

---

# Gathering Context that Supports Decisions via Entropy Search with Language Models

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Real-world decision making systems require background information about the  
2 environment to take effective actions. However, this information is frequently  
3 incomplete or costly to acquire. Rather than presuming complete context, an effective  
4 decision maker must actively gather relevant information through a sequence  
5 of targeted follow-up questions before making decisions. This paper presents a  
6 framework for adaptive information gathering using large language models (LLMs)  
7 as interactive decision-making agents. Guided by an information-theoretic objective,  
8 the LLM selects questions that minimize the entropy of the predicted optimal  
9 action distribution, effectively prioritizing information that reduces uncertainty.  
10 Our method enables instance-specific reasoning under uncertainty and improves  
11 decision quality through principled context acquisition. We evaluate our approach  
12 on modified versions of three standard benchmarks—ID-ARC, GSM8K, and  
13 Fermi—adapted to study partially observable contexts where relevant information  
14 must be actively gathered. We assess performance using state-of-the-art LLMs.  
15 Empirically, we find that our proposed Entropy Search strategy consistently out-  
16 performs strong baselines, demonstrating the effectiveness of uncertainty-guided  
17 information gathering for LLM-based decision support<sup>1</sup>.

## 18 1 Introduction

19 In everyday interactions with large language models (LLMs), users often provide incomplete problem  
20 descriptions. This is not a shortcoming of the models alone; rather, it reflects a fundamental  
21 limitation in how humans communicate—users frequently omit relevant contextual details, either due  
22 to cognitive constraints or implicit assumptions about shared knowledge. As a result, the model must  
23 reason under uncertainty, which can result in suboptimal responses [Jiang et al., 2023]. Moreover,  
24 task underspecification can pose a significant AI safety and alignment risk [Amodei et al., 2016,  
25 Ngo et al., 2024, Dalrymple et al., 2024]. These limitations and risks point to a broader challenge in  
26 human-AI collaboration: effective decision support often requires adaptive information gathering.  
27 This interactive paradigm is especially valuable in high-stakes or data-scarce domains, such as  
28 forecasting [Tetlock and Gardner, 2016, Schoenegger et al., 2025], personalized medicine [Hamburg  
29 and Collins, 2010] and education [Ericsson and Pool, 2016, Vygotsky and Cole, 1978, Robinson and  
30 Aronica, 2009], where acquiring information is costly and outcomes depend critically on a limited set  
31 of key features.

32 In this work, we formalize this interactive decision-support setting as a two-stage sequential decision-  
33 making problem: An agent observes a user with a query and can ask several follow-up questions, and  
34 after gathering the user’s answers to these questions, the agent selects an action that maximizes the  
35 expected outcome. The core research question is: *Which questions should the agent ask to inform*  
36 *its final decision best?* The key technical challenge lies in identifying informative, non-redundant

---

<sup>1</sup>Our implementation is available at <https://anonymous.4open.science/r/info-gathering-047B/>

37 questions that reduce uncertainty about the optimal action, especially when features are correlated.  
 38 Our method uses entropy minimization to determine the next most informative question, conditioned  
 39 on previous responses. By minimizing uncertainty over optimal actions, our approach allows the agent  
 40 to prioritize questions that elicit the most information about the best action. This principled strategy  
 41 enables the agent to gather just enough information to act effectively, balancing informativeness with  
 42 efficiency. Building on this foundation, we make the following contributions:

- 43 • We introduce entropy search, a method that enables language model agents to ask targeted questions  
 44 that reduce outcome entropy and improve decisions under partial information. We theoretically  
 45 characterize its optimality and identify conditions under which it ensures optimal actions.
- 46 • We experiment on four state-of-the-art language models on partially observable variants of three  
 47 popular reasoning benchmarks, showing that the entropy minimization strategy consistently out-  
 48 performs competitive baselines.

## 49 2 Background

50 Algorithm 1 defines our typical information-gathering setting where an LLM agent interacts with  
 51 an external user. At each step, the agent observes the problem and current context  $H = (q_t, a_t)_{t=1}^{t-1}$ ,  
 52 selects a follow-up question  $q$ , and receives an answer  $a$  from the user. This interaction continues  
 53 until the budget is exhausted. The agent finally responds with an action  $A$ . The different algorithms  
 54 compared in the experiments only differ in how the agent selects a follow-up question.

55 Next, we introduce some notation for a Bayesian formulation of the information-gathering prob-  
 56 lem. A language model acts as a decision-making agent that interacts with a sequence of  $T$  users  
 57  $u_1, u_2, \dots, u_T$ . Each user  $u_t$  arrives with a hidden context vector  $X_t \in \mathcal{X}^K$ , where  $K$  is the total  
 58 number of features (potentially large). This contrasts with the standard contextual bandit setting  
 59 [Neu et al., 2022, Agrawal and Goyal, 2013], in which the full context is known to the agent. At the  
 60 beginning of each interaction, the user poses a query that reveals a subset of features—modeled as a  
 61 binary vector  $S_t^0 \in \{0, 1\}^K$ , corresponding to the initially observed context  $X_{t, S_t^0}$ . The agent can  
 62 then ask a sequence of follow-up questions, with the user giving an answer that reveals additional  
 63 features one at a time. This process results in a final observed subset  $S_t \in \{0, 1\}^K$ , which defines  
 64 the gathered context  $X_{t, S_t}$ . Based on this partial context, the agent selects an action  $A_t \in \mathcal{A}$  from a  
 65 finite (potentially large) set of possible actions. After the action is taken, the agent observes a binary  
 66 outcome  $Y_{t, A_t} \in \{0, 1\}$ .

67 The agent maintains two learnable parameters:  $(\omega, \theta)$ . Let  $Y_t \in \{0, 1\}^{\mathcal{A}}$  denote the potential outcome  
 68 vector over all possible actions. In the contextual bandit literature [Neu et al., 2022, Agrawal and  
 69 Goyal, 2013],  $\theta$  typically parameterizes the reward model, encoding the agent’s belief over outcomes  
 70 conditioned on the full context:  $p_t(Y_t | X_t, \theta)$ . We generalize this definition and let  $\theta$  parameterize  
 71 the full joint distribution  $p_t(Y_t, X_t | \theta)$ , capturing both the agent’s world knowledge about user  
 72 context distributions  $p_t(X_t | \theta)$  and the conditional outcome model  $p_t(Y_t | X_t, \theta)$ .

73 The expected reward for a selected action is given by:  $r_t(\theta, A_t, X_t) = \mathbb{E}_t[Y_{t, A_t} | \theta, X_t]$ . The optimal  
 74 action under the true model  $\theta^*$  is:  $A_t^* = \arg \max_{a \in \mathcal{A}} r_t(\theta^*, a, X_t)$ . After taking action  $A_t$  and  
 75 observing outcome  $Y_{t, A_t}$  and context  $X_{t, S_t}$ , the agent updates its posterior over  $\theta$  via:

$$q_{t+1}(\theta) = \frac{q_t(\theta) p_t(Y_{t, A_t} | \theta, A_t, X_{t, S_t})}{\sum_{\theta'} q_t(\theta') p_t(Y_{t, A_t} | \theta', A_t, X_{t, S_t})}$$

76 We assume a categorical outcome distribution:  $\sum_{a \in \mathcal{A}} \mathbb{E}_t[Y_{t, a} | \theta, X_t] = \sum_{a \in \mathcal{A}} r_t(\theta, a, X_t) = 1$   
 77 which differs from the Bernoulli setting where multiple actions can yield positive outcomes. In  
 78 practice, posterior updates over  $\theta$  can be implemented via supervised fine-tuning [Ouyang et al., 2022].  
 79 As  $\theta \rightarrow \theta^*$ , the expected outcome  $\mathbb{E}_t[Y_{t, A_t} | \theta, X_t]$  converges to the true expected outcome.

80 To estimate the optimal action, the agent first engages in a sequence of  $Q_t$  follow-up questions  
 81  $\{S_t^q\}_{q \in [Q_t]}$  with the user. Each question  $S_t^q \in \{0, 1\}^K$  specifies the indices of additional user context  
 82 to be revealed at step  $q$ , resulting in observed context  $X_{t, S_t^q}$ . The agent’s policy for generating these  
 83 questions is parameterized by  $\omega$ , and is conditioned on the context revealed so far:  $p_t(S_t^q | \omega, X_{t, S_t^{<q}})$ .  
 84 The interaction terminates after  $Q_t$  follow-up queries, at which point the agent has accumulated  
 85 a final observed subset  $S_t = \bigcup_{q=1}^{Q_t} S_t^q$  corresponding to context  $X_{t, S_t}$ . Based on this gathered  
 86 context, the agent selects an action:  $A_t = \arg \max_a r_t(\theta_t, a, X_{t, S_t})$  and incurs instantaneous regret

87  $\Delta_t = r_t(\theta^*, A_t^*, X_t) - r_t(\theta^*, A_t, X_t)$ . Here,  $A_t^*$  is the optimal action under the true model  $\theta^*$ . The  
 88 agent’s objective is to minimize the expected cumulative regret over  $T$  episodes:

$$R_T = \sum_{t=1}^T \mathbb{E} [r_t(\theta^*, A_t^*, X_t) - r_t(\theta^*, A_t, X_t)] = \sum_{t=1}^T \mathbb{E} [\Delta_t]$$

89 where the expectation is taken over all randomness: environment parameters and agent policy.

90 While regret can be used as a learning signal to update both  $(\omega, \theta)$ , optimizing  $\omega$  directly from regret  
 91 suffers from the sparse rewards problem: the agent receives feedback only after completing the  
 92 entire sequence of follow-up questions, based on the final outcome for its action, as is done in prior  
 93 works [Andukuri et al., 2024]. Delayed supervision makes it difficult to attribute success or failure to  
 94 individual follow-up queries. We don’t assume gold supervision for follow-up questions.

95 To overcome this, we propose an entropy search algorithm that evaluates each potential follow-up  
 96 question before it is asked. This approach assigns a score to every candidate query based on its  
 97 expected reduction in uncertainty over the outcome. Prior works have used similar intrinsic scores  
 98 based on information-gain [Houthoof et al., 2016, Mohamed and Jimenez Rezende, 2015] to tackle  
 99 sparse rewards. These scores are then used both to select the next follow-up question and to update  
 100 the policy parameters  $\omega$ , providing a more informative and dense supervision signal during training.  
 101 For reference, we summarize all the notation in Appendix 6.

### 102 3 Entropy Search Algorithm

103 We now present the Entropy Search algorithm 2. Following this strategy at each iteration, the decision  
 104 maker asks the question with the greatest reduction in posterior uncertainty about the optimal action,  
 105 on average, i.e., the question that elicits the most information.

106 To understand why this algorithm is optimal for the information-gathering task, we begin with a  
 107 formal definition of optimal follow-up questions and an associated objective for scoring and selecting  
 108 them. We then establish theoretical guarantees for this entropy-based objective. Finally, we provide a  
 109 regret analysis for the entropy search procedure’s exact and approximate implementations.

110 **Definition 3.1** (Optimal  $q^{\text{th}}$  follow-up question  $S_t^q$  for user  $u_t$ ). Before the  $q^{\text{th}}$  follow-up question,  
 111 the agent has already observed a subset of user context indices denoted by  $S_t^{<q} = \bigcup_{j=0}^{q-1} S_t^j$ . The  
 112 ideal query  $S_t^q$  would reveal just enough additional context to fully recover the potential outcome  
 113 distribution:  $p_t(Y_t | X_{t,S_t^{<q}}, X_{t,S_t^q}) = p_t(Y_t | X_t) = p_t(Y_t | X_{t,S_t^{<q}}, X_{t,(S_t^{<q})^-})$ , where  $(S_t^{<q})^- =$   
 114  $\mathbf{1} \setminus S_t^{<q}$  denotes the unobserved context indices.

115 Among all such queries, we aim to select one that requires minimal user effort to answer—formally,  
 116 the smallest number of new context features—while avoiding redundancy. Using a statistical distance  
 117 measure  $d$ , the optimal query is defined as:

$$\begin{aligned} S_t^q &= \arg \min_{S \in \{0,1\}^K} |S| \quad \text{s.t.} \quad d \left( p_t(Y_t | X_{t,S_t^{<q}}, X_{t,S}) \parallel p_t(Y_t | X_t) \right) = 0, \\ &= \arg \min_S \left\{ |S| \mid S \in \arg \min_{S'} d \left( p_t(Y_t | X_{t,S_t^{<q}}, X_{t,S'}) \parallel p_t(Y_t | X_t) \right), S \in \{0,1\}^K \right\}. \end{aligned}$$

118 The two formulations are equivalent because the minimum distance  $d$  is zero for at least one  
 119 solution—namely,  $S_t^q = (S_t^{<q})^-$ , which reveals all remaining context. However, this may not  
 120 be the minimal-norm solution. Once the distance is reduced to zero, further follow-up queries cannot  
 121 improve the prediction, and the agent can terminate questioning and respond with its final action.

122 In practice, the posterior  $p_t(Y_t | X_t)$  is not directly observable, making it infeasible to minimize  
 123 divergence-based objectives such as total variation or KL divergence. Instead, we use conditional  
 124 entropy difference as a surrogate distance measure, which has the key advantage of not requiring  
 125 access to the full posterior:

$$S_t^q = \arg \min_{S \in \{0,1\}^K} \mathbb{H}_t(Y_t | X_{t,S_t^{<q}}, X_{t,S}) - \mathbb{H}_t(Y_t | X_t) = \arg \min_{S \in \{0,1\}^K} \mathbb{H}_t(Y_t | X_{t,S_t^{<q}}, X_{t,S})$$

126 We follow a greedy procedure that selects one context feature at a time per follow-up question,  
 127 i.e.,  $|S_t^q| = 1$ . This choice makes the optimization tractable, reducing the search complexity from

128 combinatorial  $\mathcal{O}(2^K)$  to linear  $\mathcal{O}(K)$ . To do so, we assign a score to each candidate question:  
 129  $v_t^q(S) = -\mathbb{H}_t(Y_t|X_{t,S_t^{<q}}, X_{t,S})$  and select the context feature that maximizes this score. Thus, our  
 130 optimization becomes:

$$S_t^q = \arg \min_{S \in \{0,1\}^K, |S|=1} \mathbb{H}_t(Y_t|X_{t,S_t^{<q}}, X_{t,S}) = \arg \max_{S \in \{0,1\}^K, |S|=1} v_t^q(S)$$

131 We continue this procedure until the marginal entropy reduction  $\mathbb{H}_t(Y_t|X_{t,S_t^{<q}}) -$   
 132  $\mathbb{H}_t(Y_t|X_{t,S_t^{<q}}, X_{t,S_t^q})$  is sufficiently small—that is, when the new question contributes little  
 133 additional information. We assume a budget  $B$  that bounds the total number of follow-up  
 134 questions allowed per user. The final gathered context is:  $S_t = \bigcup_{j=0}^{Q_t} S_t^j$ . The entropy score  
 135  $v_t^q(S) = -\mathbb{H}_t(Y_t | X_{t,S_t^{<q}}, X_{t,S})$  is estimated using Monte Carlo sampling, by generating potential  
 136 outcome distributions under each candidate question  $S$ . In practice, due to limited compute, we  
 137 sample only a subset of  $k \ll K$  candidate follow-up questions from the policy  $p_t(S_t^q | \omega, X_{t,S_t^{<q}})$   
 138 prior to estimating the scores.

139 In our setting, outcomes are modeled as a categorical distribution over actions, where exactly one  
 140 action is correct. This structure allows us to estimate entropy more efficiently. Rather than computing  
 141 the full distribution over potential outcomes and then selecting the best action, we directly compute  
 142 the entropy over the action distribution induced by the model. This approach avoids a separate  
 143 outcome prediction step and simplifies entropy estimation within our algorithm. The following  
 144 lemma formalizes this equivalence (proof in Appendix 7.1):

145 **Lemma 3.2** (Equivalence of Action and Outcome Entropy). *For a categorical distribution over*  
 146 *outcomes, minimization of outcome entropy is equivalent to minimization of action entropy:*

$$\arg \min_S \mathbb{H}_t(Y_t|X_{t,S}) = \arg \min_S \mathbb{H}_t(A_t|X_{t,S})$$

### 147 3.1 Optimality of Entropy Search

148 Matching entropy does not imply matching distributions in general. However, since we minimize  
 149 entropy only by expanding the conditioning set, it decreases strictly unless the conditional distributions  
 150 match. This justifies entropy as a surrogate loss for the distribution matching objective. The following  
 151 lemma formalizes this via a bound on total variation distance:

**Lemma 3.3** (Upper Bound on Total Variation). *For any  $\epsilon > 0$ , we have*

$$\mathbb{H}_t(Y_t|X_{t,S_t}) - \mathbb{H}_t(Y_t|X_t) \leq \epsilon \Rightarrow \|p_t(Y_t|X_{t,S_t}) - p_t(Y_t|X_t)\|_1 \leq \sqrt{\frac{\epsilon}{2}}$$

152 *In particular,  $\mathbb{H}_t(Y_t|X_{t,S_t}) = \mathbb{H}_t(Y_t|X_t) \Rightarrow p_t(Y_t|X_{t,S_t}) = p_t(Y_t|X_t)$*

153 The proof follows from Pinsker’s inequality, and we provide the details in the Appendix 7.2.

### 154 3.2 Convergence Analysis

155 We show that the minimal entropy-minimizing context subset  $S_t^*$  suffices for the posterior sampling  
 156 of the optimal action. In particular, conditioning on  $X_{t,S_t^*}$  is equivalent to conditioning on the full  
 157 context  $X_t$  for both the current parameters  $\theta_t$  and the true parameters  $\theta^*$ . This means no additional  
 158 context is needed once outcome entropy is minimized. The same holds for any  $S_t \supseteq S_t^*$ .

