
In-Context Learning for Pure Exploration

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this work, we study the active sequential hypothesis testing problem, also known
2 as *pure exploration*, where the goal is to actively control a data collection process
3 to efficiently identify the correct hypothesis underlying a decision problem. While
4 relevant across multiple domains, devising adaptive exploration strategies remains
5 challenging, particularly due to difficulties in encoding appropriate inductive biases.
6 To address these limitations, we introduce *In-Context Pure Exploration (ICPE)*, an
7 in-context learning approach that uses Transformers to learn exploration strategies
8 directly from experience. Numerical results across diverse benchmarks highlight
9 **ICPE**'s capability to achieve satisfactory performance in stochastic and structured
10 settings, demonstrating its ability to meta-learn exploration strategies.

11 1 Introduction

12 Modern artificial intelligence systems have achieved remarkable performance across specialized tasks
13 such as image classification [Krizhevsky et al. \[2012\]](#), Super-human board-game play [Silver et al.
14 \[2018\]](#), protein-structure prediction [Jumper et al. \[2021\]](#) and large-scale language modelling [Brown
15 et al. \[2020\]](#). Yet, there is still a lack in understanding how to autonomously discover meta-skills
16 fundamental for sequential decision making, such as active testing or active learning [Chernoff \[1992\]](#),
17 [Cohn et al. \[1996\]](#).

18 Consider an agent tasked with sequentially selecting samples to quickly improve its understanding
19 of an underlying phenomenon. When the decision maker can exert some control over the collected
20 samples' *information content*, this is a problem also known as the *active sequential hypothesis
21 testing problem* [Chernoff \[1992\]](#), [Ghosh \[1991\]](#), [Naghshvar and Javidi \[2013\]](#), [Naghshvar et al.
22 \[2012\]](#), [Mukherjee et al. \[2022\]](#) or *pure exploration problem* [Degenne and Koolen \[2019\]](#), [Degenne
23 et al. \[2019, 2020\]](#). Active hypothesis testing has become increasingly important nowadays, with
24 applications ranging from medical diagnostics [Berry et al. \[2010\]](#), image identification [Vaidhiyan et al.
25 \[2012\]](#), recommender systems [Resnick and Varian \[1997\]](#), etc. Nonetheless, devising an adaptive
26 data collection strategy is notoriously difficult and highly problem-specific.

27 In this paper, we address the question: how can sequential decision-making agents autonomously
28 discover and leverage hidden structure to enhance active exploration for hypothesis testing? We
29 introduce *In-Context Pure Explorer (ICPE)*, a novel method combining Supervised Learning and
30 Deep RL [Goodfellow et al. \[2016\]](#), [Murphy \[2023\]](#), which builds on the in-context learning and
31 sequence modeling capabilities of Transformers [Lee et al. \[2023\]](#)—a meta-learning approach that
32 uncovers underlying shared structure across a class of problems \mathcal{M} [Schaul and Schmidhuber \[2010\]](#),
33 [Bengio et al. \[1990\]](#).

34 **ICPE** operates by integrating two complementary neural networks: an inference (I) network, trained
35 via supervised learning to infer the true hypothesis given current data, and an exploration (π) network,
36 trained through reinforcement learning to select actions optimizing the inference accuracy of the I
37 network.

38 We validate **ICPE** through different benchmarks, demonstrating its ability to efficiently explore
39 in stochastic and structured environments. In particular, these results show that **ICPE** achieves
40 performance comparable to optimal instance-dependent Best Arm Identification (BAI) algorithms
41 [Garivier and Kaufmann \[2016\]](#), [Audibert and Bubeck \[2010\]](#), without requiring explicit problem-
42 specific exploration strategies that often involve solving complex optimization problems. Thanks to
43 the in-context capability of **ICPE**, it is effectively discovering active sampling techniques that at test
44 time do not need much more computation than a forward pass. Consequently, **ICPE** emerges as a
45 practical applicable method for data-efficient exploration.

46 1.1 Related Work

47 The problem of active sequential hypothesis testing [Chernoff \[1992\]](#), [Ghosh \[1991\]](#), [Lindley \[1956\]](#),
48 [Naghshvar and Javidi \[2013\]](#), [Naghshvar et al. \[2012\]](#), [Mukherjee et al. \[2022\]](#), [Gan et al. \[2021\]](#), in
49 which a learner is tasked with adaptively performing a sequence of actions to identify an unknown
50 property of the environment, is closely related to the exploration problem in Reinforcement Learning
51 (RL) [Sutton and Barto \[2018\]](#), where an agent needs to identify the optimal policy. This exploration
52 problem has long centred on regret minimisation [Sutton and Barto \[2018\]](#), with techniques based on
53 Upper-Confidence Bounds [Auer et al. \[2002, 2008\]](#), [Cappé et al. \[2013\]](#), [Lattimore and Hutter \[2012\]](#),
54 [Auer \[2002\]](#), posterior-sampling [Kaufmann et al. \[2012\]](#), [Osband et al. \[2013\]](#), [Russo and Van Roy \[2014\]](#),
55 [Gopalan et al. \[2014\]](#) and Information-Directed Sampling (IDS) [Russo et al. \[2018\]](#); yet these
56 schemes assume that minimizing regret is the sole objective and falter in identification problems.

57 A more closely related setting is that of pure exploration in bandits and Markov Decision Processes
58 (MDPs), settings known as Best Arm/Policy Identification (BAI/BPI) [Audibert and Bubeck \[2010\]](#),
59 [Garivier and Kaufmann \[2016\]](#), [Degenne and Koolen \[2019\]](#), [Al Marjani et al. \[2021\]](#), [Russo and](#)
60 [Proutiere \[2023a\]](#), [Russo et al. \[2025\]](#). In these problems the samples collected by the agent are
61 no longer perceived as rewards, and the agent must actively optimize its exploration strategy to
62 identify the optimal policy. BAI/BPI reframe the task as sequential hypothesis testing, yielding
63 instance-adaptive algorithms in fixed-confidence settings such as Track-and-Stop (TaS) [Garivier and](#)
64 [Kaufmann \[2016\]](#). However, while BAI strategy are powerful, they may be suboptimal when the
65 underlying information structure is not adequately captured within the hypothesis testing framework.
66 Although IDS and BAI offer frameworks to account for such structure, extending these approaches to
67 Deep Learning is difficult, particularly when the information structure is unknown.

68 Recently Transformers [Vaswani et al. \[2017\]](#), [Chen et al. \[2021\]](#) have demonstrated remarkable in-
69 context learning capabilities [Brown et al. \[2020\]](#), [Garg et al. \[2022\]](#). In-context learning [Moeini et al.](#)
70 [\[2025\]](#) is a form of meta-RL [Beck et al. \[2023\]](#), where agents can solve new tasks without updating any
71 parameters by simply conditioning on additional context, such as their action-observation histories.
72 Building on this ability, [Lee et al. \[2023\]](#) recently showed that Transformers can be trained in a
73 supervised manner using offline data to mimic posterior sampling in reinforcement learning. In
74 [Dai et al. \[2024\]](#) the authors present ICEE (In-Context Exploration Exploitation). ICEE uses
75 Transformer architectures to perform in-context exploration-exploration for RL. ICEE tackles this
76 challenge by expanding the framework of return conditioned RL with in-context learning [Chen et al.](#)
77 [\[2021\]](#), [Emmons et al. \[2021\]](#). Return conditioned learning is a type of technique where the agent
78 learns the return-conditional distribution of actions in each state. Actions are then sampled from the
79 distribution of actions that receive high return [Srivastava et al. \[2019\]](#), [Kumar et al. \[2019\]](#). Lastly, we
80 note the important contribution of RL² [Duan et al. \[2016\]](#), which proposes to represent an RL policy
81 as the hidden state of an RNN, whose weights are learned via RL. **ICPE** employs a similar idea, but
82 focuses on a different objective (identification), and splits the process into a supervised inference
83 network that provides rewards to an RL-trained transformer network that selects actions to maximize
84 information gain.

85 2 Learning to Explore: In-Context Pure Exploration

86 We introduce **ICPE** (In-Context Pure Exploration), a deep-learning framework that combines se-
87 quential architecture with supervised and reinforcement learning to automatically discover efficient
88 exploration policies for active sequential hypothesis testing. Instead of explicitly encoding induc-
89 tive biases, we use transformers to let the agent autonomously infer the problem structures from
90 experiences.

91 **Environment and Interaction Model.** We consider a model class of environments \mathcal{M} and a
92 distribution $\mathcal{P}(\mathcal{M}) \in \Delta(\mathcal{M})$ from which the true environment M is sampled from. We model an

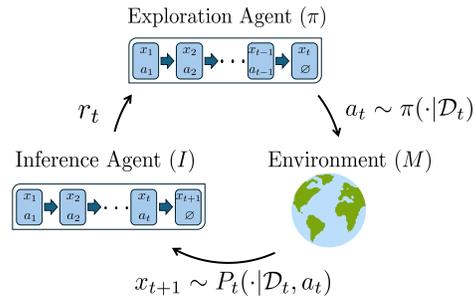
93 environment as a tuple $M = (\mathcal{X}, \mathcal{A}, P, \rho)$, where \mathcal{X} is a set of possible observations, \mathcal{A} is a finite
 94 set of actions, $P = (P_t)_{t \in \mathbb{N}}$ denotes the transition functions, with $P_t : (\mathcal{X} \times \mathcal{A})^t \rightarrow \Delta(\mathcal{X})$ and
 95 $\rho \in \Delta(\mathcal{X})$ denotes the initial observation distribution. All the environments in a class \mathcal{M} share
 96 the same set of observations \mathcal{X} and set of actions \mathcal{A} . The learner interacts with the environment in
 97 a sequential manner: (1) an initial observation $x_1 \sim \rho$ is sampled from \mathcal{X} ; (2) at time-step t , the
 98 learner chooses an action a_t and observes the next observation $x_{t+1} \sim P_t(\cdot | \mathcal{D}_t, a_t)$, meaning that
 99 x_{t+1} is drawn independently from $P_t(\cdot | \mathcal{D}_t, a_t)$ given a trajectory $\mathcal{D}_t = (x_1, a_1, \dots, x_{t-1}, a_{t-1}, x_t)$.
 100 Formally, the learner uses a *randomized policy* $\pi = (\pi_t)_{t \in \mathbb{N}}$, which is a sequence of deterministic
 101 functions, to select actions: action a_t is selected by sampling independently from $\pi_t(\mathcal{D}_t)$ (with \mathcal{D}_t
 102 being a random variable), where $\pi_t(\mathcal{D}_t)$ specifies a probability distribution over \mathcal{A} .

103 We assume a task-specific ground-truth hypothesis H_M^* from a predefined class \mathcal{H} of hypotheses for
 104 each environment, where our goal is to efficiently infer this hypothesis. Informally, we can state our
 105 objective as follows:

*Given an environment M drawn from $\mathcal{P}(\mathcal{M})$, how can we learn a sampling strategy π that
 collects data \mathcal{D} from M so the agent can reliably infer H_M^* solely from \mathcal{D} ?*

106

107 An oracle $h(\hat{H}; M) = \mathbf{1}_{\{\hat{H} = H_M^*\}}$ provides super-
 108 vised feedback at training time (not test time), indic-
 109 ating correctness without revealing hidden struc-
 110 tures. Using oracle feedback, we learn an inference
 111 mapping $I : \mathcal{D}_t \mapsto \Delta(\mathcal{H})$, yielding posterior distri-
 112 butions over hypotheses given collected data. The
 113 estimator $\hat{H}_t \sim I(\cdot | \mathcal{D}_t)$ guides exploration by pro-
 114 viding a reward signal to an RL agent collecting the
 115 data \mathcal{D}_t using an exploration policy π .



116 **Example: Best Arm Identification** A relevant ex-
 117 ample is that of Best-Arm Identification in MAB
 118 problems [Garivier and Kaufmann \[2016\]](#). Recall that in a MAB problem the decision maker can
 119 choose between K different actions a_1, \dots, a_K (we also say *arms*) at each time-step. Upon selecting
 120 an action a at time t , it observes a random reward r_t distributed according to a distribution ν_{a_t} . In
 121 BAI the goal is to identify the best action $a^* = \arg \max_a \mathbb{E}_{R \sim \nu_a} [R]$ as quickly as possible (hence
 122 $H^* = a^*$). While several algorithms have been provided for different settings [Soare et al. \[2014\]](#),
 123 [Jedra and Proutiere \[2020\]](#), [Russo and Proutiere \[2023b\]](#), [Kocák and Garivier \[2020\]](#), [Poiani et al. \[2024\]](#),
 124 a major issue is that the algorithm design can drastically change if the assumptions change.
 125 Moreover, it is difficult to design efficient techniques for more complex settings such as MDPs (in
 126 fact, the problem becomes non-convex [Marjani and Proutiere \[2021\]](#), [Russo and Pacchiano \[2025\]](#)).
 127 Therefore, in this work we address the open question of whether it is possible to learn efficient
 128 exploration strategies directly from experience, avoiding the process of designing a BAI algorithm.

129 2.1 ICPE for Fixed Confidence Problems

130 In this work, we focus on the fixed confidence setting [Garivier and Kaufmann \[2016\]](#). In this setting,
 131 the agent needs to learn to stop the data sampling process as soon as it is sufficiently confident
 132 to have correctly estimated H^* for an environment M . Let \mathbb{P}_M^π be the underlying probability
 133 measure of the process $((\mathcal{D}_t, a_t))_t$ under a sampling strategy π . In the following we also write
 134 $\mathbb{P}_{M \sim \mathcal{P}(\mathcal{M})}^\pi(\cdot) = \mathbb{E}_{M \sim \mathcal{P}(\mathcal{M})}[\mathbb{P}_M^\pi(\cdot)]$ to denote the expected probability over the prior.

135 We equip the learner with the capability to stop the sampling process at any point in time. We denote
 136 such stopping rule by τ , which is a stopping time with respect to the filtration $(\sigma(\mathcal{D}_t))_t$. Then,
 137 the learner wishes to find an optimal stopping rule τ (with $\tau < \infty$ a.s.), exploration policy π and
 138 inference network I subject to a confidence level at the stopping time τ :

$$\min_{\tau, \pi, I} \mathbb{E}_{M \sim \mathcal{P}(\mathcal{M})}[\tau] \quad \text{s.t.} \quad \mathbb{P}_{M \sim \mathcal{P}(\mathcal{M})}^\pi(h(\hat{H}_\tau; M) = 1) \geq 1 - \delta. \quad (1)$$

Algorithm 1 ICPE (In-Context Pure Exploration) - Fixed Confidence

- 1: **Input:** Tasks distribution $\mathcal{P}(\mathcal{M})$; confidence δ ; learning rates α, β ; initial λ and hyper-parameters T, N, η .
- 2: Initialize buffer \mathcal{B} , networks Q_θ, I_ϕ and set $\bar{\theta} \leftarrow \theta, \bar{\phi} \leftarrow \phi$.
- 3: **while** Training is not over **do**
- 4: Sample environment $M \sim \mathcal{P}(\mathcal{M})$ with true hypothesis H^* , observe $s_1 \sim \rho$ and set $t \leftarrow 1$.
- 5: **repeat**
- 6: Execute action $a_t = \arg \max_a Q_\theta(s_t, a)$ in M and observe next state s_{t+1} .
- 7: Add experience $z_t = (s_t, a_t, s_{t+1}, d_t = \mathbf{1}\{s_{t+1} \text{ is terminal}\}, H^*)$ to \mathcal{B} .
- 8: Set $t \leftarrow t + 1$.
- 9: **until** $a_{t-1} = a_{\text{stop}}$ or $t > N$.
- 10: Update variable λ according to

$$\lambda \leftarrow \max(0, \lambda - \beta(I_\phi(H^*|s_\tau) - 1 + \delta)). \quad (2)$$

- 11: Sample batches $B, B' \sim \mathcal{B}$ and update θ, ϕ as

$$\theta \leftarrow \theta - \alpha \nabla_\theta \frac{1}{|B|} \sum_{z \in B} \left[\mathbf{1}_{\{a \neq a_{\text{stop}}\}} (y_\lambda(z) - Q_\theta(s, a))^2 + (r_\lambda(z_{\text{stop}}) - Q_\theta(s, a_{\text{stop}}))^2 \right], \quad (3)$$

$$\phi \leftarrow \phi + \alpha \nabla_\phi \frac{1}{|B'|} \sum_{z \in B'} [\log(I_\phi(H^*|s'))]. \quad (4)$$

- 12: Update $\bar{\theta} \leftarrow (1 - \eta)\bar{\theta} + \eta\theta$ and every T steps set $\bar{\phi} \leftarrow \phi$.
 - 13: **end while**
-

139 Introducing a stopping action a_{stop} to π_t , we define $\tau = \min(N, \inf t : a_t = a_{\text{stop}})$ for a maximum
 140 horizon N (the horizon is introduced for practical reasons). We consider solving the dual formulation:

$$\min_{\lambda \geq 0} \max_{\pi, I} V_\lambda(\pi, I) = -\mathbb{E}_{M \sim \mathcal{P}(\mathcal{M})}^\pi [\tau] + \lambda \left[\mathbb{P}_{M \sim \mathcal{P}(\mathcal{M})}^\pi \left(h(\hat{H}_\tau; M) = 1 \right) - 1 + \delta \right],$$

141 with $\hat{H}_\tau \sim I(\cdot | \mathcal{D}_\tau)$. To solve this problem, **ICPE** treats each optimization separately, and optimize
 142 using a descent-ascent scheme. **ICPE** leverages transformers to encode trajectories \mathcal{D}_t as fixed-length
 143 states $s_t = (\mathcal{D}_t, \emptyset_{t:N})$ of an induced MDP M , padding with null tokens to horizon N . The resulting
 144 MDP formulation has actions $\mathcal{A} \cup a_{\text{stop}}$ and a reward structure penalizing each step until stopping,
 145 defined below.

