
Learning to Explore: An In-Context Learning Approach for Pure Exploration

Alessio Russo*
Boston University
arusso2@bu.edu

Ryan Welch*
Massachusetts Institute of Technology
Broad Institute of MIT and Harvard
rcwelch@mit.edu

Aldo Pacchiano
Boston University
Broad Institute of MIT and Harvard
pacchian@bu.edu

Abstract

In this work, we study the active sequential hypothesis testing problem, also known as *pure exploration*, where the goal is to actively control a data collection process to efficiently identify the correct hypothesis underlying a decision problem. While relevant across multiple domains, devising adaptive exploration strategies remains challenging, particularly due to difficulties in encoding appropriate inductive biases. Existing Reinforcement Learning (RL)-based methods often underperform when relevant information structures are inadequately represented, whereas more complex methods, like *Best Arm Identification* (BAI) techniques, may be difficult to devise and typically rely on explicit modeling assumptions. To address these limitations, we introduce *In-Context Pure Exploration* (ICPE), an in-context learning approach that uses Transformers to learn exploration strategies directly from experience. ICPE combines supervised learning and reinforcement learning to identify and exploit latent structure across related tasks, without requiring prior assumptions. Numerical results across diverse synthetic and semi-synthetic benchmarks highlight ICPE’s capability to achieve robust performance performance in deterministic, stochastic, and structured settings. These results demonstrate ICPE’s ability to match optimal instance-dependent algorithms using only deep learning techniques, making it a practical and general approach to data-efficient exploration.

1 Introduction

Modern artificial intelligence systems have achieved remarkable performance across specialized tasks such as image classification [59], Super-human board-game play [104], protein-structure prediction [51] and large-scale language modelling [14]. Yet, there is still a lack in understanding how to autonomously discover meta-skills fundamental for sequential decision making, such as active testing or active learning [21, 23].

Consider an agent tasked with sequentially selecting samples to quickly improve its understanding of an underlying phenomenon. When the decision maker can exert some control over the collected samples’ *information content*, this is a problem also known as the *active sequential hypothesis testing problem* [21, 37, 81, 82, 77] or *pure exploration problem* [27, 28, 29]. Active hypothesis testing has

*Equal contribution (alphabetical order).

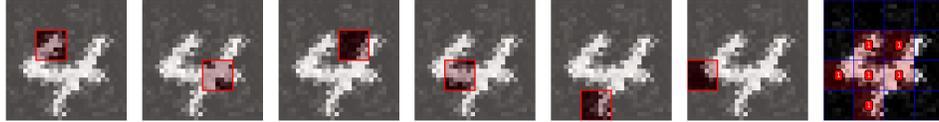


Figure 1: Starting from a blank canvas (left), **ICPE** sequentially chooses a region of pixel that is expected to maximally reduce the posterior entropy over classes, until it is confident to have inferred the true class (right).

become increasingly important nowadays, with applications ranging from medical diagnostics [12], image identification [110], recommender systems [93], etc.

Nonetheless, devising an adaptive data collection strategy is notoriously difficult and highly problem-specific. A source of issues remains the difficulty to encode the right inductive-biases in the model. For example, consider a tutor choosing which question to ask a student: each answer gives immediate feedback and hints at concepts they’ve (or haven’t) mastered. All these answers can reveal regularities such as “if a student misses x , they usually also miss y ”. But how can we *autonomously learn and encode* this information so that we can efficiently pinpoint the student’s knowledge gaps? A different, but relevant example, is that of Multi-Armed Bandit (MAB) problems [109, 62], where each action provides a random reward: how can we *autonomously learn and exploit* those situations where sampling one action provides indirect information about others, or where sequences of actions must be selected to uncover hidden information efficiently? Effective exploration requires exploiting the structural properties of the problem, but current approaches often necessitate explicitly specifying these biases, significantly limiting adaptability and efficiency.

In this paper, we address the question: how can sequential decision-making agents autonomously discover and leverage hidden structure to enhance active exploration for hypothesis testing? We introduce *In-Context Pure Explorer (ICPE)*, a novel method combining Supervised Learning and Deep RL [41, 78], which builds on the in-context learning and sequence modeling capabilities of Transformers [64]—a meta-learning approach that uncovers underlying shared structure across a class of problems \mathcal{M} [103, 11]. We hypothesize that these sequential models can learn to map histories of data to effective exploration strategies, leveraging their sequence-prediction capabilities.

ICPE operates by integrating two complementary neural networks: an inference (I) network, trained via supervised learning to infer the true hypothesis given current data, and an exploration (π) network, trained through reinforcement learning to select actions optimizing the inference accuracy of the I network. Crucially, this dual-network architecture, combined with the in-context learning abilities of Transformers, allows **ICPE** to learn from experience [105] how to autonomously identify and exploit regularities, enabling efficient exploration.

We validate **ICPE** through extensive synthetic and semi-synthetic experiments, demonstrating its ability to efficiently explore deterministic, stochastic, and structured environments. In particular, these results show that **ICPE** achieves performance comparable to optimal instance-dependent Best Arm Identification (BAI) algorithms [36, 5], without requiring explicit problem-specific exploration strategies that often involve solving complex optimization problems. Thanks to the in-context capability of **ICPE**, it is effectively discovering active sampling techniques [86], that at test time does not need much more computation than a forward pass. Consequently, **ICPE** emerges as a practical applicable method for data-efficient exploration.

1.1 Related Work

We now provide a brief overview of the related work and refer the reader to appendix A for an extended discussion. The problem of active sequential hypothesis testing [21, 37, 65, 81, 82, 77, 34], in which a learner is tasked with adaptively performing a sequence of actions to identify an unknown property of the environment, is closely related to the exploration problem in Reinforcement Learning (RL) [109], where an agent needs to identify the optimal policy. This exploration problem has long centered on regret minimization [109], with techniques based on Upper-Confidence Bounds [7, 8, 16, 61, 6], posterior-sampling [56, 88, 100, 42] and Information-Directed Sampling (IDS) [102]; yet these schemes assume that minimizing regret is the sole objective and falter in identification problems. A more closely related setting is that of pure exploration in bandits and Markov Decision

Processes (MDPs), settings known as Best Arm/Policy Identification (BAI/BPI) [5, 36, 27, 2, 96, 98]. In these problems the samples collected by the agent are no longer perceived as rewards, and the agent must actively optimize its exploration strategy to identify the optimal policy. BAI/BPI reframe the task as sequential hypothesis testing, yielding instance-adaptive algorithms in fixed-confidence settings such as Track-and-Stop (TaS) [36]. A similar setting is that of BAI/BPI with fixed horizon, which is, however, less understood [115, 4, 83] compared to the fixed-confidence one. However, while BAI strategy are powerful, they may be suboptimal when the underlying information structure is not adequately captured within the hypothesis testing framework. Hence, the issue of leveraging hidden environmental information, or problem with complex information structure remains a difficult problem. Although IDS and BAI offer frameworks to account for such structure, extending these approaches to Deep Learning is difficult, particularly when the information structure is unknown. Recently Transformers [113, 18] have demonstrated remarkable in-context learning capabilities [14, 35]. In-context learning [74] is a form of meta-RL [9], where agents can solve new tasks without updating any parameters by simply conditioning on additional context, such as their action-observation histories. Building on this ability, [64] recently showed that Transformers can be trained in a supervised manner using offline data to mimic posterior sampling in reinforcement learning. In [25] the authors present ICEE (In-Context Exploration Exploitation). ICEE uses Transformer architectures to perform in-context exploration-exploitation for RL. ICEE tackles this challenge by expanding the framework of return conditioned RL with in-context learning [18, 32]. Return conditioned learning is a type of technique where the agent learns the return-conditional distribution of actions in each state. Actions are then sampled from the distribution of actions that receive high return [107, 60]. Lastly, we note the important contribution of RL² [31], which proposes to represent an RL policy as the hidden state of an RNN, whose weights are learned via RL. ICPE employs a similar idea, but focuses on a different objective (identification), and splits the process into a supervised inference network that provides rewards to an RL-trained transformer network that selects actions to maximize information gain.

2 Learning to Explore: In-Context Pure Exploration

In this work, we introduce ICPE, an end-to-end deep-learning framework that combines sequential architecture with supervised and reinforcement learning to automatically discover efficient exploration policies for active sequential hypothesis testing. Instead of explicitly encoding inductive biases, we use transformers to let the agent autonomously infer the problem structure from experiences. We next provide an informal problem statement, followed by technical details.

Informal problem statement. Imagine a doctor sequentially ordering medical tests to produce a patient diagnosis $\hat{H} \in \text{HealthConditions}$. Each test result provides clues, helping the doctor narrow down which disease the patient truly has. Here, the doctor acts like an agent, using a sampling strategy π to choose each test possibly depending on previous tests and test results - or change to health conditions to efficiently gather information about the patient (the environment M). The patient’s test results are used by the doctor to produce a hypothesis \hat{H} . A good sampling strategy maximizes the probability that $\hat{H} = H^*$ where H^* is the true health condition of the patient. In the following, we consider a *general probabilistic*, or Bayesian, formulation in which uncertain quantities are modeled as random variables. Specifically, we adopt an in-context learning approach [64] and assume that M (e.g., the patient) is drawn from a patient environment class \mathcal{M} —a set of problems sharing some common characteristics—according to a distribution $\mathcal{P}(\mathcal{M})$.

Then, the central question we seek to answer is the following one:

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^ solely from \mathcal{D} ?*

Another relevant example is that of Best-Arm Identification in MAB problems [36]. Recall that in a MAB problem the decision maker can choose between K different actions a_1, \dots, a_K (we also say *arms*) at each time-step. Upon selecting an action a at time t , it observes a random reward r_t distributed according to a distribution ν_{a_t} . In 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 $H^* = a^*$). While several algorithms have been provided for different settings [106, 49, 97, 57, 90], a major issue is that the algorithm design can

drastically change if the assumptions change. Moreover, it is difficult to design efficient techniques for more complex settings such as MDPs (in fact, the problem becomes non-convex [70, 95]). Therefore, in this work we address the open question of whether it is possible to learn efficient exploration strategies directly from experience, avoiding the process of designing a BAI algorithm.

2.1 Problem Formulation

We now give a precise statement of the overall objective that we study in this paper. We begin with describing what is an environment, the set of assumptions, and the problems that we investigate.

Environment, sampling policy and hypotheses. We consider a model class of environments \mathcal{M} and a distribution $\mathcal{P}(\mathcal{M}) \in \Delta(\mathcal{M})$ from which the true environment M is sampled from. We model an 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 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 $\rho \in \Delta(\mathcal{X})$ denotes the initial observation distribution. All the environments in a class \mathcal{M} share the same set of observations \mathcal{X} and set of actions \mathcal{A} . The learner interacts with the environment in a sequential manner: (1) an initial observation $x_1 \sim \rho$ is sampled from \mathcal{X} ; (2) at time-step t , the 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 x_{t+1} is drawn independently from $P_t(\cdot | \mathcal{D}_t, a_t)$ given a trajectory $\mathcal{D}_t = (x_1, a_1, \dots, x_{t-1}, a_{t-1}, x_t)$. Formally, the learner uses a *randomized policy* $\pi = (\pi_t)_{t \in \mathbb{N}}$, which is a sequence of deterministic functions, to select actions: action a_t is selected by sampling independently from $\pi_t(\mathcal{D}_t)$ (with \mathcal{D}_t being a random variable), where $\pi_t(\mathcal{D}_t)$ specifies a probability distribution over \mathcal{A} . Finally, for each environment M , we assume the existence of a task-specific ground-truth hypothesis H_M^* contained in a predefined hypothesis class \mathcal{H} . For simplicity, we omit the subscript when the context is clear.

Oracle and estimator I . Our method assumes access to an oracle $h(\hat{H}; M) = \mathbf{1}_{\{\hat{H} = H_M^*\}}$ at training time (but not test time), which is a binary function that tells us whether the predicted hypothesis \hat{H} for environment M corresponds to the true hypothesis H_M^* ².

This oracle acts as a solution verifier, and, importantly, having access to this oracle is not the same as embedding an inductive bias into the model. The oracle does not reveal how to uncover hidden information; it only indicates whether our prediction \hat{H} is correct. As an example, consider the *label identification* problem. In this problem the agent seeks to correctly identify a label (e.g., the best action in a MAB problem) in a problem with K alternatives. We can set $\mathcal{H} = \{1, \dots, K\}$ and $H^* = i^*$, where i^* is the index of the correct label.

Using this oracle h , we propose to leverage supervised learning techniques and use the feedback provided by h to learn a mapping $I \in (\mathcal{D}_t \mapsto \Delta(\mathcal{H}))$ that computes the posterior distribution over the true hypothesis from trajectories of data. Then, the estimator $\hat{H}_t \sim I(\cdot | \mathcal{D}_t)$ can be thought as a Thompson sampling estimator, which can be used to provide a reward signal to an RL agent that collects the data \mathcal{D}_t using an exploration policy π (we can also use an estimator $\hat{H}_t = \arg \max_{H \in \mathcal{H}} I(H | \mathcal{D}_t)$).

Problem. Let \mathbb{P}_M^π be the underlying probability measure of the process $((\mathcal{D}_t, a_t))_t$ under a sampling strategy π . In the following we also write $\mathbb{P}_{M \sim \mathcal{P}(\mathcal{M})}^\pi(\cdot) = \mathbb{E}_{M \sim \mathcal{P}(\mathcal{M})}[\mathbb{P}_M^\pi(\cdot)]$ to denote the expected probability over the prior. Then, we consider the following two *online learning* problems:

- *Fixed confidence setting:* the agent needs to learn to stop the data sampling process as soon as it is sufficiently confident to have correctly estimated H^* for an environment M . We equip the learner with the capability to stop the sampling process at any point in time. We denote such stopping rule by τ , which is a stopping time with respect to the filtration $(\sigma(\mathcal{D}_t))_t$. Then, the learner wishes to find an optimal stopping rule τ (with $\tau < \infty$ a.s.),

²One may also consider a scenario in which, given a set of potential hypotheses $(H_i)_i$, the agent is tasked with predicting just one of them, and thus we can set $h(\hat{H}, M) = \max_i \mathbf{1}\{\hat{H} = H_i\}$.

exploration policy π and inference network I subject to a confidence level at the stopping time τ :

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

This setting is similar to Bayesian BAI with fixed confidence [47, 36], where the goal is to estimate the best action.

- *Fixed horizon (budget) setting*: in this case, for a given horizon $N \in \mathbb{N}$, we are interested in learning π, I such that

$$\max_{\pi, I} \mathbb{P}_{M \sim \mathcal{P}(\mathcal{M})}^{\pi} \left(h(\hat{H}_N; M) = 1 \right) \quad (2)$$

where the data \mathcal{D} is sampled from M using a policy π . This setting is similar to Bayesian fixed-budget BAI in the classical bandit setting [5, 54, 62, 4, 83].

To solve both problems, we consider a unified meta-RL approach that we discuss in the next section. In appendix C we also discuss more in detail the connection between our technique, and other pure exploration approaches, such as IDS [101] and Track and Stop [36].

2.2 ICPE: In-Context Pure Exploration

In this section we propose **ICPE**, a meta-RL approach for solving eqs. (1) and (2). The first aspect of our approach is the treatment of trajectories of data $\mathcal{D}_t = (x_1, a_1, \dots, x_t)$ as sequences to be given as input to sequential models, such as Transformers. This approach, inspired by sequential modeling [48], allows us to model the problem as a Markov Decision Process (MDP) [91]. An environment M can be then modeled as an MDP, which is a sequential 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 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 be defined later), $H^* \in \mathcal{H}$ is the true hypothesis in M and ρ is the initial state distribution. The model is similar to the one introduced in our problem formulation, but the state is fixed, and does not grow in time.

Nonetheless, we can use such modeling methodology for practical reasons: 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 for the remaining steps up to some pre-defined horizon N , with $s_1 = (x_1, \varnothing_{1:N})$. This limit N is a practical upper bound on the horizon that limits the dimensionality of the state, which is introduced for implementing the algorithm. The action space remains \mathcal{A} , and the transition dynamics P' are induced by (ρ, P) . By modeling the problem as an MDP, it allows us to cast the maximization over π as an RL problem. We now discuss the fixed confidence setting eq. (1) more in detail, while for the fixed horizon we refer the reader to appendix C.2.

2.2.1 Fixed Confidence Setting

In the fixed confidence setting (eq. (1)), problems terminate at some random point in time τ , chosen by the learner, or when the maximum horizon N is reached. We model this by giving π_t an additional stopping action a_{stop} such that $\pi_t : \mathcal{D}_t \rightarrow \mathcal{A} \cup \{a_{\text{stop}}\}$ so that the data collection processes terminates at the stopping-time $\tau = \min(N, \inf\{t \in \mathbb{N} : a_t = a_{\text{stop}}\})$.

To solve eq. (1) we consider the dual problem

$$\min_{\lambda \geq 0} \max_I \max_{\pi} V_{\lambda}(\pi, I) = -\mathbb{E}_{M \sim \mathcal{P}(\mathcal{M})}^{\pi} [\tau] + \lambda \left[\mathbb{P}_{M \sim \mathcal{P}(\mathcal{M})}^{\pi} \left(h(\hat{H}_{\tau}; M) = 1 \right) - 1 + \delta \right], \quad (3)$$

where $\hat{H}_{\tau} \sim I(\cdot | s_{\tau})$, with $s_{\tau} = \mathcal{D}_{\tau}$ under our MDP notation.

Then, we can use the MDP formalism to define an RL problem: we define a reward r that penalizes the agent at all time-steps, that is $r_t = -1$, while at the stopping-time we have $r_{\tau} = -1 + \lambda \mathbb{E}_{H \sim I(\cdot | s_{\tau})} [h(H; M)]$. Accordingly, one can define the Q -value of (π, I, λ) in a state-action pair (s, a) at the t -th step 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 $a_n \sim \pi_n(\cdot | s_n)$ (note that, since one can deduce the timestep from the nature of the state, due to the number of null tokens, we omit the subscript t from the Q -function). Therefore, when averaging over $\mathcal{P}(\mathcal{M})$, and the initial state, the value of (π, I, λ) is exactly $V_{\lambda}(\pi, I) = \mathbb{E}_{M \sim \mathcal{P}(\mathcal{M}), s_1 \sim \rho, a \sim \pi_1(\cdot | s_1)} [Q_{\lambda}^{\pi, I}(s_1, a)]$.

