

---

# In-Context Pure Exploration in Continuous Decision Spaces

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 In active sequential testing, also termed *pure exploration*, a learner is tasked with the  
2 goal to adaptively acquire information so as to identify an unknown ground-truth  
3 hypothesis with as few queries as possible. This problem has several motivating  
4 applications, including Best-Arm Identification (BAI) in bandits, where actions  
5 index hypotheses, and generalized search problems, where strategically chosen  
6 queries reveal partial information about a hidden label. In many modern settings,  
7 however, the hypothesis, or recommendation space, is *continuous*: for example,  
8 identifying a near optimal action in a continuous-armed bandit, localizing an  $\epsilon$ -  
9 ball contained in a target region, or estimating the minimizer of a function from  
10 noisy observations. Existing methods are predominantly frequentist and model-  
11 specific, while learned approaches have been limited to finite recommendation  
12 spaces. We introduce C-ICPE, a theory-guided learned model for Bayesian fixed-  
13 confidence pure exploration with continuous recommendations. C-ICPE meta-  
14 trains sequential architectures over a task prior to jointly learn exploration, stopping  
15 and recommendations strategies. At inference time, it actively gathers evidence on  
16 tasks and identifies an  $\epsilon$ -optimal recommendation *without* parameter updates.

## 17 1 Introduction

18 Several learning problems are inherently interactive: the learner  
19 sequentially performs interventions or stages queries, observes  
20 noisy evidence whose distribution depends on that intervention,  
21 and stops once the accumulated evidence supports a reliable  
22 conclusion. This type of interactive sequential decision-making  
23 problem goes back to Chernoff [12] and has been formalized  
24 through active sequential hypothesis testing [37] and pure ex-  
25 ploration with fixed confidence [6, 14], where the learner min-  
26 imizes the number of queries subject to returning an  $\epsilon$ -accurate  
27 recommendation with probability at least  $1 - \delta$ .

28 This regime is well understood in canonical settings with finite  
29 decision spaces, including best-arm identification in stochastic  
30 multi-armed bandit models [19] and best-policy identification  
31 in Markov Decision Processes (MDPs) [44]. In these settings  
32 the learner chooses queries (e.g., arms) and outputs an object  
33 of interest, often a best action/policy, and the theoretical guar-  
34 antees have been studied in a range of settings [14, 43, 3, 52].

35 Despite this progress, practical methods for fixed-confidence  
36 pure exploration remain limited when the learner must return a recommendation in a *continuous*

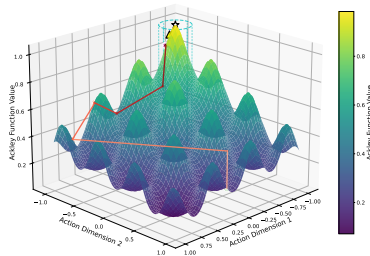


Figure 1: C-ICPE is able to identify the global maxima (with  $\epsilon$ -accuracy and  $1 - \delta$  confidence) of the inverted Ackley function (with random parameters and observation noise) without gradients while trying to use the least number of data-points.

37 space. Existing continuous methods are either frequentist and model-specific [20, 58, 42, 48] or  
 38 Bayesian but restricted to Gaussian processes and not optimizing sample complexity [65]. Even in  
 39 finite arms, the theory of Bayesian fixed-confidence pure exploration is limited, and results are known  
 40 only in the finite setting with Gaussian likelihoods/priors [29]. No analogous Bayesian theory, nor  
 41 practical methods, are known for continuous recommendation spaces under general priors.

42 Recently, [53] proposed In-Context Pure Exploration (ICPE), which meta-trains sequential neural  
 43 policies for finite active-testing problems. However, ICPE is restricted to finite hypothesis and  
 44 action sets, and does not address the continuous recommendation setting. Hence, it is currently  
 45 missing a broadly reusable learned method for Bayesian fixed-confidence pure exploration when the  
 46 recommendation itself is continuous and performance is optimized over a task prior. We introduce  
 47 C-ICPE, a theory-guided method that learns to collect data, stop, and recommend directly from  
 48 trajectories in continuous recommendation spaces. This type of  $(\epsilon, \delta)$ -PAC exploration directly  
 49 addresses problems in experimental sciences such as materials discovery [36] and dose-finding [41],  
 50 where each trial is costly, observations are noisy, and practitioners need not just a good answer but a  
 51 guarantee that the answer is  $\epsilon$ -correct with a given confidence.

52 **Contributions.** First, we formulate Bayesian fixed-confidence pure exploration with continuous  
 53 recommendations, and establish the corresponding Bellman optimality structure. Second, we prove  
 54 Bayesian  $(\epsilon, \delta)$ -correctness under a local closedness condition that is weaker than the uniqueness  
 55 assumption in [53], using a novel subdifferential argument. Third, we instantiate this framework into  
 56 a practical method, C-ICPE, to train exploration policies that deploy model-free, and evaluate it on  
 57 noisy binary search,  $\epsilon$ -best arm identification on the unit sphere, noisy Ackley minimization, value  
 58 estimation in Gaussian Processes and a real-world geochemical task where the goal is to locate peak  
 59 copper concentration from sparse soil measurements [1]. To our knowledge, this is the first practical  
 60 learned framework combining continuous recommendations and fixed-confidence stopping.

## 61 2 Problem Setting

62 We consider a Bayesian family of active sequential decision problems indexed by a latent parameter  
 63  $\theta \in \Theta$ , where  $\Theta \subset \mathbb{R}^d$  is compact and  $\theta \sim \nu$ . Each environment  $M_\theta$  specifies an initial observation  
 64 law  $\rho_\theta \in \Delta(\mathcal{Y})$  and observation kernels  $P_{\theta,t}(\cdot|h_t, a_t)$  over a compact observation space  $\mathcal{Y} \subset \mathbb{R}^n$ .

65 In sequential testing, the learner interacts with  $M_\theta$  over time: in round  $t$  it observes the history

$$H_t := (Y_1, A_1, Y_2, \dots, A_{t-1}, Y_t),$$

66 chooses a query  $A_t \in \mathcal{A} \subset \mathbb{R}^m$  (compact), and observes  $Y_{t+1} \sim P_{\theta,t}(\cdot|H_t, A_t)$ . The goal is to  
 67 collect a history that is sufficiently informative to output a high-quality recommendation  $\hat{x} \in \mathcal{X}$ ,  
 68 where  $\mathcal{X}$  is compact. We refer to  $\mathcal{X}$  as the hypothesis or recommendation space. We distinguish  
 69  $\mathcal{A}$  from  $\mathcal{X}$ :  $\mathcal{A}$  is the query space used to collect information, while  $\mathcal{X}$  is the space of objects the  
 70 learner may return. In many tasks  $\mathcal{X} = \mathcal{A}$ , but in value-identification tasks  $\mathcal{X}$  may instead be a set of  
 71 possible function values.

72 **Risk function.** Recommendation quality is measured by a task-dependent loss, or risk, function  
 73  $L_\theta : \mathcal{X} \rightarrow [0, \infty)$ , satisfying  $\inf_{x \in \mathcal{X}} L_\theta(x) = 0$ . We say that  $x$  is  $\epsilon$ -optimal for task  $\theta$  if  $L_\theta(x) \leq \epsilon$ ,  
 74 and define

$$\mathcal{X}_\epsilon(\theta) := \{x \in \mathcal{X} : L_\theta(x) \leq \epsilon\}.$$

75 In the following, we assume that  $(\theta, x) \mapsto L_\theta(x)$  is jointly lower semicontinuous (in Section B we  
 76 state the regularity assumptions used in the theoretical results). Throughout the paper,  $x_\theta^*$  denotes  
 77 a selected zero-loss target in the recommendation space  $\mathcal{X}$ , i.e.,  $L_\theta(x_\theta^*) = 0$ . Depending on the  
 78 problem, this object may be an optimizer, a best arm, a threshold, or an optimal value. When the  
 79 zero-loss set is not a singleton, we assume a fixed continuous selector  $\theta \mapsto x_\theta^*$  to ensure regularity.

80 In many examples, the loss is defined through a function  $f_\theta$  parametrized by  $\theta$ . In function optimiza-  
 81 tion problems, we set  $\mathcal{X} = \mathcal{A}$ , and the goal is to find a point  $\hat{x} \in \mathcal{X}$  that optimizes the function. In this  
 82 case, one can take the risk loss  $L_\theta$  to be value-gap loss  $L_\theta(x) := f_\theta(x) - f_\theta(x_\theta^*)$ , with  $x_\theta^* \in F^*(\theta)$ ,  
 83 where  $F^*(\theta) := \arg \min_{x \in \mathcal{X}} f_\theta(x)$ , or a distance loss in the query space  $L_\theta(x) := \|x_\theta^* - x\|$  if the  
 84 goal is to find  $x$  close to a selected optimal point  $x_\theta^* \in F^*(\theta)$ . Other problems include the  $\epsilon$ -best arm  
 85 identification problem in multi-armed bandits where  $f_\theta(x)$  is the mean reward of arm  $x$ , or noisy  
 86 binary search in a continuum, where the agent observes noisy observations of  $\text{sign}(x - x_\theta^*)$  and the  
 87 loss is defined in the query space. Problems where the recommendation space  $\mathcal{X}$  is not identical to

88 the query space  $\mathcal{A}$  include optimal-value identification, where the learner returns a scalar estimate of  
 89 the optimal value: we take  $\mathcal{X} \subset \mathbb{R}$ , and set  $L_\theta(x) := |x - x_\theta^*|$  where  $x_\theta^* := \max_{a \in \mathcal{A}} f_\theta(a)$ .

90 **Optimization objective.** We work in the fixed-confidence  $((\epsilon, \delta)$ -PAC) regime. A learner is defined  
 91 by the triplet  $(\pi, I, \tau)$ : a sampling policy  $\pi = (\pi_t)_{t \geq 1}$  such that  $A_t = \pi_t(H_t)$ ; a stopping time  
 92  $\tau$  with respect to  $\mathcal{F}_t = \sigma(H_t)$ , defining when to stop the data acquisition process; an inference  
 93 rule  $I = (I_t)_{t \geq 1}$  such that  $\hat{x}_\tau = I_\tau(H_\tau)$ . The goal of the learner is to adaptively choose queries  
 94  $A_1, A_2, \dots$  and a stopping time  $\tau$ , so that the returned  $\hat{x}_\tau$  is  $\epsilon$ -optimal, i.e.  $\hat{x}_\tau \in \mathcal{X}_\epsilon(\theta)$ , with high  
 95 probability. Hence, for a given pair  $\epsilon > 0, \delta \in (0, 1/2)$ , we seek to minimize the (expected) number of  
 96 queries while ensuring  $\delta$ -correctness: formally, we solve

$$\inf_{\tau, \pi, I} \mathbb{E}_{\theta \sim \nu}^\pi [\tau] \quad \text{s.t.} \quad \mathbb{P}_{\theta \sim \nu}^\pi (\hat{x}_\tau \in \mathcal{X}_\epsilon(\theta)) \geq 1 - \delta, \quad \mathbb{E}_{\theta \sim \nu}^\pi [\tau] < \infty. \quad (1)$$

97 Our formulation is Bayesian:  $\nu$  is both a prior over environments and the task distribution used  
 98 for training and evaluation. This enables amortized learning across tasks: the models are trained  
 99 on episodes drawn from  $\nu$  and transfer to new tasks from the same family without parameter  
 100 updates. Second, it defines the posterior success probability  $q_t(h, x)$  that drives our theory and  
 101 algorithms. Third, the average-case guarantee under  $\nu$  is the natural objective for applications where  
 102 the practitioner has domain knowledge about plausible task distributions. For this setup, we are not  
 103 aware of analogous Bayes-optimal characterizations for continuous recommendation spaces under  
 104 general priors, as results are limited to finite settings with Gaussian structure [29].

### 105 3 Theoretical Background

106 This section provides the theoretical foundation for C-ICPE, where we characterize an optimal learner.  
 107 Relative to the finite ICPE analysis of Russo et al. [53], the continuous setting introduces two technical  
 108 issues absent in finite spaces: (i) attainment of the Bellman equation over continuous actions requires  
 109 establishing a semicontinuity chain from the observation model through the posterior predictive to  
 110 the  $Q$ -function, and (ii) the  $(\epsilon, \delta)$ -correctness proof must handle non-singleton dual optima, which  
 111 we address via a weaker local closedness condition and a subdifferential argument that replaces the  
 112 uniqueness and monotonicity assumptions of the finite case. Together, these results provide the first  
 113 theoretical infrastructure for Bayesian pure exploration with continuous recommendations.

114 **Posterior success and optimal inference.** In the fixed-confidence setting with continuous  $\mathcal{X}$ , the  
 115 relevant posterior object is the posterior probability that  $x$  is  $\epsilon$ -optimal:

$$q_t(h, x) := \mathbb{P}(L_\theta(x) \leq \epsilon \mid H_t = h), \quad r_t(h) := \max_{x \in \mathcal{X}} q_t(h, x).$$

116 Here  $q_t(h, x)$  is the posterior success probability of recommending  $x$ , and  $r_t(h)$  is the best posterior  
 117 success probability achievable from history  $h$ . One can show that the maximum is attained and an  
 118 optimal inference rule is any selector (see Proposition 1 for a proof)

$$I_t^*(h) \in \arg \max_{x \in \mathcal{X}} q_t(h, x).$$

119 **Dual formulation and Bellman optimality.** We study the fixed-confidence problem in Eq. (1)  
 120 through its Lagrangian dual, following the ASHT literature [37, 53]. Introducing a multiplier  $\lambda \geq 0$   
 121 for the correctness constraint gives

$$V_\lambda(\pi, I, \tau) := -\mathbb{E}_{\theta \sim \nu}^\pi [\tau] + \lambda (\mathbb{P}_{\theta \sim \nu}^\pi (I_\tau(H_\tau) \in \mathcal{X}_\epsilon(\theta)) - 1 + \delta), \quad \inf_{\lambda \geq 0} \sup_{\pi, I, \tau} V_\lambda(\pi, I, \tau). \quad (2)$$

122 For fixed  $\pi$  and  $\tau$ , optimizing over  $I$  replaces the terminal success probability by  $\mathbb{E}^\pi [r_\tau(H_\tau)]$ . Hence,  
 123 for fixed  $\lambda$ , the inner problem is equivalent up to a constant to maximizing  $\mathbb{E}^\pi [\lambda r_\tau(H_\tau) - \tau]$ . Stopping  
 124 can also be embedded as an absorbing action  $a_{\text{stop}}$ : the learner continues with actions in  $\mathcal{A}$  until it  
 125 selects  $a_{\text{stop}}$ , at which point it stops and outputs  $I_t^*(H_t)$ ; see Lemma 6. Thus the fixed- $\lambda$  problem is  
 126 an optimal-stopping control problem on  $\bar{\mathcal{A}} = \mathcal{A} \cup \{a_{\text{stop}}\}$ .

127 For  $t \geq 1$  and  $h \in \mathcal{H}_t$ , define the optimal reward-to-collect value

$$V_t^*(h; \lambda) := \sup_{\bar{\pi}} \mathbb{E}_{\theta \sim \nu}^{\bar{\pi}} [\lambda r_{\bar{\tau}}(H_{\bar{\tau}}) - (\bar{\tau} - t) \mid H_t = h]. \quad (3)$$

128 Define

$$Q_{t, \text{stop}}^*(h; \lambda) := \lambda r_t(h), \quad Q_{t, \text{cont}}^*(h, a; \lambda) := -1 + \mathbb{E} [V_{t+1}^*(H_{t+1}; \lambda) \mid H_t = h, A_t = a],$$

129 where the expectation is under the posterior predictive kernel  $\bar{P}_t(\cdot \mid h, a)$ .

130 **Theorem 3.1** (Bellman equation and greedy deterministic optimality). *Assume the regularity condi-*  
 131 *tions of Section B.3 (Assumption 2 and Assumption 4), then*

$$V_t^*(h; \lambda) = \max \left\{ Q_{t, \text{stop}}^*(h; \lambda), \sup_{a \in \mathcal{A}} Q_{t, \text{cont}}^*(h, a; \lambda) \right\}. \quad (4)$$

132 *Moreover, let  $a_t^*(h) \in \arg \max_{a \in \mathcal{A}} Q_{t, \text{cont}}^*(h, a; \lambda)$ . The deterministic policy  $\bar{\pi}^*$  defined by*

$$\bar{\pi}^*(h) = a_{\text{stop}}^* \quad \text{if } Q_{t, \text{stop}}^*(h; \lambda) \geq Q_{t, \text{cont}}^*(h, a_t^*(h); \lambda),$$

133 *and  $\bar{\pi}^*(h) = a_t^*(h)$  otherwise, is optimal for the fixed- $\lambda$  inner problem.*

134 Unlike the finite case, where Bellman attainment is automatic, the continuous setting requires  
 135 verifying that the supremum over  $a \in \mathcal{A}$  in Eq. (4) is attained. Our proof (Section B.3.3) establishes  
 136 this through a value-iteration construction that also handles the coupling with the stopping/continue  
 137 structure, and propagates lower semicontinuity from the observation model through the posterior  
 138 predictive kernel to the  $Q$ -function. This semicontinuity chain is specific to this problem and does  
 139 not follow from any existing reference by specialization.

140 **Zero duality gap and  $(\epsilon, \delta)$ -correctness.** The Bellman theorem characterizes the inner problem  
 141 for a fixed multiplier  $\lambda$ . We now state when the Lagrangian relaxation is exact. Let  $c(\pi) := \mathbb{E}^\pi[\tau_\pi]$ ,  
 142  $\rho(\pi) := \mathbb{E}^\pi[r_{\tau_\pi}(H_{\tau_\pi})]$ , and  $\mathcal{K} := \{(c(\pi), \rho(\pi)) : \pi \in \mathcal{T}\}$ , where  $\tau_\pi$  is the stopping time of a policy  
 143 whose action space embeds the stopping decision. The time-sharing assumption, stated formally in  
 144 Assumption 7, says that ex-ante randomization between two admissible policies remains admissible  
 145 and therefore convexifies  $\mathcal{K}$ .

146 **Theorem 3.2** (Zero duality gap and  $(\epsilon, \delta)$ -correctness). *Assume time-sharing (Assumption 7), and*  
 147 *assume strict feasibility (Assumption 8), i.e., there exists a feasible  $\pi_{\text{sf}} \in \mathcal{T}$  such that  $\rho(\pi_{\text{sf}}) > 1 - \delta$ .*  
 148 *Then, the duality gap is zero. Furthermore, let  $g(\lambda) := \inf_{\pi \in \mathcal{T}} \{c(\pi) + \lambda(1 - \delta - \rho(\pi))\}$ ,  $\lambda^* \in$   
 149  $\arg \max_{\lambda \geq 0} g(\lambda)$ , and let  $\mathcal{S}(\lambda^*)$  be the set of dual minimizers at  $\lambda^*$ . If there exists  $\epsilon_0 > 0$  such that*

$$\mathcal{K}_{\epsilon_0}(\lambda^*) := \{(c, \rho) \in \mathcal{K} : c + \lambda^*(1 - \delta - \rho) \leq g(\lambda^*) + \epsilon_0\}$$

150 *is closed in  $\mathbb{R}^2$ , then there exists a dual-optimal policy  $\pi^* \in \mathcal{S}(\lambda^*)$  that is primal optimal. Conse-*  
 151 *quently, with the posterior-optimal inference rule  $I_t^*(h) \in \arg \max_{x \in \mathcal{X}} q_t(h, x)$ ,*

$$\mathbb{P}_{\theta \sim \nu}^{\pi^*} (L_\theta (I_{\tau_{\pi^*}}^*(H_{\tau_{\pi^*}})) \leq \epsilon) \geq 1 - \delta.$$

152 This theorem improves on the corresponding result in [53] in two ways. First, we replace the  
 153 assumption that the dual-optimal policy is unique with a weaker local closedness condition on  
 154 the near-optimal set  $\mathcal{K}_{\epsilon_0}(\lambda^*)$ : this allows multiple dual-optimal policies, which is more natural.  
 155 Second, while [53] derives correctness via a monotonicity argument on the optimal cost, our proof  
 156 (Section B.4) uses a direct subdifferential characterization to show that if all near-optimal policies  
 157 have  $\rho < 1 - \delta$ , then every subgradient of the dual value is strictly negative, contradicting optimality.  
 158 Zero duality gap follows from a standard perturbation argument [47, 10]; see Section B.5.

## 159 4 Continuous ICPE: C-ICPE

160 In this section we describe C-ICPE, a practical method based on the theory of the previous section.  
 161 C-ICPE has three components: inference, stopping, and exploration. Each of these are implemented  
 162 via learned sequential neural architectures, trained end-to-end from interaction data. Once trained, we  
 163 use C-ICPE at deployment time (a.k.a. test-time or inference-time) to perform pure exploration. We  
 164 now describe the models learned by C-ICPE and the training procedure. More details can be found in  
 165 Section C of the appendix.

166 **Training Protocol.** We adopt a similar meta-training protocol as the one used in [53] to train the  
 167 models. Briefly, we assume access to a simulator over  $\nu$  from which we can sample trajectories. We  
 168 use a meta-training, where we sample tasks  $\theta \sim \nu$  and assume access to a zero-loss target  $x_\theta^* \in \mathcal{X}$ ,  
 169 collect trajectories using C-ICPE, store the data in a replay buffer  $\mathcal{B}$ , and perform off-policy updates:  
 170 (1) a likelihood update of the parameters of the inference model, (2) a DQN-like update of the critic  
 171 and (3) an update of the actor based on the learned  $Q$  function. After training, we freeze all models;  
 172 deployment requires no access to  $x_\theta^*$  or the prior. We now discuss the modeling of these components  
 173 more in detail.

---

**Algorithm 1** C-ICPE
 

---

```

// Training phase
1: Initialize buffer  $\mathcal{B}$ , networks  $Q_\psi, I_\phi$ , actor  $\pi$ .
2: while Training is not over do
3:   Sample environment  $M_\theta \sim \nu$  and hypothesis  $x_\theta^*$ ; observe  $Y_1 \sim \rho$  and set  $t \leftarrow 1$ .
4:   repeat
5:     Execute action  $A_t \sim \pi(\cdot | H_t)$  according to actor  $\pi$  and observe  $Y_{t+1}$ .
6:     Add partial trajectory  $(H_t, A_t, Y_{t+1}, x_\theta^*)$  to  $\mathcal{B}$  and set  $t \leftarrow t + 1$ .
7:   until  $Q_\psi(H_t, a_{\text{stop}}) \geq Q_\psi(H_t, A_t)$ .
8:   In the fixed confidence, update  $c$  according to Eq. (9).
9:   Sample batch  $B \sim \mathcal{B}$  and update models using  $\mathcal{L}_{\text{inf}}(B; \phi)$  (Eq. (5)) and  $\mathcal{L}_{\text{critic}}(B; \psi)$  (Eq. (7)); for TD3, train  $\pi$ 
   according to  $\mathcal{L}_{\text{act}}(B; \psi)$  (Eq. (8)).
10: end while

// Inference/Deployment phase (models are fixed here)
11: Sample unknown environment  $M \sim \nu$  and collect a trajectory  $H_\tau$  using  $\pi$  (until  $Q_\psi(H_t, A_t) \leq Q_\psi(H_t, a_{\text{stop}})$ ).
12: Return  $\hat{x}_\tau = \mu_\phi(H_\tau)$  (recommendation).

```

---

174 **Gaussian inference model.** For a history  $h$ , the ideal inference rule maximizes the posterior success  
 175 probability  $q_t(h, x) := \mathbb{P}(L_\theta(x) \leq \epsilon | H_t = h)$  over recommendations  $x \in \mathcal{X}$ . However, computing  
 176  $q_t$  is not straightforward, as the posterior distribution may have a complex shape. Instead, we train  
 177 the inference model to learn the posterior law of  $x_\theta^*$  from trajectories and outputs a diagonal Gaussian  
 178 distribution characterizing the uncertainty around  $x_\theta^*$

$$I_\phi(\cdot | h) = \mathcal{N}(\mu_\phi(h), \text{diag}(\sigma_\phi^2(h))).$$

179 The recommendation at stopping is defined as the mean  $\hat{x} = \mu_\phi(h)$ . The covariance characterizes the  
 180 uncertainty around this point, and, as shown in Proposition 7, the optimal mean and covariance are  
 181 the posterior mean and covariance of  $x_\theta^*$  given  $H_t = h$ . Thus the Gaussian is a moment projection  
 182 of the posterior law of  $x_\theta^*$ . The deployed recommendation  $\hat{x} = \mu_\phi(h)$  should therefore be viewed  
 183 as a tractable approximation to  $\arg \max_{x \in \mathcal{X}} q_t(h, x)$ , rather than as an exact maximizer of  $q_t$ . In  
 184 Section B.7, we show that  $\mu_\phi$  is near-optimal when the posterior uncertainty on  $x_\theta^*$  is small relative  
 185 to the  $\epsilon$ -success margin (see Proposition 8). This justification is most direct for localization losses,  
 186 where  $\mathcal{X}_\epsilon(\theta)$  is a neighborhood of  $x_\theta^*$ .

187 We train  $\phi$  using a log-likelihood loss on a batch of partial trajectories  $B = (x_i^*, H_{t_i})_i$  sampled from  
 188 the buffer  $\mathcal{B}$ , where  $x_i^*$  is the optimal point for trajectory  $i$  and  $t_i$  is a timestep sampled uniformly at  
 189 random for that trajectory

$$\mathcal{L}_{\text{inf}}(B; \phi) = - \sum_{i=1}^{|B|} \log I_\phi(x_i^* | H_{t_i}). \quad (5)$$

190 In the following, we denote by  $\bar{\phi}$  the target parameter of the inference model, updated via a Polyak  
 191 update  $\bar{\phi} \leftarrow (1 - \tau_I)\bar{\phi} + \tau_I\phi$  with  $\tau_I \in (0, 1)$ .