146 **Learning I .** The distribution I is modeled using a transformer with parameter ϕ , and we denote it
 147 by I_ϕ . Then, considering a fixed (π, λ) , the maximization with respect to I amounts to solving

$$\max_{\phi} \mathbb{E}_{M \sim \mathcal{P}(\mathcal{M})}^\pi [h(\hat{H}_\tau; M)], \quad \hat{H}_\tau \sim I_\phi(\cdot | s_\tau).$$

148 Therefore, we can train ϕ via a cross-entropy loss $-\sum_{H'} h(H'; M) \log(I_\phi(H' | s_\tau))$.

149 **Learning π .** The policy π is learnt using RL techniques. We define a reward r that penalizes
 150 the agent at all time-steps, that is $r_t = -1$, while at the stopping-time we have $r_\tau = -1 +$
 151 $\lambda \mathbb{E}_{H \sim I(\cdot | s_\tau)} [h(H; M)]$. Accordingly, one can define the Q -value of (π, I, λ) in a state-action pair
 152 (s, a) as $Q_\lambda^{\pi, I}(s, a) = \mathbb{E}_{M \sim \mathcal{P}(\mathcal{M})}^\pi \left[\sum_{n=t}^\tau r_n \mid s_t = s, a_t = a \right]$, with $a_n \sim \pi_n(\cdot | s_n)$.

153 We model π with a transformer of parameter θ , and train it using DQN Mnih et al. [2015], Van Hasselt
 154 et al. [2016] with a replay buffer \mathcal{B} and a target network $Q_{\bar{\theta}}$ parameterized by $\bar{\theta}$. To maintain timescale
 155 separation, we introduce a separate target inference network $I_{\bar{\phi}}$, parameterized by $\bar{\phi}$, which provides
 156 feedback for training θ . Note that, as discussed earlier, we introduce a dedicated stop-action a_{stop}
 157 whose value depends solely on history. Thus, its Q -value can be updated at any time, allowing
 158 retrospective evaluation of stopping. For learning the Q -values, we define the reward for a transition
 159 $z = (s, a, s', d, H^*)$ as:

$$r_\lambda(z) := -1 + d\lambda \log I_{\bar{\phi}}(H^* | s'), \quad d = \mathbf{1}\{z \text{ terminal}\},$$

160 where we set $s' \leftarrow s$ if $a = a_{\text{stop}}$, and terminal means either $a = a_{\text{stop}}$ or the last step in the horizon.
 161 We also define the transition z_{stop} by replacing (a, s') with (a_{stop}, s) in z . Then, for $a \neq a_{\text{stop}}$, the
 162 Q -values can be learned using a target value:

$$y_\lambda(z) = r_\lambda(z) + (1 - d) \max_i Q_{\bar{\theta}}(s', a_i).$$

163 Instead, for the stopping action, we use the loss $(r_\lambda(z_{\text{stop}}) - Q_\theta(s, a_{\text{stop}}))^2$. Therefore, the overall
 164 loss used for training θ on a transition z is:

$$\mathbf{1}_{\{a \neq a_{\text{stop}}\}} (y_\lambda(z) - Q_\theta(s, a))^2 + (r_\lambda(z_{\text{stop}}) - Q_\theta(s, a_{\text{stop}}))^2,$$

165 where $\mathbf{1}_{\{a \neq a_{\text{stop}}\}}$ avoids double accounting for the stopping action.

166 **Last steps.** Then, to train (θ, ϕ) , we sample two independent batches $(B, B') \sim \mathcal{B}$ from the buffer,
 167 and compute the gradient updates as in eqs. (3) and (4) of algorithm 1. We periodically update target
 168 networks, setting $\bar{\phi} \leftarrow \phi$ every T steps and using a Polyak averaging $\bar{\theta} \leftarrow (1 - \eta)\bar{\theta} + \eta\theta$, with
 169 $\eta \in (0, 1)$.

170 Finally, we update λ by assessing the confidence of I_ϕ at the stopping time (2) for a fixed (π, I) .
 171 Thus, for sufficiently small learning rates, optimizing (λ, θ, ϕ) resembles an ascent-descent scheme.

172 3 Empirical Evaluation

173 We evaluate our approach across various tasks: stochastic bandits with or without latent structure;
 174 learning a probabilistic version of binary search. Due to space limitations, we refer the reader to
 175 appendix C for more details and more experiments on MAB problem with feedback graphs Russo
 176 et al. [2025], MDPs with hidden information and an analysis of ICPE in the broader setting of
 177 classifying images by sequentially revealing image patches.

178 **Algorithms.** In our evaluations we compare to different algorithms, depending on the problem. Some
 179 of the algorithms include: uniform sampling, TaS (Track and Stop) Garivier and Kaufmann [2016],
 180 TTPS (Top Two Sampling) Russo et al. [2018]. We also include a variant of IDS Russo and Van Roy
 181 [2018] based on the I -mapping, which uses the observation that I defines a posterior distribution over
 182 \mathcal{H} . Always based on this idea, we also introduce I -DPT, a variant of DPT Lee et al. [2023], based on
 183 the fact that I can be used to explore a problem à-la Thompson Sampling. More information about
 184 these methods, and their hyper-parameters, can be found in appendix B¹.

185 3.1 Bandit Problems

186 We now apply ICPE to the classical BAI problem within MAB tasks. For the MAB setting we have a
 187 finite number of actions $\mathcal{A} = \{1, \dots, K\}$, corresponding to the actions in the MAB problem M . For
 188 each action a , we define a corresponding reward distribution ν_a from which rewards are sampled i.i.d.
 189 Then, $\mathcal{P}(\mathcal{M})$ is a prior distribution on the actions' rewards distributions $(\nu_a)_a$ and for BAI we let
 190 $H^* = \arg \max_a \mathbb{E}_{r \sim \nu_a} [r]$, so that we need to identify the best action. Lastly, the observation at time
 191 t is $x_t = (a_t, r_t)$, where a_t is the chosen action at time t and r_t is a reward sampled from ν_{a_t} .

192 **Stochastic Bandit Problems.** We evaluate ICPE on stochastic bandit environments with $\delta = 0.1$
 193 and $N = 100$. Each action's reward distribution is normally distributed $\nu_a = \mathcal{N}(\mu_a, 0.5^2)$, with
 194 $(\mu_a)_{a \in \mathcal{A}}$ drawn from $\mathcal{P}(\mathcal{M})$. In this case $\mathcal{P}(\mathcal{M})$ is a uniform distribution over problems with
 195 minimum gap $\max_a \mu_a - \max_{b \neq a} \mu_b \geq \Delta_0$, with $\Delta_0 = 0.4$. Hence, an algorithm could exploit
 196 this property to infer H^* more quickly. For this case, we also derive some sample complexity
 bounds in appendix A. Figure 1 summarizes the results for this setting. We compare to TaS and

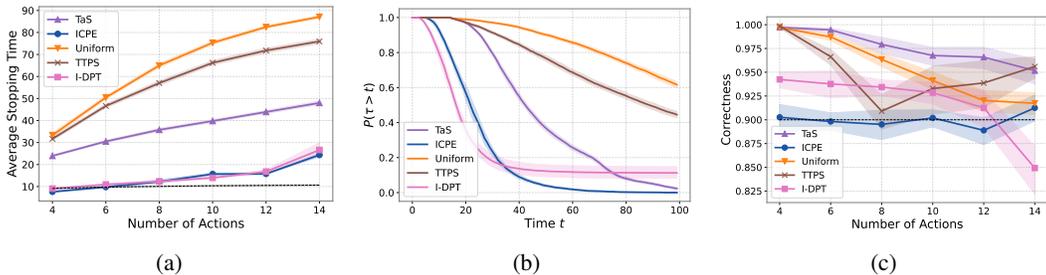


Figure 1: Results for stochastic MABs with fixed confidence $\delta = 0.1$ and $N = 100$: (a) average stopping time τ ; (b) survival function of τ ; (c) probability of correctness $\mathbb{P}_{M \sim \mathcal{P}(\mathcal{M})}^\pi (h(\hat{H}; M) = 1)$.

197 TTPS, and use the stopping rule of TaS also for Uniform and TTPS (the stopping rule is based on a
 198

¹In the results, shaded areas indicate 95% confidence intervals, computed via hierarchical bootstrapping.

199 self-normalized process, compared with a threshold function $\beta(t, \delta)$; see also appendix B for more
 200 details). Overall, we see in fig. 1a how ICPE is able to find a more efficient strategy compared to
 201 classical techniques. Interestingly, also I-DPT seems to achieve relatively small sample complexities.
 202 However, its tail distribution of τ is rather large compared to ICPE (fig. 1b) and the correctness
 203 is smaller than $1 - \delta$ for large values of K . Methods like TaS and TTPS achieve larger sample
 204 complexity, but also larger correctness values (fig. 1c). This is due to the fact that it is hard to define
 205 stopping rules. In fact, it is well known that current theoretically sound stopping rules are overly
 206 conservative Garivier and Kaufmann [2016]. Nonetheless, even using a less conservative rule such
 207 as $\beta(t, \delta) = \log((1 + \log(t))/\delta)$, which is what we use (and, yet, has not been proven to guarantee
 208 δ -correctness), is still conservative. The fact that ICPE can achieve the right value of confidence can
 209 help discover potential ways to define stopping rules. Lastly, in fig. 1a in black we show a complexity
 210 bound (proof in appendix A.1). While seemingly constant, it is actually *slowly* increasing in the
 211 number of arms.

212 **Bandit Problems with Hidden Information.** To evaluate ICPE in structured settings, we introduce
 213 bandit environments with latent informational dependencies, termed *magic actions*. In the single
 214 magic action case, the magic action a_m 's reward is distributed according to $\mathcal{N}(\mu_{a_m}, \sigma_m^2)$, where
 215 $\sigma_m \in (0, 1)$ and $\mu_{a_m} := \phi(\arg \max_{a \neq a_m} \mu_a)$ encodes information about the optimal action's identity
 216 through an invertible mapping ϕ that is unknown to the learner. The index a_m is fixed, and the mean
 217 rewards of the other actions $(\mu_a)_{a \neq a_m}$ are sampled from $\mathcal{P}(\mathcal{M})$, a uniform distribution over models
 218 guaranteeing that a_m , as defined above, is not optimal (see appendices A.2 and C.1.2 for more details).
 219 Then, we define the reward distribution of the non-magic actions as $\mathcal{N}(\mu_a, (1 - \sigma_m)^2)$.

220 In our first experiment, we vary the standard deviation σ_m in $[0, 1]$. Thus, agents must balance
 221 sampling between informative and noisy actions based on varying uncertainty levels. We evaluate
 222 ICPE in a fixed-confidence setting with error rate $\delta = 0.1$. Figure 2a compares ICPE's sample
 223 complexity against a theoretical lower bound (see appendix A) and an informed baseline, denoted as
 224 *I-IDS*, which performs standard IDS leveraging ICPE's trained inference network *I* for exploiting
 the magic action (details in Appendix B). ICPE achieves sample complexities close to the theoretical

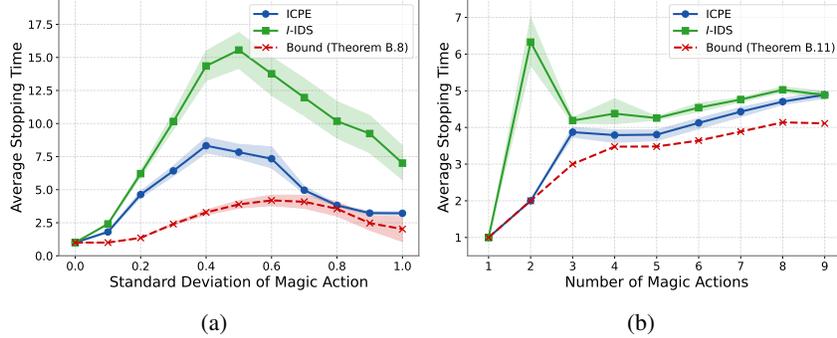


Figure 2: (a) Single magic action: average stopping time and the theoretical lower bound across varying σ_m . (b) Magic chain: average stopping time between ICPE, *I-IDS* vs. number of magic actions.

225 bound across all tested noise levels, consistently outperforming *I-IDS*. To further challenge ICPE, we
 226 introduce a multi-layered "magic chain" bandit environments, where there is a sequence of n
 227 magic actions $\mathcal{A}_m := \{a_{i_1}, \dots, a_{i_n}\} \subset \mathcal{A}$ such that $\mu_{a_{i_j}} = \phi(\mu_{a_{i_{j+1}}})$, and $\mu_{a_{i_n}} = \phi(\arg \max_{a \notin \mathcal{A}_m} \mu_a)$.
 228 The first index i_1 is known, and by following the chain, an agent can uncover the best action in n
 229 steps. However, the optimal sample complexity depends on the ratio of magic actions to non-magic
 230 arms. Varying the number of magic actions from 1 to 9 in a 10-actions environment, Figure 2b
 231 demonstrates ICPE's empirical performance, outperforming *I-IDS*.
 232

233 **Bandit Problems with Feedback Graphs.** In bandit problems, playing action u yields its reward,
 234 while full-information settings reveal all rewards. Feedback graphs interpolate between these ex-
 235 tremes: a directed graph $G \in [0, 1]^{K \times K}$ specifies that choosing u reveals the reward of v with
 236 probability $G_{u,v}$. Although feedback graphs have been extensively studied for regret minimization
 237 Mannor and Shamir [2011], their role in pure exploration remains underexplored Russo et al.
 238 [2025]; here we use them as structured testbeds, where latent relational and stochastic dependencies

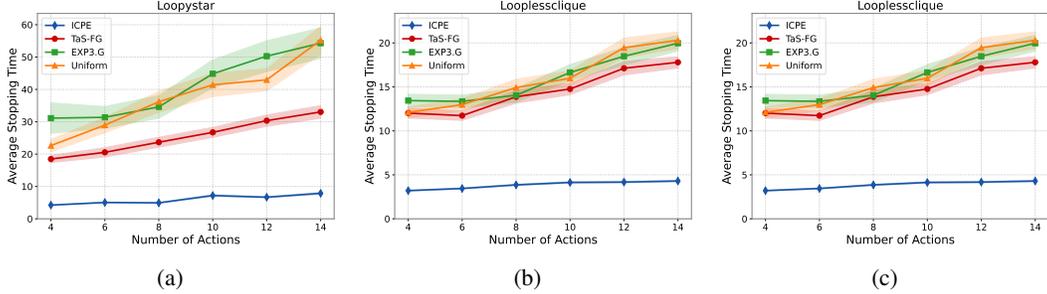


Figure 3: Sample complexity comparison under the fixed-confidence setting for: (a) Loopy Star, (b) Loopless Clique, and (c) Ring graphs.

239 must be inferred to explore efficiently. Formally, upon playing u the learner observes for each
 240 $v \in [K]$:

$$r_v \sim \begin{cases} \mathcal{N}(\mu_v, \sigma^2), & \text{with probability } G_{u,v}, \\ \text{no observation,} & \text{otherwise.} \end{cases}$$

241 We tested **ICPE** on 3 different graph families with $\delta = 0.1$: the loopy star graph, the ring graph and
 242 the loopless clique [Russo et al. \[2025\]](#). We set the optimal arm’s mean to 1 and all others to 0.5 to
 243 facilitate faster convergence. We compared it to Uniform Sampling, EXP3.G, and Tas-FG using a
 244 shared stopping rule from [Russo et al. \[2025\]](#).