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 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

$$\lambda \leftarrow \max(0, \lambda - \beta (I_\phi(H^*|s_\tau) - 1 + \delta)). \quad (4)$$

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

$$\theta \leftarrow \theta - \alpha \nabla_\theta \frac{1}{|B|} \sum_{z \in B} \left[\mathbf{1}_{\{a \neq a_{\text{stop}}\}} (y_\lambda(z) - Q_\theta(s, a))^2 + (r_\lambda(z_{\text{stop}}) - Q_\theta(s, a_{\text{stop}}))^2 \right], \quad (5)$$

$$\phi \leftarrow \phi + \alpha \nabla_\phi \frac{1}{|B'|} \sum_{z \in B'} [\log(I_\phi(H^*|s'))]. \quad (6)$$

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

Practical implementation. Optimizing $\min_{\lambda \geq 0} \max_{I, \pi} V_\lambda(\pi, I)$ can be viewed as a multi-timescale stochastic optimization problem: the slowest timescale updates the variable λ , an intermediate timescale optimizes over I , and the fastest refines the policy π .

We treat each optimization separately, and optimize using a descent-ascent scheme. The distribution I is modeled using a sequential architecture with parameter ϕ , and we denote it by I_ϕ . Then, considering a fixed (π, λ) , the maximization problem in eq. (3) with respect to I amounts to solving $\max_\phi \mathbb{E}_{M \sim \mathcal{P}(\mathcal{M})}^\pi [h(\hat{H}_\tau; M)]$ with $\hat{H}_\tau \sim I_\phi(\cdot|s_\tau)$. Therefore, we can train ϕ with a cross-entropy loss $-\sum_{H'} h(H'; M) \log(I_\phi(H'|s_\tau)) = -\log(I_\phi(H^*|s_\tau))$, averaged over different environments (to that aim, we use a replay buffer, introduced down below).

The policy π is defined by the greedy policy with respect to the Q -values defined above. Therefore, we can use RL techniques to learn the Q -function of a policy π maximizing the value in eq. (3), assuming (λ, I) fixed. We denote this Q -function by Q_θ , which is implemented using a sequential architecture with a parameter θ .

We train θ using DQN [73, 111] with a replay buffer \mathcal{B} and a target network $Q_{\bar{\theta}}$ parameterized by $\bar{\theta}$ (the replay buffer is used to collect data during training). To maintain timescale separation, we introduce a separate target inference network $I_{\bar{\phi}}$, parameterized by $\bar{\phi}$, which provides feedback for training θ . Note that, as discussed earlier, we introduce a dedicated stop-action a_{stop} whose value depends solely on history. Thus, its Q -value can be updated at any time, allowing retrospective evaluation of stopping. For learning the Q -values, we define the reward for a transition $z = (s, a, s', d, H^*)$ as:

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

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. 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 Q -values can be learned using classical TD-learning techniques [109] and the target value is:

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

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

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

where $\mathbf{1}_{\{a \neq a_{\text{stop}}\}}$ avoids double accounting for the stopping action. Then, to train (θ, ϕ) , we sample two independent batches $(B, B') \sim \mathcal{B}$ from the buffer, and compute the gradient updates as in eqs. (5)

and (6) of algorithm 1. We periodically update target networks, setting $\bar{\phi} \leftarrow \phi$ every T steps and using a Polyak averaging $\theta \leftarrow (1 - \eta)\theta + \eta\bar{\theta}$, with $\eta \in (0, 1)$.

Finally, we update λ by assessing the confidence of I_ϕ at the stopping time (4) for a fixed (π, I) . Thus, for sufficiently small learning rates, optimizing (λ, θ, ϕ) resembles an ascent-descent scheme.

3 Empirical Evaluation

We evaluate our approach across various exploration tasks, including deterministic bandits with fixed budgets and stochastic bandits with or without latent structure. Additionally, we also evaluate ICPE in the broader setting of classifying images by sequentially revealing image patches. Due to space limitations, we refer the reader to appendix D for more details and more experiments on MAB problem with feedback graphs [98], MDPs with hidden information and an example of algorithmic discovery, where ICPE learns a probabilistic version of binary search.

Algorithms. In our evaluations we compare to different algorithms, depending on the problem. Some of the algorithms include: uniform sampling, DQN [72], TaS (Track and Stop) [36], TTPS (Top Two Sampling) [102] and Deep Contextual Multi-Armed Bandits (Deep CMAB) [24]. Moreover, we also include a variant of IDS [101] based on the I -mapping, which uses the observation that I defines a posterior distribution over \mathcal{H} . Always based on this idea, we also introduce I -DPT, a variant of DPT [64], based on the fact that I can be used to explore a problem à-la Thompson Sampling. More information about these methods, and their hyper-parameters, can be found in appendix C³.

3.1 Bandit Problems

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

Deterministic Bandits with Fixed Horizon. We first evaluate ICPE in the deterministic bandit environments with a fixed horizon K , equal to the number of actions. Therefore, ICPE needs to learn to select each action only once to determine the optimal action. Each action’s reward distribution is deterministic, so that $\nu_a = \delta_{\mu_a}$ is a Dirac distribution, with $(\mu_a)_{a \in \mathcal{A}}$ drawn from $\mathcal{P}(\mathcal{M}) = \mathcal{U}([0, 1]^K)$. Figure 2 summarizes performance relative to uniform sampling, DQN [72], and I -DPT. ICPE consistently identifies optimal actions (probability of correctness ≈ 1) and learns optimal sampling strategies (fraction of unique actions ≈ 1). Without being explicitly instructed to “choose each action exactly once”, ICPE discovers on its own that sampling every action is exactly what yields enough information to identify the best one. In contrast, baseline performances degrade significantly as the number of actions increases.

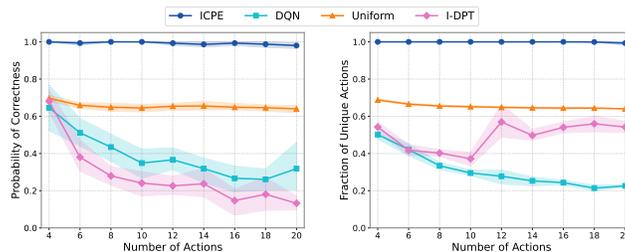


Figure 2: Deterministic bandits: (left) probability of correctly identifying the best action vs. K ; (right) average fraction of unique actions selected during exploration vs. K .

³In the results, shaded areas indicate 95% confidence intervals, computed via hierarchical bootstrapping.

Stochastic Bandit Problems. Next, we evaluate **ICPE** on stochastic bandit environments for both the fixed confidence (with $\delta = 0.1$ and $N = 100$) and fixed horizon setting (with horizon 30). Each action’s reward distribution is normally distributed $\nu_a = \mathcal{N}(\mu_a, 0.5^2)$, with $(\mu_a)_{a \in \mathcal{A}}$ drawn from $\mathcal{P}(\mathcal{M})$. In this case $\mathcal{P}(\mathcal{M})$ is a uniform distribution over problems with 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 this property to infer H^* more quickly. For this case, we also derive some sample complexity bounds in appendix B. Figures 3 and 4 summarize the numerical results for the two settings. In this case we also compare to

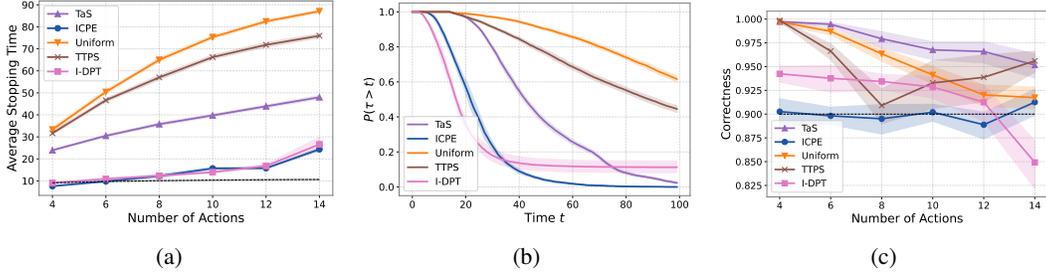


Figure 3: Results for stochastic MABs with fixed confidence $\delta = 0.1$ and $N = 100$: (a) average stopping time τ ; (b) survival function of τ ; (c) probability of correctness $\mathbb{P}_{M \sim \mathcal{P}(\mathcal{M})}^\pi(h(\hat{H}_\tau; M) = 1)$.

TaS and TTPS, and use the stopping rule of TaS also for Uniform and TTPS (the stopping rule is based on a self-normalized process, compared with a threshold function $\beta(t, \delta)$; see also appendix C for more details). Overall, we see in fig. 3a how **ICPE** is able to find a more efficient strategy compared to classical techniques. Interestingly, also **I-DPT** seems to achieve relatively small sample complexities. However, its tail distribution of τ is rather large compared to **ICPE** (fig. 3b) and the correctness is smaller than $1 - \delta$ for large values of K . Methods like TaS and TTPS achieve larger sample complexity, but also larger correctness values (fig. 3c). This is due to the fact that it is hard to define stopping rules. In fact, it is well known that current theoretically sound stopping rules are overly conservative [36]. Nonetheless, even using a less conservative rule such as $\beta(t, \delta) = \log((1 + \log(t))/\delta)$, which is what we use (and, yet, has not been proven to guarantee δ -correctness), is still conservative. The fact that **ICPE** can achieve the right value of confidence can help discover potential ways to define stopping rules. Lastly, in fig. 3a in black we show a complexity bound (proof in appendix B.2). While seemingly constant, it is actually *slowly* increasing in the number of arms.

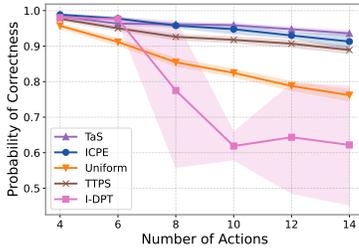


Figure 4: Correctness $\mathbb{P}_{M \sim \mathcal{P}(\mathcal{M})}^\pi(h(\hat{H}; M) = 1)$ for stochastic MABs with fixed horizon $N = 30$.

Lastly, in fig. 4, we compare the correctness of these methods in the fixed horizon setting with $N = 30$. Note that the training of **ICPE** differs from the fixed-confidence setting (i.e., there is no stopping action). While the results of **ICPE** are relatively good, the learning is more brittle compared to the fixed confidence (in fact, **I-DPT**, which is based on the I mapping learning by **ICPE**, has terrible performance). We believe that including a stopping action induces a form of curriculum learning, where the agent learns to adapt to the difficulty of the problem, greatly improving the learning process. Such effect is lost in the fixed horizon. Investigating this aspect is a future venue of research. Moreover, we believe **ICPE** can help shed a light on the fixed horizon problem, which is less studied compared to the fixed confidence one.

Bandit Problems with Hidden Information. To evaluate **ICPE** in structured settings, we introduce bandit environments with latent informational dependencies, termed *magic actions*. In the single magic action case, the magic action a_m ’s reward is distributed according to $\mathcal{N}(\mu_{a_m}, \sigma_m^2)$, where $\sigma_m \in (0, 1)$ and $\mu_{a_m} := \phi(\arg \max_{a \neq a_m} \mu_a)$ encodes information about the optimal action’s identity through an invertible mapping ϕ that is unknown to the learner. The index a_m is fixed, and the mean rewards of the other actions $(\mu_a)_{a \neq a_m}$ are sampled from $\mathcal{P}(\mathcal{M})$, a uniform distribution over models guaranteeing that a_m , as defined above, is not optimal (see appendices B.3 and D.1.3 for more details). Then, we define the reward distribution of the non-magic actions as $\mathcal{N}(\mu_a, (1 - \sigma_m)^2)$.

In our first experiment, we vary the standard deviation σ_m in $[0, 1]$. Thus, agents must balance sampling between informative and noisy actions based on varying uncertainty levels. We evaluate **ICPE** in a fixed-confidence setting with error rate $\delta = 0.1$. Figure 5a compares **ICPE**'s sample complexity against a theoretical lower bound (see appendix B) and an informed baseline, denoted as *I*-IDS, which performs standard IDS leveraging **ICPE**'s trained inference network *I* for exploiting the magic action (details in Appendix C). **ICPE** achieves sample complexities close to the theoretical

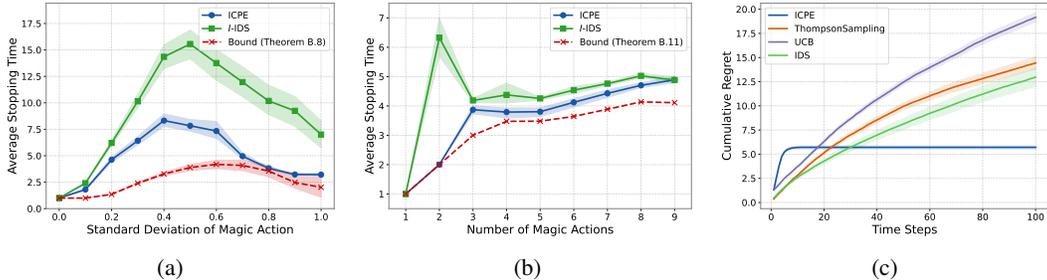


Figure 5: (a) Single magic action: average stopping time and the theoretical lower bound across varying σ_m . (b) Magic chain: average stopping time between **ICPE**, *I*-IDS vs. number of magic actions. (c) Cumulative regret minimization task with $\sigma_m = 0.1$.

bound across all tested noise levels, consistently outperforming *I*-IDS. Additionally, we evaluate **ICPE** in a cumulative regret minimization setting, despite **ICPE** not being explicitly optimized for regret minimization: at the stopping τ , **ICPE** commits to the identified best action (i.e., explore-then-commit strategy). As demonstrated in Figure 5c, **ICPE** outperforms classic algorithms such as UCB, Thompson Sampling, and standard IDS initialized with Gaussian priors.

To further challenge **ICPE**, we introduce a multi-layered "magic chain" bandit environments, where there is a sequence of n 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 $\mu_{a_{i_n}} = \phi(\arg \max_{a \notin \mathcal{A}_m} \mu_a)$. The first index i_1 is known, and by following the chain, an agent can uncover the best action in n steps. However, the optimal sample complexity depends on the ratio of magic actions to non-magic arms. Varying the number of magic actions from 1 to 9 in a 10-actions environment, Figure 5b demonstrates **ICPE**'s empirical performance, outperforming *I*-IDS.

3.2 Semi-Synthetic Experiment: Pixel Sampling for MNIST Classification

Lastly, to evaluate the applicability of **ICPE** to structured real-world decision problems, we introduce a classification task inspired by active perception settings. Formally, the environment class $\mathcal{P}(\mathcal{M})$ consists of MNIST images [63], each partitioned into a set of 36 distinct regions, corresponding to the action space $\mathcal{A} = \{1, \dots, 36\}$. Initially, all regions are masked, and at each timestep t , the agent selects an action $a_t \in \mathcal{A}$ to reveal the corresponding region in the image (example provided in fig. 1). The observation at time t is $x_t = o_t$, where o_t represents the state of the image after revealing the specified pixels. The episode concludes after a fixed budget of 12 steps. At this point, classification of the image label (or corresponding digit) H^* is performed by applying a pre-trained convolutional classifier to the partially revealed image. To promote generalization and prevent memorization of image-specific patterns, each image from the distribution $\mathcal{P}(\mathcal{M})$ undergoes several random augmentations (rotation, translation, noising, elastic deformation and contrast shifts).

For this setting we consider a slight variation of **ICPE** that may be of interest: we consider an inference net *I* that is a pre-trained classifier, trained on fully revealed images from $\mathcal{P}(\mathcal{M})$. Using this network, we benchmark **ICPE** against two baselines: standard uniform random sampling and Deep Contextual Multi-Armed Bandit (Deep CMAB) [24], which employs Bayesian neural networks to sample from a posterior distribution (Deep CMAB uses as rewards the correctness probabilities computed by *I*). Table 6b reports the classification accuracy and number of regions sampled. **ICPE** achieves substantially better performance than both baselines using fewer or comparable regions. However, to analyze whether **ICPE** learns a sampling strategy that adapts to the context of the task, we compare region selection distributions across digit classes using pairwise chi-squared tests. **ICPE** exhibits significantly more variation across classes than either baseline, as visualized in Figure 6a.

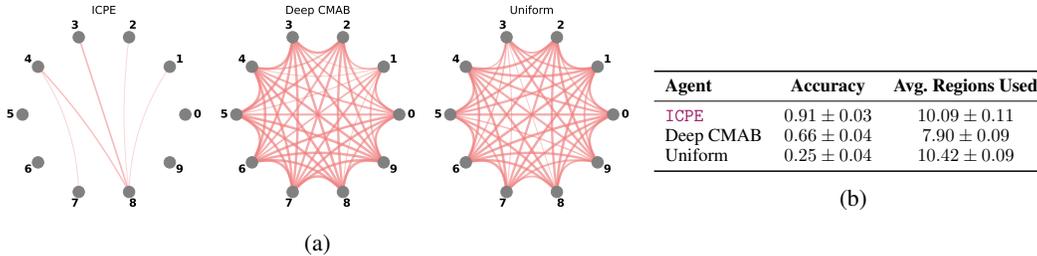


Figure 6: MNIST pixel-sampling task: (a) A chord between two digits indicates that their distributions were not significantly different (p -value > 0.05 , based on a pairwise chi-squared test), with thicker chords representing higher p -values.; (b) accuracy and performance (mean \pm 95% CI)

This suggests **ICPE** adapts its exploration to class-conditional structure, rather than applying a generic sampling policy.

4 Conclusions

In this work, we addressed the design of efficient pure-exploration strategies for the *active sequential hypothesis testing* problem, where an agent sequentially selects samples to rapidly identify the true hypothesis. While particularly relevant across different domains, it is difficult to design optimal strategies in the presence of hidden structure, and most of the existing optimal strategies are restricted to simple cases for unstructured multi-armed bandit problems. To overcome these limitations, we introduced **ICPE**, an in-context learning framework that leverages Transformers to learn exploration policies directly from experience. Our experiments span unstructured and structured stochastic/deterministic bandits, an adaptive MNIST pixel-revelation task, and complex MDP environments. Our results demonstrate that **ICPE** generalizes beyond tabular and graph-structured settings, autonomously discovering task-specific adaptive exploration strategies. We believe our work makes a fundamental contribution to active testing, and in particular to the sub-field of best-arm identification, where key questions—such as the gap between fixed-horizon and fixed-confidence settings—remain open. Future directions include: a theoretical analysis of **ICPE**’s guarantees; extending the framework to active regression and continuous-parameter hypotheses; improving the model architecture and computational efficiency to scale **ICPE** to larger, higher-dimensional problems. Lastly, designing methods capable of generalizing to unseen problems remains an important challenge.

Acknowledgments

The authors are pleased to acknowledge that the computational work reported on in this paper was performed on the Shared Computing Cluster administered by Boston University’s Research Computing Services and computing resources from the Laboratory for Information and Decision Systems at MIT. R.W. was supported by a Master of Engineering fellowship by the Eric and Wendy Schmidt Center at the Broad Institute.