192 **Critic and reward definition.** We parametrize the critic by  $\psi$ , and model it with two heads: a  
 193 continuation head  $Q_\psi(h, a)$  for  $a \in \mathcal{A}$  and a stopping head  $Q_\psi(h, a_{\text{stop}})$ . We also define the value:  
 194 let  $a_{\text{tgt}}(h')$  be the continuation target action proposed at the next history by the current policy, then  
 195 the value is defined as

$$V_\psi(h') := \max \{Q_\psi(h', a_{\text{stop}}), Q_\psi(h', a_{\text{tgt}}(h'))\},$$

196 Using the definition of the  $Q$  function from Theorem 3.1, we learn the parameter  $\psi$  using TD-learning.  
 197 While the ideal reward would be  $r_t(h) = \max_x q_t(h, x)$ , to better capture the uncertainty around the  
 198 recommendation  $\mu_t(h)$ , we use a sampled reward from the target inference model  $I_{\bar{\phi}}$ :

$$\hat{r}(h, \theta) := \frac{1}{K} \sum_{k=1}^K \mathbf{1} \{L_\theta(X^{(k)}) \leq \epsilon\}, \quad X^{(k)} \sim I_{\bar{\phi}}(\cdot | h). \quad (6)$$

199 Conditionally on  $(h, \theta)$ , this is an unbiased estimate of the probability that a sample from the inference  
 200 distribution lies in  $\mathcal{X}_\epsilon(\theta)$ . After averaging over the posterior, this reward is  $\mathbb{E}_{X \sim I_{\bar{\phi}}(\cdot | h)}[q_t(h, X)] \leq$   
 201  $r_t(h)$ , so it is a conservative version of the ideal posterior reward and allows the model to capture the  
 202 current uncertainty in the recommendation rule (Proposition 6 in the appendix makes this precise, as

203 the gap between the practical reward  $\hat{r}_t(h)$  and the ideal Bellman reward  $r_t(h)$  is controlled by the  
 204 second moment of the inference model).

205 Then, we sample a batch of partial trajectories  $B = (H_{t_i}, A_{t_i}, H_{t_i+1}, d_{t_i}, x_i^*)_i$  from the buffer, where  
 206  $t_i$  is a uniformly sampled timestep for the  $i$ -th trajectory and  $d_{t_i} = 1$  when the maximum horizon is  
 207 reached. Then, we use targets

$$y_i^{\text{stop}} = \hat{r}(H_{t_i}, \theta), \quad y_i^{\text{cont}} = -c + d_{t_i} \hat{r}(H_{t_i+1}, \theta) + (1 - d_{t_i}) V_{\bar{\psi}}(H_{t_i+1}),$$

208 where we reparametrized the Lagrange multiplier as a per-step cost  $c = 1/\lambda$  and the inner Lagrangian  
 209 objective remains the same. Then, the critic loss is

$$\mathcal{L}_{\text{critic}}(B; \psi) = \frac{1}{2|B|} \sum_{i=1}^{|B|} \left[ (Q_{\psi}(H_{t_i}, A_{t_i}) - y_i^{\text{cont}})^2 + (Q_{\psi}(H_{t_i}, a_{\text{stop}}) - y_i^{\text{stop}})^2 \right]. \quad (7)$$

210 The critic also decides when to stop: at rollout time the current policy  $\pi$  first proposes  $A_t$ , and the  
 211 learner stops iff

$$Q_{\psi}(H_t, a_{\text{stop}}) \geq Q_{\psi}(H_t, A_t).$$

212 Hence, at deployment, no access to  $\theta$  is required as the stopping comparison uses only  $H_t$ .

213 **Policy and actors.** The policy, or actor rule, decides what continuation action  $a \in \mathcal{A}$  to choose next.  
 214 Depending on the whether  $\mathcal{X} = \mathcal{A}$ , we propose three possible actor rules.

215 *Thompson Sampling (TS) rule:* this rule can be used when the recommendation space and the query  
 216 space coincide  $\mathcal{X} = \mathcal{A}$ , and eliminates the need for a separate actor. This rule draws inspiration from  
 217 classical Thompson Sampling [60]: the policy is implicitly represented by the inference model, and  
 218 the actions are sampled according to  $A_t \sim I_{\phi}(\cdot | H_t)$ . This is useful when informative queries are  
 219 themselves plausible recommendations. Early in an episode the exploration is more spread; later, as  
 220 the posterior target law contracts, TS concentrates near the current recommendation. As target action  
 221 for the critic we use the mean value of the inference model  $a_{\text{tgt}}(h) = \mu_t(h)$ .

222 *Top Two Posterior Sampling (TTPS) rule:* also this rule can be used when the recommendation space  
 223 and the query space coincide  $\mathcal{X} = \mathcal{A}$ . This rule is similar to classical TTPS [54], but we extend  
 224 it to the continuous case. This rule draws a sample  $A_t \sim I_{\phi}(\cdot | H_t)$  and, with probability 1/2, it  
 225 samples until the new sample is farther from the mean  $\mu_{\phi}(H_t)$ . The logic is that the posterior mean  
 226 is the current recommendation, while samples farther from it represent plausible alternatives. TTPS  
 227 therefore spends some probability mass checking alternatives instead of repeatedly querying near the  
 228 current mean before the stopping critic is confident. As target action for the critic we use the mean  
 229 value of the inference model  $a_{\text{tgt}}(h) = \mu_t(h)$ .

230 *TD3 rule* [18]: this rule can be used for general recommendation spaces when  $\mathcal{A} \neq \mathcal{X}$ , and formally  
 231 tries to solve the Bellman equation in Theorem 3.1. It learns a parametric actor  $\pi_{\vartheta}(h) \in \mathcal{A}$  from the  
 232 critic. The deterministic actor is trained by

$$\mathcal{L}_{\text{act}}(B; \vartheta) = -\frac{1}{|B|} \sum_i Q_{\psi}(H_{t_i}, \pi_{\vartheta}(H_{t_i})). \quad (8)$$

233 The critic target uses the usual TD3 stabilizers: a target actor, target-action smoothing, and twin  
 234 critics. In case  $\mathcal{X} = \mathcal{A}$  we can use a stochastic TD3 variants, where TD3 learns the mean  $\bar{\mu}$  and  
 235 covariance  $\bar{\Sigma}$  of a Gaussian actor. In this case the loss is augmented with an imitation learning  
 236 loss  $\text{KL}(I_{\psi}(\cdot | h) || \pi_{\vartheta}(\cdot | h))$  that provides a rich signal: the actor can increase variance when sampled  
 237 actions have higher continuation value and shrink it when exploration is no longer useful.

238 **Cost update.** We update the per-step cost  $c$  by simply performing a gradient step on the  
 239 dual variable. We sample a fresh batch of trajectories, and estimate the success rate  $\hat{p} =$   
 240  $\frac{1}{|B|} \sum_{i=1}^{|B|} \mathbf{1} \left\{ \mu_{\phi}(H_{t_i}^{(i)}) \in \mathcal{X}_c(\theta^{(i)}) \right\}$ , and update the cost as follows

$$c \leftarrow \text{Proj}_{[0,1]} (c - \eta_c ((1 - \delta) - \hat{p})). \quad (9)$$

241 If empirical success is below  $1 - \delta$ , the cost decreases and trajectories become longer; otherwise, if  
 242 above the target, the cost increases and stopping becomes more aggressive.

243 **Correctness certification.** The zero-duality result from the previous section justifies the ideal  
 244 Lagrangian objective, but a trained model is still approximate. In Section B.6 we outline how to  
 245 obtain formal  $(\epsilon, \delta)$ -guarantees on the trained model (see Proposition 5).

## 246 5 Empirical Evaluation

247 We evaluate C-ICPE on various benchmarks: noisy binary search,  $\epsilon$ -best arm identification on the  
 248 unit sphere, Ackley minimization, and GP max-value estimation. We also validate on a real-world  
 249 geochemical exploration task [1]. We compare four exploration rules: TS, TTPS, TD3, and uniform  
 250 sampling. For all experiments we set a maximum sample complexity  $t_{\max}$  (details in Section D) and  
 251 report 95% confidence intervals using bootstrap. We also compare against Bayesian optimization  
 252 baselines whenever possible: Tree-structured Parzen Estimator (TPE) [8], Gaussian Process (GP)  
 253 with UCB or Expected Improvement [57, 5], and CMA-ES [24]. Each is given a fixed sample budget  
 254 larger than the median stopping time of C-ICPE. These methods are not  $(\epsilon, \delta)$ -correct competitors,  
 255 and optimize a fixed-budget objective. We include them to test whether standard methods already  
 256 attains the target correctness at comparable budgets. For  $\epsilon$ -best arm on the sphere, we additionally  
 257 compare against Lazy Track-and-Stop [30], an optimal frequentist fixed-confidence baseline.

### 258 5.1 Synthetic Benchmarks

259 We now provide a brief description of the benchmarks, and then discuss the results.

260 **Noisy binary search.** The environment parameter  $\theta$  is drawn uniformly from  $[-1, 1]^d$ , with selector  
 261  $x^*(\theta) = \theta$ . The agent queries  $a \in [-1, 1]^d$  and observes, per coordinate,  $y_i = \xi_i \cdot \text{sign}(\theta_i - a_i)$ ,  
 262 where  $\xi_i \in \{-1, +1\}$  are i.i.d. Rademacher random variables with  $\mathbb{P}(\xi_i = +1) = 1 - p$ . The loss is  
 263  $L_\theta(x) = \|x - \theta\|_2$ , so the  $\epsilon$ -optimal set is  $\mathcal{X}_\epsilon(\theta) = \{x : \|x - \theta\|_2 \leq \epsilon\}$ . Here  $\mathcal{X} = \mathcal{A} = [-1, 1]^d$ .  
 264 The difficulty is controlled by the noise rate  $p$  and the dimension  $d$ : each coordinate provides one bit  
 265 of corrupted information per query, and the agent must simultaneously localize all  $d$  coordinates.

266  **$\epsilon$ -best arm on the sphere.** The environment parameter  $\theta$  is drawn uniformly on the unit sphere  
 267  $\mathbb{S}^{d-1}$ , with selector  $x^*(\theta) = \theta$ . The agent queries  $a \in [-1, 1]^d$  and observes a noisy linear reward  
 268  $y = f_\theta(a) + \xi$ , where  $f_\theta(a) = \theta^\top a$  and  $\xi \sim \mathcal{N}(0, \sigma^2)$ . The loss is defined via the inner product:  
 269  $L_\theta(x) = 1 - f_\theta(x)$ , so the  $\epsilon$ -optimal set is  $\mathcal{X}_\epsilon(\theta) = \{x : f_\theta(x) \geq 1 - \epsilon\}$ . Here  $\mathcal{X} = \mathcal{A}$ , and the  
 270 difficulty lies in estimating a direction from scalar projections.

271 **Ackley minimization.** The agent must locate the global minimizer of a randomly parametrized  
 272 Ackley function [40], a standard multimodal benchmark for global optimization. The parameter  
 273 is  $\theta = (\alpha, \beta, \gamma, \theta^*)$ , where  $(\alpha, \beta, \gamma)$  control the function shape and  $\theta^* \in [-1, 1]^d$  is the global  
 274 minimizer (selector  $x^*(\theta) = \theta^*$ ). The agent queries  $a \in [-1, 1]^d$  and observes a normalized  
 275 function evaluation  $y = \tilde{f}_{\alpha, \beta, \gamma}(a - \theta^*) + \xi$ ,  $\xi \sim \mathcal{N}(0, \sigma^2)$ . The loss is  $L_\theta(x) = \|x - \theta^*\|_2$ . Here  
 276  $\mathcal{X} = \mathcal{A} = [-1, 1]^d$ . The difficulty arises from the function’s many local optima and nearly flat outer  
 277 region, which can trap greedy strategies; the agent must explore globally before converging.

278 **GP max-value estimation.** A function  $f$  is sampled from a Gaussian process  $\text{GP}(0, k_{\text{RBF}}(\ell, \sigma_f))$   
 279 on  $[0, 1]^d$ , with lengthscale  $\ell \sim \text{Unif}[0.05, 0.2]$  and output scale  $\sigma_f = 1$ . The target is the scalar  
 280 maximum value  $\theta^* = \max_x f(x)$ , with selector  $x^*(\theta) = \theta^* \in \mathbb{R}$ . The agent queries  $a \in [0, 1]^d$  and  
 281 observes  $y = f(a) + \xi$ ,  $\xi \sim \mathcal{N}(0, \sigma^2)$ . The loss is  $L_\theta(x) = |x - \theta^*|$ . This is the  $\mathcal{X} \neq \mathcal{A}$  setting:  
 282 the recommendation space  $\mathcal{X} \subseteq \mathbb{R}$  is scalar while the action space  $\mathcal{A} = [0, 1]^d$  is  $d$ -dimensional,

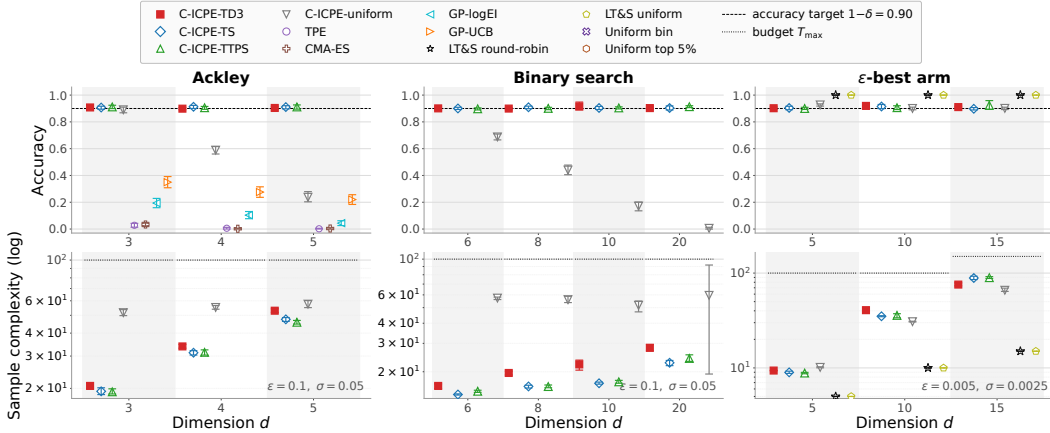


Figure 2: Accuracy (top) and sample complexity (bottom) at the hardest  $(\epsilon, \sigma)$  per benchmark.

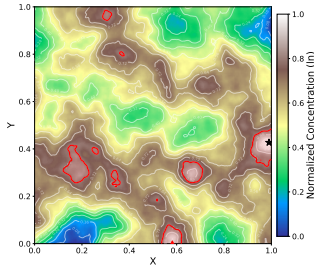


Figure 3: copper concentration in a 2D region in the geochemical exploration task. Red regions indicate concentration of copper within  $\epsilon$  of the maximum value.

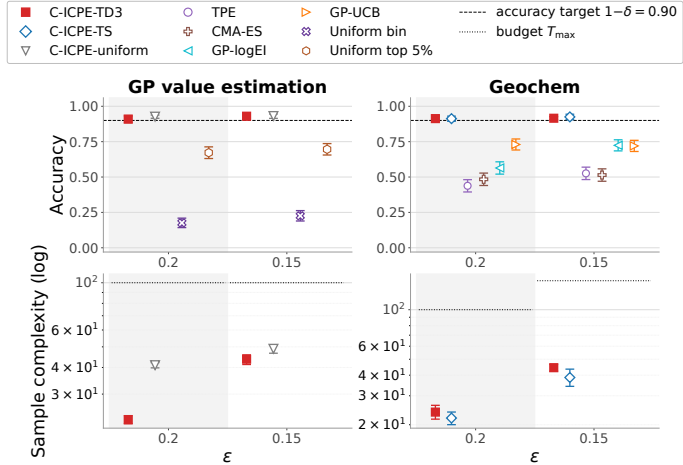


Figure 4: Accuracy (top) and sample complexity (bottom) at the hardest  $\sigma$  per benchmark.

283 requiring the TD3 actor to learn an exploration policy decoupled from the inference model. The  
 284 difficulty is twofold: the agent must both explore to find the region of high function values and  
 285 estimate the peak value to within  $\epsilon$ , without knowing the function’s lengthscale in advance. For this  
 286 problem we compare against two non-parametric baselines: (1) uniform sampling over the domain,  
 287 reporting the trimmed mean of the top-5% observed values; (2) partitioning the domain into uniform  
 288 bins, sampling uniformly within each bin, and reporting the highest bin average as the value estimate.

289 **Results.** Figs. 2-4 report accuracy and sample complexity (and their 95% confidence intervals) at  
 290 the hardest  $(\epsilon, \sigma)$  per benchmark; full sweeps over  $(\epsilon, \sigma, d)$ , experimental details, and robustness  
 291 to prior misspecification are in Section D. Across all four tasks, C-ICPE with learned exploration  
 292 (TS, TTPS, or TD3) consistently meets the  $1 - \delta$  accuracy target, while BO baselines and non-  
 293 parametric estimators fall well below, confirming that fixed-budget optimization does not yield  
 294  $(\epsilon, \delta)$ -correctness. C-ICPE-uniform is competitive at low  $d$  but degrades as dimension increases on  
 295 Ackley and binary search, where directed exploration matters. On  $\epsilon$ -best arm, uniform exploration is  
 296 optimal by rotational symmetry [30], and Lazy Track-and-Stop achieves optimal sample complexity  
 297 by exploiting the linear structure; C-ICPE matches the accuracy target but uses more samples,  
 298 reflecting the cost of a model-agnostic stopping rule. On GP max-value estimation ( $\mathcal{X} \neq \mathcal{A}$ ), C-ICPE-  
 299 TD3 meets the target while non-parametric baselines fall short; BO methods are inapplicable here as  
 300 they return locations rather than values. Particularly, in Section B.9 we prove that (Theorem B.5),  
 301 under an RBF-GP prior with interior regularity, max-value estimation is asymptotically harder than  
 302 argmax localization: we establish a sample complexity lower bound for value estimation and an upper  
 303 bound for argmax identification via a two-stage algorithm (T-BAL; Algorithm 2), showing that the  
 304 geometric structure of the function helps localization but not value estimation.

## 305 5.2 Geochemical Exploration

306 Lastly, we construct a realistic task using data from the USGS Geochemical Survey [1], which  
 307 provides measurements of copper concentration across the United States. The goal is to identify the  
 308 location of peak copper concentration in an unknown region with  $(\epsilon, \delta)$ -guarantees. We partition the  
 309 data into geographic regions, fit a sparse variational Gaussian process to each, and split regions into  
 310 training and evaluation. Training regions are used to meta-train C-ICPE; we evaluate on held-out  
 311 regions whose spatial structure was not seen during training. See also Section D.5 for more details.

312 **Results.** In Fig. 4 we present the results. In this problem, for sake of simplicity we only test the  
 313 TD3 and TS actors for ICPE, and compare with respect to classical Bayesian baselines. The results  
 314 show accuracy and sample complexity on held-out regions: C-ICPE-TD3, C-ICPE-TS achieve the  
 315  $1 - \delta$  accuracy target, while BO baselines fall below the target while using a larger sample budget,  
 316 demonstrating that both learned exploration and learned stopping contribute to this real-data task.  
 317 This shows evidence that C-ICPE transfers across tasks with genuine distribution shift. In Section D.5  
 318 we report the experimental details, and results for different values of  $\epsilon$ .

## 319 6 Discussion, Related Work and Conclusions

320 Active sequential hypothesis testing (ASHT) provides the broad conceptual umbrella for this paper.  
321 In ASHT, a learner adaptively selects experiments and decides when to stop and declare a hypothesis,  
322 with the objective of minimizing expected sample size subject to a correctness constraint [12, 62, 22,  
323 37, 38]. A key methodological theme in this literature is that fixed-confidence constraints can be  
324 handled via Lagrangian duality. Closely related line of works include Bayesian experimental design  
325 and Bayesian active learning, which study adaptive data acquisition when the unknown is drawn  
326 from a known prior, typically optimizing expected utility or information gain [35, 23, 45], and active  
327 learning, which emphasizes selecting informative queries/labels to reduce uncertainty efficiently [13].  
328 Despite the shared emphasis on adaptive measurement and sequential stopping, most classical ASHT  
329 results assume substantial knowledge of the observation model: one typically has access to likelihoods  
330 (or at least to a parametric family) for every experiment under every hypothesis, enabling explicit  
331 likelihood-ratio statistics and model-based allocation [37]. While there are efforts toward relaxing  
332 this assumption to partial model knowledge [11], the need for explicit likelihood structure remains  
333 a limiting factor for modern continuous environments with complex, history-dependent feedback.  
334 In many practical settings, the learner must instead infer both (i) which latent environment/task it is  
335 facing and (ii) which actions are informative, using only interaction data and function approximation.  
336 This motivates data-driven approaches that preserve the ASHT objective while reducing dependence  
337 on fully specified likelihoods.

338 The most developed special case of ASHT is fixed-confidence pure exploration in bandits, where the  
339 hypothesis is the identity of an optimal action. In finite-armed bandits [34], best-arm identification  
340 (BAI) at  $\epsilon = 0$  is characterized by a mature theory: instance-dependent lower bounds quantify the  
341 intrinsic complexity of identifying the best arm [19, 14, 63, 30, 33, 50, 43, 52], and a family of  
342 algorithms achieves near-optimal sample complexity by coupling adaptive allocation with statistically  
343 valid stopping rules [6, 54, 19, 64, 31]. These results provide both sharp guidance and strong baselines,  
344 but they rely on a finite decision set and model-specific likelihood constructions. Similar themes arise  
345 in pure exploration for Markov decision processes, where the goal is to identify an optimal policy  
346 with probability at least  $1 - \delta$  [3, 59, 4, 49, 51]. This literature yields sharp insights into exploration  
347 complexity but is likewise developed for finite state-action structure and frequentist guarantees.

348 Even under well-specified models, moving from  $\epsilon = 0$  to  $(\epsilon, \delta)$ -PAC identification in continuous  
349 decision spaces can be technically demanding. Several recent works address continuous pure  
350 exploration in bandit models [20]. Takemori et al. [58] give a tractable algorithm for continuous-arm  
351 linear bandits and Poiani et al. [42] derive lower bounds and a Track-and-Stop framework for infinite-  
352 answer problems; both are frequentist, model-specific, and require explicit likelihood structure.  
353 In MDPs, some work begun to treat  $(\epsilon, \delta)$ -PAC objectives in finite MDP settings for best policy  
354 identification [61] and optimal data-collection for policy evaluation [48], both in the frequentist setting.  
355 To our knowledge, no general Bayesian theory is known for continuous recommendation spaces under  
356 general priors; current Bayesian fixed-confidence results are limited to finite bandits with Gaussian  
357 likelihoods and Gaussian priors [29]. While Bayesian ideas drive exploration in finite bandit (with  
358 frequentist guarantees), e.g. posterior sampling and top-two methods [55, 54, 56], the most developed  
359 Bayesian framework in continuous spaces is Bayesian optimization (BO) [21], which maintains a  
360 posterior over an unknown objective and selects queries via acquisition functions [28, 26]. BO aims  
361 to identify an optimizer, but is typically posed as fixed-budget optimization without a correctness  
362 constraint. Wilson [65] recently introduced a Bayesian  $(\epsilon, \delta)$ -stopping rule for BO, but it is restricted  
363 to GP surrogates and does not optimize sample complexity. Despite this progress, no existing method  
364 combines three elements: fixed-confidence  $(\epsilon, \delta)$  stopping, continuous recommendations, and a  
365 learned exploration and recommender procedure under a Bayesian formulation. C-ICPE addresses  
366 this gap.

367 **Conclusions.** C-ICPE is a theory-inspired method for Bayesian fixed-confidence pure exploration  
368 with continuous recommendations. On the theoretical side, we establish that the Lagrangian duality  
369 and Bellman optimality structure of finite ASHT carries over to continuous spaces under regularity  
370 conditions, and prove  $(\epsilon, \delta)$ -correctness under a local closedness assumption that is weaker than the  
371 uniqueness condition required in [53]. On the algorithmic side, we show that C-ICPE achieves the  
372 desired guarantees across different tasks, including a real-world geochemical exploration problem,  
373 while using fewer samples than standard optimization baselines. To our knowledge, no prior method  
374 combines continuous recommendations, fixed-confidence stopping in a single practical framework.  
375 Limitations and broader impact are discussed in Section A.