159 **Lemma 3.4.** *Following the information-gathering procedure and posterior sampling for action*  
 160 *selection, the following equalities hold under exact entropy search:*

$$\mathbb{P}_t(A_t^*, \theta^* | X_{t,S_t^*}) = \mathbb{P}_t(A_t, \theta_t | X_{t,S_t^*}) = \mathbb{P}_t(A_t^*, \theta^* | X_t) = \mathbb{P}_t(A_t, \theta_t | X_t)$$

161 The proof in Appendix 7.3 follows from posterior sampling and Lemma 3.3.

162 We state a key lemma extending the information ratio framework of Neu et al. [2022] to partial  
 163 context settings, yielding an analogous regret bound:

164 **Lemma 3.5** (Bounded Information Ratio). *The information ratio for the observed set  $S_t$  is defined*  
 165 *as follows:*

$$\rho_t(S_t) = \frac{(\mathbb{E}_t[r_t(\theta_t, A_t, X_{t,S_t}) - r_t(\theta^*, A_t, X_{t,S_t}) | X_{t,S_t}])^2}{\mathbb{I}_t(\theta^*; Y_{t,A_t} | A_t, X_{t,S_t})}$$

166 *admits the bound  $\rho_t(S_t) \leq 2 \sum_{a \in \mathcal{A}} \mathbb{E}[\bar{r}_t(a, X_{t,S_t})]$*

167 The proof technique is similar to [Neu et al., 2022], and more details are in Appendix 7.4.

168 **Theorem 3.6** (Regret of Exact Entropy Search). *Exact entropy search achieves the regret bound:*  
 169  $R_T \leq \sqrt{2T\mathbb{H}(\theta^*)}$

170 We provide a detailed proof in Appendix 7.7.

171 At each turn  $t$ , we update the  $\omega$  parameters using the sampled entropy scores. This process can be  
 172 viewed as transferring knowledge from System-2 (deliberate, exhaustive evaluation) to System-1 (fast,  
 173 heuristic-based selection) over time, consistent with the dual-process theory of cognition [Booch  
 174 et al., 2021]. The key idea is that we are learning a policy that learns how to gather information  
 175 efficiently. Since evaluating all  $K$  follow-up questions (a full System-2 search) is computationally  
 176 infeasible, we instead train  $\omega$  to generate a smaller, high-quality subset of  $k \ll K$  candidate questions  
 177 using System-1-style inference. We then apply System-2 reasoning within this restricted subset by  
 178 selecting the candidate with the lowest estimated entropy. As  $\omega$  improves, the quality of this subset  
 179 increases, yielding better candidates more often. This motivates our assumption that the entropy gap  
 180  $\epsilon_t$  diminishes over time—a property we exploit in Theorem 3.7 to establish sublinear regret despite  
 181 constrained computational resources.

182 The main contributor to non-zero gap  $\epsilon_t$  is failing to include the optimal question  $\bar{S}_t^q =$   
 183  $\arg \min_{S_t^q} \mathbb{H}_t(Y_t|X_t, S_t^{<q}, X_t, S_t^q)$  in the candidate set for entropy estimation, not poor outcome  
 184 estimation. While updates to  $\theta$  do refine entropy estimates through calibrated outcome prediction,  
 185 they are unlikely to significantly reduce  $\epsilon_t$ , as most language models have a strong enough prior to  
 186 choosing the optimal question if all options were considered.

187 **Theorem 3.7** (Regret). *Approximate entropy search achieves the regret bound:  $R_T \leq \sqrt{2T\mathbb{H}(\theta^*)} +$   
 188  $\mathcal{O}(T^{1-\alpha/2})$  in the one-hot outcome setting.*

189 *We assume that Entropy Search learns to select the optimal subset with diminishing error  $\epsilon_t =$   
 190  $\mathbb{H}_t(Y_t|X_t, S_t, \theta) - \mathbb{H}_t(Y_t|X_t, \theta)$ ,  $\mathbb{E}[\epsilon_t] = \mathcal{O}(\frac{1}{t^\alpha})$ . In particular, for  $\alpha = 1$ , we obtain  $R_T = \mathcal{O}(\sqrt{T})$ ,  
 191 identical to the exact entropy search setting.*

192 We provide a detailed proof in Appendix 7.8.

**Lemma 3.8** (Entropy Search Lowers Per-Question Suboptimality over ReAct). *We define the per-  
 question suboptimality of an algorithm as  $\epsilon_t^q = \mathbb{H}_t(Y_t|X_t, S_t^{<q}, X_t, S_t^q) - \mathbb{H}_t(Y_t|X_t, S_t^{<q}, X_t, \bar{S}_t^q)$  where  
 $\bar{S}_t^q$  is the optimal  $q^{\text{th}}$  question that minimizes the outcome entropy. Let  $p_{t,q}^* = p_t(S_t^q = \bar{S}_t^q | \omega, X_t, S_t^{<q})$   
 be the probability that the policy selects the optimal  $q^{\text{th}}$  follow-up question. For a budget of  $k \ll K$   
 candidate questions in Entropy Search (ES), we achieve lower suboptimality than the ReAct baseline:*

$$(\epsilon_t^q)_{ES} \leq (1 - p_{t,q}^*)^{k-1} (\epsilon_t^q)_{ReAct}$$

193 We defer the proof to Appendix 7.9 along with a discussion of this result.

## 194 4 Experiments

195 We have presented a theoretical framework for deriving a training objective that encourages a language  
 196 model to seek information for decision-making. In this section, we experiment with three benchmarks,  
 197 each comprising 1000 data points, to facilitate empirical studies. We go into detail on each of the  
 198 benchmarks GSM8K-Q, 1D-ARC, and Fermi Problems, in Appendix 9.

### 199 4.1 Algorithms

200 We compare ENTROPY SEARCH to two other agent baselines:

- 201 1. REACT: the agent first produces a chain-of-thought conditioned on  $(problem, H)$  and emits the  
 202 next question greedily at the point where reasoning requires a missing fact.
- 203 2. REFLEXION: The reflexion agent [Shinn et al., 2023] includes a dedicated reflection component,  
 204 where after each interaction, the agent evaluates what it has learned, assesses the effectiveness of its  
 205 previous queries, identifies missing information, and adjusts its questioning strategy accordingly.  
 206 The prompt structure incorporates a dedicated “Reflection” section, includes past reflections for  
 207 continuity, and instructs the agent explicitly to reflect before reasoning. Memory stores and  
 208 displays reflections in subsequent interactions, enabling learning and strategic adjustments.

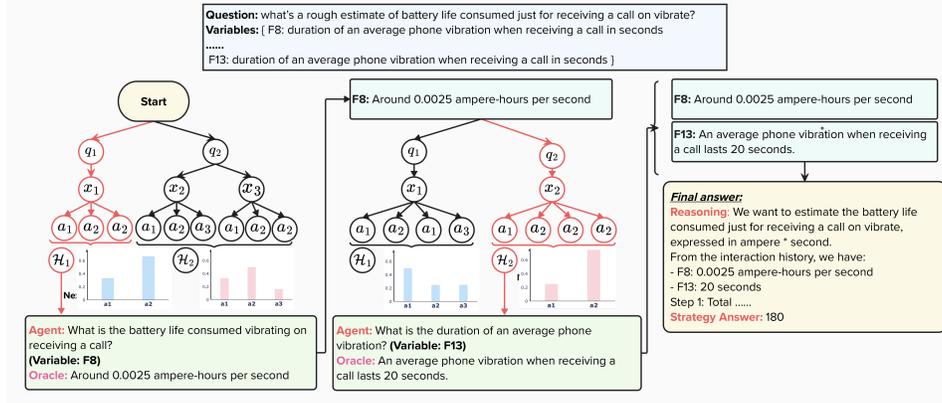


Figure 2: Example of the ENTROPY SEARCH agent interacting on a Fermi estimation task. The agent selects targeted follow-up questions (e.g., requesting values for Variables F8 and F13) to reduce uncertainty and gather key facts incrementally. Using the retrieved values, the agent composes a final answer by chaining units and performing multi-step reasoning, highlighting the interpretability and structured nature of the entropy-guided decision process.

209 We also compare against a FULL INFORMATION setting with  
 210 access to the entire user context.

211 We primarily used the Gemma-3 12B [Team et al., 2025] model  
 212 to balance performance and computational efficiency. We also  
 213 ran evaluations using multiple state-of-the-art language models  
 214 to assess robustness across architectures. These additional details  
 215 and results are in Appendix 14.1.

## 216 4.2 Results

217 Figure 12 presents the comparative evaluation of method perfor-  
 218 mance across three challenging datasets: 1D-ARC, GSM8K, and  
 219 Fermi. As expected, the FULL INFORMATION baseline serves  
 220 as a practical performance ceiling due to its complete contex-  
 221 tual details. Crucially, ENTROPY SEARCH substantially narrows  
 222 this gap more effectively than agentic reasoning methods like  
 223 REFLEXION and REACT, highlighting the value of strategically  
 224 prioritizing uncertainty reduction.

225 In principle, one might expect sufficiently intelligent agents could  
 226 infer complete symbolic solutions directly from GSM and Fermi  
 227 question statements, reducing the need for additional information.  
 228 However, our empirical findings robustly contradict this hypothe-  
 229 sis, showing that entropy-guided information gathering significantly enhances performance even  
 230 when symbolic solutions might seem inferable initially.

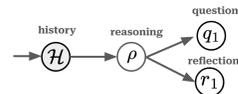
231 Moreover, the performance on Fermi problems underscores another dimension of our method’s  
 232 strength: its ability to handle real-world estimation tasks that require nuanced common-sense reason-  
 233 ing. Remarkably, ENTROPY SEARCH not only outperformed the competing methods but matched the  
 234 performance of the FULL INFORMATION scenario. This is particularly compelling given the inherent  
 235 uncertainty and complexity involved in Fermi estimations.

236 Further granularity is provided by the turn-based evaluation depicted in Figure 4, confirming that  
 237 ENTROPY SEARCH consistently improves with more turns and outperforms baselines at every stage.  
 238 These findings collectively underscore ENTROPY SEARCH’s effectiveness in strategically acquiring  
 239 context under uncertainty, thus significantly enhancing decision-making quality in partially observable  
 240 environments. Additional detailed turn-based analyses for the 1D-ARC and GSM8K benchmarks are  
 241 available in Appendix 14.2, further validating these insights.

### REACT procedure



### REFLEXION procedure



### ENTROPY SEARCH procedure

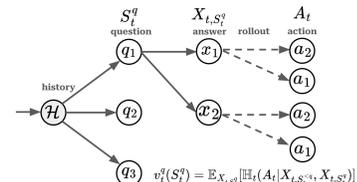


Figure 1: Comparison of question selection mechanisms

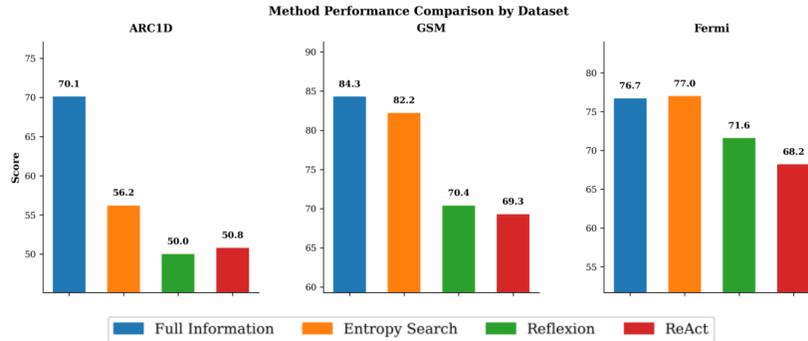


Figure 3: Method performance comparison across three benchmarks. FULL INFORMATION represents the upper bound with complete context. ENTROPY SEARCH consistently outperforms REACT and REFLEXION across all datasets and approaches the performance of the FULL INFORMATION baseline, demonstrating the value of principled, entropy-guided question selection.

## 242 5 Conclusion, Limitations, and Future Work

243 We introduce an entropy search for context  
 244 gathering in LLMs agents, enabling the se-  
 245 lection of questions that reduce uncertainty  
 246 and improve downstream decision-making.  
 247 Our framework is model-agnostic and may  
 248 be paired with a variety of models. By  
 249 applying this method across various rea-  
 250 soning tasks, we demonstrate that entropy  
 251 search leads to superior performance over  
 252 state-of-the-art baselines. One limitation of  
 253 our work is the computational cost associ-  
 254 ated with the Monte Carlo estimation of en-  
 255 tropy. Future work could address this by us-  
 256 ing value network estimation to reduce the  
 257 computational burden. Further, our work  
 258 did not consider information gathering in  
 259 the context of safe exploration [Ray et al.,  
 260 2019, Garcia and Fernández, 2015, Wachi  
 261 et al., 2023], which is an important direc-  
 262 tion for future work. This work also opens  
 263 several other promising directions for fu-  
 264 ture research. Applying this framework to  
 265 interactive human-AI collaboration in do-  
 266 mains such as scientific discovery, diagnos-  
 267 tic reasoning, and judgemental forecasting  
 268 could yield both practical benefits and  
 269 deeper insights into the role of adaptive  
 270 information gathering in complex decision  
 271 systems. Information seeking is also a key  
 capability that should be monitored for AI  
 safety and alignment; future frontier mod-  
 els can be benchmarked against our work  
 to monitor information seeking capability  
 development. By endorsing language mod-  
 els with information seeking capabilities,  
 our work contributes a framework for a  
 cooperative, interactive intelligent system,  
 highlighting the broader potential of LLMs  
 as strategic agents in sequential decision-  
 making under uncertainty.

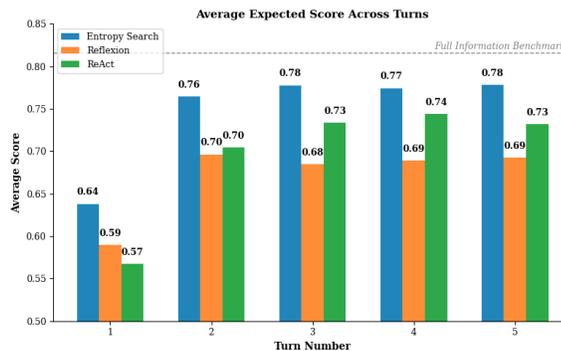


Figure 4: Turn-based evaluation for Fermi Problems. We conduct detailed comparisons across methods (ENTROPY SEARCH, REFLEXION, REACT) and turns, revealing that ES consistently closes the performance gap with the full information setting with each additional question. This highlights both the generality and efficiency of entropy-guided information acquisition.

## 272 References

- 273 Shipra Agrawal and Navin Goyal. Thompson sampling for contextual bandits with linear payoffs. In  
 274 *International conference on machine learning*, pages 127–135. PMLR, 2013.
- 275 Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané.  
 276 Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.
- 277 Chinmaya Andukuri, Jan-Philipp Fränken, Tobias Gerstenberg, and Noah Goodman. STar-GATE:  
 278 Teaching language models to ask clarifying questions. In *First Conference on Language Modeling*,  
 279 2024. URL <https://openreview.net/forum?id=CrzAj0kZjR>.

- 280 Grady Booch, Francesco Fabiano, Lior Horesh, Kiran Kate, Jonathan Lenchner, Nick Linck, Andreas  
 281 Loreggia, Keerthiram Murgesan, Nicholas Mattei, Francesca Rossi, et al. Thinking fast and  
 282 slow in ai. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages  
 283 15042–15046, 2021.
- 284 Djallel Bouneffouf, Irina Rish, Guillermo A. Cecchi, and Raphael Feraud. Context attentive bandits:  
 285 Contextual bandit with restricted context, 2017. URL <https://arxiv.org/abs/1705.03821>.
- 286 Kathryn Chaloner and Isabella Verdinelli. Bayesian experimental design: A review. *Statistical Science*,  
 287 10(3):273–304, 1995. ISSN 08834237. URL <http://www.jstor.org/stable/2246015>.
- 288 David Dalrymple, Joar Skalse, Yoshua Bengio, Stuart Russell, Max Tegmark, Sanjit Seshia, Steve  
 289 Omohundro, Christian Szegedy, Ben Goldhaber, Nora Ammann, et al. Towards guaranteed safe ai:  
 290 A framework for ensuring robust and reliable ai systems. *arXiv preprint arXiv:2405.06624*, 2024.
- 291 M. H. DeGroot. Uncertainty, Information, and Sequential Experiments. *The Annals of Mathematical*  
 292 *Statistics*, 33(2):404 – 419, 1962. doi: 10.1214/aoms/1177704567. URL [https://doi.org/10.](https://doi.org/10.1214/aoms/1177704567)  
 293 [1214/aoms/1177704567](https://doi.org/10.1214/aoms/1177704567).
- 294 Audrey Durand and Christian Gagné. Thompson sampling for combinatorial bandits and its applica-  
 295 tion to online feature selection. In *Workshops at the Twenty-Eighth AAAI Conference on Artificial*  
 296 *Intelligence*, 2014.
- 297 Anders Ericsson and Robert Pool. *Peak: Secrets from the new science of expertise*. Random House,  
 298 2016.
- 299 Javier Garcia and Fernando Fernández. A comprehensive survey on safe reinforcement learning.  
 300 *Journal of Machine Learning Research*, 16(1):1437–1480, 2015.
- 301 Margaret A Hamburg and Francis S Collins. The path to personalized medicine. *New England*  
 302 *Journal of Medicine*, 363(4):301–304, 2010.
- 303 Rein Houthoofd, Xi Chen, Yan Duan, John Schulman, Filip De Turck, and Pieter Abbeel. Vime:  
 304 Variational information maximizing exploration. *Advances in neural information processing*  
 305 *systems*, 29, 2016.
- 306 Mingjian Jiang, Yangjun Ruan, Sicong Huang, Saifei Liao, Silviu Pitis, Roger Baker Grosse,  
 307 and Jimmy Ba. Calibrating language models via augmented prompt ensembles. [https:](https://openreview.net/pdf?id=L0dc4wqbNs)  
 308 [//openreview.net/pdf?id=L0dc4wqbNs](https://openreview.net/pdf?id=L0dc4wqbNs), 2023.
- 309 Ashwin Kalyan, Abhinav Kumar, Arjun Chandrasekaran, Ashish Sabharwal, and Peter Clark. How  
 310 much coffee was consumed during emnlp 2019? fermi problems: A new reasoning challenge for  
 311 ai. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*,  
 312 pages 7318–7328, 2021.
- 313 Sanath Kumar Krishnamurthy and Susan Athey. Survey bandits with regret guarantees, 2020. URL  
 314 <https://arxiv.org/abs/2002.09814>.
- 315 Belinda Z Li, Been Kim, and Zi Wang. Questbench: Can llms ask the right question to acquire  
 316 information in reasoning tasks? *arXiv preprint arXiv:2503.22674*, 2025.
- 317 Shakir Mohamed and Danilo Jimenez Rezende. Variational information maximisation for intrinsically  
 318 motivated reinforcement learning. *Advances in neural information processing systems*, 28, 2015.
- 319 Gergely Neu, Iuliia Olkhovskaia, Matteo Papini, and Ludovic Schwartz. Lifting the information  
 320 ratio: An information-theoretic analysis of thompson sampling for contextual bandits. *Advances in*  
 321 *Neural Information Processing Systems*, 35:9486–9498, 2022.
- 322 Richard Ngo, Lawrence Chan, and Sören Mindermann. The alignment problem from a deep learning  
 323 perspective. In *The Twelfth International Conference on Learning Representations*, 2024. URL  
 324 <https://openreview.net/forum?id=fh8EYKFKns>.