References

- [1] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [2] A. Al Marjani, A. Garivier, and A. Proutiere. Navigating to the best policy in markov decision processes. *Advances in Neural Information Processing Systems*, 34:25852–25864, 2021.
- [3] D. Arumugam and T. L. Griffiths. Toward efficient exploration by large language model agents. *arXiv preprint arXiv:2504.20997*, 2025.
- [4] A. Atsidakou, S. Katariya, S. Sanghavi, and B. Kveton. Bayesian fixed-budget best-arm identification. *arXiv preprint arXiv:2211.08572*, 2022.
- [5] J.-Y. Audibert and S. Bubeck. Best arm identification in multi-armed bandits. In *COLT-23th Conference on learning theory-2010*, pages 13–p, 2010.
- [6] P. Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, 3(Nov):397–422, 2002.
- [7] P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47:235–256, 2002.
- [8] P. Auer, T. Jaksch, and R. Ortner. Near-optimal regret bounds for reinforcement learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 21, 2008.
- [9] J. Beck, R. Vuorio, E. Z. Liu, Z. Xiong, L. Zintgraf, C. Finn, and S. Whiteson. A survey of meta-reinforcement learning. *arXiv preprint arXiv:2301.08028*, 2023.
- [10] G. Bellala, S. Bhavnani, and C. Scott. Extensions of generalized binary search to group identification and exponential costs. *Advances in Neural Information Processing Systems*, 23, 2010.
- [11] Y. Bengio, S. Bengio, and J. Cloutier. *Learning a synaptic learning rule*. Citeseer, 1990.
- [12] S. M. Berry, B. P. Carlin, J. J. Lee, and P. Muller. *Bayesian adaptive methods for clinical trials*. CRC press, 2010.
- [13] L. Biewald. Experiment tracking with weights and biases, 2020. Software available from wandb.com.
- [14] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901, 2020.
- [15] S. Bubeck, R. Munos, and G. Stoltz. Pure exploration in finitely-armed and continuous-armed bandits. *Theoretical Computer Science*, 412(19):1832–1852, 2011.
- [16] O. Cappé, A. Garivier, O.-A. Maillard, R. Munos, and G. Stoltz. Kullback-leibler upper confidence bounds for optimal sequential allocation. *The Annals of Statistics*, pages 1516–1541, 2013.
- [17] F. Cecchi and N. Hegde. Adaptive active hypothesis testing under limited information. *Advances in Neural Information Processing Systems*, 30, 2017.
- [18] L. Chen, K. Lu, A. Rajeswaran, K. Lee, A. Grover, M. Laskin, P. Abbeel, A. Srinivas, and I. Mordatch. Decision transformer: Reinforcement learning via sequence modeling. In *Advances in Neural Information Processing Systems*, volume 34, pages 15084–15097, 2021.
- [19] X. Chen, Q. Liu, and Y. Wang. Active learning for contextual search with binary feedback. *Management Science*, 69(4):2165–2181, 2023.
- [20] H. Chernoff. Sequential design of experiments. *The Annals of Mathematical Statistics*, 30(3):755–770, 1959.

- [21] H. Chernoff. *Sequential design of experiments*. Springer, 1992.
- [22] J. Coda-Forno, M. Binz, Z. Akata, M. Botvinick, J. Wang, and E. Schulz. Meta-in-context learning in large language models. *Advances in Neural Information Processing Systems*, 36:65189–65201, 2023.
- [23] D. A. Cohn, Z. Ghahramani, and M. I. Jordan. Active learning with statistical models. *Journal of artificial intelligence research*, 4:129–145, 1996.
- [24] M. Collier and H. U. Llorens. Deep contextual multi-armed bandits. *arXiv preprint arXiv:1807.09809*, 2018.
- [25] Z. Dai, F. Tomasi, and S. Ghiassian. In-context exploration-exploitation for reinforcement learning. *arXiv preprint arXiv:2403.06826*, 2024.
- [26] S. Dasgupta. Analysis of a greedy active learning strategy. In *Advances in Neural Information Processing Systems*, volume 17, 2004.
- [27] R. Degenne and W. M. Koolen. Pure exploration with multiple correct answers. *Advances in Neural Information Processing Systems*, 32, 2019.
- [28] R. Degenne, W. M. Koolen, and P. Ménard. Non-asymptotic pure exploration by solving games. *Advances in Neural Information Processing Systems*, 32, 2019.
- [29] R. Degenne, P. Ménard, X. Shang, and M. Valko. Gamification of pure exploration for linear bandits, 2020.
- [30] S. Diamond and S. Boyd. CVXPY: A Python-embedded modeling language for convex optimization. *Journal of Machine Learning Research*, 17(83):1–5, 2016.
- [31] Y. Duan, J. Schulman, X. Chen, P. L. Bartlett, I. Sutskever, and P. Abbeel. RI^2 : Fast reinforcement learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*, 2016.
- [32] S. Emmons, B. Eysenbach, I. Kostrikov, and S. Levine. Rvs: What is essential for offline rl via supervised learning? *arXiv preprint arXiv:2112.10751*, 2021.
- [33] E. Even-Dar, S. Mannor, Y. Mansour, and S. Mahadevan. Action elimination and stopping conditions for the multi-armed bandit and reinforcement learning problems. *Journal of machine learning research*, 7(6), 2006.
- [34] K. Gan, S. Jia, and A. Li. Greedy approximation algorithms for active sequential hypothesis testing. *Advances in Neural Information Processing Systems*, 34:5012–5024, 2021.
- [35] S. Garg, D. Tsipras, P. S. Liang, and G. Valiant. What can transformers learn in-context? a case study of simple function classes. *Advances in Neural Information Processing Systems*, 35:30583–30598, 2022.
- [36] A. Garivier and E. Kaufmann. Optimal best arm identification with fixed confidence. In *Conference on Learning Theory*, pages 998–1027. PMLR, 2016.
- [37] B. K. Ghosh. A brief history of sequential analysis. *Handbook of sequential analysis*, 1, 1991.
- [38] D. Ghosh, M. K. Hanawal, and N. Zlatanov. Fixed budget best arm identification in unimodal bandits. *Transactions on Machine Learning Research*, 2024.
- [39] D. Golovin and A. Krause. Adaptive submodularity: Theory and applications in active learning and stochastic optimization. *Journal of Artificial Intelligence Research*, 42:427–486, 2011.
- [40] D. Golovin, A. Krause, and D. Ray. Near-optimal bayesian active learning with noisy observations. *Advances in Neural Information Processing Systems*, 23, 2010.
- [41] I. Goodfellow, Y. Bengio, A. Courville, and Y. Bengio. *Deep learning*, volume 1. MIT press Cambridge, 2016.

- [42] A. Gopalan, S. Mannor, and Y. Mansour. Thompson sampling for complex online problems. In *International conference on machine learning*, pages 100–108. PMLR, 2014.
- [43] Gurobi Optimization, LLC. Gurobi Optimizer Reference Manual, 2024.
- [44] C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau, E. Wieser, J. Taylor, S. Berg, N. J. Smith, R. Kern, M. Picus, S. Hoyer, M. H. van Kerkwijk, M. Brett, A. Haldane, J. F. del Río, M. Wiebe, P. Peterson, P. Gérard-Marchant, K. Sheppard, T. Reddy, W. Weckesser, H. Abbasi, C. Gohlke, and T. E. Oliphant. Array programming with NumPy. *Nature*, 585(7825):357–362, Sept. 2020.
- [45] K. Harris and A. Slivkins. Should you use your large language model to explore or exploit? *arXiv preprint arXiv:2502.00225*, 2025.
- [46] J. D. Hunter. Matplotlib: A 2d graphics environment. *Computing in science & engineering*, 9(3):90–95, 2007.
- [47] K. Jang, J. Komiyama, and K. Yamazaki. Fixed confidence best arm identification in the bayesian setting. *Advances in Neural Information Processing Systems*, 37, 2024.
- [48] M. Janner, Q. Li, and S. Levine. Offline reinforcement learning as one big sequence modeling problem. In *Advances in Neural Information Processing Systems*, volume 34, 2021.
- [49] Y. Jedra and A. Proutiere. Optimal best-arm identification in linear bandits. *Advances in Neural Information Processing Systems*, 33:10007–10017, 2020.
- [50] M. Jourdan, R. Degenne, D. Baudry, R. de Heide, and E. Kaufmann. Top two algorithms revisited. *Advances in Neural Information Processing Systems*, 35:26791–26803, 2022.
- [51] J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvunakool, R. Bates, A. Žídek, A. Potapenko, et al. Highly accurate protein structure prediction with alphafold. *nature*, 596(7873):583–589, 2021.
- [52] J. Kagan. Motives and development. *Journal of personality and social psychology*, 22(1):51, 1972.
- [53] Z. Karnin, T. Koren, and O. Somekh. Almost optimal exploration in multi-armed bandits. In *International conference on machine learning*, pages 1238–1246. PMLR, 2013.
- [54] E. Kaufmann, O. Cappé, and A. Garivier. On the complexity of best-arm identification in multi-armed bandit models. *The Journal of Machine Learning Research*, 17(1):1–42, 2016.
- [55] E. Kaufmann and W. M. Koolen. Mixture martingales revisited with applications to sequential tests and confidence intervals. *Journal of Machine Learning Research*, 22(246):1–44, 2021.
- [56] E. Kaufmann, N. Korda, and R. Munos. Thompson sampling: An asymptotically optimal finite-time analysis. In *International conference on algorithmic learning theory*, pages 199–213. Springer, 2012.
- [57] T. Kocák and A. Garivier. Best arm identification in spectral bandits. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI’20*, Yokohama, Yokohama, Japan, 2021.
- [58] A. Krishnamurthy, K. Harris, D. J. Foster, C. Zhang, and A. Slivkins. Can large language models explore in-context? In *Advances in Neural Information Processing Systems*, volume 37, 2024.
- [59] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, volume 25, 2012.
- [60] A. Kumar, X. B. Peng, and S. Levine. Reward-conditioned policies. *arXiv preprint arXiv:1912.13465*, 2019.

- [61] T. Lattimore and M. Hutter. Pac bounds for discounted mdps. In *Algorithmic Learning Theory: 23rd International Conference, ALT 2012, Lyon, France, October 29-31, 2012. Proceedings 23*, pages 320–334. Springer, 2012.
- [62] T. Lattimore and C. Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- [63] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [64] J. Lee, A. Xie, A. Pacchiano, Y. Chandak, C. Finn, O. Nachum, and E. Brunskill. Supervised pretraining can learn in-context reinforcement learning. *Advances in Neural Information Processing Systems*, 36:43057–43083, 2023.
- [65] D. V. Lindley. On a measure of the information provided by an experiment. *The Annals of Mathematical Statistics*, 27(4):986–1005, 1956.
- [66] J. Lisman, G. Buzsáki, H. Eichenbaum, L. Nadel, C. Ranganath, and A. D. Redish. Viewpoints: how the hippocampus contributes to memory, navigation and cognition. *Nature neuroscience*, 20(11):1434–1447, 2017.
- [67] G. Liu, M. Tang, and B. Eysenbach. A single goal is all you need: Skills and exploration emerge from contrastive rl without rewards, demonstrations, or subgoals. *arXiv preprint arXiv:2408.05804*, 2024.
- [68] S. Mannor and O. Shamir. From bandits to experts: On the value of side-observations. In *Advances in Neural Information Processing Systems*, volume 24, 2011.
- [69] A. A. Marjani, A. Garivier, and A. Proutiere. Navigating to the best policy in markov decision processes. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [70] A. A. Marjani and A. Proutiere. Adaptive sampling for best policy identification in markov decision processes. In *International Conference on Machine Learning*, pages 7459–7468. PMLR, 2021.
- [71] W. McKinney et al. Data structures for statistical computing in python. In *Proceedings of the 9th Python in Science Conference*, volume 445, pages 51–56. Austin, TX, 2010.
- [72] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- [73] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- [74] A. Moeini, J. Wang, J. Beck, E. Blaser, S. Whiteson, R. Chandra, and S. Zhang. A survey of in-context reinforcement learning. *arXiv preprint arXiv:2502.07978*, 2025.
- [75] G. Monea, A. Bosselut, K. Brantley, and Y. Artzi. Llms are in-context reinforcement learners. *arXiv preprint arXiv:2410.05362*, 2024.
- [76] L. S. Morris, M. M. Grehl, S. B. Rutter, M. Mehta, and M. L. Westwater. On what motivates us: a detailed review of intrinsic v. extrinsic motivation. *Psychological medicine*, 52(10):1801–1816, 2022.
- [77] S. Mukherjee, A. S. Tripathy, and R. Nowak. Chernoff sampling for active testing and extension to active regression. In *International Conference on Artificial Intelligence and Statistics*, pages 7384–7432. PMLR, 2022.
- [78] K. P. Murphy. *Probabilistic machine learning: Advanced topics*. MIT press, 2023.
- [79] L. Nadel. The hippocampus and space revisited. *Hippocampus*, 1(3):221–229, 1991.
- [80] L. Nadel and M. A. Peterson. The hippocampus: part of an interactive posterior representational system spanning perceptual and memorial systems. *Journal of Experimental Psychology: General*, 142(4):1242, 2013.

- [81] M. Naghshvar and T. Javidi. Active sequential hypothesis testing. *The Annals of Statistics*, 41(6):2703–2738, 2013.
- [82] M. Naghshvar, T. Javidi, and K. Chaudhuri. Noisy bayesian active learning. In *2012 50th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, pages 1626–1633. IEEE, 2012.
- [83] N. Nguyen, I. Aouali, A. György, and C. Vernade. Prior-Dependent Allocations for Bayesian Fixed-Budget Best-Arm Identification in Structured Bandits. In *Proceedings of The 28th International Conference on Artificial Intelligence and Statistics*, pages 379–387. PMLR, Apr. 2025.
- [84] A. Nie, Y. Su, B. Chang, J. N. Lee, E. H. Chi, Q. V. Le, and M. Chen. Evolve: Evaluating and optimizing llms for exploration. *arXiv preprint arXiv:2410.06238*, 2024.
- [85] R. Nowak. Generalized binary search. In *2008 46th annual Allerton conference on communication, control, and computing*, pages 568–574. IEEE, 2008.
- [86] J. Oh, M. Hessel, W. M. Czarnecki, Z. Xu, H. P. van Hasselt, S. Singh, and D. Silver. Discovering reinforcement learning algorithms. *Advances in Neural Information Processing Systems*, 33:1060–1070, 2020.
- [87] J. O’keefe and L. Nadel. Précis of o’keefe & nadel’s the hippocampus as a cognitive map. *Behavioral and Brain Sciences*, 2(4):487–494, 1979.
- [88] I. Osband, D. Russo, and B. Van Roy. (more) efficient reinforcement learning via posterior sampling. *Advances in Neural Information Processing Systems (NeurIPS)*, 26, 2013.
- [89] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in Neural Information Processing Systems (NeurIPS)*, 32, 2019.
- [90] R. Poiani, M. Jourdan, E. Kaufmann, and R. Degenne. Best-Arm Identification in Unimodal Bandits. In *Proceedings of The 28th International Conference on Artificial Intelligence and Statistics*, pages 2233–2241. PMLR, Apr. 2025.
- [91] M. L. Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- [92] T. Rainforth, A. Foster, D. R. Ivanova, and F. Bickford Smith. Modern bayesian experimental design. *Statistical Science*, 39(1):100–114, 2024.
- [93] P. Resnick and H. R. Varian. Recommender systems. *Communications of the ACM*, 40(3):56–58, 1997.
- [94] C. Rouyer, D. van der Hoeven, N. Cesa-Bianchi, and Y. Seldin. A near-optimal best-of-both-worlds algorithm for online learning with feedback graphs. *Advances in Neural Information Processing Systems*, 35:35035–35048, 2022.
- [95] A. Russo and A. Pacchiano. Adaptive exploration for multi-reward multi-policy evaluation. *arXiv preprint arXiv:2502.02516*, 2025.
- [96] A. Russo and A. Proutiere. Model-free active exploration in reinforcement learning. *Advances in Neural Information Processing Systems*, 36:54740–54753, 2023.
- [97] A. Russo and A. Proutiere. On the sample complexity of representation learning in multi-task bandits with global and local structure. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 9658–9667, 2023.
- [98] A. Russo, Y. Song, and A. Pacchiano. Pure exploration with feedback graphs. In *Proceedings of The 28th International Conference on Artificial Intelligence and Statistics*, Proceedings of Machine Learning Research. PMLR, 2025.
- [99] A. Russo and F. Vannella. Multi-reward best policy identification. *Advances in Neural Information Processing Systems*, 37:105583–105662, 2025.

- [100] D. Russo and B. Van Roy. Learning to optimize via posterior sampling. *Mathematics of Operations Research*, 39(4):1221–1243, 2014.
- [101] D. Russo and B. Van Roy. Learning to optimize via information-directed sampling. *Operations Research*, 66(1):230–252, 2018.
- [102] D. J. Russo, B. Van Roy, A. Kazerouni, I. Osband, Z. Wen, et al. A tutorial on thompson sampling. *Foundations and Trends® in Machine Learning*, 11(1):1–96, 2018.
- [103] T. Schaul and J. Schmidhuber. Metalearning. *Scholarpedia*, 5(6):4650, 2010.
- [104] D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre, D. Kumaran, T. Graepel, et al. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*, 362(6419):1140–1144, 2018.
- [105] D. Silver and R. S. Sutton. Welcome to the era of experience. *Google AI*, 2025.
- [106] M. Soare, A. Lazaric, and R. Munos. Best-arm identification in linear bandits. In *Advances in Neural Information Processing Systems*, volume 27, 2014.
- [107] R. K. Srivastava, P. Shyam, F. Mutz, W. Jaskowski, and J. Schmidhuber. Training agents using upside-down reinforcement learning. *CoRR*, abs/1912.02877, 2019.
- [108] J. Sun, Z. Wang, R. Yang, C. Xiao, J. Lui, and Z. Dai. Large language model-enhanced multi-armed bandits. *arXiv preprint arXiv:2502.01118*, 2025.
- [109] R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [110] N. K. Vaidhiyan, S. Arun, and R. Sundaresan. Active sequential hypothesis testing with application to a visual search problem. In *2012 IEEE International Symposium on Information Theory Proceedings*, pages 2201–2205. IEEE, 2012.
- [111] H. Van Hasselt, A. Guez, and D. Silver. Deep reinforcement learning with double q-learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 30, 2016.
- [112] G. Van Rossum and F. L. Drake Jr. *Python reference manual*. Centrum voor Wiskunde en Informatica Amsterdam, 1995.
- [113] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- [114] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, S. J. van der Walt, M. Brett, J. Wilson, K. J. Millman, N. Mayorov, A. R. J. Nelson, E. Jones, R. Kern, E. Larson, C. J. Carey, Í. Polat, Y. Feng, E. W. Moore, J. VanderPlas, D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen, E. A. Quintero, C. R. Harris, A. M. Archibald, A. H. Ribeiro, F. Pedregosa, P. van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:261–272, 2020.
- [115] P.-A. Wang, R.-C. Tzeng, and A. Proutiere. Best arm identification with fixed budget: A large deviation perspective. *Advances in Neural Information Processing Systems*, 36:16804–16815, 2023.
- [116] M. Waskom, O. Botvinnik, D. O’Kane, P. Hobson, S. Lukauskas, D. C. Gemperline, T. Augspurger, Y. Halchenko, J. B. Cole, J. Warmenhoven, J. de Ruiter, C. Pye, S. Hoyer, J. Vanderplas, S. Villalba, G. Kunter, E. Quintero, P. Bachant, M. Martin, K. Meyer, A. Miles, Y. Ram, T. Yarkoni, M. L. Williams, C. Evans, C. Fitzgerald, Brian, C. Fonnesbeck, A. Lee, and A. Qalieh. mwaskom/seaborn: v0.8.1 (september 2017), Sept. 2017.
- [117] A. X. Zheng, I. Rish, and A. Beygelzimer. Efficient test selection in active diagnosis via entropy approximation. In *Conference on Uncertainty in Artificial Intelligence*, 2005.