## References

- 376 [1] The National Geochemical Survey: Database and documentation. Report 2004-1001, 2004.
- 377 [2] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama. Optuna: A next-generation hyper-  
378 parameter optimization framework. In *Proceedings of the 25th ACM SIGKDD international*  
379 *conference on knowledge discovery & data mining*, pages 2623–2631, 2019.
- 380 [3] A. Al Marjani, A. Garivier, and A. Proutiere. Navigating to the best policy in markov decision  
381 processes. In *Advances in neural information processing systems*, volume 34, pages 25852–  
382 25864, 2021.
- 383 [4] A. Al Marjani, T. Kocak, and A. Garivier. On the Complexity of All  $\epsilon$ -Best Arms Identifi-  
384 cation. In M.-R. Amini, S. Canu, A. Fischer, T. Guns, P. Kralj Novak, and G. Tsoumakas,  
385 editors, *Machine Learning and Knowledge Discovery in Databases*, volume 13716, pages  
386 317–332, Cham, 2023. Springer Nature Switzerland. ISBN 978-3-031-26411-5 978-3-031-  
387 26412-2. doi: 10.1007/978-3-031-26412-2\_20. URL [https://link.springer.com/10.](https://link.springer.com/10.1007/978-3-031-26412-2_20)  
388 [1007/978-3-031-26412-2\\_20](https://link.springer.com/10.1007/978-3-031-26412-2_20). Series Title: Lecture Notes in Computer Science.
- 389 [5] S. Ament, S. Daulton, D. Eriksson, M. Balandat, and E. Bakshy. Unexpected improvements to  
390 expected improvement for bayesian optimization. *Advances in neural information processing*  
391 *systems*, 36:20577–20612, 2023.
- 392 [6] J.-Y. Audibert and S. Bubeck. Best arm identification in multi-armed bandits. In *COLT-23th*  
393 *Conference on learning theory-2010*, pages 13–p, 2010.
- 394 [7] M. Balandat, B. Karrer, D. Jiang, S. Daulton, B. Letham, A. G. Wilson, and E. Bakshy. Botorch:  
395 A framework for efficient monte-carlo bayesian optimization. *Advances in neural information*  
396 *processing systems*, 33:21524–21538, 2020.
- 397 [8] J. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl. Algorithms for hyper-parameter optimization.  
398 *Advances in neural information processing systems*, 24, 2011.
- 399 [9] D. Bertsekas and S. E. Shreve. *Stochastic optimal control: the discrete-time case*, volume 5.  
400 Athena Scientific, 1996.
- 401 [10] J. Borwein and A. Lewis. *Convex Analysis and Nonlinear Optimization: Theory and Examples*.  
402 Springer, 2006.
- 403 [11] F. Cecchi and N. Hegde. Adaptive active hypothesis testing under limited information. *Advances*  
404 *in Neural Information Processing Systems*, 30, 2017.
- 405 [12] H. Chernoff. Sequential design of experiments. *The Annals of Mathematical Statistics*, 30(3):  
406 755 – 770, 1959. doi: 10.1214/aoms/1177706205. URL [https://doi.org/10.1214/aoms/](https://doi.org/10.1214/aoms/1177706205)  
407 [1177706205](https://doi.org/10.1214/aoms/1177706205). Publisher: Institute of Mathematical Statistics.
- 408 [13] D. A. Cohn, Z. Ghahramani, and M. I. Jordan. Active learning with statistical models. *Journal*  
409 *of artificial intelligence research*, 4:129–145, 1996.
- 410 [14] R. Degenne, W. M. Koolen, and P. Ménard. Non-asymptotic pure exploration by solving  
411 games. In H. Wallach, H. Larochelle, A. Beygelzimer, F. dAlché Buc, E. Fox, and R. Garnett,  
412 editors, *Advances in neural information processing systems*, volume 32. Curran Associates,  
413 Inc., 2019. URL [https://proceedings.neurips.cc/paper\\_files/paper/2019/file/](https://proceedings.neurips.cc/paper_files/paper/2019/file/8d1de7457fa769ece8d93a13a59c8552-Paper.pdf)  
414 [8d1de7457fa769ece8d93a13a59c8552-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/8d1de7457fa769ece8d93a13a59c8552-Paper.pdf).
- 415 [15] B. Efron. Bootstrap methods: another look at the jackknife. In *Breakthroughs in statistics:*  
416 *Methodology and distribution*, pages 569–593. Springer, 1992.
- 417 [16] I. Ekeland and R. Temam. *Convex analysis and variational problems*. SIAM, 1999.
- 418 [17] E. A. Feinberg and P. O. Kasyanov. MDPs with setwise continuous transition probabilities.  
419 *Operations Research Letters*, 49(5):734–740, 2021.
- 420 [18] S. Fujimoto, H. Hoof, and D. Meger. Addressing function approximation error in actor-critic  
421 methods. In *International conference on machine learning*, pages 1587–1596. PMLR, 2018.
- 422

- 423 [19] A. Garivier and E. Kaufmann. Optimal best arm identification with fixed confidence. *Proceed-*  
424 *ings of the 29th Conference on Learning Theory*, 49:998–1027, 2016.
- 425 [20] A. Garivier and E. Kaufmann. Nonasymptotic sequential tests for overlapping hypotheses  
426 applied to near-optimal arm identification in bandit models. *Sequential Analysis*, 40(1):61–96,  
427 2021. Publisher: Taylor & Francis.
- 428 [21] R. Garnett. *Bayesian Optimization*. Cambridge University Press, 2023.
- 429 [22] B. K. Ghosh. A brief history of sequential analysis. *Handbook of sequential analysis*, 1, 1991.
- 430 [23] D. Golovin and A. Krause. Adaptive submodularity: Theory and applications in active learning  
431 and stochastic optimization. *Journal of Artificial Intelligence Research*, 42:427–486, 2011.
- 432 [24] N. Hansen. The cma evolution strategy: A tutorial. *arXiv preprint arXiv:1604.00772*, 2016.
- 433 [25] A. Hantoute and MA. López. Characterizations of the subdifferential of the supremum of  
434 convex functions. *Journal of Convex Analysis*, 15:831–858, 2008.
- 435 [26] P. Hennig and C. J. Schuler. Entropy Search for Information-Efficient Global Optimization.  
436 *Journal of Machine Learning Research*, 13(57):1809–1837, 2012. ISSN 1533-7928.
- 437 [27] O. Hernández-Lerma and J. B. Lasserre. *Discrete-Time Markov Control Processes*. Springer,  
438 New York, NY, 1996. ISBN 978-1-4612-6884-0 978-1-4612-0729-0. doi: 10.1007/  
439 978-1-4612-0729-0. URL <http://link.springer.com/10.1007/978-1-4612-0729-0>.
- 440 [28] J. M. Hernández-Lobato, M. W. Hoffman, and Z. Ghahramani. Predictive Ent-  
441 *ropy Search for Efficient Global Optimization of Black-box Functions*. In *Ad-*  
442 *vances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc.,  
443 2014. URL [https://proceedings.neurips.cc/paper\\_files/paper/2014/hash/](https://proceedings.neurips.cc/paper_files/paper/2014/hash/6488484c982e9af5c35689523ba1abfe-Abstract.html)  
444 [6488484c982e9af5c35689523ba1abfe-Abstract.html](https://proceedings.neurips.cc/paper_files/paper/2014/hash/6488484c982e9af5c35689523ba1abfe-Abstract.html).
- 445 [29] K. Jang, J. Komiyama, and K. Yamazaki. Fixed Confidence Best Arm Iden-  
446 *tification in the Bayesian Setting*. In A. Globerson, L. Mackey, D. Belgrave,  
447 A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors, *Advances in Neural Infor-*  
448 *mation Processing Systems*, volume 37, pages 17789–17829. Curran Associates, Inc.,  
449 2024. URL [https://proceedings.neurips.cc/paper\\_files/paper/2024/file/](https://proceedings.neurips.cc/paper_files/paper/2024/file/1fb0a4de9c14f5557eeea886e22569cd-Paper-Conference.pdf)  
450 [1fb0a4de9c14f5557eeea886e22569cd-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/1fb0a4de9c14f5557eeea886e22569cd-Paper-Conference.pdf).
- 451 [30] Y. Jedra and A. Proutiere. Optimal best-arm identification in linear bandits. In *Advances in*  
452 *neural information processing systems*, volume 33, pages 10007–10017, 2020.
- 453 [31] M. Jourdan, R. Degenne, D. Baudry, R. de Heide, and E. Kaufmann. Top Two Algorithms Revis-  
454 *ited*. In *Advances in Neural Information Processing Systems*, volume 35, pages 26791–26803,  
455 Dec. 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/hash/](https://proceedings.neurips.cc/paper_files/paper/2022/hash/ab5f5f22e3e09f4424592ffb06840ab0-Abstract-Conference.html)  
456 [ab5f5f22e3e09f4424592ffb06840ab0-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2022/hash/ab5f5f22e3e09f4424592ffb06840ab0-Abstract-Conference.html).
- 457 [32] E. Kaufmann and W. M. Koolen. Mixture martingales revisited with applications to sequential  
458 tests and confidence intervals. *Journal of Machine Learning Research*, 22(246):1–44, 2021.
- 459 [33] T. Kocák and A. Garivier. Best arm identification in spectral bandits. In *Proceedings of the*  
460 *twenty-ninth international joint conference on artificial intelligence, IJCAI’20*, Yokohama,  
461 Yokohama, Japan, 2021. ISBN 978-0-9992411-6-5. Number of pages: 7 tex.articleno: 307.
- 462 [34] T. Lattimore and C. Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- 463 [35] D. V. Lindley. On a measure of the information provided by an experiment. *The Annals of*  
464 *Mathematical Statistics*, 27(4):986–1005, 1956.
- 465 [36] Y. Liu, T. Zhao, W. Ju, and S. Shi. Materials discovery and design using machine learning.  
466 *Journal of Materiomics*, 3(3):159–177, 2017.
- 467 [37] M. Naghshvar and T. Javidi. Active Sequential Hypothesis Testing. *The Annals of Statistics*, 41  
468 (6):2703–2738, 2013. ISSN 0090-5364.

- 469 [38] M. Naghshvar, T. Javidi, and K. Chaudhuri. Noisy bayesian active learning. In *2012 50th annual*  
470 *allerton conference on communication, control, and computing (allerton)*, pages 1626–1633.  
471 IEEE, 2012.
- 472 [39] NASA/GSFC/METI/ERSDAC/JAROS, and U.S./Japan ASTER Science Team. Mountain  
473 pass mine, california, 2010. URL [https://www.jpl.nasa.gov/images/  
474 pia13979-mountain-pass-mine-california](https://www.jpl.nasa.gov/images/pia13979-mountain-pass-mine-california). Image PIA13979, acquired March 28,  
475 2010.
- 476 [40] M. Naser, M. K. Al-Bashiti, A. T. G. Tapeh, A. Naser, V. Kodur, R. Hawileh, J. Abdalla,  
477 N. Khodadadi, A. H. Gandomi, and A. D. Eslamlou. A review of benchmark and test func-  
478 tions for global optimization algorithms and metaheuristics. *Wiley Interdisciplinary Reviews:  
479 Computational Statistics*, 17(2):e70028, 2025.
- 480 [41] J. O’Quigley, M. Pepe, and L. Fisher. Continual reassessment method: a practical design for  
481 phase I clinical trials in cancer. *Biometrics*, 46(1):33–48, Mar. 1990. ISSN 0006-341X.
- 482 [42] R. Poiani, M. Bernasconi, and A. Celli. Pure Exploration with Infinite Answers, May 2025.
- 483 [43] R. Poiani, M. Jourdan, E. Kaufmann, and R. Degenne. Best-Arm Identification in Unimodal  
484 Bandits. In *Proceedings of The 28th International Conference on Artificial Intelligence and  
485 Statistics*, pages 2233–2241. PMLR, Apr. 2025. URL [https://proceedings.mlr.press/  
486 v258/poiani25a.html](https://proceedings.mlr.press/v258/poiani25a.html). ISSN: 2640-3498.
- 487 [44] M. L. Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John  
488 Wiley & Sons, 2014.
- 489 [45] T. Rainforth, A. Foster, D. R. Ivanova, and F. Bickford Smith. Modern Bayesian experimental  
490 design. *Statistical Science*, 39(1):100–114, 2024. Publisher: Institute of Mathematical Statistics.
- 491 [46] R. T. Rockafellar. *Conjugate duality and optimization*. Society for Industrial and Applied  
492 Mathematics, 1974.
- 493 [47] R. T. Rockafellar and R. J. Wets. *Variational analysis*. Springer, 1998.
- 494 [48] A. Russo and A. Pacchiano. Adaptive exploration for multi-reward multi-policy evaluation. In  
495 A. Singh, M. Fazel, D. Hsu, S. Lacoste-Julien, F. Berkenkamp, T. Maharaj, K. Wagstaff, and  
496 J. Zhu, editors, *Proceedings of the 42nd international conference on machine learning*, volume  
497 267 of *Proceedings of machine learning research*, pages 52382–52421. PMLR, July 2025. URL  
498 <https://proceedings.mlr.press/v267/russo25a.html>.
- 499 [49] A. Russo and A. Proutiere. Model-free active exploration in reinforcement learning. In *Advances  
500 in neural information processing systems*, volume 36, pages 54740–54753, 2023.
- 501 [50] A. Russo and A. Proutiere. On the sample complexity of representation learning in multi-task  
502 bandits with global and local structure. In *Proceedings of the AAAI conference on artificial  
503 intelligence*, volume 37, pages 9658–9667, 2023.
- 504 [51] A. Russo and F. Vannella. Multi-reward best policy identification. In *Advances in neural  
505 information processing systems*, volume 37, pages 105583–105662, 2024.
- 506 [52] A. Russo, Y. Song, and A. Pacchiano. Pure exploration with feedback graphs. In *Proceedings  
507 of the 28th international conference on artificial intelligence and statistics*, volume 258 of  
508 *Proceedings of machine learning research*, pages 1810–1818. PMLR, 2025.
- 509 [53] A. Russo, R. Welch, and A. Pacchiano. Learning to Explore: An In-Context Learning  
510 Approach for Pure Exploration, June 2025. URL <http://arxiv.org/abs/2506.01876>.  
511 arXiv:2506.01876 [cs].
- 512 [54] D. Russo. Simple Bayesian Algorithms for Best Arm Identification. In *Conference on Learning  
513 Theory*, pages 1417–1418. PMLR, June 2016. URL [https://proceedings.mlr.press/  
514 v49/russo16.html](https://proceedings.mlr.press/v49/russo16.html).
- 515 [55] D. Russo and B. Van Roy. Learning to optimize via posterior sampling. *Mathematics of  
516 Operations Research*, 39(4):1221–1243, 2014.

- 517 [56] X. Shang, R. Heide, P. Menard, E. Kaufmann, and M. Valko. Fixed-confidence guarantees for  
518 bayesian best-arm identification. In *International Conference on Artificial Intelligence and*  
519 *Statistics*, pages 1823–1832. PMLR, 2020.
- 520 [57] N. Srinivas, A. Krause, S. Kakade, and M. Seeger. Gaussian process optimization in the bandit  
521 setting: no regret and experimental design. In *Proceedings of the 27th International Conference*  
522 *on International Conference on Machine Learning*, pages 1015–1022, 2010.
- 523 [58] S. Takemori, Y. Umeda, and A. Gopalan. Instance-optimal pure exploration for linear bandits  
524 on continuous arms. In *Forty-second International Conference on Machine Learning*, 2025.
- 525 [59] J. Taupin, Y. Jedra, and A. Proutiere. Best policy identification in discounted linear MDPs. In  
526 *Sixteenth european workshop on reinforcement learning*, 2023.
- 527 [60] W. R. Thompson. On the likelihood that one unknown probability exceeds another in view of  
528 the evidence of two samples. *Biometrika*, 25(3-4):285–294, 1933.
- 529 [61] A. Tirinzoni, A. Al Marjani, and E. Kaufmann. Near instance-optimal pac reinforcement learning  
530 for deterministic mdps. In *Advances in neural information processing systems*, volume 35,  
531 pages 8785–8798, 2022.
- 532 [62] A. Wald and J. Wolfowitz. Optimum character of the sequential probability ratio test. *The*  
533 *Annals of Mathematical Statistics*, 19(3):326 – 339, 1948. doi: 10.1214/aoms/1177730197.  
534 URL <https://doi.org/10.1214/aoms/1177730197>.
- 535 [63] P.-A. Wang, A. Proutiere, K. Ariu, Y. Jedra, and A. Russo. Optimal algorithms for mul-  
536 tiplayer multi-armed bandits. In S. Chiappa and R. Calandra, editors, *Proceedings of the*  
537 *twenty third international conference on artificial intelligence and statistics*, volume 108  
538 of *Proceedings of machine learning research*, pages 4120–4129. PMLR, Aug. 2020. URL  
539 <https://proceedings.mlr.press/v108/wang20m.html>.
- 540 [64] P.-A. Wang, R.-C. Tzeng, and A. Proutiere. Fast Pure Exploration via Frank-Wolfe. In *Advances*  
541 *in Neural Information Processing Systems*, volume 34, pages 5810–5821. Curran Associates,  
542 Inc., 2021.
- 543 [65] J. T. Wilson. Stopping bayesian optimization with probabilistic regret bounds. *Advances in*  
544 *Neural Information Processing Systems*, 37:98264–98296, 2024.

545	<b>Contents</b>	
546	<b>1 Introduction</b>	<b>1</b>
547	<b>2 Problem Setting</b>	<b>2</b>
548	<b>3 Theoretical Background</b>	<b>3</b>
549	<b>4 Continuous ICPE: C-ICPE</b>	<b>4</b>
550	<b>5 Empirical Evaluation</b>	<b>7</b>
551	5.1 Synthetic Benchmarks . . . . .	7
552	5.2 Geochemical Exploration . . . . .	8
553	<b>6 Discussion, Related Work and Conclusions</b>	<b>9</b>
554	<b>Appendix</b>	<b>15</b>
555	<b>A Limitations and Broader Impact</b>	<b>16</b>
556	<b>B Appendix: Theoretical Results</b>	<b>17</b>
557	B.1 Problem Modeling . . . . .	17
558	B.2 Posterior distribution over the true hypothesis and inference rule optimality . . . . .	19
559	B.3 Fixed-confidence setting: formulation and optimal rules . . . . .	22
560	B.3.1 Optimal Inference Rule . . . . .	23
561	B.3.2 Stopping as an action (equivalence) . . . . .	23
562	B.3.3 Optimal Policy . . . . .	23
563	B.3.4 Reward Shaping and Removal of the Stop Action . . . . .	28
564	B.4 Fixed-confidence setting: $(\epsilon, \delta)$ -correctness of dual-optimal points . . . . .	30
565	B.5 Fixed-confidence setting: zero duality gap via perturbation values . . . . .	33
566	B.5.1 Zero Duality by Perturbation . . . . .	34
567	B.5.2 Primal attainment and KKT . . . . .	35
568	B.6 Training-time certification of $(\epsilon, \delta)$ -correctness . . . . .	37
569	B.7 Choice of inference model and reward modeling . . . . .	40
570	B.8 Robustness to prior misspecification . . . . .	44
571	B.9 Sample Complexity: Value Estimation vs Argmax Localization . . . . .	47
572	B.9.1 Max-value estimation complexity . . . . .	48
573	B.9.2 Argmax localization complexity . . . . .	51
574	<b>C Appendix: Algorithms</b>	<b>58</b>
575	C.1 History Encoder and Time Pooling Layer . . . . .	58
576	C.2 Replay Buffer and Prefix Sampling . . . . .	58
577	C.3 Inference Update and Its Regularization . . . . .	59

578	C.4	Reward, Critic Update, and Critic Regularization . . . . .	59
579	C.5	Actor Rules: TS, TTPS, and TD3 . . . . .	60
580	C.6	Cost Update for Fixed Confidence . . . . .	61
581	<b>D</b>	<b>Appendix: Numerical Results</b>	<b>63</b>
582	D.1	Synthetic Benchmarks: description . . . . .	63
583	D.1.1	Noisy binary search. . . . .	64
584	D.1.2	$\epsilon$ -best-arm identification on the sphere. . . . .	64
585	D.1.3	Ackley minimizer identification. . . . .	65
586	D.1.4	GP max-value estimation . . . . .	66
587	D.2	Synthetic Benchmarks: baselines . . . . .	66
588	D.3	Synthetic Benchmarks: numerical results . . . . .	68
589	D.3.1	Noisy Binary Search . . . . .	68
590	D.3.2	$\epsilon$ -best arm problem . . . . .	70
591	D.3.3	Ackley minimization . . . . .	72
592	D.3.4	GP max-value estimation . . . . .	74
593	D.4	Synthetic Benchmarks: robustness . . . . .	76
594	D.4.1	Noisy Binary Search . . . . .	76
595	D.4.2	$\epsilon$ -best arm problem . . . . .	76
596	D.4.3	Ackley minimization . . . . .	77
597	D.4.4	GP max-value estimation . . . . .	78
598	D.5	Geochemical Exploration: Experimental Details and Numerical Results . . . . .	80
599	D.5.1	Dataset and Motivation . . . . .	80
600	D.5.2	Region Partitioning and GP Fitting . . . . .	80
601	D.5.3	Ground Truth Construction . . . . .	81
602	D.5.4	Task Prior and Train/Test Split . . . . .	81
603	D.5.5	Numerical Results . . . . .	85

## 604 A Limitations and Broader Impact

605 **Limitations.** *Parametric inference model.* C-ICPE models the posterior law of the target  $x_\theta^*$  with  
606 a diagonal Gaussian. In the ideal NLL objective, this corresponds to a moment projection: the  
607 mean matches  $\mathbb{E}[x_\theta^*|H_t]$  and the diagonal covariance matches the posterior coordinate variances.  
608 This is only a surrogate for the Bayes  $\epsilon$ -rule  $\arg \max_{x \in \mathcal{X}} q_t(H_t, x)$ . When the posterior over  $x_\theta^*$  is  
609 multimodal, or poorly summarized by first and second moments, the mean recommendation may  
610 lie between plausible targets and the covariance may misrepresent uncertainty. This can affect both  
611 exploration and stopping, since the critic evaluates samples from the same inference distribution.  
612 Richer posterior families, such as mixtures or normalizing flows, or a direct model of  $q_t(h, x)$ , could  
613 reduce this mismatch at the cost of additional optimization and training complexity.

614 Note that this limitation is most benign in localization benchmarks, where  $\mathcal{X}_\epsilon(\theta)$  is a ball or interval  
615 around  $x_\theta^*$ . It is more pronounced in value-gap tasks such as Ackley or geochemical optimization,  
616 where the  $\epsilon$ -optimal set can be anisotropic, nonconvex, or multimodal. These tasks are therefore  
617 useful stress tests for the Gaussian inference model.

618 *Selector availability.* Training the inference model requires access to the selector  $x^*(\theta)$  for each  
619 sampled task  $\theta \sim \nu$ . In our benchmarks this is available in closed form (the shifted minimizer for  
620 Ackley, the parameter itself for binary search, the GP maximum for max-value estimation). In general,  
621 computing  $x^*(\theta)$  may require numerical optimization, introducing approximation error in the NLL  
622 targets. If the selector can only be evaluated approximately or with noise, the inference model may  
623 learn a biased posterior, potentially affecting both recommendation quality and stopping calibration.  
624 Extending C-ICPE to settings where only noisy or approximate selectors are available is an important  
625 direction for future work.

626 *Cost calibration.* The dual variable  $c$ , which controls the exploration–stopping tradeoff, is updated  
627 online during training via primal feedback on the empirical success rate. In practice, the convergence  
628 of  $c$  and the sensitivity of stopping behavior to its value require careful tuning of learning rates and  
629 update schedules.

630 *Decoupled action and recommendation spaces.* When  $\mathcal{X} \neq \mathcal{A}$ , C-ICPE requires a separate TD3  
631 actor to learn the exploration policy. This adds architectural complexity and an additional source of  
632 approximation error. Our GP max-value experiment exercises this setting.

633 *Bayesian guarantee and prior dependence.* The  $(\epsilon, \delta)$ -correctness guarantee is average-case under  
634 the task prior  $\nu$ . When the deployment distribution differs substantially from  $\nu$ , correctness may  
635 degrade. Our robustness experiments and the prior-shift bounds in Section B.8 provide some  
636 quantitative control, but worst-case guarantees for individual task instances are not provided. This  
637 limitation is shared by all Bayesian methods and is analogous to the dependence of frequentist  
638 methods on their parametric assumptions.

639 *Scalability.* The method’s behavior in higher dimensions (larger than  $d \geq 50$ ), where posterior  
640 concentration is slower and exploration is harder, remains to be investigated. The LSTM (or trans-  
641 former) architecture scales with the maximum horizon  $t_{\max}$ , which may need to grow with dimension,  
642 increasing both training and inference cost.

643 **Broader impact.** C-ICPE is a general-purpose tool for adaptive experimentation with correct-  
644 ness guarantees. Potential applications include materials discovery, dose-finding in clinical trials,  
645 environmental monitoring, and any setting where sequential experiments are costly and the practi-  
646 tioner requires a principled stopping criterion. In such settings, reducing sample complexity directly  
647 translates to reduced cost, time, and resource consumption.

648 We do not foresee direct negative societal impacts from the method itself. However, as with any  
649 system that automates experimental decisions, users should be aware that the  $(\epsilon, \delta)$  guarantee is  
650 conditional on the modeling assumptions (the task prior  $\nu$  and the observation model). Deploying  
651 C-ICPE in safety-critical domains, such as clinical dose-finding, would require careful validation of  
652 these assumptions and, where appropriate, additional safeguards beyond the Bayesian guarantee.

## 653 B Appendix: Theoretical Results

654 **Roadmap and novelty guide.** The theoretical analysis proceeds in four stages. We summarize  
655 what is standard and what is new relative to the finite ICPE framework of Russo et al. [53].

- 656 • **§B.2: Posterior success probability.** We define  $q_t(h, x) = \mathbb{P}(L_\theta(x) \leq \epsilon \mid H_t = h)$  and  
657 establish its regularity properties. The proofs use standard tools (Radon–Nikodym, reverse  
658 Fatou, Portmanteau); the object  $q_t(h, x)$  itself, the natural continuous analogue of posterior  
659 mass, has not previously been studied in the pure exploration literature.
- 660 • **§B.3: Bellman optimality.** We prove that the optimal value satisfies a stop/continue Bellman  
661 equation and that the supremum over continuation actions is attained by a measurable  
662 selector. The proof adapts the value-iteration framework of Bertsekas and Shreve [9] to  
663 our setting, with the main technical content being a semicontinuity induction that threads  
664 likelihood continuity through posterior weak continuity, predictive weak continuity,  $Q$ -  
665 function lower semicontinuity, and a sup–inf interchange. None of these steps appear in  
666 Russo et al. [53], where attainment is automatic for finite action spaces. The resulting  
667 dependency chain and its coupling with the posterior success payoff are specific to this  
668 problem.
- 669 • **§B.4–B.5: Correctness and zero duality gap.** We prove  $(\epsilon, \delta)$ -correctness under a local  
670 closedness condition that is strictly weaker than the uniqueness assumption in Russo et al.  
671 [53]. Our proof uses a subdifferential characterization from Hantoute and López [25]  
672 to derive a contradiction without the monotonicity argument of Russo et al. [53]; this is  
673 the strongest theorem-level novelty in the appendix. Zero duality gap follows using a  
674 perturbation argument [46].
- 675 • **§B.6: Model certification.** We introduce a checkpointwise certification protocol based on  
676 mixture supermartingales [32]. Unlike the pooled approach in Russo et al. [53], it tests each  
677 frozen checkpoint independently and requires no monotonicity assumption on the training  
678 trajectory.
- 679 • **§B.7: Inference model and sampled reward.** We relate the implemented Gaussian  
680 inference model to the ideal posterior quantities used in the Bellman characterization. We  
681 show that the sampled reward is an unbiased estimate of the success probability of the  
682 stochastic selector and is conservative relative to the ideal reward  $r_t(h)$ ; the gap is controlled  
683 by second moments of the inference distribution (Proposition 6). We also characterize the  
684 Gaussian NLL as a moment projection of the posterior law of  $x_\theta^*$  and give conditions under  
685 which the NLL mean is near-optimal (Propositions 7 and 8).
- 686 • **§B.8: Robustness to prior misspecification.** We analyse the robustness to prior misspecifi-  
687 cation, and what is the predicted impact on sample complexity and accuracy.
- 688 • **§B.9: Sample Complexity of Value Estimation vs Argmax Localization in Gaussian**  
689 **Processes.** In this section we prove that, under a hierarchical RBF-GP prior with a high-  
690 probability interior regularity condition, max-value estimation is asymptotically harder than  
691 argmax localization. Specifically, under this regularity assumption, we establish that the  
692 value estimation problem cannot be circumvented by the geometric structure of the function.
  - 693 – The analysis is based on showing a lower bound on the sample complexity of estimating  
694 the max-value, and an upper bound on estimating the argmax.
  - 695 – To this aim, we introduce an algorithm, Two-Stage Bayesian Argmax Localization  
696 (T-BAL) Algorithm 2, for locating the argmax of a Gaussian Process.
  - 697 – T-BAL works in two phases: first, searches the domain for a region where  $X^*$  may be  
698 located, and then perform gradient ascent using noisy finite differences to approximate  
699 the gradients.
  - 700 – We provide a sample complexity upper bound of T-BAL and provide  $(\epsilon, \delta)$ -guarantees.