245 As shown in Figure 3, **ICPE** consistently achieves significantly lower sample complexity, suggesting
 246 that that **ICPE** is able to meta-learn and leverage the underlying structure of the graph.

247 3.2 Algorithm Discovery: Meta-Learning Binary Search

248 To test **ICPE**’s ability to recover classical exploration algorithms, we evaluate whether it can au-
 249 tonomously meta-learn binary search. We define an action space of $\mathcal{A} = \{1, \dots, K\}$, where K is the
 250 upper bound on the possible location of the hidden target $H^* \sim \mathcal{A}$. Pulling an arm above or below
 251 H^* yields a observation $x_t = -1$ or $x_t = +1$, respectively—providing directional feedback. We
 252 train **ICPE** under the fixed-confidence setting for $K = 2^3, \dots, 2^8$ using a target error rate of $\delta = 0.01$.
 253 In table 1 we report results on 100 held-out tasks per setting. **ICPE** consistently achieves perfect
 254 accuracy with worst-case stopping times that match the optimal $\log_2(K)$ rate, demonstrating that it
 255 has successfully rediscovered binary search purely from experience. While simple, this task illustrates
 256 **ICPE**’s broader potential to learn efficient search strategies in domains where no hand-designed
 257 algorithm is available.

K (Actions)	Min Accuracy	Mean Stop Time	Max Stop Time	$\log_2 K$
8	1.00	2.13 ± 0.12	3	3
16	1.00	2.93 ± 0.12	4	4
32	1.00	3.71 ± 0.15	5	5
64	1.00	4.50 ± 0.21	6	6
128	1.00	5.49 ± 0.23	7	7
256	1.00	6.61 ± 0.26	8	8

Table 1: **ICPE** performance on the binary search task as the number of actions K increases.

258 4 Conclusions

259 In this work, we addressed the design of efficient pure-exploration strategies for the *active sequential*
 260 *hypothesis testing* problem, where an agent sequentially selects samples to rapidly identify the true
 261 hypothesis. While particularly relevant across different domains, it is difficult to design optimal
 262 strategies in the presence of hidden structure, and most of the existing optimal strategies are restricted
 263 to simple cases for unstructured multi-armed bandit problems. To overcome these limitations, we
 264 introduced **ICPE**, an in-context learning framework that leverages Transformers to learn exploration
 265 policies directly from experience. Our results demonstrate that **ICPE** is able to autonomously
 266 discovering task-specific adaptive exploration strategies. We believe our work makes a fundamental
 267 contribution to active testing, and in particular to the sub-field of best-arm identification. Future
 268 directions include several directions, including a theoretical analysis of **ICPE**’s guarantees and scaling
 269 **ICPE** to larger, higher-dimensional problems.

270 References

- 271 Aymen Al Marjani, Aurélien Garivier, and Alexandre Proutiere. Navigating to the best policy in
272 markov decision processes. *Advances in Neural Information Processing Systems*, 34:25852–25864,
273 2021.
- 274 Jean-Yves Audibert and Sébastien Bubeck. Best arm identification in multi-armed bandits. In
275 *COLT-23th Conference on learning theory-2010*, pages 13–p, 2010.
- 276 Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine*
277 *Learning Research*, 3(Nov):397–422, 2002.
- 278 Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit
279 problem. *Machine learning*, 47:235–256, 2002.
- 280 Peter Auer, Thomas Jaksch, and Ronald Ortner. Near-optimal regret bounds for reinforcement
281 learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 21, 2008.
- 282 Jacob Beck, Risto Vuorio, Evan Zheran Liu, Zheng Xiong, Luisa Zintgraf, Chelsea Finn, and Shimon
283 Whiteson. A survey of meta-reinforcement learning. *arXiv preprint arXiv:2301.08028*, 2023.
- 284 Yoshua Bengio, Samy Bengio, and Jocelyn Cloutier. *Learning a synaptic learning rule*. Citeseer,
285 1990.
- 286 Scott M Berry, Bradley P Carlin, J Jack Lee, and Peter Muller. *Bayesian adaptive methods for clinical*
287 *trials*. CRC press, 2010.
- 288 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
289 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
290 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 291 Olivier Cappé, Aurélien Garivier, Odalric-Ambrym Maillard, Rémi Munos, and Gilles Stoltz.
292 Kullback-leibler upper confidence bounds for optimal sequential allocation. *The Annals of Statistics*,
293 pages 1516–1541, 2013.
- 294 Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel,
295 Arvind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence
296 modeling. *Advances in neural information processing systems*, 34:15084–15097, 2021.
- 297 Herman Chernoff. Sequential design of experiments. *The Annals of Mathematical Statistics*, 30(3):
298 755–770, 1959.
- 299 Herman Chernoff. *Sequential design of experiments*. Springer, 1992.
- 300 David A Cohn, Zoubin Ghahramani, and Michael I Jordan. Active learning with statistical models.
301 *Journal of artificial intelligence research*, 4:129–145, 1996.
- 302 Zhenwen Dai, Federico Tomasi, and Sina Ghiassian. In-context exploration-exploitation for rein-
303 forcement learning. *arXiv preprint arXiv:2403.06826*, 2024.
- 304 Rémy Degenne and Wouter M Koolen. Pure exploration with multiple correct answers. *Advances in*
305 *Neural Information Processing Systems*, 32, 2019.
- 306 Rémy Degenne, Wouter M Koolen, and Pierre Ménard. Non-asymptotic pure exploration by solving
307 games. *Advances in Neural Information Processing Systems*, 32, 2019.
- 308 Rémy Degenne, Pierre Ménard, Xuedong Shang, and Michal Valko. Gamification of pure exploration
309 for linear bandits, 2020.
- 310 Yan Duan, John Schulman, Xi Chen, Peter L Bartlett, Ilya Sutskever, and Pieter Abbeel. Rl^2 : Fast
311 reinforcement learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*, 2016.
- 312 Scott Emmons, Benjamin Eysenbach, Ilya Kostrikov, and Sergey Levine. Rvs: What is essential for
313 offline rl via supervised learning? *arXiv preprint arXiv:2112.10751*, 2021.

- 314 Kyra Gan, Su Jia, and Andrew Li. Greedy approximation algorithms for active sequential hypothesis
315 testing. *Advances in Neural Information Processing Systems*, 34:5012–5024, 2021.
- 316 Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers learn
317 in-context? a case study of simple function classes. *Advances in Neural Information Processing*
318 *Systems*, 35:30583–30598, 2022.
- 319 Aurélien Garivier and Emilie Kaufmann. Optimal best arm identification with fixed confidence. In
320 *Conference on Learning Theory*, pages 998–1027. PMLR, 2016.
- 321 Bashkar K Ghosh. A brief history of sequential analysis. *Handbook of sequential analysis*, 1, 1991.
- 322 Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1.
323 MIT press Cambridge, 2016.
- 324 Aditya Gopalan, Shie Mannor, and Yishay Mansour. Thompson sampling for complex online
325 problems. In *International conference on machine learning*, pages 100–108. PMLR, 2014.
- 326 Yassir Jedra and Alexandre Proutiere. Optimal best-arm identification in linear bandits. *Advances in*
327 *Neural Information Processing Systems*, 33:10007–10017, 2020.
- 328 Marc Jourdan, Rémy Degenne, Dorian Baudry, Rianne de Heide, and Emilie Kaufmann. Top two
329 algorithms revisited. *Advances in Neural Information Processing Systems*, 35:26791–26803, 2022.
- 330 John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger,
331 Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate
332 protein structure prediction with alphafold. *nature*, 596(7873):583–589, 2021.
- 333 Emilie Kaufmann and Wouter M Koolen. Mixture martingales revisited with applications to sequential
334 tests and confidence intervals. *Journal of Machine Learning Research*, 22(246):1–44, 2021.
- 335 Emilie Kaufmann, Nathaniel Korda, and Rémi Munos. Thompson sampling: An asymptotically
336 optimal finite-time analysis. In *International conference on algorithmic learning theory*, pages
337 199–213. Springer, 2012.
- 338 Emilie Kaufmann, Olivier Cappé, and Aurélien Garivier. On the complexity of best-arm identification
339 in multi-armed bandit models. *The Journal of Machine Learning Research*, 17(1):1–42, 2016.
- 340 Tomáš Kocák and Aurélien Garivier. Best arm identification in spectral bandits. *arXiv preprint*
341 *arXiv:2005.09841*, 2020.
- 342 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolu-
343 tional neural networks. *Advances in neural information processing systems*, 25, 2012.
- 344 Aviral Kumar, Xue Bin Peng, and Sergey Levine. Reward-conditioned policies. *arXiv preprint*
345 *arXiv:1912.13465*, 2019.
- 346 Tor Lattimore and Marcus Hutter. Pac bounds for discounted mdps. In *Algorithmic Learning Theory:*
347 *23rd International Conference, ALT 2012, Lyon, France, October 29-31, 2012. Proceedings 23*,
348 pages 320–334. Springer, 2012.
- 349 Jonathan Lee, Annie Xie, Aldo Pacchiano, Yash Chandak, Chelsea Finn, Ofir Nachum, and Emma
350 Brunskill. Supervised pretraining can learn in-context reinforcement learning. *Advances in Neural*
351 *Information Processing Systems*, 36:43057–43083, 2023.
- 352 Dennis V Lindley. On a measure of the information provided by an experiment. *The Annals of*
353 *Mathematical Statistics*, 27(4):986–1005, 1956.
- 354 Shie Mannor and Ohad Shamir. From bandits to experts: On the value of side-observations. *Advances*
355 *in neural information processing systems*, 24, 2011.
- 356 Aymen Al Marjani and Alexandre Proutiere. Adaptive sampling for best policy identification in
357 markov decision processes. In *International Conference on Machine Learning*, pages 7459–7468.
358 PMLR, 2021.

- 359 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare,
360 Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control
361 through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- 362 Amir Moeini, Jiuqi Wang, Jacob Beck, Ethan Blaser, Shimon Whiteson, Rohan Chandra, and
363 Shangtong Zhang. A survey of in-context reinforcement learning. *arXiv preprint arXiv:2502.07978*,
364 2025.
- 365 Subhojyoti Mukherjee, Ardhendu S Tripathy, and Robert Nowak. Chernoff sampling for active
366 testing and extension to active regression. In *International Conference on Artificial Intelligence
367 and Statistics*, pages 7384–7432. PMLR, 2022.
- 368 Kevin P Murphy. *Probabilistic machine learning: Advanced topics*. MIT press, 2023.
- 369 Mohammad Naghshvar and Tara Javidi. Active sequential hypothesis testing. *The Annals of Statistics*,
370 41(6):2703–2738, 2013.
- 371 Mohammad Naghshvar, Tara Javidi, and Kamalika Chaudhuri. Noisy bayesian active learning. In
372 *2012 50th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*,
373 pages 1626–1633. IEEE, 2012.
- 374 Ian Osband, Daniel Russo, and Benjamin Van Roy. (more) efficient reinforcement learning via
375 posterior sampling. *Advances in Neural Information Processing Systems (NeurIPS)*, 26, 2013.
- 376 Riccardo Poiani, Marc Jourdan, Emilie Kaufmann, and Rémy Degenne. Best-arm identification in
377 unimodal bandits. *arXiv preprint arXiv:2411.01898*, 2024.
- 378 Paul Resnick and Hal R Varian. Recommender systems. *Communications of the ACM*, 40(3):56–58,
379 1997.
- 380 Alessio Russo and Aldo Pacchiano. Adaptive exploration for multi-reward multi-policy evaluation.
381 *arXiv preprint arXiv:2502.02516*, 2025.
- 382 Alessio Russo and Alexandre Proutiere. Model-free active exploration in reinforcement learning.
383 *Advances in Neural Information Processing Systems*, 36:54740–54753, 2023a.
- 384 Alessio Russo and Alexandre Proutiere. On the sample complexity of representation learning in
385 multi-task bandits with global and local structure. In *Proceedings of the AAAI Conference on
386 Artificial Intelligence*, volume 37, pages 9658–9667, 2023b.
- 387 Alessio Russo and Filippo Vannella. Multi-reward best policy identification. *Advances in Neural
388 Information Processing Systems*, 37:105583–105662, 2025.
- 389 Alessio Russo, Yichen Song, and Aldo Pacchiano. Pure exploration with feedback graphs. In *Pro-
390 ceedings of The 28th International Conference on Artificial Intelligence and Statistics*, Proceedings
391 of Machine Learning Research. PMLR, 2025.
- 392 Daniel Russo and Benjamin Van Roy. Learning to optimize via posterior sampling. *Mathematics of
393 Operations Research*, 39(4):1221–1243, 2014.
- 394 Daniel Russo and Benjamin Van Roy. Learning to optimize via information-directed sampling.
395 *Operations Research*, 66(1):230–252, 2018.
- 396 Daniel J Russo, Benjamin Van Roy, Abbas Kazerouni, Ian Osband, Zheng Wen, et al. A tutorial on
397 thompson sampling. *Foundations and Trends® in Machine Learning*, 11(1):1–96, 2018.
- 398 Tom Schaul and Jürgen Schmidhuber. Metalearning. *Scholarpedia*, 5(6):4650, 2010.
- 399 David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez,
400 Marc Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore Graepel, et al. A general reinforcement
401 learning algorithm that masters chess, shogi, and go through self-play. *Science*, 362(6419):
402 1140–1144, 2018.
- 403 Marta Soare, Alessandro Lazaric, and Rémi Munos. Best-arm identification in linear bandits.
404 *Advances in neural information processing systems*, 27, 2014.

- 405 Rupesh Kumar Srivastava, Pranav Shyam, Filipe Mutz, Wojciech Jaskowski, and Jürgen Schmidhuber.
406 Training agents using upside-down reinforcement learning. *CoRR*, abs/1912.02877, 2019. URL
407 <http://arxiv.org/abs/1912.02877>.
- 408 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- 409 Nidhin Koshy Vaidhiyan, SP Arun, and Rajesh Sundaresan. Active sequential hypothesis testing with
410 application to a visual search problem. In *2012 IEEE International Symposium on Information*
411 *Theory Proceedings*, pages 2201–2205. IEEE, 2012.
- 412 Hado Van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-
413 learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 30, 2016.
- 414 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
415 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*
416 *systems*, 30, 2017.

417 Appendix

418 Contents

419	1 Introduction	1
420	1.1 Related Work	2
421	2 Learning to Explore: In-Context Pure Exploration	2
422	2.1 ICPE for Fixed Confidence Problems	3
423	3 Empirical Evaluation	5
424	3.1 Bandit Problems	5
425	3.2 Algorithm Discovery: Meta-Learning Binary Search	7
426	4 Conclusions	7
427	A Theoretical Results	13
428	A.1 Sample Complexity Bounds for MAB Problems with Fixed Minimum Gap	13
429	A.2 Sample Complexity Lower Bound for the Magic Action MAB Problem	15
430	A.3 Sample Complexity Bound for the Multiple Magic Actions MAB Problem	19
431	B Algorithms	23
432	B.1 ICPE with Fixed Confidence	23
433	B.2 Other Algorithms	25
434	B.2.1 Track and Stop	25
435	B.2.2 I -IDS	26
436	B.2.3 I -DPT	26
437	B.3 Transformer Architecture	27
438	C Experiments	28
439	C.1 Bandit Problems	28
440	C.1.1 Stochastic Bandits Problems	28
441	C.1.2 Bandit Problems with Hidden Information	29
442	C.2 Exploration on Feedback Graphs	31
443	C.3 Meta-Learning Binary Search	32

444 **Appendix**

445 **A Theoretical Results**

446 In this section we provide different theoretical results, mainly for the sample complexity of different
447 MAB problems with structure.

448 **A.1 Sample Complexity Bounds for MAB Problems with Fixed Minimum Gap**

449 We now derive a sample complexity lower bound for a MAB problem where the minimum gap is
450 known and the rewards are normally distributed.

451 Consider a MAB problem with K arms $\{1, \dots, K\}$. To each arm a is associated a reward distribution
452 $\nu_a = \mathcal{N}(\mu_a, \sigma^2)$ that is simply a Gaussian distribution. Let $a^*(\mu) = \arg \max_a \mu_a$, and define the
453 gap in arm a to be $\Delta_a(\mu) = \mu_{a^*(\mu)} - \mu_a$. In the following, without loss of generality, we assume
454 that $a^*(\mu) = 1$.

455 We define the minimum gap to be $\Delta_{\min}(\mu) = \min_{a \neq a^*(\mu)} \Delta_a(\mu)$. Assume now to know that
456 $\Delta_{\min} \geq \Delta_0 > 0$.

457 Then, for any δ -correct algorithm, guaranteeing that at some stopping time τ the estimated optimal
458 arm \hat{a}_τ is δ -correct, i.e., $\mathbb{P}_\mu(\hat{a}_\tau \neq a^*(\mu)) \leq \delta$, we have the following result.

