After observing  $y_t$ , only the selected action is updated:

$$A_{t,a_t} = A_{t-1,a_t} + z_{t,a_t} z_{t,a_t}^\top, \quad b_{t,a_t} = b_{t-1,a_t} + y_t z_{t,a_t}.$$

131 For all  $a \neq a_t$ , set  $A_{t,a} = A_{t-1,a}$  and  $b_{t,a} = b_{t-1,a}$ .

132 **5 Theoretical Derivation**

133 **Theorem 5.1** (Cumulative regret analysis). *Suppose Assumptions A.1, A.2, A.3, A.4, and A.5 hold.*  
 134 *Suppose that at every round  $t = 1, \dots, T$ , the selected arm  $a_t$  satisfies*

$$a_t \in \arg \max_{a \in \mathcal{A}_t} \left[ z_{t,a}^\top \hat{\theta}_{t-1} + \beta_{t-1}(\delta) \sqrt{z_{t,a}^\top A_{t-1}^{-1} z_{t,a}} \right].$$

135 *Then, with probability at least  $1 - \delta$ ,*

$$R_T := \sum_{t=1}^T r_t \leq 2\beta_T(\delta) \sqrt{2T c_\lambda \log \frac{\det(A_T)}{\det(\lambda I_D)}},$$

136 where

$$\beta_t(\delta) := R\sqrt{2\log\frac{\det(A_t)^{1/2}}{\det(\lambda I_D)^{1/2}\delta}} + \sqrt{\lambda}S_\Delta$$

137 and

$$c_\lambda := \max\left\{1, \frac{L^2}{\lambda}\right\}.$$

138 The above theorem can be specialized to LLM-PriorCB as follows.

139 **Corollary 5.2** (Cumulative regret analysis for LLM-PriorCB). *Suppose the assumptions in Theorem 5.1 hold. Consider the disjoint linear model*

$$\mu(z_{t,a}) = z_{t,a}^\top \theta_a^*, \quad z_{t,a} \in \mathbb{R}^D, \quad a \in [K],$$

141 with arm-specific prior parameters  $\theta_{0,a} \in \mathbb{R}^D$ . Then the arm-wise LLM-PriorCB rule induced by  
142 this block representation satisfies, with probability at least  $1 - \delta$ ,

$$R_T \leq 2\beta_T^{\text{dis}}(\delta)\sqrt{2Tc_\lambda \log\frac{\prod_{a=1}^K \det(A_{T,a})}{\lambda^{Kd}}},$$

143 where

$$A_{t,a} := \lambda I_d + \sum_{i=1}^t \mathbf{1}\{a_i = a\} z_{i,a} z_{i,a}^\top$$

144 and

$$\beta_t^{\text{dis}}(\delta) := R\sqrt{2\log\frac{\prod_{a=1}^K \det(A_{t,a})^{1/2}}{\lambda^{Kd/2}\delta}} + \sqrt{\lambda}S_\Delta.$$

## 145 6 Experiments

146 We evaluate our method on two benchmark tasks constructed from the MovieLens dataset [10] and  
147 the GAIA dataset [16].

### 148 6.1 MovieLens

149 We adopt a semi-synthetic contextual bandit task based on the MovieLens dataset [10], following an  
150 experimental design introduced in prior work, most notably EVOLvE [17]. This semi-synthetic setup  
151 provides controlled access to key components of the environment, such as oracle context features and  
152 ground-truth rewards, while preserving realistic user and movie information derived from real-world  
153 data.

154 **Dataset** We use the MovieLens-100k dataset [10], a widely adopted benchmark for recommendation  
155 system evaluation. It consists of 100,000 explicit ratings on a 5-point scale, collected from  $N = 943$   
156 users on 1,682 movies. Each user has provided at least 20 ratings and is associated with demographic  
157 attributes, including age, gender, occupation, and zip code. Each movie is described by metadata  
158 such as its title, release year, and genre tags. Together, these features provide a natural testbed for  
159 studying contextual decision-making with textual side information. Following prior work, we restrict  
160 the action set  $\mathcal{A}$  to the top- $K$  most-watched movies, where  $K = 10$  corresponds to a relatively easy  
161 setting and  $K = 30$  induces a more challenging problem.

162 **Environment** To construct the ground-truth reward function, we apply singular value decomposition  
163 (SVD) to the user–movie rating matrix  $P \in \mathbb{R}^{N \times K}$  and obtain a low-rank approximation [11].  
164 Specifically, we approximate  $P$  as  $\tilde{P} = U\Sigma V^\top$ , where  $U$  and  $V$  denote user and movie embedding  
165 matrices, respectively. The reward for recommending movie  $a$  to user  $i$  is computed as  $\mu(z_{i,a}) =$   
166  $u_i^\top \Sigma v_a$ .

167 At each time step, a user is sampled, and the agent receives contextual information according to the  
168 current information regime. The agent selects one movie from a finite action set and receives a scalar  
169 reward corresponding to the user’s preference.

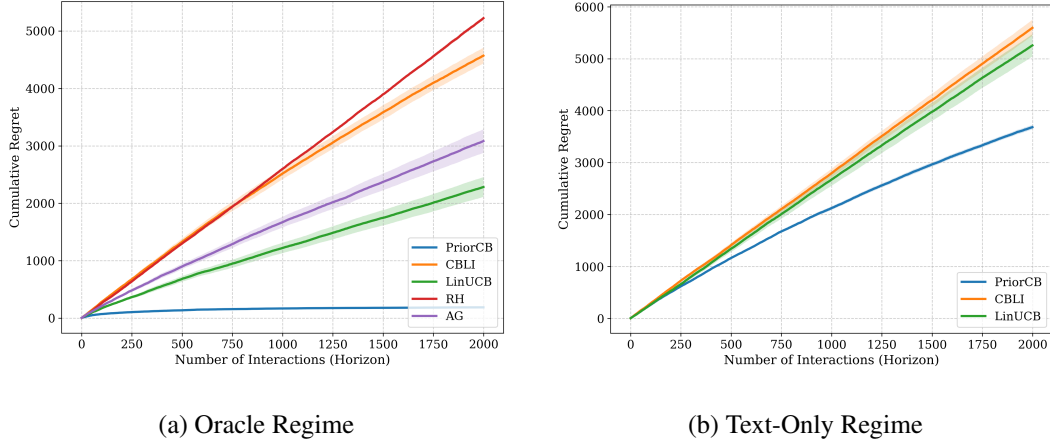


Figure 1: Performance as a function of the interaction horizon for  $K = 10$ . PriorCB (ours) consistently achieves significantly lower cumulative regret compared to baseline methods, demonstrating strong sample efficiency in the oracle regime. We observe a similar qualitative trend in the text-only regime, indicating that the advantage of PriorCB persists even when only textual information is provided. Shaded regions indicate  $\pm$  one standard error over 10 independent runs.

170 **Information Regimes** We consider two information regimes. (i) *Oracle Context Regime*: the agent  
 171 has direct access to latent user context features  $u_i$  obtained from the low-rank approximation, which  
 172 are also used to generate rewards. (ii) *Text-Only Context Regime*: the latent user context features are  
 173 hidden, and the agent must infer effective context features solely from textual descriptions of user  
 174 demographic information. The Oracle Context Regime closely resembles the classical contextual  
 175 bandit setting, whereas the Text-Only Context Regime more closely reflects practical real-world  
 176 scenarios.

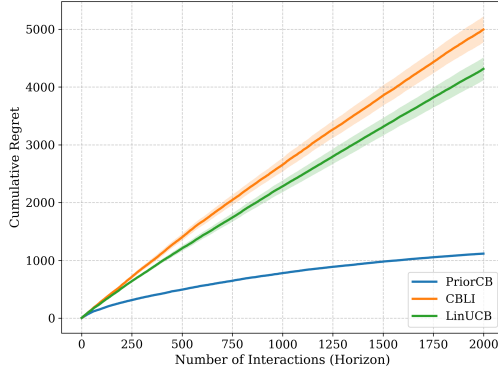
177 **Baselines and Experimental Setup** In the Oracle Context Regime, we compare our method  
 178 against CBLI [3], standard LinUCB [15], and two prompted LLM policy baselines, RH and AG,  
 179 both of which prompt an LLM to directly select actions [17]. RH represents past interactions as raw  
 180 (context, action, reward) tuples and provides them, together with the current context and available  
 181 actions, as the LLM prompt. AG augments the prompt with LinUCB-derived exploitation values and  
 182 exploration bonuses, so that the LLM receives explicit algorithmic guidance for action selection. In  
 183 the Text-Only Context Regime, we focus on baselines whose action-selection rules are fully specified  
 184 under the same text-derived feature representation. We therefore report CBLI [3] and LinUCB [15].  
 185 We omit RH and AG because they are prompted LLM policy baselines requiring LLM inference  
 186 at every round, and adapting AG’s oracle-context LinUCB guidance to text-only features would  
 187 introduce an additional implementation variant beyond the original baseline.

188 For prior construction, we use QWEN3-4B-INSTRUCT-2507 [24] as the LLM  $\mathcal{M}$  and QWEN3-  
 189 EMBEDDING-8B as the text encoder  $\psi$ . We select a subset of 100 users from the dataset as  $\mathcal{D}_0$ . For  
 190 each user  $i$  and movie  $a$  in  $\mathcal{D}_0$ , we estimate the corresponding rating by prompting the LLM  $\mathcal{M}$  with  
 191 the user’s demographic attributes  $c_i$  and the movie’s metadata  $s_a$ . Each user–movie pair is queried 10  
 192 times, and the resulting predictions are averaged to obtain prior reward estimates. These estimated  
 193 ratings are used both for prior construction in our algorithm and for CBLI pretraining.

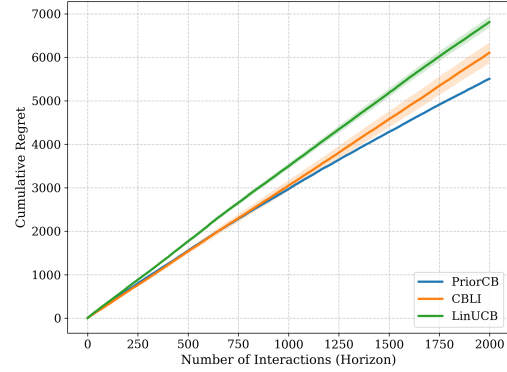
194 For the action encoder, the mapping function  $f_\phi$  is implemented as a two-layer multilayer perceptron  
 195 (MLP). We train  $f_\phi$  by minimizing the MSE objective defined in Equation 1. During online learning,  
 196 users are sampled from a dataset disjoint from  $\mathcal{D}_0$ . All online experiments are conducted over 10  
 197 random seeds, and performance is evaluated using cumulative regret. Details of the prompting  
 198 templates and training are provided in Appendix D.1.

## 199 6.2 Experimental results

200 **Experiments with  $K = 10$**  We first evaluate all methods in an oracle regime, where the agent is  
 201 provided with oracle-quality action representations. This setting isolates the effect of the learning al-

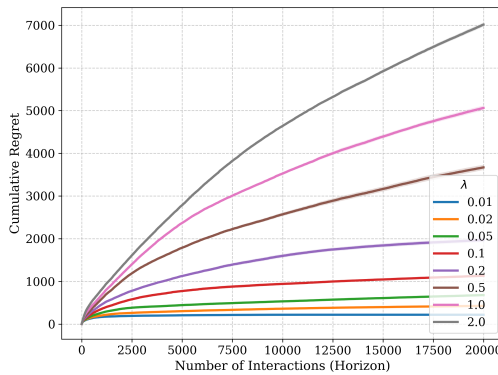


(a) Oracle Regime

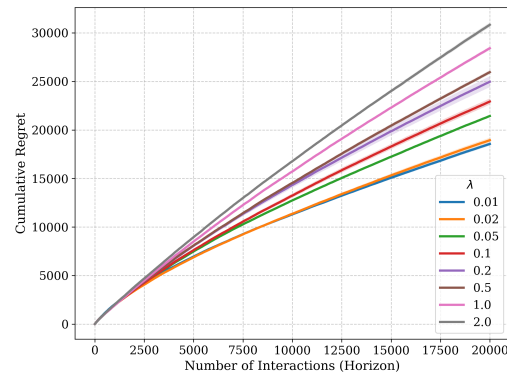


(b) Text-Only Regime

Figure 2: Performance as a function of the interaction horizon when  $K = 10$  in prior construction and  $K = 30$  when online learning. Shaded regions indicate  $\pm$  one standard error over 10 independent runs.