### 701 B.1 Problem Modeling

702 We specialize to the fixed-confidence  $((\epsilon, \delta)$ -PAC) setting introduced in Section 2, and provide a  
703 self-contained definition of the induced probability measures.

704 We now provide a formal definition of the underlying probability measures of the problem we consider.  
 705 To that aim, it is important to formally define what a model  $M$  is, as well as the definition of policy  $\pi$   
 706 and inference rule  $I$  (inference rules are also known as recommendation rules).

707 **Spaces and histories.** Let  $\Theta \subset \mathbb{R}^d$  be compact. Let  $\mathcal{A} \subset \mathbb{R}^m$  be a compact action (query) space  
 708 and  $\mathcal{Y} \subset \mathbb{R}^n$  a compact observation space, each endowed with the Borel  $\sigma$ -algebra. Let  $\mathcal{X}$  be a  
 709 compact hypothesis/decision space (in our experiments  $\mathcal{X} = \mathcal{A}$ ). For  $t \in \mathbb{N}$ , define the history space

$$\mathcal{H}_t := (\mathcal{Y} \times \mathcal{A})^{t-1} \times \mathcal{Y}, \quad h_t = (y_1, a_1, \dots, a_{t-1}, y_t),$$

710 with its product Borel  $\sigma$ -algebra. We also write  $\mathcal{H}_\infty := \mathcal{Y} \times (\mathcal{A} \times \mathcal{Y})^\mathbb{N}$  for infinite histories. Since  
 711  $\mathcal{A}, \mathcal{Y}$  are compact metric spaces,  $\mathcal{H}_t$  and  $\mathcal{H}_\infty$  are standard Borel spaces.

712 **Environment (observation model).** An environment is indexed by  $\theta \in \Theta$  and specified by an  
 713 initial observation law  $\rho_\theta \in \Delta(\mathcal{Y})$  and a sequence of (possibly history-dependent) observation kernels

$$P_{\theta,t}(\cdot|h_t, a_t) \in \Delta(\mathcal{Y}), \quad t \geq 1,$$

714 such that for every Borel  $C \subset \mathcal{Y}$  the map  $(h_t, a) \mapsto P_{\theta,t}(C|h_t, a)$  is measurable. Optionally, one  
 715 may assume weak continuity in  $\theta$ . However, we do assume weak continuity in  $a$ , as this is later used  
 716 to prove optimality.

717 **Assumption 1** (Weak continuity of the transition). *For all  $\theta \in \Theta$  we assume  $a \mapsto P_{\theta,t}(\cdot|h_t, a)$  to be*  
 718 *weakly continuous.*<sup>1</sup>

719 **Learner: policy, stopping time, inference rule.** A (possibly randomized) sampling policy is a  
 720 sequence of probability kernels

$$\pi_t(\cdot|h_t) \in \Delta(\mathcal{A}), \quad t \geq 1,$$

721 measurable as maps  $\mathcal{H}_t \rightarrow \Delta(\mathcal{A})$ . Let  $H_t = (Y_1, A_1, \dots, A_{t-1}, Y_t)$  be the random history and  
 722  $\mathcal{F}_t = \sigma(H_t)$ . A stopping time  $\tau$  is defined w.r.t.  $(\mathcal{F}_t)_{t \geq 1}$ . An inference rule is a sequence of  
 723 measurable maps  $I_t : \mathcal{H}_t \rightarrow \mathcal{X}$ , and the learner outputs

$$\hat{x}_\tau := I_\tau(H_\tau).$$

724 **Loss and  $\epsilon$ -optimal set.** For each  $\theta \in \Theta$ , the environment induces a loss function  $L_\theta : \mathcal{X} \rightarrow [0, \infty)$   
 725 with  $\inf_{x \in \mathcal{X}} L_\theta(x) = 0$ . Define the  $\epsilon$ -optimal set

$$\mathcal{X}_\epsilon(\theta) := \{x \in \mathcal{X} : L_\theta(x) \leq \epsilon\}.$$

726 In the following we the following conditions on  $L_\theta(x)$ .

727 **Assumption 2.** *We assume joint lower-semicontinuity of  $(x, \theta) \mapsto L_\theta(x)$  and joint Borel measurabil-*  
 728 *ity.*

729 **Remark 1** (Regularity of selector-based localization losses). *Some of our localization examples*  
 730 *define the loss through a selected target  $x_{\text{sel}}^*(\theta) \in \mathcal{X}$ , for instance  $L_\theta(x) = \|x - x_{\text{sel}}^*(\theta)\|$ . In this*  
 731 *case, the standing lower-semicontinuity assumption on  $L(\theta, x) = L_\theta(x)$  is satisfied whenever the*  
 732 *selector  $\theta \mapsto x_{\text{sel}}^*(\theta)$  is continuous. There are two standard ways to obtain such a selector:*

733 1. *First, if  $F(\theta) := \arg \max_{x \in \mathcal{X}} f_\theta(x) = \{x^*(\theta)\}$  is singleton for every  $\theta$ ,  $f(\theta, x)$  is jointly*  
 734 *continuous, and  $\mathcal{X}$  is compact, then the maximum theorem implies that the unique optimizer*  
 735 *is continuous, and thus  $x_{\text{sel}}^*(\theta) = x^*(\theta)$  is continuous. This covers the shifted Ackley*  
 736 *benchmark when the shift determines a unique optimizer continuously, the linear bandit*  
 737 *benchmark on the unit sphere where  $x^*(\theta) = \theta$ , and noisy binary search when the target*  
 738 *map is continuous.*

739 2. *Second, if  $F(\theta)$  is not singleton but has nonempty compact convex values and is Hausdorff-*  
 740 *continuous in  $\theta$ , then the minimum-norm selector  $x_{\text{sel}}^*(\theta) := \arg \min_{x \in F(\theta)} \|x\|^2$  is well-*  
 741 *defined and continuous: closed convex values give uniqueness of the minimum-norm point,*  
 742 *while Hausdorff continuity gives stability of this point as  $\theta$  varies.*

743 *Thus the selector-based distance loss is jointly continuous in these cases. This selector also gives a*  
 744 *canonical target for the Gaussian inference model: the NLL objective learns a moment projection of*  
 745 *the posterior law of  $x_{\text{sel}}^*(\theta)$ , while correctness remains defined through the loss  $L_\theta$  and the set  $\mathcal{X}_\epsilon(\theta)$ .*

<sup>1</sup>That is, for all continuous bounded functions  $f$  we have that  $a \mapsto \int_{\mathcal{Y}} f(y) P_{\theta,t}(dy|h_t, a)$  is continuous.

746 **Path measures (Ionescu–Tulcea).** Fix  $\theta \in \Theta$  and a policy  $\pi$ . By the Ionescu–Tulcea theorem,  
 747 there exists a unique probability measure  $\mathbb{P}_{\theta,t}^\pi$  on  $(\mathcal{H}_t, \mathcal{B}(\mathcal{H}_t))$  such that for all cylinder sets  $C =$   
 748  $C_1 \times B_1 \times \cdots \times B_{t-1} \times C_t$  (with  $C_i \in \mathcal{B}(\mathcal{Y})$  and  $B_i \in \mathcal{B}(\mathcal{A})$ ),

$$\mathbb{P}_{\theta,t}^\pi(C) = \int_{C_1} \rho_\theta(dy_1) \prod_{s=1}^{t-1} \left[ \int_{B_s} \pi_s(da_s|h_s) \int_{C_{s+1}} P_{\theta,s}(dy_{s+1}|h_s, a_s) \right].$$

749 Analogously, one obtains a unique path measure  $\mathbb{P}_\theta^\pi$  on  $(\mathcal{H}_\infty, \mathcal{B}(\mathcal{H}_\infty))$ .

750 **Mixture law over tasks.** Given a prior  $\nu$  on  $\Theta$ , define the joint law on  $\Theta \times \mathcal{H}_t$  by

$$\mathbf{P}_t^\pi(d\theta, dh_t) := \nu(d\theta) \mathbb{P}_{\theta,t}^\pi(dh_t),$$

751 and the trajectory marginal  $\mathbb{P}_t^\pi(\cdot) = \int \mathbb{P}_{\theta,t}^\pi(\cdot) \nu(d\theta)$ . We use  $\mathbb{E}_{\theta \sim \nu}^\pi[\cdot]$  and  $\mathbb{P}_{\theta \sim \nu}^\pi(\cdot)$  for expectation-  
 752 s/probabilities under this mixture.

753 **Fixed-confidence objective.** The learner is  $(\epsilon, \delta)$ -correct (under  $\nu$ ) if

$$\mathbb{P}_{\theta \sim \nu}^\pi(\hat{x}_\tau \in \mathcal{X}_\epsilon(\theta)) \geq 1 - \delta.$$

754 In the fixed-confidence regime, we seek to minimize the expected number of queries subject to  
 755  $(\epsilon, \delta)$ -correctness:

$$\inf_{\pi, I, \tau} \mathbb{E}_{\theta \sim \nu}^\pi[\tau] \quad \text{s.t.} \quad \mathbb{P}_{\theta \sim \nu}^\pi(\hat{x}_\tau \in \mathcal{X}_\epsilon(\theta)) \geq 1 - \delta.$$

## 756 B.2 Posterior distribution over the true hypothesis and inference rule optimality

757 We first record a domination assumption that allows us to express likelihoods w.r.t. fixed reference  
 758 measures.

759 **Assumption 3 (Domination).** *There exist probability measures  $\lambda_0, \lambda$  on  $(\mathcal{Y}, \mathcal{B}(\mathcal{Y}))$  such that, for all*  
 760  *$\theta \in \Theta$ , all  $t \geq 1$ , and all  $(h_t, a) \in \mathcal{H}_t \times \mathcal{A}$ ,*

$$\rho_\theta(\cdot) \ll \lambda_0(\cdot) \quad \text{and} \quad P_{\theta,t}(\cdot|h_t, a) \ll \lambda(\cdot).$$

761 Let  $p_{\theta,0}(y) := \frac{d\rho_\theta}{d\lambda_0}(y)$  and  $p_{\theta,t}(y'|h_t, a) := \frac{dP_{\theta,t}(\cdot|h_t, a)}{d\lambda}(y')$  be versions of the corresponding  
 762 densities, chosen jointly measurable in their arguments.

763 *Remark.* The assumption holds, for instance, when all  $\rho_\theta$  and  $P_{\theta,t}(\cdot|h_t, a)$  admit densities w.r.t. a  
 764 common reference measure (e.g., Lebesgue on  $\mathcal{Y} \subset \mathbb{R}^n$  or counting measure when  $\mathcal{Y}$  is finite).

765 Under Assumption 3, define the (policy-independent) likelihood of a realized history  $h_t =$   
 766  $(y_1, a_1, \dots, a_{t-1}, y_t) \in \mathcal{H}_t$  under parameter  $\theta$ :

$$\ell_t(\theta, h_t) := p_{\theta,0}(y_1) \prod_{s=1}^{t-1} p_{\theta,s}(y_{s+1}|h_s, a_s), \quad h_s = (y_1, a_1, \dots, a_{s-1}, y_s).$$

767 We now give a posterior kernel representation that is independent of  $\pi$ .

768 **Lemma 1 (Posterior kernel over  $\Theta$ ).** *Consider Assumption 3. For each  $t \in \mathbb{N}$  there exists a probability*  
 769 *kernel  $R_t : \mathcal{H}_t \times \mathcal{B}(\Theta) \rightarrow [0, 1]$ , independent of  $\pi$ , such that for every policy  $\pi$ , all  $A \in \mathcal{B}(\Theta)$  and*  
 770  *$Z \in \mathcal{B}(\mathcal{H}_t)$ ,*

$$\mathbf{P}_t^\pi(\theta \in A, H_t \in Z) = \int_Z R_t(A|h) \mathbb{P}_t^\pi(dh),$$

771 where  $\mathbf{P}_t^\pi(d\theta, dh) = \nu(d\theta) \mathbb{P}_{\theta,t}^\pi(dh)$  and  $\mathbb{P}_t^\pi$  is its  $\mathcal{H}_t$ -marginal. Moreover, for  $\mathbb{P}_t^\pi$ -a.e.  $h \in \mathcal{H}_t$ ,

$$R_t(A|h) = \frac{\int_A \ell_t(\theta, h) \nu(d\theta)}{\int_\Theta \ell_t(\theta, h) \nu(d\theta)}.$$

772 Consequently, for any measurable map  $g : \Theta \rightarrow \mathcal{S}$  into a standard Borel space  $\mathcal{S}$  and any  $B \in \mathcal{B}(\mathcal{S})$ ,

$$\mathbb{P}(g(\theta) \in B | H_t = h) = R_t(\{\theta : g(\theta) \in B\} | h) \quad \text{for } \mathbb{P}_t^\pi\text{-a.e. } h.$$

773 *Proof.* Fix  $\pi$  and  $t$ . Define the reference measure on  $\mathcal{H}_t$  (depending on  $\pi$ )

$$\nu_t^\pi(dh_t) := \lambda_0(dy_1) \prod_{s=1}^{t-1} [\pi_s(da_s|h_s) \lambda(dy_{s+1})].$$

774 By construction and Assumption 3,  $\mathbb{P}_{\theta,t}^\pi \ll \nu_t^\pi$  for every  $\theta$ , with Radon–Nikodym density

$$\frac{d\mathbb{P}_{\theta,t}^\pi}{d\nu_t^\pi}(h_t) = \ell_t(\theta, h_t),$$

775 which does not depend on  $\pi$ . Therefore, for  $A \in \mathcal{B}(\Theta)$  and  $Z \in \mathcal{B}(\mathcal{H}_t)$ ,

$$\mathbf{P}_t^\pi(\theta \in A, H_t \in Z) = \int_A \int_Z \ell_t(\theta, h) \nu_t^\pi(dh) \nu(d\theta),$$

776 and

$$\mathbb{P}_t^\pi(Z) = \int_Z \int_\Theta \ell_t(\theta, h) \nu(d\theta) \nu_t^\pi(dh).$$

777 Hence  $\mathbf{P}_t^\pi(\theta \in A, \cdot) \ll \mathbb{P}_t^\pi(\cdot)$  and the Radon–Nikodym derivative is the displayed Bayes ratio, which  
 778 defines the kernel  $R_t(A|h)$ . Measurability and the fact that  $R_t(\cdot|h)$  is a probability measure follow  
 779 from standard properties of Radon–Nikodym derivatives. Independence of  $\pi$  is immediate from the  
 780 explicit formula.  $\square$

781 In the following we also need to consider in what cases the mapping  $h \mapsto R_t(\cdot|h)$  is weakly  
 782 continuous. To that aim, we require a further assumption.

783 **Assumption 4** (Likelihood continuity). *For each state  $s < t$ , the kernel  $p_{\theta,s}(y|h_s, a_s)$  is jointly con-*  
 784 *tinuous in  $(\theta, y, h_s, a_s)$  with strictly positive density. Hence,  $(\theta, h_t) \mapsto \ell(\theta, h_t)$  is jointly continuous.*

785 Under this assumption, we have the following.

786 **Lemma 2** (Weak continuity of the posterior). *Consider Assumption 3 and Assumption 4. For each  $t$*   
 787 *we have that the mapping  $h \mapsto R_t(\cdot|h)$ ,  $h \in \mathcal{H}_t$ , is weakly continuous.*

788 *Proof.* Consider any sequence  $(h_n)_n \in \mathcal{H}_t$  such that  $h_n \rightarrow h$ ,  $h \in \mathcal{H}_t$ . Fix  $f \in C_b(\Theta)$  (continuous  
 789 and bounded). Since for  $\mathbb{P}_t^\pi$ -a.e.  $h' \in \mathcal{H}_t$  we have  $R_t(d\theta|h') = \frac{\ell_t(\theta, h') \nu(d\theta)}{\int_\Theta \ell_t(\theta, h') \nu(d\theta)}$ , we have that

$$\int_\Theta f(\theta) R_t(d\theta|h_n) = \frac{\int_\Theta f(\theta) \ell_t(\theta, h_n) \nu(d\theta)}{\int_\Theta \ell_t(\theta, h_n) \nu(d\theta)}.$$

790 Now, since  $\Theta$  is compact and  $\theta \mapsto f$  is continuous and  $\theta \mapsto \ell_t(\theta, h')$  is continuous for each  $h' \in \mathcal{H}_t$ ,  
 791 we have that  $\sup_{\theta \in \Theta} |f(\theta) \ell_t(\theta, h_n)| < \infty$ , therefore by dominated convergence we have

$$\int_\Theta f(\theta) \ell_t(\theta, h_n) \nu(d\theta) \rightarrow \int_\Theta f(\theta) \ell_t(\theta, h) \nu(d\theta) \quad \text{and} \quad \int_\Theta \ell_t(\theta, h_n) \nu(d\theta) \rightarrow \int_\Theta \ell_t(\theta, h) \nu(d\theta).$$

792 Since  $\int_\Theta \ell_t(\theta, h) \nu(d\theta) > 0$  by strict positivity of the density  $p_s$ , we have that

$$\int_\Theta f(\theta) \ell_t(\theta, h_n) \nu(d\theta) \rightarrow \int_\Theta f(\theta) \ell_t(\theta, h) \nu(d\theta),$$

793 so  $h \mapsto R_t(\cdot|h)$  is weakly continuous.  $\square$

794 **Optimal inference rule.** Fix  $\epsilon > 0$  and  $t \in \mathbb{N}$ . Recall that, for each  $\theta \in \Theta$ , the  $\epsilon$ -optimal set  
 795 is  $\mathcal{X}_\epsilon(\theta) = \{x \in \mathcal{X} : L_\theta(x) \leq \epsilon\}$ . Given a realized history  $h \in \mathcal{H}_t$ , define the *posterior success*  
 796 *probability* of recommending  $x \in \mathcal{X}$  as

$$q_t(h, x) := \mathbb{P}_{\theta \sim \nu}^\pi(x \in \mathcal{X}_\epsilon(\theta) | H_t = h) = R_t(\{\theta \in \Theta : L_\theta(x) \leq \epsilon\} | h), \quad (10)$$

797 where  $R_t(\cdot|h)$  is the posterior kernel from Lemma 1. We also define

$$r_t(h) := \sup_{x \in \mathcal{X}} q_t(h, x). \quad (11)$$

798 **Lemma 3.** *Under Assumption 2, we have that  $q_t(h, x)$  is jointly Borel measurable in  $(h, x)$  and upper*  
 799 *semicontinuous in  $x$  for each fixed  $h$ . If in addition to Assumption 2 we also assume Assumption 4,*  
 800 *then  $q_t(h, x)$  is jointly upper semicontinuous.*

801 *Proof.* We prove the 3 properties separately.

- 802 • Regarding measurability, note that  $(x, \theta) \mapsto \mathbf{1}\{L_\theta(x) \leq \epsilon\}$  is Boreal measurable by As-  
 803 sumption 2. Using that for every  $A \in \mathcal{B}(\Theta)$  we have that  $h \mapsto R_t(A|h)$  is Borel-measurable,  
 804 then  $q_t(h, x) = \int_\Theta \mathbf{1}\{L_\theta(x) \leq \epsilon\} R_t(d\theta|h)$  is jointly measurable since integration preserves  
 805 measurability.
- 806 • Consider now the u.s.c. property of  $x \mapsto q_t(h, x)$  for each  $h$ . If, for every  $\theta$ , the map  $x \mapsto$   
 807  $L_\theta(x)$  is continuous on the compact set  $\mathcal{X}$ , then  $\mathcal{X}_\epsilon(\theta)$  is closed and  $x \mapsto \mathbf{1}\{x \in \mathcal{X}_\epsilon(\theta)\}$  is  
 808 upper semicontinuous. Consequently,  $x \mapsto q_t(h, x)$  is upper semicontinuous  $\mathbb{P}_t^\pi$ -a.s.. To see  
 809 this, let  $(x_n)_n$  be a sequence in  $\mathcal{X}$  such that  $x_n \rightarrow x^*$ . Define  $y_n = \mathbf{1}\{x_n \in \mathcal{X}_\epsilon(\theta)\}$ . By  
 810 Fatou's reverse lemma we have

$$\limsup_n \mathbb{E}_t[y_n | H_t = h] \leq \mathbb{E}_t[\limsup_n y_n | H_t = h] \leq \mathbb{P}_t(x \in \mathcal{X}_\epsilon(\theta) | H_t = h).$$

811 where the last inequality follows from the fact that  $\limsup_n y_n \leq \mathbf{1}\{x \in \mathcal{X}_\epsilon(\theta)\}$  from the  
 812 upper semicontinuity. Thus the posterior is upper semicontinuous on  $\mathbb{P}_t^\pi$ -a.s.

- 813 • Define  $C_\epsilon = \{(\theta, x) \in \Theta \times \mathcal{X} : L_\theta(x) \leq \epsilon\}$ . By assumption on  $L$  (joint continuity), we  
 814 have that  $C_\epsilon$  is closed. Consider any sequence  $(h_n, x_n) \rightarrow (h, x)$  and define the distribution  
 815  $\mu_n(\cdot) := R_t(\cdot|h_n) \otimes \delta_{x_n}$  and  $\mu(\cdot) := R_t(\cdot|h) \otimes \delta_x$ . Since  $h \mapsto R_t(\cdot|h)$  is weakly  
 816 continuous (Lemma 2, follows from Assumption 4), by Portmanteau's lemma we have

$$\limsup_{n \rightarrow \infty} \mu_n(C_\epsilon) \leq \mu(C_\epsilon).$$

817 But  $\mu_n(C_\epsilon) = \int_\Theta \mathbf{1}\{L_\theta(x_n) \leq \epsilon\} R_t(d\theta|h_n) = q_t(h_n, x_n)$  and  $\mu(C_\epsilon) = \int_\Theta \mathbf{1}\{L_\theta(x) \leq$   
 818  $\epsilon\} R_t(d\theta|h) = q_t(h, x)$ , hence  $(h, x) \mapsto q_t(h, x)$  is jointly upper semicontinuous.

819 □

820 Since  $\mathcal{X}$  is compact, by the Extreme Value theorem we have that the supremum in (11) is attained (so  
 821 one may replace sup by max).

822 **Proposition 1** (Optimal inference). *Consider a fixed policy  $\pi$  and a fixed time  $t \in \mathbb{N}$ . Assume*  
 823 *Assumption 3 and Assumption 2. Among measurable inference rules  $I_t : \mathcal{H}_t \rightarrow \mathcal{X}$ , the maximal value*  
 824 *of  $\mathbb{P}_{\theta \sim \nu}^\pi(I_t(H_t) \in \mathcal{X}_\epsilon(\theta))$  is achieved by any rule satisfying, for  $\mathbb{P}_t^\pi$ -a.e.  $h \in \mathcal{H}_t$ ,*

$$I_t(h) \in \arg \max_{x \in \mathcal{X}} q_t(h, x).$$

825 *Furthermore, a measurable arg max selector exists and  $I_t(h)$  is measurable.*

826 *Proof.* Fix  $\pi$  and  $t$ , and let  $\hat{x}_t := I_t(H_t)$ . Using the posterior kernel,

$$\begin{aligned} \mathbb{P}_{\theta \sim \nu}^\pi(\hat{x}_t \in \mathcal{X}_\epsilon(\theta)) &= \int \mathbf{1}\{\hat{x}_t \in \mathcal{X}_\epsilon(\theta)\} \mathbf{P}_t^\pi(d\theta, dh) \\ &= \int_{\mathcal{H}_t} \left[ \int_\Theta \mathbf{1}\{I_t(h) \in \mathcal{X}_\epsilon(\theta)\} R_t(d\theta|h) \right] \mathbb{P}_t^\pi(dh) \\ &= \int_{\mathcal{H}_t} q_t(h, I_t(h)) \mathbb{P}_t^\pi(dh) \\ &\leq \int_{\mathcal{H}_t} \sup_{x \in \mathcal{X}} q_t(h, x) \mathbb{P}_t^\pi(dh) = \int_{\mathcal{H}_t} r_t(h) \mathbb{P}_t^\pi(dh). \end{aligned}$$

827 If  $I_t(h) \in \arg \max_{x \in \mathcal{X}} q_t(h, x)$  for  $\mathbb{P}_t^\pi$ -a.e.  $h$ , then the inequality holds with equality, yielding the  
 828 optimal value.

829 The existence of an arg max rule, and measurability of  $I_t$ , follows from an application of [27,  
 830 Proposition D.5]. □

831 We also note the following lower bound on the posterior  $q_t$ .

832 **Lemma 4** (Markov lower bound on posterior success). *Fix  $\epsilon > 0$ ,  $t \in \mathbb{N}$ , and a realized history*  
 833  *$h \in \mathcal{H}_t$ . For any decision  $x \in \mathcal{X}$ , define the posterior success probability*

$$q_t(h, x) := \mathbb{P}_{\theta \sim \nu}^{\pi}(L_{\theta}(x) \leq \epsilon | H_t = h),$$

834 *and the posterior mean loss*

$$\bar{L}_t(h, x) := \mathbb{E}_{\theta \sim \nu}^{\pi}[L_{\theta}(x) | H_t = h].$$

835 *Then*

$$q_t(h, x) \geq 1 - \frac{\bar{L}_t(h, x)}{\epsilon}.$$

836 *Equivalently, with the (clipped) shaped reward  $r_{\epsilon}(\theta, x) := [1 - L_{\theta}(x)/\epsilon]_{+}$ ,*

$$q_t(h, x) \geq \mathbb{E}_{\theta \sim \nu}^{\pi}[r_{\epsilon}(\theta, x) | H_t = h].$$

837 *Proof.* Since  $L_{\theta}(x) \geq 0$ , Markov's inequality yields

$$\mathbb{P}_{\theta \sim \nu}^{\pi}(L_{\theta}(x) > \epsilon | H_t = h) \leq \frac{\bar{L}_t(h, x)}{\epsilon}.$$

838 Taking complements gives the first claim. For the second, note that  $\mathbf{1}\{L_{\theta}(x) \leq \epsilon\} \geq [1 - L_{\theta}(x)/\epsilon]_{+}$   
 839 pointwise, and take conditional expectations.  $\square$

### 840 B.3 Fixed-confidence setting: formulation and optimal rules

841 We consider the fixed-confidence problem from Section 2 in its Bayesian (task-averaged) form. A  
 842 learner is a triplet  $(\pi, I, \tau)$  with sampling policy  $\pi$ , inference rule  $I = (I_t)_{t \geq 1}$ , and stopping time  $\tau$ .  
 843 The objective is

$$\inf_{\pi, I, \tau} \mathbb{E}_{\theta \sim \nu}^{\pi}[\tau] \quad \text{s.t.} \quad \mathbb{P}_{\theta \sim \nu}^{\pi}(I_{\tau}(H_{\tau}) \in \mathcal{X}_{\epsilon}(\theta)) \geq 1 - \delta, \quad \mathbb{E}_{\theta \sim \nu}^{\pi}[\tau] < \infty. \quad (12)$$

844 Throughout this section,  $H_t = (Y_1, A_1, \dots, A_{t-1}, Y_t)$  is the history,  $\hat{x}_{\tau} = I_{\tau}(H_{\tau})$  and  $\mathcal{F}_t = \sigma(H_t)$ .