Appendix

Contents

1	Introduction	1
1.1	Related Work	2
2	Learning to Explore: In-Context Pure Exploration	3
2.1	Problem Formulation	4
2.2	ICPE: In-Context Pure Exploration	5
2.2.1	Fixed Confidence Setting	5
3	Empirical Evaluation	7
3.1	Bandit Problems	7
3.2	Semi-Synthetic Experiment: Pixel Sampling for MNIST Classification	9
4	Conclusions	10
	Acknowledgments	10
	Limitations and Broader Impact	19
A	Extended Related Work	20
B	Theoretical Results	22
B.1	Suboptimality of IDS	22
B.2	Sample Complexity Bounds for MAB Problems with Fixed Minimum Gap	25
B.3	Sample Complexity Lower Bound for the Magic Action MAB Problem	27
B.4	Sample Complexity Bound for the Multiple Magic Actions MAB Problem	32
C	Algorithms	36
C.1	ICPE with Fixed Confidence	36
C.2	ICPE with Fixed Horizon	38
C.3	Other Algorithms	39
C.3.1	Track and Stop	39
C.3.2	I -IDS	40
C.3.3	In-Context Explore-then-Commit	41
C.3.4	I -DPT	41
C.4	Transformer Architecture	41
D	Experiments	43
D.1	Bandit Problems	43
D.1.1	Deterministic Bandits with Fixed Horizon	43

D.1.2	Stochastic Bandits Problems	44
D.1.3	Bandit Problems with Hidden Information	45
D.2	Semi-Synthetic Pixel Sampling	48
D.3	MDP Problems: Magic Room	49
D.4	Exploration on Feedback Graphs	52
D.5	Meta-Learning Binary Search	54

Appendix

Limitations and Broader Impact

Finite vs continuous sets of hypotheses. One of the main limitations of this work is the assumption that \mathcal{H} is finite. This is a common assumption in active sequential hypothesis testing, and the continuous case is also referred to as *active regression* [77]. We believe our framework can be extended to this case with a proper parametrization of the inference mapping I that allows to sample from a continuous set.

On the prior set of problems $\mathcal{P}(\mathcal{M})$. One limitation of our approach is the assumption of access to a prior set of problems $\mathcal{P}(\mathcal{F})$. In reality, such set may lack a common structure and need not be stationary. Nonetheless, we view this as a useful starting point for developing more sophisticated methods. A natural direction for future work is to extend our framework to an adversarial setting, in which problem instances can evolve or even be chosen to thwart the learner.

Oracle assumption. Another limitation arises from assuming access to a *perfect* oracle $h(\hat{H}; M) = \mathbf{1}_{\{\hat{H}=H_M^*\}}$. In practice, computing such an oracle may be infeasible, and noise-free feedback is rarely available. We make this assumption because, in many settings, one can instantly verify the correct hypothesis (for example, classifying an image as a dog or selecting the best arm from a vector of averages). Furthermore, our main focus is in the sequential process of starting from "no data", to being able to predict the right hypothesis as quickly as possible (see the MNIST example). We believe this framework to be valuable when one can build verifiable simulations to train policies that transfer to real-world problems.

Practical limit on the horizon of ICPE. A limitation of ICPE is the current limit N on the horizon of the trajectory. This is due to the computation cost of training and using transformer architectures. Future work could investigate how to extend this limit, or completely remove it.

Scaling to larger problems with transformers. Another technical limitation of ICPE is the hardness to scaling to larger problems. This is closely related to the above limitation, and it is mainly an issue of investigating how to improve the current architecture of ICPE and/or distribute training.

On the training of transformers. Lastly, we believe that ICPE does not use the full capabilities of transformer architectures. For example, during training and evaluation, we always use the last hidden state of the transformer to make prediction, while the other hidden states are left untouched.

Bayesian BAI. Some of our work falls within the Bayesian Best Arm Identification theoretical framework. However, the Bayesian setting is less known compared to the frequentist one, and only recently some work [47] studied the unstructured Gaussian case. Future work should compare ICPE more thoroughly with Bayesian techniques once the Bayesian setting is more developed.

Broader impact . This paper primarily focuses on foundational research in pure exploration problems. Although we do not directly address societal impacts, we recognize their importance. The methods proposed here improve the sample efficiency of active sequential hypothesis testing procedures, and could be applied in various contexts with societal implications. For instance, our technique could be used in decision-making systems in healthcare, finance, and autonomous vehicles, where biases or errors could have significant consequences. Therefore, while the immediate societal impact of our work may not be evident, we urge future researchers and practitioners to carefully consider the ethical implications and potential negative impacts in their specific applications

A Extended Related Work

Exploration for Regret Minimization. The problem of exploration is particularly relevant in RL [109], and many strategies have been introduced, often with the goal of minimizing regret. Notably, approaches based on Posterior Sampling [56, 88, 100, 42] and Upper Confidence Bounds [7, 8, 16, 61, 6] have received significant attention. However, the problem of minimizing regret is a relevant objective only when one cares about the rewards accumulated so far, and does not answer the problem of how to efficiently gather data to reach some desired goal. In this context, *Information-Directed Sampling* (IDS) [100, 102] has been proposed to strike a balance between minimizing regret and maximizing information gain, where the latter is quantified as the mutual information between the true optimal action and the subsequent observation. However, when the information structure is unknown, it effectively becomes a significant challenge to exploit it. Importantly, if the state does not encode the structure of the problem, RL techniques may not be able to exploit hidden information.

In-Context Learning, LLMs and Return Conditioned Learning. Recently, Transformers [113, 18] have demonstrated remarkable in-context learning capabilities [14, 35]. In-context learning [74] is a form of meta-RL [9], where agents can solve new tasks without updating any parameters by simply conditioning on additional context such as their action-observation histories. When provided with a few supervised input-output examples, a pretrained model can predict the most likely next token [64]. Building on this ability, [64] recently showed that Transformers can be trained in a supervised manner using offline data to mimic posterior sampling in reinforcement learning. In [58] the authors investigate the extent to which LLMs [1] can perform in-context exploration in multi-armed bandit problems. Similarly, other works [22, 75, 84, 45, 108] evaluate the in-context learning capabilities of LLMs in sequential decision making problems, with [45] showing that LLMs can help at exploring large action spaces with inherent semantics. On a different note, in [3] investigate how to use LLMs to implement PSRL, leveraging the full expressivity and fluidity of natural language to express the prior and current knowledge about the problem.

In [25] the authors present ICEE (In-Context Exploration Exploitation), a method closely related to ICPE. ICEE uses Transformer architectures to perform in-context exploration-exploitation for RL. ICEE tackles this challenge by expanding the framework of return conditioned RL with in-context learning [18, 32]. Return conditioned learning is a type of technique where the agent learns the return-conditional distribution of actions in each state. Actions are then sampled from the distribution of actions that receive high return. This methodology was first proposed for the online RL setting by work on Upside Down RL [107] and Reward Conditioned Policies [60]. Lastly, we note the important contribution of RL² [31], which proposes to represent an RL policy as the hidden state of an RNN, whose weights are learned via RL. ICPE employs a similar idea, but focuses on a different objective (identification), and splits the process into a supervised inference network that provides rewards to an RL-trained transformer network that selects actions to maximize information gain.

Active Pure Exploration in Bandit and RL Problems. Other strategies consider the *pure exploration problem* [33, 5, 15, 54], or Best Arm Identification (BAI), in which the samples collected by the agent are no longer perceived as rewards, and the agent must actively optimize its exploration strategy to identify the optimal action. In this pure exploration framework, the task is typically formulated as a hypothesis testing problem: given a desired goal, the agent must reject the hypothesis that the observed data could have been generated by any environment whose behavior is fundamentally inconsistent with the true environment [36]. This approach leads to instance-dependent exploration strategies that adapt to the difficulty of the environment and has been extensively studied in the context of bandit problems under the fixed confidence setting [33, 36, 28, 47, 98], where the objective is to identify the optimal policy using the fewest number of samples while maintaining a specified level of confidence. Similar ideas have been applied to Markov Decision Processes for identifying the best policy [70, 69, 96, 99] or rapidly estimating the value of a given policy [95]. Another setting is that of identifying the best arm in MAB problems with a fixed horizon. In this case characterizing the complexity of the problem is challenging, and this is an area of work that is less developed compared to the fixed confidence one [115, 53, 5, 4, 83, 38]. Because of this reason, we believe ICPE can help better understand the nuances of this specific setting.

However, while BAI strategies are powerful, they may be suboptimal when the underlying information structure is not adequately captured within the hypothesis testing framework. Hence, the issue of leveraging hidden environmental information, or problem with complex information structure remains

a difficult problem. Although IDS and BAI techniques offer frameworks to account for such structure, extending these approaches to Deep Learning is difficult, particularly when the information structure is unknown to the learner.

A closely related work is that of [67]. In [67] the authors present empirical evidence of skills and directed exploration emerging from using RL with a sparse reward and a contrastive loss. They define a goal state, and encode a sparse reward using that goal state. Their objective, which maximizes the probability of reaching the goal state, is similar to ours, where in our framework the goal state would be a hypothesis. Note, however, that they do not learn an inference network, and we do not assume the observations to possess the Markov property.

Active Learning and Active Sequential Hypothesis Testing In the problem of active sequential hypothesis testing [21, 37, 65, 81, 82, 77, 34], a learner is tasked with adaptively performing a sequence of actions to identify an unknown property of the environment. Each action yields noisy feedback about the true hypothesis, and the goal is to minimize the number of samples required to make a confident and correct decision. Similarly, active learning [23, 19] studies the problem of data selection, and, closely related, Bayesian Active Learning [39], or Bayesian experimental design [92], studies how to adaptively select from a number of expensive tests in order to identify an unknown hypothesis sampled from a known prior distribution.

Active sequential hypothesis testing generalizes the pure exploration setting in bandits and RL by allowing for the identification of arbitrary hypotheses, rather than just the optimal action. However, most existing approaches assume full knowledge of the observation model [81], which is the distribution of responses for each action under each hypothesis. While some work has attempted to relax this assumption to partial knowledge [17], it remains highly restrictive in practice. As in bandit settings, real-world exploration and hypothesis testing often proceed without access to the true observation model, requiring strategies that can learn both the structure and the hypothesis from interaction alone.

Algorithm Discovery. Our method is also closely related to the problem of discovering algorithms [86]. In fact, one can argue that ICPE is effectively discovering active sampling techniques. This is particularly important for BAI and Best Policy Identification (BPI) problems, where often one needs to solve a computationally expensive optimization technique numerous times. For BPI the problem is even more exacerbated, since the optimization problem is usually non-convex [70, 95].

Cognitive Theories of Exploration. Our approach draws inspiration from cognitive theories of exploration. Indeed, in animals, exploration arises naturally from detecting mismatches between sensory experiences and internal cognitive maps—mental representations encoding episodes and regularities within environments [87, 80]. Detection of novelty prompts updates of these cognitive maps, a function strongly associated with the hippocampus [79, 66]. Conversely, exploration can also be explicitly goal-directed: psychological theories posit that an internal representation of goals, combined with cognitive maps formed through experience, guides adaptive action selection [52, 76]. ICPE embodies these cognitive principles computationally: the exploration (π) network learns an internal model (analogous to a cognitive map), while the inference (I) network encodes goal-directed evaluation. This interplay enables ICPE to effectively manage exploration as an adaptive, structure-sensitive behavior.

B Theoretical Results

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

B.1 Suboptimality of IDS

In pure exploration IDS [101] the main objective is to maximize the *information gain*. For example, consider the BAI problem: we set $\alpha_t(a) = \mathbb{P}(\hat{H} = a | \mathcal{D}_t)$ to be the posterior distribution of the optimal arm. Then, the information gain is defined through the following quantity

$$g_t(a) = \mathbb{E}[H(\alpha_t) - H(\alpha_{t+1}) | \mathcal{D}_t, a_t = a],$$

which measures the expected reduction in entropy of the posterior distribution of the best arm due to selecting arm a at time t .

For the BAI problem, the authors in [101] propose a *myopic* sampling policy $a_t \in \arg \max_a g_t(a)$, which only considers the information gain from the next sample. The reason for using a greedy policy stems from the fact that such a strategy is competitive with the optimal policy in problems where the information gain satisfies a property named *adaptive submodularity* [39], a generalization of submodular set functions to adaptive policies. For example, in the noiseless Optimal Decision Tree problem, it is known [117] that a greedy strategy based on the information gain is equivalent to a nearly-optimal [26, 40, 39] strategy named *generalized binary search* (GBS) [85, 10], which maximizes the expected reduction of the *version space* (the space of hypotheses consistent with the data observed so far). However, for the noisy case both strategies perform poorly [40].

The myopic pure exploration IDS strategy $a_t \in \arg \max_a g_t(a)$ can perform poorly in environments where the sampling decisions influence the observation distributions, or where an action taken at time t can greatly affect the complexity of the problem at a later stage (hence, IDS can perform poorly on credit assignments problems).

First example. As a first example, consider a bandit problem with K arms, where the reward for each arm a_i is distributed according to $\mathcal{N}(\mu_i, 1)$, with priors $\mu_1 = \delta_0$ and $\mu_i \sim \mathcal{U}([0, 1])$ independently for each $i \in \{2, \dots, K\}$. Thus, almost surely, the optimal arm a^* lies within $\{2, \dots, K\}$, and the goal is to estimate a^* .

We introduce the following twist: if arm a_1 is sampled exactly twice, its reward distribution changes permanently to a Dirac delta distribution $\delta_{\phi(a^*)}$, where ϕ is a known invertible mapping. Consequently, sampling arm a_1 twice fully reveals the identity of a^* . However, if arm a_1 has not yet been sampled, the expected immediate information gain at any step t is zero, i.e., $g_t(a_1) = 0$, since arm a_1 is already known to be suboptimal. In contrast, the immediate information gain for any other arm remains strictly positive. Therefore, under this setting and for nontrivial values of (σ, K) , the myopic IDS strategy cannot achieve the optimal constant sample complexity, and instead scales linearly in K .

Second example. Another example is a bandit environment containing a chain of two magic actions $\{1, m\}$, where the index of the first magic action (1) is known. Action 1 reveals the index m , and pulling arm m subsequently identifies the best arm with certainty. In this scenario, IDS is myopic and typically neglects arm 1 because of its inability to plan more than 1-step ahead in the future. However, depending on the total number of arms and reward variances, IDS may still select arm 1 if doing so significantly reduces the set of candidate best arms faster than pulling other arms (e.g., if the variance is significantly large). The following theorem illustrates the sub-optimality of IDS.

Theorem B.1. *Consider a bandit environment with a chain of 2 magic actions. The reward of the regular arms is $\mathcal{N}(\mu_a, 1)$ with $\mu_a \sim \mathcal{U}([0, 1])$, $a \neq 1, m$. For $K \geq 5$ there exists $\delta_0 \in (0, 1/2)$ such that for any $\delta \leq \delta_0$, we have that IDS is not sample optimal in the fixed confidence setting.*

Proof of theorem B.1. Let $Y_{t,a}$ be the random reward observed upon selecting arm $a_t = a$. We use that $g_t(a) = I_t(A^*; Y_{t,a}) = \text{KL} \left(\frac{\mathbb{P}(A^*, Y_{t,a} | \mathcal{D}_t)}{\mathbb{P}(A^* | \mathcal{D}_t) \mathbb{P}(Y_{t,a} | \mathcal{D}_t)} \right)$, with \mathcal{D}_1 containing an empty observation. The proof relies on showing that action a_1 is not chosen during the first two rounds for large values of K .

In the proofs, for brevity, we write $\mathbb{P}_t(\cdot) = \mathbb{P}(\cdot | \mathcal{D}_t)$. Observe the following lemmas.

Lemma B.2. Let $Y_{t,a}$ be the random reward observed upon selecting arm $a_t = a$ and let $S_{t,a} = \mathbf{1}\{a_t = a \text{ is magic}\}$. Under the assumption that the agent knows with absolute certainty that a is magic after observing $Y_{t,a}$, we have that $I_t(A^*; Y_{t,a}) = I_t(A^*; Y_{t,a}, S_{t,a})$.

Proof. Note that

$$I_t(A^*; Y_{t,a}, S_{t,a}) = H_t(Y_{t,a}, S_{t,a}) - H_t(Y_{t,a}, S_{t,a} | A^*).$$

Note that by assumption we have that $H_t(S_{t,a} | Y_{t,a}) = 0$. Then, the first term can also be rewritten as

$$H_t(Y_{t,a}, S_{t,a}) = H_t(S_{t,a} | Y_{t,a}) + H_t(Y_{t,a}) = H_t(Y_{t,a}).$$

Similarly, we also have $H_t(Y_{t,a}, S_{t,a} | A^*) = H_t(S_{t,a} | Y_{t,a}, A^*) + H_t(Y_{t,a} | A^*) = H_t(Y_{t,a} | A^*)$. Henceforth

$$I_t(A^*; Y_{t,a}, S_{t,a}) = H_t(Y_{t,a}) - H_t(Y_{t,a} | A^*) = I_t(A^*; Y_{t,a}).$$

□

Using the decomposition from the previous lemma we can rewrite the mutual information between A^* and $Y_{t,a}$ as

$$I_t(A^*; Y_{t,a}) = I_t(A^*; Y_{t,a}, S_{t,a}) = I_t(A^*; Y_{t,a} | S_{t,a}) + I_t(A^*; S_{t,a}).$$

Lemma B.3. Let $\mathcal{E}_t = \{(a_1, \dots, a_{t-1}) \text{ are not magic actions}\}$, with $\mathcal{E}_1 = \emptyset$. Under $a_t = 1$ we have that $I_t(A^*; Y_{t,1} | \mathcal{E}_t) = \log \left(\frac{K - |\mathcal{A}_t|}{K - |\mathcal{A}_t| - 1} \right)$ where $\mathcal{A}_t = \{a | \exists i < t : a_t = a\}$ is the unique number of actions chosen in $t \in \{1, \dots, t-1\}$.