(a) Oracle Regime



(b) Text-Only Regime

Figure 3: We report cumulative regret of LLM-PriorCB for  $\lambda \in [0.01, 2.0]$  under (a) the oracle regime and (b) the text-only regime. Smaller  $\lambda$  consistently yields lower regret, while larger  $\lambda$  degrades performance, indicating that overly confident priors can suppress exploration and hinder effective online adaptation.

202 algorithm itself by removing confounding factors arising from noisy or mis-specified textual encodings  
 203 of the context. As shown in Figure 1(a), LLM-PriorCB (abbreviated as PriorCB in figures) achieves  
 204 significantly lower cumulative regret than all baseline methods across the entire horizon, demon-  
 205 strating strong sample efficiency. This result indicates that LLM-PriorCB can effectively exploit  
 206 informative representations while maintaining sufficient exploration. In contrast, LinUCB performs  
 207 poorly in this setting, as it learns action representations entirely from scratch and therefore requires a  
 208 substantial amount of exploration. CBLI also exhibits limited exploration due to its inherent design,  
 209 leading to suboptimal performance even under oracle access. As a result, its convergence is slower  
 210 than that of the naive LinUCB. Furthermore, RH and AG fail to achieve competitive performance,  
 211 suggesting that relying solely on LLM-based decision making without principled online learning and  
 212 exploration mechanisms is insufficient.

213 We then evaluate the methods in a text-only regime (Figure 1(b)), where context and action repre-  
 214 sentations are derived solely from natural language descriptions. This setting is closer to real-world  
 215 problem scenarios, and despite the additional noise and ambiguity introduced by textual representa-  
 216 tions, LLM-PriorCB maintains a clear advantage over competing methods, indicating that its benefits  
 217 are not limited to idealized oracle settings. In Appendix B, we report additional results over an  
 218 extended horizon of up to 20,000 steps. Together, these results demonstrate that LLM-PriorCB

219 effectively leverages prior information both when representations are reliable and when they must be  
220 inferred from noisy textual descriptions.

221 **Experiments with dynamic action set** Furthermore, we evaluate the methods in an arm-evolving  
222 setting, where the available action set changes over time. Specifically, the offline textual dataset  $\mathcal{D}_0$   
223 contains prior information for only  $K = 10$  actions, while online learning is conducted with an  
224 expanded action set of  $K = 30$  actions. Thus, many actions appearing during online interaction are  
225 not included in the offline prior construction stage. As shown in Figure 2, LLM-PriorCB consistently  
226 outperforms baseline methods despite this mismatch between the offline prior set and the online  
227 action set. While this experiment represents only a limited evaluation of evolving-action scenarios,  
228 the results suggest that LLM-PriorCB can effectively accommodate newly introduced actions by  
229 leveraging textual prior information.

230 **Ablation studies** We further conduct an ablation study on the exploration parameter  $\lambda$ , which  
231 controls the strength of the prior in LLM-PriorCB. As shown in Figure 3 smaller values of  $\lambda$  con-  
232 sistently lead to lower cumulative regret in this experiments. This observation suggests that overly  
233 trusting the prior can hinder effective exploration, preventing the agent from adequately incorporating  
234 online feedback. This behavior is consistent with the poor performance of CBLI observed in earlier  
235 experiments, where limited exploration due to strong reliance on prior estimates leads to suboptimal  
236 outcomes. Additional ablation studies on the feature dimension, rank, and the number of actions are  
237 provided in Appendix B.

### 238 6.3 GAIA

239 **Dataset** We conduct experiments on the GAIA [16] benchmark, a dataset designed to evaluate  
240 General AI Assistants on realistic user queries that require a combination of fundamental abilities.  
241 The questions are formulated to have a brief, unambiguous correct answer, making verification  
242 straightforward. We construct a text-only subset of GAIA by filtering out instances that require visual  
243 understanding. After filtering, the subset contains 127 questions from validation set and 230 questions  
244 from test set.

245 **Environment** We formulate tool selection as a contextual bandit with a set of  $K = 3$  tools. The  
246 three arms correspond to reason-only (using only internal knowledge without any web search; low  
247 cost), search-light (using single web search with one query; mid cost), and search-deep (using  
248 three web search with three distinct queries; high cost). Detailed descriptions are provided in the  
249 Appendix D.2.2. At each time step, the agent receives a question, and the agent selects one tool  
250 from a set of  $K$  tools, after which it receives a scalar reward defined by a pre-designed function.  
251 Specifically, we set the expected reward for tool  $a$  to question  $c_i$  as  $\mu(z_{t,a}) = accuracy - cost$ ,  
252 where  $accuracy \in \{0, 1\}$  is computed via quasi-exact match against the ground-truth answer, and  
253 the  $cost$  is a fixed tool-dependent penalty.

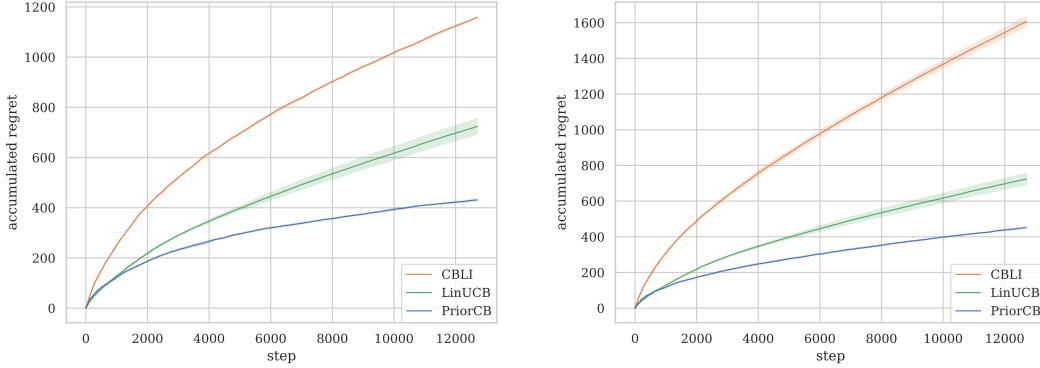
254 **Baelines** We consider two baselines: naive LinUCB [15] with froze embedding backbone, CBLI [3].

255 **Experimental setup** We evaluate gpt-5-mini [23] and gpt-5.2 on GAIA benchmark. These models  
256 are used as prior LLMs for offline prior construction, while online learning and evaluation are  
257 performed on the validation set under a contextual bandit.

258 For prior construction, we use the GAIA benchmark test set as  $\mathcal{D}_0$ . For each question  $c_i$  and tool  $a$ ,  
259 we obtain an estimated score by prompting an LLM with the question together with the descriptions  
260 of all available tools. The full prompt template is provided in Appendix D.2.2. Subsequently, we train  
261 a two-layer MLP action encoder on the generated dataset, using QWEN/QWEN3-EMBEDDING-8B as  
262 a frozen embedding backbone. The experimental details are explained in Appendix D.2.3.

263 For online learning, we use the GAIA benchmark validation set. Due to the limited number of  
264 validation instances, we repeat the online evaluation for 100 iterations (with reshuffling and different  
265 random seeds).

266 **Experimental results** Figure 4 reports the cumulative regret of PriorCB, LinUCB, and CBLI in  
267 the main GAIA experiment. Overall, PriorCB achieves substantially lower cumulative regret than



(a) Cumulative regret of Prior from different eval contexts      (b) Cumulative regret of Prior from same eval contexts

Figure 4: Performance as a function of the interaction horizon when  $K = 3$  using gpt-5.2 with medium reasoning effort. Shaded regions indicate  $\pm$  one standard error over 3 independent runs. We use the validation split for evaluation due to answer-label availability. Priors are constructed either from the validation set (same eval contexts) or from the test set (different eval contexts), using no labels.

268 the baselines, indicating that the LLM-induced prior provides useful initialization for subsequent  
 269 online learning. Figure 4(a) and Figure 4(b) correspond to priors constructed from a different context  
 270 set and from the evaluation context set, respectively. In both cases, prior construction uses only  
 271 reward-free textual information and does not access environment rewards, labels, or answers. CBLI  
 272 shows substantial variation depending on the dataset used for prior construction. By contrast, PriorCB  
 273 remains stable across the two prior-construction settings, while LinUCB is also stable but lacks  
 274 the benefit of LLM-induced initialization. Additional results for the prompted LLM policy baseline  
 275 RH [17] are provided in Appendix C.1. Since RH requires LLM inference at every decision step, we  
 276 report it separately from the main long-horizon comparison. Appendix C.1 also includes RH results  
 277 with GPT-5.2 and GPT-5-mini, together with the long-horizon GPT-5-mini results for LinUCB, CBLI,  
 278 and PriorCB. The average reward results are reported in Appendix C.2.

## 279 7 Conclusion

280 We studied *textual contextual bandits*, where both contexts and actions are described in natural  
 281 language and the action set may evolve over time. To bridge LLM-based decision making and  
 282 principled bandit learning, we proposed LLM-PriorCB, a two-stage framework that uses an LLM only  
 283 to construct reward-free offline priors from textual information, and then performs online learning  
 284 solely through a standard contextual bandit algorithm. Theoretically, we establish prior-dependent  
 285 cumulative regret guarantees for disjoint linear contextual bandits, showing that well-aligned offline  
 286 priors can reduce the regret bound, while prior misspecification enters only through explicit arm-wise  
 287 bias terms. Across experiments on MovieLens and GAIA, LLM-PriorCB consistently achieved lower  
 288 regret than baselines under both oracle and text-only regimes, degraded gracefully under mis-specified  
 289 priors, and showed promising performance when the action set expanded between prior construction  
 290 and online learning. Our results highlight the importance of separating prior knowledge extraction  
 291 from online adaptation, enabling stable exploration and reliable improvement from sparse bandit  
 292 feedback.

293 There are several directions for future work. First, while our evolving-arm experiment provides initial  
 294 evidence of robustness, broader evaluations under more diverse non-stationary action dynamics would  
 295 further clarify the limits of the approach. Second, extending the prior construction stage to incorporate  
 296 richer uncertainty estimates (e.g., calibrated prior covariances) and exploring tighter integrations  
 297 with contextual bandit variants beyond linear models are promising avenues. Finally, applying the  
 298 framework to real-world deployments with latency and cost constraints remains an important step  
 299 toward practical LLM-assisted online decision-making systems.

## References

- [1] Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. *Advances in neural information processing systems*, 24, 2011.
- [2] Shipra Agrawal and Navin Goyal. Thompson sampling for contextual bandits with linear payoffs. In *International conference on machine learning*, pages 127–135. PMLR, 2013.
- [3] Parand Alamdari, Yanshuai Cao, and Kevin Wilson. Jump starting bandits with llm-generated prior knowledge. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 19821–19833, 2024.
- [4] Anthropic. Introducing the model context protocol. Anthropic News, 2024.
- [5] Anthropic. Agent skills. Claude API Documentation, 2025.
- [6] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [7] Shuhao Chen, Weisen Jiang, Baijiong Lin, James Kwok, and Yu Zhang. Routerdc: Query-based router by dual contrastive learning for assembling large language models. *Advances in Neural Information Processing Systems*, 37:66305–66328, 2024.
- [8] Wei Chu, Lihong Li, Lev Reyzin, and Robert 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.
- [9] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pages 4171–4186, 2019.
- [10] F Maxwell Harper and Joseph A Konstan. The movielens datasets: History and context. *Acm transactions on interactive intelligent systems (tiis)*, 5(4):1–19, 2015.
- [11] Yehuda Koren, Robert Bell, and Chris Volinsky. Matrix factorization techniques for recommender systems. *Computer*, 42(8):30–37, 2009.
- [12] Akshay Krishnamurthy, Keegan Harris, Dylan J Foster, Cyril Zhang, and Aleksandrs Slivkins. Can large language models explore in-context? *Advances in Neural Information Processing Systems*, 37:120124–120158, 2024.
- [13] Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- [14] Lei Li, Yongfeng Zhang, Dugang Liu, and Li Chen. Large language models for generative recommendation: A survey and visionary discussions. In *Proceedings of the 2024 joint international conference on computational linguistics, language resources and evaluation (LREC-COLING 2024)*, pages 10146–10159, 2024.
- [15] Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th international conference on World wide web*, pages 661–670, 2010.
- [16] Grégoire Mialon, Clémentine Fourier, Thomas Wolf, Yann LeCun, and Thomas Scialom. Gaia: a benchmark for general ai assistants. In *The Twelfth International Conference on Learning Representations*, 2023.
- [17] Allen Nie, Yi Su, Bo Chang, Jonathan N Lee, Ed H Chi, Quoc V Le, and Minmin Chen. Evolve: Evaluating and optimizing llms for exploration. *arXiv preprint arXiv:2410.06238*, 2024.
- [18] Isaac Ong, Amjad Almahairi, Vincent Wu, Wei-Lin Chiang, Tianhao Wu, Joseph E Gonzalez, M Waleed Kadous, and Ion Stoica. Routellm: Learning to route llms with preference data. *arXiv preprint arXiv:2406.18665*, 2024.