845 **Posterior success.** For each  $t$  and realized history  $h \in \mathcal{H}_t$ , define the posterior success probability  
 846 of recommending  $x \in \mathcal{X}$  as

$$q_t(h, x) := \mathbb{P}_{\theta \sim \nu}^{\pi}(x \in \mathcal{X}_{\epsilon}(\theta) | H_t = h),$$

847 and recall  $r_t(h) := \sup_{x \in \mathcal{X}} q_t(h, x)$ . We have the following lemma that relates the success probability  
 848 to the expected posterior success.

849 **Lemma 5** (Stopped success as expected posterior success). *For any policy  $\pi$ , stopping time  $\tau$ , and*  
 850 *inference rule  $I$ ,*

$$\mathbb{P}_{\theta \sim \nu}^{\pi}(\hat{x}_{\tau} \in \mathcal{X}_{\epsilon}(\theta)) = \mathbb{E}^{\pi}[q_{\tau}(H_{\tau}, \hat{x}_{\tau})].$$

851 *Proof.* By the tower rule and  $\hat{x}_{\tau}$  being  $\sigma(H_{\tau})$ -measurable,

$$\mathbb{P}_{\theta \sim \nu}^{\pi}(\hat{x}_{\tau} \in \mathcal{X}_{\epsilon}(\theta)) = \mathbb{E}^{\pi}[\mathbb{E}^{\pi}[\mathbf{1}\{\hat{x}_{\tau} \in \mathcal{X}_{\epsilon}(\theta)\} | H_{\tau}]] = \mathbb{E}^{\pi}[q_{\tau}(H_{\tau}, \hat{x}_{\tau})].$$

852  $\square$

853 Therefore, we have that  $\mathbb{P}_{\theta \sim \nu}^{\pi}(I_{\tau}(H_{\tau}) \in \mathcal{X}_{\epsilon}(\theta)) = \mathbb{E}^{\pi}[q_{\tau}(H_{\tau}, I_{\tau}(H_{\tau}))]$ .

854 **Lagrangian dual.** Define, for  $\lambda \geq 0$ , the Lagrangian value

$$V_{\lambda}(\pi, I, \tau) := \mathbb{E}_{\theta \sim \nu}^{\pi}[\tau] + \lambda \left( (1 - \delta) - \mathbb{P}_{\theta \sim \nu}^{\pi}(I_{\tau}(H_{\tau}) \in \mathcal{X}_{\epsilon}(\theta)) \right).$$

855 Using Lemma 5, this can be written as

$$V_{\lambda}(\pi, I, \tau) = \lambda(1 - \delta) + \mathbb{E}^{\pi}[\tau - \lambda q_{\tau}(H_{\tau}, I_{\tau}(H_{\tau}))]. \quad (13)$$

856 The Lagrangian dual of (12) is then

$$\sup_{\lambda \geq 0} \inf_{\pi, I, \tau} V_{\lambda}(\pi, I, \tau). \quad (14)$$

857 **B.3.1 Optimal Inference Rule**

858 Fix  $\pi, \lambda$  and  $t$ . Among all measurable inference rules  $I_t : \mathcal{H}_t \rightarrow \mathcal{X}$ , the maximal probability of  
 859  $\epsilon$ -success at time  $t$  is achieved by any

$$I_t(h) \in \arg \max_{x \in \mathcal{X}} q_t(h, x),$$

860 equivalently,  $q_t(h, I_t(h)) = r_t(h)$  for  $\mathbb{P}_t^\pi$ -a.e.  $h$ . (see Proposition 1.)

861 Since  $\tau$  is adapted, plugging the optimal inference rule into Eq. (13) yields the simplified dual  
 862 objective

$$\sup_{\lambda \geq 0} \inf_{\pi, \tau} \lambda(1 - \delta) + \mathbb{E}^\pi [\tau - \lambda r_\tau(H_\tau)]. \quad (15)$$

863 **B.3.2 Stopping as an action (equivalence)**

864 The additional optimization over stopping rules can be avoided by introducing an additional stopping  
 865 action  $a_{\text{stop}}$ . Introduce an augmented action space  $\bar{\mathcal{A}} := \mathcal{A} \cup \{a_{\text{stop}}\}$ , where choosing  $a_{\text{stop}}$  terminates  
 866 interaction (no new observation is collected). Let  $\bar{\tau} := \inf\{t \geq 1 : A_t = a_{\text{stop}}\}$ .

867 **Lemma 6** (Embedding stopping times as a stop action). *For every triplet  $(\pi, I, \tau)$  with  $\tau < \infty$  a.s.,  
 868 there exists a policy  $\bar{\pi}$  on  $\bar{\mathcal{A}}$  such that, under  $\bar{\pi}$ , (i)  $\bar{\tau} = \tau$  a.s., and (ii) the stopped history  $H_{\bar{\tau}}$  has the  
 869 same distribution as  $H_\tau$  under  $\pi$ . In particular, for every  $\lambda \geq 0$ ,  $V_\lambda(\pi, I, \tau) = V_\lambda(\bar{\pi}, I, \bar{\tau})$ .*

870 *Proof.* Since  $\tau$  is a stopping time w.r.t.  $\mathcal{F}_t = \sigma(H_t)$ , the event  $\{\tau = t\}$  belongs to  $\mathcal{F}_t$ ; hence there  
 871 exists a measurable set  $S_t \subset \mathcal{H}_t$  such that  $\{\tau = t\} = \{H_t \in S_t\}$ . Define  $\bar{\pi}$  as follows: at time  $t$ ,  
 872 given history  $h \in \mathcal{H}_t$ ,

$$\bar{\pi}_t(a_{\text{stop}}|h) = \mathbf{1}\{h \in S_t\}, \quad \bar{\pi}_t(\cdot|h) = \pi_t(\cdot|h) \text{ on } \mathcal{A} \text{ when } h \notin S_t.$$

873 Then  $\{A_t = a_{\text{stop}}\} = \{H_t \in S_t\} = \{\tau = t\}$ , so  $\bar{\tau} = \tau$  a.s. Moreover, on the event  $\{\tau > t\}$  the  
 874 action distribution and observation kernel coincide with those under  $\pi$ , so the induced law of  $(H_t)_{t \leq \tau}$   
 875 is the same; in particular  $H_{\bar{\tau}}$  under  $\bar{\pi}$  has the same distribution as  $H_\tau$  under  $\pi$ . Hence, one can easily  
 876 show that the equality  $V_\lambda(\pi, I, \tau) = V_\lambda(\bar{\pi}, I, \bar{\tau})$  follows.  $\square$

877 **B.3.3 Optimal Policy**

878 Lemma 6 shows that (for fixed  $\lambda$ ) the inner problem in Eq. (15) can be viewed as an optimal-stopping  
 879 control problem on the augmented action space: each continuation step incurs unit cost, while  
 880 stopping at history  $h \in \mathcal{H}_t$  incurs terminal cost  $-\lambda r_t(h)$ .

881 Define the optimal cost-to-go (for fixed  $\lambda$ ) from a history  $h \in \mathcal{H}_t$  as

$$V_t^*(h; \lambda) := \inf_{\bar{\pi} = (\bar{\pi}_i)_{i \geq t}} \mathbb{E}_{\bar{\theta} \sim \nu} \left[ \sum_{s=t}^{\bar{\tau}-1} 1 - \lambda r_{\bar{\tau}}(H_{\bar{\tau}}) \mid H_t = h \right], \quad (16)$$

882 where the infimum is over policies on  $\bar{\mathcal{A}}$  and  $\bar{\tau}$  is the first time  $a_{\text{stop}}$  is chosen. Similarly to [53], we  
 883 can define the following optimal  $Q$ -functions

$$Q_{t, \text{stop}}^*(h; \lambda) := -\lambda r_t(h), \quad Q_{t, \text{cont}}^*(h, a; \lambda) := 1 + \mathbb{E} [V_{t+1}^*(H_{t+1}; \lambda) \mid H_t = h, A_t = a].$$

884 where the latter expectation is over the posterior mixture, defined as

$$\bar{P}_t(y' \in Y \mid H_t = h, A_t = a) = \int P_{\theta, t}(y' \in Y \mid H_t, A_t = a) R_t(d\theta \mid H_t = h), \quad \forall Y \in \mathcal{B}(\mathcal{Y}).$$

885 Furthermore, similarly to [53], a standard decomposition yields the Bellman optimality relation

$$V_t^*(h; \lambda) = \min \left\{ Q_{t, \text{stop}}^*(h; \lambda), \inf_{a \in \mathcal{A}} Q_{t, \text{cont}}^*(h, a; \lambda) \right\}. \quad (17)$$

886 However, in order to guarantee that the infimum,  $\inf_{a \in \mathcal{A}} Q_{t, \text{cont}}^*(h, a; \lambda)$  is attained, since  $\mathcal{A}$  is  
 887 compact, we need to guarantee that the  $Q$ -value is lower semicontinuous. We begin by showing that  
 888  $V_t^*$  is lower semicontinuous. To that aim, we need some results first. We begin by showing that the  
 889 mixture posterior is weakly continuous.

890 **Lemma 7** (Weak continuity of the mixture posterior). Fix  $t$  and  $h \in \mathcal{H}_t$ . Let  $R_t(\cdot|h)$  be the posterior  
 891 on  $\Theta$  and define the posterior predictive kernel

$$\bar{P}_t(\cdot|h, a) := \int_{\Theta} P_{\theta,t}(\cdot|h, a) R_t(d\theta|h).$$

892 Under Assumption 1,  $a \mapsto \bar{P}_t(\cdot|h, a)$  is weakly continuous. If in addition we assume Assumption 4,  
 893 then  $(h, a) \mapsto \bar{P}_t(\cdot|h, a)$  is jointly weakly continuous.

894 *Proof.* We prove weak continuity in the action first. Fix  $f \in C_b(\mathcal{Y})$  (continuous and bounded) and a  
 895 sequence  $(a_n)_n$  such that  $a_n \rightarrow a$ . For each  $\theta$ , define

$$g_n(\theta) := \int f(y) P_{\theta,t}(dy|h, a_n), \quad g(\theta) := \int f(y) P_{\theta,t}(dy|h, a).$$

896 By Assumption 1,

$$g_n(\theta) \rightarrow g(\theta).$$

897 Moreover,  $|\int f(y) P_{\theta,t}(dy|h, a_n)| \leq \|f\|_{\infty} < \infty$  for all  $\theta, n$ , and thus is also bounded. By  
 898 dominated convergence,

$$\begin{aligned} \int f(y) \bar{P}_t(dy|h, a_n) &= \int_{\Theta} g_n(\theta) R_t(d\theta | h), \\ &\rightarrow \int_{\Theta} \underbrace{\left( \int f(y) P_{\theta,t}(dy|h, a) \right)}_{=g(\theta)} R_t(d\theta | h), \\ &= \int f(y) \bar{P}_t(dy|h, a). \end{aligned} \quad (\text{Fubini-Tonelli})$$

899 Regarding the second part, it is less straightforward to prove. We assume Assumption 4. Recall that  
 900  $R_t(d\theta | h) = \frac{\ell_t(\theta, h)\nu(d\theta)}{\int_{\Theta} \ell_t(\theta, h)\nu(d\theta)}$ . We omit the normalization constant for simplicity, and simply include  
 901 it in  $\ell_t$ .

902 Define some sequence  $(h_n, a_n) \rightarrow (h, a)$ , and consider

$$\int f(y) \bar{P}_t(dy | h_n, a_n) = \int_{\Theta} \underbrace{\int f(y) P_{\theta,t}(dy | h_n, a_n) \ell_t(\theta, h_n)}_{=: G_f(h_n, a_n, \theta)} \nu(d\theta) = \int_{\Theta} G_f(h_n, a_n, \theta) \ell_t(\theta, h_n) \nu(d\theta).$$

903 Clearly  $G_f$  is bounded and  $\ell_t$  is jointly continuous by assumption. Since  $\Theta$  is compact, we have that  
 904 along any convergence sequence  $h_n \rightarrow h$ , the likelihood  $\ell_t$  is uniformly bounded. If  $G_f$  is jointly  
 905 continuous in  $(h, a)$ , then, by dominated convergence, we obtain

$$\int_{\Theta} G_f(h_n, a_n, \theta) \ell_t(\theta, h_n) \nu(d\theta) \rightarrow \int_{\Theta} G_f(h, a, \theta) \ell_t(\theta, h) \nu(d\theta) = \int f(y) \bar{P}_t(dy | h, a),$$

906 which shows the claim. Hence, we need to show that  $(h, a) \mapsto P_{\theta,t}(\cdot | h, a)$  is weakly continuous for  
 907 each  $\theta$ .

908 Using that  $dP_{\theta,t}(\cdot | h, a) = p_{\theta,t}(y | h, a)d\lambda$ , and recalling that by Assumption 4  $p_{\theta,t}$  is jointly  
 909 continuous. We have

$$\int f(y) P_{\theta,t}(dy | h_n, a_n) = \int f(y) p_{\theta,t}(y | h_n, a_n) d\lambda$$

910 using again compactness (of  $\mathcal{Y}$ ) and continuity of the arguments, we derive weak continuity of  
 911  $(h, a) \mapsto P_{\theta,t}(\cdot | h, a)$ , which concludes the proof.

912 □

913 The main technical challenge in the continuous case is showing that the infimum  $\inf_{a \in \mathcal{A}} Q_{t, \text{cont}}^*(h, a)$   
 914 in the Bellman equation is attained. In finite  $\mathcal{A}$  this is trivial; in compact continuous  $\mathcal{A}$  it requires  
 915 lower semicontinuity of  $a \mapsto Q_{t, \text{cont}}^*(h, a)$ . We establish this through a value-iteration construction

916 that propagates lower semicontinuity from the observation model through the posterior predictive  
 917 to the  $Q$ -function. The proof adapts the general template of negative dynamic programming [9, 27]  
 918 to our stop/continue structure, where the stopping payoff involves  $r_t(h) = \sup_x q_t(h, x)$  — itself a  
 919 supremum over a continuous set whose upper semicontinuity must be established first (Lemma 3).  
 920 We build a sequence of values  $W_t^{(n)}(h)$  and show that these are l.s.c. in  $h$ .

921 Starting from  $V_t^*$ , we construct  $W_t^*$  and build an increasing sequence  $W_t^{(n)}$  from below. We show  
 922 that each  $W_t^{(n)}$  is l.s.c. in  $h$  and that for each fixed  $h$ , the map  $a \mapsto Q_t^{(n)}(h, a)$  is l.s.c. on  $\mathcal{A}$ . We  
 923 then show that  $W_t^{(n)}$  approaches  $W_t^*$ . We conclude by showing that  $W_t^* = V^* + \lambda$ , proving that  $V_t^*$   
 924 is l.s.c.

925 Define  $W_t^{(0)}(h) := 0$  for all  $t, h$ . Recursively for  $n \geq 0$ , define

$$\begin{aligned} Q_t^{(n)}(h, a) &:= 1 + \int_{\mathcal{Y}} W_{t+1}^{(n)}(h, a, y) \bar{P}_t(dy | h, a), \\ C_t^{(n)}(h) &:= \inf_{a \in \mathcal{A}} Q_t^{(n)}(h, a), \\ W_t^{(n+1)}(h) &:= \min\{\lambda - \lambda r_t(h); C_t^{(n)}(h)\}. \end{aligned}$$

926 Also define  $W_t^*(h) = V_t^*(h; \lambda) + \lambda$ . We also define the operator  $\mathcal{T}_t$  as follows:

$$(\mathcal{T}_t u)(h) := \min \left\{ \lambda(1 - r_t(h)), 1 + \inf_{a \in \mathcal{A}} \int_{\mathcal{Y}} u(h, a, y) \bar{P}_t(dy | h, a) \right\},$$

927 therefore  $W_t^{(n+1)} = \mathcal{T}_t W_{t+1}^{(n)}$ . One can clearly show that the operator is monotone, as for  $u(h) \leq v(h)$   
 928 we have  $(\mathcal{T}_t u)(h) \leq (\mathcal{T}_t v)(h)$ . Therefore  $U_t(h) := \sup_{n \geq 0} W_t^{(n)}(h)$  exists pointwise, and satisfies  
 929  $0 \leq U_t(h) \leq \lambda$ .

930 **Lemma 8.** *Assume Assumption 2 and Assumption 4. For every  $t \in \mathbb{N}$ , the iterates  $W_t^{(n)}$  and  $C_t^{(n)}$   
 931 are lower semicontinuous, and for each fixed  $h$ ,  $a \mapsto Q_t^{(n)}(h, a)$  is lower semicontinuous on  $\mathcal{A}$  and  
 932 attains its minimum.*

933 *Proof.* Let  $g_t(h) = \lambda(1 - r_t(h))$  and fix  $t$ . We prove the argument by induction on  $n$ : for a fixed  $n$ ,  
 934  $W_s^{(n)}$  is l.s.c. for every  $s$ .

935 Base Step. Clearly, for all  $t$ ,  $W_t^{(0)}$  is l.s.c. , and since  $Q_t^{(0)}(h, a) = 1$  and  $C_t^{(0)}(h) = 1$ , we have  
 936  $W_t^{(1)}(h) = \min(g_t(h), 1)$  which is l.s.c. if  $-r_t$  is l.s.c. (which is, by Lemma 3).

937 Then, assume the induction hypothesis. We use that  $W_{t+1}^{(n)}$  is l.s.c. and prove that  $W_t^{(n+1)}$  is l.s.c.  
 938 We do so in 3 steps: first we prove that  $a \mapsto Q_t^{(n)}(h, a)$  is l.s.c. Then, we prove that  $h \mapsto C_t^{(n)}(h)$  is  
 939 l.s.c. Lastly we show that  $W_t^{(n+1)}$  is l.s.c.

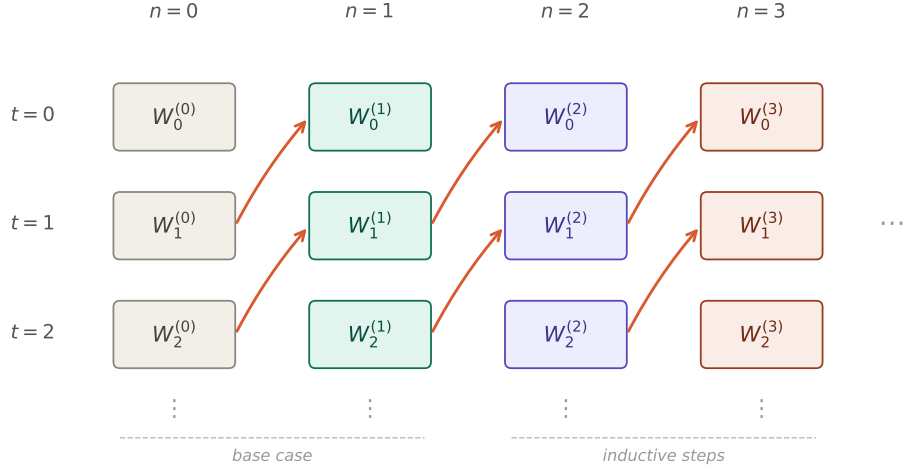
940 Step 1. Fix  $h$  and consider a sequence  $a_m \rightarrow a$ . Define  $\mu_m := \delta_{a_m} \otimes \bar{P}_t(\cdot | h, a_m)$  and  $\mu :=$   
 941  $\delta_a \otimes \bar{P}_t(\cdot | h, a)$ . By Lemma 7 we have that  $a \rightarrow \bar{P}_t(\cdot | h, a)$  is weakly continuous, and thus  
 942  $\mu_m \Rightarrow \mu$ . Since  $W_{t+1}^{(n)}$  is l.s.c. and bounded, the Portmanteau theorem yields

$$\liminf_{m \rightarrow \infty} \int_{\mathcal{Y} \times \mathcal{A}} W_{t+1}^{(n)}(h, a', y) d\mu_m \geq \int_{\mathcal{Y} \times \mathcal{A}} W_{t+1}^{(n)}(h, a', y) d\mu.$$

943 therefore  $a \mapsto Q_t^{(n)}(h, a)$  is l.s.c. on compact  $\mathcal{A}$ , and the infimum is attained.

944 Step 2. We now show that  $h \mapsto C_t^{(n)}(h)$  is l.s.c. Define a sequence  $(h_m)_m$  such that  $h_m \rightarrow h$ .  
 945 Choose a subsequence  $(h_{m_k})_k$  such that  $C_t^{(n)}(h_{m_k}) \rightarrow \liminf_{m \rightarrow \infty} C_t^{(n)}(h_m)$ . For each  $k$  choose  
 946  $a_{m_k} \in \arg \min_{a \in \mathcal{A}} Q_t^{(n)}(h_{m_k}, a)$ . Since  $\mathcal{A}$  is compact, by passing to a further subsequence if  
 947 necessary, we may assume  $a_{m_k} \rightarrow a$ . Define  $\mu_{m_k} := \delta_{(h_{m_k}, a_{m_k})} \otimes \bar{P}_t(\cdot | h_{m_k}, a_{m_k})$  and  $\mu :=$   
 948  $\delta_{(h, a)} \otimes \bar{P}_t(\cdot | h, a)$ . Since by assumption we have that  $(h, a) \mapsto \bar{P}_t(\cdot | h, a)$  is jointly weakly  
 949 continuous, then  $\mu_{m_k} \Rightarrow \mu$ . Similarly to before, we obtain

$$\liminf_{k \rightarrow \infty} Q_t^{(n)}(h_{m_k}, a_{m_k}) \geq Q_t^{(n)}(h, a) \geq C_t^{(n)}(h).$$



Each diagonal arrow (3-step chain)

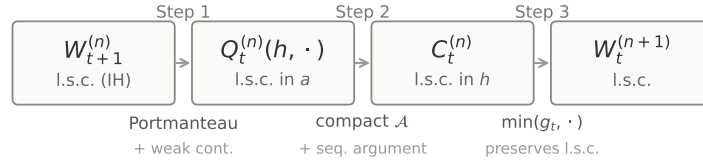


Figure 5: Induction diagram used in the proof of Lemma 8.

950 But  $Q_t^{(n)}(h_{m_k}, a_{m_k}) = C_t^{(n)}(h_{m_k})$ , therefore  $\liminf_{m \rightarrow \infty} C_t^{(n)}(h_m) = \lim_{k \rightarrow \infty} C_t^{(n)}(h_{m_k}) \geq$   
 951  $C_t^{(n)}(h)$ . Therefore  $C_t^{(n)}$  is l.s.c. in  $h$ .

952 Step 3. Lastly, consider  $W_t^{(n+1)} = \min(g_t(h), C_t^{(n)}(h))$ . Since both arguments are l.s.c., then  
 953  $W_t^{(n+1)}$  is l.s.c. Since  $t$  was arbitrary, the statements holds for all  $t$ .

954

□

955 Then, since each  $W_t^{(n)}$  is l.s.c., we get that  $U_t(h)$  is l.s.c. (arbitrary suprema of l.s.c. functions are  
 956 l.s.c.).

957 **Lemma 9.** Assume Assumption 2 and Assumption 4. For every  $t \in \mathbb{N}$ , define

$$Q_t^U(h, a) := 1 + \int_{\mathcal{Y}} U_{t+1}(h, a, y) \bar{P}_t(dy | h, a), \quad C_t^U(h) := \inf_{a \in \mathcal{A}} Q_t^U(h, a),$$

958 then  $U_t(h) = \min\{g_t(h), C_t^U(h)\}$ . Furthermore, we have that  $a \mapsto Q_t^U(h, \cdot)$  is lower semicontinu-  
 959 ous and  $Q_t^U(h, a)$  is jointly Borel measurable.

960 *Proof.* Let  $g_t(h) = \lambda(1 - r_t(h))$  and fix  $t$ . Since  $W_{t+1}^{(n)} \uparrow U_{t+1}$  and  $W_{t+1}^{(n)} \geq 0$ , by monotone  
 961 convergence for each  $(h, a)$  we have  $Q_t^{(n)} \uparrow Q_t^U$ . By Lemma 8, each  $Q_t^{(n)}(h, \cdot)$  is l.s.c., thus  $Q_t^U$  is  
 962 also l.s.c. being the suprema of l.s.c. functions. Furthermore, since each  $U_{t+1}(h, a, y)$  is l.s.c., and  
 963  $(h, a) \mapsto \bar{P}_t(\cdot | h, a)$  is a Borel kernel, then  $Q_t^U(h, a)$  is jointly Borel measurable.

964 Next, we show that  $\sup_n \inf_a Q_t^{(n)}(h, a) = \inf_a \sup_n Q_t^{(n)}(h, a) = C_t^U(h)$ . To show this, let  
 965  $m_n = \inf_a Q_t^{(n)}(h, a)$  and  $m = \sup_n m_n$ . First, note that by monotonicity for all  $N$  we have

966  $m_n = \inf_a Q_t^{(n)}(h, a) \leq \inf_a Q_t^U(h, a) = C_t^U(h)$ , and thus  $m \leq C_t^U(h)$ . Now, choose some  
 967 minimizers  $a_n \in \arg \min_a Q_t^{(n)}(h, a)$ . By compactness, there is some subsequence satisfying  
 968  $a_{n_k} \rightarrow a$ . Choose some integers  $i_0$ , then by monotonicity for all large  $k$  we have

$$Q_t^{(i_0)}(h, a_{n_k}) \leq Q_t^{(n_k)}(h, a_{n_k}) = m_{n_k}.$$

969 Therefore  $\liminf_k m_{n_k} = m \geq Q_t^{(i_0)}(h, a)$ . Since this holds for any  $i_0$ , we have  
 970  $m \geq \sup_{i_0} Q_t^{(i_0)}(h, a) = Q_t^U(h, a) \geq C_t^U(h)$ , which shows that  $\sup_n \inf_a Q_t^{(n)}(h, a) =$   
 971  $\inf_a \sup_n Q_t^{(n)}(h, a) = C_t^U(h)$ .