Proof. We use that $\mathbb{P}_t(S_{t,1} = 1 | a_t = 1) = 1$. Hence, for arm 1 we have

$$\begin{aligned} I_t(A^*; S_{t,1} | \mathcal{E}_t) &= \text{KL}(\mathbb{P}_t(A^*, S_{t,1} | \mathcal{E}_t) || \mathbb{P}_t(A^* | \mathcal{E}_t) \mathbb{P}_t(S_{t,1} | \mathcal{E}_t)), \\ &= \sum_{a \neq 1} \mathbb{P}(A^* = a | S_{t,1} = 1, \mathcal{E}_t) \log \left(\frac{\mathbb{P}(S_{t,1} = 1 | A^* = a, \mathcal{E}_t)}{\mathbb{P}_t(S_{t,1} = 1 | \mathcal{E}_t)} \right), \\ &= 0. \end{aligned}$$

Then, we have

$$\begin{aligned} Y_t(A^*; Y_{t,1} | S_{t,1}, \mathcal{E}_t) &= Y_1(A^*; Y_{t,1} | S_{t,1} = 1, \mathcal{E}_t), \\ &= \text{KL}(\mathbb{P}_t(A^*, Y_{t,1} | S_{t,1} = 1, \mathcal{E}_t) || \mathbb{P}_t(A^* | S_{t,1} = 1, \mathcal{E}_t) \mathbb{P}_t(Y_{t,1} | S_{t,1} = 1, \mathcal{E}_t)), \\ &= \text{KL}(\mathbb{P}_t(Y_{t,1} | A^*, S_{t,1} = 1, \mathcal{E}_t) || \mathbb{P}_t(Y_{t,1} | S_{t,1} = 1, \mathcal{E}_t)), \\ &= \log \left(\frac{1/(K - |\mathcal{A}_t| - 1)}{1/(K - |\mathcal{A}_t|)} \right), \end{aligned}$$

where we used that under \mathcal{E}_t exactly \mathcal{A}_t regular arms have been pulled and recognised as regular; the still-unrevealed set of candidates for the second magic arm has therefore size $K - |\mathcal{A}_t| - 1$ (since arm 1 is known to be magic). Thus the result follows from applying the previous lemma. □

Lemma B.4. For any un-pulled arm $a \neq 1$ at time t we have that $I_t(A^*; Y_{t,a} | \mathcal{E}_t) \geq \frac{1}{K - |\mathcal{A}_t| - 1} \log(K - |\mathcal{A}_t| - 1)$.

Proof. To compute the mutual information we use that $I_t(A^*; Y_{t,a} | \mathcal{E}_t) = I_t(A^*; Y_{t,a}, S_{t,a} | \mathcal{E}_t) = I_t(A^*; Y_{t,a} | S_{t,a}, \mathcal{E}_t) + I_t(A^*; S_{t,a} | \mathcal{E}_t)$. We start by computing the first term of this expression, and finding a non-trivial lower bound.

First term $I_t(A^*; Y_{t,a} | S_{t,a}, \mathcal{E}_t)$. Note that for $a \neq 1$ we have

$$\begin{aligned} I_t(A^*; Y_{t,a} | S_{t,a}, \mathcal{E}_t) &= \mathbb{P}_t(S_{t,a} = 0 | \mathcal{E}_t) I_t(A^*; Y_{t,a} | S_{t,a} = 0, \mathcal{E}_t) \\ &\quad + \mathbb{P}_t(S_{t,a} = 1 | \mathcal{E}_t) I_t(A^*; Y_{t,a} | S_{t,a} = 1, \mathcal{E}_t), \\ &\geq \frac{1}{K - |\mathcal{A}_t| - 1} I_t(A^*; Y_{t,a} | S_{t,a} = 1, \mathcal{E}_t), \end{aligned}$$

where we used that under \mathcal{E}_t , we have a uniform prior over the remaining $K - (t - 1)$ un-pulled arms, and the agent knows that arm 1 is magic.

If $a \neq 1$ and $S_{t,a} = 1$, then a is the second magic arm. Therefore we have $\mathbb{P}_t(Y_{t,a}|A^*, S_{t,a} = 1, \mathcal{E}_t) = 1$. Hence $I_t(A^*; Y_{t,a}|S_{t,a} = 1, \mathcal{E}_t) = \log(K - |\mathcal{A}_t| - 1)$ since $Y_{t,a}$ can only take values uniformly over $K - |\mathcal{A}_t| - 1$ arms under the event $\{S_{t,a} = 1, \mathcal{E}_t\}$.

Second term $I_t(A^*; S_{t,a}|\mathcal{E}_t)$. We are only left with computing the following term

$$\begin{aligned}
I_t(A^*; S_{t,a}|\mathcal{E}_t) &= \mathbb{E}_{t, A^*, S_{t,a}} \left[\log \frac{\mathbb{P}(A^*, S_{t,a}|\mathcal{E}_t)}{\mathbb{P}(A^*|\mathcal{E}_t)\mathbb{P}(S_{t,a}|\mathcal{E}_t)} \right], \\
&= \sum_{b \in \{0,1\}} \mathbb{P}_t(S_{t,a} = b|\mathcal{E}_t) \text{KL}(\mathbb{P}_t(S_{t,a} = b|A^*, \mathcal{E}_t) || \mathbb{P}_t(S_{t,a} = b|\mathcal{E}_t)), \\
&\geq \mathbb{P}_t(S_{t,a} = 1|\mathcal{E}_t) \text{KL}(\mathbb{P}_t(S_{t,a} = 1|A^*, \mathcal{E}_t) || \mathbb{P}_t(S_{t,a} = 1|\mathcal{E}_t)), \\
&= \frac{1}{K - |\mathcal{A}_t| - 1} \text{KL}(\mathbb{P}_t(S_{t,a} = 1|A^*, \mathcal{E}_t) || \mathbb{P}_t(S_{t,a} = 1|\mathcal{E}_t)), \\
&= \frac{1}{K - |\mathcal{A}_t| - 1} \sum_{j \neq \{1,a\}} \mathbb{P}_t(S_{t,a} = 1|A^* = j, \mathcal{E}_t) \log \frac{\mathbb{P}_t(S_{t,a} = 1|A^* = j, \mathcal{E}_t)}{1/(K - |\mathcal{A}_t| - 1)}, \\
&= \frac{1}{K - |\mathcal{A}_t| - 1} \sum_{j \neq \{1,a\}} \frac{1}{K - t - 1} \log \frac{1/(K - |\mathcal{A}_t| - 2)}{1/(K - |\mathcal{A}_t| - 1)}, \\
&= \frac{1}{K - |\mathcal{A}_t| - 1} \frac{K - 2}{K - |\mathcal{A}_t| - 2} \log \frac{1/(K - |\mathcal{A}_t| - 2)}{1/(K - |\mathcal{A}_t| - 1)}, \\
&\geq \frac{1}{K - |\mathcal{A}_t| - 1} \log \frac{K - |\mathcal{A}_t| - 1}{K - |\mathcal{A}_t| - 2}.
\end{aligned}$$

Conclusion. Finally, we can derive that

$$\begin{aligned}
I_t(A^*; Y_{t,a}|\mathcal{E}_t) &= I_1(A^*; Y_{t,a}|S_{t,a}, \mathcal{E}_t) + I_t(A^*; S_{t,a}|\mathcal{E}_t), \\
&\geq \frac{1}{K - |\mathcal{A}_t| - 1} \log(K - |\mathcal{A}_t| - 1) + \frac{1}{K - |\mathcal{A}_t| - 1} \log \left(\frac{K - |\mathcal{A}_t| - 1}{K - |\mathcal{A}_t| - 2} \right), \\
&= \frac{1}{K - |\mathcal{A}_t| - 1} \log(K - |\mathcal{A}_t| - 1).
\end{aligned}$$

□

Lemma B.5. Assume $a_1 = j$ is a regular arm, pulled at the first time-step. Then $I_2(A^*; Y_{2,j}|a_1 = j) \leq \frac{1}{2} \ln(1 + \frac{1}{12\sigma^2})$.

Proof. First, note that

$$I_2(A^*; Y_{2,j} | a_1 = j) \leq I_2(\mu_j; Y_{2,j} | a_1 = j) = H_2(Y_{2,j}|a_1 = j) - \frac{1}{2} \ln(2\pi e\sigma^2)$$

Then, since $\text{Var}_2(Y_{2,j}|a_1 = j) = \text{Var}_2(\mu_j|a_1 = j) + \sigma^2 \leq 1/12 + \sigma^2$. Therefore $H_2(Y_{2,j}|a_1 = j) \leq \frac{1}{2} \ln(2\pi e(1/12 + \sigma^2))$. Hence $I_2(A^*; Y_{2,j}|a_1 = j) \leq \frac{1}{2} \ln(1 + \frac{1}{12\sigma^2})$. □

Hence, one can verify that for $K \geq 4$ the first magic arm will never be chosen at the first time-step. Similarly, at the second time-step the first magic arm will not be chosen if $K \geq 5$.

Consider the fixed-confidence setting with some confidence level $\delta < 1/2$. Let $\mathcal{A}_1 = \{\text{second magic arm sampled at } t = 1\}$ and $\mathcal{A}_2 = \{\text{second magic arm sampled at } t = 2\}$. Then, the sample complexity of IDS satisfies $\mathbb{E}[\tau_{IDS} | \mathcal{A}_1^c, \mathcal{A}_2^c] \geq 3$ for δ sufficiently small (since the sample complexity scales as $\log(1/(2.4\delta))$).

We also have that at the first time-step the decision is uniform over $\{2, \dots, K\}$. Lastly, if the first sampled arm is not magic, then it's a regular arm, and by the previous lemmas the information gain of such arm will be smaller than the information gain of another un-pulled arm. In fact the inequality

$$\frac{\log(x-2)}{x-2} > \frac{1}{2} \ln\left(1 + \frac{1}{12}\right)$$

it satisfied over $x \in \{4, \dots, 120\}$. Since it is sub-optimal to sample again the same regular arm Since the information gain on all the other arms remains the same, we have that the decision at the second time-step is again uniform over the remaining unchosen arms. Therefore

$$\begin{aligned}
\mathbb{E}[\tau_{IDS}] &= \mathbb{E}[\tau_{IDS}|\mathcal{A}_1]\mathbb{P}(\mathcal{A}_1) + \mathbb{E}[\tau_{IDS}|\mathcal{A}_1^c]\mathbb{P}(\mathcal{A}_1^c), \\
&= \frac{1}{K-1} + \frac{K-2}{K-1}\mathbb{E}[\tau_{IDS}|\mathcal{A}_1^c], \\
&= \frac{1}{K-1} + \frac{K-2}{K-1}\left(\mathbb{E}[\tau_{IDS}|\mathcal{A}_1^c, \mathcal{A}_2]\frac{1}{K-2} + \mathbb{E}[\tau_{IDS}|\mathcal{A}_1^c, \mathcal{A}_2^c]\frac{K-3}{K-2}\right), \\
&\geq \frac{1}{K-1} + \frac{K-2}{K-1}\left(2\frac{1}{K-2} + 3\frac{K-3}{K-2}\right), \\
&= \frac{3}{K-1} + 3\frac{K-3}{K-1},
\end{aligned}$$

which is larger than 2 for $K > 4$. Since there is a policy with sample complexity 2, we have that IDS cannot be sample optimal for $K \in \{5, \dots, 120\}$.

Similarly, for large values of $K > 120$, resampling the same regular arm at the second timestep leads IDS to a sample complexity larger than 2. And therefore cannot be sample optimal. \square

B.2 Sample Complexity Bounds for MAB Problems with Fixed Minimum Gap

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

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

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

Then, for any δ -correct algorithm, guaranteeing that at some stopping time τ the estimated optimal 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.

Theorem B.6. *Consider a model μ satisfying $\Delta_{\min} \geq \Delta_0 > 0$. Then, for any δ -probably correct method Alg, with $\delta \in (0, 1/2)$, we have that the optimal sample complexity is bounded as*

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

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

Proof. Step 1: Log-likelihood ratio. The initial part of the proof is rather standard, and follows the same argument used in the Best Arm Identification and Best Policy Identification literature [36, 99].

Define the set of models

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

and the set of alternative models

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

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

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

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

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

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

We also know from the information processing inequality [54] that $\mathbb{E}_\mu[\Lambda_\tau] \geq \sup_{\mathcal{E} \in \mathcal{M}_\tau} \text{kl}(\mathbb{P}_\mu(\mathcal{E}), \mathbb{P}_{\mu'}(\mathcal{E}))$, where $\mathcal{M}_t = \sigma(A_1, R_1, \dots, A_t, R_t)$. We use the fact that the algorithm 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 $\mathbb{P}_\mu(\mathcal{E}) \geq 1 - \delta$ 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 the monotonicity properties of the Bernoulli KL divergence). Hence

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

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

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

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

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

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

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

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

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 decomposed as

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

Hence, the optimization problem over the alternative models becomes

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

The inner infimum over μ' can then be written as

$$\begin{aligned} P_{\bar{a}}^*(\omega) &:= \inf_{\mu' \in \mathbb{R}^K} \sum_a \omega_a \frac{(\mu_a - \mu'_a)^2}{2\sigma^2}. \\ \text{s.t. } &\mu'_{\bar{a}} - \mu'_b \geq \Delta_0 \quad \forall b \neq \bar{a}. \end{aligned} \tag{7}$$

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

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

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 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 arms. Clearly We have that $\mu'_{\bar{a}} - \mu'_b \geq \Delta_0$ for all $b \neq \bar{a}$. Hence, we obtain

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

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 $\sup_{\omega \in \Delta(K)} \min_{\bar{a} \neq 1} P_{\bar{a}}^*(\omega) \leq \frac{\Delta_0^2}{2\sigma^2}$. This corresponds to the case where all the mass is given to $\omega_1 = 1$. Otherwise, the solution is to set $\omega_1 = 0$ and $\omega_a = \frac{1/\Delta_a^2}{\sum_b 1/\Delta_b^2}$ for $a \neq 1$.

Hence, we conclude that

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

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

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

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

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

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

Therefore

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

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

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

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

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

□

B.3 Sample Complexity Lower Bound for the Magic Action MAB Problem

We now consider a special class of models that embeds information about the optimal arm in the mean reward of some of the arms. Let $\phi : \mathbb{R} \rightarrow \mathbb{R}$ be a strictly decreasing function over $\{2, \dots, K\}$ ⁴.

Particularly, we make the following assumptions:

⁴One could also consider strictly increasing functions.

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

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

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

Hence, define the set of models

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

and the set of alternative models

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

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

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

Theorem B.7. *For any δ -correct algorithm, the sample complexity lower bound on the magic action problem is*

$$\mathbb{E}_\mu[\tau] \geq T^*(\mu) \text{kl}(1 - \delta, \delta), \quad (8)$$

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

$$(T^*(\mu))^{-1} = \max_{\omega \in \Delta(K)} \min_{a \neq 1, a^*} \omega_1 \frac{(\phi(a^*) - \phi(a))^2}{2\sigma_1^2} + \sum_{b \in \mathcal{K}_a(\omega)} \omega_b \frac{(\mu_b - m(\omega; \mathcal{C}_a(\omega)))^2}{2\sigma^2}, \quad (9)$$

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

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

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

Proof. Step 1: Log-likelihood ratio. The initial part of the proof is rather standard, and follows the same argument used in the Best Arm Identification and Best Policy Identification literature [36, 99].

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

$$\Lambda_\tau = \sum_a \sum_{n=1}^{\tau} \mathbf{1}_{\{A_n=a\}} \log \frac{f_a(R_n)}{f'_a(R_n)}$$

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

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

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

We also know from the information processing inequality [54] that $\mathbb{E}_\mu[\Lambda_\tau] \geq \sup_{\mathcal{E} \in \mathcal{M}_\tau} \text{kl}(\mathbb{P}_\mu(\mathcal{E}), \mathbb{P}_{\mu'}(\mathcal{E}))$, where $\mathcal{M}_t = \sigma(A_1, R_1, \dots, A_t, R_t)$. We use the fact that the algorithm 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

$\mathbb{P}_\mu(\mathcal{E}) \geq 1 - \delta$ 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 the monotonicity properties of the Bernoulli KL divergence). Hence

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

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

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

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

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

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

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

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

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

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

$$\begin{aligned} \inf_{\mu' \in \mathbb{R}^K} \quad & \omega_1 \frac{(\phi(a^*) - \phi(a))^2}{2\sigma_1^2} + \sum_{a \neq 1} \omega_a \frac{(\mu_a - \mu'_a)^2}{2\sigma^2} \\ \text{s.t.} \quad & \mu'_a \geq \max(\phi(2), \mu'_b) \quad \forall b, \\ & \mu'_1 = \phi(a). \end{aligned} \tag{10}$$

To solve it, we construct the following Lagrangian

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

where $\theta \in \mathbb{R}_+^K$ is the multiplier vector. From the KKT conditions we already know that $\theta_1 = 0, \theta_a = 0$ 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 $\mu'_b = \mu'_a$ 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 solution depends also on ω .

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 necessarily include arm a . Then, note that

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

and

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

Then, using the observations derived above, we conclude that

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

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} ?

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

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

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

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

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

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

As a corollary, we have the following upper bound on $T^*(\mu)$.

Corollary B.8. *We have that*

$$T^*(\mu) \leq \min_{\omega \in \Delta(K)} \max_{a \neq 1, a^*} \frac{2\sigma_1^2}{\omega_1 (\phi(a^*) - \phi(a))^2}.$$

In particular, for $\phi(x) = 1/x$ and $a^* < K$ we have

$$T^*(\mu) \leq 2\sigma^2 (a^*(a^* + 1))^2,$$

while for $a^* = K$ we get $T^*(\mu) \leq 2\sigma^2 (a^*(a^* - 1))^2$.