459 **Theorem A.1.** *Consider a model μ satisfying $\Delta_{\min} \geq \Delta_0 > 0$. Then, for any δ -probably correct
460 method Alg, with $\delta \in (0, 1/2)$, we have that the optimal sample complexity is bounded as*

$$\frac{1}{\max\left(\Delta_0^2, \frac{1}{\sum_{a \neq 1} 1/\Delta_a^2}\right)} \leq \inf_{\tau: \text{Alg is } \delta\text{-correct}} \frac{\mathbb{E}_\mu[\tau]}{2\sigma^2 \text{kl}(1 - \delta, \delta)} \leq 2 \sum_a \frac{1}{(\Delta_a + \Delta_0)^2},$$

461 with $\Delta_1 = 0$ and $\text{kl}(x, y) = x \log(x/y) + (1 - x) \log((1 - x)/(1 - y))$. In particular, the solution
462 $\omega_a \propto 1/(\Delta_a + \Delta_0)^2$ (up to a normalization constant) achieves the upper bound.

463 *Proof. Step 1: Log-likelihood ratio.* The initial part of the proof is rather standard, and follows the
464 same argument used in the Best Arm Identification and Best Policy Identification literature [Garivier](#)
465 [and Kaufmann \[2016\]](#), [Russo and Vannella \[2025\]](#).

466 Define the set of models

$$\mathcal{S} = \{\mu' \in \mathbb{R}^K : \Delta_{\min}(\mu') \geq \Delta_0\},$$

467 and the set of alternative models

$$\text{Alt}(\mu) = \left\{ \mu' \in \mathcal{S} : \arg \max_a \mu'_a \neq 1 \right\}.$$

468 Take the expected log-likelihood ratio between μ and $\mu' \in \text{Alt}(\mu)$ of the data observed up to τ
469 $\Lambda_\tau = \log \frac{d\mathbb{P}_\mu(A_1, R_1, \dots, A_\tau, R_\tau)}{d\mathbb{P}_{\mu'}(A_1, R_1, \dots, A_\tau, R_\tau)}$, where A_t is the action taken in round t , and R_t is the reward observed
470 upon selecting A_t . Then, we can write

$$\Lambda_t = \sum_a \sum_{n=1}^t \mathbf{1}_{\{A_n=a\}} \log \frac{f_a(R_n)}{f'_a(R_n)}$$

471 where f_a, f'_a , are, respectively, the reward density for action a in the two models μ, μ' with respect to
472 the Lebesgue measure. Letting $N_a(t)$ denote the number of times action a has been selected up to
473 round t , by an application of Wald's lemma the expected log-likelihood ratio can be shown to be

$$\mathbb{E}_\mu[\Lambda_\tau] = \sum_a \mathbb{E}_\mu[N_a(\tau)] \text{KL}(\mu_a, \mu'_a)$$

474 where $\text{KL}(\mu_a, \mu'_a)$ is the KL divergence between two Gaussian distributions $\mathcal{N}(\mu_a, \sigma)$ and $\mathcal{N}(\mu'_a, \sigma)$
475 (note that we have σ_1 instead of σ for $a = 1$).

476 We also know from the information processing inequality [Kaufmann et al. \[2016\]](#) that $\mathbb{E}_\mu[\Lambda_\tau] \geq$
477 $\sup_{\mathcal{E} \in \mathcal{M}_\tau} \text{kl}(\mathbb{P}_\mu(\mathcal{E}), \mathbb{P}_{\mu'}(\mathcal{E}))$, where $\mathcal{M}_t = \sigma(A_1, R_1, \dots, A_t, R_t)$. We use the fact that the algo-
478 rithm is δ -correct: by choosing $\mathcal{E} = \{\hat{a}_\tau = a^*\}$ we obtain that $\mathbb{E}_\mu[\Lambda_\tau] \geq \text{kl}(1 - \delta, \delta)$, since

479 $\mathbb{P}_\mu(\mathcal{E}) \geq 1 - \delta$ and $\mathbb{P}_{\mu'}(\mathcal{E}) = 1 - \mathbb{P}_{\mu'}(\hat{a}_\tau \neq a^*) \leq 1 - \mathbb{P}_{\mu'}(\hat{a}_\tau = \arg \max_a \mu'_a) \leq \delta$ (we also used
480 the monotonicity properties of the Bernoulli KL divergence). Hence

$$\sum_a \mathbb{E}_\mu[N_a(\tau)] \text{KL}(\mu_a, \mu'_a) \geq \text{kl}(1 - \delta, \delta).$$

481 Letting $\omega_a = \mathbb{E}_\mu[N_a(\tau)]/\mathbb{E}_\mu[\tau]$, we have that

$$\mathbb{E}_\mu[\tau] \sum_a \omega_a \text{KL}(\mu_a, \mu'_a) \geq \text{kl}(1 - \delta, \delta).$$

482 Lastly, optimizing over $\mu' \in \text{Alt}(\mu)$ and $\omega \in \Delta(K)$ yields the bound:

$$\mathbb{E}_\mu[\tau] \geq T^*(\mu) \text{kl}(1 - \delta, \delta),$$

483 where $T^*(\mu)$ is defined as

$$(T^*(\mu))^{-1} = \sup_{\omega \in \Delta(K)} \inf_{\mu' \in \text{Alt}(\mu)} \sum_a \omega_a \text{KL}(\mu_a, \mu'_a).$$

484 **Step 2: Optimization over the set of alternative models.** We now face the problem of optimizing
485 over the set of alternative models.

486 Defining $\text{Alt}_a = \{\mu' \in \mathbb{R}^K : \mu'_a - \mu'_b \geq \Delta_0 \forall b \neq a\}$, the set of alternative models can be decom-
487 posed as

$$\begin{aligned} \text{Alt}(\mu) &= \left\{ \mu' \in \mathbb{R}^K : \arg \max_a \mu'_a \neq 1, \Delta_{\min}(\mu') \geq \Delta_0 \right\}, \\ &= \cup_{a \neq 1} \text{Alt}_a. \end{aligned}$$

488 Hence, the optimization problem over the alternative models becomes

$$\inf_{\mu' \in \text{Alt}(\mu)} \sum_a \omega_a \text{KL}(\mu_a, \mu'_a) = \min_{\bar{a} \neq 1} \inf_{\mu' \in \text{Alt}_{\bar{a}}} \sum_a \omega_a \frac{(\mu_a - \mu'_a)^2}{2\sigma^2}.$$

489 The inner infimum over μ' can then be written as

$$\begin{aligned} P_{\bar{a}}^*(\omega) &:= \inf_{\mu' \in \mathbb{R}^K} \sum_a \omega_a \frac{(\mu_a - \mu'_a)^2}{2\sigma^2}. \\ &\text{s.t. } \mu'_{\bar{a}} - \mu'_b \geq \Delta_0 \quad \forall b \neq \bar{a}. \end{aligned} \tag{5}$$

490 While the problem is clearly convex, it does not yield an immediate closed form solution.

491 To that aim, we try to derive a lower bound and an upper bound of the value of this minimization
492 problem.

493 **Step 3: Upper bound on $P_{\bar{a}}^*$.** Note that an upper bound on $\min_{\bar{a} \neq 1} P_{\bar{a}}^*(\omega)$ can be found by finding a
494 feasible solution μ' . Consider then the solution $\mu'_1 = \mu_1 - \Delta$, $\mu'_{\bar{a}} = \mu_1$ and $\mu'_b = \mu_b$ for all other
495 arms. Clearly We have that $\mu'_{\bar{a}} - \mu'_b \geq \Delta_0$ for all $b \neq \bar{a}$. Hence, we obtain

$$\min_{\bar{a} \neq 1} P_{\bar{a}}^*(\omega) \leq \omega_1 \frac{\Delta_0^2}{2\sigma^2} + \min_{\bar{a} \neq 1} \omega_{\bar{a}} \frac{\Delta_{\bar{a}}^2}{2\sigma^2}.$$

496 At this point, one can easily note that if $\frac{\Delta_0^2}{2\sigma^2} \geq \frac{1}{2\sigma^2 \sum_{a \neq 1} \frac{1}{\Delta_a^2}}$, then $\sup_{\omega \in \Delta(K)} \min_{\bar{a} \neq 1} P_{\bar{a}}^*(\omega) \leq \frac{\Delta_0^2}{2\sigma^2}$.

497 This corresponds to the case where all the mass is given to $\omega_1 = 1$. Otherwise, the solution is to set
498 $\omega_1 = 0$ and $\omega_a = \frac{1/\Delta_a^2}{\sum_b 1/\Delta_b^2}$ for $a \neq 1$.

499 Hence, we conclude that

$$(T^*(\mu))^{-1} = \sup_{\omega \in \Delta(K)} \min_{\bar{a} \neq 1} P_{\bar{a}}^*(\omega) \leq \frac{1}{2\sigma^2} \max \left(\Delta_0^2, \frac{1}{\sum_{a \neq 1} 1/\Delta_a^2} \right).$$

500 **Step 4: Lower bound on $P_{\bar{a}}^*$.** For the lower bound, note that we can relax the constraint to only
 501 consider $\mu'_{\bar{a}} - \mu'_1 \geq \Delta_0$. This relaxation enlarges the feasible set, and thus the infimum of this new
 502 problem lower bounds $P_{\bar{a}}^*(\omega)$.

503 By doing so, since the other arms are not constrained, by convexity of the KL divergence at the
 504 infimum we have $\mu'_b = \mu_b$ for all $b \notin \{1, \bar{a}\}$. Therefore

$$P_{\bar{a}}^*(\omega) \geq \inf_{\mu': \mu'_{\bar{a}} - \mu'_1 \geq \Delta_0} \sum_a \omega_a \frac{(\mu_a - \mu'_a)^2}{2\sigma^2} = \inf_{\mu': \mu'_{\bar{a}} - \mu'_1 \geq \Delta_0} \omega_1 \frac{(\mu_1 - \mu'_1)^2}{2\sigma^2} + \omega_{\bar{a}} \frac{(\mu_{\bar{a}} - \mu'_{\bar{a}})^2}{2\sigma^2}.$$

505 Solving the KKT conditions we find the equivalent conditions $\mu'_{\bar{a}} = \mu'_1 + \Delta_0$ and

$$\omega_1(\mu_1 - \mu'_1) + \omega_{\bar{a}}(\mu_{\bar{a}} - \mu'_1 - \Delta_0) = 0 \Rightarrow \mu'_1 = \frac{\omega_1\mu_1 + \omega_{\bar{a}}\mu_{\bar{a}} - \omega_{\bar{a}}\Delta_0}{\omega_1 + \omega_{\bar{a}}}.$$

506 Therefore

$$\mu'_{\bar{a}} = \frac{\omega_1\mu_1 + \omega_{\bar{a}}\mu_{\bar{a}} - \omega_{\bar{a}}\Delta_0}{\omega_1 + \omega_{\bar{a}}} + \Delta_0 = \frac{\omega_1\mu_1 + \omega_{\bar{a}}\mu_{\bar{a}} + \omega_1\Delta_0}{\omega_1 + \omega_{\bar{a}}}.$$

507 Plugging these solutions back in the value of the problem, we obtain

$$\begin{aligned} P_{\bar{a}}^*(\omega) &\geq \frac{\omega_1\omega_{\bar{a}}^2}{(\omega_1 + \omega_{\bar{a}})^2} \frac{(\mu_1 - \mu_{\bar{a}} + \Delta_0)^2}{2\sigma^2} + \frac{\omega_{\bar{a}}\omega_1^2}{(\omega_1 + \omega_{\bar{a}})^2} \frac{(\mu_{\bar{a}} - \mu_1 - \Delta_0)^2}{2\sigma^2}, \\ &= \frac{\omega_1\omega_{\bar{a}}}{\omega_1 + \omega_{\bar{a}}} \frac{(\mu_1 - \mu_{\bar{a}} + \Delta_0)^2}{2\sigma^2}, \\ &= \frac{\omega_1\omega_{\bar{a}}}{\omega_1 + \omega_{\bar{a}}} \frac{(\Delta_{\bar{a}} + \Delta_0)^2}{2\sigma^2}. \end{aligned}$$

508 Let $\theta_a = \Delta_a + \Delta_0$, with $\theta_1 = \Delta_0$. We plug in a feasible solution $\omega_a = \frac{1/\theta_a^2}{\sum_b 1/\theta_b^2}$, yielding

$$\begin{aligned} (T^*(\mu))^{-1} &= \sup_{\omega \in \Delta(K)} \min_{\bar{a} \neq 1} P_{\bar{a}}^*(\omega) \geq \min_{\bar{a} \neq 1} \frac{1/(\theta_1\theta_{\bar{a}})^2}{\sum_b 1/\theta_b^2(1/\theta_1^2 + 1/\theta_{\bar{a}}^2)} \frac{\theta_{\bar{a}}^2}{2\sigma^2}, \\ &= \min_{\bar{a} \neq 1} \frac{1}{\sum_b 1/\theta_b^2(1 + \theta_1^2/\theta_{\bar{a}}^2)} \frac{1}{2\sigma^2}, \\ &= \frac{1}{2\sigma^2 \sum_b 1/\theta_b^2} \min_{\bar{a} \neq 1} \frac{1}{1 + \theta_1^2/\theta_{\bar{a}}^2}, \\ &\geq \frac{1}{2\sigma^2 \sum_b 1/\theta_b^2} \frac{1}{1 + \theta_1^2/\Delta_0^2}, \\ &= \frac{1}{4\sigma^2 \sum_b 1/\theta_b^2}. \end{aligned}$$

509

□

510 A.2 Sample Complexity Lower Bound for the Magic Action MAB Problem

511 We now consider a special class of models that embeds information about the optimal arm in the
 512 mean reward of some of the arms. Let $\phi : \mathbb{R} \rightarrow \mathbb{R}$ be a strictly decreasing function over $\{2, \dots, K\}$ ².

513 Particularly, we make the following assumptions:

514 1. We consider mean rewards μ satisfying $\mu_1 = \phi(\arg \max_{a \neq 1} \mu_a)$, and $\mu^* = \max_a \mu_a >$
 515 $\phi(2)$. Arm 1 is called "magic action", and with this assumption we are guaranteed that the
 516 magic arm is not optimal, since

$$\mu_1 \frac{1}{\max_a \mu_a} = \phi(\arg \max_{a \neq 1} \mu_a) \frac{1}{\max_a \mu_a} \leq \phi(2) \frac{1}{\max_a \mu_a} < 1 \Rightarrow \max_a \mu_a > \mu_1.$$

517 2. The rewards are normally distributed, with a fixed known standard deviation σ_1 for the
 518 magic arm, and fixed standard deviation σ for all the other arms.

²One could also consider strictly increasing functions.

519 Hence, define the set of models

$$\mathcal{S} = \left\{ \mu \in \mathbb{R}^K : \mu_1 = \phi(\arg \max_{a \neq 1} \mu_a), \max_a \mu_a > \phi(2) \right\},$$

520 and the set of alternative models

$$\text{Alt}(\mu) = \left\{ \mu' \in \mathcal{S} : \arg \max_a \mu'_a \neq a^* \right\},$$

521 where $a^* = \arg \max_a \mu_a$.

522 Then, for any δ -correct algorithm, guaranteeing that at some stopping time τ the estimated optimal
523 arm \hat{a}_τ is δ -correct, i.e., $\mathbb{P}_\mu(\hat{a}_\tau \neq a^*) \leq \delta$, we have the following result.

524 **Theorem A.2.** *For any δ -correct algorithm, the sample complexity lower bound on the magic action
525 problem is*

$$\mathbb{E}_\mu[\tau] \geq T^*(\mu) \text{kl}(1 - \delta, \delta), \quad (6)$$

526 where $\text{kl}(x, y) = x \log(x/y) + (1 - x) \log((1 - x)/(1 - y))$ and $T^*(\mu)$ is the characteristic time of
527 μ , defined as

$$(T^*(\mu))^{-1} = \max_{\omega \in \Delta(K)} \min_{a \neq 1, a^*} \omega_1 \frac{(\phi(a^*) - \phi(a))^2}{2\sigma_1^2} + \sum_{b \in \mathcal{K}_a(\omega)} \omega_b \frac{(\mu_b - m(\omega; \mathcal{K}_a(\omega)))^2}{2\sigma^2}, \quad (7)$$

528 where $m(\omega; \mathcal{C}) = \frac{\sum_{a \in \mathcal{C}} \omega_a \mu_a}{\sum_{a \in \mathcal{C}} \omega_a}$ and the set $\mathcal{K}_a(\omega)$ is defined as

$$\mathcal{K}_a(\omega) = \{a\} \cup \{b \in \{2, \dots, K\} : \mu_b \geq m(\omega; \mathcal{C}_b \cup \{a\}) \text{ and } \mu_b \geq \phi(2)\}.$$

529 with $\mathcal{C}_x = \{b \in \{2, \dots, K\} : \mu_b \geq \mu_x\}$ for $x \in [K]$.