- 346 [19] Chanwoo Park, Xiangyu Liu, Asuman Ozdaglar, and Kaiqing Zhang. Do llm agents have regret?  
347 a case study in online learning and games. *arXiv preprint arXiv:2403.16843*, 2024.
- 348 [20] Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu,  
349 and Alexander Miller. Language models as knowledge bases? In *Proceedings of the 2019*  
350 *conference on empirical methods in natural language processing and the 9th international joint*  
351 *conference on natural language processing (EMNLP-IJCNLP)*, pages 2463–2473, 2019.
- 352 [21] Scott Sanner, Krisztian Balog, Filip Radlinski, Ben Wedin, and Lucas Dixon. Large language  
353 models are competitive near cold-start recommenders for language-and item-based preferences.  
354 In *Proceedings of the 17th ACM conference on recommender systems*, pages 890–896, 2023.
- 355 [22] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro,  
356 Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models  
357 can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36:  
358 68539–68551, 2023.
- 359 [23] Aaditya Singh, Adam Fry, Adam Perelman, Adam Tart, Adi Ganesh, Ahmed El-Kishky, Aidan  
360 McLaughlin, Aiden Low, AJ Ostrow, Akhila Ananthram, Akshay Nathan, Alan Luo, Alec Helyar,  
361 Aleksander Madry, Aleksandr Efremov, Aleksandra Spyra, Alex Baker-Whitcomb, Alex Beutel,  
362 Alex Karpenko, Alex Makelov, Alex Neitz, Alex Wei, Alexandra Barr, Alexandre Kirchmeyer,  
363 Alexey Ivanov, Alexi Christakis, Alistair Gillespie, Allison Tam, Ally Bennett, Alvin Wan,  
364 Alyssa Huang, Amy McDonald Sandjideh, Amy Yang, Ananya Kumar, Andre Saraiva, Andrea  
365 Vallone, Andrei Gheorghe, Andres Garcia Garcia, Andrew Braunstein, Andrew Liu, Andrew  
366 Schmidt, Andrey Mereskin, Andrey Mishchenko, Andy Applebaum, Andy Rogerson, Ann  
367 Rajan, Annie Wei, Anoop Kotha, Anubha Srivastava, Anushree Agrawal, Arun Vijayvergiya,  
368 Ashley Tyra, Ashvin Nair, Avi Nayak, Ben Eggers, Bessie Ji, Beth Hoover, Bill Chen, Blair Chen,  
369 Boaz Barak, Borys Minaiev, Botao Hao, Bowen Baker, Brad Lightcap, Brandon McKinzie,  
370 Brandon Wang, Brendan Quinn, Brian Fioca, Brian Hsu, Brian Yang, Brian Yu, Brian Zhang,  
371 Brittany Brenner, Callie Riggins Zetino, Cameron Raymond, Camillo Lugaresi, Carolina Paz,  
372 Cary Hudson, Cedric Whitney, Chak Li, Charles Chen, Charlotte Cole, Chelsea Voss, Chen  
373 Ding, Chen Shen, Chengdu Huang, Chris Colby, Chris Hallacy, Chris Koch, Chris Lu, Christina  
374 Kaplan, Christina Kim, CJ Minott-Henriques, Cliff Frey, Cody Yu, Coley Czarnecki, Colin Reid,  
375 Colin Wei, Cory Decareaux, Cristina Scheau, Cyril Zhang, Cyrus Forbes, Da Tang, Dakota  
376 Goldberg, Dan Roberts, Dana Palmie, Daniel Kappler, Daniel Levine, Daniel Wright, Dave  
377 Leo, David Lin, David Robinson, Declan Grabb, Derek Chen, Derek Lim, Derek Salama,  
378 Dibya Bhattacharjee, Dimitris Tsipras, Dinghua Li, Dingli Yu, DJ Strouse, Drew Williams,  
379 Dylan Hunn, Ed Bayes, Edwin Arbus, Ekin Akyurek, Elaine Ya Le, Elana Widmann, Eli Yani,  
380 Elizabeth Proehl, Enis Sert, Enoch Cheung, Eri Schwartz, Eric Han, Eric Jiang, Eric Mitchell,  
381 Eric Sigler, Eric Wallace, Erik Ritter, Erin Kavanaugh, Evan Mays, Evgenii Nikishin, Fangyuan  
382 Li, Felipe Petroski Such, Filipe de Avila Belbute Peres, Filippo Raso, Florent Bekerman, Foivos  
383 Tsimplouras, Fotis Chantzis, Francis Song, Francis Zhang, Gaby Raila, Garrett McGrath, Gary  
384 Briggs, Gary Yang, Giambattista Parascandolo, Gildas Chabot, Grace Kim, Grace Zhao, Gregory  
385 Valiant, Guillaume Leclerc, Hadi Salman, Hanson Wang, Hao Sheng, Haoming Jiang, Haoyu  
386 Wang, Haozhun Jin, Harshit Sikchi, Heather Schmidt, Henry Aspegren, Honglin Chen, Huida  
387 Qiu, Hunter Lightman, Ian Covert, Ian Kivlichan, Ian Silber, Ian Sohl, Ibrahim Hammoud, Ignasi  
388 Clavera, Ikai Lan, Ilge Akkaya, Ilya Kostrikov, Irina Kofman, Isak Etinger, Ishaan Singal, Jackie  
389 Hehir, Jacob Huh, Jacqueline Pan, Jake Wilczynski, Jakub Pachocki, James Lee, James Quinn,  
390 Jamie Kiros, Janvi Kalra, Jasmyr Samaroo, Jason Wang, Jason Wolfe, Jay Chen, Jay Wang, Jean  
391 Harb, Jeffrey Han, Jeffrey Wang, Jennifer Zhao, Jeremy Chen, Jerene Yang, Jerry Tworek, Jesse  
392 Chand, Jessica Landon, Jessica Liang, Ji Lin, Jiancheng Liu, Jianfeng Wang, Jie Tang, Jihan  
393 Yin, Joanne Jang, Joel Morris, Joey Flynn, Johannes Ferstad, Johannes Heidecke, John Fishbein,  
394 John Hallman, Jonah Grant, Jonathan Chien, Jonathan Gordon, Jongsoo Park, Jordan Liss, Jos  
395 Kraaijeveld, Joseph Guay, Joseph Mo, Josh Lawson, Josh McGrath, Joshua Vendrow, Joy Jiao,  
396 Julian Lee, Julie Steele, Julie Wang, Junhua Mao, Kai Chen, Kai Hayashi, Kai Xiao, Kamyar  
397 Salahi, Kan Wu, Karan Sekhri, Karan Sharma, Karan Singhal, Karen Li, Kenny Nguyen, Keren  
398 Gu-Lemberg, Kevin King, Kevin Liu, Kevin Stone, Kevin Yu, Kristen Ying, Kristian Georgiev,  
399 Kristie Lim, Kushal Tirumala, Kyle Miller, Lama Ahmad, Larry Lv, Laura Clare, Lorraine  
400 Fauconnet, Lauren Itow, Lauren Yang, Laurentia Romaniuk, Leah Anise, Lee Byron, Leher  
401 Pathak, Leon Maksin, Leyan Lo, Leyton Ho, Li Jing, Liang Wu, Liang Xiong, Lien Mamitsuka,

402 Lin Yang, Lindsay McCallum, Lindsey Held, Liz Bourgeois, Logan Engstrom, Lorenz Kuhn,  
403 Louis Feuvrier, Lu Zhang, Lucas Switzer, Lukas Kondraciuk, Lukasz Kaiser, Manas Joglekar,  
404 Mandeep Singh, Mandip Shah, Manuka Stratta, Marcus Williams, Mark Chen, Mark Sun,  
405 Marselus Cayton, Martin Li, Marvin Zhang, Marwan Aljube, Matt Nichols, Matthew Haines,  
406 Max Schwarzer, Mayank Gupta, Meghan Shah, Melody Huang, Meng Dong, Mengqing Wang,  
407 Mia Glaese, Micah Carroll, Michael Lampe, Michael Malek, Michael Sharman, Michael Zhang,  
408 Michele Wang, Michelle Pokrass, Mihai Florian, Mikhail Pavlov, Miles Wang, Ming Chen,  
409 Mingxuan Wang, Minnia Feng, Mo Bavarian, Molly Lin, Moose Abdo, Mostafa Rohaninejad,  
410 Nacho Soto, Natalie Staudacher, Natan LaFontaine, Nathan Marwell, Nelson Liu, Nick Preston,  
411 Nick Turley, Nicklas Ansmann, Nicole Blades, Nikil Pancha, Nikita Mikhaylin, Niko Felix,  
412 Nikunj Handa, Nishant Rai, Nitish Keskar, Noam Brown, Ofir Nachum, Oleg Boiko, Oleg Murk,  
413 Olivia Watkins, Oona Gleeson, Pamela Mishkin, Patryk Lesiewicz, Paul Baltescu, Pavel Belov,  
414 Peter Zhokhov, Philip Pronin, Phillip Guo, Phoebe Thacker, Qi Liu, Qiming Yuan, Qinghua  
415 Liu, Rachel Dias, Rachel Puckett, Rahul Arora, Ravi Teja Mullapudi, Raz Gaon, Reah Miyara,  
416 Rennie Song, Rishabh Aggarwal, RJ Marsan, Robel Yemiru, Robert Xiong, Rohan Kshirsagar,  
417 Rohan Nuttall, Roman Tsiupa, Ronen Eldan, Rose Wang, Roshan James, Roy Ziv, Rui Shu,  
418 Ruslan Nigmatullin, Saachi Jain, Saam Talaie, Sam Altman, Sam Arnesen, Sam Toizer, Sam  
419 Toyer, Samuel Miserendino, Sandhini Agarwal, Sarah Yoo, Savannah Heon, Scott Ethersmith,  
420 Sean Grove, Sean Taylor, Sebastien Bubeck, Sever Banesiu, Shaokyi Amdo, Shengjia Zhao,  
421 Sherwin Wu, Shibani Santurkar, Shiyu Zhao, Shraman Ray Chaudhuri, Shreyas Krishnaswamy,  
422 Shuaiqi, Xia, Shuyang Cheng, Shyamal Anadkat, Simón Posada Fishman, Simon Tobin, Siyuan  
423 Fu, Somay Jain, Song Mei, Sonya Egoian, Spencer Kim, Spug Golden, SQ Mah, Steph Lin,  
424 Stephen Imm, Steve Sharpe, Steve Yadlowsky, Sulman Choudhry, Sungwon Eum, Suvansh  
425 Sanjeev, Tabarak Khan, Tal Stramer, Tao Wang, Tao Xin, Tarun Gogineni, Taya Christianson,  
426 Ted Sanders, Tejal Patwardhan, Thomas Degry, Thomas Shadwell, Tianfu Fu, Tianshi Gao,  
427 Timur Garipov, Tina Sriskandarajah, Toki Sherbakov, Tomer Kaftan, Tomo Hiratsuka, Tongzhou  
428 Wang, Tony Song, Tony Zhao, Troy Peterson, Val Kharitonov, Victoria Chernova, Vineet  
429 Kosaraju, Vishal Kuo, Vitchyr Pong, Vivek Verma, Vlad Petrov, Wanning Jiang, Weixing Zhang,  
430 Wenda Zhou, Wenlei Xie, Wenting Zhan, Wes McCabe, Will DePue, Will Ellsworth, Wulfie  
431 Bain, Wyatt Thompson, Xiangning Chen, Xiangyu Qi, Xin Xiang, Xinwei Shi, Yann Dubois,  
432 Yaodong Yu, Yara Khakbaz, Yifan Wu, Yilei Qian, Yin Tat Lee, Yinbo Chen, Yizhen Zhang,  
433 Yizhong Xiong, Yonglong Tian, Young Cha, Yu Bai, Yu Yang, Yuan Yuan, Yuanzhi Li, Yufeng  
434 Zhang, Yuguang Yang, Yujia Jin, Yun Jiang, Yunyun Wang, Yushi Wang, Yutian Liu, Zach  
435 Stubenvoll, Zehao Dou, Zheng Wu, and Zhigang Wang. Openai gpt-5 system card, 2025. URL  
436 <https://arxiv.org/abs/2601.03267>.