- 325 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong  
326 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow  
327 instructions with human feedback. *Advances in neural information processing systems*, 35:27730–  
328 27744, 2022.
- 329 Alex Ray, Joshua Achiam, and Dario Amodei. Benchmarking safe exploration in deep reinforcement  
330 learning. *arXiv preprint arXiv:1910.01708*, 7(1):2, 2019.
- 331 Ken Robinson and Lou Aronica. *The element: How finding your passion changes everything*. Penguin,  
332 2009.
- 333 Daniel Russo and Benjamin Van Roy. An information-theoretic analysis of thompson sampling.  
334 *Journal of Machine Learning Research*, 17(68):1–30, 2016.
- 335 Philipp Schoenegger, Peter S Park, Ezra Karger, Sean Trott, and Philip E Tetlock. Ai-augmented  
336 predictions: Llm assistants improve human forecasting accuracy. *ACM Transactions on Interactive*  
337 *Intelligent Systems*, 15(1):1–25, 2025.
- 338 Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion:  
339 Language agents with verbal reinforcement learning. *Advances in Neural Information Processing*  
340 *Systems*, 36:8634–8652, 2023.
- 341 Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej,  
342 Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, et al. Gemma 3 technical  
343 report. *arXiv preprint arXiv:2503.19786*, 2025.
- 344 Philip E Tetlock and Dan Gardner. *Superforecasting: The art and science of prediction*. Random  
345 House, 2016.
- 346 Sabina Tomkins, Peng Liao, Predrag Klasnja, and Susan Murphy. Intelligentpooling: practical  
347 thompson sampling for mhealth. *Machine Learning*, 110(9):2685–2727, 2021. doi: 10.1007/  
348 s10994-021-05995-8. URL <https://doi.org/10.1007/s10994-021-05995-8>.
- 349 Lev Semenovich Vygotsky and Michael Cole. *Mind in society: Development of higher psychological*  
350 *processes*. Harvard university press, 1978.
- 351 Akifumi Wachi, Wataru Hashimoto, Xun Shen, and Kazumune Hashimoto. Safe exploration in  
352 reinforcement learning: A generalized formulation and algorithms. *Advances in Neural Information*  
353 *Processing Systems*, 36:29252–29272, 2023.
- 354 Yudong Xu, Wenhao Li, Pashootan Vaezipoor, Scott Sanner, and Elias Boutros Khalil. Llms and  
355 the abstraction and reasoning corpus: Successes, failures, and the importance of object-based  
356 representations. *Transactions on Machine Learning Research*, 2024.
- 357 Jiayu Yao, Emma Brunskill, Weiwei Pan, Susan Murphy, and Finale Doshi-Velez. Power constrained  
358 bandits. In Ken Jung, Serena Yeung, Mark Sendak, Michael Sjoding, and Rajesh Ranganath,  
359 editors, *Proceedings of the 6th Machine Learning for Healthcare Conference*, volume 149 of  
360 *Proceedings of Machine Learning Research*, pages 209–259. PMLR, 06–07 Aug 2021. URL  
361 <https://proceedings.mlr.press/v149/yao21a.html>.

Table 1: Notation

$t$	Index of user interacting with the agent
$T$	Total number of users the agent interacts with
$u_t$	The $t^{\text{th}}$ user
$X_t$	Random variable representing full context for user $u_t$
$K$	Total number of features in a user’s context (potentially large)
$\mathcal{X}$	Set of values from which a user’s context feature is drawn, $X_t \in \mathcal{X}^K$
$A_t$	Final action performed by agent in response to user $u_t$
$\mathcal{A}$	Finite set of actions that agent can perform
$A$	$ \mathcal{A} $ , number of actions in agent’s action space (potentially large)
$Y_{t,A_t}$	Outcome in $\{0, 1\}$ for user $u_t$ given action $A_t$
$Y_t$	Potential outcomes under different actions, in $\{0, 1\}^{\mathcal{A}}$
$S_t^0$	Binary random variable representing the set of context feature indices initially revealed by the user, in $\{0, 1\}^K$
$S_t^q$	Binary random variable representing the set of additional context feature indices gathered by asking the $q^{\text{th}}$ follow-up question
$X_{t,S_t^q}$	Partial context of $u_t$ revealed by question asking about features in $S_t^q$
$S_t$	Binary random vector representing the entire set of user feature indices that are gathered by asking all the follow-up questions $\bigcup_{q \geq 0} S_t^q$
$X_{t,S_t}$	Random variable representing the partial user context that is observed by the agent at the time of performing an action
$Q_t$	Total number of follow-up questions that the agent asks the user $u_t$
$B$	Budget for maximum number of follow-up questions that the agent can ask the user $u_t$
$S_t^{<q}$	$\bigcup_{j=0}^{q-1} S_t^j$ , set of user feature indices observed before asking the $q^{\text{th}}$ question
$\mathbb{1}$	Set of all user feature indices $\{1\}^K$
$(S)^-$	$\mathbb{1} \setminus S$ , complement of user context features in set $S$ , possibly unobserved feature indices
$\theta$	Random variable representing agent parameters for the joint distribution of user context $X$ and potential outcomes $Y$
$\omega$	Random variable representing agent parameters for follow-up question generation policy
$\mathcal{F}_t$	Filtration / History prior to user $t$ , equals $\{X_{1,S_1}, A_1, Y_{1,A_1}, \dots, X_{t,S_t}, A_t, Y_{t,A_t}\}$
$p_t$	Probability distribution conditioned on $\mathcal{F}_t$
$r_t$	Reward / Expected outcome model for user $u_t$ , $r_t(\theta, A_t, X_{t,S}) = \mathbb{E}_t[Y_{t,A_t}   \theta, X_{t,S}]$
$\theta^*$	Random variable representing optimal parameters for reward / expected outcome model
$A_t^*$	Random variable representing optimal action that maximizes expected outcome $r_t(\theta^*, A, X_t)$ given full context
$q_t(\theta)$	Posterior distribution over parameter $\theta$ after observing $\mathcal{F}_t$
$\Delta_t$	Instantaneous regret for action selected for user $u_t$
$R_T$	Expected cumulative regret
$d(\cdot    \cdot)$	Statistical distance measure on two distributions
$\mathbb{H}_t$	Entropy of distribution conditional on $\mathcal{F}_t$
$v_t^q(S)$	Entropy score for candidate $q^{\text{th}}$ question, negative of expected outcome entropy conditional on observations before $q$
$k$	Budget for number of candidate follow-up questions that are evaluated with the entropy score prior to selecting one to ask the user
$\rho_t(S_t)$	Generalization of the information ratio of Neu et al. [2022] to partial context in $S_t$
$\mathbb{I}_t$	Mutual information conditional on $\mathcal{F}_t$
$\epsilon_t$	Gap between the outcome entropy for the selected subset $S_t$ and the full context
$\alpha$	Rate of exponential decay of $\epsilon_t$
$\epsilon_t^q$	Gap between outcome entropy for selected $q^{\text{th}}$ follow-up question for user $u_t$ and the optimal entropy minimizing question
$p_{t,q}$	Probability distribution conditional on $\mathcal{F}_t$ and user context observed prior to the $q^{\text{th}}$ question
$\mathbb{P}$	Probability measure

363 **7 Theoretical Results**

364 **Lemma 7.1** (Equivalence of Action and Outcome Entropy). *For a categorical distribution over*  
 365 *outcomes:*

$$\arg \min_S \mathbb{H}_t(Y_t|X_{t,S}) = \arg \min_S \mathbb{H}_t(A_t|X_{t,S})$$

366 *Proof.* The model selects actions with the same likelihood as the expected outcome

$$p_t(A_t = a|X_{t,S}) = p_t(Y_{t,a} = 1|X_{t,S})$$

367 Let  $e_a \in \{0, 1\}^{\mathcal{A}}$  denote the one-hot vector that is one at index  $a$ .

$$\begin{aligned} \mathbb{H}_t(Y_t|X_{t,S}) &= - \sum_{y \in \{0,1\}^{\mathcal{A}}} p_t(Y_t = y|X_{t,S}) \log(p_t(Y_t = y|X_{t,S})) \\ &= - \sum_{y \in \{e_a | a \in \mathcal{A}\}} p_t(Y_t = y|X_{t,S}) \log(p_t(Y_t = y|X_{t,S})) \quad (\text{one only correct action}) \\ &= - \sum_{a \in \mathcal{A}} p_t(Y_{t,a} = 1|X_{t,S}) \log(p_t(Y_{t,a} = 1|X_{t,S})) \\ &= \sum_{a \in \mathcal{A}} p_t(A_t = a|X_{t,S}) \log(p_t(A_t = a|X_{t,S})) \\ &= \mathbb{H}_t(A_t|X_{t,S}) \end{aligned}$$

368

□

**Lemma 7.2** (Upper Bound on Total Variation). *For any  $\epsilon > 0$ , we have*

$$\mathbb{H}_t(Y_t|X_{t,S_t}) - \mathbb{H}_t(Y_t|X_t) \leq \epsilon \Rightarrow \|p_t(Y_t|X_{t,S_t}) - p_t(Y_t|X_t)\|_1 \leq \sqrt{\frac{\epsilon}{2}}$$

369 *In particular,  $\mathbb{H}_t(Y_t|X_{t,S_t}) = \mathbb{H}_t(Y_t|X_t) \Rightarrow p_t(Y_t|X_{t,S_t}) = p_t(Y_t|X_t)$*

*Proof.*

$$\begin{aligned} \|p_t(Y_t|X_{t,S_t}) - p_t(Y_t|X_t)\|_1 &\leq \sqrt{\frac{1}{2} D_{KL}(p_t(Y_t|X_{t,S_t}) || p_t(Y_t|X_t))} \quad (\text{Pinsker's Inequality}) \\ &= \sqrt{\frac{1}{2} D_{KL}(p_t(Y_t|X_{t,S_t}) || p_t(Y_t|X_{t,S_t}, X_{t,S_t}^-))} \\ &= \sqrt{\frac{1}{2} \mathbb{I}_t(Y_t; X_{t,S_t}^- | X_{t,S_t})} = \sqrt{\frac{1}{2} (\mathbb{H}_t(Y_t|X_{t,S_t}) - \mathbb{H}_t(Y_t|X_{t,S_t}, X_{t,S_t}^-))} \\ &= \sqrt{\frac{1}{2} (\mathbb{H}_t(Y_t|X_{t,S_t}) - \mathbb{H}_t(Y_t|X_t))} \leq \sqrt{\frac{\epsilon}{2}} \end{aligned}$$

370 For  $\epsilon = 0$ , we have a simpler argument:

$$\begin{aligned} S_t &= \arg \min_S \mathbb{H}_t(Y_t|X_{t,S}) \Rightarrow \mathbb{H}_t(Y_t|X_{t,S_t}) = \mathbb{H}_t(Y_t|X_t) \\ &\Rightarrow Y_t \perp X_t | X_{t,S_t} \Rightarrow p_t(Y_t|X_{t,S_t}) = p_t(Y_t|X_t) \end{aligned}$$

371 In particular, this also holds for the smallest norm set of user features  $S_t^*$  that satisfies this condition.

372

□

373 **Lemma 7.3.** *Following the information-gathering procedure and posterior sampling for action*  
 374 *selection, the following equalities hold under exact entropy minimization:*

$$\mathbb{P}_t(A_t^*, \theta^* | X_{t,S_t^*}) = \mathbb{P}_t(A_t, \theta_t | X_{t,S_t^*}) = \mathbb{P}_t(A_t^*, \theta^* | X_t) = \mathbb{P}_t(A_t, \theta_t | X_t)$$

375 *Proof.* Exact entropy minimization implies that the outcome distributions conditioned on the minimal  
 376 sufficient context  $S_t^*$  match those conditioned on the full context by Lemma 3.3:

$$\mathbb{P}_t(Y_t | \theta_t, X_{t,S_t^*}) = \mathbb{P}_t(Y_t | \theta_t, X_t)$$

377 Since posterior sampling draws  $\theta_t$  from the same distribution as  $\theta^*$ , we also have:

$$\mathbb{P}_t(Y_t | \theta^*, X_{t,S_t^*}) = \mathbb{P}_t(Y_t | \theta^*, X_t)$$

378 Therefore, under both  $\theta_t$  and  $\theta^*$ , conditioning on  $X_{t,S_t^*}$  is equivalent to conditioning on  $X_t$  for  
 379 predicting  $Y_t$ .

380 Because the action is selected deterministically to maximize the expected outcome, the distribution  
 381 over actions must also be identical:

$$\begin{aligned} \mathbb{P}_t(A_t | \theta_t, X_{t,S_t^*}) &= \mathbb{P}_t(A_t | \theta_t, X_t) \\ \mathbb{P}_t(A_t^* | \theta^*, X_{t,S_t^*}) &= \mathbb{P}_t(A_t^* | \theta^*, X_t) \end{aligned}$$

382 Finally, combining the equivalence of action distributions with that of  $\theta$  under posterior sampling  
 383 gives:

$$\begin{aligned} \mathbb{P}_t(A_t, \theta_t | X_{t,S_t^*}) &= \mathbb{P}_t(A_t, \theta_t | X_t) \\ \mathbb{P}_t(A_t^*, \theta^* | X_{t,S_t^*}) &= \mathbb{P}_t(A_t^*, \theta^* | X_t) \end{aligned}$$

384 which proves the lemma. □

385 **Lemma 7.4** (Bounded Information Ratio). *Information ratio for observed set  $S_t$  defined as follows:*

$$\rho_t(S_t) = \frac{(\mathbb{E}_t[r_t(\theta_t, A_t, X_{t,S_t}) - r_t(\theta^*, A_t, X_{t,S_t}) | X_{t,S_t}])^2}{\mathbb{I}_t(\theta^*; Y_{t,A_t} | A_t, X_{t,S_t})}$$

386 admits the bound  $\rho_t(S_t) \leq 2 \sum_{a \in \mathcal{A}} \mathbb{E}[\bar{r}_t(a, X_{t,S_t})]$

387 *Proof.* Following the proof technique in Neu et al. [2022], we bound a quantity analogous to the  
 388 instantaneous regret  $\Delta_t = r_t(\theta_t, A_t, X_{t,S_t}) - r_t(\theta^*, A_t, X_{t,S_t})$ .

389 It is easy to see that for  $S_t \supseteq S_t^*$ , this simplifies to the instantaneous regret:

$$\begin{aligned} \Delta_t &= r_t(\theta^*, A_t^*, X_{t,S_t^*}) - r_t(\theta^*, A_t, X_{t,S_t^*}) \\ &\quad (\text{under posterior sampling conditional distribution } (\theta^*, A_t^*) \text{ matches } (\theta_t, A_t)) \\ &= r_t(\theta_t^*, A_t^*, X_t) - r_t(\theta^*, A_t, X_t) \\ &\quad (\text{Using } \mathbb{P}_t(Y_{t,A_t^*} | \theta^*, X_{t,S_t^*}) = \mathbb{P}_t(Y_{t,A_t^*} | \theta^*, X_t)) \end{aligned}$$

390 Define  $\bar{r}_t(a_t, X_{t,S_t^*}) = \mathbb{E}_t[r_t(\theta^*, a_t, X_{t,S_t^*})]$  and  $g(p|q) = p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q}$ .

391 The Fenchel-Young inequality states:  $f(x) + f^*(y) \geq \langle x, y \rangle$ , where  $f^*$  is the convex conjugate of  $f$ :  
 392  $f^*(a) = \sup_{x \in \mathcal{X}} \{ \langle a, x \rangle - f(x) \}$ . Hence,  $g^*(u|q) = \sup_{p \in [0,1]} \{ pu - g(p|q) \} = \log(1 + q(e^u -$   
 393  $1)) \leq q \left( u + \frac{u^2}{2} \right)$ .

$$\begin{aligned} \mathbb{E}_t[\Delta_t | X_{t,S_t}] &= \mathbb{E}_{t, X_{t,S_t}} [r_t(\theta_t, A_t, X_{t,S_t}) - r_t(\theta^*, A_t, X_{t,S_t})] \\ &= \mathbb{E}_{t, X_{t,S_t}} [r_t(\theta_t, A_t, X_{t,S_t}) - \bar{r}_t(A_t, X_{t,S_t})] \\ &= \mathbb{E}_{t, X_{t,S_t}} \left[ \sum_{a \in \mathcal{A}} \mathbb{I}[A_t = a] \frac{\eta p_t(A_t = a | X_{t,S_t})}{\eta p_t(A_t = a | X_{t,S_t})} r_t(\theta_t, a, X_{t,S_t}) - \bar{r}_t(A_t, X_{t,S_t}) \right] \\ &= \mathbb{E}_{t, X_{t,S_t}} \left[ \eta \sum_{a \in \mathcal{A}} p_t(A_t = a | X_{t,S_t}) \left( \frac{\mathbb{I}[A_t = a]}{\eta p_t(A_t = a | X_{t,S_t})} r_t(\theta_t, a, X_{t,S_t}) \right) - \bar{r}_t(A_t, X_{t,S_t}) \right] \\ &\leq \mathbb{E}_{t, X_{t,S_t}} \left[ \eta \sum_{a \in \mathcal{A}} p_t(A_t = a | X_{t,S_t}) \left( g(r_t(\theta_t, a, X_{t,S_t}) | \bar{r}_t(A_t, X_{t,S_t})) \right) \right] \end{aligned}$$

$$\begin{aligned}
& + g^* \left( \frac{\mathbb{1}[A_t = a]}{\eta p_t(A_t = a | X_t, S_t)} \|\bar{r}_t(A_t, X_t, S_t)\| \right) - \bar{r}_t(A_t, X_t, S_t) \Big] \\
& \leq \mathbb{E}_{t, X_t, S_t} \left[ \eta \sum_{a \in \mathcal{A}} p_t(A_t = a | X_t, S_t) \left( g(r_t(\theta_t, a, X_t, S_t) | \bar{r}_t(A_t, X_t, S_t)) \right. \right. \\
& \quad \left. \left. + \frac{\mathbb{1}[A_t = a]}{\eta p_t(A_t = a | X_t, S_t)} \bar{r}_t(A_t, X_t, S_t) + \frac{\mathbb{1}[A_t = a]}{2\eta^2 p_t^2(A_t = a | X_t, S_t)} \bar{r}_t(A_t, X_t, S_t) \right) - \bar{r}_t(A_t, X_t, S_t) \right] \\
& = \mathbb{E}_{t, X_t, S_t} \left[ \eta \sum_{a \in \mathcal{A}} p_t(A_t = a | X_t, S_t) g(r_t(\theta_t, a, X_t, S_t) | \bar{r}_t(A_t, X_t, S_t)) + \frac{1}{2\eta} \sum_{a \in \mathcal{A}} \bar{r}_t(a, X_t, S_t) \right] \\
& = \mathbb{E}_{t, X_t, S_t} \left[ \eta \sum_{a \in \mathcal{A}} p_t(A_t = a | X_t, S_t) g(r_t(\theta^*, a, X_t, S_t) | \bar{r}_t(A_t, X_t, S_t)) + \frac{1}{2\eta} \sum_{a \in \mathcal{A}} \bar{r}_t(a, X_t, S_t) \right] \\
& = \eta \mathbb{I}_t(\theta^*; Y_{t, A_t} | A_t, X_t, S_t) + \frac{1}{2\eta} \sum_{a \in \mathcal{A}} \mathbb{E}_t[\bar{r}_t(a, X_t, S_t)]
\end{aligned}$$