530 *Proof. Step 1: Log-likelihood ratio.* The initial part of the proof is rather standard, and follows the
531 same argument used in the Best Arm Identification and Best Policy Identification literature [Garivier](#)
532 [and Kaufmann \[2016\]](#), [Russo and Vannella \[2025\]](#).

533 Take the expected log-likelihood ratio between μ and $\mu' \in \text{Alt}(\mu)$ of the data observed up to τ
534 $\Lambda_\tau = \log \frac{d\mathbb{P}_\mu(A_1, R_1, \dots, A_\tau, R_\tau)}{d\mathbb{P}_{\mu'}(A_1, R_1, \dots, A_\tau, R_\tau)}$, where A_t is the action taken in round t , and R_t is the reward observed
535 upon selecting A_t . Then, we can write

$$\Lambda_t = \sum_a \sum_{n=1}^t \mathbf{1}_{\{A_n=a\}} \log \frac{f_a(R_n)}{f'_a(R_n)}$$

536 where f_a, f'_a are, respectively, the reward density for action a in the two models μ, μ' with respect to
537 the Lebesgue measure. Letting $N_a(t)$ denote the number of times action a has been selected up to
538 round t , by an application of Wald's lemma the expected log-likelihood ratio can be shown to be

$$\mathbb{E}_\mu[\Lambda_\tau] = \sum_a \mathbb{E}_\mu[N_a(\tau)] \text{KL}(\mu_a, \mu'_a)$$

539 where $\text{KL}(\mu_a, \mu'_a)$ is the KL divergence between two Gaussian distributions $\mathcal{N}(\mu_a, \sigma)$ and $\mathcal{N}(\mu'_a, \sigma)$
540 (note that we have σ_1 instead of σ for $a = 1$).

541 We also know from the information processing inequality [Kaufmann et al. \[2016\]](#) that $\mathbb{E}_\mu[\Lambda_\tau] \geq$
542 $\sup_{\mathcal{E} \in \mathcal{M}_\tau} \text{kl}(\mathbb{P}_\mu(\mathcal{E}), \mathbb{P}_{\mu'}(\mathcal{E}))$, where $\mathcal{M}_t = \sigma(A_1, R_1, \dots, A_t, R_t)$. We use the fact that the algo-
543 rithm is δ -correct: by choosing $\mathcal{E} = \{\hat{a}_\tau = a^*\}$ we obtain that $\mathbb{E}_\mu[\Lambda_\tau] \geq \text{kl}(1 - \delta, \delta)$, since
544 $\mathbb{P}_\mu(\mathcal{E}) \geq 1 - \delta$ and $\mathbb{P}_{\mu'}(\mathcal{E}) = 1 - \mathbb{P}_{\mu'}(\hat{a}_\tau \neq a^*) \leq 1 - \mathbb{P}_{\mu'}(\hat{a}_\tau = \arg \max_a \mu'_a) \leq \delta$ (we also used
545 the monotonicity properties of the Bernoulli KL divergence). Hence

$$\sum_a \mathbb{E}_\mu[N_a(\tau)] \text{KL}(\mu_a, \mu'_a) \geq \text{kl}(1 - \delta, \delta).$$

546 Letting $\omega_a = \mathbb{E}_\mu[N_a(\tau)] / \mathbb{E}_\mu[\tau]$, we have that

$$\mathbb{E}_\mu[\tau] \sum_a \omega_a \text{KL}(\mu_a, \mu'_a) \geq \text{kl}(1 - \delta, \delta).$$

547 Lastly, optimizing over $\mu' \in \text{Alt}(\mu)$ and $\omega \in \Delta(K)$ yields the bound:

$$\mathbb{E}_\mu[\tau] \geq T^*(\mu) \text{kl}(1 - \delta, \delta),$$

548 where $T^*(\mu)$ is defined as

$$(T^*(\mu))^{-1} = \sup_{\omega \in \Delta(K)} \inf_{\mu' \in \text{Alt}(\mu)} \sum_a \omega_a \text{KL}(\mu_a, \mu'_a).$$

549 **Step 2: Optimization over the set of alternative models.** We now face the problem of optimizing
 550 over the set of alternative models. First, we observe that $\mathcal{S} = \cup_{a \neq a^*} \{\mu : \mu_1 = \phi(a), \mu_a > \phi(2)\}$.
 551 Therefore, we can write

$$\text{Alt}(\mu) = \cup_{a \notin \{1, a^*\}} \{\mu' : \mu'_1 = \phi(a), \mu'_a > \max(\phi(2), \mu'_b) \forall b \neq a\}.$$

552 Hence, for a fixed $a \notin \{1, a^*\}$, the inner infimum becomes

$$\begin{aligned} \inf_{\mu' \in \mathbb{R}^K} \quad & \omega_1 \frac{(\phi(a^*) - \phi(a))^2}{2\sigma_1^2} + \sum_{a \neq 1} \omega_a \frac{(\mu_a - \mu'_a)^2}{2\sigma^2} \\ \text{s.t.} \quad & \mu'_a \geq \max(\phi(2), \mu'_b) \quad \forall b, \\ & \mu'_1 = \phi(a). \end{aligned} \tag{8}$$

553 To solve it, we construct the following Lagrangian

$$\ell(\mu', \theta) = \omega_1 \frac{(\phi(a^*) - \phi(a))^2}{2\sigma_1^2} + \sum_{b \neq 1} \omega_b \frac{(\mu_b - \mu'_b)^2}{2\sigma^2} + \sum_b \theta_b (\max(\phi(2), \mu'_b) - \mu'_a),$$

554 where $\theta \in \mathbb{R}_+^K$ is the multiplier vector. From the KKT conditions we already know that $\theta_1 = 0, \theta_a = 0$
 555 and $\theta_b = 0$ if $\mu'_b \leq \phi(2)$, with $b \in \{2, \dots, K\}$. In particular, we also know that either we have
 556 $\mu'_b = \mu'_a$ or $\mu'_b = \mu_b$. Therefore, for $\mu_b \leq \phi(2)$ the solution is $\mu'_b = \mu_b$, while for $\mu_b > \phi(2)$ the
 557 solution depends also on ω .

558 To fix the ideas, let \mathcal{K} be the set of arms for which $\mu'_b = \mu'_a$ at the optimal solution. Such set must
 559 necessarily include arm a . Then, note that

$$\frac{\partial \ell}{\partial \mu'_a} = \omega_a \frac{\mu'_a - \mu_a}{\sigma^2} - \sum_{b \in [K]} \theta_b = 0.$$

560 and

$$\frac{\partial \ell}{\partial \mu'_b} = \omega_b \frac{\mu'_b - \mu_b}{\sigma^2} + \theta_b = 0 \quad \text{for } b \neq (1, a).$$

561 Then, using the observations derived above, we conclude that

$$\mu'_a = \frac{\sum_{b \in \mathcal{K}} \omega_b \mu_b}{\sum_{b \in \mathcal{K}} \omega_b},$$

562 with $\mu'_b = \mu'_a$ if $b \in \mathcal{K}$, and $\mu'_b = \mu_b$ otherwise. However, how do we compute such set \mathcal{K} ?

563 First, \mathcal{K} includes arm a . However, in general we have $\mathcal{K} \neq \{a\}$: if that were not true we would have
 564 $\mu'_a = \mu_a$ and $\mu'_b = \mu_b$ for the other arms – but if any μ_b is greater than μ_a , then a is not optimal,
 565 which is a contradiction. Therefore, also arm a^* is included in \mathcal{K} , since any convex combination of
 566 $\{\mu_a\}$ is necessarily smaller than μ_{a^*} . We apply this argument repeatedly for every arm b to obtain \mathcal{K} .

567 Hence, for some set $\mathcal{C} \subseteq [K]$ define the average reward

$$m(\omega; \mathcal{C}) = \frac{\sum_{a \in \mathcal{C}} \omega_a \mu_a}{\sum_{a \in \mathcal{C}} \omega_a},$$

568 and the set $\mathcal{C}_x = \{a\} \cup \{b \in \{2, \dots, K\} : \mu_b \geq \mu_x\}$ for $x \in [K]$. Then,

$$\mathcal{K} := \mathcal{K}(\omega) = \{a\} \cup \{b \in \{2, \dots, K\} : \mu_b \geq m(\omega; \mathcal{C}_b) \text{ and } \mu_b \geq \phi(2)\}.$$

569 In other words, \mathcal{K} is the set of *confusing arms* for which the mean reward in the alternative model
 570 changes. An arm b is *confusing* if the average reward m , taking into account b , is smaller than μ_b . If
 571 this holds for b , then it must also hold all the arms b' such that $\mu_{b'} \geq \mu_b$. \square

572 Finally, to get a better intuition of the main result, we can look at the 3-arms case: it is optimal to
 573 only sample the magic arm iff $|\phi(a^*) - \phi(a)| > \frac{\sigma_1(\mu_{a^*} - \mu_a)}{2\sigma}$.

574 **Lemma A.3.** *With $K = 3$ we have that $\omega_1 = 1$ if and only if*

$$|\phi(a^*) - \phi(a)| > \frac{\sigma_1(\mu_{a^*} - \mu_a)}{2\sigma},$$

575 and $\omega_1 = 0$ if the reverse inequality holds.

576 *Proof.* With 3 arms, from the proof of the theorem we know that $\mathcal{K}_a(\omega) = \{a, a^*\}$ for all ω . Letting
 577 $m(\omega) = \frac{\omega_a \mu_a + \omega_{a^*} \mu_{a^*}}{\omega_a + \omega_{a^*}}$, we obtain

$$(T^*(\mu))^{-1} = \max_{\omega \in \Delta(3)} \omega_1 \frac{(\phi(a^*) - \phi(a))^2}{2\sigma_1^2} + \frac{\omega_a(\mu_a - m(\omega))^2 + \omega_{a^*}(\mu_{a^*} - m(\omega))^2}{2\sigma^2}.$$

578 Clearly the solution is $\omega_1 = 1$ as long as

$$\frac{(\phi(a^*) - \phi(a))^2}{2\sigma_1^2} > \max_{\omega: \omega_a + \omega_{a^*} = 1} \frac{\omega_a(\mu_a - m(\omega))^2 + \omega_{a^*}(\mu_{a^*} - m(\omega))^2}{2\sigma^2}.$$

579 To see why this is the case, let $f_1 = \frac{(\phi(a^*) - \phi(a))^2}{2\sigma_1^2}$, $f_2(\omega_a, \omega_{a^*}) = \frac{\omega_a(\mu_a - m(\omega))^2}{2\sigma^2}$ and $f_3(\omega_a, \omega_{a^*}) =$
 580 $\frac{\omega_{a^*}(\mu_{a^*} - m(\omega))^2}{2\sigma^2}$. Then, we can write

$$\omega_1 f_1 + \omega_a f_2(\omega_a, \omega_{a^*}) + \omega_{a^*} f_3(\omega_a, \omega_{a^*}) = \omega_1 f_1 + (1 - \omega_1) \left[\frac{\omega_a f_2}{1 - \omega_1} + \frac{\omega_{a^*} f_3}{1 - \omega_1} \right].$$

581 Being a convex combination, this last term can be upper bounded as

$$\omega_1 f_1 + \omega_a f_2(\omega_a, \omega_{a^*}) + \omega_{a^*} f_3(\omega_a, \omega_{a^*}) \leq \max \left(f_1, \frac{\omega_a f_2}{1 - \omega_1} + \frac{\omega_{a^*} f_3}{1 - \omega_1} \right).$$

582 Now, note that also the term inside the bracket is a convex combination. Therefore, let $\omega_a = (1 - \omega_1)\alpha$
 583 and $\omega_{a^*} = (1 - \omega_1)(1 - \alpha)$ for some $\alpha \in [0, 1]$. We have that

$$m(\omega) = \frac{(1 - \omega_1)\alpha \mu_a + (1 - \omega_1)(1 - \alpha)\mu_{a^*}}{1 - \omega_1} = \alpha \mu_a + (1 - \alpha)\mu_{a^*}.$$

584 Hence, we obtain that

$$\begin{aligned} \frac{\omega_a(\mu_a - m(\omega))^2 + \omega_{a^*}(\mu_{a^*} - m(\omega))^2}{2(1 - \omega_1)\sigma^2} &= \frac{\omega_a f_2 + \omega_{a^*} f_3}{1 - \omega_1}, \\ &= \frac{\alpha(1 - \alpha)^2(\mu_a - \mu_{a^*})^2 + (1 - \alpha)\alpha^2(\mu_{a^*} - \mu_a)^2}{2\sigma^2}, \\ &= \alpha(1 - \alpha) \frac{(1 - \alpha)(\mu_a - \mu_{a^*})^2 + \alpha(\mu_{a^*} - \mu_a)^2}{2\sigma^2}, \\ &= \alpha(1 - \alpha) \frac{(\mu_a - \mu_{a^*})^2}{2\sigma^2}. \end{aligned}$$

585 Since this last term is maximized for $\alpha = 1/2$, we obtain

$$\omega_1 f_1 + \omega_a f_2(\omega_a, \omega_{a^*}) + \omega_{a^*} f_3(\omega_a, \omega_{a^*}) \leq \max \left(f_1, \frac{(\mu_a - \mu_{a^*})^2}{8\sigma^2} \right).$$

586 Since f_1 is attained for $\omega_1 = 1$, we have that as long as $f_1 > \frac{(\mu_a - \mu_{a^*})^2}{8\sigma^2}$, then the solution is $\omega_1 = 1$.

587 On the other hand, if $\frac{(\mu_a - \mu_{a^*})^2}{8\sigma^2} > f_1$, then we can set $\omega_a = (1 - \omega_1)/2$ and $\omega_{a^*} = (1 - \omega_1)/2$,
 588 leading to

$$\omega_1 f_1 + \omega_a f_2(\omega_a, \omega_{a^*}) + \omega_{a^*} f_3(\omega_a, \omega_{a^*}) = \omega_1 f_1 + (1 - \omega_1) \frac{(\mu_a - \mu_{a^*})^2}{8\sigma^2},$$

589 which is maximized at $\omega_1 = 0$. □

590 **A.3 Sample Complexity Bound for the Multiple Magic Actions MAB Problem**

591 We now extend our analysis to the case where multiple magic actions can be present in the environment.
 592 In contrast to the single magic action setting, here a *chain* of magic actions sequentially reveals
 593 information about the location of the optimal action. Without loss of generality, assume that the first
 594 n arms (with indices $1, \dots, n$) are the magic actions, and the remaining $K - n$ arms are non-magic.
 595 The chain structure is such that pulling magic arm j (with $1 \leq j < n$) yields information about only
 596 the location of the next magic arm $j + 1$, while pulling the final magic action (arm n) reveals the
 597 identity of the optimal action. As before, we assume that the magic actions are informational only
 598 and are never optimal.

599 To formalize the model, let $\phi : \{1, \dots, n\} \rightarrow \mathbb{R}$ be a strictly decreasing function. We assume that the
 600 magic actions have fixed means given by

$$\mu_j = \begin{cases} \phi(j + 1), & \text{if } j = 1, \dots, n - 1, \\ \phi\left(\arg \max_{a \notin \{1, \dots, n\}} \mu_a\right), & \text{if } j = n. \end{cases}$$

601 and that the non-magic arms satisfy

$$\mu^* = \max_{a \notin \{1, \dots, n\}} \mu_a > \phi(n).$$

602 Thus, the optimal arm lies among the non-magic actions. Considering the noiseless case where the
 603 rewards of all actions are fixed and the case where we can identify if an action is magic once revealed,
 604 we have the following result.

605 **Theorem A.4.** *Consider noiseless magic bandit problem with K arms and n magic actions. The*
 606 *optimal sample complexity is upper bounded as*

$$\inf_{\text{Alg}} \mathbb{E}_{\text{Alg}}[\tau] \leq \min \left(n, \sum_{j=1}^{K-n} \left(\prod_{i=j+1}^{K-n} \frac{i}{n-1+i} \right) \left(1 + \frac{n-1}{n-1+j} \min \left(\frac{n-2}{2}, \frac{j(n-1+j)}{j+1} \right) \right) \right).$$

607 *Proof.* In the proof we derive a sample complexity bound for a policy based on some insights. We
 608 use the assumption that upon observing a reward from a magic arm, the learner can almost surely
 609 identify that the pulled arm is a magic arm.