437 [24] An Yang, Anpeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang  
438 Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu,  
439 Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin  
440 Yang, Jiayi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang,  
441 Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui  
442 Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang  
443 Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger  
444 Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan  
445 Qiu. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.

446 [25] Chicheng Zhang, Alekh Agarwal, Hal Daumé III, John Langford, and Sahand N Negahban.  
447 Warm-starting contextual bandits: Robustly combining supervised and bandit feedback. *arXiv*  
448 *preprint arXiv:1901.00301*, 2019.

449 [26] Wenlin Zhang, Chuhan Wu, Xiangyang Li, Yuhao Wang, Kuicai Dong, Yichao Wang, Xinyi  
450 Dai, Xiangyu Zhao, Huifeng Guo, and Ruiming Tang. Llmtrerec: Unleashing the power of  
451 large language models for cold-start recommendations. In *Proceedings of the 31st International*  
452 *Conference on Computational Linguistics*, pages 886–896, 2025.

453 [27] Dongruo Zhou, Lihong Li, and Quanquan Gu. Neural contextual bandits with ucb-based  
454 exploration. In *International conference on machine learning*, pages 11492–11502. PMLR,  
455 2020.

- 456 [28] Yinglun Zhu, Dylan J Foster, John Langford, and Paul Mineiro. Contextual bandits with  
457 large action spaces: Made practical. In *International Conference on Machine Learning*, pages  
458 27428–27453. PMLR, 2022.

---

**Algorithm 1** Prior-centered LinUCB

---

**Require:** Prior parameter  $\theta_0 \in \mathbb{R}^D$ , regularization  $\lambda > 0$ , confidence level  $\delta \in (0, 1)$

```
1: Initialize  $A_0 \leftarrow \lambda I_D$ 
2: Initialize  $b_0 \leftarrow \lambda \theta_0$ 
3: for  $t = 1, 2, \dots, T$  do
4:   Observe context  $c_t$  and available arm set  $\mathcal{A}_t$ 
5:   for all  $a \in \mathcal{A}_t$  do
6:     Construct joint feature  $z_{t,a} \leftarrow \Psi(\psi(c_t), \psi(s_a))$ 
7:   end for
8:   Compute  $\hat{\theta}_{t-1} \leftarrow A_{t-1}^{-1} b_{t-1}$ 
9:   Select
      
$$a_t \leftarrow \arg \max_{a \in \mathcal{A}_t} \left[ z_{t,a}^\top \hat{\theta}_{t-1} + \beta_{t-1}(\delta) \sqrt{z_{t,a}^\top A_{t-1}^{-1} z_{t,a}} \right]$$

10:  Observe reward  $y_t$ 
11:  Update
      
$$A_t \leftarrow A_{t-1} + z_{t,a_t} z_{t,a_t}^\top$$

12:  Update
      
$$b_t \leftarrow b_{t-1} + y_t z_{t,a_t}$$

13: end for
```

---

## 459 A Theoretical Derivation

460 For the theoretical derivation, we analyze the sequence generated by an arbitrary adaptive policy.  
461 The proof below does not depend on the internal implementation of Algorithm 1, except when we  
462 explicitly impose the UCB selection rule for the instantaneous regret bound.

463 **Assumption A.1** (Adaptive bandit protocol). At each round  $t$ , before observing the reward  $y_t$ , the  
464 learner observes the context  $c_t$  and a nonempty finite arm set  $\mathcal{A}_t$ . For each  $a \in \mathcal{A}_t$ , the joint feature is

$$z_{t,a} = \Psi(\psi(c_t), \psi(s_a)) \in \mathbb{R}^D.$$

465 The learner selects  $a_t \in \mathcal{A}_t$  before observing  $y_t$ .

466 Let  $\mathcal{H}_t$  denote the pre-reward sigma-field at round  $t$ , containing all past observations and the current  
467 context, arm set, and selected arm:

$$\mathcal{H}_t = \sigma(\theta_0, c_1, \mathcal{A}_1, a_1, y_1, \dots, c_{t-1}, \mathcal{A}_{t-1}, a_{t-1}, y_{t-1}, c_t, \mathcal{A}_t, a_t).$$

468 Then  $x_t := z_{t,a_t}$  is  $\mathcal{H}_t$ -measurable.

469 **Assumption A.2** (Linear reward model). There exists a fixed but unknown parameter

$$\theta^* \in \mathbb{R}^D$$

470 such that for every context-arm pair,

$$\mu(z_{t,a}) = z_{t,a}^\top \theta^*.$$

471 The optimal arm at round  $t$  is

$$a_t^* \in \arg \max_{a \in \mathcal{A}_t} \mu(z_{t,a}),$$

472 and the optimal expected reward is

$$\mu_t^* = \mu(z_{t,a_t^*}).$$

473 **Assumption A.3** (Bounded features). There exists  $L > 0$  such that

$$\|z_{t,a}\|_2 \leq L$$

474 for all  $t$  and all  $a \in \mathcal{A}_t$ . Equivalently,

$$\|z_{t,a}\|_{(\lambda I_D)^{-1}}^2 \leq \frac{L^2}{\lambda}.$$

---

**Algorithm 2** LLM-PriorCB

---

**Require:** Offline textual dataset  $\mathcal{D}_0 = (\mathcal{C}_0, \mathcal{A}_0, \{s_a\}_{a \in \mathcal{A}_0})$ , LLM  $\mathcal{M}$ , number of LLM samples  $M$ , embedding map  $\psi$ , joint feature map  $\Psi$ , prior encoder  $f_\phi$ , regularization  $\lambda > 0$ , confidence level  $\delta \in (0, 1)$ , horizon  $T$

**Stage I: Offline prior construction**

```
1: for all  $c_i \in \mathcal{C}_0$  do
2:    $z_i \leftarrow \psi(c_i)$ 
3:   for all  $a \in \mathcal{A}_0$  do
4:      $z_a \leftarrow \psi(s_a), z_{i,a} \leftarrow \Psi(z_i, z_a)$ 
5:     Query  $\mathcal{M}$   $M$  times with  $(c_i, s_a)$  to obtain  $\{\hat{\mu}_{i,a}^{(m)}\}_{m=1}^M$ 
6:      $\hat{\mu}_{i,a} \leftarrow \frac{1}{M} \sum_{m=1}^M \hat{\mu}_{i,a}^{(m)}$ 
7:   end for
8: end for
9: Train  $f_\phi$  by minimizing Eq. equation 1
10: for all  $a \in \mathcal{A}_0$  do
11:    $\theta_{0,a} \leftarrow f_\phi(\psi(s_a))$ 
12: end for
```

**Stage II: Online learning**

```
13: Initialize the set of initialized arms  $\mathcal{I} \leftarrow \emptyset$ 
14: for  $t = 1, 2, \dots, T$  do
15:   Observe context  $c_t$  and available arm set  $\mathcal{A}_t$ 
16:    $z_t \leftarrow \psi(c_t)$ 
17:   for all  $a \in \mathcal{A}_t$  do
18:      $z_a \leftarrow \psi(s_a), z_{t,a} \leftarrow \Psi(z_t, z_a)$ 
19:     if  $a \notin \mathcal{I}$  then
20:        $\theta_{0,a} \leftarrow f_\phi(z_a)$ 
21:        $A_a \leftarrow \lambda I_D$ 
22:        $b_a \leftarrow \lambda \theta_{0,a}$ 
23:        $\mathcal{I} \leftarrow \mathcal{I} \cup \{a\}$ 
24:     end if
25:      $\hat{\theta}_a \leftarrow A_a^{-1} b_a$ 
26:      $U_t(a) \leftarrow z_{t,a}^\top \hat{\theta}_a + \beta_{t-1}^{\text{dis}}(\delta) \sqrt{z_{t,a}^\top A_a^{-1} z_{t,a}}$ 
27:   end for
28:   Select  $a_t \leftarrow \arg \max_{a \in \mathcal{A}_t} U_t(a)$ 
29:   Observe reward  $y_t$ 
30:   Update  $A_{a_t} \leftarrow A_{a_t} + z_{t,a_t} z_{t,a_t}^\top$ 
31:   Update  $b_{a_t} \leftarrow b_{a_t} + y_t z_{t,a_t}$ 
32: end for
```

---

475 **Assumption A.4** (Conditional sub-Gaussian noise). The observed reward satisfies

$$y_t = x_t^\top \theta^* + \eta_t.$$

476 The noise satisfies

$$\mathbb{E}[\eta_t \mid \mathcal{H}_t] = 0,$$

477 and for all  $s \in \mathbb{R}$ ,

$$\mathbb{E}[\exp(s\eta_t) \mid \mathcal{H}_t] \leq \exp\left(\frac{s^2 R^2}{2}\right).$$

478 **Assumption A.5** (Bounded prior misspecification). There exists a known constant  $S_\Delta > 0$  such that

$$\|\theta^* - \theta_0\|_2 \leq S_\Delta.$$

479 We write

$$\Delta^* := \theta^* - \theta_0.$$

480 **Definition A.6.** For any adaptively chosen sequence  $\{a_t\}_{t \geq 1}$ , define

$$A_t = \lambda I_D + \sum_{i=1}^t x_i x_i^\top, \quad b_t = \lambda \theta_0 + \sum_{i=1}^t y_i x_i,$$

481 where  $x_i = z_{i,a_i}$ .

482 The prior-centered ridge estimator is

$$\hat{\theta}_t = A_t^{-1} b_t.$$

483 Since  $\lambda > 0$ ,

$$A_t \succeq \lambda I_D \succ 0,$$

484 so  $A_t$  is symmetric positive definite and  $A_t^{-1}$  is well-defined.