394 Choosing the value of  $\eta > 0$  which minimizes the above expression, we obtain

$$\mathbb{E}_t[\Delta_t | X_t, S_t] \leq \sqrt{2\mathbb{I}_t(\theta^*; Y_{t, A_t} | A_t, X_t, S_t) \sum_{a \in \mathcal{A}} \mathbb{E}_t[\bar{r}_t(a, X_t, S_t)]}$$

395 Rearranging terms, we get  $\rho_t(S_t) = \frac{(\mathbb{E}_t[\Delta_t | X_t, S_t])^2}{\mathbb{I}_t(\theta^*; Y_{t, A_t} | A_t, X_t, S_t)} \leq 2 \sum_{a \in \mathcal{A}} \mathbb{E}_t[\bar{r}_t(a, X_t, S_t)] \quad \square$

396 The following Lemma is a precursor to the bound on conditional mutual information  
397 (Lemma 7.6)

398 **Lemma 7.5** (Recursive Property of Posterior Sampling).  $\prod_{t=1}^T \sum_{\theta} q_t(\theta) p_t(Y_{t, A_t} | \theta, A_t, X_t, S_t) =$   
399  $\sum_{\theta} q_1(\theta) \prod_{t=1}^T p_t(Y_{t, A_t} | \theta, A_t, X_t, S_t)$

400 *Proof.* We repeat the proof in Neu et al. [2022] for easy reference:

401 From posterior sampling,  $q_{t+1}(\theta) = \frac{q_t(\theta) p_t(Y_{t, A_t} | \theta, A_t, X_t, S_t)}{\sum_{\theta'} q_t(\theta') p_t(Y_{t, A_t} | \theta', A_t, X_t, S_t)}$

$$\begin{aligned}
\sum_{\theta} q_1(\theta) \prod_{t=1}^T p_t(Y_{t, A_t} | \theta, A_t, X_t, S_t) &= \prod_{t=1}^T \frac{\sum_{\theta} q_1(\theta) \prod_{k=1}^t p_k(Y_{k, A_k} | \theta, A_k, X_{k, S_k})}{\sum_{\theta'} q_1(\theta') \prod_{k=1}^{t-1} p_k(Y_{k, A_k} | \theta', A_k, X_{k, S_k})} \\
&= \prod_{t=1}^T \sum_{\theta} q_t(\theta) p_t(Y_{t, A_t} | \theta, A_t, X_t, S_t)
\end{aligned}$$

402 □

403 **Lemma 7.6** (Bounded Conditional Mutual Information). *The cumulative mutual information for any*  
404 *observation set  $S_t$  admits the following bound:*  $\mathbb{E} \left[ \sum_{t=1}^T \mathbb{I}_t(\theta^*; Y_{t, A_t} | A_t, X_t, S_t) \right] \leq \mathbb{H}(\theta^*)$

405 *Proof.* We generalize the proof of Neu et al. [2022] to an arbitrary subset of features  $S_t$ . Let  
406  $q_{t+1}(\theta) = \frac{q_t(\theta) p_t(Y_{t, A_t} | \theta, A_t, X_{t, S_t}^*)}{\sum_{\theta'} q_t(\theta') p_t(Y_{t, A_t} | \theta', A_t, X_{t, S_t}^*)}$  be the posterior sampling update of our belief of the value  
407 of optimal parameter  $\theta^*$ .  $q_1(\theta)$  is our initial prior before any interactions.

$$\mathbb{E} \left[ \sum_{t=1}^T \mathbb{I}_t(\theta^*; Y_{t, A_t} | A_t, X_t, S_t) \right] = \mathbb{E} \left[ \sum_{t=1}^T \mathbb{E}_t D_{KL} [p_t(Y_{t, A_t} | \theta^*, A_t, X_{t, S_t}) | p_t(Y_{t, A_t} | A_t, X_{t, S_t})] \right]$$

$$\begin{aligned}
&= \mathbb{E} \left[ \sum_{t=1}^T \mathbb{E}_t \left[ \log \frac{p_t(Y_{t,A_t} | \theta^*, A_t, X_{t,S_t})}{p_t(Y_{t,A_t} | A_t, X_{t,S_t})} \right] \right] \\
&= \mathbb{E} \left[ \sum_{t=1}^T \mathbb{E}_t \left[ \log \frac{p_t(Y_{t,A_t} | \theta^*, A_t, X_{t,S_t})}{\sum_{\theta} q_t(\theta) p_t(Y_{t,A_t} | \theta, A_t, X_{t,S_t})} \right] \right] \\
&= \mathbb{E} \left[ \log \frac{\prod_{t=1}^T p_t(Y_{t,A_t} | \theta^*, A_t, X_{t,S_t})}{\sum_{\theta} q_1(\theta) \prod_{t=1}^T p_t(Y_{t,A_t} | \theta, A_t, X_{t,S_t})} \right] \text{ (by Lemma 7.5)} \\
&\leq \mathbb{E} \left[ \log \frac{\prod_{t=1}^T p_t(Y_{t,A_t} | \theta^*, A_t, X_{t,S_t})}{q_1(\theta^*) \prod_{t=1}^T p_t(Y_{t,A_t} | \theta^*, A_t, X_{t,S_t})} \right] \\
&= \mathbb{E}[-\log(q_1(\theta^*))] = \mathbb{H}(\theta^*)
\end{aligned}$$

408

□

409 **Theorem 7.7** (Regret of Exact Entropy Minimization). *Exact entropy minimization achieves the*  
410 *regret bound:  $R_T \leq \sqrt{2T\mathbb{H}(\theta^*)}$*

411 *Proof.* In the exact entropy minimization setting we assume  $S_t \supseteq S_t^*$ , which is the smallest entropy-  
412 minimizing subset of user features, and therefore, the information ratio  $\rho_t(S_t) = \rho_t(S_t^*)$  can simply  
413 be denoted as  $\rho_t$ .

$$\rho_t = \frac{(\mathbb{E}_t [r_t(\theta^*, A_t^*, X_t) - r_t(\theta^*, A_t, X_{t,S_t^*})])^2}{\mathbb{I}_t(\theta^*; Y_{t,A_t} | A_t, X_{t,S_t^*})}$$

414 The cumulative regret of the decision maker can be bounded

$$\begin{aligned}
R_T &= \mathbb{E} \left[ \sum_{t=1}^T (r_t(\theta^*, A_t^*, X_t) - r_t(\theta^*, A_t, X_{t,S_t^*})) \right] = \mathbb{E} \left[ \sum_{t=1}^T \mathbb{E}_t [r_t(\theta^*, A_t^*, X_t) - r_t(\theta^*, A_t, X_{t,S_t^*})] \right] \\
&= \mathbb{E} \left[ \sum_{t=1}^T \sqrt{\rho_t \mathbb{I}_t(\theta^*; Y_{t,A_t} | A_t, X_{t,S_t^*})} \right] \leq \sqrt{\mathbb{E} \left[ \sum_{t=1}^T \rho_t \right] \mathbb{E} \left[ \sum_{t=1}^T \mathbb{I}_t(\theta^*; Y_{t,A_t} | A_t, X_{t,S_t^*}) \right]}
\end{aligned}$$

415 The last step follows by Cauchy-Schwarz Inequality. We proceed by deriving bounds for each of  
416  $\mathbb{E} \left[ \sum_{t=1}^T \rho_t \right]$  and  $\mathbb{E} \left[ \sum_{t=1}^T \mathbb{I}_t(\theta^*; Y_{t,A_t} | A_t, X_{t,S_t^*}) \right]$ .

417 To bound  $\mathbb{E} \left[ \sum_{t=1}^T \rho_t \right]$ , we invoke Lemma 3.5:

$$\rho_t \leq 2\mathbb{E} \left[ \sum_{a \in \mathcal{A}} \bar{r}_t(a, X_{t,S_t}) \right] = 2 \quad (\text{categorical distribution over actions})$$

418 Next, we invoke Lemma 7.6 which says  $\mathbb{E} \left[ \sum_{t=1}^T \mathbb{I}_t(\theta^*; Y_{t,A_t} | A_t, X_{t,S_t^*}) \right] \leq \mathbb{H}(\theta^*)$ .

419 Combining these two bounds we obtain:

$$R_T \leq \sqrt{\mathbb{E} \left[ \sum_{t=1}^T \rho_t \right] \mathbb{E} \left[ \sum_{t=1}^T \mathbb{I}_t(\theta^*; Y_{t,A_t} | A_t, X_{t,S_t^*}) \right]} \leq \sqrt{2T\mathbb{H}(\theta^*)}$$

420

□

421 **Theorem 7.8** (Regret). *Approximate entropy minimization achieves the regret bound:  $R_T \leq$*   
422  *$\sqrt{2T\mathbb{H}(\theta^*)} + \mathcal{O}(T^{1-\alpha/2})$  in the one-hot outcome setting.*

423 *We assume that entropy minimization learns to select the optimal subset with diminishing error  $\epsilon_t =$*   
424  *$\mathbb{H}_t(Y_t | X_{t,S_t}, \theta) - \mathbb{H}_t(Y_t | X_t, \theta)$ ,  $\mathbb{E}[\epsilon_t] = \mathcal{O}(\frac{1}{t^\alpha})$ . In particular, for  $\alpha = 1$ , we obtain  $R_T = \mathcal{O}(\sqrt{T})$ ,*  
425 *identical to the exact entropy minimization setting.*

426 *Proof.* Let us define  $\Delta_t = r_t(\theta^*, A_t^*, X_t) - r_t(\theta^*, A_t, X_t)$  and  $\tilde{\Delta}_t = r_t(\theta_t, A_t, X_{t,S_t}) -$   
 427  $r_t(\theta^*, A_t, X_{t,S_t})$ . We use the definition of information ratio from Lemma 3.5 and define  $\tilde{\rho}_t = \rho_t(S_t)$ :

$$\tilde{\rho}_t = \frac{(\mathbb{E}_t [\tilde{\Delta}_t | X_{t,S_t}])^2}{\mathbb{I}_t(\theta^*; Y_{t,A_t} | A_t, X_{t,S_t})}$$

428 Define  $\gamma_t = \mathbb{E}_t [\Delta_t | X_{t,S_t^*}] - \mathbb{E}_t [\tilde{\Delta}_t | X_{t,S_t}]$ . Regret  $R_T$  can be decomposed as follows:

$$\begin{aligned} R_T &= \mathbb{E} \left[ \sum_{t=1}^T \Delta_t \right] = \mathbb{E} \left[ \sum_{t=1}^T \mathbb{E}_t [\Delta_t | X_{t,S_t^*}] \right] = \mathbb{E} \left[ \sum_{t=1}^T \mathbb{E}_t [\Delta_t | X_{t,S_t^*}] - \mathbb{E}_t [\tilde{\Delta}_t | X_{t,S_t}] + \mathbb{E}_t [\tilde{\Delta}_t | X_{t,S_t}] \right] \\ &= \mathbb{E} \left[ \sum_{t=1}^T \gamma_t + \sqrt{\tilde{\rho}_t \mathbb{I}_t(\theta^*; Y_{t,A_t} | A_t, X_{t,S_t})} \right] \\ &\leq \mathbb{E} \left[ \sum_{t=1}^T \gamma_t \right] + \sqrt{\mathbb{E} \left[ \sum_{t=1}^T \tilde{\rho}_t \right] \mathbb{E} \left[ \sum_{t=1}^T \mathbb{I}_t(\theta^*; Y_{t,A_t} | A_t, X_{t,S_t}) \right]} \\ &\leq \mathbb{E} \left[ \sum_{t=1}^T \gamma_t \right] + \sqrt{2T\mathbb{H}(\theta^*)} \end{aligned}$$

429 Where we bound the second term in a similar manner to Theorem 3.6.  $\mathbb{E} \left[ \sum_{t=1}^T \tilde{\rho}_t \right] \leq 2T$  from

430 Lemma 3.5.  $\mathbb{E} \left[ \sum_{t=1}^T \mathbb{I}_t(\theta^*; Y_{t,A_t} | A_t, X_{t,S_t}) \right] \leq \mathbb{H}(\theta^*)$  from Lemma 7.6.

431 Now we bound  $\gamma_t = \mathbb{E}_t [\Delta_t | X_{t,S_t^*}] - \mathbb{E}_t [\tilde{\Delta}_t | X_{t,S_t}]$ :

432 Note that, under the exact entropy search setting, we have  $S_t = S_t^*$  which yields  $\gamma_t = 0$ .

$$\begin{aligned} \gamma_t &= \mathbb{E}_t [\Delta_t | X_{t,S_t^*}] - \mathbb{E}_t [\tilde{\Delta}_t | X_{t,S_t}] \\ &= r_t(\theta^*, A_t^*, X_{t,S_t^*}) - r_t(\theta^*, A_t, X_{t,S_t^*}) - (r_t(\theta_t, A_t, X_{t,S_t}) - r_t(\theta^*, A_t, X_{t,S_t})) \\ &= r_t(\theta_t, A_t, X_{t,S_t^*}) - r_t(\theta^*, A_t, X_{t,S_t^*}) - (r_t(\theta_t, A_t, X_{t,S_t}) - r_t(\theta^*, A_t, X_{t,S_t})) \\ &\text{(Conditional distribution of } (\theta^*, A_t^*) \text{ and } (\theta_t, A_t) \text{ are same)} \\ &= r_t(\theta_t, A_t, X_{t,S_t^*}) - r_t(\theta_t, A_t, X_{t,S_t}) - (r_t(\theta^*, A_t, X_{t,S_t^*}) - r_t(\theta^*, A_t, X_{t,S_t})) \end{aligned}$$

433 Let  $\delta_t(\theta) = r_t(\theta, A_t, X_{t,S_t^*}) - r_t(\theta, A_t, X_{t,S_t}) = p_t(Y_{t,A_t} = 1 | X_{t,S_t^*}, \theta) - p_t(Y_{t,A_t} = 1 | X_{t,S_t}, \theta)$ .

434 Then, we have  $\gamma_t = \delta_t(\theta_t) - \delta_t(\theta^*)$ .

$$\begin{aligned} |\delta_t(\theta)| &= |p_t(Y_{t,A_t} = 1 | X_{t,S_t^*}, \theta) - p_t(Y_{t,A_t} = 1 | X_{t,S_t}, \theta)| \\ &\leq \sup_{a \in \mathcal{A}} |p_t(Y_{t,a} = 1 | X_{t,S_t^*}, \theta) - p_t(Y_{t,a} = 1 | X_{t,S_t}, \theta)| \\ &= \sup_{Y_t} |p_t(Y_t | X_{t,S_t^*}, \theta) - p_t(Y_t | X_{t,S_t}, \theta)| \quad \text{(Total Variation (TV) distance)} \\ &\leq \sqrt{\frac{1}{2}} \epsilon_t \quad \text{(by Lemma 3.3)} \end{aligned}$$

435 Where we assume that  $\epsilon_t = \mathbb{H}_t(Y_t | X_{t,S_t}, \theta) - \mathbb{H}_t(Y_t | X_{t,S_t^*}, \theta)$  is bounded.

436 Now we derive a bound on  $\mathbb{E} \left[ \sum_{t=1}^T \gamma_t \right]$ .

$$\mathbb{E} \left[ \sum_{t=1}^T \gamma_t \right] = \mathbb{E} \left[ \sum_{t=1}^T \delta_t(\theta_t) - \delta_t(\theta^*) \right] \leq \mathbb{E} \left[ \sum_{t=1}^T |\delta_t(\theta_t)| + |\delta_t(\theta^*)| \right]$$

$$\leq \sqrt{2} \mathbb{E} \left[ \sum_{t=1}^T \sqrt{\epsilon_t} \right]$$

437 Assuming that the entropy gap decays as  $\mathbb{E}[\epsilon_t] = \mathcal{O}\left(\frac{1}{t^\alpha}\right)$ , we achieve the following bound:

$$\mathbb{E} \left[ \sum_{t=1}^T \gamma_t \right] = \mathcal{O}\left(\sum_{t=1}^T \frac{1}{t^{\alpha/2}}\right) = \mathcal{O}(T^{1-\alpha/2})$$

438 Replacing this in the regret bound we obtain:

$$R_T \leq \mathcal{O}(T^{1-\alpha/2}) + \sqrt{2T\mathbb{H}(\theta^*)} = \mathcal{O}(\sqrt{T}) \quad (\text{for } \alpha = 1)$$

439

□

**Lemma 7.9** (Entropy Minimization Lowers Per-Question Suboptimality over ReAct). *We define the per-question suboptimality of an algorithm as  $\epsilon_t^q = \mathbb{H}_t(Y_t|X_{t,S_t^{<q}}, X_{t,S_t^q}) - \mathbb{H}_t(Y_t|X_{t,S_t^{<q}}, X_{t,\bar{S}_t^q})$  where  $\bar{S}_t^q$  is the optimal  $q^{\text{th}}$  question that minimizes the outcome entropy. Let  $p_{t,q}^* = p_t(S_t^q = \bar{S}_t^q|\omega, X_{t,S_t^{<q}})$  be the probability that the policy selects the optimal  $q^{\text{th}}$  follow-up question. For a budget of  $k \ll K$  candidate questions in Entropy-Search (ES), we achieve lower suboptimality than the ReAct baseline:*

$$(\epsilon_t^q)_{\text{ES}} \leq (1 - p_{t,q}^*)^{k-1} (\epsilon_t^q)_{\text{ReAct}}$$

440 *Proof.* Let  $\bar{S}_t^q = \arg \min_{S_t^q} \mathbb{H}_t(Y_t|X_{t,S_t^{<q}}, X_{t,S_t^q})$  be the optimal next follow-up question to the  
 441 user for a given observed context  $X_{t,S_t^{<q}}$ . The per-question suboptimality of an algorithm  $\epsilon_t^q =$   
 442  $\mathbb{H}_t(Y_t|X_{t,S_t^{<q}}, X_{t,S_t^q}) - \mathbb{H}_t(Y_t|X_{t,S_t^{<q}}, X_{t,\bar{S}_t^q})$  measures how much more entropy could have been  
 443 minimized if the optimal question was selected.