610 Let us define the state (m, r, l) , where m denotes the number of remaining unrevealed magic actions
 611 ($m_0 = n - 1$), r denotes the number of remaining unrevealed non-magic actions ($r_0 = K - n$), and
 612 l is the binary indicator with value 1 if we have revealed any hidden magic action and 0 otherwise.

613 Before any observation the learner has no information about which $n - 1$ indices among $\{2, \dots, K\}$
 614 form the chain of intermediate magic arms. Hence, one can argue that at the first time-step is optimal
 615 to sample uniformly at random an action in $\{2, \dots, K\}$.

616 Upon observing a magic action, and thus we are in state $(m, r, 1)$, we consider the following candidate
 617 policies: (1) start from the revealed action and follow the chain, or (2) keep sampling unrevealed
 618 actions uniformly at random until all non-magic actions are revealed. As previously discussed,
 619 starting the chain from the initial magic action would be suboptimal and we do not consider it.

620 Upon drawing a hidden magic arm, let its chain index be $j \in \{2, \dots, n\}$ (which is uniformly
 621 distributed). The remaining cost to complete the chain is $n - j$, and hence its expected value is

$$\mathbb{E}[n - j] = \frac{n - 2}{2}.$$

622 Therefore, the total expected cost for strategy (1) is

$$T_1 = \frac{n - 2}{2}.$$

623 We can additionally compute the expected cost for strategy (2) as follows: if the last non-magic action
 624 is revealed at step i , then among the first $i - 1$ draws there are exactly $r - 1$ non-magic arms. Since

625 there are $\binom{m+r}{r}$ ways to place all r non-magic arms $m+r$ slots, we have

$$\begin{aligned}
T_2 &= \mathbb{E}[\text{Draws until all non-magic revealed}] \\
&= \sum_{i=r}^{m+r} i \cdot \mathbb{P}[\text{Last non-magic revealed at step } i] \\
&= \sum_{i=r}^{m+r} i \cdot \frac{\binom{i-1}{r-1}}{\binom{m+r}{r}} \\
&= \frac{r! \cdot m!}{(m+r)!} \sum_{i=r}^{m+r} i \binom{i-1}{r-1} \\
&= \frac{r! \cdot m!}{(m+r)!} \sum_{i=r}^{m+r} \frac{i!}{(r-1)!(i-r)!} \\
&= \frac{r! \cdot m!}{(m+r)!} \sum_{i=r}^{m+r} r \binom{i}{r} \\
&= \frac{r \cdot r! \cdot m!}{(m+r)!} \binom{m+r+1}{r+1} \\
&= \frac{r \cdot r! \cdot m!}{(m+r)!} \cdot \frac{(m+r+1) \cdot (m+r)!}{(r+1) \cdot r! \cdot m!} \\
&= \frac{r(m+r+1)}{r+1}
\end{aligned}$$

626 Finally, we define a policy in $(m, r, 1)$ as the one choosing between strategy 1 and strategy 2,
627 depending on which one achieves the minimum cost. Hence, the complexity of this policy is

$$V(m, r, 1) = \min \left(\frac{n-2}{2}, \frac{r(m+r+1)}{r+1} \right).$$

628 Now, before finding a magic arm, consider a policy that uniformly samples between the non-revealed
629 arms. Therefore, in $(m, r, 0)$ we can achieve a complexity of $1 + \frac{m}{m+r}V(m-1, r, 1) + \frac{r}{m+r}V(m, r-1, 0)$.
630 Since we can always achieve a sample complexity of n , we can find a policy with the following
631 complexity:

$$\begin{aligned}
V(m, r, 0) &= \min \left(n, 1 + \frac{m}{m+r}V(m-1, r, 1) + \frac{r}{m+r}V(m, r-1, 0) \right) \\
&= \min \left(n, 1 + \frac{m}{m+r} \min \left(\frac{n-2}{2}, \frac{r(m+r)}{r+1} \right) + \frac{r}{m+r}V(m, r-1, 0) \right)
\end{aligned}$$

632 Given we always start with $n-1$ hidden magic actions we can define a recursion in terms of just the
633 variable r as follows:

$$V(r) = 1 + \frac{n-1}{n-1+r}T(r) + \frac{r}{n-1+r}V(r-1),$$

634 where $T(r) = \min \left(\frac{n-2}{2}, \frac{r(n-1+r)}{r+1} \right)$. Letting $A(r) = \frac{r}{n-1+r}$ and $B(r) = 1 + \frac{n-1}{n-1+r}T(r)$, we can
635 write

$$V(r) = B(r) + A(r)V(r-1),$$

636 Clearly $V(0) = 0$ since if all non-magic actions are revealed, then we know the optimal action
 637 deterministically. Unrolling the recursion we get

$$\begin{aligned} V(1) &= B(1), \\ V(2) &= B(2) + A(2)B(1), \\ V(3) &= B(3) + A(3)B(2) + A(3)A(2)B(1), \\ &\dots \end{aligned}$$

$$V(r) = \sum_{j=1}^r \left(\prod_{i=j+1}^r A(i) \right) B(j).$$

638 Substituting back in our expression, we get

$$V(r) = \sum_{j=1}^r \left(\prod_{i=j+1}^r \frac{i}{n-1+i} \right) \left(1 + \frac{n-1}{n-1+j} T(j) \right).$$

639 Thus starting at $r = K - n$ we get the following expression:

$$\min \left(n, \sum_{j=1}^{K-n} \left(\prod_{i=j+1}^{K-n} \frac{i}{n-1+i} \right) \left(1 + \frac{n-1}{n-1+j} \min \left(\frac{n-2}{2}, \frac{j(n-1+j)}{j+1} \right) \right) \right),$$

640 which is also an upper bound on the optimal sample complexity.

641

□

642 To get a better intuition of the result, we also have the following corollary, which shows that we
 643 should expect a scaling linear in n for small values of n (for large values the complexity tends instead
 644 to "flatten").

645 **Corollary A.5.** *Let T be the scaling in theorem A.4. We have that*

$$\min(n, (K - n)/2) \lesssim T \lesssim C \min(n, K/2).$$

646 *Proof.* First, observe the scaling

$$\left(1 + \frac{n-1}{n-1+j} \min \left(\frac{n-2}{2}, \frac{j(n-1+j)}{j+1} \right) \right) = O(n/2).$$

647 At this point, note that

$$\prod_{i=j+1}^{K-n} \frac{i}{n-1+i} = \prod_{i=j+1}^{K-n} \left(1 + \frac{n-1}{i} \right)^{-1}.$$

648 Using that $\frac{x}{1+x} \leq \log(1+x) \leq x$, we have

$$\log \prod_{i=j+1}^{K-n} \frac{i}{n-1+i} = \sum_{i=j+1}^{K-n} -\log \left(1 + \frac{n-1}{i} \right) \geq -(n-1) \sum_{i=j+1}^{K-n} \frac{1}{i}.$$

649 and

$$\log \prod_{i=j+1}^{K-n} \frac{i}{n-1+i} = \sum_{i=j+1}^{K-n} -\log \left(1 + \frac{n-1}{i} \right) \leq -(n-1) \sum_{i=j+1}^{K-n} \frac{1}{n-1+i}.$$

650 Define $H_n = \sum_{i=1}^n 1/i$ to be the n -th Harmonic number, we also have

$$\sum_{i=j+1}^{K-n} \frac{1}{i} = H_{K-n} - H_j.$$

651 Therefore

$$-(n-1)(H_{K-n} - H_j) \leq \log \prod_{i=j+1}^{K-n} \frac{i}{n-1+i} \leq -(n-1)(H_{K-1} - H_{n+j-1})$$

652 Using that $H_\ell \sim \log(\ell) + \gamma + O(1/\ell)$, where γ is the Euler–Mascheroni constant, we get

$$\left(\frac{j}{K-n}\right)^{n-1} \lesssim \prod_{i=j+1}^{K-n} \frac{i}{n-1+i} \lesssim \left(\frac{n+j-1}{K-1}\right)^{n-1}.$$

653 Therefore, we can bound $\sum_{j=1}^{K-n} \left(\frac{n+j-1}{K-1}\right)^{n-1}$ using an integral bound

$$\sum_{j=1}^{K-n} \left(\frac{n+j-1}{K-1}\right)^{n-1} \leq \int_0^{K-n} \left(\frac{n+x}{K-1}\right)^{n-1} dx \leq \frac{e(K-1)}{n}.$$

654 From which follows that the original expression can be upper bounded by an expression scaling as
 655 $O(\min(n, (K-1)/2))$.

656 Similarly, using that $\sum_{j=1}^{K-n} \left(\frac{j}{K-n}\right)^{n-1} \geq (K-n)/n$, we have that the lower bound scales as
 657 $\min(n, (K-n)/2)$. □

658 B Algorithms

659 In this section we present some of the algorithms more in detail. These includes: **ICPE**, TaS, *I-DPT*
660 and *I-IDS*.

661 **MDP Formulation for ICPE.** Recall that in **ICPE** we treat trajectories of data $\mathcal{D}_t = (x_1, a_1, \dots, x_t)$
662 as sequences to be given as input to sequential models, such as Transformers. We treat trajectories
663 as states of an MDP M . An environment M can be then modeled as an MDP, which is a sequential
664 model characterized by a tuple $M = (\mathcal{S}, \mathcal{A}, P', r, H_M^*, \rho)$, where \mathcal{S} is the state space, \mathcal{A} the action
665 space, $P' : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ is the transition function, $r : \mathcal{S} \rightarrow [0, 1]$ defines the reward function (to
666 be defined later), $H^* \in \mathcal{H}$ is the true hypothesis in M and ρ is the initial state distribution.

667 We define the state at time-step t as $s_t = (\mathcal{D}_t, \varnothing_{t:N})$, with $\varnothing_{t:N}$ indicating a null sequence of tokens
668 for the remaining steps up to some pre-defined horizon N , with $s_1 = (x_1, \varnothing_{1:N})$.

669 To be more precise, letting $(s_t^\varnothing, a_t^\varnothing)$ denote, respectively, the null elements in the state and action at
670 time-step t , we have $\varnothing_{t:t+k} = \{s_t^\varnothing, a_{t+1}^\varnothing, s_{t+1}^\varnothing, \dots, a_{t+k-1}^\varnothing, s_{t+k}^\varnothing\}$.

671 The limit N is a practical upper bound on the horizon that limits the dimensionality of the state,
672 which is introduced for implementing the algorithm. The action space remains \mathcal{A} , and the transition
673 dynamics P' are induced by (ρ, P) .

674 B.1 ICPE with Fixed Confidence

675 Recall that $\mathcal{D}_t = (x_1, a_1, \dots, x_{t-1}, a_{t-1}, x_t)$ and $\hat{H}_\tau \sim I(\cdot | D_\tau)$. In the fixed confidence setting,
676 problems terminate at some random point in time τ , chosen by the learner, or when the maximum
677 horizon N is reached. We model this by giving π_t an additional stopping action a_{stop} such that
678 $\pi_t : \mathcal{D}_t \rightarrow \mathcal{A} \cup \{a_{\text{stop}}\}$ so that the data collection processes terminates at the stopping-time
679 $\tau = \min(N, t_{\text{stop}})$, with $t_{\text{stop}} := \inf\{t \in \mathbb{N} : a_t = a_{\text{stop}}\}$.

680 Optimizing the dual formulation

$$\min_{\lambda \geq 0} \max_{I, \pi} V_\lambda(\pi, I)$$

681 can be viewed as a multi-timescale stochastic optimization problem: the slowest timescale updates
682 the variable λ , an intermediate timescale optimizes over I , and the fastest refines the policy π .

Algorithm 2 **ICPE** (In-Context Pure Exploration) - Fixed Confidence

- 1: **Input:** Tasks distribution $\mathcal{P}(\mathcal{M})$; confidence δ ; learning rates α, β ; initial λ and hyper-parameters T, N, η .
- 2: Initialize buffer \mathcal{B} , networks Q_θ, I_ϕ and set $\bar{\theta} \leftarrow \theta, \bar{\phi} \leftarrow \phi$.
- 3: **while** Training is not over **do**
- 4: Sample environment $M \sim \mathcal{P}(\mathcal{M})$ with hypothesis H^* , observe $s_1 \sim \rho$ and set $t \leftarrow 1$.
- 5: **for** $t = 1, \dots, N - 1$ **do**
- 6: Execute action $a_t = \arg \max_a Q_\theta(s_t, a)$ in M and observe next state s_{t+1} .
- 7: Add experience $z_t = (s_t, a_t, s_{t+1}, d_t = \mathbf{1}\{s_{t+1} \text{ is terminal}\}, H^*)$ to \mathcal{B} .
- 8: If $a_t = a_{\text{stop}}$, break the loop.
- 9: **end for**
- 10: Update variable λ according to

$$\lambda \leftarrow \max(0, \lambda - \beta (I_\phi(H^* | s_{\tau+1}) - 1 + \delta). \quad (9)$$

- 11: Sample batches $B, B' \sim \mathcal{B}$ and update θ, ϕ as

$$\theta \leftarrow \theta - \alpha \nabla_\theta \frac{1}{|B|} \sum_{z \in B} \left[\mathbf{1}_{\{a \neq a_{\text{stop}}\}} (y_\lambda(z) - Q_\theta(s, a))^2 + (r_\lambda(z_{\text{stop}}) - Q_\theta(s, a_{\text{stop}}))^2 \right], \quad (10)$$

$$\phi \leftarrow \phi + \alpha \nabla_\phi \frac{1}{|B'|} \sum_{z \in B'} [\log(I_\phi(H^* | s))]. \quad (11)$$

- 12: Update $\bar{\theta} \leftarrow (1 - \eta)\bar{\theta} + \eta\theta$ and every T steps set $\bar{\phi} \leftarrow \phi$.
 - 13: **end while**
-

683 **MDP Formulation.** We can use the MDP formalism to define an RL problem: we define a reward
 684 r that penalizes the agent at all time-steps, that is $r_t = -1$, while at the stopping-time we have
 685 $r_\tau = -1 + \lambda \mathbb{E}_{H \sim I(\cdot|s_\tau)}[h(H; M)]$. Hence, a trajectory’s return can be written as

$$G_\tau = \sum_{t=1}^{\tau} r_t = -\tau + 1 + \underbrace{r(s_\tau, a_\tau)}_{r_\tau} = -\tau + \lambda I(H^*|s_\tau).$$

686 Accordingly, one can define the Q -value of (π, I, λ) in a state-action pair (s, a) at the t -th step as

$$687 \quad Q_\lambda^{\pi, I}(s, a) = \mathbb{E}_{M \sim \mathcal{P}(\mathcal{M})}^{\pi} \left[\sum_{n=t}^{\tau} r_n \mid s_t = s, a_t = a \right], \text{ with } a_n \sim \pi_n(\cdot|s_n)$$

688 **Optimization over ϕ .** We treat each optimization separately, employing a descent-ascent scheme.
 689 The distribution I is modeled using a sequential architecture parameterized by ϕ , denoted by I_ϕ .
 690 Fixing (π, λ) , the inner maximization in eq. (1) corresponds to

$$\max_{\phi} \mathbb{E}_{M \sim \mathcal{P}(\mathcal{M})}^{\pi} [h(\hat{H}_\tau; M)], \quad \text{with } \hat{H}_\tau \sim I_\phi(\cdot|s_\tau).$$

691 We train ϕ via cross-entropy loss:

$$-\sum_{H'} h(H'; M) \log I_\phi(H'|s_\tau) = -\log I_\phi(H^*|s_\tau),$$

692 averaged across environments. Alternatively, a MAP estimator may be used with the same loss.

693 **Optimization over π .** The policy π is defined as the greedy policy with respect to learned Q -values.
 694 Therefore, standard RL techniques can learn the Q -function that maximizes the value in eq. (1)
 695 given (λ, I) . Denoting this function by Q_θ , it is parameterized using a sequential architecture with
 696 parameters θ .

697 We train Q_θ using DQN Mnih et al. [2015], Van Hasselt et al. [2016], employing a replay buffer
 698 \mathcal{B} and a target network $Q_{\bar{\theta}}$ parameterized by $\bar{\theta}$. To maintain timescale separation, we introduce an
 699 additional inference target network $I_{\bar{\phi}}$, parameterized by $\bar{\phi}$, which provides stable training feedback
 700 for θ . When (I, λ) are fixed, optimizing π reduces to maximizing:

$$-\tau + \lambda \log I_\phi(H^*|s_\tau).$$

701 Hence, we define the reward at the transition $z = (s, a, s', d, H^*)$ (with the convention that $s' \leftarrow s$ if
 702 $a = a_{\text{stop}}$) as:

$$r_\lambda(z) := -1 + d\lambda \log I_{\bar{\phi}}(H^*|s'),$$

703 where $d = \mathbf{1}\{z \text{ is terminal}\}$ (z is terminal if the transition corresponds to the last time-step in
 704 a horizon, or $a = a_{\text{stop}}$). Furthermore, for a transition $z = (s, a, s', d, H^*)$ we define $z_{\text{stop}} :=$
 705 $z|_{(a, s') \leftarrow (a_{\text{stop}}, s)}$ as the same transition z with $a \leftarrow a_{\text{stop}}$ and $s' \leftarrow s$.