485 **Lemma A.7** (Equivalence to residual ridge regression). *Define*

$$\hat{\Delta}_t := \hat{\theta}_t - \theta_0, \quad \Delta^* := \theta^* - \theta_0.$$

486 *Then*

$$\hat{\Delta}_t = A_t^{-1} \sum_{i=1}^t x_i (y_i - x_i^\top \theta_0).$$

487 *Equivalently,  $\hat{\Delta}_t$  is the ridge estimator obtained by regressing the residualized rewards*

$$\tilde{y}_i := y_i - x_i^\top \theta_0$$

488 *on  $x_i$ , with regularization centered at zero:*

$$\hat{\Delta}_t = \arg \min_{\Delta \in \mathbb{R}^D} \left\{ \sum_{i=1}^t (\tilde{y}_i - x_i^\top \Delta)^2 + \lambda \|\Delta\|_2^2 \right\}.$$

489 *Proof.* By definition,

$$\hat{\Delta}_t = \hat{\theta}_t - \theta_0 = A_t^{-1} b_t - \theta_0 = A_t^{-1} (b_t - A_t \theta_0).$$

490 Now,

$$\begin{aligned} b_t - A_t \theta_0 &= \lambda \theta_0 + \sum_{i=1}^t y_i x_i - \left( \lambda I_D + \sum_{i=1}^t x_i x_i^\top \right) \theta_0 \\ &= \sum_{i=1}^t y_i x_i - \sum_{i=1}^t x_i x_i^\top \theta_0 \\ &= \sum_{i=1}^t x_i (y_i - x_i^\top \theta_0). \end{aligned}$$

491 Therefore,

$$\hat{\Delta}_t = A_t^{-1} \sum_{i=1}^t x_i (y_i - x_i^\top \theta_0).$$

492 The optimization form follows from the normal equation of ridge regression:

$$\left( \lambda I_D + \sum_{i=1}^t x_i x_i^\top \right) \Delta = \sum_{i=1}^t x_i \tilde{y}_i.$$

493 The left-hand matrix is  $A_t$ , so the unique minimizer is exactly

$$\hat{\Delta}_t = A_t^{-1} \sum_{i=1}^t x_i \tilde{y}_i.$$

494 □

495 **Lemma A.8** (Self-normalized bound). *Under Assumptions A.1 and A.4, with probability at least*  
496  *$1 - \delta$ , simultaneously for all  $t \geq 0$ ,*

$$\left\| \sum_{i=1}^t x_i \eta_i \right\|_{A_t^{-1}} \leq R \sqrt{2 \log \frac{\det(A_t)^{1/2}}{\det(\lambda I_D)^{1/2} \delta}}.$$

497 *Proof.* This is the standard self-normalized concentration inequality for vector-valued martingales.  
 498 In our notation,  $x_t$  is measurable before observing  $y_t$ , and  $\eta_t$  is conditionally  $R$ -sub-Gaussian given  
 499 the pre-reward sigma-field  $\mathcal{H}_t$ . Therefore, the process

$$\sum_{i=1}^t x_i \eta_i$$

500 satisfies the conditions of the self-normalized martingale inequality with regularization matrix  $\lambda I_D$ .  
 501 For more details, see Theorem 1 in [1].  $\square$

502 **Lemma A.9** (Confidence ellipsoid). *Suppose Assumptions A.1, A.2, A.4, and A.5 hold. Define*

$$\beta_t(\delta) = R \sqrt{2 \log \frac{\det(A_t)^{1/2}}{\det(\lambda I_D)^{1/2} \delta}} + \sqrt{\lambda} S_\Delta.$$

503 *Then, with probability at least  $1 - \delta$ , simultaneously for all  $t \geq 0$ ,*

$$\|\hat{\theta}_t - \theta^*\|_{A_t} \leq \beta_t(\delta).$$

504 *Proof.* By Lemma A.7,

$$\hat{\Delta}_t = A_t^{-1} \sum_{i=1}^t x_i (y_i - x_i^\top \theta_0).$$

505 Using the linear model,

$$y_i = x_i^\top \theta^* + \eta_i,$$

506 we have

$$y_i - x_i^\top \theta_0 = x_i^\top (\theta^* - \theta_0) + \eta_i = x_i^\top \Delta^* + \eta_i.$$

507 Therefore,

$$\begin{aligned} \hat{\Delta}_t &= A_t^{-1} \sum_{i=1}^t x_i (x_i^\top \Delta^* + \eta_i) \\ &= A_t^{-1} \left( \sum_{i=1}^t x_i x_i^\top \Delta^* + \sum_{i=1}^t x_i \eta_i \right). \end{aligned}$$

508 Since

$$A_t = \lambda I_D + \sum_{i=1}^t x_i x_i^\top,$$

509 we have

$$\sum_{i=1}^t x_i x_i^\top = A_t - \lambda I_D.$$

510 Thus,

$$\begin{aligned} \hat{\Delta}_t &= A_t^{-1} \left( (A_t - \lambda I_D) \Delta^* + \sum_{i=1}^t x_i \eta_i \right) \\ &= \Delta^* - \lambda A_t^{-1} \Delta^* + A_t^{-1} \sum_{i=1}^t x_i \eta_i. \end{aligned}$$

511 Hence,

$$\hat{\Delta}_t - \Delta^* = A_t^{-1} \left( \sum_{i=1}^t x_i \eta_i - \lambda \Delta^* \right).$$

512 Since

$$\hat{\Delta}_t - \Delta^* = \hat{\theta}_t - \theta^*,$$

513 we obtain

$$\begin{aligned}
\|\hat{\theta}_t - \theta^*\|_{A_t} &= \left\| A_t^{-1} \left( \sum_{i=1}^t x_i \eta_i - \lambda \Delta^* \right) \right\|_{A_t} \\
&= \left\| \sum_{i=1}^t x_i \eta_i - \lambda \Delta^* \right\|_{A_t^{-1}} \\
&\leq \left\| \sum_{i=1}^t x_i \eta_i \right\|_{A_t^{-1}} + \|\lambda \Delta^*\|_{A_t^{-1}}.
\end{aligned}$$

514 For the second term, since  $A_t \succeq \lambda I_D$ , we have

$$A_t^{-1} \preceq \frac{1}{\lambda} I_D.$$

515 Therefore,

$$\begin{aligned}
\|\lambda \Delta^*\|_{A_t^{-1}}^2 &= \lambda^2 \Delta^{*\top} A_t^{-1} \Delta^* \\
&\leq \lambda^2 \Delta^{*\top} \left( \frac{1}{\lambda} I_D \right) \Delta^* \\
&= \lambda \|\Delta^*\|_2^2 \\
&\leq \lambda S_\Delta^2.
\end{aligned}$$

516 Thus,

$$\|\lambda \Delta^*\|_{A_t^{-1}} \leq \sqrt{\lambda} S_\Delta.$$

517 By Lemma A.8, with probability at least  $1 - \delta$ , simultaneously for all  $t \geq 0$ ,

$$\left\| \sum_{i=1}^t x_i \eta_i \right\|_{A_t^{-1}} \leq R \sqrt{2 \log \frac{\det(A_t)^{1/2}}{\det(\lambda I_D)^{1/2} \delta}}.$$

518 Combining the two inequalities gives

$$\|\hat{\theta}_t - \theta^*\|_{A_t} \leq R \sqrt{2 \log \frac{\det(A_t)^{1/2}}{\det(\lambda I_D)^{1/2} \delta}} + \sqrt{\lambda} S_\Delta = \beta_t(\delta).$$

519

□

520 **Definition A.10.** For a fixed horizon  $T$ , define the confidence event

$$\mathcal{E}_T(\delta) := \left\{ \forall t \in \{0, \dots, T\}, \|\hat{\theta}_t - \theta^*\|_{A_t} \leq \beta_t(\delta) \right\}.$$

521 By Lemma A.9,

$$\Pr(\mathcal{E}_T(\delta)) \geq 1 - \delta.$$

522 **Lemma A.11** (Instantaneous regret under UCB selection). *On the event  $\mathcal{E}_T(\delta)$ , fix a round  $t$ . Suppose*  
523 *the selected arm  $a_t$  satisfies*

$$a_t \in \arg \max_{a \in \mathcal{A}_t} \left[ z_{t,a}^\top \hat{\theta}_{t-1} + \beta_{t-1}(\delta) \sqrt{z_{t,a}^\top A_{t-1}^{-1} z_{t,a}} \right].$$

524 *Then the instantaneous regret*

$$r_t := \mu_t^* - \mu(z_{t,a_t})$$

525 *satisfies*

$$r_t \leq 2\beta_{t-1}(\delta) \sqrt{z_{t,a_t}^\top A_{t-1}^{-1} z_{t,a_t}}.$$

526 *Consequently, with probability at least  $1 - \delta$ , the above bound holds simultaneously for all rounds  $t$*   
527 *in which the UCB selection rule is used.*

528 *Proof.* For any arm  $a \in \mathcal{A}_t$ , define

$$w_{t,a} := \sqrt{z_{t,a}^\top A_{t-1}^{-1} z_{t,a}}.$$

529 By Cauchy's inequality under the  $A_{t-1}$ -norm,

$$\begin{aligned} \left| z_{t,a}^\top (\hat{\theta}_{t-1} - \theta^*) \right| &= \left| \left( A_{t-1}^{-1/2} z_{t,a} \right)^\top \left( A_{t-1}^{1/2} (\hat{\theta}_{t-1} - \theta^*) \right) \right| \\ &\leq \left\| A_{t-1}^{-1/2} z_{t,a} \right\|_2 \left\| A_{t-1}^{1/2} (\hat{\theta}_{t-1} - \theta^*) \right\|_2 \\ &= w_{t,a} \|\hat{\theta}_{t-1} - \theta^*\|_{A_{t-1}} \\ &\leq \beta_{t-1}(\delta) w_{t,a}. \end{aligned}$$

530 Therefore, for every  $a \in \mathcal{A}_t$ , we obtain two upper-confidence inequalities:

$$z_{t,a}^\top \theta^* \leq z_{t,a}^\top \hat{\theta}_{t-1} + \beta_{t-1}(\delta) w_{t,a},$$

531 and

$$z_{t,a}^\top \hat{\theta}_{t-1} \leq z_{t,a}^\top \theta^* + \beta_{t-1}(\delta) w_{t,a}.$$

532 Let

$$U_t(a) := z_{t,a}^\top \hat{\theta}_{t-1} + \beta_{t-1}(\delta) w_{t,a}.$$

533 By the first upper-confidence inequality,

$$\mu(z_{t,a_t^*}) = z_{t,a_t^*}^\top \theta^* \leq U_t(a_t^*).$$

534 Since  $a_t$  maximizes the UCB score,

$$U_t(a_t^*) \leq U_t(a_t).$$

535 Thus,

$$\mu(z_{t,a_t^*}) \leq U_t(a_t).$$

536 Expanding  $U_t(a_t)$  gives

$$\mu(z_{t,a_t^*}) \leq z_{t,a_t}^\top \hat{\theta}_{t-1} + \beta_{t-1}(\delta) w_{t,a_t}.$$

537 Using the second confidence inequality for the selected arm  $a_t$ ,

$$z_{t,a_t}^\top \hat{\theta}_{t-1} \leq z_{t,a_t}^\top \theta^* + \beta_{t-1}(\delta) w_{t,a_t}.$$

538 Therefore,

$$\begin{aligned} \mu(z_{t,a_t^*}) &\leq z_{t,a_t}^\top \theta^* + 2\beta_{t-1}(\delta) w_{t,a_t} \\ &= \mu(z_{t,a_t}) + 2\beta_{t-1}(\delta) w_{t,a_t}. \end{aligned}$$

539 Rearranging,

$$r_t = \mu(z_{t,a_t^*}) - \mu(z_{t,a_t}) \leq 2\beta_{t-1}(\delta) w_{t,a_t}.$$

540 Substituting the definition of  $w_{t,a_t}$  yields

$$r_t \leq 2\beta_{t-1}(\delta) \sqrt{z_{t,a_t}^\top A_{t-1}^{-1} z_{t,a_t}},$$

541 □

542 **Theorem 5.1** (Cumulative regret analysis). *Suppose Assumptions A.1, A.2, A.3, A.4, and A.5 hold.*

543 *Suppose that at every round  $t = 1, \dots, T$ , the selected arm  $a_t$  satisfies*

$$a_t \in \arg \max_{a \in \mathcal{A}_t} \left[ z_{t,a}^\top \hat{\theta}_{t-1} + \beta_{t-1}(\delta) \sqrt{z_{t,a}^\top A_{t-1}^{-1} z_{t,a}} \right].$$

544 *Then, with probability at least  $1 - \delta$ ,*

$$R_T := \sum_{t=1}^T r_t \leq 2\beta_T(\delta) \sqrt{2T c_\lambda \log \frac{\det(A_T)}{\det(\lambda I_D)}},$$

545 *where*

$$\beta_t(\delta) := R \sqrt{2 \log \frac{\det(A_t)^{1/2}}{\det(\lambda I_D)^{1/2} \delta}} + \sqrt{\lambda} S_\Delta$$

546 *and*

$$c_\lambda := \max \left\{ 1, \frac{L^2}{\lambda} \right\}.$$

547 *Proof.* We prove the bound on the event  $\mathcal{E}_T(\delta)$ , which holds with probability at least  $1 - \delta$ . Define

$$w_t := x_t^\top A_{t-1}^{-1} x_t = \|x_t\|_{A_{t-1}^{-1}}^2.$$

548 We prove the regret bound on the event  $\mathcal{E}_T(\delta)$ . By Lemma A.11, for every  $t = 1, \dots, T$ ,

$$r_t \leq 2\beta_{t-1}(\delta) \sqrt{x_t^\top A_{t-1}^{-1} x_t} = 2\beta_{t-1}(\delta) \sqrt{w_t}.$$

549 Since

$$A_t = A_{t-1} + x_t x_t^\top,$$

550 we have

$$A_t \succeq A_{t-1}.$$

551 Moreover, by the matrix determinant lemma,

$$\det(A_t) = \det(A_{t-1}) (1 + x_t^\top A_{t-1}^{-1} x_t) = \det(A_{t-1})(1 + w_t).$$

552 Hence  $\det(A_t)$  is nondecreasing in  $t$ , and therefore  $\beta_t$  is also nondecreasing. Thus,

$$\beta_{t-1}(\delta) \leq \beta_T(\delta)$$

553 for all  $t \leq T$ .