972 Hence, we have that  $U_t(h) = \sup_n W_t^{(n+1)}(h) = \sup_n \min\{g_t(h), C_t^{(n)}(h)\}$ . Since  $C_t^{(n)}(h) \uparrow$   
 973  $C_t^U(h)$ , we have  $U_t(h) = \min\{g_t(h), C_t^U(h)\}$ .  $\square$

974 Finally, we show that  $U_t$  is actually the true optimal value  $W_t^*$ .

975 **Theorem B.1.** *Assume Assumption 2 and Assumption 4. For every  $t \in \mathbb{N}$ , we have that  $U_t(h) =$*   
 976  *$V_t^*(h) + \lambda$ , and  $V_t^*$  is l.s.c.*

977 *Proof.* Let  $g_t(h) = \lambda(1 - r_t(h))$ . We first show that  $U_t$  is actually a lower bound on the cost of any  
 978 policy. Then we construct a policy that achieves equality.

979 *Step 1 (lower bound).* Consider any admissible policy  $\pi$ . At any history  $h$  at time  $t$ ,

- 980 1. if  $\pi$  stops, then  $U_t(h) \leq g_t(h)$  since  $U_t(h) = \min\{g_t(h), C_t^U(h)\}$ .
- 981 2. if  $\pi$  continues with some action  $a$ , then

$$U_t(h) \leq 1 + \int_{\mathcal{Y}} U_{t+1}(h, a, y) \bar{P}_t(dy | h, a),$$

982 since  $C_t^U(h) \leq Q_t^U(h, a)$ .

983 Let  $\tau$  be the first timestep the policy  $\pi$  decides to stop. Then, iterating up to  $\min(\tau, t + N)$

$$U_t(h) \leq \mathbb{E}^\pi \left[ \sum_{s=t}^{\min(\tau, t+N)-1} 1 + \mathbf{1}_{\{\tau \leq t+N\}} g_\tau(H_\tau) + \mathbf{1}_{\{\tau > t+N\}} U_{t+N}(H_{t+N}) \middle| H_t = h \right] =: J_N^\pi(h).$$

984 If the policy achieves infinite cost, then the inequality is true for all  $N$  since  $U_t(h)$  is bounded.  
 985 Then, consider the case where the cost is finite as  $N \rightarrow \infty$ . Since  $g_\tau \geq 0$ , this implies  $\mathbb{E}^\pi[\tau - t |$   
 986  $H_t = h] < \infty$ , therefore  $\tau < \infty$  almost surely. Consequently, by dominated convergence the  
 987 remainder term  $+\mathbf{1}_{\{\tau > t+N\}} U_{t+N}(H_{t+N})$  vanishes as  $N \rightarrow \infty$ . Therefore, letting  $N \rightarrow \infty$  we  
 988 obtain  $U_t(h) \leq \lim_{N \rightarrow \infty} J_N^\pi(h) := J^\pi(h)$ , and taking infimum over  $\pi$  we get  $U_t(h) \leq W_t^*(h)$ .

989 *Step 2 (equality).* From Lemma 9,  $Q_t^U(h, a)$  is l.s.c. on  $\mathcal{A}$  for each  $h$ , and  $\mathcal{A}$  is compact.  
 990 Furthermore  $Q_t^U(h, a)$  is jointly Borel-measurable. Then, there exists a measurable selector  
 991  $a_t^*(h) \in \arg \min_a Q_t^U(h, a)$  [17]. Then both  $C_t^U(h) = Q_t^U(h, a_t^*(h))$  is Borel measurable. Consider  
 992 then a policy  $\pi^*$  that stops at  $h$  if  $g_t(h) \leq C_t^U(h)$ , and otherwise continue with  $a_t^*(h)$ . Denote by  $\tau^*$   
 993 this stopping rule. Along this policy, the Bellman minimum is attained with equality at every step,  
 994 therefore

$$U_t(h) = \mathbb{E}^{\pi^*} \left[ \sum_{s=t}^{\min(\tau^*, t+N)-1} 1 + \mathbf{1}_{\{\tau^* \leq t+N\}} g_{\tau^*}(H_{\tau^*}) + \mathbf{1}_{\{\tau^* > t+N\}} U_{t+N}(H_{t+N}) \middle| H_t = h \right].$$

995 Since  $U_t(h)$  is bounded, we have

$$\lambda \geq U_t(h) \geq \mathbb{E}^{\pi^*}[\min(\tau^*, t+N) - t | H_t = h].$$

996 Letting  $N \rightarrow \infty$ , we obtain  $\mathbb{E}^{\pi^*}[\tau - t | H_t = h] \leq \lambda$ . Therefore the remainder term  $\mathbf{1}_{\{\tau^* > t+N\}} U_{t+N}$   
 997 vanishes as  $N \rightarrow \infty$  by dominated convergence since  $U_{t+N} \leq \lambda$  and  $\tau^* < \infty$  almost surely.  
 998 Therefore we obtain  $U_t(h) = J^{\pi^*}(h)$ . Since  $W^*(h)$  is the minimal cost we obtain  $W_t^*(h) \leq U_t(h)$ ,  
 999 but from the previous step we also have  $U_t(h) \leq W_t^*(h)$ . Therefore  $U_t(h) = W_t^*(h) = V_t^*(h; \lambda) + \lambda$ .  
 1000 Finally,  $V_t^*$  is l.s.c. since  $U_t$  is l.s.c.  $\square$

1001 Hence, we conclude with the following result

1002 **Proposition 2** (Actor over  $\mathcal{A}$  and stopping action). *Assume Assumption 2 and Assumption 4. and*  
 1003 *let  $a^*(h) \in \arg \min_{a \in \mathcal{A}} Q_{t,\text{cont}}^*(h, a)$ . Then an optimal policy can be implemented by: (i) selecting*  
 1004  *$a^*(h)$  as the continuation action, and (ii) stopping iff  $Q_{t,\text{stop}}^*(h) \leq Q_{t,\text{cont}}^*(h, a^*(h))$ .*

1005 *Proof.* By (17), at each history  $h$  the optimal action is whichever attains the minimum between  
 1006  $Q_{t,\text{stop}}^*(h)$  and  $\inf_{a \in \mathcal{A}} Q_{t,\text{cont}}^*(h, a)$ . If  $a^*(h)$  attains the infimum over  $\mathcal{A}$ , then the comparison  
 1007  $Q_{t,\text{stop}}^*(h) \leq Q_{t,\text{cont}}^*(h, a^*(h))$  is equivalent to  $Q_{t,\text{stop}}^*(h) \leq \inf_a Q_{t,\text{cont}}^*(h, a)$ , i.e., stopping is  
 1008 optimal. Otherwise continuing with  $a^*(h)$  is optimal.

1009 □

1010 This last result shows that an algorithm can act in the augmented space  $\bar{\mathcal{A}}$  while never explicitly  
 1011 parameterizing a “stop action” in the actor.

1012 Hence an optimal policy can be implemented by:

- 1013 1. choose a continuation action  $a^*(h) \in \arg \min_{a \in \mathcal{A}} Q_{t,\text{cont}}^*(h, a)$  (this is the actor over  $\mathcal{A}$   
 1014 only);
- 1015 2. stop if  $Q_{t,\text{stop}}^*(h) \leq Q_{t,\text{cont}}^*(h, a^*(h))$ , otherwise continue with  $a^*(h)$ .

1016 This is mathematically equivalent to having a single policy over  $\bar{\mathcal{A}}$  that selects  
 1017  $\arg \min\{Q_{t,\text{stop}}^*(h), \min_a Q_{t,\text{cont}}^*(h, a)\}$ , but it decomposes the decision into (i) a continu-  
 1018 ous control over  $\mathcal{A}$  and (ii) an optimal-stopping comparison against a scalar stop value. Therefore,  
 1019 one can learn a  $Q$ -value for the stop decision and comparing it to the  $Q$ -value of the actor-chosen  
 1020 action: this is consistent with the optimal control structure of (17), provided the actor approximates  
 1021 the minimizer of  $Q_{t,\text{cont}}^*(h, \cdot)$  and the critic estimates are calibrated.

### 1022 B.3.4 Reward Shaping and Removal of the Stop Action

1023 The stop action need not be explicitly parameterized by the actor. Indeed, for fixed  $\lambda > 0$ , the  
 1024 Lagrangian objective

$$\mathbb{E}^\pi \left[ \sum_{s=t}^{\tau-1} 1 - \lambda r_\tau(H_\tau) \mid H_t = h \right]$$

1025 is equivalent, up to an additive constant independent of the policy, to maximizing the shaped return  
 1026 obtained from the one-step reward

$$R_s^{\text{sh}} := -1 + \lambda (r_{s+1}(H_{s+1}) - r_s(H_s)).$$

1027 Thus the actor may be restricted to continuation actions  $a \in \mathcal{A}$ , while the stopping decision is  
 1028 implemented by comparing the best continuation advantage against the stop value, which is zero in  
 1029 the shaped formulation. Equivalently, using the rescaled reward

$$-\frac{1}{\lambda} + r_{s+1}(H_{s+1}) - r_s(H_s)$$

1030 gives the same optimal policies.

1031 **Lemma 10** (Reward shaping). *Fix  $\lambda > 0$ . For a continuation policy  $\pi$  and stopping time  $\tau$ , define*  
 1032 *the original Lagrangian cost*

$$J_t^\lambda(h; \pi, \tau) := \mathbb{E}^\pi \left[ \sum_{s=t}^{\tau-1} 1 - \lambda r_\tau(H_\tau) \mid H_t = h \right],$$

1033 *and the shaped return*

$$G_t^\lambda(h; \pi, \tau) := \mathbb{E}^\pi \left[ \sum_{s=t}^{\tau-1} R_s^{\text{sh}}(H_s, A_s, H_{s+1}) \mid H_t = h \right].$$

1034 *Then*

$$\arg \min_{\pi, \tau} J_t^\lambda(h; \pi, \tau) = \arg \max_{\pi, \tau} G_t^\lambda(h; \pi, \tau).$$

1035 *Proof.* First, the shaped reward telescopes:

$$\sum_{s=t}^{\tau-1} R_s^{\text{sh}}(H_s, A_s, H_{s+1}) = \sum_{s=t}^{\tau-1} [\lambda(r_{s+1}(H_{s+1}) - r_s(H_s)) - 1] = \lambda r_\tau(H_\tau) - \lambda r_t(H_t) - \sum_{s=t}^{\tau-1} 1.$$

1036 Conditioning on  $H_t = h$  gives

$$G_t^\lambda(h; \pi, \tau) = \mathbb{E}^\pi \left[ \lambda r_\tau(H_\tau) - \sum_{s=t}^{\tau-1} 1 \mid H_t = h \right] - \lambda r_t(h).$$

1037 Since

$$J_t^\lambda(h; \pi, \tau) = \mathbb{E}^\pi \left[ \sum_{s=t}^{\tau-1} 1 - \lambda r_\tau(H_\tau) \mid H_t = h \right],$$

1038 we obtain

$$G_t^\lambda(h; \pi, \tau) = -J_t^\lambda(h; \pi, \tau) - \lambda r_t(h).$$

1039 The term  $-\lambda r_t(h)$  does not depend on  $(\pi, \tau)$ , so maximizing  $G_t^\lambda$  is equivalent to minimizing  $J_t^\lambda$ .

1040  $\square$

1041 Using this reward shaping, we do not need to learn the stopping  $Q$ -value, as shown in the next  
1042 proposition. The actor only needs to output a continuation action  $a \in \mathcal{A}$ . The stopping rule is a scalar  
1043 gate:

$$\text{continue iff } Q_{t,\text{cont}}^{\text{sh}}(h, a_t(h)) > 0.$$

1044 Equivalently, the stop value in the shaped formulation is identically zero, and the state-value target is

$$S_t(h) = \max\{0, Q_{t,\text{cont}}^{\text{sh}}(h, a_t(h))\}.$$

1045 **Proposition 3** (Reward shaping does not require learning a stopping  $Q$ -value). *Let*

$$S_t^*(h) := \sup_{\pi, \tau} G_t^\lambda(h; \pi, \tau)$$

1046 *be the optimal shaped value. Then*

$$S_t^*(h) = \max \left\{ 0, \sup_{a \in \mathcal{A}} \mathbb{E} \left[ -1 + \lambda(r_{t+1}(H_{t+1}) - r_t(h)) + S_{t+1}^*(H_{t+1}) \mid H_t = h, A_t = a \right] \right\}.$$

1047 *Define the shaped continuation  $Q$ -value*

$$Q_{t,\text{cont}}^{\text{sh}}(h, a) := -1 + \mathbb{E} \left[ \lambda(r_{t+1}(H_{t+1}) - r_t(h)) + S_{t+1}^*(H_{t+1}) \mid H_t = h, A_t = a \right].$$

1048 *If the supremum over  $a \in \mathcal{A}$  is attained by a measurable selector  $a_t^*(h) \in \arg \max_{a \in \mathcal{A}} Q_{t,\text{cont}}^{\text{sh}}(h, a)$ ,*  
1049 *then an optimal policy is implemented by*

$$\text{stop at } h \iff Q_{t,\text{cont}}^{\text{sh}}(h, a_t^*(h)) \leq 0,$$

1050 *and otherwise continuing with  $a_t^*(h)$ .*

1051 *Proof.* Next, in the shaped formulation, stopping immediately produces the empty sum and hence  
1052 shaped return 0. If instead the learner continues with action  $a$ , it receives the immediate shaped  
1053 reward

$$\lambda(r_{t+1}(H_{t+1}) - r_t(h)) - 1$$

1054 and then obtains the optimal future shaped value  $S_{t+1}^*(H_{t+1})$ . Therefore

$$S_t^*(h) = \max \left\{ 0, \sup_{a \in \mathcal{A}} Q_{t,\text{cont}}^{\text{sh}}(h, a) \right\}.$$

1055 Thus stopping is optimal exactly when

$$\sup_{a \in \mathcal{A}} Q_{t,\text{cont}}^{\text{sh}}(h, a) \leq 0.$$

1056 If  $a_t^*(h)$  attains the supremum, this is equivalent to

$$Q_{t,\text{cont}}^{\text{sh}}(h, a_t^*(h)) \leq 0.$$

1057 Otherwise, if

$$Q_{t,\text{cont}}^{\text{sh}}(h, a_t^*(h)) > 0,$$

1058 continuing with  $a_t^*(h)$  attains the continuation branch and is optimal.  $\square$

1059 **B.4 Fixed-confidence setting:  $(\epsilon, \delta)$ -correctness of dual-optimal points**

1060 In this section we provide fixed-confidence guarantees for the continuous ICPE setting. Compared to  
 1061 the correctness argument in [53], our proof differs in three aspects, each addressing a limitation of  
 1062 the finite-space analysis.

- 1063 1. **Posterior object.** The finite ICPE framework uses  $\mathbb{P}(x^* = x \mid H_t)$ , which is ill-  
 1064 defined in continuous  $\mathcal{X}$ . We replace it with the posterior success probability  $r_t(h) =$   
 1065  $\sup_{x \in \mathcal{X}} \mathbb{P}(L_\theta(x) \leq \epsilon \mid H_t = h)$ . This is the natural continuous analogue and the only  
 1066 change forced by the setting.
- 1067 2. **Non-singleton dual optima (main technical improvement).** Russo et al. [53] assume the  
 1068 dual-optimal policy set  $\mathcal{S}(\lambda)$  is a singleton for each  $\lambda$ . We relax this to local closedness of  
 1069 the near-optimal cost-reward set  $\mathcal{K}_{\epsilon_0}(\lambda^*)$  (Assumption 6). This is meaningful in continuous  
 1070 spaces where the policy class is richer and uniqueness is harder to verify. The key proof step  
 1071 is showing that if all policies in  $\mathcal{K}_{\epsilon_0}(\lambda^*)$  have  $\rho < 1 - \delta$ , then a uniform gap propagates to  
 1072 all near-optimal policies (Step 2 below), which requires the closedness assumption in an  
 1073 essential way.
- 1074 3. **Subdifferential argument (alternative proof technique).** Rather than deriving a contra-  
 1075 diction from monotonicity of the optimal cost in  $\lambda$  (as in [53]), we use the subdifferential  
 1076 characterization of Hantoute and López [25] directly. This avoids the monotonicity argument  
 1077 entirely: we show that every subgradient of  $f(\lambda) = -g(\lambda)$  at  $\lambda^*$  is strictly negative, contra-  
 1078 dicting  $0 \in \partial f(\lambda^*)$ . This argument works with multiple dual optima without modification.

1079 We begin by stating our assumptions. To do so, let  $\Pi$  denote the class of admissible policies on the  
 1080 augmented action space  $\bar{\mathcal{A}}$ , and for each  $\pi$  let  $\tau_\pi$  be the corresponding stopping time, i.e., the first  
 1081 time the policy chooses the stop action. Define

$$\mathcal{T} := \{\pi \in \Pi : \mathbb{E}^\pi[\tau_\pi] < \infty\},$$

1082 be the set of admissible policies with finite expected sample complexity.

1083 **Assumption 5** (Strict feasibility). *We assume there exists a policy  $\pi \in \mathcal{T}$  such that*

$$\mathbb{E}^\pi[r_{\tau_\pi}(H_{\tau_\pi})] > 1 - \delta.$$

1084 We also make the following assumption of closedness on the attainable cost-reward set. For any  
 1085  $\pi \in \mathcal{T}$  define  $c(\pi) := \mathbb{E}^\pi[\tau_\pi]$  and  $\rho(\pi) := \mathbb{E}^\pi[r_{\tau_\pi}(H_{\tau_\pi})]$ . Define also

$$g(\lambda) = \inf_{\pi \in \Pi: c(\pi) < \infty} c(\pi) + \lambda(1 - \delta - \rho(\pi)), \quad \mathcal{S}(\lambda) = \arg \min_{\pi \in \Pi: c(\pi) < \infty} c(\pi) + \lambda(1 - \delta - \rho(\pi)),$$

1086 and the set of attainable cost-rewards:

$$\mathcal{K} := \{(c(\pi), \rho(\pi)) : \pi \in \mathcal{T}\} \subset [0, \infty) \times [0, 1].$$

1087 We have the following properties on  $g(\lambda)$ .

1088 **Lemma 11.** *Under Assumption 5,  $g(\lambda)$  is finite, concave on  $[0, \infty)$  and attains a maximum on  $[0, \infty)$ .*

1089 *Proof.* For simplicity, we write the dual value directly on the attainable cost-reward set:

$$g(\lambda) = \inf_{(c, \rho) \in \mathcal{K}} c + \lambda(1 - \delta - \rho),$$

1090 Then,  $g(\lambda) \geq -\delta\lambda$ . Let  $(c, \rho)$  be a point satisfying Assumption 5. Then  $\rho > 1 - \delta$ , and thus  
 1091  $g(\lambda) \leq c - \lambda\eta$  for  $\eta = -(1 - \delta - \rho)$ . Since  $\eta > 0$ , and  $-\lambda\delta \leq g(\lambda) \leq c - \lambda\eta$ , then  $g(\lambda) \rightarrow -\infty$   
 1092 as  $\lambda \rightarrow \infty$ .

1093 Finally, since  $\lambda \mapsto c + \lambda(1 - \delta - \rho)$  is affine, the infimum of any family of concave functions is  
 1094 concave, therefore  $g(\lambda)$  is concave. Then, since  $g$  is finite and concave, it's continuous on  $(0, \infty)$   
 1095 and u.s.c. in 0. Therefore  $g$  is u.s.c. on  $[0, \infty)$ , and since  $\lim_{\lambda \rightarrow \infty} g(\lambda) \rightarrow -\infty$ , there exists  $\lambda^* \in [0, \infty)$   
 1096 such that  $g(\lambda^*) = \max_{\lambda \geq 0} g(\lambda)$ .  $\square$

1097 We impose then the following assumption.

1098 **Assumption 6** (Closed optimal cost-reward set). Let  $\mathcal{G} = \arg \max_{\lambda} g(\lambda)$ . We assume that for every  
 1099 maximizer  $\lambda^* \in \mathcal{G}$  there exists  $\epsilon_0 > 0$  such that

$$\mathcal{K}_{\epsilon}(\lambda^*) := \{(c, \rho) \in \mathcal{K} : c + \lambda^*(1 - \delta - \rho) \leq g(\lambda^*) + \epsilon\}$$

1100 is closed for all  $\epsilon \in (0, \epsilon_0]$

1101 Essentially, this geometric assumption is used to relax the assumption used in [53] that  $\mathcal{S}(\lambda)$  is a  
 1102 singleton for each  $\lambda$ , and allows us to generalize the argument to multiple optimal dual policies.

1103 To verify correctness, we use the fact that the sub-gradient of the optimal value of the dual problem is  
 1104 non-decreasing. To show this result, we employ the following proposition from [25] (see Prop. 3.1  
 1105 therein), which characterizes the subdifferential of the supremum of a family of affine functions.

1106 **Proposition 4** (Subdifferential of the supremum of affine functions [25]). *Given a non-empty set*  
 1107  *$\{(a_t, b_t) : t \in \mathcal{T}\} \subset \mathbb{R}^2$ , and the supremum function  $f(x) : \mathbb{R} \rightarrow \mathbb{R} \cup \{\infty\}$*

$$f(x) = \sup\{a_t x - b_t : t \in \mathcal{T}\},$$

1108 for every  $x \in \text{dom} f$  we have

$$\partial f(x) = \bigcap_{\epsilon > 0} \text{cl}(\text{conv}\{a_t : t \in \mathcal{T}_{\epsilon}(x)\} + B(x))$$

1109 with

$$\mathcal{T}_{\epsilon}(x) := \{t \in \mathcal{T} : a_t x - b_t \geq f(x) - \epsilon\},$$

1110 and

$$B(x) := \{y \in \mathbb{R} : (y, yx) \in (\overline{\text{conv}}\{(a_t, b_t) : t \in \mathcal{T}\})_{\infty}\},$$

1111 where  $C_{\infty}$  is the recession cone of a set  $C$  and  $\overline{\text{conv}}(\cdot)$  denotes the closed convex hull of a set. In  
 1112 particular, if  $x \in \text{int}(\text{dom} f)$  we have

$$\partial f(x) = \bigcap_{\epsilon > 0} \overline{\text{conv}}\{a_t : t \in \mathcal{T}_{\epsilon}(x)\}.$$

1113 This last proposition permits us to define the subdifferential of the supremum of affine functions, and,  
 1114 as we see now, we can also find a lower bound on any subdifferential  $d \in \partial f(x)$ .

1115

1116 We are now ready to prove  $(\epsilon, \delta)$ -PAC guarantees for ICPE in the continuous setting.

1117 **Theorem B.2** ( $(\epsilon, \delta)$ -correctness in the continuous ICPE case). *Under Assumption 5 and Assumption 6*  
 1118 *for all  $\lambda^* \in \mathcal{G} := \arg \max_{\lambda} g(\lambda)$  there exists a dual-optimal policy  $\pi^* \in \mathcal{S}(\lambda^*)$  such that*

$$\mathbb{E}^{\pi^*}[r_{\tau_{\pi^*}}(H_{\tau_{\pi^*}})] \geq 1 - \delta, \quad \text{and} \quad \mathbb{P}^{\pi^*}(L_{\theta}(I_{\tau^*}^*(H_{\tau^*})) \leq \epsilon) \geq 1 - \delta.$$

1119 where  $I_t^*(h) \in \arg \max_h q_t(h, x)$  is the optimal inference selector.

1120 *Proof.* For simplicity, we write the dual value directly on the attainable cost-reward set:

$$g(\lambda) = \inf_{(c, \rho) \in \mathcal{K}} c + \lambda(1 - \delta - \rho),$$

1121 and the active set  $M(\lambda) := \arg \min_{(c, \rho) \in \mathcal{K}} c + \lambda(1 - \delta - \rho)$ .

1122 Consider Lemma 11, then there exists  $\lambda^* \in \arg \max_{\lambda} g(\lambda)$  such that  $\lambda \in [0, \infty)$ . We now prove that  
 1123  $M(\lambda^*)$  contains at-least a point  $(c, \rho)$  with reward at-least  $1 - \delta$ .

1124 By contradiction, **(HYP)** assume that all points  $(c, \rho) \in M(\lambda^*)$  satisfy  $\rho < 1 - \delta$ .

1125

1126 *First step: there exists  $\alpha > 0$  such that  $\rho \leq 1 - \delta - \alpha$  for optimal policies.* Under **(HYP)**, we claim  
 1127 that there exists  $\alpha > 0$  such that  $\rho \leq 1 - \delta - \alpha$  for all  $(c, \rho) \in M(\lambda^*)$ . First, note that by  
 1128 Assumption 6 there exists  $\epsilon_0$  such that  $M(\lambda^*) = \bigcap_{\epsilon \in (0, \epsilon_0]} \mathcal{K}_{\epsilon}(\lambda^*)$ . Note that each  $\mathcal{K}_{\epsilon}(\lambda^*)$  is closed,  
 1129 non-empty (since  $g$  is finite), and since  $\rho \in [0, 1]$  and  $c \leq g(\lambda^*) + \epsilon + \lambda^* \delta$ , it's bounded. Thus each  
 1130  $\mathcal{K}_{\epsilon}(\lambda^*)$  is compact: therefore, also the intersection  $M(\lambda^*)$  is compact.