706 There is one thing to note: the logarithm in the reward is justified since the original problem can be
 707 equivalently written as:

$$\min_{\lambda \geq 0} \max_{I, \pi} -\mathbb{E}_{M \sim \mathcal{P}(\mathcal{M})}^{\pi} [\tau] + \lambda \left[\log \left(\mathbb{P}_{M \sim \mathcal{P}(\mathcal{M})}^{\pi} (h(\hat{H}_\tau; M) = 1) \right) - \log(1 - \delta) \right],$$

708 after noting that we can apply the logarithm to the constraint in eq. (1), before considering the dual.
 709 Thus the optimal solutions (I, π) remain the same.

710 Then, using classical TD-learning Sutton and Barto [2018], the training target for a transition
 711 $z = (s, a, s', d, H^*)$ can be defined as:

$$y_\lambda(z) = r_\lambda(z) + (1 - d)\gamma \max_{a'} Q_{\bar{\theta}}(s', a'),$$

712 where $\gamma \in (0, 1]$ is the discount factor.

713 As discussed earlier, we have a dedicated stopping action a_{stop} , whose value depends solely on history.
 714 Thus, its Q -value is updated retrospectively at any state s using an additional loss:

$$(r_\lambda(z_{\text{stop}}) - Q_\theta(s, a_{\text{stop}}))^2.$$

715 Therefore, the overall loss that we consider for θ for a single transition z can be written as

$$\mathbf{1}_{\{a \neq a_{\text{stop}}\}} (y_\lambda(z) - Q_\theta(s, a))^2 + (r_\lambda(z_{\text{stop}}) - Q_\theta(s, a_{\text{stop}}))^2,$$

716 where $\mathbf{1}_{\{a \neq a_{\text{stop}}\}}$ avoids double accounting for the stopping action.

717 To update parameters (θ, ϕ) , we sample independent batches $(B, B') \sim \mathcal{B}$ from the replay buffer and
 718 apply gradient updates as specified in eqs. (3) and (4) of algorithm 1. Target networks are periodically
 719 updated, with $\bar{\phi} \leftarrow \phi$ every M steps, and θ using Polyak averaging: $\bar{\theta} \leftarrow (1 - \eta)\bar{\theta} + \eta\theta$, $\eta \in (0, 1)$.

720 **Optimization over λ .** Finally, we update λ by assessing the confidence of I_ϕ at the stopping time
 721 according to eq. (2), maintaining a slow ascent-descent optimization schedule for sufficiently small
 722 learning rates.

723 **Implementation with the MAP estimator.** A practical implementation may consider to use the
 724 MAP estimator $\hat{H}_\tau = \arg \max_H I_\phi(H|s_\tau)$, which is what we do in practice, since it results in a
 725 lower variance estimator. We note that the loss function for I_ϕ , and the reward for Q_θ , as defined
 726 above, still yield the same optimal solution.

727 **Cost implementation.** Lastly, in practice, we optimize a reward $r_\lambda(z) = -c + dI_{\bar{\phi}}(H^*|s')$, by
 728 setting $c = 1/\lambda$, and noting that for a fixed λ the RL optimization remains the same. The reason why
 729 we do so is due to the fact that with this expression we do not have the product $\lambda \mathbb{E}_{H' \sim I_\phi} [h(H'; M)]$,
 730 which makes the descent-ascent process more difficult.

731 We also use the following cost update

$$c_{t+1} = c_t - \beta(1 - \delta - I_\phi(H_M^*|s_{\tau+1})),$$

732 or $c_{t+1} = c_t - \beta(1 - \delta - h(\hat{H}_\tau; M))$ if one uses the MAP estimator. To see why the cost can be
 733 updated in this way, define the parametrization $\lambda = e^{-x}$. Then the optimization problem becomes

$$\min_x \max_I \min_\pi -\mathbb{E}_{M \sim \mathcal{P}(\mathcal{M})}^\pi [\tau] + e^{-x} \left[\mathbb{P}_{M \sim \mathcal{P}(\mathcal{M})}^\pi \left(h(\hat{H}_\tau; M) = 1 \right) - 1 + \delta \right],$$

734 Letting $\rho = \mathbb{P}_{M \sim \mathcal{P}(\mathcal{M})}^\pi \left(h(\hat{H}_\tau; M) = 1 \right) - 1 + \delta$, the gradient update for x with a learning rate β
 735 simply is

$$x_{t+1} = x_t - \beta e^{-x_t} \rho,$$

736 implying that

$$-\log(\lambda_{t+1}) = -\log(\lambda_t) - \beta \lambda_t \rho.$$

737 Defining $c_t = 1/\lambda_t$, we have that

$$\log(c_{t+1}) = \log(c_t) - (\beta\rho/c_t) \Rightarrow c_{t+1} = c_t e^{\beta\rho/c_t}.$$

738 Using then the approximation $e^x \approx 1 + x$, we find $c_{t+1} = c_t + \beta\rho = c_t - \beta(1 - \delta - I_\phi(H_M^*|s_{\tau+1}))$.

739 **Training vs Deployment.** Thus far, our discussion of **ICPE** has focused on the training phase. After
 740 training completes, the learned policy π and inference network I can be deployed directly: during
 741 deployment, π both collects data and determines when to stop—either by triggering its stopping
 742 action or upon reaching the horizon N .

743 B.2 Other Algorithms

744 In this section we describe Track and Stop (TaS) [Garivier and Kaufmann \[2016\]](#), and some variants
 745 such as *I*-IDS, *I*-DPT and the explore then commit variant of **ICPE**.

746 B.2.1 Track and Stop

747 Track and Stop (TaS, [Garivier and Kaufmann \[2016\]](#)) is an asymptotically optimal as $\delta \rightarrow 0$ for MAB
 748 problems. For simplicity, we consider a Gaussian MAB problem with K actions, where the reward
 749 of each action is normally distributed according to $\mathcal{N}(\mu_a, \sigma^2)$, and let $\mu = (\mu_a)_{a \in [K]}$ denote the
 750 model. The TaS algorithm consists of: (1) the model estimation procedure and recommender rule; (2)
 751 the sampling rule, dictating which action to select at each time-step; (3) the stopping rule, defining
 752 when enough evidence has been collected to identify the best action with sufficient confidence, and
 753 therefore to stop the algorithm.

754 **Estimation Procedure and Recommender Rule** The algorithm maintains a maximum likelihood
755 estimate $\hat{\mu}_a(t)$ of the average reward for each arm based on the observations up to time t . Then, the
756 recommender rule is defined as $\hat{a}_t = \arg \max_a \hat{\mu}_a(t)$.

757 **Sampling Rule.** The sampling rule is based on the observation that any δ -correct algorithm, that is
758 an algorithm satisfying $\mathbb{P}(\hat{a}_\tau = a^*) \geq 1 - \delta$, with $a^* = \arg \max_a \mu_a$, satisfies the following sample
759 complexity

$$\mathbb{E}[\tau] \geq T^*(\mu) \text{kl}(1 - \delta, \delta),$$

760 where $\text{kl}(x, y) = x \log(x/y) + (1 - x) \log((1 - x)/(1 - y))$ and

$$(T^*(\mu))^{-1} = \sup_{\omega \in \Delta(K)} \min_{a \neq a^*} \frac{\omega_{a^*} \omega_a}{\omega_a + \omega_{a^*}} \frac{\Delta_a^2}{2\sigma^2},$$

761 with $\Delta_a = \mu_{a^*} - \max_{a \neq a^*} \mu_a$. Interestingly, to design an algorithm with minimal sample complexity,
762 we can look at the solution $\omega^* = \arg \inf_{\omega \in \Delta(K)} T(\omega; \mu)$, with $(T(\omega))^{-1} = \min_{a \neq a^*} \frac{\omega_{a^*} \omega_a}{\omega_a + \omega_{a^*}} \frac{\Delta_a^2}{2\sigma^2}$.

763 The solution ω^* provides the best proportion of draws, that is, an algorithm selecting an action a with
764 probability ω_a^* matches the lower bound and is therefore optimal with respect to T^* . Therefore, an idea
765 is to ensure that N_t/t tracks ω^* , where N_t is the visitation vector $N(t) := [N_1(t) \ \dots \ N_K(t)]^\top$.

766 However, the average rewards $(\mu_a)_a$ are initially unknown. A commonly employed idea [Garivier
767 and Kaufmann, 2016, Kaufmann et al., 2016] is to track an estimated optimal allocation $\omega^*(t) =$
768 $\arg \inf_{\omega \in \Delta(K)} T(\omega; \hat{\mu}(t))$ using the current estimate of the model $\hat{\mu}(t)$.

769 However, we still need to ensure that $\hat{\mu}(t) \rightarrow \mu$. To that aim, we employ a D-tracking rule Garivier
770 and Kaufmann [2016], which guarantees that arms are sampled at a rate of \sqrt{t} . If there is an
771 action a with $N_a(t) \leq \sqrt{t} - K/2$ then we choose $a_t = a$. Otherwise, choose the action $a_t =$
772 $\arg \min_a N_a(t) - t\omega_a^*(t)$.

773 **Stopping rule.** The stopping rule determines when enough evidence has been collected to determine
774 the optimal action with a prescribed confidence level. The problem of determining when to stop can
775 be framed as a statistical hypothesis testing problem [Chernoff, 1959], where we are testing between
776 K different hypotheses $(\mathcal{H}_a : (\mu_a > \max_{b \neq a} \mu_b))_a$.

777 We consider the following statistic $L(t) = tT(N(t)/t; \hat{\mu}(t))^{-1}$, which is a Generalized Likelihood
778 Ratio Test (GLRT), similarly as in [Garivier and Kaufmann, 2016]. Comparing with the lower bound,
779 one needs to stop as soon as $L(t) \geq \text{kl}(\delta, 1 - \delta) \sim \ln(1/\delta)$. However, to account for the random
780 fluctuations, a more natural threshold is $\beta(t, \delta) = \ln((1 + \ln(t))/\delta)$, thus we use $L(t) \geq \beta(t, \delta)$ for
781 stochastic MAB problems. We also refer the reader to Kaufmann and Koolen [2021] for more details.

782 B.2.2 I-IDS

783 We implement a variant of Information Directed Sampling (IDS) Russo and Van Roy [2018], where
784 we use the inference network I_ϕ learned during ICPE training as a posterior over optimal arms. This
785 approach enables IDS to exploit latent structure in the environment without explicitly modeling it via
786 a probabilistic model; instead, the learned I -network implicitly captures such structure.

787 By using the same inference network in both ICPE and I -IDS, we directly compare the exploration
788 policy learned by ICPE to the IDS heuristic applied on top of the same posterior distribution. While
789 computing the expected information gain in IDS requires intractable integrals, we approximate them
790 using a Monte Carlo grid of 30 candidate reward values per action. The full pseudocode for I -IDS is
791 given in Algorithm 3.

792 B.2.3 I-DPT

793 We implement a variant of DPT Lee et al. [2023] using the inference network. The idea is to act
794 greedily with respect to the posterior distribution I at inference time.

795 First, we train I using ICPE. Then, at deployment we act with respect to I : in round t we selection
796 action $a_t = \arg \max_H I(H|D_t)$. Upon observing x_{t+1} , we update D_{t+1} and stop as soon as
797 $\arg \max_H I(H|D_t) \geq 1 - \delta$.

798 **B.3 Transformer Architecture**

799 Here we briefly describe the architecture of the inference network I and of the network Q .

800 Both networks are implemented using a Transformer architecture. For the inference network, it is
 801 designed to predict a hypothesis H given a sequence of observations. Let the input tensor be denoted
 802 by $X \in \mathbb{R}^{B \times H \times m}$, where:

- 803 • B is the batch size,
- 804 • H is the sequence length (horizon), and
- 805 • $m = (d + |\mathcal{A}|)$, where d is the dimensionality of each observation x_t .

806 The inference network operates as follows:

- 807 1. **Embedding Layer:** Each observation vector $m_t = (x_t, a_t)$ is first embedded into a higher-
 808 dimensional space of size d_e using a linear transformation followed by a GELU activation:
 809 $h_t = \text{GELU}(W_{\text{embed}}m_t + b_{\text{embed}})$, $h_t \in \mathbb{R}^{d_e}$.
- 810 2. **Transformer Layers:** The embedded sequence $h \in \mathbb{R}^{B \times H \times d_e}$ is then passed through
 811 multiple Transformer layers (specifically, a GPT-2 model configuration). The Transformer
 812 computes self-attention over the embedded input to model dependencies among observations:

$$h' = \text{Transformer}(h), \quad h' \in \mathbb{R}^{B \times H \times d_e}.$$

- 813 3. **Output Layer:** The final hidden state corresponding to the last element of the sequence
 814 $(h'_{:, -1, :})$ is fed into a linear output layer that projects it to logits representing the hypotheses:

$$o = W_{\text{out}}h'_{:, -1, :} + b_{\text{out}}, \quad o \in \mathbb{R}^{B \times |\mathcal{H}|}.$$

- 815 4. **Probability Estimation:** The output logits are transformed into log-probabilities via a
 816 log-softmax operation along the last dimension

$$\log p(H|X) = \log_{\text{softmax}}(o).$$

817 For Q , we use the same architecture, but do not take a log-softmax at the final step.

Algorithm 3 I -IDS

- 1: **Input:** Pre-trained inference network I_ϕ ; prior means and variances μ_a, σ_a^2 for all $a \in \mathcal{A}$; target error threshold δ
 - 2: **Initialize:** $f_a(x) = \mathcal{N}(x \mid \mu_a, \sigma_a^2)$ for each a
 - 3: **for** $t = 1, 2, \dots$ **do**
 - 4: **if** $\max_a I_\phi(a \mid \mathcal{D}_{t-1}) \geq 1 - \delta$ **then**
 - 5: **return** $\arg \max_a I_\phi(a \mid \mathcal{D}_{t-1})$
 - 6: **end if**
 - 7: **for** each arm $a \in \mathcal{A}$ **do**
 - 8: Approximate information gain:

$$g_t(a) = H(I_\phi(\cdot \mid \mathcal{D}_{t-1})) - \mathbb{E}_{r \sim p(r \mid a, \mathcal{D}_{t-1})} [H(I_\phi(\cdot \mid \mathcal{D}_{t-1}, a, r))]$$
 - 9: **end for**
 - 10: Select action $a_t = \arg \max_a g_t(a)$
 - 11: Observe reward r_t
 - 12: Update dataset $\mathcal{D}_t = \mathcal{D}_{t-1} \cup \{(a_t, r_t)\}$
 - 13: **end for**
-

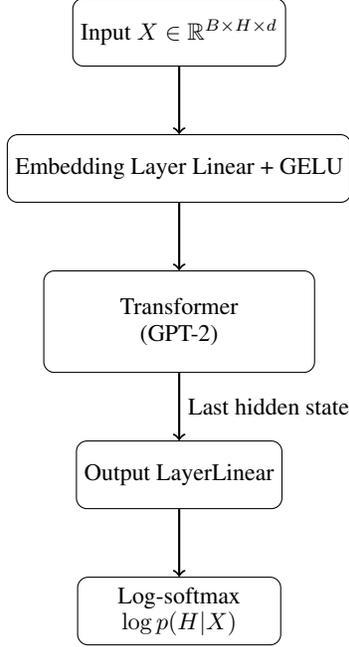


Figure 4: Model architecture for the inference network I (similarly for Q).

818 C Experiments

819 This section provides additional experimental results, along with detailed training and evaluation
 820 protocols to ensure clarity and reproducibility. All experiments were conducted using four NVIDIA
 821 A100 GPUs.

822 C.1 Bandit Problems

823 Here, we provide the implementation and evaluation details for the bandit experiments reported in
 824 Section 3.1, covering deterministic, stochastic, and structured settings.

825 **Model Architecture and Optimization.** For all bandit tasks, ICPE uses a Transformer with 3
 826 layers, 2 attention heads, hidden dimension 256, GELU activations, and dropout of 0.1 applied
 827 to attention, embeddings, and residuals. Layer normalization uses $\epsilon = 10^{-5}$. Inputs are one-hot
 828 action-reward pairs with positional encodings. Models are trained using Adam with learning rates
 829 between 1×10^{-4} and 1×10^{-6} , and batch sizes from 128 to 1024 depending on task complexity.