554 Therefore,

$$R_T = \sum_{t=1}^T r_t \leq 2\beta_T(\delta) \sum_{t=1}^T \sqrt{w_t}.$$

555 By Cauchy's inequality,

$$\sum_{t=1}^T \sqrt{w_t} \leq \sqrt{T \sum_{t=1}^T w_t}.$$

556 Thus,

$$R_T \leq 2\beta_T(\delta) \sqrt{T \sum_{t=1}^T w_t}.$$

557 It remains to upper bound  $\sum_{t=1}^T w_t$ . By Assumption A.3,

$$\|x_t\|_2 \leq L.$$

558 Since

$$A_{t-1} \succeq \lambda I_D,$$

559 we have

$$A_{t-1}^{-1} \preceq \frac{1}{\lambda} I_D.$$

560 Therefore,

$$w_t = x_t^\top A_{t-1}^{-1} x_t \leq \frac{\|x_t\|_2^2}{\lambda} \leq \frac{L^2}{\lambda}.$$

561 Since

$$c_\lambda = \max \left\{ 1, \frac{L^2}{\lambda} \right\},$$

562 we have for every  $t$ ,

$$w_t \leq c_\lambda \min\{1, w_t\}.$$

563 Indeed, if  $w_t \leq 1$ , then

$$w_t \leq c_\lambda w_t = c_\lambda \min\{1, w_t\},$$

564 and if  $w_t > 1$ , then

$$w_t \leq \frac{L^2}{\lambda} \leq c_\lambda = c_\lambda \min\{1, w_t\}.$$

565 Now use the standard elliptical potential argument. For every  $u \geq 0$ ,

$$\min\{1, u\} \leq 2 \log(1 + u).$$

566 Therefore,

$$\sum_{t=1}^T w_t \leq c_\lambda \sum_{t=1}^T \min\{1, w_t\} \leq 2c_\lambda \sum_{t=1}^T \log(1 + w_t).$$

567 Using the determinant identity above,

$$\sum_{t=1}^T \log(1 + w_t) = \sum_{t=1}^T \log \frac{\det(A_t)}{\det(A_{t-1})} = \log \frac{\det(A_T)}{\det(A_0)}.$$

568 Since

$$A_0 = \lambda I_D,$$

569 we obtain

$$\sum_{t=1}^T w_t \leq 2c_\lambda \log \frac{\det(A_T)}{\det(\lambda I_D)}.$$

570 Substituting this into the previous regret bound gives

$$R_T \leq 2\beta_T \sqrt{T \cdot 2c_\lambda \log \frac{\det(A_T)}{\det(\lambda I_D)}}.$$

571 Hence,

$$R_T \leq 2\beta_T \sqrt{2Tc_\lambda \log \frac{\det(A_T)}{\det(\lambda I_D)}}.$$

572 Since the event  $\mathcal{E}_T(\delta)$  holds with probability at least  $1 - \delta$ , the theorem follows.  $\square$

573 **Corollary 5.2** (Cumulative regret analysis for LLM-PriorCB). *Suppose the assumptions in Theorem 5.1 hold. Consider the disjoint linear model*

$$\mu(z_{t,a}) = z_{t,a}^\top \theta_a^*, \quad z_{t,a} \in \mathbb{R}^D, \quad a \in [K],$$

575 *with arm-specific prior parameters  $\theta_{0,a} \in \mathbb{R}^D$ . Then the arm-wise LLM-PriorCB rule induced by this block representation satisfies, with probability at least  $1 - \delta$ ,*

$$R_T \leq 2\beta_T^{\text{dis}}(\delta) \sqrt{2Tc_\lambda \log \frac{\prod_{a=1}^K \det(A_{T,a})}{\lambda^{Kd}}},$$

577 *where*

$$A_{t,a} := \lambda I_d + \sum_{i=1}^t \mathbf{1}\{a_i = a\} z_{i,a} z_{i,a}^\top$$

578 *and*

$$\beta_t^{\text{dis}}(\delta) := R \sqrt{2 \log \frac{\prod_{a=1}^K \det(A_{t,a})^{1/2}}{\lambda^{Kd/2\delta}}} + \sqrt{\lambda} S_\Delta.$$

579 *Proof.* Let  $e_a \in \mathbb{R}^K$  denote the  $a$ -th standard basis vector and define the block-lifted feature

$$\tilde{z}_{t,a} = e_a \otimes z_{t,a} \in \mathbb{R}^{Kd}.$$

580 Also define the stacked true and prior parameters

$$\Theta^* = (\theta_1^{*\top}, \dots, \theta_K^{*\top})^\top, \quad \Theta_0 = (\theta_{0,1}^\top, \dots, \theta_{0,K}^\top)^\top.$$

581 Then, for every arm  $a$ ,

$$\tilde{z}_{t,a}^\top \Theta^* = z_{t,a}^\top \theta_a^* = \mu(z_{t,a}).$$

582 Thus the disjoint model is a special case of the linear model in Theorem 5.1, with feature  $\tilde{z}_{t,a}$  and  
 583 parameter  $\Theta^*$ .

584 Define the global prior-centered design matrix and response vector in the block representation by

$$\tilde{A}_t = \lambda I_{Kd} + \sum_{i=1}^t \tilde{z}_{i,a_i} \tilde{z}_{i,a_i}^\top,$$

585

$$\tilde{b}_t = \lambda \Theta_0 + \sum_{i=1}^t y_i \tilde{z}_{i,a_i}.$$

586 Since  $\tilde{z}_{i,a_i}$  has nonzero entries only in the block corresponding to arm  $a_i$ , we have

$$\tilde{A}_t = \text{diag}(A_{t,1}, \dots, A_{t,K}),$$

587 and

$$\tilde{b}_t = (b_{t,1}^\top, \dots, b_{t,K}^\top)^\top.$$

588 Therefore,

$$\tilde{A}_t^{-1} \tilde{b}_t = (\hat{\theta}_{t,1}^\top, \dots, \hat{\theta}_{t,K}^\top)^\top.$$

589 Moreover, for any candidate arm  $a$ ,

$$\tilde{z}_{t,a}^\top \tilde{A}_{t-1}^{-1} \tilde{z}_{t,a} = z_{t,a}^\top A_{t-1,a}^{-1} z_{t,a},$$

590 and

$$\tilde{z}_{t,a}^\top \tilde{A}_{t-1}^{-1} \tilde{b}_{t-1} = z_{t,a}^\top \hat{\theta}_{t-1,a}.$$

591 Hence the UCB rule in Theorem 5.1 coincides with the arm-wise PriorCB-Disjoint rule.

592 Finally,

$$\det(\tilde{A}_t) = \prod_{a=1}^K \det(A_{t,a}), \quad \det(\lambda I_{Kd}) = \lambda^{Kd}.$$

593 Also,

$$\|\Theta^* - \Theta_0\|_2 = \left( \sum_{a=1}^K \|\theta_a^* - \theta_{0,a}\|_2^2 \right)^{1/2}.$$

594 Substituting these identities into Theorem 5.1 gives

$$R_T \leq 2\beta_T^{\text{dis}}(\delta) \sqrt{2T c_\lambda \log \frac{\prod_{a=1}^K \det(A_{T,a})}{\lambda^{Kd}}},$$

595 which proves the claim. □

596 **B Ablation Studies on MovieLens**

597 **B.1 CBLI**

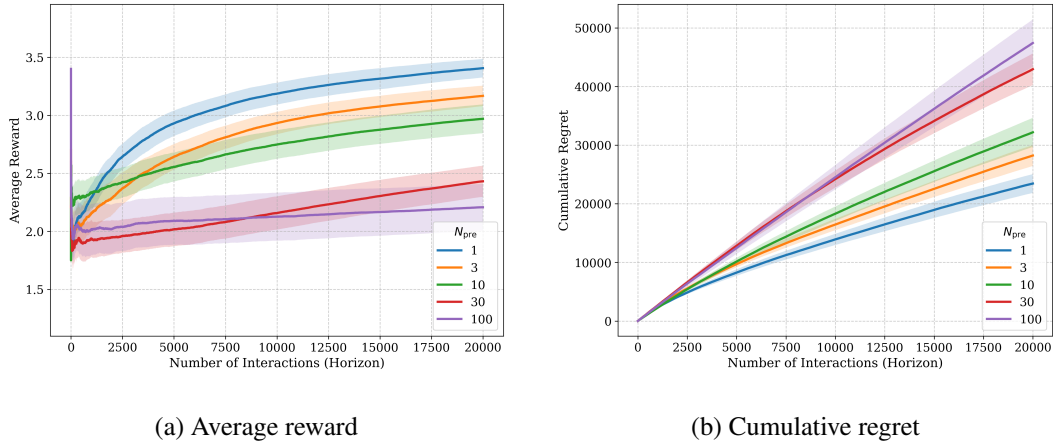


Figure 5: Average reward and cumulative regret of CBLI under oracle regime with varying numbers of pretraining steps. Here, the total pretraining budget scales as  $N_{\text{pre}} \times |\mathcal{C}_0| \times |\mathcal{A}_0|$ . Interestingly, increasing  $N_{\text{pre}}$  leads to degraded online performance. This is because excessive pretraining induces an overconfident prior, causing the UCB confidence intervals to collapse and the agent to overly exploit estimated rewards as if they were ground truth, thereby suppressing exploration.

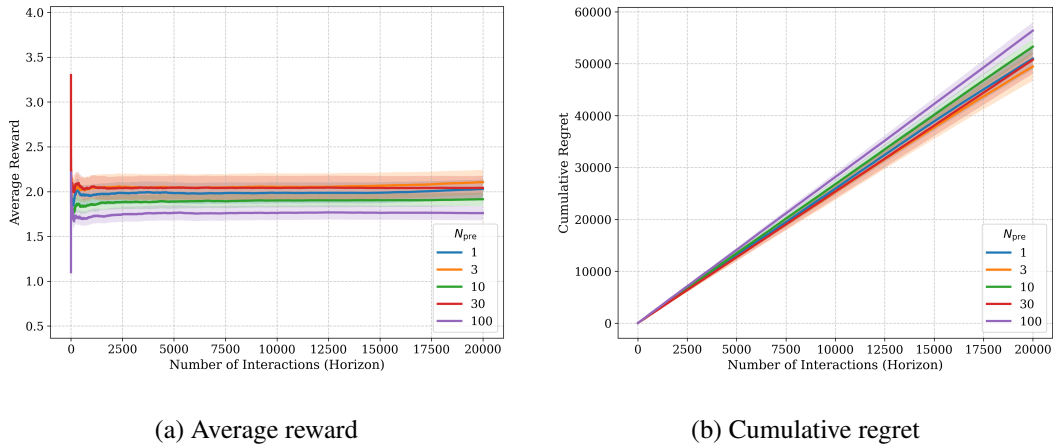
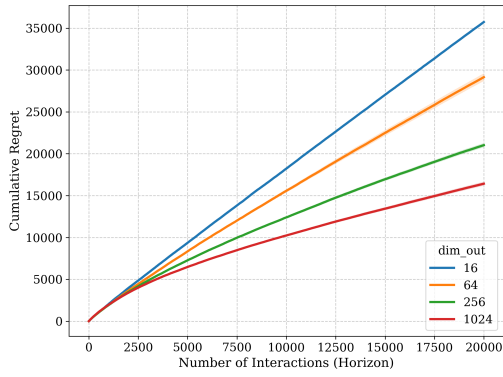
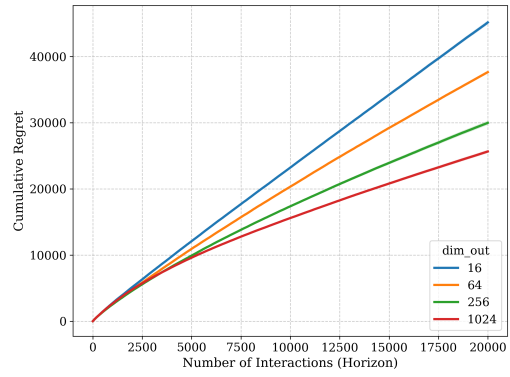


Figure 6: Average reward and cumulative regret of CBLI under text-only regime with varying numbers of pretraining steps. Here, the total pretraining budget scales as  $N_{\text{pre}} \times |\mathcal{C}_0| \times |\mathcal{A}_0|$ . Interestingly, increasing  $N_{\text{pre}}$  cannot lead to improve online performance. This is because excessive pretraining induces an overconfident prior, causing the UCB confidence intervals to collapse and the agent to overly exploit estimated rewards as if they were ground truth, thereby suppressing exploration.