Proof. Let $f_1(a) = \frac{(\phi(a^*) - \phi(a))^2}{2\sigma_1^2}$. For every weight vector $\omega \in \Delta(K)$ and every $a \neq 1, a^*$, the quantity

$$g_a(\omega) = \omega_1 f_1(a) + \sum_{b \in \mathcal{K}_a} \omega_b \frac{(\mu_b - m(\omega; \mathcal{K}_a))^2}{2\sigma^2}$$

satisfies $g_a(\omega) \geq \omega_1 f_1(a)$ because the variance term is non-negative. Hence

$$(T^*(\mu))^{-1} = \max_{\omega \in \Delta(K)} \min_{a \neq 1, a^*} g_a(\omega) \geq \max_{\omega \in \Delta(K)} \omega_1 \min_{a \neq 1, a^*} f_1(a).$$

Since $\omega_1 \leq 1$, the right-hand side is lower bounded by $\omega_1 = 1$, giving

$$(T^*(\mu))^{-1} \geq \min_{a \neq 1, a^*} f_1(a) = \frac{1}{2\sigma_1^2} \min_{a \neq 1, a^*} (\phi(a^*) - \phi(a))^2.$$

Taking reciprocals yields

$$T^*(\mu) \leq \frac{2\sigma_1^2}{\min_{a \neq 1, a^*} (\phi(a^*) - \phi(a))^2} = \min_{\omega \in \Delta(K)} \max_{a \neq 1, a^*} \frac{2\sigma_1^2}{\omega_1 (\phi(a^*) - \phi(a))^2},$$

because the minimisation over ω clearly selects $\omega_1 = 1$. (This justifies the form stated in the corollary.)

Specialising to $\phi(x) = 1/x$. With $\phi(x) = 1/x$ the difference $\phi(a^*) - \phi(a) = \frac{1}{a^*} - \frac{1}{a}$ is positive for all $a > a^*$ and negative otherwise; its smallest non-zero magnitude is obtained for the *closest* index to a^* :

- If $a^* < K$, that index is $a^* + 1$, giving

$$\min_{a \neq 1, a^*} (\phi(a^*) - \phi(a))^2 = \left(\frac{1}{a^*} - \frac{1}{a^* + 1} \right)^2 = \frac{1}{[a^*(a^* + 1)]^2}.$$

- If $a^* = K$, the closest index is $K - 1$, leading to

$$\min_{a \neq 1, a^*} (\phi(a^*) - \phi(a))^2 = \left(\frac{1}{K - 1} - \frac{1}{K} \right)^2 = \frac{1}{[a^*(a^* - 1)]^2}.$$

Plugging each expression in the general upper bound above concludes the proof. \square

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

Lemma B.9. *With $K = 3$ we have that $\omega_1 = 1$ if and only if*

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

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

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

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

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

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

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 $f_3(\omega_a, \omega_{a^*}) = \frac{\omega_{a^*}(\mu_{a^*} - m(\omega))^2}{2\sigma^2}$. Then, we can write

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

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

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

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

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

Hence, we obtain that

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

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

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

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 $\omega_1 = 1$.

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 $\omega_{a^*} = (1 - \omega_1)/2$, leading to

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

which is maximized at $\omega_1 = 0$. \square

B.4 Sample Complexity Bound for the Multiple Magic Actions MAB Problem

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

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

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

and that the non-magic arms satisfy

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

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

Theorem B.10. *Consider noiseless magic bandit problem with K arms and n magic actions. The optimal sample complexity is upper bounded as*

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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 write

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

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

$$\begin{aligned} V(1) &= B(1), \\ V(2) &= B(2) + A(2)B(1), \\ V(3) &= B(3) + A(3)B(2) + A(3)A(2)B(1), \\ &\dots \\ V(r) &= \sum_{j=1}^r \left(\prod_{i=j+1}^r A(i) \right) B(j). \end{aligned}$$

Substituting back in our expression, we get

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

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

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

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

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

Corollary B.11. *Let T be the scaling in theorem B.10. We have that*

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

Proof. First, observe the scaling

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

At this point, note that

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

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

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

and

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

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

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

Therefore

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

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

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

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

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

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

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 $\min(n, (K-n)/2)$. \square

C Algorithms

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

MDP Formulation for ICPE. Recall that in **ICPE** we treat trajectories of data $\mathcal{D}_t = (x_1, a_1, \dots, x_t)$ as sequences to be given as input to sequential models, such as Transformers. We treat trajectories as states of an MDP M . An environment M can be then modeled as an MDP, which is a sequential 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 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 be defined later), $H^* \in \mathcal{H}$ is the true hypothesis in M and ρ is the initial state distribution.

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 for the remaining steps up to some pre-defined horizon N , with $s_1 = (x_1, \varnothing_{1:N})$.

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

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

C.1 ICPE with Fixed Confidence

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 (eq. (1)), problems terminate at some random point in time τ , chosen by the learner, or when the maximum horizon N is reached. We model this by giving π_t an additional stopping action a_{stop} such that $\pi_t : \mathcal{D}_t \rightarrow \mathcal{A} \cup \{a_{\text{stop}}\}$ so that the data collection processes terminates at the stopping-time $\tau = \min(N, t_{\text{stop}})$, with $t_{\text{stop}} := \inf\{t \in \mathbb{N} : a_t = a_{\text{stop}}\}$.

Optimizing the dual formulation

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

can be viewed as a multi-timescale stochastic optimization problem: the slowest timescale updates 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

$$\lambda \leftarrow \max(0, \lambda - \beta(I_\phi(H^*|s_{\tau+1}) - 1 + \delta). \quad (11)$$

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

$$\theta \leftarrow \theta - \alpha \nabla_\theta \frac{1}{|B|} \sum_{z \in B} \left[\mathbf{1}_{\{a \neq a_{\text{stop}}\}} (y_\lambda(z) - Q_\theta(s, a))^2 + (r_\lambda(z_{\text{stop}}) - Q_\theta(s, a_{\text{stop}}))^2 \right], \quad (12)$$

$$\phi \leftarrow \phi + \alpha \nabla_\phi \frac{1}{|B'|} \sum_{z \in B'} [\log(I_\phi(H^*|s))]. \quad (13)$$

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

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

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

Accordingly, one can define the Q -value of (π, I, λ) in a state-action pair (s, a) at the t -th step as $Q_\lambda^{\pi, I}(s, a) = \mathbb{E}_{M \sim \mathcal{P}(\mathcal{M})} \left[\sum_{n=t}^{\tau} r_n \mid s_t = s, a_t = a \right]$, with $a_n \sim \pi_n(\cdot|s_n)$

Optimization over ϕ . We treat each optimization separately, employing a descent-ascent scheme. The distribution I is modeled using a sequential architecture parameterized by ϕ , denoted by I_ϕ . Fixing (π, λ) , the inner maximization in eq. (3) corresponds to

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

We train ϕ via cross-entropy loss:

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

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

Optimization over π . The policy π is defined as the greedy policy with respect to learned Q -values. Therefore, standard RL techniques can learn the Q -function that maximizes the value in eq. (3) given (λ, I) . Denoting this function by Q_θ , it is parameterized using a sequential architecture with parameters θ .

We train Q_θ using DQN [73, 111], employing a replay buffer \mathcal{B} and a target network $Q_{\bar{\theta}}$ parameterized by θ . To maintain timescale separation, we introduce an additional inference target network $I_{\bar{\phi}}$, parameterized by $\bar{\phi}$, which provides stable training feedback for θ . When (I, λ) are fixed, optimizing π reduces to maximizing:

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

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

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

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

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

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

after noting that we can apply the logarithm to the constraint in eq. (3), before considering the dual. Thus the optimal solutions (I, π) remain the same.

Then, using classical TD-learning [109], the training target for a transition $z = (s, a, s', d, H^*)$ can be defined as:

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

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

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

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

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

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

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

To update parameters (θ, ϕ) , we sample independent batches $(B, B') \sim \mathcal{B}$ from the replay buffer and apply gradient updates as specified in eqs. (5) and (6) of algorithm 1. Target networks are periodically updated, with $\bar{\phi} \leftarrow \phi$ every M steps, and $\bar{\theta}$ using Polyak averaging: $\bar{\theta} \leftarrow (1 - \eta)\bar{\theta} + \eta\theta$, $\eta \in (0, 1)$.

Optimization over λ . Finally, we update λ by assessing the confidence of I_ϕ at the stopping time according to eq. (4), maintaining a slow ascent-descent optimization schedule for sufficiently small learning rates.

Implementation with the MAP estimator. A practical implementation may consider to use the 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 lower variance estimator. We note that the loss function for I_ϕ , and the reward for Q_θ , as defined above, still yield the same optimal solution.

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

We also use the following cost update

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

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 updated in this way, define the parametrization $\lambda = e^{-x}$. Then the optimization problem becomes

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

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 simply is

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

implying that

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

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

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

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

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

C.2 ICPE with Fixed Horizon

In the fixed horizon setting (problem in eq. (2)) the MDP terminates at time-step N , and we set the reward to be $r_t = 0$ for $t < N$ and $r_N = h(\hat{H}_N; M)$, where $\hat{H}_N \sim I(\cdot|\mathcal{D}_N)$ (or $\hat{H}_n = \arg \max_{H \in \mathcal{H}} I(H|\mathcal{D}_N)$) is the inferred hypothesis. Accordingly, one can define the value of (π, I) in a state-action pair (s, a) (at time-step t) as before, and therefore we obtain that $Q_t^{\pi, I}(s, a) = \mathbb{P}_{M \sim \mathcal{P}(\mathcal{M})}^\pi \left(h(\hat{H}_N; M) = 1 \mid s_t = s, a_t = a \right)$. Accordingly, we also define the value of (π, I) at the first time-step in state s as $V^{\pi, I}(s) = \mathbb{E}_{a \sim \pi_1(\cdot|s)}[Q_1(s, a)]$.

Therefore, averaging over the initial state $s \sim \rho$, we find the overall value of (π, I) , which corresponds to

$$V(\pi, I) = \mathbb{E}_{s \sim \rho}[V_M^{\pi, I}(s)] = \mathbb{P}_{M \sim \mathcal{P}(\mathcal{M})}^\pi \left(h(I(\hat{H}_N; M) = 1) \right). \quad (14)$$

This term $V(\pi, I)$ is exactly the overall probability of correctness, the quantity maximized in eq. (2).

Algorithm 3 ICPE (In-Context Pure Explorer) - Fixed Horizon

- 1: **Input:** Tasks distribution $\mathcal{P}(\mathcal{M})$; horizon N ; learning rate α and hyper-parameters (γ, M, η) .
- 2: Initialize buffer \mathcal{B} , networks Q_θ, E_ϕ and set $Q_{\bar{\theta}} \leftarrow Q_\theta, E_{\bar{\phi}} \leftarrow E_\phi$.
- 3: **while** Training is not over **do**
- 4: Sample environment $M \sim \mathcal{P}(\mathcal{M})$ and observe $s_1 \sim \rho$.
- 5: **for** $t = 1, \dots, N - 1$ **do**
- 6: Execute action $a_t = \arg \max_a Q_\theta(s_t, a)$ 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: **end for**
- 9: Sample batches $B, B' \sim \mathcal{B}$ and update θ, ϕ as

$$\begin{aligned}\theta &\leftarrow \theta - \alpha \nabla_\theta \frac{1}{|B|} \sum_{z \in B} \left[(y_\lambda(z) - Q_\theta(s, a))^2 \right], \\ \phi &\leftarrow \phi + \alpha \nabla_\phi \frac{1}{|B'|} \sum_{z \in B'} [\log(I_\phi(H^*|s'))].\end{aligned}$$

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

Practical implementation. The practical implementation for the fixed horizon follows closely that of the fixed confidence setting, and we refer the reader to that section for most of the details. In this case the reward in a transition $z = (s, a, s', d, H^*)$ is defined as as:

$$r_\lambda(z) := d \log \mathbb{E}_{H' \sim I_{\bar{\phi}}(\cdot|s')} [h(H'; M)] = d \log I_{\bar{\phi}}(H^*|s'), \quad (15)$$

where $d = \mathbf{1}\{s' \text{ terminal}\}$ (i.e., the last state observed in a trajectory). Note that we can use the logarithm, since solving the original problem is also equivalent to solving But note that the original problem is also equivalent to solving

$$\max_I \max_\pi \log \left(\mathbb{P}_{M \sim \mathcal{P}(\mathcal{M})}^\pi \left(h(\hat{H}_T; M) = 1 \right) \right), \quad (16)$$

due to monotonicity of the logarithm.

The Q -values can be learned using classical TD-learning techniques [109]: to that aim, for a transition $z = (s, a, s', d, H^*)$, we define the target:

$$y_\lambda(z) = r_\lambda(z) + (1 - d) \max_{a'} Q_{\bar{\theta}}(s', a'). \quad (17)$$

Then, the gradient updates are the same as for the fixed confidence setting.

C.3 Other Algorithms

In this section we describe Track and Stop (TaS) [36], and some variants such as I -IDS, I -DPT and the explore then commit variant of ICPE.

C.3.1 Track and Stop

Track and Stop (TaS, [36]) is an asymptotically optimal as $\delta \rightarrow 0$ for MAB problems. For simplicity, we consider a Gaussian MAB problem with K actions, where the reward of each action is normally distributed according to $\mathcal{N}(\mu_a, \sigma^2)$, and let $\mu = (\mu_a)_{a \in [K]}$ denote the model. The TaS algorithm consists of: (1) the model estimation procedure and recommender rule; (2) the sampling rule, dictating which action to select at each time-step; (3) the stopping rule, defining when enough evidence has been collected to identify the best action with sufficient confidence, and therefore to stop the algorithm.

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

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

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

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

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

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

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

However, the average rewards $(\mu_a)_a$ are initially unknown. A commonly employed idea [36, 54] is to track an estimated optimal allocation $\omega^*(t) = \arg \inf_{\omega \in \Delta(K)} T(\omega; \hat{\mu}(t))$ using the current estimate of the model $\hat{\mu}(t)$.

However, we still need to ensure that $\hat{\mu}(t) \rightarrow \mu$. To that aim, we employ a D-tracking rule [36], which guarantees that arms are sampled at a rate of \sqrt{t} . If there is an action a with $N_a(t) \leq \sqrt{t} - K/2$ then we choose $a_t = a$. Otherwise, choose the action $a_t = \arg \min_a N_a(t) - t\omega_a^*(t)$.

Stopping rule. The stopping rule determines when enough evidence has been collected to determine the optimal action with a prescribed confidence level. The problem of determining when to stop can be framed as a statistical hypothesis testing problem [20], where we are testing between K different hypotheses $(\mathcal{H}_a : (\mu_a > \max_{b \neq a} \mu_b))_a$.

We consider the following statistic $L(t) = tT(N(t)/t; \hat{\mu}(t))^{-1}$, which is a Generalized Likelihood Ratio Test (GLRT), similarly as in [36]. Comparing with the lower bound, 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 fluctuations, a more natural threshold is $\beta(t, \delta) = \ln((1 + \ln(t))/\delta)$, thus we use $L(t) \geq \beta(t, \delta)$ for stochastic MAB problems. We also refer the reader to [55] for more details.

C.3.2 I-IDS

We implement a variant of Information Directed Sampling (IDS) [101], where we use the inference network I_ϕ learned during ICPE training as a posterior over optimal arms. This approach enables IDS to exploit latent structure in the environment without explicitly modeling it via a probabilistic model; instead, the learned I -network implicitly captures such structure.

Algorithm 4 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**
-

By using the same inference network in both ICPE and I -IDS, we directly compare the exploration policy learned by ICPE to the IDS heuristic applied on top of the same posterior distribution. While computing the expected information gain in IDS requires intractable integrals, we approximate them using a Monte Carlo grid of 30 candidate reward values per action. The full pseudocode for I -IDS is given in Algorithm 4.

C.3.3 In-Context Explore-then-Commit

We implement an ICPE variant for regret minimization via an *explore-then-commit* framework. This method reuses the exploration policy and inference network learned during fixed-confidence training. The agent interacts with the environment using the learned exploration policy until it selects the stopping action. At that point, it commits to the arm predicted to be optimal by the I -network and repeatedly pulls that arm for the remainder of the episode. The full pseudo-code is provided in Algorithm 5.

Algorithm 5 In-Context Explore-then-Commit

```

1: Input: Environment  $M \sim \mathcal{P}(\mathcal{M})$ ; pre-trained critic network  $Q_\theta$ ; pre-trained inference network  $I_\phi$ 
2: Initialize  $stopped \leftarrow \text{False}$ 
3: Observe initial state  $s_1 \sim \rho$ 
4: for  $t = 1$  to  $N$  do
5:   if  $stopped = \text{False}$  and  $a_{\text{stop}} \neq \arg \max_a Q_\theta(s_t, a)$  then
6:     Execute  $a_t = \arg \max_a Q_\theta(s_t, a)$  and observe  $s_{t+1}$ 
7:   else if  $stopped = \text{False}$  and  $a_{\text{stop}} = \arg \max_a Q_\theta(s_t, a)$  then
8:     Set  $stopped \leftarrow \text{True}$ 
9:     Execute  $a_t = \arg \max_a I_\phi(s_t)$  and observe  $s_{t+1}$ 
10:  else
11:    Execute  $a_t = \arg \max_a I_\phi(s_t)$  and observe  $s_{t+1}$ 
12:  end if
13: end for

```

C.3.4 I -DPT

We implement a variant of DPT [64] using the inference network. The idea is to act greedily with respect to the posterior distribution I at inference time.

First, we train I using ICPE. Then, at deployment we act with respect to I : in round t we selection 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 $\arg \max_H I(H|D_t) \geq 1 - \delta$.

C.4 Transformer Architecture

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

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

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

The inference network operates as follows:

1. **Embedding Layer:** Each observation vector $m_t = (x_t, a_t)$ is first embedded into a higher-dimensional space of size d_e using a linear transformation followed by a GELU activation: $h_t = \text{GELU}(W_{\text{embed}}m_t + b_{\text{embed}})$, $h_t \in \mathbb{R}^{d_e}$.

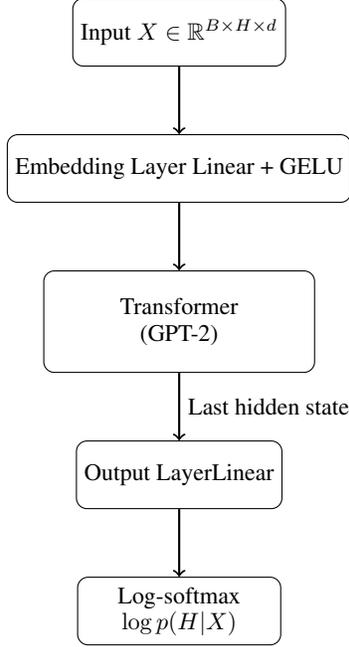


Figure 7: Model architecture for the inference network I (similarly for Q).