830 C.1.1 Stochastic Bandits Problems

831 In the stochastic Gaussian bandit setting, we evaluate ICPE on best-arm identification tasks with
 832 $K \in \{4, 6, 8, \dots, 14\}$. Arm means are sampled uniformly from $[0, 0.4K]$, with a guaranteed
 833 minimum gap of $1/K$ between the top two arms. All arms have a fixed reward standard deviation of
 834 0.5. The target confidence level is set to $\delta = 0.1$.

835 We compare ICPE against several widely used baselines: *Top-Two Probability Sampling (TTPS)* Jour-
 836 dan et al. [2022], *Track-and-Stop (TaS)* Garivier and Kaufmann [2016], *Uniform Sampling*, and
 837 *I-DPT*. For *I-DPT*, stopping occurs when the predicted optimal arm has estimated confidence at least
 838 $1 - \delta$. For *TTPS* and *TaS*, we apply the classical stopping rule based on the characteristic time $T^*(\hat{\mu}_t)$:

$$t \cdot T^*(\hat{\mu}_t) \geq \log \left(\frac{1 + \log t}{\delta} \right).$$

839 Each method is evaluated over three seeds, with 30 environments, and 30 trajectories per environment.
 840 95% confidence intervals were computed with hierarchical bootstrapping.

841 **C.1.2 Bandit Problems with Hidden Information**

842 **Magic Action Environments** We evaluate ICPE in bandit environments where certain actions
 843 reveal information about the identity of the optimal arm, testing its ability to uncover and exploit
 844 latent structure under the fixed-confidence setting.

845 Each environment contains $K = 5$ arms. Non-magic arms have mean rewards sampled uniformly
 846 from $[1, 5]$, while the mean reward of the designated *magic action* (always arm 1) is defined as
 847 $\mu_1 = \phi(\arg \max_{a \neq 1} \mu_a)$ with $\phi(i) = i/K$. The magic action is not the optimal arm, but it encodes
 848 information about which of the other arms is. To control the informativeness of this signal, we vary
 849 the standard deviation of the magic arm $\sigma_1 \in \{0.0, 0.1, \dots, 1.0\}$, while fixing the standard deviation
 850 of all other arms to $\sigma = 1 - \sigma_1$.

851 ICPE is trained under the fixed-confidence setting with a target confidence level of 0.9. For each
 852 σ_1 , we compare ICPE’s sample complexity to two baselines: (1) the average theoretical lower
 853 bound computed for the problem computed via averaging the result of Theorem A.2 over numerous
 854 environmental mean rewards, and (2) *I-IDS*, a pure-exploration information-directed sampling
 855 algorithm that uses ICPE’s *I*-network for posterior estimation. All methods are over 500 environments,
 856 with 10 trajectories per environment. 95% confidence intervals are computed using hierarchical
 857 bootstrapping with two levels.

858 Beyond the exploration efficiency analysis shown in Figure 2a, we also assess the correctness of
 859 each method’s final prediction and its usage of the magic action. As shown in Figure 5a, both
 860 ICPE and *I-IDS* consistently achieve the target accuracy of 0.9, validating their reliability under the
 861 fixed-confidence formulation.

862 Figure 5b plots the proportion of total actions that were allocated to the magic arm across different
 863 values of σ_1 . While both methods adapt their reliance on the magic action as its informativeness
 864 degrades, *I-IDS* tends to abandon it earlier. This behavior suggests that ICPE is better able to retain
 865 and exploit structured latent information beyond what is captured by simple heuristics for expected
 866 information gain.

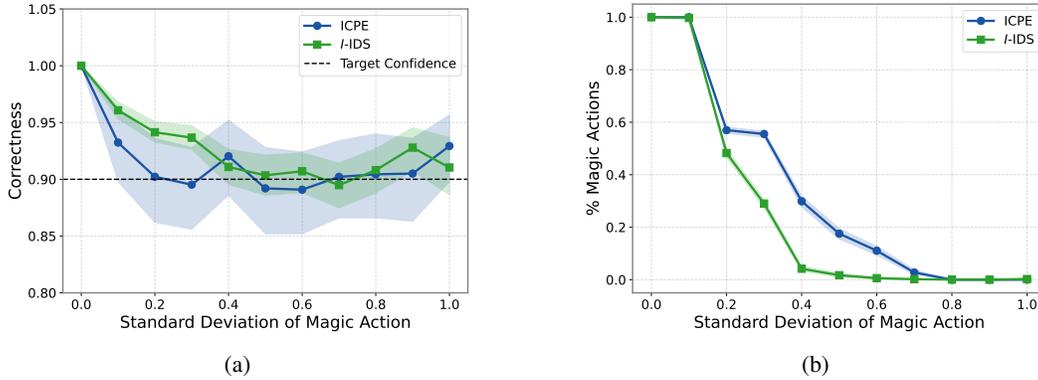


Figure 5: (a) Final prediction accuracy across varying levels of noise in the magic action (σ_1). Both ICPE and *I-IDS* consistently achieve the target confidence threshold of 0.9. (b) Percentage of actions allocated to the magic arm as a function of σ_1 . ICPE continues to exploit the magic action longer than *I-IDS*, suggesting more robust use of latent structure.

867 **Magic Chain Environments** To assess ICPE’s ability to perform multi-step reasoning over latent
 868 structure, we evaluate it in environments where identifying the optimal arm requires sequentially
 869 uncovering a chain of informative actions. In these *magic chain* environments, each magic action
 870 reveals partial information about the next, culminating in identification of the best arm.

871 We use $K = 10$ arms and vary the number of magic actions $n \in \{1, 2, \dots, 9\}$. Mean rewards for
 872 magic actions are defined recursively as:

$$\mu_{i_j} = \begin{cases} \phi(i_{j+1}), & \text{if } j = 1, \dots, n - 1, \\ \phi(\arg \max_{a \notin \{i_1, \dots, i_n\}} \mu_a), & \text{if } j = n, \end{cases}$$

873 where $\phi(i) = i/K$, and the remaining arms have mean rewards sampled uniformly from $[1, 2]$. All
 874 rewards are deterministic (zero variance).

875 ICPE is trained under the fixed-confidence setting with $\delta = 0.99$. For each n , five models are trained
 876 across five seeds. We compare ICPE’s average stopping time to the theoretical optimum (computed
 877 via Theorem A.4) and to the *I-IDS* baseline equipped with access to the *I*-network. Each model
 878 is evaluated over 1000 test environments per seed. 95% confidence intervals are computed using
 879 hierarchical bootstrapping.

880 In interpreting the results from Figure 2b, we observe that for environments with one or two magic
 881 actions, ICPE reliably learns the optimal policy of following the magic chain to completion. In these
 882 cases, the agent is able to identify the optimal arm without ever directly sampling it, by exploiting the
 883 structured dependencies embedded in the reward signals of the magic actions. Figure 6 illustrates a
 884 representative trajectory from the two-magic-arm setting, where the first magic action reveals the
 885 location of the second, which in turn identifies the optimal arm. The episode terminates without
 886 requiring the agent to explicitly sample the best arm itself.

Initial State										
?	?	?	?	?	?	?	?	?	?	Stop
Step 1: Selected Action 0										
0.400	?	?	?	?	?	?	?	?	?	Stop
Step 2: Selected Action 4										
0.400	?	?	?	0.900	?	?	?	?	?	Stop
Step 3: Selected STOP										
0.400	?	?	?	0.900	?	?	?	?	?	STOP

Figure 6: Example trajectory in the 2-magic-arm environment. ICPE follows the magic chain: the first magic action reveals the second, which identifies the optimal arm.

887 For environments with more than two magic actions, however, ICPE learns a different strategy. As the
 888 length of the magic chain increases, the expected sample complexity of following the chain from the
 889 start becomes suboptimal. Instead, ICPE learns to randomly sample actions until it encounters one of
 890 the magic arms and then proceeds to follow the chain from that point onward. This behavior represents
 891 an efficient, learned compromise between exploration cost and informativeness. Figure 7 shows an
 892 example trajectory from the six-magic-arm setting, where the agent initiates random sampling until it
 893 lands on a magic action, then successfully follows the remaining chain to identify the optimal arm.

Initial State										
?	?	?	?	?	?	?	?	?	?	Stop
Step 1: Selected Action 7										
?	?	?	?	?	?	?	1.299	?	?	Stop
Step 2: Selected Action 4										
?	?	?	?	0.600	?	?	1.299	?	?	Stop
Step 3: Selected Action 6										
?	?	?	?	0.600	?	0.500	1.299	?	?	Stop
Step 4: Selected Action 5										
?	?	?	?	0.600	0.200	0.500	1.299	?	?	Stop
Step 5: Selected Action 2										
?	?	0.900	?	0.600	0.200	0.500	1.299	?	?	Stop
Step 6: Selected Action 9										
?	?	0.900	?	0.600	0.200	0.500	1.299	?	1.916	Stop
Step 7: Selected STOP										
?	?	0.900	?	0.600	0.200	0.500	1.299	?	1.916	STOP

Figure 7: Example trajectory in the 6-magic-arm environment. Rather than starting from the first magic action, ICPE samples randomly until finding a magic action and then follows the chain to the optimal arm.

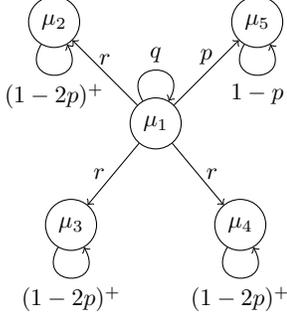


Figure 8: Loopy star graph.

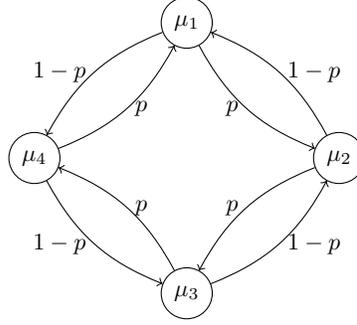


Figure 9: Ring graph.

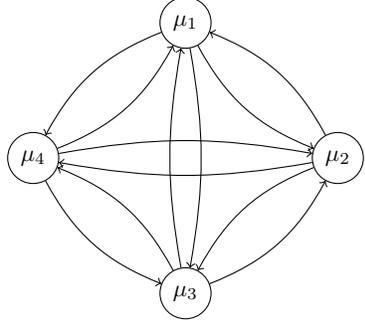


Figure 10: Loopless clique graph.

894 C.2 Exploration on Feedback Graphs

895 In the standard bandits setting we studied in Section 3.1, the learner observes the reward of the selected
 896 action, while in full-information settings, all rewards are revealed. Feedback graphs generalize this
 897 spectrum by specifying, via a directed graph G which additional rewards are observed when a
 898 particular action is chosen. Each node corresponds to an action, and an edge from u to v means that
 899 playing u may reveal feedback about v .

900 While feedback graphs have been widely studied for regret minimization [Mannor and Shamir \[2011\]](#),
 901 their use in pure exploration remains relatively underexplored [Russo et al. \[2025\]](#). We study them
 902 here as a challenging and structured testbed for in-context exploration. Unlike unstructured bandits,
 903 these environments contain latent relational structure and stochastic feedback dependencies that must
 904 be inferred and exploited to explore efficiently.

905 Formally, we define a feedback graph as an adjacency matrix $G \in [0, 1]^{K \times K}$, where $G_{u,v}$ denotes
 906 the probability that playing action u reveals the reward of action v . The learner observes a feedback
 907 vector $r \in \mathbb{R}^K$, where each coordinate is revealed independently with probability $G_{u,v}$:

$$r_v \sim \begin{cases} \mathcal{N}(\mu_v, \sigma^2), & \text{with probability } G_{u,v}, \\ \text{no observation,} & \text{otherwise.} \end{cases}$$

908 This setting allows us to test whether ICPE can learn to uncover and leverage latent graph struc-
 909 ture to guide exploration. We evaluate performance on best-arm identification tasks across three
 910 representative feedback graph families:

- 911 • **Loopy Star Graph** (Figure 8): A star-shaped graph with self-loops, parameterized by
 912 (p, q, r) . The central node observes itself with probability q , one neighboring node with
 913 probability p , and all others with probability r . When p is small, it may be suboptimal to
 914 pull the central node, requiring the agent to adapt its strategy accordingly.
- 915 • **Ring Graph** (Figure 9): A cyclic graph where each node observes its right neighbor with
 916 probability p and its left neighbor with probability $1 - p$. Effective exploration requires
 917 reasoning about which neighbors provide more informative feedback.
- 918 • **Loopless Clique Graph** (Figure 10): A fully connected graph with no self-loops. Edge
 919 probabilities are defined as:

$$G_{u,v} = \begin{cases} 0 & \text{if } u = v, \\ \frac{p}{u} & \text{if } v \neq u \text{ and } v \text{ is odd,} \\ 1 - \frac{p}{u} & \text{otherwise.} \end{cases}$$

920 Here, informativeness varies systematically with action index, requiring the learner to infer
 921 which actions are most useful.

922 These environments offer a diverse testbed for evaluating whether ICPE can uncover and exploit
 923 complex feedback structures without direct access to the underlying graph.

924 We tested ICPE in a fixed-confidence setting, using the same graph families but setting the optimal
 925 arm’s mean to 1 and all others to 0.5 to facilitate faster convergence. ICPE was trained for $K =$
 926 $4, 6, \dots, 14$ with a target error rate of $\delta = 0.1$. We compared it to Uniform Sampling, EXP3.G, and
 927 Tas-FG using a shared stopping rule from Russo et al. [2025].

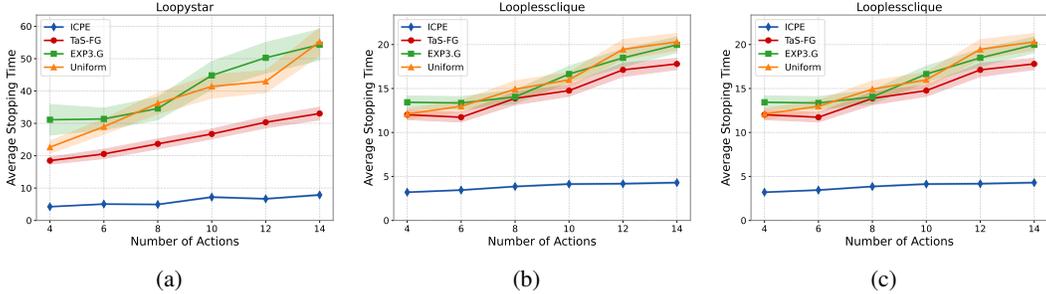


Figure 11: Sample complexity comparison under the fixed-confidence setting for: (a) Loopy Star, (b) Loopless Clique, and (c) Ring graphs.

928 As shown in Figure 11, ICPE consistently achieves significantly lower sample complexity than all
 929 baselines. This suggests that ICPE is able to meta-learn the underlying structure of the feedback
 930 graphs and leverage this knowledge to explore more efficiently than *uninformed* strategies. These
 931 results align with expectations: when environments share latent structure, learning to explore from
 932 experience offers a substantial advantage over fixed heuristics that cannot adapt across tasks.

933 C.3 Meta-Learning Binary Search

934 To test ICPE’s ability to recover classical exploration algorithms, we evaluate whether it can au-
 935 tonomously meta-learn binary search.

936 We frame the task as a structured multi-armed bandit problem where the optimal arm (i.e., the
 937 target number) is uniformly drawn from $1, \dots, K$. Pulling the correct arm yields a reward of +10,
 938 while pulling an arm above or below the target yields -1 or $+1$, respectively—providing directional
 939 feedback. The agent must learn to interpret and exploit this structure to efficiently locate the target.

940 We train ICPE under the fixed-confidence setting for $K = 2^3, \dots, 2^8$, using 150,000 in-context
 941 episodes and a target error rate of $\delta = 0.01$. Evaluation was conducted on 100 held-out tasks per
 942 setting. We report the minimum accuracy, mean stopping time, and worst-case stopping time, and
 943 compare against the theoretical binary search bound $O(\log_2 K)$.

Number of Actions (K)	Minimum Accuracy	Mean Stopping Time	Max Stopping Time	$\log_2 K$
8	1.00	2.13 ± 0.12	3	3
16	1.00	2.93 ± 0.12	4	4
32	1.00	3.71 ± 0.15	5	5
64	1.00	4.50 ± 0.21	6	6
128	1.00	5.49 ± 0.23	7	7
256	1.00	6.61 ± 0.26	8	8

Table 2: ICPE performance on the binary search task as the number of actions K increases.

944 As shown in Table 2, ICPE consistently achieves perfect accuracy with worst-case stopping times that
 945 match the optimal $\log_2(K)$ rate, demonstrating that it has successfully rediscovered binary search
 946 purely from data. While simple, this task illustrates ICPE’s broader potential to learn efficient search
 947 strategies in domains where no hand-designed algorithm is available.