598 **B.2 Output Dimension**



(a) LLM-PriorCB with  $K = 10$



(b) LLM-PriorCB with  $K = 30$

Figure 7: Cumulative regret of our algorithm with varying feature vector dimensions. We empirically find that higher-dimensional representations achieve lower cumulative regret across the horizon regardless of the number of arms. We conjecture that increased dimensionality provides greater representational capacity to capture relevant contextual information, thereby improving reward estimation during online learning.

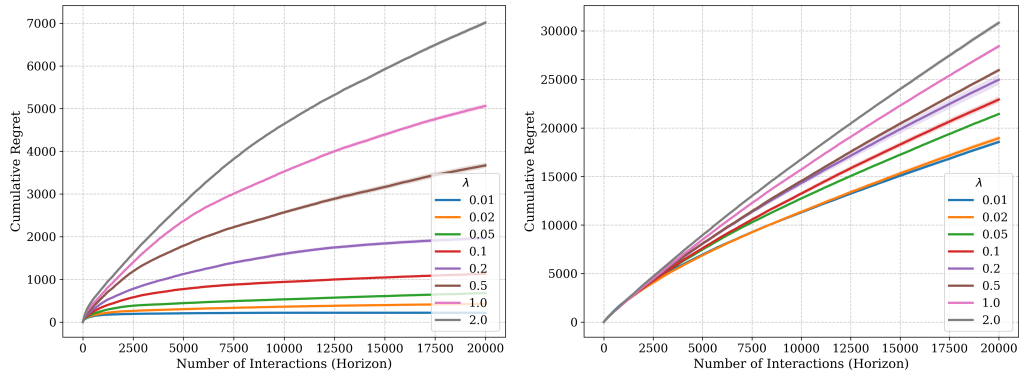
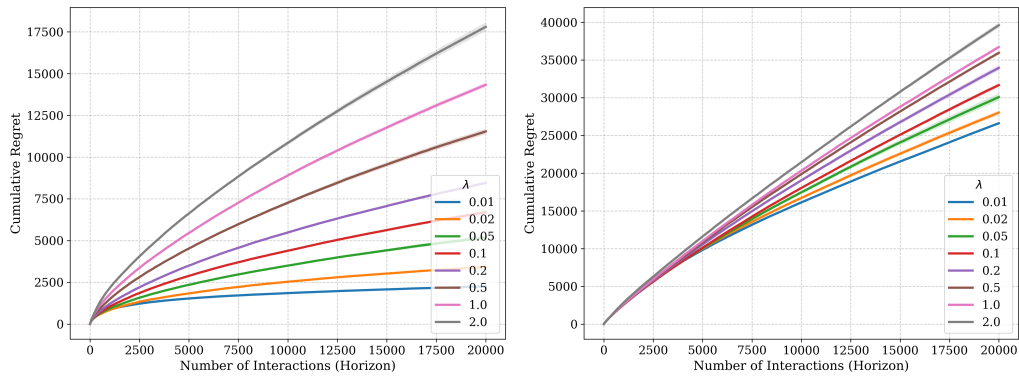
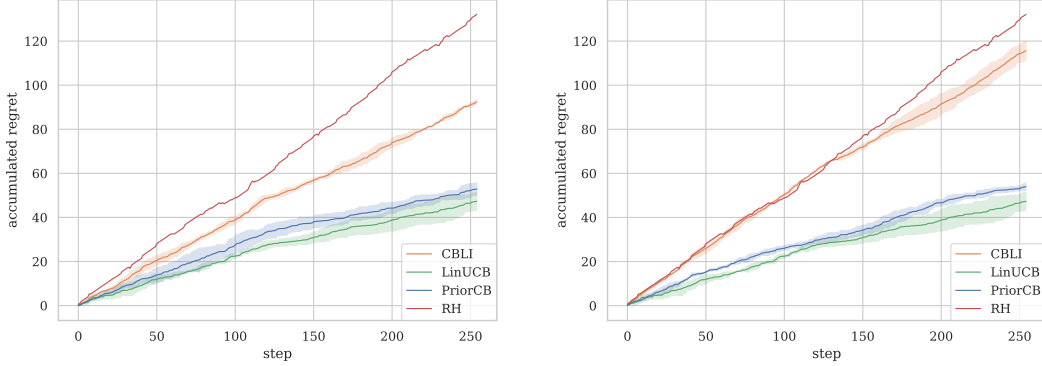
(a) Our algorithm with  $K = 10$ (b) Our algorithm with  $K = 30$ 

Figure 8: We report cumulative regret of LLM-PriorCB for  $\lambda \in [0.01, 2.0]$  under (a) the oracle regime and (b) the text-only regime. Smaller  $\lambda$  consistently yields lower regret, while larger  $\lambda$  degrades performance, indicating that overly confident priors can suppress exploration and hinder effective online adaptation.

600 **C Ablation Studies on GAIA**

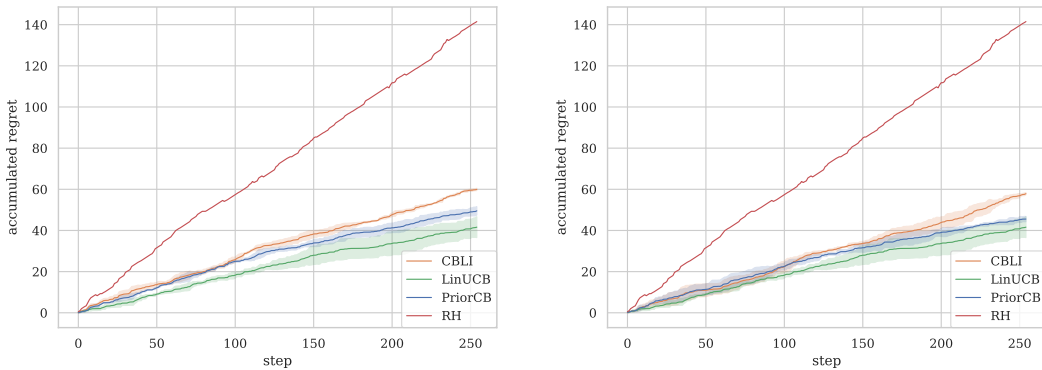
601 **C.1 Cumulative regret**



(a) Cumulative regret of Prior from different eval contexts      (b) Cumulative regret of Prior from same eval contexts

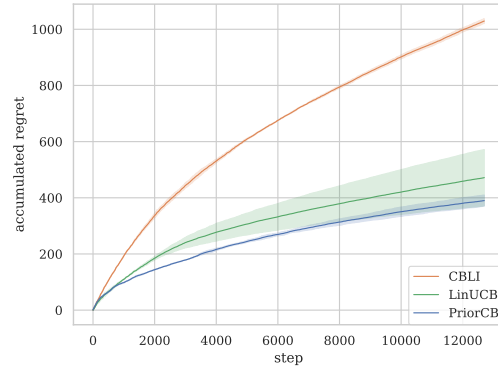
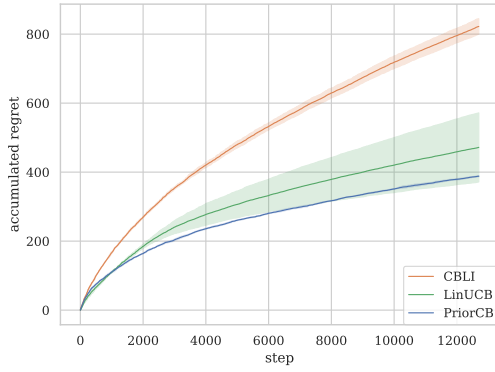
Figure 9: Performance as a function of the interaction horizon when  $K = 3$  using gpt-5.2. Shaded regions indicate  $\pm$  one standard error over 3 independent runs. We use the validation split for evaluation due to answer-label availability. Priors are constructed either from the validation set (same eval contexts) or from the test set (different eval contexts), using no labels.

602 Especially, Figure 9 presents cumulative regret in the early stage of online learning for PriorCB and  
 603 baselines—including RH, where RH results are shown only for the first three epochs due to cost and  
 604 time constraints. The near-linear growth of the cumulative regret of LLM clearly indicates that an LLM  
 605 alone struggles to effectively solve the contextual bandit problem. In contrast, PriorCB—leveraging  
 606 the LLM as a prior—achieves cumulative regret in the early stage that is comparable to LinUCB,  
 607 indicating effective learning. Additional results using GPT-5-mini are provided in the Appendix C.1.  
 608 The average reward in the early stage of online learning is reported in the Appendix C.2.



(a) Cumulative regret of Prior from different eval contexts      (b) Cumulative regret of Prior from same eval contexts

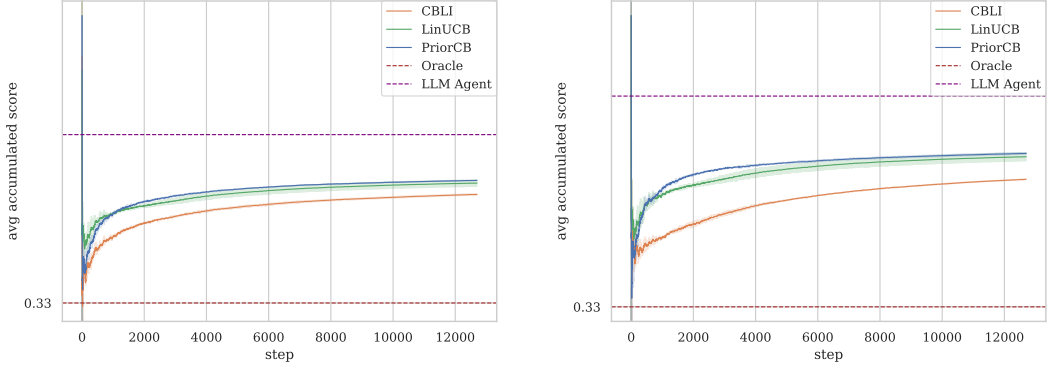
Figure 10: Performance as a function of the interaction horizon when  $K = 3$  using gpt-5-mini. Shaded regions indicate  $\pm$  one standard error over 3 independent runs. We use the validation split for evaluation due to answer-label availability. Priors are constructed either from the validation set (same eval contexts) or from the test set (different eval contexts), using no labels.