444 Under the ReAct baseline, we sample a follow-up question  $S_t^q \sim p_t(S_t^q|\omega, X_{t,S_t^{<q}})$ . Let  $p_{t,q}^* =$   
 445  $p_t(S_t^q = \bar{S}_t^q|\omega, X_{t,S_t^{<q}})$  be the probability that the policy selects the optimal  $q^{\text{th}}$  follow-up question.

446 For compactness, let us use subscript  $(t, q)$  to denote conditioning over  $(\mathcal{F}_t, X_{t,S_t^{<q}})$ . Then the per-  
 447 question suboptimality of ReAct can be re-written as  $(\epsilon_t^q)_{\text{ReAct}} = \mathbb{H}_{t,q}(Y_t|X_{t,S_t^q}) - \mathbb{H}_{t,q}(Y_t|X_{t,\bar{S}_t^q})$ .

$$\begin{aligned} (\epsilon_t^q)_{\text{ReAct}} &= \mathbb{H}_{t,q}(Y_t|X_{t,S_t^q}) - \mathbb{H}_{t,q}(Y_t|X_{t,\bar{S}_t^q}) \\ &= \left( \sum_{S \neq \bar{S}_t^q} p_{t,q}(S_t^q = S|\omega) \mathbb{H}_{t,q}(Y_t|X_{t,S}) + (p_{t,q}^*) \mathbb{H}_{t,q}(Y_t|X_{t,\bar{S}_t^q}) \right) - \mathbb{H}_{t,q}(Y_t|X_{t,\bar{S}_t^q}) \\ &= (1 - p_{t,q}^*) \frac{\sum_{S \neq \bar{S}_t^q} p_{t,q}(S_t^q = S|\omega) \mathbb{H}_{t,q}(Y_t|X_{t,S})}{1 - p_{t,q}^*} + (p_{t,q}^* - 1) \mathbb{H}_{t,q}(Y_t|X_{t,\bar{S}_t^q}) \\ &= (1 - p_{t,q}^*) \mathbb{H}_{t,q}^{\text{sub}} + (p_{t,q}^* - 1) \mathbb{H}_{t,q}^{\text{min}} = (1 - p_{t,q}^*) (\mathbb{H}_{t,q}^{\text{sub}} - \mathbb{H}_{t,q}^{\text{min}}) \end{aligned}$$

448 where  $\mathbb{H}_{t,q}^{\text{sub}}$  is the expected outcome entropy in suboptimal cases,  $\mathbb{H}_{t,q}^{\text{min}}$  is the minimum outcome  
 449 entropy possible after selecting the  $q^{\text{th}}$  follow-up question.

450 Now let us compare this with the per-question suboptimality for Entropy Search (ES). Under ES,  
 451 we first sample  $k$  candidate follow-up questions  $S_1, S_2, \dots, S_k \sim p_{t,q}(S_t^q|\omega)$ . Then amongst these  
 452 candidates, we select the entropy minimizing question  $S_t^q = \arg \min_{S \in \{S_1, S_2, \dots, S_k\}} \mathbb{H}_{t,q}(Y_t|X_{t,S})$ .  
 453 Assume that the primary limiting factor for ES is the budget  $k \ll K$ : with probability  $(1 - p_{t,q}^*)^k$ ,  
 454 the candidate set does not contain  $\bar{S}_t^q$  which is the entropy minimizing solution. Then we have:

$$\begin{aligned} (\epsilon_t^q)_{\text{ES}} &= \mathbb{H}_{t,q}(Y_t|X_{t,S_t^q}) - \mathbb{H}_{t,q}(Y_t|X_{t,\bar{S}_t^q}) \\ &\leq \left( (1 - p_{t,q}^*)^k \mathbb{H}_{t,q}^{\text{sub}} + (1 - (1 - p_{t,q}^*)^k) \mathbb{H}_{t,q}(Y_t|X_{t,\bar{S}_t^q}) \right) - \mathbb{H}_{t,q}(Y_t|X_{t,\bar{S}_t^q}) \end{aligned}$$

(In suboptimal cases, ES still achieves lower entropy than ReAct as it minimizes entropy over the  $k$  subset)

$$= (1 - p_{t,q}^*)^k (\mathbb{H}_{t,q}^{\text{sub}} - \mathbb{H}_{t,q}^{\text{min}})$$

Combining the results for per-turn suboptimality for ReAct and ES, we obtain:

$$(\epsilon_t^q)_{\text{ES}} \leq (1 - p_{t,q}^*)^{k-1} (\epsilon_t^q)_{\text{ReAct}}$$

455

□

456 *Discussion:* Selecting optimal follow-up questions is crucial because observing the answer to these  
 457 questions can help the model predict the correct action in fewer questions and more accurately. We  
 458 observe in Lemma 3.8 that the per-question suboptimality of Entropy Search depends on the compute  
 459 budget  $k$  in the exponent. Therefore, the more candidate questions we consider for inference-time  
 460 search, the more likely we are to select the optimal follow-up question, as we would expect. Moreover,  
 461 this bound assumes that both ES and ReAct share the same question-selection policy  $\omega$ . However, we  
 462 expect  $p_{t,q}^*$  to grow much faster for ES than for ReAct since (1) inference time selection based on  
 463 the entropy score  $v_t^q(S) = -\mathbb{H}_t(Y_t | X_{t,S_t^{<q}}, X_{t,S})$  increases the frequency of selecting the correct  
 464 follow-up question, and (2)  $v_t^q(S)$  serves as a dense, per-step training signal to optimize  $\omega$ , that  
 465 directly maximizes  $p_{t,q}^*$  whenever a correct follow-up question is selected in the candidate set. As a  
 466 result, we expect ES to select optimal queries much more often than ReAct, leading to better actions  
 467 and, therefore, better user outcomes.

## 468 8 Baseline Methods

469 Now we compare several baseline strategies for information gathering with language models, some  
 470 of which we compare to ENTROPY SEARCH in our experiments (Section 4).

### 471 8.1 Random Sampling

472 A naive alternative is to start with random question selection and rely on observed outcomes to guide  
 473 future decisions. This approach assumes that information-seeking behavior will emerge naturally  
 474 through reward maximization over a sufficiently large number of interactions. While conceptually  
 475 simple, such a strategy has notable drawbacks, particularly in its inefficiency due to the lack of  
 476 prioritization of informative questions. To illustrate the limitations of random sampling, consider the  
 477 following example: Let  $B \ll K$  be the budget for total questions the agent is allowed to ask per user  
 478 before performing an action, where  $K$  is the total number of features in the user’s context. Assume  
 479 that the user’s reward model depends on a sparse subset of  $d \ll K$  relevant features. The probability  
 480 that random sampling selects all  $d$  relevant features in a single round is:

$$\mathbb{P}(|S_t \cap S_t^*| = d) = \frac{\binom{d}{d} \binom{K-d}{B-d}}{\binom{K}{B}},$$

481 where  $S_t^*$  is the minimal set of relevant context features. This probability becomes exponentially  
 482 small for large  $K$  and small  $d$ . Consequently, the agent is unlikely to collect sufficient context to  
 483 make optimal decisions, leading to a high likelihood of suboptimal outcomes, hence incurring linear  
 484 regret. This sparse signal provided by random sampling exacerbates the challenges for reinforcement  
 485 learning algorithms, which depend on meaningful feedback to improve their policies. Without  
 486 effective strategies for sampling follow-up questions, the agent may fail to uncover the structure of  
 487 the reward model, limiting the potential for learning and adaptation.

### 488 8.2 REACT

489 A natural improvement to random follow-up questions is to utilize the LLM agent’s world knowledge  
 490 to come up with meaningful and relevant questions. REACT is a popular agentic framework that  
 491 prompts the model to *Reason then Act*. The agent is prompted with the user’s questions and the  
 492 context gathered through previous follow-up questions and is asked to produce Chain-of-Thought  
 493 reasoning prior to selecting the next best question. This approach is better than random sampling,  
 494 since the model can reason about what facts and features are relevant to the problem and selectively  
 495 query those from the user. However, a key limitation of REACT is that the model is unable to  
 496 look far ahead into the future and anticipate the user’s potential responses to subsequent follow-up  
 497 questions  $S_t^{>q}$  while selecting the current follow-up question  $S_t^q$ . However, as we show in Section 3,  
 498 such System-2 thinking is crucial for making optimal decisions and reliably estimating the outcome  
 499 distribution. This is particularly important for problem domains where there are complex conditional

500 interdependencies between different follow-up questions where some questions are strictly more  
 501 informative than others and the answer to these questions can make other questions irrelevant. We  
 502 demonstrate that explicitly searching over candidate questions significantly lowers the per-question  
 503 suboptimality of ENTROPY SEARCH when compared to REACT in Lemma 7.9.

504 Now we discuss another Lemma that gives a lower bound on the suboptimality  $\epsilon_t$  for REACT  
 505 and ENTROPY SEARCH and we show that this gap decreases significantly faster for ENTROPY  
 506 SEARCH.

507 **Lemma 8.1** (Lower Bound on Suboptimality Gap). *The entropy gap (suboptimality)  $\epsilon_t =$   
 508  $\mathbb{H}_t(Y_t|X_{t,S_t}) - \mathbb{H}_t(Y_t|X_t)$  admits the following lower bound for REACT and ENTROPY SEARCH  
 509 under a budget of  $B$  follow-up questions with  $R$  underlying entropy-minimizing questions.*

$$\mathbb{E}[\epsilon_t^{\text{React}}] \geq \mathbb{H}_{\text{sub}} \sum_{r=0}^{R-1} \binom{B}{r} p^r (1-p)^{B-r}$$

$$\mathbb{E}[\epsilon_t^{\text{ES}}] \geq \mathbb{H}_{\text{sub}} \sum_{r=0}^{R-1} \binom{B}{r} \bar{p}^r (1-\bar{p})^{B-r}$$

510  $p$  denotes the maximum probability that REACT selects a good follow-up question and  $\bar{p} = 1 - (1-p)^k$   
 511 is the corresponding upper bound on probability for ENTROPY SEARCH under a search budget of  $k$   
 512 candidate questions.  $\mathbb{H}_{\text{sub}}$  denotes the minimum entropy gap over suboptimal trajectories that are  
 513 missing one or more key questions about the user.

514 *Proof.* Let  $B$  be the budget for maximum number of follow-up questions that can be asked to  
 515 the user. Let  $R$  be the minimum number of queries needed to minimize the outcome entropy  
 516  $\mathbb{H}_t(Y_t|X_{t,S_t})$ . In general, there may be multiple queries that are equivalent in the sense that the  
 517 corresponding context features are highly correlated and knowing one is as good as knowing another.  
 518 Let  $S_t^* = \arg \min_S \mathbb{H}_t(Y_t|X_{t,S})$  be an entropy-minimizing subset of questions. At any given turn  $q$ ,  
 519 the agent picks a question  $S_t^q$ , while  $q \leq B$ . A query  $S_t^q$  is entropy-reducing if it is equivalent to any  
 520 of the queries in  $S_t^*$  and is not equivalent to any previously asked queries  $S_t^{<q}$ .

521 Let the probability that  $S_t^q$  is an entropy-minimizing query be  $\leq p$ . Then the probability that the  
 522 set of questions asked by the agent  $S_t \supseteq S_t^*$  (upto any equivalent questions) is  $p_t(\epsilon_t = 0) \leq$   
 523  $\sum_{r=R}^B \binom{B}{r} p^r (1-p)^{B-r}$ . This follows from the fact that this is an increasing function of  $p$  and the  
 524 actual probabilities of selecting an entropy reducing question at any given turn  $q$  are lower than  $p$ ,  
 525 which means  $p_t(\epsilon_t = 0)$  would be lower.

526 Now let us turn our attention to the case where a suboptimal set of questions is asked within the  
 527 budget  $B$ . This means that atleast one question from  $S_t^*$  or equivalent has not been asked by the  
 528 agent. Let  $\mathbb{H}_{\text{sub}} = \min_{S_t \not\supseteq S_t^*} \mathbb{H}_t(Y_t|X_{t,S_t}) = \min_{S_t \not\supseteq S_t^*} \epsilon_t$  be the minimum suboptimality gap in  
 529 the instances where the agent does not cover all entropy minimizing questions. This happens with  
 530 probability  $\geq \sum_{r=0}^{R-1} \binom{B}{r} p^r (1-p)^{B-r}$  because this quantity is decreasing in  $p$  (intuitively, this is the  
 531 CDF of a binomial distribution, and as  $p$  increases, the mass of the distribution shifts right, reducing  
 532 the CDF). Thus, it follows that  $\mathbb{E}[\epsilon_t^{\text{React}}] \geq \mathbb{H}_{\text{sub}} \sum_{r=0}^{R-1} \binom{B}{r} p^r (1-p)^{B-r}$ .

533 In contrast to REACT, the probability of selecting an entropy minimizing question is directly linked to  
 534 the computational budget for search which determines how many candidates questions  $k$  are sampled  
 535 from the agent policy  $p_t(S_t^q|\omega, X_{t,S_t^{<q}})$ . The probability that  $S_t^q$  is an entropy-minimizing query is  
 536  $\leq \bar{p} = 1 - (1-p)^k$ . The corresponding bound on suboptimality for ENTROPY SEARCH is given by  
 537  $\mathbb{E}_t[\epsilon_t^{\text{ES}}] \geq \mathbb{H}_{\text{sub}} \sum_{r=0}^{R-1} \binom{B}{r} \bar{p}^r (1-\bar{p})^{B-r}$ .

538 We consider a simple case where  $B = 50, R = 25$  and plot the upper bound on  $p_t(\epsilon_t = 0)$  against  
 539  $p$  for different values of  $k$ . We can see in Figure 5 that even for moderate values of  $k$ , ENTROPY  
 540 SEARCH achieves significantly lower  $\epsilon_t$  (practically zero) even when REACT selects an entropy-  
 541 minimizing question  $S_t^q$  with  $p = 0.5$  probability. This shows that ENTROPY SEARCH effectively  
 542 lowers  $\mathbb{E}_t[\epsilon_t]$  by utilizing test-time compute. Moreover, we expect  $\mathbb{E}_t[\epsilon_t]$  to approach zero much  
 543 faster even with marginal improvements in  $p$  as we update  $\omega$  which can be seen from Figure 5.  $\square$

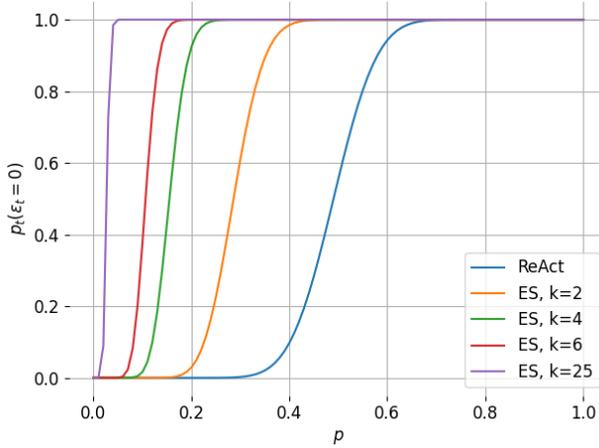


Figure 5: Upper bound,  $p_t(\epsilon_t = 0) \leq \sum_{r=R}^B \binom{B}{r} p^r (1-p)^{B-r}$  vs  $p$ : REACT is  $k = 1$ , ENTROPY SEARCH is  $k > 1$

544 **8.3 REFLEXION**

545 REFLEXION Shinn et al. [2023] builds on top of REACT and adds a dedicated reflection component  
 546 before the reasoning block in REACT. The agent is explicitly prompted to reflect on what it has  
 547 learned from the previous interactions and assess the effectiveness of its information gathering  
 548 strategy. The prompt also includes past reflections for improved continuity between questions. As  
 549 seen in Section 4, REFLEXION does slightly better than REACT due to added backward reflection.  
 550 However, as we demonstrate with ENTROPY SEARCH, looking forward is crucial for optimal question  
 551 selection.

552 **9 Benchmarks**

553 **GSM8K-Q** We utilized the GSM8K-Q benchmark from QuestBench [Li et al., 2025], which  
 554 contains human-annotated grade school math problems designed with one missing variable assign-  
 555 ment necessary for solving each problem. This benchmark provides clear, structured scenarios for  
 556 evaluating the agent’s ability to identify and acquire the minimal information required to resolve  
 557 underspecified reasoning tasks. The controlled yet challenging nature of GSM8K-Q makes it ideal  
 558 for rigorously assessing our algorithm’s capability to ask targeted, informative questions.

559 **1D-ARC** The 1D-ARC benchmark [Xu et al., 2024], a widely studied adaptation of the Abstraction  
 560 and Reasoning Corpus (ARC), is designed to evaluate abstract reasoning in language models. We  
 561 further adapt it to test sequential information gathering by masking the output sequence. Each  
 562 task presents an input sequence and a partially hidden output generated via transformations such as  
 563 recoloring based on token frequency, repeating patterns, dynamic shifts, and conditional replacements.  
 564 To uncover the transformation rule, the agent must query specific output indices, receiving exact  
 565 values from an oracle. This setup provides a controlled and interpretable environment for evaluating  
 566 how effectively an agent can reduce uncertainty and reason under partial information.