1131 Since  $M(\lambda^*)$  is compact, we have that the maximum is achieved, and since  $(c, \rho) \mapsto \rho$  is continuous,  
 1132 we set  $\rho_{\max} = \max\{\rho : (c, \rho) \in M(\lambda^*)\}$  and  $\alpha = 1 - \delta - \rho_{\max} > 0$ .  
 1133

1134 *Second step: near optimal policies satisfy  $\rho \leq 1 - \delta - \alpha/2$ .* Always under **(HYP)**, we now claim  
 1135 that for near-optimal points we actually have  $\rho \leq 1 - \delta - \alpha/2$ . Suppose that this is not true: then, by  
 1136 Assumption 6 there exists  $\epsilon_0$  such that for every  $\epsilon \in (0, \epsilon_0]$  there exists a point  $(c, \rho)$  in  $\mathcal{K}_\epsilon(\lambda^*)$  such  
 1137 that  $\rho > 1 - \delta - \alpha/2$ . Let  $\epsilon_n = 1/n$  and consider such a sequence  $(c_n, \rho_n) \in \mathcal{K}_{\epsilon_n}(\lambda^*)$  satisfying  
 1138  $\rho_n > 1 - \delta - \alpha/2$  for all  $n$ . Therefore, for each  $n$  we have that

$$c_n + \lambda^*(1 - \delta - \rho_n) \leq g(\lambda^*) + \frac{1}{n}.$$

1139 Since  $\rho_n \in [0, 1]$  and  $c_n \leq g(\lambda^*) + 1/n + \lambda^*\delta$ , thus bounded, we have that the sequence  $(c_n, \rho_n)$  is  
 1140 bounded. Since for each  $\epsilon_n$  the set  $\mathcal{K}_{\epsilon_n}(\lambda^*)$  is closed, we can take a convergent subsequence that  
 1141 satisfies  $(c_n, \rho_n) \rightarrow (c, \rho)$  for some  $(c, \rho) \in \mathcal{K}_{\epsilon_0}$ . Therefore, we have that

$$\lim_{n \rightarrow \infty} c_n + \lambda^*(1 - \delta - \rho_n) \leq g(\lambda^*) + 1/n \implies c + \lambda^*(1 - \delta - \rho) \leq g(\lambda^*).$$

1142 However, for any  $(c, \rho)$  we also have  $g(\lambda^*) \leq c + \lambda^*(1 - \delta - \rho)$ , thus  $c + \lambda^*(1 - \delta - \rho) = g(\lambda^*)$ ,  
 1143 while also having  $\rho \geq 1 - \delta - \alpha/2$ , which contradicts  $\rho \leq 1 - \delta - \alpha$ . Therefore, for all  $\epsilon_0$  close  
 1144 points we have  $\rho \leq 1 - \delta - \alpha/2$ .

1145 We now distinguish the cases  $\lambda^* = 0$  and  $\lambda^* > 0$ .  
 1146

1147 *Third step: case  $\lambda^* = 0$  is not optimal.* Under **(HYP)**, assume  $\lambda^* = 0$  is optimal. We show that this  
 1148 is not the case, and there exists  $\lambda : g(\lambda) > g(0)$ . We proceed by showing there exists  $\lambda > 0$  small  
 1149 such that for any  $(c, \rho) \in \mathcal{K}$  we have  $c + \lambda(1 - \delta - \rho) > g(0)$ , and thus there exists  $\lambda > 0$  such that  
 1150  $g(\lambda) > g(0)$ .

1151 Let  $\lambda > 0$  to be chosen. By Assumption 6 and the previous step there exists  $\epsilon_0$  such that for all  
 1152  $(c, \rho) \in \mathcal{K}_\epsilon(0)$  with  $\epsilon \in (0, \epsilon_0]$  we have  $1 - \delta - \alpha/2 \geq \rho$  and  $c \leq g(0) + \epsilon$ . Consider then the  
 1153 following two cases:

1154 1. Case where  $(c, \rho)$  satisfies  $c \geq g(0) + \epsilon_0$  (far away point). Since  $\rho \in [0, 1]$ , we have  
 1155  $1 - \delta - \rho \in [-1, 1]$ , and thus  $\lambda(1 - \delta - \rho) \geq -\lambda$ . Hence

$$g(0) + \epsilon_0 - \lambda \leq c + \lambda(1 - \delta - \rho).$$

1156 2. Case where  $(c, \rho)$  satisfies  $c \leq g(0) + \epsilon_0$  (near optimal point). Then, in this case we have  
 1157 that  $(1 - \delta - \rho) \geq \alpha/2$ , thus, using that  $g(0) \leq c$ , combining the two inequalities we have

$$g(0) + \lambda\alpha/2 \leq c + \lambda(1 - \delta - \rho).$$

1158 Then, choose  $\lambda \in (0, \epsilon_0/2)$ : we obtain  $-\lambda > -\epsilon_0/2$ , and  $\lambda > 0$ , thus

$$g(0) < g(0) + \frac{\epsilon_0}{2} < c + \lambda(1 - \delta - \rho), \quad g(0) < g(0) + \lambda\frac{\alpha}{2} \leq c + \lambda(1 - \delta - \rho).$$

1159 Therefore  $g(0) < c + \lambda(1 - \delta - \rho)$  for all  $(c, \rho) \in \mathcal{K}$ . This shows that  $g(\lambda) > g(0)$ , and contradicts  
 1160 the optimality of  $\lambda^* = 0$ .

1161 *Fourth step: case  $\lambda^* > 0$  is not optimal.* Define the function  $f(\lambda) = -g(\lambda)$ . From Lemma 11 then  
 1162  $f$  is convex. Since  $\lambda^* > 0$ , it lies in the interior of  $\text{dom}(f)$ . Then, by Proposition 4, we have

$$\partial f(\lambda^*) = \bigcap_{\epsilon \in (0, \epsilon_0]} \overline{\text{conv}}\{\rho - (1 - \delta) : (c, \rho) \in \mathcal{K}_\epsilon(\lambda^*)\}$$

1163 By step 2 we have  $1 - \delta - \alpha/2 \geq \rho$ , therefore  $\rho - (1 - \delta) \leq -\alpha/2$ , implying that

$$d \leq -\frac{\alpha}{2} \quad \forall d \in \partial f(\lambda^*).$$

1164 But  $\lambda^*$  minimizes the convex function  $f$  and is an interior point of the domain, so necessarily  
 1165  $0 \in \partial f(\lambda^*)$ , which is a contradiction.

1166 *Last step.* Since both cases are impossible, we conclude there exists  $(c, \rho) \in M(\lambda^*)$  with  $\rho \geq 1 - \delta$ .  
 1167 Hence, there exists a dual optimal policy  $\pi \in \mathcal{S}(\lambda^*)$  such that  $c(\pi) = c$  and  $\rho(\pi) \geq 1 - \delta$ .  
 1168 Lastly, since  $I_t^*$  attains the supremum in the definition of  $r_t(h)$ , we have  $r_t(h) = q_t(h, I_t^*(h)) =$   
 1169  $\mathbb{P}(L_\theta(I_t^*(h)) \leq \epsilon \mid H_t = h)$ . Thus, we get

$$1 - \delta \leq \rho(\pi) = \mathbb{E}^\pi[r_{\tau_\pi}(H_{\tau_\pi})] = \mathbb{E}_{H_{\tau_\pi}}^\pi \left[ \mathbb{P} \left( L_\theta \left( I_{\tau_\pi}^* \left( H_{\tau_\pi} \right) \right) \leq \epsilon \mid H_{\tau_\pi} \right) \right] = \mathbb{P} \left( L_\theta \left( I_{\tau_\pi}^* \left( H_{\tau_\pi} \right) \right) \leq \epsilon \right).$$

1170

□

1171 **Remark 2.** *The above result yields a Bayesian fixed-confidence guarantee under the prior  $\nu$ . It*  
 1172 *is therefore the natural continuous counterpart of the fixed-confidence correctness statement in the*  
 1173 *discrete ICPE setting.*

## 1174 B.5 Fixed-confidence setting: zero duality gap via perturbation values

1175 In the previous sections we reduced the Bayesian fixed-confidence problem to an optimal-stopping  
 1176 problem with terminal posterior-success reward. We now study when the corresponding Lagrangian  
 1177 relaxation is exact. The main point of this section is that, although the policy space is infinite-  
 1178 dimensional, the duality question can be analyzed through the two-dimensional attainable cost-reward  
 1179 set

$$\mathcal{K} := \{(c(\pi), \rho(\pi)) : \pi \in \mathcal{T}\} \subset [0, \infty) \times [0, 1],$$

1180 where

$$c(\pi) := \mathbb{E}^\pi[\tau_\pi], \quad \rho(\pi) := \mathbb{E}^\pi[r_{\tau_\pi}(H_{\tau_\pi})].$$

1181 The zero duality gap result uses a standard one-dimensional perturbation argument from convex  
 1182 duality [46, 16, 47, 10]. While the general technique idea is not new, what is new is its application to  
 1183 this specific problem studied in this manuscript, where the attainable cost-reward set  $\mathcal{K}$  arises from  
 1184 an optimal-stopping control problem over continuous spaces.

1185 After optimizing over the inference rule, the reduced primal problem is

$$P^* := \inf_{\pi \in \mathcal{T}} c(\pi) \quad \text{s.t.} \quad \rho(\pi) \geq 1 - \delta. \quad (18)$$

1186 Equivalently,  $P^* = \inf\{c : (c, \rho) \in \mathcal{K}, \rho \geq 1 - \delta\}$ . The associated Lagrangian dual value is

$$D^* := \sup_{\lambda \geq 0} g(\lambda), \quad g(\lambda) := \inf_{\pi \in \mathcal{T}} \{c(\pi) + \lambda(1 - \delta - \rho(\pi))\}. \quad (19)$$

1187 Equivalently,  $g(\lambda) = \inf_{(c, \rho) \in \mathcal{K}} \{c + \lambda(1 - \delta - \rho)\}$ , and the multiplier satisfies  $\lambda \geq 0$ . Clearly, by  
 1188 Lagrangian duality, we have that the weak duality easily holds  $D^* \leq P^*$ .

1189 **Assumptions.** We state some assumptions. The first one is a convexification assumption on the  
 1190 policy class.

1191 **Assumption 7 (Time-sharing).** *For every  $\pi_0, \pi_1 \in \mathcal{T}$  and every  $\alpha \in [0, 1]$ , there exists a policy*  
 1192  *$\pi_\alpha \in \mathcal{T}$  such that*

$$c(\pi_\alpha) = \alpha c(\pi_1) + (1 - \alpha)c(\pi_0), \quad \rho(\pi_\alpha) = \alpha \rho(\pi_1) + (1 - \alpha)\rho(\pi_0).$$

1193 This assumption is natural for randomized sequential policies: before the episode starts, the learner  
 1194 samples an independent Bernoulli random variable and then follows either  $\pi_1$  or  $\pi_0$  for the entire  
 1195 episode. We use a-priori randomization over complete policies, rather than only randomized actions  
 1196 at each history, because the latter does not automatically convexify the full stopping-time law. Hence,  
 1197 under this assumption one can easily show that the cost-reward set  $\mathcal{K}$  is convex.

1198 **Assumption 8 (Strict feasibility).** *There exists  $\pi_{\text{sf}} \in \mathcal{T}$  such that  $\rho(\pi_{\text{sf}}) > 1 - \delta$ .*

1199 This is the same assumption as in Assumption 5.

1200 The next assumption is not needed for zero duality gap. It is only needed when we want to extract an  
 1201 actual primal-dual saddle point from the zero-gap identity. Unlike global closedness of all bounded  
 1202 slices of  $\mathcal{K}$ , it only asks for closedness of one near-optimal dual level set around a dual maximizer.

1203 **Assumption 9 (Local closedness of near-optimal dual level sets).** *Let  $\mathcal{G} := \arg \max_{\lambda \geq 0} g(\lambda)$ . For*  
 1204 *every  $\lambda^* \in \mathcal{G}$ , there exists  $\epsilon_0 > 0$  such that*

$$\mathcal{K}_{\epsilon_0}(\lambda^*) := \{(c, \rho) \in \mathcal{K} : c + \lambda^*(1 - \delta - \rho) \leq g(\lambda^*) + \epsilon_0\}$$

1205 *is closed in  $\mathbb{R}^2$ .*

1206 **B.5.1 Zero Duality by Perturbation**

1207 We now show how Assumption 7 and Assumption 8 can be used to prove zero duality gap via a  
 1208 perturbed value. Define the scalar perturbation value function

$$\varphi(u) := \inf \{c : (c, \rho) \in \mathcal{K}, 1 - \delta - \rho \leq u\}, \quad u \in \mathbb{R}, \quad (20)$$

1209 with the convention  $\inf \emptyset = +\infty$ . The original primal value is

$$P^* = \varphi(0).$$

1210 The variable  $u$  is the allowed violation of the confidence constraint. If  $u > 0$ , the constraint is relaxed;  
 1211 if  $u < 0$ , the constraint is strengthened.

1212 **Lemma 12** (Basic properties of the perturbation value). *Consider Assumption 7 and Assumption 8.*  
 1213 *Then  $\varphi$  is convex, nonincreasing, and finite on an open interval containing 0. In particular,  $\varphi$  is*  
 1214 *continuous at 0, and  $\partial\varphi(0) \neq \emptyset$ .*

1215 *Proof.* We prove the claims one by one.

1216 Step 1: monotonicity. If  $u_1 \leq u_2$ , then

$$\{(c, \rho) \in \mathcal{K} : 1 - \delta - \rho \leq u_1\} \subseteq \{(c, \rho) \in \mathcal{K} : 1 - \delta - \rho \leq u_2\}.$$

1217 Therefore

$$\varphi(u_2) \leq \varphi(u_1).$$

1218 Thus  $\varphi$  is nonincreasing.

1219 Step 2: convexity. Using Assumption 7 one can easily show that  $\mathcal{K}$  is convex. Take  $u_0, u_1 \in \mathbb{R}$ ,  
 1220  $\alpha \in [0, 1]$ , and  $\eta > 0$ . If  $\varphi(u_i) < \infty$ , choose  $(c_i, \rho_i) \in \mathcal{K}$  such that

$$1 - \delta - \rho_i \leq u_i, \quad c_i \leq \varphi(u_i) + \eta, \quad i \in \{0, 1\}.$$

1221 If one of the two values is  $+\infty$ , the convexity inequality is trivial. By convexity of  $\mathcal{K}$ ,

$$(c_\alpha, \rho_\alpha) := \alpha(c_1, \rho_1) + (1 - \alpha)(c_0, \rho_0) \in \mathcal{K}.$$

1222 Moreover,

$$1 - \delta - \rho_\alpha = \alpha(1 - \delta - \rho_1) + (1 - \alpha)(1 - \delta - \rho_0) \leq \alpha u_1 + (1 - \alpha)u_0.$$

1223 Hence  $(c_\alpha, \rho_\alpha)$  is feasible for  $\varphi(\alpha u_1 + (1 - \alpha)u_0)$ , and so

$$\varphi(\alpha u_1 + (1 - \alpha)u_0) \leq c_\alpha \leq \alpha\varphi(u_1) + (1 - \alpha)\varphi(u_0) + \eta.$$

1224 Letting  $\eta \downarrow 0$  gives convexity.

1225 Step 3: finiteness near zero. By strict feasibility, there exists  $(c_{\text{sf}}, \rho_{\text{sf}}) \in \mathcal{K}$  such that

$$\rho_{\text{sf}} > 1 - \delta.$$

1226 Let

$$\eta_{\text{sf}} := \rho_{\text{sf}} - (1 - \delta) > 0.$$

1227 Then

$$1 - \delta - \rho_{\text{sf}} = -\eta_{\text{sf}}.$$

1228 Therefore, for every  $u \geq -\eta_{\text{sf}}$ , the same point  $(c_{\text{sf}}, \rho_{\text{sf}})$  is feasible for  $\varphi(u)$ . Hence

$$\varphi(u) \leq c_{\text{sf}} < \infty, \quad \forall u \geq -\eta_{\text{sf}}.$$

1229 Also, since  $c \geq 0$  for all  $(c, \rho) \in \mathcal{K}$ , we have  $\varphi(u) \geq 0$  whenever  $\varphi(u) < \infty$ . Thus  $\varphi$  is finite on the  
 1230 open interval  $(-\eta_{\text{sf}}, \infty)$ , which contains 0.

1231 Step 4: continuity and existence of a subgradient. A proper convex function that is finite on an open  
 1232 interval is continuous on that interval. Since  $0 \in \text{int}(\text{dom } \varphi)$ , the one-dimensional subdifferential  
 1233  $\partial\varphi(0)$  is nonempty. Equivalently, one may take any slope between the left and right derivatives of  $\varphi$   
 1234 at zero.  $\square$

1235 **Theorem B.3** (Zero duality gap by perturbation). Assume Assumption 7 and Assumption 8. Then  
 1236 there exists  $\lambda^* \geq 0$  such that

$$g(\lambda^*) = P^*.$$

1237 Consequently,

$$\inf_{\pi \in \mathcal{T}: \rho(\pi) \geq 1 - \delta} c(\pi) = \sup_{\lambda \geq 0} \inf_{\pi \in \mathcal{T}} \{c(\pi) + \lambda(1 - \delta - \rho(\pi))\}. \quad (21)$$

1238 In particular, the Lagrangian dual has no duality gap, and the dual supremum is attained.

1239 *Proof.* By Lemma 12, choose  $s^* \in \partial\varphi(0)$ . Since  $\varphi$  is nonincreasing, every subgradient at zero is  
 1240 nonpositive. Indeed, for  $h > 0$ , the subgradient inequality gives

$$\varphi(h) \geq \varphi(0) + s^*h,$$

1241 while monotonicity gives  $\varphi(h) \leq \varphi(0)$ . Therefore  $s^*h \leq 0$ , and hence  $s^* \leq 0$ .

1242 Define  $\lambda^* := -s^* \geq 0$ . The subgradient inequality gives, for every  $u \in \mathbb{R}$ ,

$$\varphi(u) \geq \varphi(0) + s^*u = P^* - \lambda^*u.$$

1243 Now fix any  $(c, \rho) \in \mathcal{K}$  and set  $u = 1 - \delta - \rho$ . Since  $(c, \rho)$  is feasible for  $\varphi(u)$ , we have  $c \geq \varphi(u)$ .  
 1244 Thus

$$c \geq \varphi(u) \geq P^* - \lambda^*u = P^* - \lambda^*(1 - \delta - \rho).$$

1245 Equivalently,  $c + \lambda^*(1 - \delta - \rho) \geq P^*$ . Taking the infimum over all  $(c, \rho) \in \mathcal{K}$  gives

$$g(\lambda^*) = \inf_{(c, \rho) \in \mathcal{K}} \{c + \lambda^*(1 - \delta - \rho)\} \geq P^*.$$

1246 By weak duality, we also have  $g(\lambda^*) \leq P^*$ . Hence  $g(\lambda^*) = P^*$ . Taking the supremum over  $\lambda \geq 0$   
 1247 proves Eq. (21).  $\square$

1248 **Remark 3.** Strict feasibility (Assumption 8) is a simple sufficient condition. The exact perturbation  
 1249 condition is lower semicontinuity of  $\varphi$  at the origin by the Fenchel-Moreau theorem. In fact,  $\varphi^{**}$  is the  
 1250 lower semicontinuous convex envelope of  $\varphi$ . Since  $\varphi$  is already convex, equality  $\varphi^{**}(0) = \varphi(0)$  holds  
 1251 exactly when  $\varphi$  is lower semicontinuous at 0. Finally, Lemma 12 shows that strict feasibility implies  
 1252 continuity, hence lower semicontinuity, at 0. We state Theorem B.3 under strict feasibility because it  
 1253 is easier to interpret and check. Strict feasibility places 0 in the interior of  $\text{dom } \varphi$ , and finite convex  
 1254 functions are continuous on the interior of their domain. The price is that strict feasibility is stronger  
 1255 than necessary.

### 1256 B.5.2 Primal attainment and KKT

1257 The zero-gap result above does not require any closedness assumption on  $\mathcal{K}$ . However, zero duality  
 1258 gap is only a statement about values. To obtain an actual policy that is both primal feasible and  
 1259 dual optimal, we need an attainment condition. The local closedness assumption above is enough:  
 1260 it compactifies a near-optimal dual level set around a dual maximizer, and a feasible minimizing  
 1261 sequence can be extracted inside this compact set.

1262 **Lemma 13** (Primal-dual attainment from local closedness). Assume 7, 8, and 9. Let  $\lambda^* \in \mathcal{G}$  be any  
 1263 dual maximizer. Then there exists  $(c^*, \rho^*) \in \mathcal{K}$  such that

$$\rho^* \geq 1 - \delta, \quad c^* = P^*, \quad c^* + \lambda^*(1 - \delta - \rho^*) = g(\lambda^*).$$

1264 Equivalently, there exists a policy  $\pi^* \in \mathcal{T}$  such that

$$c(\pi^*) = P^*, \quad \rho(\pi^*) \geq 1 - \delta, \quad \pi^* \in \mathcal{S}(\lambda^*).$$

1265 Moreover,  $\lambda^*(1 - \delta - \rho(\pi^*)) = 0$ .

1266 *Proof.* By Theorem B.3,  $g(\lambda^*) = P^*$ . Let  $(c_n, \rho_n) \in \mathcal{K}$  be a feasible minimizing sequence for the  
 1267 primal problem, so that

$$\rho_n \geq 1 - \delta, \quad c_n \rightarrow P^*.$$

1268 Define  $L_{\lambda^*}(c, \rho) := c + \lambda^*(1 - \delta - \rho)$ . Since  $\rho_n \geq 1 - \delta$  and  $\lambda^* \geq 0$ , we have

$$L_{\lambda^*}(c_n, \rho_n) \leq c_n.$$

1269 On the other hand, by definition of  $g(\lambda^*)$ ,

$$L_{\lambda^*}(c_n, \rho_n) \geq g(\lambda^*) = P^*.$$

1270 Therefore,

$$P^* \leq L_{\lambda^*}(c_n, \rho_n) \leq c_n \rightarrow P^*,$$

1271 and hence  $L_{\lambda^*}(c_n, \rho_n) \rightarrow P^* = g(\lambda^*)$ .

1272 Let  $\epsilon_0 > 0$  be given by Assumption 9. For all sufficiently large  $n$ ,

$$L_{\lambda^*}(c_n, \rho_n) \leq g(\lambda^*) + \epsilon_0,$$

1273 so  $(c_n, \rho_n) \in \mathcal{K}_{\epsilon_0}(\lambda^*)$ .

1274 We now show that  $\mathcal{K}_{\epsilon_0}(\lambda^*)$  is compact. It is closed by assumption. It is also bounded: for every  
1275  $(c, \rho) \in \mathcal{K}_{\epsilon_0}(\lambda^*)$ , since  $\rho \in [0, 1]$ , we have  $1 - \delta - \rho \geq -\delta$ , and thus

$$c \leq g(\lambda^*) + \epsilon_0 + \lambda^* \delta.$$

1276 Also  $c \geq 0$  and  $\rho \in [0, 1]$ . Hence  $\mathcal{K}_{\epsilon_0}(\lambda^*)$  is closed and bounded in  $\mathbb{R}^2$ , and therefore compact.

1277 Passing to a subsequence, we may assume  $(c_n, \rho_n) \rightarrow (c^*, \rho^*)$  for some  $(c^*, \rho^*) \in \mathcal{K}_{\epsilon_0}(\lambda^*) \subseteq \mathcal{K}$ .  
1278 Since  $c_n \rightarrow P^*$ , we have  $c^* = P^*$ . Since  $\rho_n \geq 1 - \delta$  for all  $n$ , we have  $\rho^* \geq 1 - \delta$ . By continuity  
1279 of  $L_{\lambda^*}$ ,

$$L_{\lambda^*}(c^*, \rho^*) = \lim_n L_{\lambda^*}(c_n, \rho_n) = P^* = g(\lambda^*).$$

1280 Because  $(c^*, \rho^*) \in \mathcal{K}$ , there exists  $\pi^* \in \mathcal{T}$  such that  $c(\pi^*) = c^*$  and  $\rho(\pi^*) = \rho^*$ . The equality  
1281  $L_{\lambda^*}(c^*, \rho^*) = g(\lambda^*)$  implies  $\pi^* \in \mathcal{S}(\lambda^*)$ . Finally, since  $c^* = P^*$  and  $L_{\lambda^*}(c^*, \rho^*) = P^*$ , we obtain

$$\lambda^*(1 - \delta - \rho^*) = 0.$$

1282 This proves the claim.  $\square$

1283 Define the dual-optimal policy set

$$\mathcal{S}(\lambda) := \arg \min_{\pi \in \mathcal{T}} \{c(\pi) + \lambda(1 - \delta - \rho(\pi))\}.$$

1284 **Theorem B.4** (KKT and Bayesian correctness). Assume 7, 8, and 9. Let  $\lambda^* \in \mathcal{G}$  be any dual  
1285 maximizer. Then there exists  $\pi^* \in \mathcal{S}(\lambda^*)$  such that

$$c(\pi^*) = P^*, \quad \rho(\pi^*) \geq 1 - \delta, \quad \lambda^*(1 - \delta - \rho(\pi^*)) = 0.$$

1286 Consequently, if  $I_t^*(h) \in \arg \max_{x \in \mathcal{X}} q_t(h, x)$  is a measurable posterior-optimal inference selector,  
1287 then

$$\mathbb{P}_{\theta \sim \nu}^{\pi^*} (L_{\theta} (I_{\tau_{\pi^*}}^* (H_{\tau_{\pi^*}})) \leq \epsilon) \geq 1 - \delta.$$

1288 *Proof.* By Lemma 13, there exists  $\pi^* \in \mathcal{T}$  such that

$$c(\pi^*) = P^*, \quad \rho(\pi^*) \geq 1 - \delta, \quad \pi^* \in \mathcal{S}(\lambda^*),$$

1289 and

$$\lambda^*(1 - \delta - \rho(\pi^*)) = 0.$$

1290 It remains only to translate  $\rho(\pi^*) \geq 1 - \delta$  into the desired correctness statement.

1291 By definition,

$$\rho(\pi^*) = \mathbb{E}^{\pi^*} [r_{\tau_{\pi^*}} (H_{\tau_{\pi^*}})].$$

1292 Since  $I_t^*$  is posterior-optimal,  $r_t(h) = q_t(h, I_t^*(h))$ . Hence

$$\rho(\pi^*) = \mathbb{E}^{\pi^*} [q_{\tau_{\pi^*}} (H_{\tau_{\pi^*}}, I_{\tau_{\pi^*}}^* (H_{\tau_{\pi^*}}))].$$

1293 By Lemma 5, the last display equals

$$\mathbb{P}_{\theta \sim \nu}^{\pi^*} (L_{\theta} (I_{\tau_{\pi^*}}^* (H_{\tau_{\pi^*}})) \leq \epsilon).$$

1294 Since  $\rho(\pi^*) \geq 1 - \delta$ , the result follows.  $\square$

1295 **Remark 4** (What is stronger than the previous correctness statement). *The result above separates*  
1296 *three issues that were previously entangled: (i) zero duality gap, (ii) attainment of a primal-dual*  
1297 *saddle point, (iii) Bayesian-  $(\epsilon, \delta)$ -correctness. Time-sharing plus strict feasibility gives zero duality*  
1298 *gap and dual attainment. The local closedness assumption is only used to extract an actual feasible*  
1299 *dual-optimal policy from the zero-gap identity. Finally, Bayesian correctness follows from the*  
1300 *posterior-optimal inference identity.*

1301 **Weaknesses and limitations.** There are several limitations to keep in mind.