(a) Cumulative regret of Prior from different eval contexts

(b) Cumulative regret of Prior from same eval contexts

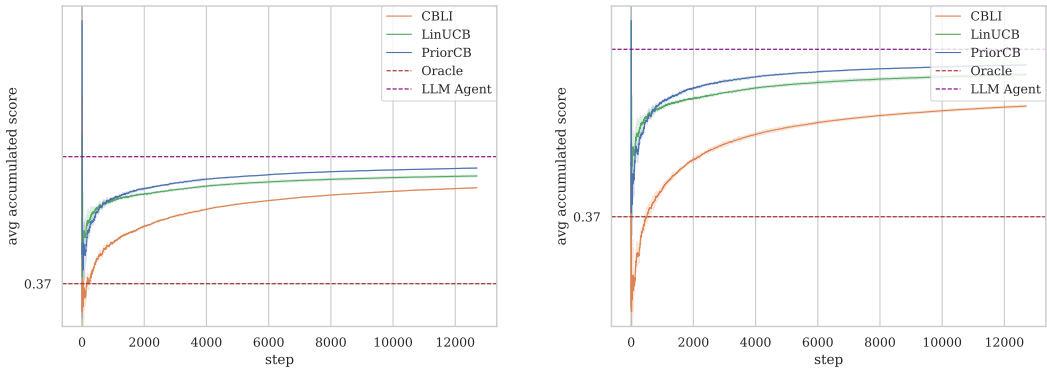
Figure 11: Performance as a function of the interaction horizon when  $K = 3$  using gpt-5-mini. Shaded regions indicate  $\pm$  one standard error over 3 independent runs. We use the validation split for evaluation due to answer-label availability. Priors are constructed either from the validation set (same eval contexts) or from the test set (different eval contexts), using no labels.



(a) Average reward of Prior from different eval contexts

(b) Average reward of Prior from same eval contexts

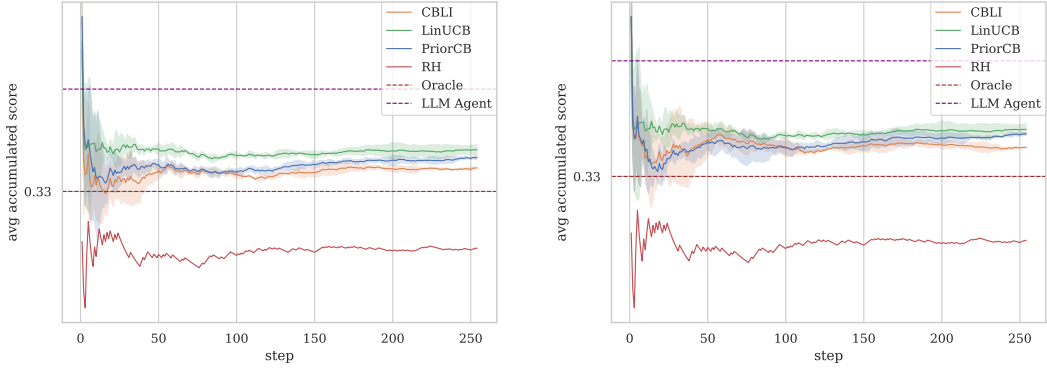
Figure 12: Performance as a function of the interaction horizon when  $K = 3$  using gpt-5-mini with medium reasoning effort. Shaded regions indicate  $\pm$  one standard error over 3 independent runs. We use the validation split for evaluation due to answer-label availability. Priors are constructed either from the validation set (same eval contexts) or from the test set (different eval contexts), using no labels. Oracle denotes the tool choice that achieves best reward, whereas LLM Agent denotes the tool choice selected by the LLM agent.



(a) Average reward of Prior from different eval contexts

(b) Average reward of Prior from same eval contexts

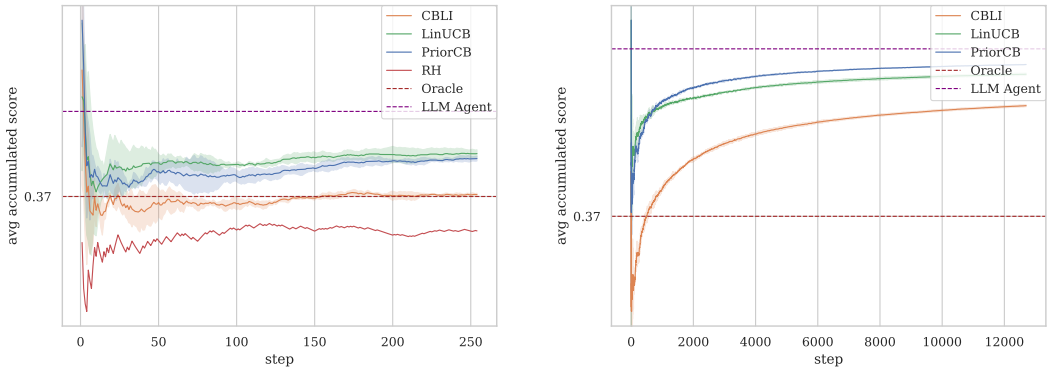
Figure 13: Performance as a function of the interaction horizon when  $K = 3$  using gpt-5.2 with medium reasoning effort. Shaded regions indicate  $\pm$  one standard error over 3 independent runs. We use the validation split for evaluation due to answer-label availability. Priors are constructed either from the validation set (same eval contexts) or from the test set (different eval contexts), using no labels. Oracle denotes the tool choice that achieves best reward, whereas LLM agent denotes the tool choice selected by the LLM agent.



(a) Average reward of Prior from different eval contexts

(b) Average reward of Prior from same eval contexts

Figure 14: Performance as a function of the interaction horizon when  $K = 3$  using gpt-5-mini with medium reasoning effort. Shaded regions indicate  $\pm$  one standard error over 3 independent runs. We use the validation split for evaluation due to answer-label availability. Priors are constructed either from the validation set (same eval contexts) or from the test set (different eval contexts), using no labels. Oracle denotes the tool choice that achieves best reward, whereas LLM agent denotes the tool choice selected by the LLM agent.



(a) Average reward of Prior from different eval contexts

(b) Average reward of Prior from same eval contexts

Figure 15: Performance as a function of the interaction horizon when  $K = 3$  using gpt-5.2 with medium reasoning effort. Shaded regions indicate  $\pm$  one standard error over 3 independent runs. We use the validation split for evaluation due to answer-label availability. Priors are constructed either from the validation set (same eval contexts) or from the test set (different eval contexts), using no labels. Oracle denotes the tool choice that achieves best reward, whereas LLM Agent denotes the tool choice selected by the LLM agent.

## 610 **D Experimental Details**

### 611 **D.1 Movie Lens**

Tool Description: search\_deep

**Purpose:** Run a deeper web search to answer questions that require multiple independent lookups and cross-checking.

**What it does:**

- Generates **three distinct, high-precision search queries**.
- Retrieves results for **all three** queries.
- **Synthesizes** the evidence to produce a final answer.

**Use when (triggers):**

- The task is **multi-hop** (needs 2+ facts or subquestions).
- The answer benefits from **cross-validation** across sources.
- Entities are **ambiguous**, niche, or likely to have changed.
- You need **citations** from multiple sources for reliability.

**Do not use when (anti-triggers):**

- A **single lookup** is enough (use search\_light).
- The question can be answered reliably from **parametric knowledge** alone.

**Cost: High** (3 web searches).

612

613 **D.2 GAIA**

614 **D.2.1 Tool Descriptions**

Tool Description: reason\_only

Purpose: answer using only internal (parametric) knowledge without any web search.

What it does:

- Produces an answer using only internal knowledge and reasoning.
- Performs no web search.
- Lowest cost.

Use when (triggers):

- The question is reasoning-heavy (math/logic/proof/code reading/conceptual explanation).
- Facts are stable and unlikely to have changed (basic definitions, well-known principles).
- The task is to transform/rewrite/summarize provided text (no external facts needed).
- You are confident the answer does not depend on the latest information or specific citations.

Do not use when (anti-triggers):

- The question requires up-to-date facts (current events, latest versions, prices, active roles/titles).
- The question requests sources, citations, or “according to” style evidence.
- High risk of hallucination: niche proper nouns, many numbers, exact dates/titles.

Output expectation:

- Provide the best possible answer from internal knowledge.
- If uncertainty is high or factuality is critical, recommend escalation to search\_lite/search\_heavy.

Common failure modes:

- Outdated facts.
- Overconfidence on niche or recent details.

615

#### Tool Description: search\_light

Purpose: quickly verify or look up a single key fact via the web.

What it does:

- Runs EXACTLY ONE web search with ONE query.
- Uses the retrieved results as evidence for the final answer.
- Medium cost (1 search).

Use when (triggers):

- The question can be resolved by ONE primary lookup (one fact / definition / spec / date).
- Internal knowledge may be outdated or uncertain, but only a single check is needed.
- The task benefits from a source but does not require multiple independent sources.
- A high-precision query is obvious from the question.

Do not use when (anti-triggers):

- The question is multi-hop (requires 2+ independent facts or subquestions).
- You need cross-validation across sources (high risk of misinformation or ambiguous entities).
- The answer requires aggregation/comparison across multiple entities (list/compare/rank).
- The problem is pure reasoning/math/logic that does not benefit from web facts.

Output expectation:

- The final answer must be supported by evidence from the search results.
- If evidence is insufficient, state so and escalate to search\_heavy if needed.

Common failure modes:

- Wrong or overly broad query → irrelevant results (one-shot failure).
- A single source may be insufficient for disambiguation.

Query writing hint (1-shot):

- Include unique identifiers (proper nouns, version/year) and intent words such as “official”, “documentation”, “paper pdf”.

616

#### Tool Description: search\_deep

**Purpose:** solve questions that require multiple lookups or stronger evidence via multi-query web search.

**What it does:**

- Generates **THREE** distinct web-search queries.
- Retrieves results for all three queries.
- Combines and cross-checks the evidence to support the final answer.
- Highest cost (3 searches).

**Use when (triggers):**

- The question is multi-hop (requires 2+ subfacts; e.g., find a value, verify a condition, then produce a final result).
- The question demands high precision (exact number/date/quote/title) where small errors cause failure.
- The question implies freshness or current status (e.g., “latest”, “current”, “updated”, “as of <date>”).
- The answer benefits from cross-validation (multiple independent sources or official + secondary confirmation).
- Ambiguity is likely (multiple entities with the same name) and disambiguation is needed.

**Do not use when (anti-triggers):**

- A single lookup is sufficient (use `search_lite`).
- The task is pure reasoning with stable facts (use `knowledge_only`).
- You have no clear way to form three meaningfully distinct queries (risk of wasted cost).

**Output expectation:**

- Synthesize evidence across queries.
- If sources conflict, explicitly note the conflict and prefer more authoritative/primary sources.
- The final answer must be traceable to the retrieved evidence.

**Suggested 3-query decomposition:**

1. Core fact/definition (what is X?).
2. Exact detail (number/date/spec/quote).
3. Verification (official document or independent confirmation).

**Common failure modes:**

- Query diversity too low → redundant evidence.
- Over-searching simple questions → unnecessary cost.

618 **D.2.2 Reward Generation Prompt**

```

Reward Generation Prompt

REWARD_W_COST_SYSTEM_PROMPT = (
    "You are an evaluation engine for cost-aware tool selection.\n"
    "Your objective is to maximize expected utility: higher correctness with
    lower cost.\n"
    "Choose the single tool with the best cost--accuracy tradeoff.\n"
    "Think carefully, but do NOT output any reasoning.\n"
    "Output MUST be strict JSON only."
)

REWARD_W_COST_USER_PROMPT = (
    "Select exactly ONE tool to answer the question.\n"
    "Objective: maximize cost-aware correctness (expected utility).\n"
    "Decision rule:\n"
    "- Prefer the lowest-cost tool that is likely to answer correctly.\n"
    "- Choose a higher-cost tool only if it meaningfully increases the
    probability of correctness.\n\n"
    "Question:\n"
    "{question}\n\n"
    "Available tools (name: description; cost is implied by the tool
    choice):\n"
    "{tool_descriptions_in_text}\n\n"
    "Strict output format (JSON only; exactly one key):\n"
    "{\"answer\": \"<tool_name>\"}"
)

```

619

620 **D.2.3 Hyperparameters**

Method	Setting
CBLI	epoch = 10
PriorCB	epoch = 10, learning_rate = $1 \times 10^{-3}$ , weight_decay = $1 \times 10^{-4}$

Table 1: Training settings.