567 **Fermi Problems** The Fermi problems [Kalyan et al., 2021] evaluate algorithms through  
 568 commonsense-based estimation tasks, such as estimating coffee consumption at EMNLP 2019  
 569 or the potential rise in sea levels from melting ice. Problems are designed with an oracle that provides  
 570 answers from a mix of relevant and distractor facts, totaling 20 variables, but without indicating which  
 571 facts are necessary. Solving Fermi problems requires recursively breaking down complex queries,  
 572 abstracting details, and synthesizing multiple reasoning steps, aligning well with our interactive  
 573 decision-making framework. These problems challenge even advanced fine-tuned models, which  
 574 typically exhibit substantial errors; the best fine-tuned T5 model only reached a score of 0.23 [Kalyan  
 575 et al., 2021]. To ensure the necessity for interactive information gathering, we filtered out easily  
 576 solvable questions and manually removed inaccurately annotated ones, resulting in 1000 high-quality  
 577 tasks from an initial pool of 5000. Figure 2 illustrates an example interaction, where the ENTROPY  
 578 SEARCH agent selects informative questions to arrive at an accurate estimate.

**Algorithm 1** Interactive Question Answering

---

```

1: procedure INTERACTIVEQA(Problem, User, LLMAgent)
2:    $H \leftarrow \{\}$  ▷ Instantiate context
3:   for  $b$  to budget do
4:     question  $q = \text{LLMAgent}(\text{Problem}, H)$ 
5:     answer  $a = \text{User}(q)$  ▷ Query user
6:      $H \leftarrow H \cup \{(q, a)\}$  ▷ Append to context
7:   end for
8:   return action  $A = \text{LLMAgent}(\text{Problem}, H)$ 
9: end procedure

```

---

**Algorithm 2** Entropy Search

---

```

1: procedure ENTROPYSEARCH(problem, history  $H$ )
2:    $\mathcal{H} \leftarrow \{\}$ 
3:    $\mathcal{Q} = \text{LLMAgent.sample}(\text{problem}, \text{sample size} = k, H)$  ▷ Sample  $k$  questions
4:   for  $q$  in  $\mathcal{Q}$  do
5:     answers =  $\text{LLMAgent.sample}(q, \text{sample size}, \text{history})$  ▷ Sample answers
6:     for answer in answers do
7:       actions =  $\text{LLMAgent.sample}((q, \text{answer}), \text{sample size}, \text{history})$ 
8:     end for
9:      $\mathcal{H} \leftarrow \mathcal{H} \cup \{\text{ComputeAvgEntropy}(\text{actions} \mid q, \text{answer})\}$  ▷ Compute using Lemma 3.2
10:  end for
11:  return  $\arg \min_{\mathcal{Q}} \mathcal{H}$ 
12: end procedure

```

---

580 **11 Other Related Works**

581 Our work integrates ideas from contextual bandits and experiment design in a novel manner. In the  
582 following subsections, we examine relevant threads in prior art, including derived settings such as  
583 survey bandits and combinatorial bandits. We conclude by explaining how these ideas come together  
584 in ENTROPY SEARCH.

585 **11.1 Experiment Design**

586 Experimental design addresses the problem of selecting the most informative set of measurements to  
587 reveal a hidden quantity of interest. Prior work has approached this task from a decision-theoretic  
588 perspective [Chaloner and Verdinelli, 1995, DeGroot, 1962]. A central tool in this literature is  
589 Shannon information (entropy), which is equivalent to the expected KL divergence between the prior  
590 and the posterior after incorporating new information from the experiment. This criterion, also known  
591 as Expected Information Gain (EIG), has been widely adopted to select optimal experiments. It  
592 motivates our use of entropy minimization as the objective for selecting the next follow-up question  
593 in ENTROPY SEARCH.

594 **11.2 Contextual Bandits**

595 Contextual bandits are a class of Multi-Armed Bandit (MAB) problems where the decision maker  
596 observes a context  $X_t$  before selecting an action  $A_t$  and receiving the corresponding outcome  $Y_t$   
597 [Agrawal and Goyal, 2013]. The decision maker leverages both the current context and the history of  
598 past interactions—i.e., tuples  $(X_j, A_j, Y_j)_{j < t}$ —to estimate the optimal action. Thompson sampling  
599 is a popular strategy in this setting, known for its strong theoretical guarantees [Agrawal and Goyal,  
600 2013]. Neu et al. [2022] extend the information-theoretic analysis of regret developed in Russo and  
601 Van Roy [2016] to the contextual bandit case. Their key contribution is the notion of an information  
602 ratio, quantifying the trade-off between incurred regret and knowledge gained about the optimal action.  
603 Our work builds on this idea by introducing an analogous quantity in our setting (Lemma 3.5).

### 604 11.3 Combinatorial Bandits

605 Combinatorial bandit problems involve selecting combinations of arms, where the reward depends  
606 on the joint set of chosen arms [Durand and Gagné, 2014]. This is a more general setting than  
607 our information-gathering problem, as the outcome can be an arbitrary function of the selected arm  
608 features. In contrast, our setting has a key property: outcome uncertainty necessarily decreases as more  
609 features are added to the observed context subset. This enables the use of greedy selection strategies,  
610 such as ENTROPY SEARCH, which are significantly more tractable than general combinatorial  
611 optimization.

### 612 11.4 Survey Bandits

613 The survey bandit framework introduced by Krishnamurthy and Athey [2020] studies how to gather  
614 subsets of context before identifying the best arm in contextual bandits. Earlier, Bouneffouf et al.  
615 [2017] proposed a related problem—contextual bandits with restricted context—where the agent  
616 selects a fixed-size subset of the user context prior to choosing an arm. While inspired by combi-  
617 natorial bandits, their setting imposes additional assumptions on the reward distribution to make  
618 subset selection tractable. Related efforts in personalization have also explored feature relevance  
619 in contextual bandits [Tomkins et al., 2021, Yao et al., 2021]. However, a key limitation of these  
620 approaches is their reliance on strong structural assumptions about the context or reward, which  
621 restricts general applicability.

### 622 11.5 Combining Ideas

623 While prior work in sequential experimental design has focused on selecting sets of experiments to  
624 maximize information about outcomes, this approach has not been rigorously applied to the survey  
625 bandit setting. In reinforcement learning, information gain has been used as an intrinsic reward  
626 to encourage exploration [Houthoof et al., 2016, Mohamed and Jimenez Rezende, 2015]. To our  
627 knowledge, ENTROPY SEARCH is the first method to formalize entropy-based question selection in  
628 online reinforcement learning, particularly in the survey bandit context.

## 629 12 Chain-of-Thought & Min-Entropy Search Prompt

### 630 12.1 ARC-1D:

631 Task description:

```
632 You are solving a binary sequence puzzle. Problem: You have a {n}-bit input sequence: {  
633     ↪ input-seq}. A hidden output sequence was created by applying an unknown number of  
634     ↪ right circular shifts. A right circular shift moves the last bit to front and  
635     ↪ shifts all others right.
```

636 The previously gathered information follows this segment of the prompt. Here, we choose to include  
637 the previously gathered context in the initial prompt rather than in the chat history. We include this as  
638 follows:

```
639 Known output bits: {"(new line)".join(f"- Index {q} -> Value {v}" for q, v in seen)}
```

640 For Chain-of-Thought:

```
641 Instructions: Do not query the indices for which output is known.  
642 You have EXACTLY two options:  
643  
644 If you have sufficient information to determine the number of right circular shifts:  
645 1) Submit your final answer in this format:  
646 Reasoning: <step-by-step reasoning>  
647 Answer: <number of right circular shifts to transform the input into the output>  
648  
649 Else if you need to gather more information to find the number of right circular shifts:  
650 2) Ask a query in this format:  
651 Reasoning: <step-by-step reasoning>  
652 Query: <index 0-{n-1}>  
653  
654 Use EXACTLY one of these two formats, and NOTHING ELSE (no markdown, no extra text).
```

655 For MES:

```
656 Instructions: Do not query the indices for which output is known.  
657 You have EXACTLY two options:  
658  
659 If you have sufficient information to determine the number of right circular shifts:
```

660 1) Submit your final answer in this format:  
 661 Reasoning: <step-by-step reasoning>  
 662 Answer: <number of right circular shifts to transform the input into the output>  
 663  
 664 Else if you need to gather more information to find the number of right circular shifts:  
 665 2) Ask a query in this format:  
 666 Reasoning: <step-by-step reasoning>  
 667 Query: <index 0-{n-1}>  
 668 Value: <possible output value at the queried index>  
 669 Answer: <number of right circular shifts to transform the input into the output>  
 670  
 671 Use EXACTLY one of these two formats, and NOTHING ELSE (no markdown, no extra text).

## 672 12.2 Fermi

673 Both approaches have shared an initial prompt segment describing the task

674 You are a methodical math problem solver. Your goal is to solve the word problem below.  
 675 ↪ You have the problem text, a history of questions you previously asked and the  
 676 ↪ answers you received, and a list of variables that can be related to the problem.  
 677  
 678 **\*\*Problem Text\*\***  
 679 Question: How many individual life forms are there on Earth?  
 680  
 681 **\*\*Variable list\*\***  
 682 Variables:{"F1": "number of timezones where it is meal time at any given point.", "F2": "  
 683 ↪ total number of animals living in inhabitable areas.", "F3": "average number of  
 684 ↪ people living in a single timezone."}  
 685  
 686 **\*\*History of Questions Asked and Answers Received:**  
 687 --- Interaction 1 ---  
 688 Your Question: What is the total animal population (F13) on Earth?  
 689 Answer Received: The total animal population on Earth is 2e+19.  
 690 --- End of History ---\*\*

691 Following the problem context and past history, the agent needs to decide which variable to query  
 692 based on CoT and MES Algorithm

693 For CoT:

694 **\*\*Your Task:\*\***  
 695 1. Analyze the problem text and the interaction history carefully.  
 696 2. Reason step-by-step towards solving the problem. Show your thinking process clearly.  
 697 3. Use the information from the problem text and any answers received from previous  
 698 ↪ questions.  
 699 4. Determine if you can calculate the final answer based on your current understanding.  
 700  
 701 **\*\*Choose EXACTLY ONE of the following actions:\*\***  
 702  
 703 A) If you can calculate the final answer:  
 704 Provide your step-by-step Chain-of-Thought reasoning, explaining how you derived the  
 705 ↪ answer using the problem text and previous answers received.  
 706 **\*\*Output Format (Required):\*\***  
 707 Reasoning: <Your step-by-step derivation>  
 708 FinalAnswer: <The final numerical answer>  
 709  
 710 B) If you are stuck or need more specific information to proceed:  
 711 Identify the *\*single, specific question\** you need to ask to get the necessary  
 712 ↪ clarification or missing piece of information required for your reasoning.  
 713 Provide reasoning explaining *\*why\** you need to ask this question based on your current  
 714 ↪ reasoning progress.  
 715 **\*\*Output Format (Required):\*\***  
 716 Reasoning: <Explain why this question is necessary for your next step>  
 717 QueryOracle: <Your single, specific question>  
 718  
 719 **\*\*Important:\*\***  
 720 - Perform calculations step-by-step in your reasoning.  
 721 - Ask focused, specific questions based on gaps in your understanding or missing  
 722 ↪ information identified in your reasoning.  
 723 - Do not ask for information you can deduce yourself from the text or history.  
 724 - Use only one output format (Reasoning + FinalAnswer OR Reasoning + QueryOracle).  
 725 - for large number that needs scientific notation, use scientific notation in format Xe+Y  
 726 ↪ , DO NOT DO 10^{-10}  
 727  
 728 e.g. \{(8.26 \times 10^{6})\} meters should be written as 8.26e+6  
 729  
 730 Begin your Chain-of-Thought reasoning and determine your next action.

731 For MES, We have different prompt for two steps: Query step and rollout step:

## 732 Query step:

```
733 **Instructions:**
734 1. Review the problem context (text, variables) and interaction history.
735 2. think of what is the formula involving the variables to compute the end result
736 3. Identify variables whose values are *missing* but are *required* to solve the problem
737   ↪ based on the equations.
738 3. **Do NOT ask about any variable that has already been asked in the interaction history
739   ↪ - regardless of whether it was answered**
740 4. For each new question:
741   * Provide brief reasoning explaining *why* this specific value is needed for the
742     ↪ calculation based on the equations.
743   * Avoid asking for the same variable again.
744   * Make a reasonable estimation for the variable
745   * Ask concise, clear natural language questions (e.g., "What is the value of N?", "
746     ↪ How many customers were there?").
747 5. Output *only* the JSON object and conform exactly to the provided schema.
748
749 **IMPORTANT: Your response MUST be valid parseable JSON that exactly matches the schema
750   ↪ below. DO NOT include any explanations or text outside the JSON structure.**
751 - The response must start with {{ and end with }}
752 **Required JSON Schema:**
753 {schema_json}
754 Generate the JSON object:
```

## 755 Rollout step:

```
756 **Your Task & Decision Logic:**
757
758 1. **Analyze:** Review all provided information (Problem, History, Hypothetical Answer)
759 2. **Decide:**
760   * **Scenario 1: Numeric Answer Possible**
761     * If you *can* calculate a final numeric answer:
762       * Perform the calculation step-by-step in your reasoning.
763       * Your primary output is the **numeric value**.
764   * **Scenario 2: Numeric Answer NOT Possible**
765     * If you *cannot* calculate a final numeric answer (because some values are
766       ↪ still unknown even with the hypothetical answer):
767       * Explain *why* a numeric calculation isn't possible in your reasoning.
768       * Derive the most simplified **symbolic expression** for the answer based
769       ↪ on knowns and unknowns.
770       * Your primary output is the **symbolic expression**.
771
772 **Output Format (Based on Decision):**
773
774 * **'simulated_reasoning' (Required for BOTH scenarios):**
775   * Provide clear, step-by-step reasoning.
776   * If Scenario 1: Show the derivation and calculation leading to the numeric answer.
777   * If Scenario 2: Explain the derivation of the symbolic expression and state why a
778     ↪ numeric value cannot be determined.
779
780 * **'symbolic_expression' (Required for BOTH scenarios):**
781   * Provide the final symbolic formula representing the answer.
782   * Must include only variables and Python/SymPy-style operators ('+', '-', '*', '/',
783     ↪ '**').
784   * No assignment ('='), units, or descriptive text. Just the pure expression.
785   * *Even if you provide a 'numeric_value', this field should contain the underlying
786     ↪ symbolic formula used.*
787   * Correct Examples: 'N * A', '(X + Y) / Z', '1800 - 900 - C - 2*R'
788   * Wrong Examples: 'Cost = X + Y', 'X = 5 + N', 'Final = 40'
789
790 * **'numeric_value' (Conditional):**
791   * **Scenario 1:** Provide the **single calculated number** (int or float).
792     * Correct Examples: '40', '12.5', '-100', '40.0'
793     * **CRITICALLY WRONG Examples (NEVER USE THESE):** '"40"', '7 + 3*7 + 12"', '7
794       ↪ + 3*7 + 12', '[40]', '10 * 4'
795     * for large number that needs scientific notation, use scientific notation in
796       ↪ format Xe+Y , DO NOT DO 10^{-10}
797       e.g. \{(8.26 \times 10^{\{6\}}\}) meters should be written as 8.26e+6
798     * **Scenario 2:** Set this field to **'null'**.
799
800
801 **IMPORTANT: Your response MUST be valid parseable JSON that exactly matches the schema
802   ↪ below. DO NOT include any explanations or text outside the JSON structure.**
803 The response must start with {{ and end with }}
804
805 **Required JSON Schema:**
806 {schema_json}
807 Generate the JSON object:
```

808 **12.3 GSM8K**

809 **Stop condition:**

```
810 You are a methodical math problem solver using Chain-of-Thought reasoning. Your goal is
811 ↪ to solve the word problem below. You have the problem text and a history of
812 ↪ questions you previously asked and the answers you received. You have access to
813 ↪ the original problem, all variable definitions, all equations, and the history of
814 ↪ your interactions with an Oracle (questions you asked and answers you received).
815
816 **Problem Context:**
817 {problem_context}
818
819 **Full Interaction History (including latest Oracle answer):**
820 {history_string}
821
822 **Your Task:**
823 Based on ALL the information currently available (problem text, variable definitions,
824 ↪ equations, and the entire interaction history), critically assess if you have
825 ↪ sufficient information to calculate a definitive, single numerical answer to the
826 ↪ problem.
827
828 1. **Reasoning:** Explain your thought process.
829 * Consider the goal of the problem.
830 * Review the necessary equations.
831 * Check if all variables in those equations either have known values (from context
832 ↪ or Oracle answers) or can be derived from other knowns.
833 2. **Decision ('can_solve_now'):**
834 * Set to 'true' if you are confident you can now proceed to calculate the final
835 ↪ numerical answer without needing further Oracle input.
836 * Set to 'false' if you believe crucial numerical information is still missing, or
837 ↪ if there's critical ambiguity in definitions/equations that prevents a final
838 ↪ numeric solution.
839
840 Output your decision strictly as a JSON object conforming to the schema below.
```

841 **Query Step:**

```
842 # ROLE: Strategic Question Generator for Math Problems
843 You are analyzing a math word problem to identify the *most useful questions* an agent
844 ↪ could ask an Oracle to obtain **missing numerical values** required to solve the
845 ↪ problem based on the provided equations and variables.
846 **Problem Context:**
847 {problem_context}
848
849 **Interaction History (Up to Query):**
850 {history_string}
851
852 **Instructions:**
853 1. Review the problem context (text, variables, equations) and interaction history.
854 2. Identify variables whose values are *missing* but are *required* to solve the problem
855 ↪ based on the equations.
856 3. **Do NOT ask about any variable that has already been asked in the interaction history
857 ↪ - regardless of whether it was answered - with only one exception:**
858 - You may ask again about a previously asked variable **only if you are certain that
859 ↪ obtaining this one value will allow you to complete the problem.** This
860 ↪ exception should be used sparingly and only when absolutely justified.
861 4. For each new question:
862 * Provide brief reasoning explaining **why** this specific value is needed for the
863 ↪ calculation based on the equations.
864 * Use 'value' as the query type unless clarification of a symbol is essential (in
865 ↪ that case, use 'definition' or 'equation').
866 * Avoid asking for the same variable again.
867 * Ask concise, clear natural language questions (e.g., "What is the value of N?", "
868 ↪ How many customers were there?").
869 5. Output **only** the JSON object and conform exactly to the provided schema.
```