2. **Transformer Layers:** The embedded sequence $h \in \mathbb{R}^{B \times H \times d_e}$ is then passed through multiple Transformer layers (specifically, a GPT-2 model configuration). The Transformer computes self-attention over the embedded input to model dependencies among observations:

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

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

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

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

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

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

D Experiments

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

For more informations about the hyper-parameters, we also refer the reader to the `README.md` file in the code, as well as the training configurations in the `configs/experiments` folder.

Libraries used in the experiments. We set up our experiments using Python 3.10.12 [112] (For more information, please refer to the following link <http://www.python.org>), and made use of the following libraries: NumPy [44], SciPy [114], PyTorch [89], Seaborn [116], Pandas [71], Matplotlib [46], CVXPY [30], Wandb [13], Gurobi [43]. Changes, and new code, are published under the MIT license. To run the code, please, read the attached README file for instructions.

D.1 Bandit Problems

Here, we provide the implementation and evaluation details for the bandit experiments reported in Section 3.1, covering deterministic, stochastic, and structured settings. Note that for this setting the observations are simply the observed rewards, i.e., $x_t = r_t$.

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

D.1.1 Deterministic Bandits with Fixed Horizon

Each environment consists of K arms, where $K \in \{4, 6, 8, \dots, 20\}$. Mean rewards for each arm are sampled uniformly from $[0, 1]$, and rewards are deterministic (i.e., zero variance). Agents interact with the environment for exactly K steps and are then required to predict the optimal arm. Success is measured by the probability of correctly identifying the best arm. We also compute the average number of unique arms selected during training episodes as a proxy for exploration diversity.

ICPE is compared against three baselines in the deterministic setting: *Uniform Sampling*, which selects arms uniformly at random; *DQN*, a deep Q-network trained directly on environmental rewards [72]; and *I-DPT*, which performs posterior sampling using ICPE’s *I*-network [64]. All methods were evaluated over five seeds, with 900 environments per seed. 95% confidence intervals were computed with hierarchical bootstrapping.

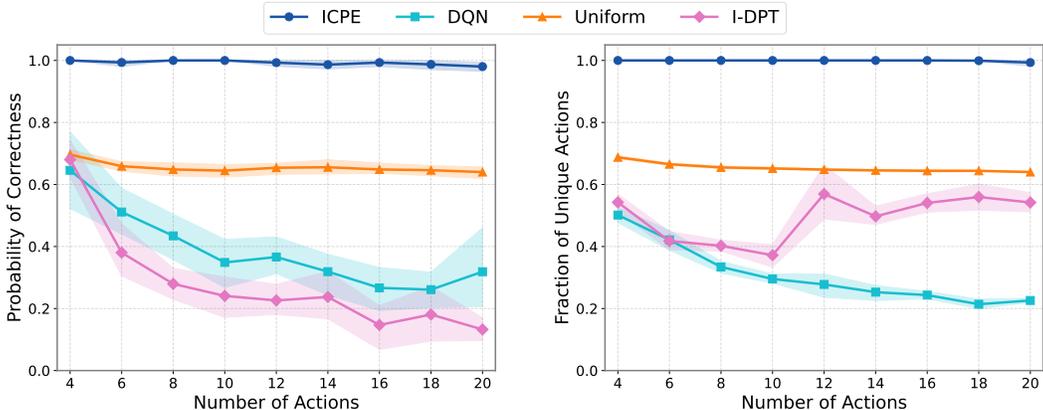


Figure 8: Deterministic bandits: (left) probability of correctly identifying the best action vs. K ; (right) average fraction of unique actions selected during exploration vs. K .

D.1.2 Stochastic Bandits Problems

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

We compare **ICPE** against several widely used baselines: *Top-Two Probability Sampling (TTPS)* [50], *Track-and-Stop (TaS)* [36], *Uniform Sampling*, and *I-DPT*. For *I-DPT*, stopping occurs when the predicted optimal arm has estimated confidence at least $1 - \delta$. For *TTPS* and *TaS*, we apply the classical stopping rule based on the characteristic time $T^*(N_t/t; \hat{\mu}_t)$ (explained in appendix C.3.1):

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

Each method is evaluated over three seeds, with 300 environments, and 15 trajectories per environment. 95% confidence intervals were computed with hierarchical bootstrapping.

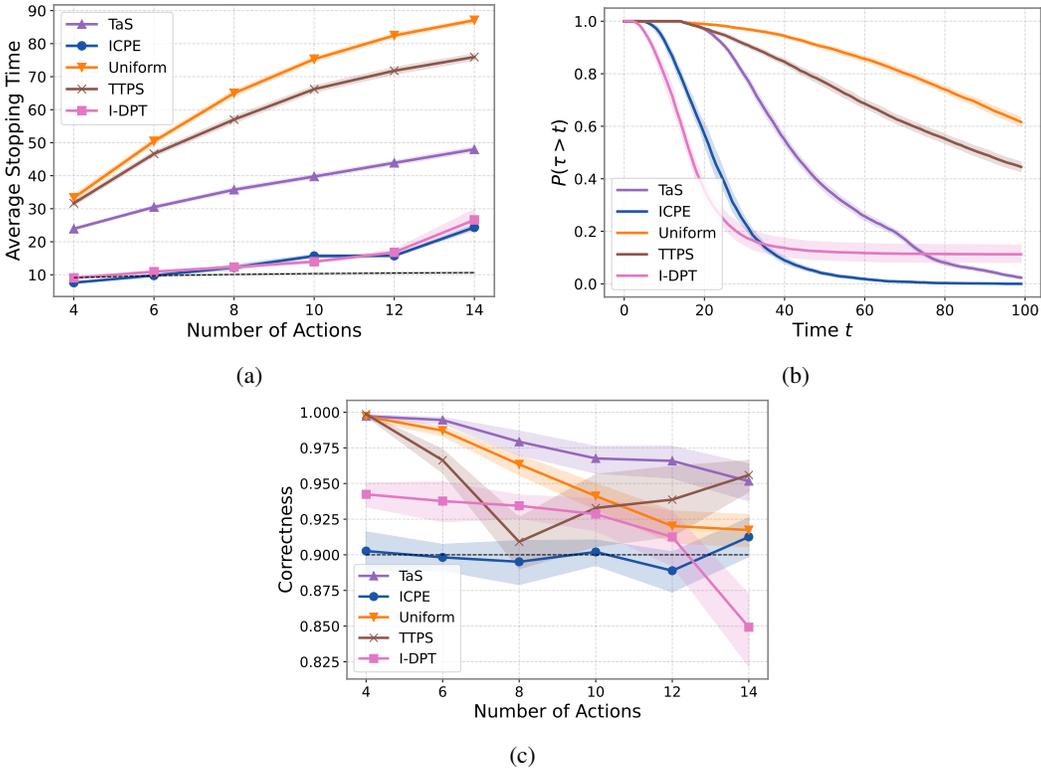


Figure 9: Results for stochastic MABs with fixed confidence $\delta = 0.1$ and $N = 100$: (a) average stopping time τ ; (b) survival function of τ ; (c) probability of correctness $\mathbb{P}_{M \sim \mathcal{P}(M)}^\pi(h(\hat{H}_\tau; M) = 1)$.

Does ICPE learn randomized policies? An intriguing question is whether **ICPE** is capable of learning randomized policies. Intuitively, one might expect randomized methods, such as actor-critic algorithms, to perform better. However, we observe that this is not the case for **ICPE**. Crucially, the inherent randomness of the environment, when passed as input to the transformer architecture, already serves as a source of stochasticity. Thus, although **ICPE** employs a deterministic mapping (via DQN) from observed trajectories, these trajectories themselves constitute random variables, rendering the policy’s output effectively stochastic. To illustrate this, we examine an **ICPE** policy trained with fixed confidence ($\delta = 0.1$) in a setting with $K = 14$ actions (see the two rightmost plots in fig. 10). By analyzing 100 trajectories from this environment and computing an averaged policy, we clearly observe how trajectory randomness influences the policy’s outputs. Specifically, exploration intensity peaks around the middle of the horizon and diminishes as the confidence level increases.

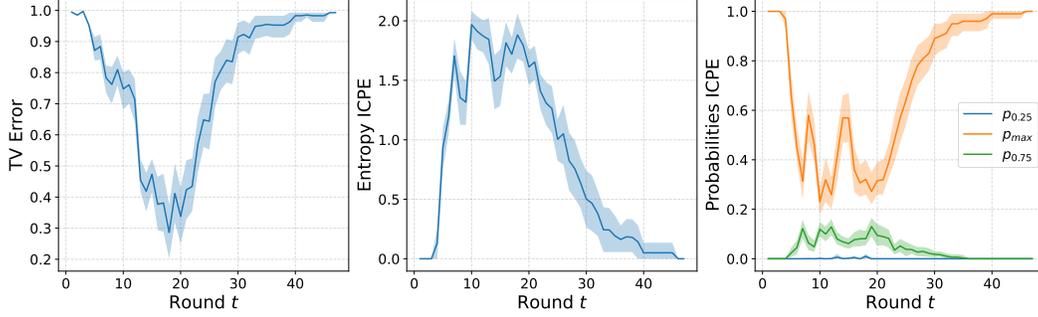


Figure 10: Statistics of **ICPE** with fixed confidence on 100 trajectories from a single environment, with $K = 14$. From left to right: Total variation error between the average **ICPE** policy and the approximate Track and Stop policy; entropy of the average policy of **ICPE**; probabilities of the average **ICPE** policy, with p_{max} representing the maximum probability and p_α the α -quantile.

Does ICPE resembles Track and Stop? In fig. 10 (left figure) we compare an **ICPE** policy trained in the fixed confidence setting ($\delta = 0.1$) with an almost optimal version of TaS, that can be easily computed without solving any optimization problem. Let $\hat{\Delta}_t(a) = \hat{\mu}_{\hat{a}_t}(t) - \max_{a \neq \hat{a}_t} \hat{\mu}_a$, where $\hat{\mu}_a(t)$ is the empirical reward of arm a in round t and $\hat{a}_t = \arg \max_a \hat{\mu}_a(t)$ is the estimated optimal arm. Then, the approximate TaS policy is defined as

$$\pi_t(a) = \frac{1/\hat{\Delta}_a(t)}{\sum_b 1/\hat{\Delta}_b(t)},$$

with $\hat{\Delta}_{\hat{a}_t}(t) = \min_{a \neq \hat{a}_t} \hat{\Delta}_a(t)$. In the figure we sampled 100 trajectories from a single environment with $K = 14$, and computed an average **ICPE** policy. Then, we compared this policy to the approximate TaS policy, and computed the total variation. We can see that the two policies are not always similar. We believe this is due to the fact that **ICPE** is exploiting prior information on the environment, including the minimum gap assumption, and the fact that the average rewards are bounded in $[0, 0.4K]$.

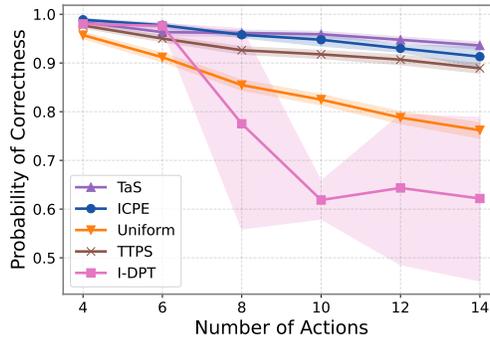


Figure 11: Correctness $\mathbb{P}_{M \sim \mathcal{P}(\mathcal{M})}^\pi(h(\hat{H}_\tau; M) = 1)$ for stochastic MABs with fixed horizon $N = 30$.

D.1.3 Bandit Problems with Hidden Information

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

Each environment contains $K = 5$ arms. Non-magic arms have mean rewards sampled uniformly from $[1, 5]$, while the mean reward of the designated *magic action* (always arm 1) is defined as $\mu_m = \phi(\arg \max_{a \neq a_m} \mu_a)$ with $\phi(i) = i/K$. The magic action is not the optimal arm, but it encodes information about which of the other arms is. To control the informativeness of this signal,

we vary the standard deviation of the magic arm $\sigma_m \in \{0.0, 0.1, \dots, 1.0\}$, while fixing the standard deviation of all other arms to $\sigma = 1 - \sigma_m$.

ICPE is trained under the fixed-confidence setting with a target confidence level of 0.9. For each σ_m , we compare **ICPE**'s sample complexity to two baselines: (1) the average theoretical lower bound computed for the problem computed via averaging the result of Theorem B.7 over numerous environmental mean rewards, and (2) *I-IDS*, a pure-exploration information-directed sampling algorithm that uses **ICPE**'s *I*-network for posterior estimation. All methods are over 500 environments, with 10 trajectories per environment. 95% confidence intervals are computed using hierarchical bootstrapping with two levels.

Beyond the exploration efficiency analysis shown in Figure 5a, we also assess the correctness of each method's final prediction and its usage of the magic action. As shown in Figure 12a, both **ICPE** and *I-IDS* consistently achieve the target accuracy of 0.9, validating their reliability under the fixed-confidence formulation.

Figure 12b plots the proportion of total actions that were allocated to the magic arm across different values of σ_m . While both methods adapt their reliance on the magic action as its informativeness degrades, *I-IDS* tends to abandon it earlier. This behavior suggests that **ICPE** is better able to retain and exploit structured latent information beyond what is captured by simple heuristics for expected information gain.

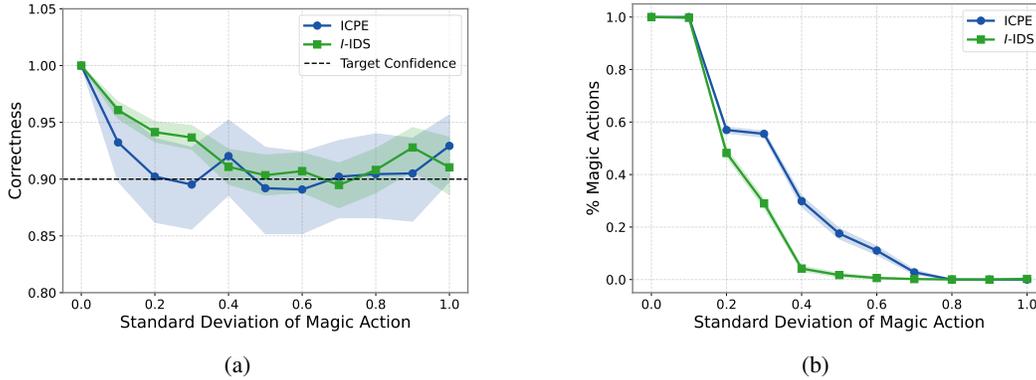


Figure 12: (a) Final prediction accuracy across varying levels of noise in the magic action (σ_m). Both **ICPE** and *I-IDS* consistently achieve the target confidence threshold of 0.9. (b) Percentage of actions allocated to the magic arm as a function of σ_m . **ICPE** continues to exploit the magic action longer than *I-IDS*, suggesting more robust use of latent structure.

We also assess **ICPE**'s performance in a regret minimization setting. We define an *In-Context Explore-then-Commit* variant of **ICPE**, which explores until the *I*-network reaches confidence $1 - \delta$, then repeatedly selects the estimated optimal action. We compare this policy's cumulative regret to that of three standard algorithms: *UCB*, *Thompson Sampling*, and *IDS*, each initialized with Gaussian priors. For this evaluation, we fix $\sigma_m = 0.1$, $\sigma = 0.9$, and $\delta = 0.01$.

Implementation details for *I-IDS* and *In-Context Explore-then-Commit* are provided in Sections C.3.2 and C.3.3 respectively.

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

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

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

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

ICPE is trained under the fixed-confidence setting with $\delta = 0.99$. For each n , five models are trained across five seeds. We compare ICPE’s average stopping time to the theoretical optimum (computed via Theorem B.10) and to the *I-IDS* baseline equipped with access to the *I*-network. Each model is evaluated over 1000 test environments per seed. 95% confidence intervals are computed using hierarchical bootstrapping.

In interpreting the results from Figure 5b, we observe that for environments with one or two magic actions, ICPE reliably learns the optimal policy of following the magic chain to completion. In these cases, the agent is able to identify the optimal arm without ever directly sampling it, by exploiting the structured dependencies embedded in the reward signals of the magic actions. Figure 13 illustrates a representative trajectory from the two-magic-arm setting, where the first magic action reveals the location of the second, which in turn identifies the optimal arm. The episode terminates without 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 13: 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.

For environments with more than two magic actions, however, ICPE learns a different strategy. As the length of the magic chain increases, the expected sample complexity of following the chain from the start becomes suboptimal. Instead, ICPE learns to randomly sample actions until it encounters one of the magic arms and then proceeds to follow the chain from that point onward. This behavior represents an efficient, learned compromise between exploration cost and informativeness. Figure 14 shows an example trajectory from the six-magic-arm setting, where the agent initiates random sampling until it 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 14: 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.

D.2 Semi-Synthetic Pixel Sampling

To evaluate **ICPE** in a setting that more closely resembles real-world decision-making tasks, we designed a semi-synthetic environment based on the MNIST dataset [63], where the agent must adaptively reveal image regions to classify a digit while minimizing the number of pixels observed. This experiment serves as a proof-of-concept for using **ICPE** in perceptual tasks where observations are costly and information must be acquired efficiently.

Environment Details. Each MNIST image is partitioned into 36 non-overlapping 5×5 pixel regions, defining an action space of size $K = 36$. At each timestep, the agent selects a region to reveal, progressively unmasking the image. The agent begins with a fully masked image and has a fixed horizon of $H = 12$ steps to acquire information and make a prediction.

To prevent overfitting and encourage generalizable policies, we apply strong augmentations at each episode: random rotations ($\pm 30^\circ$), translations (up to 2 pixels), Gaussian noise ($\mathcal{N}(0, 0.1)$), elastic deformations, and random contrast adjustments. These augmentations ensure that agents cannot memorize specific pixel layouts and must instead learn adaptive exploration strategies.

Model Architecture and Optimization. Due to the visual nature of the task, we use a convolutional encoder rather than a transformer. The **ICPE** critic network combines a CNN image encoder with a separate action-count encoder. The CNN consists of 3 convolutional blocks with 16 base channels, followed by max pooling and global average pooling. The action counts (which track how often each region has been sampled) are passed through a linear embedding layer with 32 output units, followed by ReLU activation and LayerNorm. The image and action embeddings are concatenated and processed through two residual MLP layers, producing Q -values over actions. The I -network shares the same architecture but outputs logits over 10 digit classes.

All models are implemented in PyTorch and trained with Adam using a learning rate of 1×10^{-4} . Training is performed over 500,000 episodes using 40 parallel environment instances. We use a batch size of 128, a replay buffer of size 100,000, and a discount factor $\gamma = 0.999$. The Q -network is updated using Polyak averaging with coefficient 0.01, and the I -network is updated every two steps using 30 bootstrap batches. To populate the buffer initially, we perform 10 batches of bootstrap updates before standard training begins. Gradients are clipped to a maximum norm of 2.