870 **Symbolic equation generation step:**

```
871 # ROLE: Symbolic Problem Solver (Initial Expression)
872
873 You are tasked with deriving an initial symbolic expression to solve a math word problem.
874 You have been given a problem context, the actual interaction history with an Oracle so
875 ↪ far,
876 and a hypothetical first query to the Oracle along with its hypothetical answer for THIS
877 ↪ simulation.
878
879 **Problem Context:**
880 {problem_context}
881
```

```

882 **Actual Interaction History (Overall Agent-Oracle):**
883 {base_interaction_history}
884
885 Question Asked: {initial_query_text}
886
887 **Instructions:**
888 1. Incorporate the hypothetical Oracle answer into your understanding, along with the
889     ↪ actual history.
890 2. Perform step-by-step reasoning to derive a **single, final symbolic mathematical
891     ↪ expression** that represents the solution to the problem.
892 3. This expression may contain known numbers and variable names from the problem context
893     ↪ or history.
894 4. Use standard Python/SymPy math notation ('+', '-', '*', '/', '**').
895 5. Use the variable in the initial_query_text to create a symbolic expression.
896 6. **DO NOT evaluate to a number yet, even if you think you can.** The goal is the
897     ↪ symbolic form.
898 7. Output ONLY the JSON object strictly following the schema.
899
900 Example Output: '(A * B) + C - D/2'
901
902 Generate the Symbolic Expression:

```

903 Rollout step:

```

904 # ROLE: Math Problem Solver (Internal Hallucination Step)
905
906 You are in a multi-step simulation to solve a math problem.
907 This simulation started after a hypothetical query ('{initial_query_text}') was
908     ↪ hypothetically answered ('{initial_hypothetical_answer}').
909 You are now at an internal step of this simulation.
910
911 **Original Problem Context:**
912 {problem_context}
913
914 **Actual Interaction History (before this simulation started):**
915 {base_interaction_history}
916
917 **Simulated History (This Rollout: Initial Query + Internal Hallucinations so far):**
918 {internal_hallucination_history_string}
919
920 **Current State of Symbolic Expression (if any):** {current_symbolic_expression}
921
922 **Instructions (Current Internal Step: {current_internal_step}/{max_internal_steps}):**
923 1. Review ALL information. Perform step-by-step reasoning based on the current state.
924 2. Update 'current_symbolic_expression_after_step' with the symbolic expression after
925     ↪ your reasoning for this step.
926 3. **Decision Point:**
927     * **IF you can calculate a final, single numeric answer NOW:**
928       * Set 'is_final_numeric_answer_reached' to 'true'.
929       * Provide the number in 'final_numeric_answer'.
930       * Set 'next_variable_to_hallucinate' and 'hallucinated_value_for_next_variable'
931         ↪ to 'null'.
932     * **ELSE IF this IS THE FINAL internal step ({current_internal_step} == {
933         ↪ max_internal_steps}):**
934       * **YOU MUST ATTEMPT TO EVALUATE the 'current_symbolic_expression_after_step'
935         ↪ to a number.**
936       * If evaluation is successful, set 'is_final_numeric_answer_reached' to 'true'
937         ↪ and provide it in 'final_numeric_answer'.
938       * If evaluation fails (still symbolic or error), set '
939         ↪ is_final_numeric_answer_reached' to 'false' and 'final_numeric_answer' to
940         ↪ 'null'.
941       * Set 'next_variable_to_hallucinate' and 'hallucinated_value_for_next_variable'
942         ↪ to 'null'.
943     * **ELSE IF this is NOT the final internal step AND you are still blocked by a
944         ↪ missing numerical value:**
945       * Identify the 'next_variable_to_hallucinate' (the symbol of the variable).
946       * Propose a PLAUSIBLE, SIMPLE 'hallucinated_value_for_next_variable' (e.g.,
947         ↪ small integer).
948       * Set 'is_final_numeric_answer_reached' to 'false' and 'final_numeric_answer'
949         ↪ to 'null'.
950     * **OTHERWISE (e.g., cannot identify a clear next variable to hallucinate before
951         ↪ max steps):**
952       * Set 'is_final_numeric_answer_reached' to 'false', 'final_numeric_answer' to '
953         ↪ null'.
954       * Set 'next_variable_to_hallucinate' and 'hallucinated_value_for_next_variable'
955         ↪ to 'null'.

```

956 Final answer step:

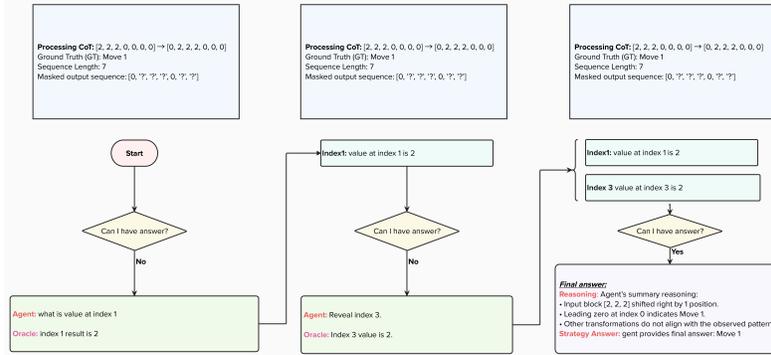


Figure 6: arc-1d CoT

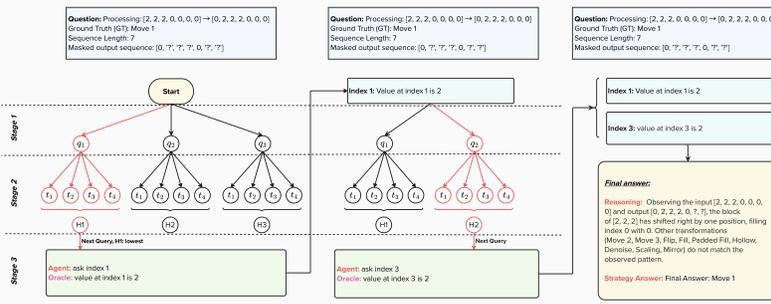


Figure 7: arc-1d MES

```

957 You are a methodical math problem solver using Chain-of-Thought reasoning. Your goal is
958     ↪ to solve the word problem below. You have the problem text and a history of
959     ↪ questions you previously asked and the answers you received.
960 **This is your final attempt. You MUST provide a final answer.**
961
962 You are in the **final internal step** of a multi-step to solve a math problem.
963 You have derived a symbolic expression and made several internal value.
964
965 **Original Problem Context:**
966 {problem_context}
967
968 **Interaction History (Overall Agent-Oracle):**
969 {history_string}
970
971 **Instructions:**
972 1. Review ALL known values from the original context, actual history, and ALL values
973     ↪ from this.
974 2. **YOU MUST ATTEMPT TO EVALUATE the 'current_symbolic_expression' to a single numeric
975     ↪ value.**
976 3. Provide step-by-step 'final_reasoning' showing this evaluation.
977 4. You MUST provide a final answer now. Make your best calculation based on the
978     ↪ information you have, even if some parts were uncertain or missing.
979 5. Output ONLY the JSON object strictly following the schema.
980
981 **Output Format (Required):**
982 Reasoning: <Your final step-by-step derivation, justifying your answer based on all
983     ↪ information>
984 FinalAnswer: <Your best calculated numerical answer>

```

## 985 13 Benchmark Examples

### 986 13.1 ARC-1D:

987 Example: Reasoning output after it has queried index 0.

**Agent loaded problem:**

Processing CoT: [2, 2, 2, 0, 0, 0, 0] → [0, 2, 2, 2, 0, 0, 0]

Ground Truth (GT): Move 1

Sequence Length: 7

Masked output sequence: [0, '?', '?', '?', 0, '?', '?']

— Turn 1 —

**Agent's thought:** Let's analyze the input [2, 2, 2, 0, 0, 0, 0] and output [0, '?', '?', '?', 0, '?', '?']. The first element is 2, but the output at index 0 is 0, indicating the beginning of the sequence has changed. The last element is 0, and the output at index 6 is '?', giving limited clues.

Possible transformations considered:

- **Flip:** Unlikely (first element would stay 2).
- **Move 1, 2, 3:** Shifts could explain index 0 changing to 0.
- **Fill / Padded Fill:** No repeated or gapped markers, unlikely.
- **Hollow / Denoise / Scaling / Mirror:** No match to observed pattern.

**Agent query:** Reveal index 1 to check if shift occurred.

**Oracle's answer:** Index 1 value is 2.

— Turn 2 —

**Agent's thought:** Now the output is [0, 2, '?', '?', 0, '?', '?']. Considering the non-zero block, only a shift explains index 0 becoming 0. Reveal index 3 to confirm the shift.

**Agent query:** Reveal index 3.

**Oracle's answer:** Index 3 value is 2.

— Turn 3 —

**Agent's thought:** Output is now [0, 2, 2, 2, 0, '?', '?']. This suggests a shift of the [2, 2, 2] block. To narrow down the move type, reveal index 2.

**Agent query:** Reveal index 2.

**Oracle's answer:** Index 2 value is 2.

— Turn 4 —

**Agent's thought:** With [0, 2, 2, 2, 0, '?', '?'], to distinguish between Move 1/2/3, reveal index 5.

**Agent query:** Reveal index 5.

**Oracle's answer:** Index 5 value is 0.

**Final Decision**

**Agent's summary reasoning:**

- Input block [2, 2, 2] shifted right by 1 position.
- Leading zero at index 0 indicates Move 1.
- Other transformations do not align with the observed pattern.

Predicted final output: [0, 2, 2, 2, 0, 0, 0].

**Agent Final Answer: Move 1**

**Correct Answer?: True (Expected: Move 1)**

**Agent loaded problem:**

Processing: [2, 2, 2, 0, 0, 0, 0] → [0, 2, 2, 2, 0, 0, 0]

Ground Truth (GT): Move 1

Sequence Length: 7

Masked output sequence: [0, '?', '?', '?', 0, '?', '?']

— **Turn 1** —

— **Entropy Strategy: Evaluating Queries** —

Query candidates: [1]

Evaluating query 1:

Predicted value probabilities: {0: 0.375, 2: 0.625}

- If value 0 (prob 0.38): Trans\_Probs: {'Move 1': 0.125, 'Move 3': 0.375, 'Move 2': 0.125, 'Fill': 0.125, 'Mirror': 0.125, 'Flip': 0.125}, Entropy: 2.41

- If value 2 (prob 0.62): Trans\_Probs: {'Move 1': 0.86, 'Move 3': 0.14}, Entropy: 0.59

Expected entropy: 1.27

Best query chosen: 1 (Entropy: 1.27)

Queried index 1, got value 2

Updated masked sequence: [0, 2, '?', '?', 0, '?', '?']

Votes: [False, False, False, False, True]

— **Turn 2** —

Query candidates: [2, 3]

Evaluating query 2:

Predicted value probabilities: {2: 0.625, 0: 0.375}

- Value 2 (prob 0.62): Trans\_Probs: {'Move 1': 1.0}, Entropy: 0.00

- Value 0 (prob 0.38): Trans\_Probs: {'Move 1': 0.86, 'Hollow': 0.14}, Entropy: 0.59

Expected entropy: 0.22

Evaluating query 3:

Predicted value probabilities: {0: 0.625, 2: 0.375}

- Value 0 (prob 0.62): Trans\_Probs: {'Move 1': 1.0}, Entropy: 0.00

- Value 2 (prob 0.38): Trans\_Probs: {'Move 1': 1.0}, Entropy: 0.00

Expected entropy: 0.00

Best query chosen: 3 (Entropy: 0.00)

Queried index 3, got value 2

Updated masked sequence: [0, 2, '?', 2, 0, '?', '?']

Votes: [False, True, False, False, False]

— **Turn 3** —

Query candidates: [2, 5]

Evaluating query 2:

Predicted value probabilities: {2: 0.875, 0: 0.125}

- Value 2 (prob 0.88): Trans\_Probs: {'Move 1': 1.0}, Entropy: 0.00

- Value 0 (prob 0.12): Trans\_Probs: {'Move 1': 1.0}, Entropy: 0.00

Expected entropy: 0.00

Evaluating query 5:

Predicted value probabilities: {0: 0.71, 2: 0.29}

- Value 0 (prob 0.71): Trans\_Probs: {'Move 1': 0.88, 'Move 2': 0.12}, Entropy: 0.54

- Value 2 (prob 0.29): Trans\_Probs: {'Move 1': 0.88, 'Denoise': 0.12}, Entropy: 0.54

Expected entropy: 0.54

Best query chosen: 2 (Entropy: 0.00)

Queried index 2, got value 2

Updated masked sequence: [0, 2, 2, 2, 0, '?', '?']

Votes: [False, True, False, True, True]

— **Final Reasoning** —

**Agent's thought:** Observing the input [2, 2, 2, 0, 0, 0, 0] and output [0, 2, 2, 2, 0, ?, ?], the block of [2, 2, 2] has shifted right by one position, filling index 0 with 0. Other transformations (Move 2, Move 3, Flip, Fill, Padded Fill, Hollow, Denoise, Scaling, Mirror) do not match the observed pattern.

**Agent Decision:** Final transformation selected: **Move 1**

**Agent Log:** Final Answer: Move 1

Votes: ['Move 1', 'Move 1', 'Move 1', 'Move 1', 'Move 1']

**Correct Answer?:** True

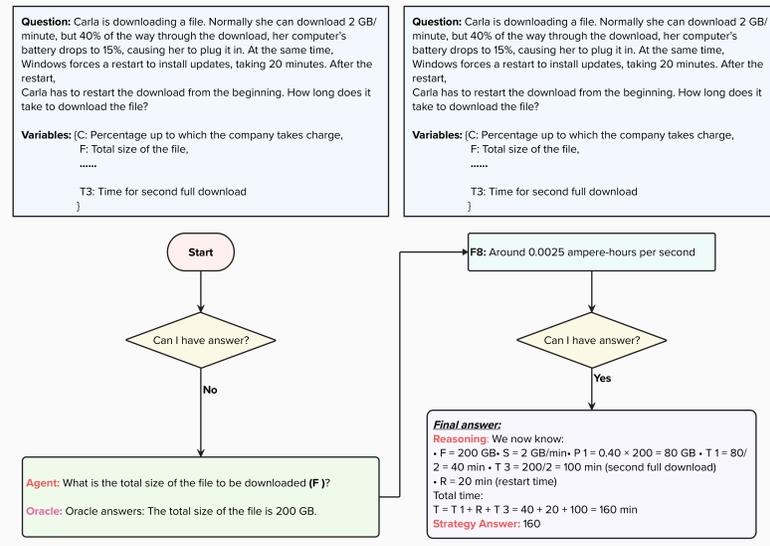


Figure 8: GSM CoT

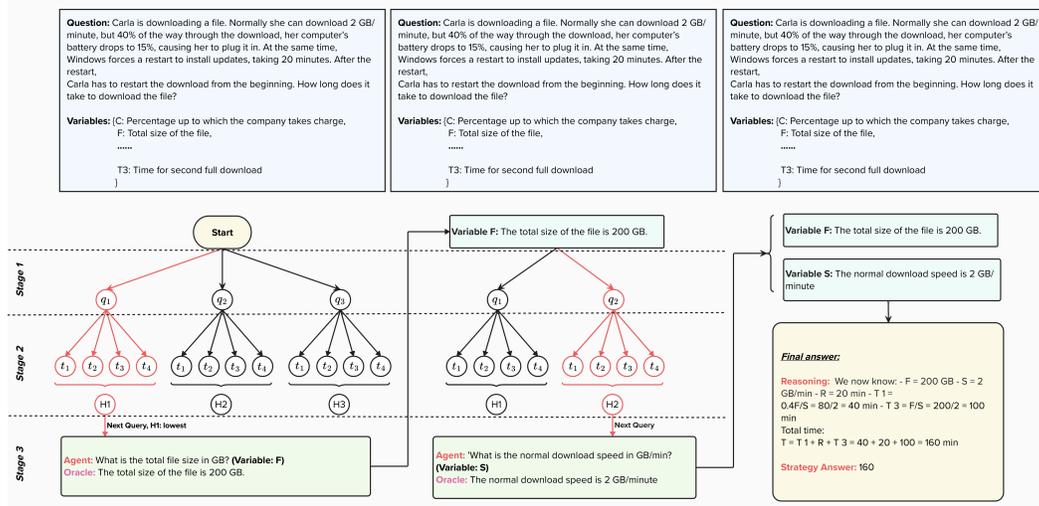


Figure 9: GSM MES

## Agent Problem Log - GSM - CoT

**Agent loaded problem:**

Carla is downloading a file. Normally she can download 2 GB/minute, but 40% of the way through the download, her computer's battery level drops to 15%, causing her to plug it into a power source. At the same time, Windows forces a restart to install updates, which takes 20 minutes. During this time, her computer charges but this doesn't affect the download speed. After the restart, Carla has to restart the download from the beginning. How long does it take to download the file?

**— Problem ID 0 | Turn 1/2 —**

**Agent's thought:** The problem states Carla downloads at 2 GB/minute. She gets 40% of the way through, then has a 20-minute interruption, and then has to restart the download from the beginning. We need to find the total time to download the file. Let's define the variables:

- $S = 2$  GB/min (Normal download speed)
- $P1 = 40\%$  (Percentage downloaded before the issue)
- $R = 20$  min (Time taken for restart)

We are missing the total file size  $F$ . Without  $F$ , we cannot calculate the time for the first 40% or the total time.

**Agent's query:** What is the total size of the file to be downloaded ( $F$ )?

**Oracle's answer:** The total size of the file is 200 GB.

**— Problem ID 0 | Turn 2/2 —**

**Agent's thought:** We now know:

- $F = 200$  GB
- $S = 2$  GB/min
- $P1 = 0.40 \times 200 = 80$  GB
- $T1 = 80/2 = 40$  min
- $T3 = 200/2 = 100$  min (second full download)
- $R = 20$  min (restart time)

Total time:

$$T = T1 + R + T3 = 40 + 20 + 100 = 160 \text{ min}$$

**Agent Final Answer: 160.0**

**Correct Answer? :** Correct Answer? True (Expected: 160.0)

**Agent loaded problem:**

Carla is downloading a file. Normally she can download 2 GB/minute, but 40% of the way through the download, her computer's battery drops to 15%, causing her to plug it in. At the same time, Windows forces a restart to install updates, taking 20 minutes. After the restart, Carla has to restart the download from the beginning. How long does it take to download the file?

**Agent Log:**

=== Variables: {C: Percentage up to which the company takes charge, F: Total size of the file, P1: Percentage downloaded before issue, R: Time for restart and updates, S: Download speed, T: Total download time, T1: Time spent downloading first 40%, T3: Time for second full download} ===

**Agent Log:**

=== Missing Info: F = 200 GB, C = 0.35 ===

**Agent Log:**

=== Solution: 160.0 minutes ===

**Agent Log:**

=== Starting Problem ID: 0 ===

— Turn 1/5 —

— Entropy Strategy: Stage 1: Generating 32 Queries —

— Strategy: Successfully generated/parsed 32/32 queries —

— Strategy: Clustered 32 queries into 1 group —

— Entropy Strategy: Stage 2: Simulating Rollouts —

Simulating 32 rollouts for cluster: ('F', 'value') (Rep Query: 'F')

Cluster ('F', 'value'): Outcomes={0.7\*F + 20: 30, 27.0: 2}, Entropy=0.3373

— Entropy Strategy: Stage 3: Selecting Action —

Min Entropy Cluster: ('F', 'value') (Entropy=0.3373)

Most Frequent Outcome: '0.7\*F + 20' (Count=30/32, Ratio=0.94)

**Agent's thought:** We know: - Download speed  $S = 2$  GB/min. - File size  $F$ . - First 40% downloaded:  $P1 = 0.4 \times F$ . - Time for first 40%:  $T1 = P1/S = 0.4F/2 = 0.2F$ . - Restart time  $R = 20$  min. - Full second download:  $T3 = F/S = F/2$ .  
Total time  $T = T1 + R + T3 = 0.2F + 20 + 0.5F = 0.7F + 20$ .

**Agent Decision:** Query (Most frequent outcome is symbolic: '0.7\*F + 20')

**Agent Log:**

Strategy selected Query: 'What is the total file size in GB? (Variable: F)'

**Agent Log:**

Strategy intermediate answer state: '0.7\*F + 20'

**Oracle's answer:** 'The total size of the file is 200 GB.'

— Turn 2/5 —

— Entropy Strategy: Stage 1: Generating 32 Queries —

— Strategy: Successfully generated/parsed 32/32 queries —

— Strategy: Clustered 32 queries into 3 groups —

— Entropy Strategy: Stage 2: Simulating Rollouts —

Simulating 32 rollouts for cluster: ('C', 'value') (Rep Query: 'C')

Cluster ('C', 'value'): Outcomes={160.0: 15, 100.0: 13, 120.0: 4}, Entropy=1.4153

Simulating 32 rollouts for cluster: ('S', 'value') (Rep Query: 'S')

Cluster ('S', 'value'): Outcomes={20 + 280/S: 29, 160.0: 1, 20 + 200/S: 1}, Entropy=0.4096

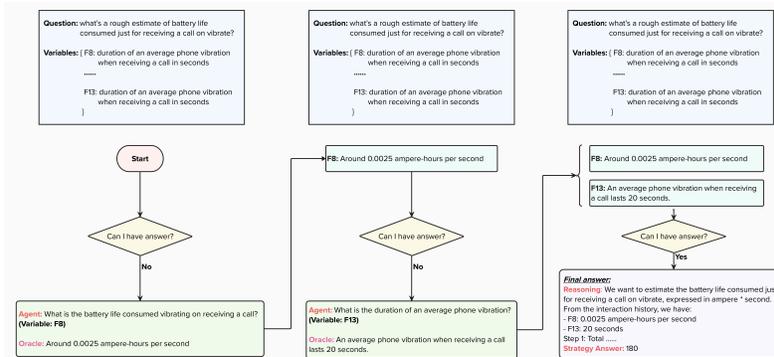


Figure 10: Fermi CoT

### Agent Execution Log (Continued)

#### — Entropy Strategy: Stage 3: Selecting Action —

Min Entropy Cluster: ('S', 'value') (Entropy=0.4096)

Most Frequent Outcome: '20 + 280/S' (Count=29/31, Ratio=0.94)

Agent's thought: With  $F = 200$  GB, total time:

$$T = \frac{1.4F}{S} + R = \frac{280}{S} + 20$$

We need  $S$  to compute numerically.

Agent Decision: Query (Most frequent outcome is symbolic: '20 + 280/S')

Agent Log: Strategy selected Query: 'What is the normal download speed in GB/min? (Variable: S)'

Agent Log: Strategy intermediate answer state: '20 + 280/S'

Oracle's answer: 'The normal download speed is 2 GB/minute.'

#### — Turn 3/5 —

#### — Entropy Strategy: Stage 1: Generating 32 Queries —

— Strategy: Successfully generated/parsed 32/32 queries —

— Strategy: Clustered 32 queries into 2 groups —

#### — Entropy Strategy: Stage 2: Simulating Rollouts —

Simulating 32 rollouts for cluster: ('P1', 'value') (Rep Query: 'P1')

Cluster ('P1', 'value'): Various symbolic outcomes, Entropy=3.0097

Simulating 32 rollouts for cluster: ('C', 'value') (Rep Query: 'C')

Cluster ('C', 'value'): Various symbolic outcomes, Entropy=2.1437

#### — Entropy Strategy: Stage 3: Selecting Action —

Min Entropy Cluster: ('C', 'value') (Entropy=2.1437)

Most Frequent Outcome: '160.0' (Count=10/32, Ratio=0.31)

Agent's thought: We now know: -  $F = 200$  GB -  $S = 2$  GB/min -  $R = 20$  min -  $T1 = 0.4F/S = 80/2 = 40$  min -  $T3 = F/S = 200/2 = 100$  min

Total time:

$$T = T1 + R + T3 = 40 + 20 + 100 = 160 \text{ min}$$

Agent Log: Strategy decided to Answer Numerically.

Agent Log: Strategy Reasoning: We calculated  $T = 160$  min.

Agent Log: Strategy Answer: 160.0

Agent Log: Final Status: solved, Correct: 1

993

### 994 13.3 Fermi

995 Example Fermi Problem:

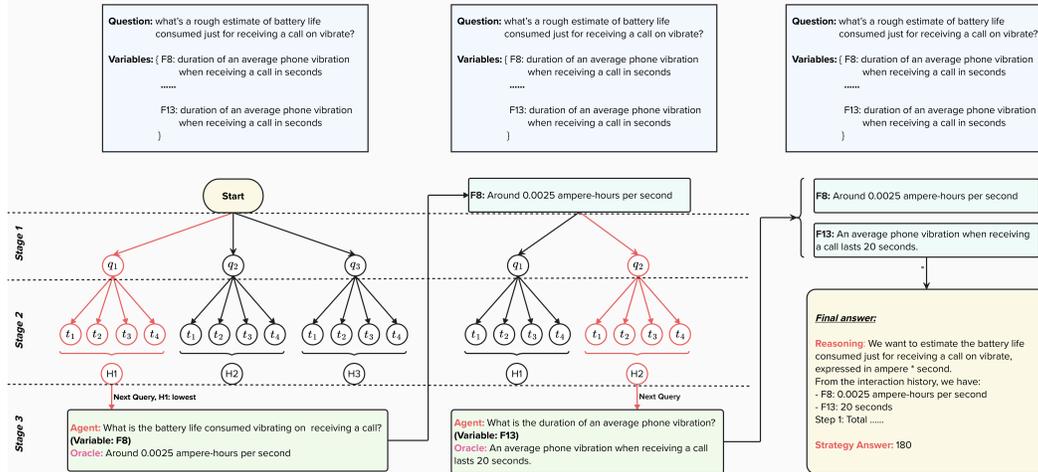


Figure 11: fermi MES

### Agent Problem Log - Fermi - CoT

#### Agent loaded problem:

On average, across all phones, what's a rough estimate of battery life consumed just for receiving a call on vibrate? Regardless of previous unit, **YOU MUST** answer in terms of Unit: ampere \* second

Below are Variables you should ask info within this list:

- F1: number of fermi problems solved by an average enthusiast each day
- F2: pressure equivalent of 1 atmosphere in pascals
- F3: speed of sound at sea level in kilometers per day
- F4: circumference of the earth in kilometers
- F5: number of days in a year
- F6: calories burned per hour when the phone heats up in the pocket
- F7: energy liberated when the Hiroshima bomb detonated in joules
- F8: battery life consumed per second while vibrating on receiving a call in ampere-hours
- F9: number of fermi problem enthusiasts in the world
- F10: number of hours in a day
- F11: energy liberated when a pop-it is burst in joules
- F12: average American lifespan in years
- F13: duration of an average phone vibration when receiving a call in seconds
- F14: force contained in a fart in newtons
- F15: average duration a phone stays in a pocket during the day in hours
- F16: normal heartbeats per minute
- F17: number of minutes in an hour
- F18: average life expectancy in years
- F19: impact area of a fart in square meters
- F20: reduction in average lifespan caused by cancer in years

#### — Problem ID 4 | Turn 1/5 —

**Agent Log:** Agent asks Oracle: What is the battery life consumed per second while vibrating on receiving a call in ampere-hours (F8)?

**Oracle's answer:** Around 0.0025 ampere-hours of battery life is consumed per second while vibrating on receiving a call.

#### — Problem ID 4 | Turn 2/5 —

**Agent Log:** Agent asks Oracle: What is the duration of an average phone vibration when receiving a call in seconds (F13)?

**Oracle's answer:** An average phone vibration while receiving a call lasts 20 seconds.

### Agent Execution Log (Continued)

#### — Problem ID 4 | Turn 3/5 —

**Agent's thought:** We want to estimate the battery life consumed just for receiving a call on vibrate, expressed in ampere \* second.

From the history, we have two key pieces of information:

- F8: 0.0025 ampere-hours per second
- F13: 20 seconds

Step 1: Calculate total battery consumption (in ampere-hours): Battery consumption = consumption rate  $\times$  duration =  $0.0025 \text{ Ah/s} \times 20 \text{ s} = 0.05 \text{ Ah}$

Step 2: Convert ampere-hours to ampere-seconds:  $1 \text{ Ah} = 3600 \text{ A} \cdot \text{s}$  Battery consumption =  $0.05 \text{ Ah} \times 3600 \text{ s/Ah} = 180 \text{ A} \cdot \text{s}$

**Agent Final Answer: 180.0**

**Correct Answer?: True (Expected: 180.0)**

997

**Agent loaded problem:** On average, across all phones, what's a rough estimate of battery life consumed just for receiving a call on vibrate? **Regardless of what unit previously asked, the answer MUST BE in terms of Unit: ampere \* second**

Below are Variables you should ask info within this list:

- F1: number of fermi problems solved by an average enthusiast each day
- F2: pressure equivalent of 1 atmosphere in pascals
- F3: speed of sound at sea level in kilometers per day
- F4: circumference of the earth in kilometers
- F5: number of days in a year
- F6: calories burned per hour when the phone heats up in the pocket
- F7: energy liberated when the Hiroshima bomb detonated in joules
- F8: battery life consumed per second while vibrating on receiving a call in ampere-hours
- F9: number of fermi problem enthusiasts in the world
- F10: number of hours in a day
- F11: energy liberated when a pop-it is burst in joules
- F12: average American lifespan in years
- F13: duration of an average phone vibration when receiving a call in seconds
- F14: force contained in a fart in newtons
- F15: average duration a phone stays in a pocket during the day in hours
- F16: normal heartbeats per minute
- F17: number of minutes in an hour
- F18: average life expectancy in years
- F19: impact area of a fart in square meters
- F20: reduction in average lifespan caused by cancer in years

Agent Log:

=== Variables: ... ===

Agent Log:

=== Solution: 180.0 ===

Agent Log:

=== Starting Problem ID: 4 ===

— Turn 1/5 —

— Entropy Strategy: Stage 1: Generating 32 Queries —

— Strategy: Successfully generated/parsed 32/32 queries —

— Strategy: Clustered 32 queries into 2 groups —

— Entropy Strategy: Stage 2: Simulating Rollouts in Parallel —

Simulating 32 rollouts for cluster: ('F8', 'value') (Rep Query: 'F8')

Simulating 32 rollouts for cluster: ('F13', 'value') (Rep Query: 'F13')

Cluster ('F8', 'value'): Outcomes=3600\*F13\*F8: 32, Entropy=0.0000

Cluster ('F13', 'value'): Outcomes=3600\*F13\*F8: 32, Entropy=0.0000

— Entropy Strategy: Stage 3: Selecting Action —

Min Entropy Cluster: ('F8', 'value') (Entropy=0.0000)

Most Frequent Outcome: '3600\*F13\*F8' (Count=32/32, Ratio=1.00)

Agent Decision: Query (Most frequent outcome is symbolic: '3600\*F13\*F8')

Agent Log:

Strategy selected Query: 'What is the battery life consumed per second while vibrating on receiving a call in ampere-hours? (Variable: F8)'

Agent Log:

Strategy intermediate answer state: '3600\*F13\*F8'

Oracle's answer: 'Around 0.0025 ampere-hours per second while vibrating on receiving a call.'

## Agent Execution Log (Continued)

— Turn 2/5 —

— Entropy Strategy: Stage 1: Generating 32 Queries —

— Strategy: Successfully generated/parsed 32/32 queries —

— Strategy: Clustered 32 queries into 1 group —

— Entropy Strategy: Stage 2: Simulating Rollouts in Parallel —

Simulating 32 rollouts for cluster: ('F13', 'value') (Rep Query: 'F13')

Cluster ('F13', 'value'): Outcomes=9.0: 3, 3600\*F13\*F8: 22, 9\*F13: 1, 9.0\*F13: 2, 27.0: 1, 18.0: 1, Entropy=1.4115

— Entropy Strategy: Stage 3: Selecting Action —

Min Entropy Cluster: ('F13', 'value') (Entropy=1.4115)

Most Frequent Outcome: '3600\*F13\*F8' (Count=22/30, Ratio=0.73)

**Agent's thought:** We want to estimate the battery life consumed (in ampere-seconds) just for receiving a call on vibrate.

**Agent Decision:** Query (Most frequent outcome is symbolic: '3600\*F13\*F8')

**Agent Log:**

Strategy selected Query: 'What is the duration of an average phone vibration when receiving a call in seconds? (Variable: F13)'

**Agent Log:**

Strategy intermediate answer state: '3600\*F13\*F8'

**Oracle's answer:** 'An average phone vibration when receiving a call lasts 20 seconds.'

— Turn 3/5 —

— Entropy Strategy: Stage 1: Generating 32 Queries —

— Strategy: Successfully generated/parsed 32/32 queries —

— Strategy: Clustered 32 queries into 2 groups —

— Entropy Strategy: Stage 2: Simulating Rollouts in Parallel —

Simulating 32 rollouts for cluster: ('F17', 'value') (Rep Query: 'F17')

Simulating 32 rollouts for cluster: ('F13', 'value') (Rep Query: 'F13')

Cluster ('F13', 'value'): Outcomes=180.0: 32, Entropy=0.0000

Cluster ('F17', 'value'): Outcomes=180.0: 32, Entropy=0.0000

— Entropy Strategy: Stage 3: Selecting Action —

Min Entropy Cluster: ('F13', 'value') (Entropy=0.0000)

Most Frequent Outcome: '180.0' (Count=32/32, Ratio=1.00)

**Agent's thought:** We want to estimate the battery life consumed just for receiving a call on vibrate, expressed in ampere \* second.

From the interaction history, we have: - F8: 0.0025 ampere-hours per second - F13: 20 seconds

Step 1: Total ampere-hours =  $0.0025 \times 20 = 0.05$  Ah

Step 2: Convert to ampere-seconds =  $0.05 \times 3600 = 180$  A·s

Therefore, the battery life consumed is 180 ampere-seconds.

Symbolic expression:  $F8 \times F13 \times 3600$

**Agent Log:**

Strategy decided to Answer Numerically.

**Agent Log:**

Strategy Reasoning: As above.

**Agent Log:**

Strategy Answer: 180.0

**Agent Log:**

Final Status: solved, Correct: 1

## 1000 14 More Experimental Details

### 1001 14.1 Multi-Model Evaluation

### 1002 14.2 More Turn Based Evaluation

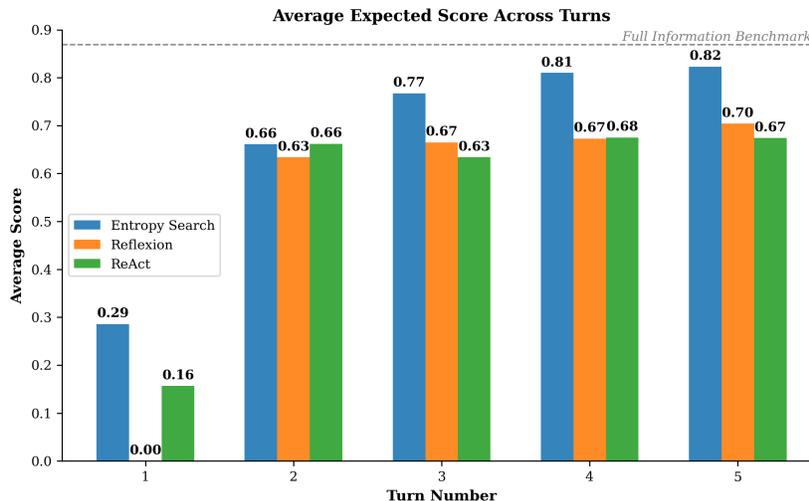


Figure 12: Turn-based evaluation on the **GSM8K** benchmark. We compare ENTROPY SEARCH (ES) with REFLEXION and REACT across five question asking turns. ES rapidly reduces the performance gap to the *Full-Information* upper bound, already matching the baselines after two questions and pulling decisively ahead thereafter. This illustrates the generality and efficiency of entropy-guided information acquisition.

1003 Following from the turned based evaluation for the Fermi dataset in Figure 4, we now present a  
1004 turn based analysis for the GSM8K dataset in Figure 4. On GSM8K, ENTROPY SEARCH takes  
1005 an early lead after the very first question (0.29 vs. 0.16 for REACT and 0.00 for REFLEXION) and  
1006 continues to improve with each turn. By the third question it attains an average score of 0.77—already  
1007 closing more than 80% of the gap to the full-information ceiling—while REACT and REFLEXION lag  
1008 behind at 0.63 and 0.66, respectively. After five questions ES reaches 0.82, whereas REACT plateaus  
1009 around 0.67 and REFLEXION tops out at 0.70. These results confirm that ES not only benefits from  
1010 additional turns but also outperforms both baselines at every stage, consistently acquiring the most  
1011 decision-critical context under uncertainty.

1012 Collectively, these findings underscore ENTROPY SEARCH’s effectiveness in strategically gathering  
1013 information on GSM8K, thereby enhancing decision quality in partially observable settings. Detailed  
1014 turn-based analyses for the 1D-ARC and Fermi benchmarks are provided in Appendix 14.2, further  
1015 corroborating these insights.

## 1016 15 Extended Results