Pretraining the Inference Network. To provide stable reward signals and ensure consistency with baselines, we pretrain a separate CNN classifier to predict digit labels from fully revealed images. This classifier consists of three convolutional layers with max pooling, followed by two linear layers and a softmax head. The classifier is trained on the same augmented data used during **ICPE** training and is frozen during exploration learning. Its softmax confidence for the correct digit is used as the reward signal. This setup simulates realistic scenarios in which high-quality predictive models already exist for fully observed data (e.g., in clinical diagnosis).

Evaluation. We compare **ICPE** to two baselines: *Uniform Sampling*, which selects image regions uniformly at random at each timestep, and *Deep CMAB* [24], a contextual bandit algorithm that uses a Bayesian neural network to model $p(r | x, a)$ and performs posterior sampling via dropout.

The Deep CMAB model uses a convolutional encoder to extract image features, which are concatenated with a learned embedding of the action count vector. The combined representation is passed through a multilayer perceptron with dropout applied to each hidden layer. At each decision point, the agent samples a dropout rate from a uniform distribution over $(0, 1)$ and uses the resulting forward pass as a sample from the posterior over rewards (Thompson sampling). The reward signal for each action is computed using the pretrained MNIST classifier: specifically, the increase in softmax probability for the correct digit class after a new region is revealed.

We train Deep CMAB for 100,000 episodes using Adam optimization. During training, the agent interacts with multiple MNIST instances in parallel, and updates its model based on the marginal improvement in confidence after each action. The model learns to maximize this incremental reward signal by associating particular visual contexts with the most informative actions.

For each trained model, we sample 1000 test environments and report on (1) the average final classification accuracy by the pretrained classifier at the end of trajectory, and (2) the average number of regions used before prediction. Confidence intervals are computed via bootstrapping.

Adaptive Sampling Analysis. To assess whether agents learn digit-specific exploration strategies, we analyze the distribution of selected image regions across digit classes. For each agent, we compute pairwise chi-squared tests between all digit pairs, testing whether the distributions of selected regions are statistically distinguishable.

To ensure sufficient support for the test, we only compare digit pairs that each have at least five trajectories and remove actions that appear in fewer than five total samples across the two classes. For each qualifying digit pair, we construct a $2 \times \tilde{K}$ contingency table, where \tilde{K} is the number of region indices that are meaningfully used by either digit. The rows correspond to digit classes, and each column counts how many samples from each class selected the corresponding region at least once.

We apply the chi-squared test of independence to each contingency table. A low p-value indicates that the region selection distributions for the two digits are significantly different, suggesting digit-specific adaptation. By comparing the number and strength of significant differences across agents, we evaluate the extent to which each method tailors its exploration policy to the structure of the input class.

We visualize the resulting pairwise p-values in Figure 15 using a heatmap. Each cell shows the chi-squared test p-value between a pair of digits. Lower values (blue cells) indicate greater divergence in sampling behavior, and thus more adaptive and digit-specific strategies.

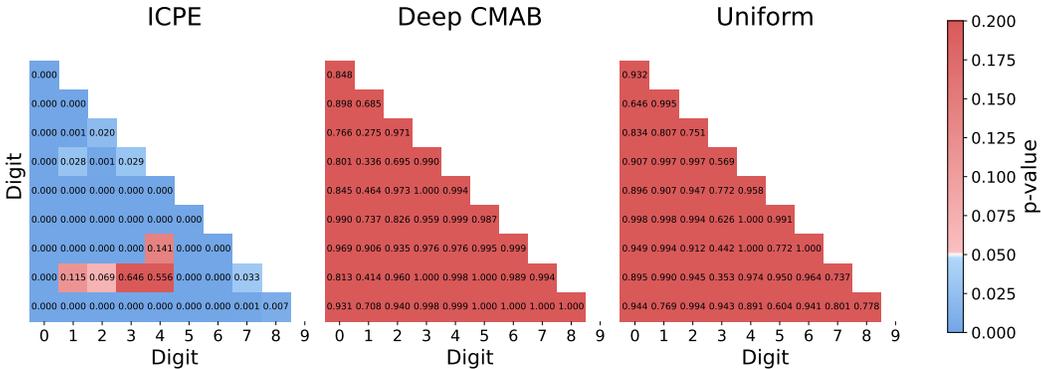


Figure 15: Pairwise chi-squared test p-values for region selection distributions across digit classes. Lower values indicate more statistically distinct exploration behaviors.

For further intuition into the sampling process, Figure 16 shows a representative example of the ICPE pipeline progressively revealing image regions and correctly classifying the digit ‘2’. This highlights the interplay between exploration and inference as the agent strategically uncovers informative regions to guide its decision.

To illustrate the impact of input corruption, Figure 17 presents an example where ICPE fails to correctly classify the digit. Although the agent successfully reveals the central body of the digit, the applied augmentations distort the image to the extent that the digit becomes visually ambiguous. In this case, the agent incorrectly predicts an ‘8’ when the true label is a ‘9’, underscoring the challenge introduced by realistic image corruptions in this setting.

D.3 MDP Problems: Magic Room

The Magic Room is a sequential decision-making environment structured as a $K \times K$ grid-shaped room containing four doors, each positioned at the midpoint of one of the four walls (top, bottom, left, right). At the beginning of each episode, exactly one of these doors is randomly chosen to be the correct door (H^*), unknown to the agent.

The agent’s goal is to identify and pass through the correct door. Each episode lasts for a maximum of $N = K^2$ time steps, during which the agent navigates the grid, observes clues, and attempts

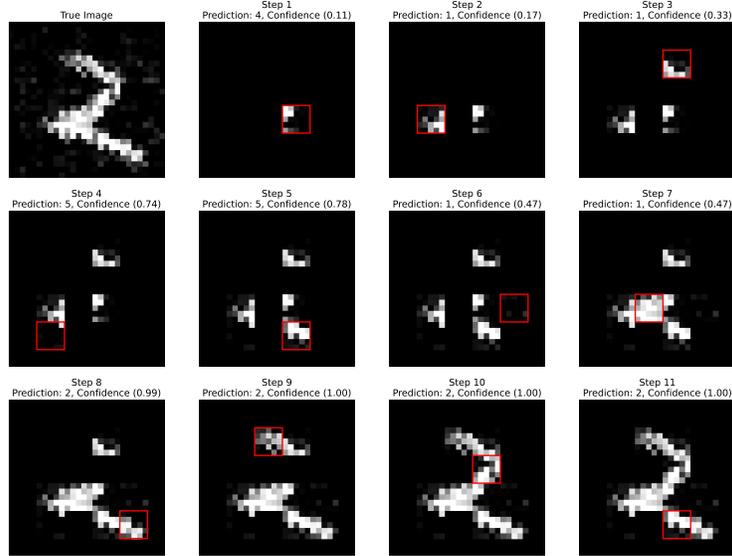


Figure 16: Illustrative example of the ICPE agent revealing regions of an MNIST digit and correctly classifying it as a ‘2’. The sequence shows the intermediate revealed image and predicted label at each timestep.

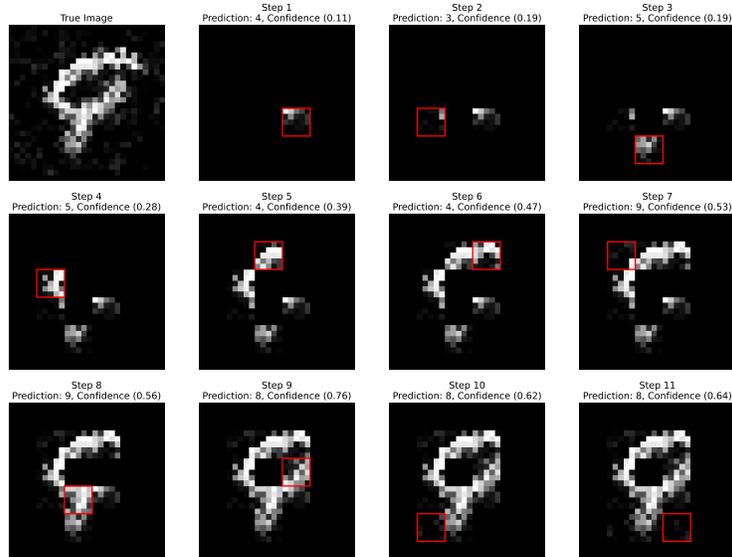


Figure 17: Example of an incorrect classification due to aggressive data augmentations. Although the agent reveals the central region of the digit, the distortions cause it to misclassify a ‘9’ as an ‘8’.

to determine the correct door. Two binary clues, each randomly assigned a location within the sub-grid $[1, 1] \times [K - 1, K - 1]$, are placed in the room at the start of each episode. Each clue has a binary value, randomly set to either -1 or 1 . Collecting both clues provides sufficient information to unambiguously determine the correct door, given that the agent has learned the mapping from clue configurations to door identity.

At each time step t , the agent observes the state vector:

$$x_t = (z_t, y_t, c_{1,t}, c_{2,t}, r_t),$$

where:

- (z_t, y_t) are the agent’s current coordinates on the grid.

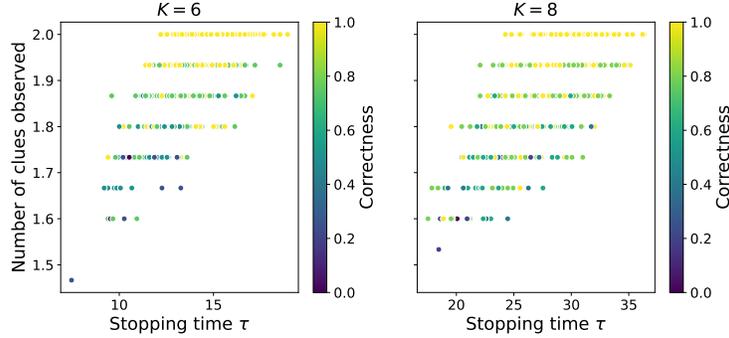


Figure 20: Magic room environment. Relationship among agent correctness, the number of clues observed, and the stopping time.

We trained **ICPE** on 3 seeds, using the fixed confidence setting (disabling the stopping action) using $\delta = 0.05$ and evaluated the policies on 4500 episodes for $K = 6$ and $K = 8$. In table 1 are shown the statistics of the average correctness and of the stopping time.

In fig. 18 we can see a sample trajectory taken by **ICPE**. Starting from the middle of the room, **ICPE** follows a path that allows to find the clues C_1, C_2 in the green area. As soon as the second clue is found, it goes through the closest door.

In fig. 19, we present the survival functions of the stopping time τ for environments with grid sizes $K = 6, 8$, alongside the corresponding correctness densities. Lastly, fig. 20 illustrates the relationship among agent correctness, the number of clues observed, and the stopping time. Specifically, smaller stopping times correlate with fewer observed clues, leading to lower correctness. Conversely, when the agent observes both clues, it consistently selects the correct door, demonstrating that it has effectively learned the association between the clues and the correct hypothesis.

D.4 Exploration on Feedback Graphs

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

While feedback graphs have been widely studied for regret minimization [68], their use in pure exploration remains relatively underexplored [98]. We study them here as a challenging and structured testbed for in-context exploration. Unlike unstructured bandits, these environments contain latent relational structure and stochastic feedback dependencies that must be inferred and exploited to explore efficiently.

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

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

This setting allows us to test whether **ICPE** can learn to uncover and leverage latent graph structure to guide exploration. As in the bandits setting, we have a finite number of actions $\mathcal{A} = \{1, \dots, K\}$, corresponding to the actions (or vertices) in a feedback graph G . The learner’s goal is to identify the best action, where $H^* = \arg \max_a \mu_a$. At each time step t , the observation is the partially observed reward vector $x_t = r_t$.

We evaluate performance on best-arm identification tasks across three representative feedback graph families:

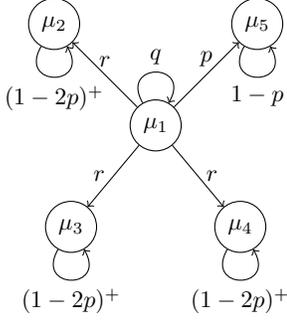


Figure 21: Loopy star graph.

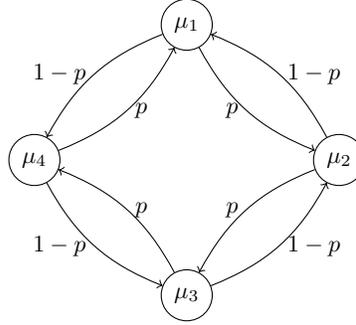


Figure 22: Ring graph.

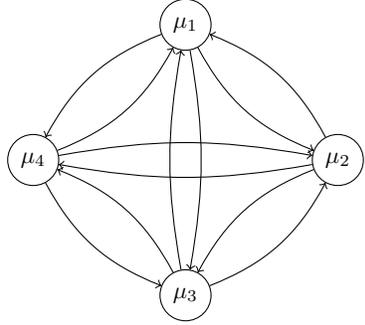


Figure 23: Loopless clique graph.

- **Loopy Star Graph** (Figure 21): A star-shaped graph with self-loops, parameterized by (p, q, r) . The central node observes itself with probability q , one neighboring node with probability p , and all others with probability r . When p is small, it may be suboptimal to pull the central node, requiring the agent to adapt its strategy accordingly.
- **Ring Graph** (Figure 22): A cyclic graph where each node observes its right neighbor with probability p and its left neighbor with probability $1 - p$. Effective exploration requires reasoning about which neighbors provide more informative feedback.
- **Loopless Clique Graph** (Figure 23): A fully connected graph with no self-loops. Edge probabilities are defined as:

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

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

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

Fixed-Horizon. For each graph family, mean rewards were sampled uniformly from $[0, 1]$ with fixed variance 0.2, using hyperparameters: $(p, q, r) = (0.25, 0.3, 0.35)$ for the loopy star graph, $p = 0.3$ for the ring, and $p = 0.5$ for the loopless clique. We considered both small ($K = 5, H = 25$) and large ($K = 10, H = 50$) environments.

ICPE was compared to three baselines: Uniform Sampling, EXP3.G [94], and Tas-FG [98]. All methods performed maximum likelihood inference at the end of the trajectory. Table 2 reports the average probability of correctly identifying the best arm.

Algorithm	Loopy Star		Loopless Clique		Ring	
	Small	Large	Small	Large	Small	Large
ICPE	0.88 ± 0.01	0.59 ± 0.02	0.95 ± 0.01	0.79 ± 0.04	0.79 ± 0.01	0.51 ± 0.03
TasFG	0.82 ± 0.01	0.73 ± 0.02	0.84 ± 0.01	0.83 ± 0.01	0.70 ± 0.02	0.56 ± 0.02
EXP3.G	0.66 ± 0.02	0.40 ± 0.01	0.84 ± 0.01	0.78 ± 0.02	0.77 ± 0.02	0.52 ± 0.02
Uniform	0.73 ± 0.02	0.60 ± 0.02	0.86 ± 0.01	0.79 ± 0.02	0.78 ± 0.02	0.62 ± 0.02

Table 2: Probability of correctly identifying the best arm. Small environments: $K = 5, H = 25$; Large: $K = 10, H = 50$. Results reported as mean ± 95% CI.

ICPE outperforms all baselines in small environments across all graph families, highlighting its ability to learn efficient strategies from experience. Performance slightly degrades in larger environments, likely due to difficulty in credit assignment over long horizons. Still, **ICPE** remains competitive, validating its capacity to generalize across graph-structured settings.

Fixed-Confidence. We next tested **ICPE** in a fixed-confidence setting, using the same graph families but setting the optimal arm’s mean to 1 and all others to 0.5 to facilitate faster convergence. **ICPE** was trained for $K = 4, 6, \dots, 14$ with a target error rate of $\delta = 0.1$. We compared it to Uniform Sampling, EXP3.G, and Tas-FG using a shared stopping rule from [98].

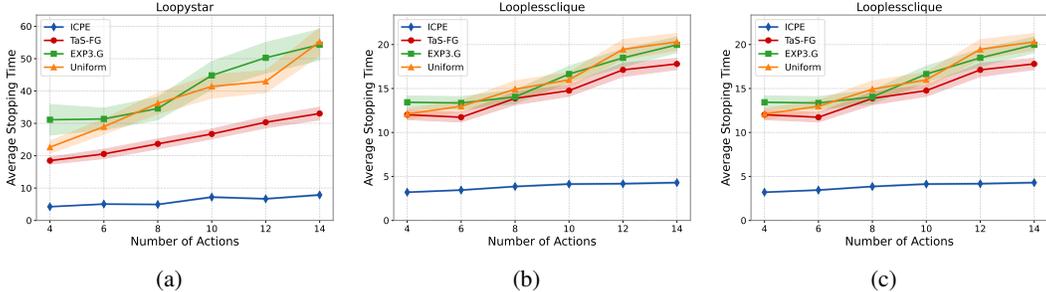


Figure 24: Sample complexity comparison under the fixed-confidence setting for: (a) Loopy Star, (b) Loopless Clique, and (c) Ring graphs.

As shown in Figure 24, **ICPE** consistently achieves significantly lower sample complexity than all baselines. This suggests that **ICPE** is able to meta-learn the underlying structure of the feedback graphs and leverage this knowledge to explore more efficiently than *uninformed* strategies. These results align with expectations: when environments share latent structure, learning to explore from experience offers a substantial advantage over fixed heuristics that cannot adapt across tasks.

D.5 Meta-Learning Binary Search

To test **ICPE**’s ability to recover classical exploration algorithms, we evaluate whether it can autonomously meta-learn binary search.

We define an action space of $\mathcal{A} = \{1, \dots, K\}$, where K is the upper bound on the possible location of the hidden target $H^* \sim \mathcal{A}$. Pulling an arm above or below H^* yields a observation $x_t = -1$ or $x_t = +1$, respectively—providing directional feedback.

We train **ICPE** under the fixed-confidence setting for $K = 2^3, \dots, 2^8$, using 150,000 in-context episodes and a target error rate of $\delta = 0.01$. Evaluation was conducted on 100 held-out tasks per setting. We report the minimum accuracy, mean stopping time, and worst-case stopping time, and 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 3: **ICPE** performance on the binary search task as the number of actions K increases.

As shown in Table 3, **ICPE** consistently achieves perfect accuracy with worst-case stopping times that match the optimal $\log_2(K)$ rate, demonstrating that it has successfully rediscovered binary search purely from experience. While simple, this task illustrates **ICPE**’s broader potential to learn efficient search strategies in domains where no hand-designed algorithm is available.