- 1302 1. First, the theorem is a statement about the randomized, time-sharing closure of the policy  
1303 class. If one restricts to deterministic policies without ex-ante randomization,  $\mathcal{K}$  need not be  
1304 convex and a Lagrangian duality gap can occur.
- 1305 2. Second, zero duality gap does not imply that every dual-optimal policy is feasible. The  
1306 theorem guarantees the existence of a feasible dual-optimal policy under the local closedness  
1307 assumption. A learned policy that approximately minimizes the Lagrangian at  $\lambda^*$  is not  
1308 automatically certified by this theorem.
- 1309 3. Third, local closedness is still an attainment assumption. It is weaker than requiring all  
1310 bounded slices of  $\mathcal{K}$  to be closed, but it is not automatic. Without some local closedness, the  
1311 primal value may be approached by policies whose cost-reward pairs converge to a boundary  
1312 point outside  $\mathcal{K}$ , so no exact optimal policy need exist.

### 1313 **B.6 Training-time certification of $(\epsilon, \delta)$ -correctness**

1314 We now give a training-time correctness guarantee for ICPE that treats the learned policy, stopping  
1315 rule, and recommender as a black box. The guarantee is designed for the practical workflow used in  
1316 our experiments: every few training epochs we freeze the current model, evaluate it on a finite batch  
1317 of fresh task realizations, and certify whether the checkpoint satisfies the  $(\epsilon, \delta)$ -guarantees.

1318 The key point is that certification must account for repeated testing. A fixed-level certificate applied  
1319 repeatedly at many checkpoints is not valid without correction. We therefore assign a certification  
1320 budget  $\alpha_m$  to checkpoint  $m$ , with total budget at most  $\alpha \in (0, 1)$ , and certify each checkpoint using  
1321 a checkpointwise mixture martingale.

1322 Fix a target success threshold

$$q^* = 1 - \delta,$$

1323 and an additional parameter  $\alpha \in (0, 1)$  that controls the probability of ever falsely certifying a  
1324 checkpoint.

1325 **Protocol.** We index certification checkpoints by  $m \geq 1$ . In the experiments, for example, check-  
1326 point  $m$  may correspond to every 2500 training epochs. At checkpoint  $m$ , let

$$w_m = (\phi_m, \psi_m, \vartheta_m)$$

1327 denote the frozen parameters of the inference model  $I_\phi$ , the critic  $Q_\psi$  and the actor  $\pi_\vartheta$ . We denote by  
1328  $I_m$ ,  $\pi_m$ , and  $\tau_m$ , respectively, the induced inference rule, sampling policy, and stopping time. The  
1329 probability of success of checkpoint  $m$  is

$$p_m := \mathbb{P}_{\theta \sim \nu}^{\pi_m} (I_m(H_{\tau_m}) \in \mathcal{X}_\epsilon(\theta)).$$

1330 Thus  $p_m$  is the probability of success of the actual recommender rule used by the frozen checkpoint.

1331 For each checkpoint  $m$ , after  $w_m$  is fixed, we draw i.i.d. a batch of  $B$  evaluation episodes. For  
1332  $j = 1, \dots, B$ , let  $\theta^{(m,j)} \sim \nu$  be the task in the  $j$ -th evaluation episode for checkpoint  $m$ , and let  
1333  $H_{\tau_m}^{(m,j)}$  be the stopped history obtained by running the frozen snapshot  $w_m$  on that task. Define the  
1334 success indicator

$$Z_{m,j} := \mathbf{1} \left\{ I_m \left( H_{\tau_m}^{(m,j)} \right) \in \mathcal{X}_\epsilon(\theta^{(m,j)}) \right\}.$$

1335 We write

$$S_{m,t} := \sum_{j=1}^t Z_{m,j}, \quad t = 0, 1, \dots, B.$$

1336 Hence, For each checkpoint  $m$ , conditional on the training history before evaluating checkpoint  $m$   
1337 and conditional on the frozen checkpoint  $w_m$ , the variables

$$Z_{m,1}, \dots, Z_{m,B}$$

1338 are independent and identically distributed as  $\text{Ber}(p_m)$ . Moreover, the certification batch at check-  
1339 point  $m$  is not reused for training updates before the certification decision for checkpoint  $m$  is  
1340 made.

1341 **Checkpointwise mixture martingale.** At checkpoint  $m$ , we test the null hypothesis

$$H_{0,m} : p_m \leq q^*.$$

1342 Let  $\Pi_m$  be a probability distribution on  $[q^*, 1]$ , chosen before observing the certification batch at  
1343 checkpoint  $m$ . For  $r \in [q^*, 1]$ , define

$$L_{m,t}(r) := \left(\frac{r}{q^*}\right)^{S_{m,t}} \left(\frac{1-r}{1-q^*}\right)^{t-S_{m,t}}, \quad t = 0, 1, \dots, B,$$

1344 with the usual convention  $0^0 = 1$ . The checkpointwise mixture martingale is

$$M_{m,t} := \int_{q^*}^1 L_{m,t}(r) \Pi_m(dr), \quad M_{m,0} = 1.$$

1345 In implementation,  $\Pi_m$  may be a finite grid distribution, in which case no numerical integration is  
1346 needed. If

$$\Pi_m = \sum_{\ell=1}^L \omega_{m,\ell} \delta_{r_{m,\ell}}, \quad \sum_{\ell=1}^L \omega_{m,\ell} = 1,$$

1347 then

$$M_{m,t} = \sum_{\ell=1}^L \omega_{m,\ell} \left(\frac{r_{m,\ell}}{q^*}\right)^{S_{m,t}} \left(\frac{1-r_{m,\ell}}{1-q^*}\right)^{t-S_{m,t}}.$$

1348 **Certification budgets.** Let  $(\alpha_m)_{m \geq 1}$  be nonnegative certification budgets, chosen predictably  
1349 before observing the certification batch at checkpoint  $m$ , and satisfying

$$\sum_{m \geq 1} \alpha_m \leq \alpha.$$

1350 We declare checkpoint  $m$  certified if

$$\max_{1 \leq t \leq B} M_{m,t} \geq \alpha_m^{-1}.$$

1351 Let

$$\mathcal{C} := \left\{ m \geq 1 : \max_{1 \leq t \leq B} M_{m,t} \geq \alpha_m^{-1} \right\}$$

1352 be the set of certified checkpoints.

1353 **Proposition 5** (Training-time certification by checkpointwise mixture martingales). *Suppose that the*  
1354 *certification budgets satisfy  $\sum_{m \geq 1} \alpha_m \leq \alpha$  almost surely. Then*

$$\mathbb{P}(\forall m \in \mathcal{C}, p_m > q^*) \geq 1 - \alpha.$$

1355 *Equivalently, every certified checkpoint is  $(\epsilon, \delta)$ -correct, and, after training, any adaptively selected*  
1356 *checkpoint*

$$\hat{m} \in \mathcal{C}$$

1357 *is  $(\epsilon, \delta)$ -correct with confidence at least  $1 - \alpha$ .*

1358 *Proof.* Fix a checkpoint  $m$ . We condition on the training history before the certification batch for  
1359 checkpoint  $m$ , and on the frozen checkpoint  $w_m$ . Under this conditioning,  $p_m$ ,  $\Pi_m$ , and  $\alpha_m$  are fixed,  
1360 and

$$Z_{m,1}, \dots, Z_{m,B} \stackrel{\text{i.i.d.}}{\sim} \text{Bernoulli}(p_m).$$

1361 Assume the null hypothesis

$$H_{0,m} : p_m \leq q^*$$

1362 holds. Fix any  $r \in [q^*, 1]$ . We show that  $\{L_{m,t}(r)\}_{t=0}^B$  is a nonnegative supermartingale under  $H_{0,m}$ .

1363 Let  $\mathcal{F}_{m,t}$  be the sigma-field generated by the training history, the frozen checkpoint  $w_m$ , and the first  
1364  $t$  certification outcomes

$$Z_{m,1}, \dots, Z_{m,t}.$$

1365 For  $t < B$ ,

$$\frac{L_{m,t+1}(r)}{L_{m,t}(r)} = \begin{cases} r/q^*, & Z_{m,t+1} = 1, \\ (1-r)/(1-q^*), & Z_{m,t+1} = 0. \end{cases}$$

1366 Therefore,

$$\mathbb{E}\left[\frac{L_{m,t+1}(r)}{L_{m,t}(r)} \mid \mathcal{F}_{m,t}\right] = p_m \frac{r}{q^*} + (1-p_m) \frac{1-r}{1-q^*}.$$

1367 A direct calculation gives

$$p_m \frac{r}{q^*} + (1-p_m) \frac{1-r}{1-q^*} = 1 + \frac{(p_m - q^*)(r - q^*)}{q^*(1-q^*)}.$$

1368 Since  $p_m \leq q^*$  and  $r \geq q^*$ , the final term is nonpositive. Hence

$$\mathbb{E}\left[\frac{L_{m,t+1}(r)}{L_{m,t}(r)} \mid \mathcal{F}_{m,t}\right] \leq 1.$$

1369 Thus

$$\mathbb{E}[L_{m,t+1}(r) \mid \mathcal{F}_{m,t}] \leq L_{m,t}(r),$$

1370 so  $\{L_{m,t}(r)\}_{t=0}^B$  is a nonnegative supermartingale.

1371 Because  $M_{m,t}$  is a mixture of the nonnegative supermartingales  $L_{m,t}(r)$ , Tonelli's theorem gives

$$\mathbb{E}[M_{m,t+1} \mid \mathcal{F}_{m,t}] = \int_{q^*}^1 \mathbb{E}[L_{m,t+1}(r) \mid \mathcal{F}_{m,t}] \Pi_m(dr) \leq \int_{q^*}^1 L_{m,t}(r) \Pi_m(dr) = M_{m,t}.$$

1372 Therefore  $\{M_{m,t}\}_{t=0}^B$  is a nonnegative supermartingale with  $M_{m,0} = 1$ . By Ville's inequality,

$$\mathbb{P}\left(\max_{1 \leq t \leq B} M_{m,t} \geq \alpha_m^{-1} \mid \text{past}, w_m, H_{0,m}\right) \leq \alpha_m.$$

1373 Let

$$A_m := \left\{ p_m \leq q^* \text{ and } \max_{1 \leq t \leq B} M_{m,t} \geq \alpha_m^{-1} \right\}$$

1374 be the event that checkpoint  $m$  is falsely certified. From the conditional bound above,

$$\mathbb{P}(A_m \mid \text{past}) \leq \alpha_m.$$

1375 Taking expectations,

$$\mathbb{P}(A_m) \leq \mathbb{E}[\alpha_m].$$

1376 By the union bound,

$$\mathbb{P}(\exists m \geq 1 : A_m) \leq \sum_{m \geq 1} \mathbb{P}(A_m) \leq \mathbb{E}\left[\sum_{m \geq 1} \alpha_m\right] \leq \alpha.$$

1377 Thus, with probability at least  $1 - \alpha$ , no checkpoint with  $p_m \leq q^*$  is certified. Equivalently,

$$p_m > q^* \quad \forall m \in \mathcal{C}.$$

1378 Taking  $q^* = 1 - \delta$  gives the stated  $(\epsilon, \delta)$ -correctness guarantee.  $\square$

1379 **Remark 5** (What the guarantee certifies). *The guarantee is simultaneous over all certified check-*  
 1380 *points:*

$$\mathbb{P}(\forall m \in \mathcal{C}, p_m > 1 - \delta) \geq 1 - \alpha.$$

1381 *Thus the final returned model need not be the first certified model. It may be any checkpoint selected*  
 1382 *after training, as long as it belongs to  $\mathcal{C}$ . The result does not certify checkpoints that were never*  
 1383 *certified, nor does it imply that all checkpoints after the first certified checkpoint are correct.*

1384 **Remark 6** (Relation to the original ICPE training-time argument). *The original finite-hypothesis*  
 1385 *ICPE argument pools evidence across checkpoints and tests a global null of the form*

$$\sup_m p_m \leq q^*.$$

1386 *That pooled martingale can show that training has produced evidence against the hypothesis that all*  
 1387 *checkpoints are bad. However, by itself it does not certify an arbitrary later checkpoint unless one*  
 1388 *adds a persistence or monotonicity assumption on the sequence  $(p_m)$ , as done in [53].*

1389 *The checkpointwise martingale above is different. It tests each frozen checkpoint separately, spends*  
 1390 *error probability across checkpoints, and therefore certifies the checkpoints that actually pass the*  
 1391 *test. This matches the practical workflow in which we may check a model, fail to certify it, continue*  
 1392 *training, and later certify a different checkpoint. No monotonicity or persistence assumption on the*  
 1393 *training trajectory is required.*

1394 **Remark 7** (Choice of certification budgets). *If all checkpoints are treated symmetrically and a*  
 1395 *maximum number  $M_{\max}$  of checks is fixed in advance, the uniform allocation*

$$\alpha_m = \frac{\alpha}{M_{\max}}$$

1396 *is the simplest choice. If later checkpoints are expected to be better, one can use the exponentially*  
 1397 *back-loaded allocation*

$$\alpha_m = \alpha \frac{(\gamma - 1)\gamma^{m-1}}{\gamma^{M_{\max}} - 1}, \quad \gamma > 1.$$

1398 *This preserves the same correctness theorem because the proof only uses*

$$\sum_m \alpha_m \leq \alpha.$$

1399 *The choice of schedule affects power, not validity. Larger  $\alpha_m$  makes checkpoint  $m$  easier to certify,*  
 1400 *so back-loading the budget gives more power to later checkpoints at the expense of earlier ones.*

1401 **Remark 8** (Finite-grid implementation). *In practice we use a finite grid*

$$q^* \leq r_1 < \dots < r_L \leq 1$$

1402 *with weights  $\omega_1, \dots, \omega_L$ . Then*

$$M_{m,t} = \sum_{\ell=1}^L \omega_{\ell} \left( \frac{r_{\ell}}{q^*} \right)^{S_{m,t}} \left( \frac{1 - r_{\ell}}{1 - q^*} \right)^{t - S_{m,t}},$$

1403 *which is easy to compute and remains a valid mixture martingale. The grid should place mass on*  
 1404 *plausible alternatives above  $q^*$ , for example  $r_{\ell} \in \{0.91, 0.92, 0.93, 0.94, 0.95, 0.97\}$  when  $q^* = 0.9$ .*  
 1405 *Placing mass closer to  $q^*$  improves power for small margins but requires larger batches; placing*  
 1406 *mass farther above  $q^*$  improves power when the checkpoint is substantially better than the target.*

## 1407 B.7 Choice of inference model and reward modeling

1408 This subsection clarifies the relation between the ideal posterior quantities used in the Bellman  
 1409 characterization and the inference model used in the implementation. The ideal inference rule  
 1410 maximizes

$$q_t(h, x) := \mathbb{P}(L_{\theta}(x) \leq \epsilon | H_t = h), \quad r_t(h) := \sup_{x \in \mathcal{X}} q_t(h, x),$$

1411 and any measurable maximizer is denoted by

$$I_t^*(h) \in \arg \max_{x \in \mathcal{X}} q_t(h, x).$$

1412 The implementation does not parameterize  $I_t^*$  directly. Instead, it uses a stochastic selector  $\hat{I}_t(\cdot|h)$ ,  
 1413 represented by a diagonal Gaussian base distribution

$$\hat{I}_t(\cdot|h) = \mathcal{N}(\mu_t(h), \Sigma_t(h)).$$

1414 The recommendation is the mean  $\mu_t(h)$ , while the critic evaluates samples from  $\hat{I}_t(\cdot|h)$ .

1415 **Sampled reward.** For fixed  $(h, \theta)$ , define the success probability of the implemented stochastic  
 1416 selector by

$$\hat{r}_t(h, \theta) := \mathbb{E}_{x \sim \hat{I}_t(\cdot|h)}[\mathbf{1}\{x \in \mathcal{X}_\epsilon(\theta)\}].$$

1417 Its posterior average is

$$\hat{r}_t(h) := \mathbb{E}[\hat{r}_t(h, \theta) | H_t = h].$$

1418 Given  $M$  Monte Carlo samples  $X_t^{(1)}, \dots, X_t^{(M)} \sim \hat{I}_t(\cdot|h)$ , we use

$$\hat{r}_{t,M}(h, \theta) := \frac{1}{M} \sum_{m=1}^M \mathbf{1}\{X_t^{(m)} \in \mathcal{X}_\epsilon(\theta)\}.$$

1419 **Lemma 14** (Sampled stochastic-selector reward). *Conditionally on  $(H_t = h, \theta)$ ,*

$$\mathbb{E}[\hat{r}_{t,M}(h, \theta) | H_t = h, \theta] = \hat{r}_t(h, \theta), \quad \text{Var}(\hat{r}_{t,M}(h, \theta) | H_t = h, \theta) \leq \frac{1}{4M}.$$

1420 *Furthermore,*

$$\hat{r}_t(h) = \mathbb{E}[q_t(h, \hat{X}_t) | H_t = h] \leq r_t(h).$$

1421 *Proof.* The first claim follows because  $\hat{r}_{t,M}$  is the average of  $M$  Bernoulli random variables with  
 1422 success probability  $\hat{r}_t(h, \theta)$ . The variance bound follows from  $p(1-p) \leq 1/4$ . For the posterior  
 1423 identity,  $\hat{X}_t$  is conditional independent of  $\theta$  given  $H_t = h$ , thus

$$\hat{r}_t(h) = \int_{\mathcal{X}} \mathbb{P}(x \in \mathcal{X}_\epsilon(\theta) | H_t = h) \hat{I}_t(dx|h) = \int_{\mathcal{X}} q_t(h, x) \hat{I}_t(dx|h).$$

1424 The last display is bounded by  $\sup_{x \in \mathcal{X}} q_t(h, x) = r_t(h)$ . □

1425 The reward  $\hat{r}_t(h)$  evaluates the stochastic selector we actually use. It is therefore conservative relative  
 1426 to the ideal deterministic reward  $r_t(h)$ , which assumes access to the best posterior recommendation.

1427 **Second-moment robustness.** The next proposition gives two sufficient conditions under which the  
 1428 sampled reward is close to the ideal one. The first bound uses concentration around the task-level  
 1429 zero-loss target  $x_\theta^*$ . The second uses concentration around the posterior-optimal rule  $I_t^*(h)$ .

1430 For  $z \in \mathcal{X}$ , define

$$D_t(h, z) := \mathbb{E}[\|\hat{X}_t - z\|^2 | H_t = h], \quad \hat{X}_t \sim \hat{I}_t(\cdot|h).$$

1431 **Proposition 6** (Second-moment robustness and ideal-reward gap). *Assume that, for posterior-a.e.  $\theta$ ,*  
 1432 *there exists  $\rho(\theta) > 0$  such that*

$$B(x_\theta^*, \rho(\theta)) \cap \mathcal{X} \subseteq \mathcal{X}_\epsilon(\theta).$$

1433 *where  $B(x, r)$  is an euclidean ball of radius  $r$  around  $x$ . Then*

$$\hat{r}_t(h, \theta) \geq 1 - \frac{D_t(h, x_\theta^*)}{\rho(\theta)^2}.$$

1434 *Moreover, assume  $\hat{X}_t$  is conditionally independent of  $\theta$  given  $H_t = h$ . If  $I_t^*(h) \in$   
 1435  $\arg \max_{x \in \mathcal{X}} q_t(h, x)$  and  $q_t(h, \cdot)$  is  $L_t(h)$ -Lipschitz on  $B(I_t^*(h), R_t(h)) \cap \mathcal{X}$ , then*

$$0 \leq r_t(h) - \hat{r}_t(h) \leq L_t(h) \sqrt{D_t(h, I_t^*(h))} + \frac{D_t(h, I_t^*(h))}{R_t(h)^2}.$$

1436 *If  $q_t(h, \cdot)$  is globally  $L_t(h)$ -Lipschitz on  $\mathcal{X}$ , the second term is unnecessary.*

1437 *Proof.* For the first claim, the margin assumption gives

$$\{\|\hat{X}_t - x_\theta^*\| \leq \rho(\theta)\} \subseteq \{\hat{X}_t \in \mathcal{X}_\epsilon(\theta)\}.$$

1438 Therefore

$$1 - \hat{r}_t(h, \theta) \leq \mathbb{P}(\|\hat{X}_t - x_\theta^*\| > \rho(\theta) | H_t = h, \theta).$$

1439 Since the law of  $\hat{X}_t$  depends on  $h$  but not on  $\theta$  conditional on  $h$ , Markov's inequality gives

$$1 - \hat{r}_t(h, \theta) \leq \frac{D_t(h, x_\theta^*)}{\rho(\theta)^2}.$$

1440 For the second claim, by Lemma 14,

$$\hat{r}_t(h) = \mathbb{E}[q_t(h, \hat{X}_t) | H_t = h].$$

1441 Since  $I_t^*(h)$  maximizes  $q_t(h, \cdot)$ ,

$$r_t(h) - \hat{r}_t(h) = \mathbb{E}[q_t(h, I_t^*(h)) - q_t(h, \hat{X}_t) | H_t = h] \geq 0.$$

1442 Let

$$E_h := \{\|\hat{X}_t - I_t^*(h)\| \leq R_t(h)\}.$$

1443 On  $E_h$ , Lipschitzness gives

$$q_t(h, I_t^*(h)) - q_t(h, \hat{X}_t) \leq L_t(h) \|\hat{X}_t - I_t^*(h)\|.$$

1444 On  $E_h^c$ , the same difference is at most 1. Hence

$$r_t(h) - \hat{r}_t(h) \leq L_t(h) \mathbb{E}[\|\hat{X}_t - I_t^*(h)\| | H_t = h] + \mathbb{P}(E_h^c | H_t = h).$$

1445 Jensen's inequality gives

$$\mathbb{E}[\|\hat{X}_t - I_t^*(h)\| | H_t = h] \leq \sqrt{D_t(h, I_t^*(h))},$$

1446 and Markov's inequality gives

$$\mathbb{P}(E_h^c | H_t = h) \leq \frac{D_t(h, I_t^*(h))}{R_t(h)^2}.$$

1447 If  $q_t(h, \cdot)$  is globally Lipschitz, take  $E_h = \mathcal{X}$  and remove the last term.  $\square$

1448 The first bound gives the margin interpretation of the sampled reward: if the stochastic selector has  
 1449 small second moment around  $x_\theta^*$ , then its samples are likely to be  $\epsilon$ -optimal. The second bound  
 1450 is different: it compares the stochastic selector to the ideal posterior reward and is small when the  
 1451 selector concentrates around  $I_t^*(h)$  and  $q_t(h, \cdot)$  is locally regular. Thus the sampled reward evaluates  
 1452 the implemented selector through the same event used for correctness; it does not assume that the  
 1453 NLL mean is exactly optimal.

1454 **Gaussian NLL.** We next characterize the population Gaussian NLL objective. Let

$$Z := x_\theta^*$$

1455 be the selected zero-loss target, viewed as a random variable under the posterior law  $\theta | H_t = h$ . For  
 1456  $\mu \in \mathbb{R}^d$  and  $\Sigma \in \mathbb{S}_{++}^d$ , define

$$\mathcal{L}_h(\mu, \Sigma) := \mathbb{E}[-\log \mathcal{N}(Z; \mu, \Sigma) | H_t = h].$$

1457 Let

$$\mu_t^{\text{NLL}}(h) := \mathbb{E}[Z | H_t = h], \quad \Sigma_t^{\text{NLL}}(h) := \text{Cov}(Z | H_t = h).$$

1458 **Proposition 7** (Gaussian NLL moment projection). *Assume  $Z | H_t = h$  has finite second moment and*  
 1459 *positive definite covariance. Then the unique minimizer of  $\mathcal{L}_h(\mu, \Sigma)$  over  $\mu \in \mathbb{R}^d$  and  $\Sigma \in \mathbb{S}_{++}^d$  is*

$$\mu = \mu_t^{\text{NLL}}(h), \quad \Sigma = \Sigma_t^{\text{NLL}}(h).$$

1460 *If the covariance is restricted to be diagonal, the optimal mean is still  $\mu_t^{\text{NLL}}(h)$  and the optimal*  
 1461 *diagonal entries are the posterior coordinate variances of  $Z$ .*

1462 *Proof.* Up to an additive constant,

$$\mathcal{L}_h(\mu, \Sigma) = \frac{1}{2} \log \det \Sigma + \frac{1}{2} \mathbb{E}[(Z - \mu)^\top \Sigma^{-1} (Z - \mu) | H_t = h].$$

1463 For fixed  $\Sigma$ , the second term is minimized at  $\mu = \mathbb{E}[Z|H_t = h]$ . With this choice, the objective  
 1464 becomes

$$\frac{1}{2} \log \det \Sigma + \frac{1}{2} \text{tr}(\Sigma^{-1} \Sigma_t^{\text{NLL}}(h)).$$

1465 The first-order condition in  $\Sigma$  gives

$$\Sigma = \Sigma_t^{\text{NLL}}(h),$$

1466 and strict convexity in the natural parameters gives uniqueness. The diagonal case follows by the  
 1467 same calculation coordinate-wise.  $\square$

1468 Thus Gaussian NLL performs a moment projection of the posterior law of  $x_\theta^*$ : it matches the first  
 1469 two posterior moments, or the coordinate variances in the diagonal case. This does not imply that  
 1470  $\mu_t^{\text{NLL}}(h)$  maximizes  $q_t(h, \cdot)$  in general. The NLL objective learns the posterior target law, while  
 1471 optimality for stopping is defined by the  $\epsilon$ -success probability.

1472 **Near-optimality of the NLL mean.** The NLL objective does not directly maximize  $q_t(h, x)$ . It  
 1473 learns the posterior mean of the selected target  $x_\theta^*$  under the Gaussian moment projection. The  
 1474 next result gives a simple condition under which this mean is nevertheless close to the ideal rule in  
 1475 posterior success probability. The condition is posterior concentration relative to the margin of the  
 1476 success set, not exact equality between the posterior mean and the maximizer of  $q_t(h, \cdot)$ .

1477 **Proposition 8** (Near-optimality of the NLL mean). *Fix a history  $h$  and write*

$$Z := x_\theta^*, \quad \mu_t^{\text{NLL}}(h) := \mathbb{E}[Z|H_t = h], \quad V_t(h) := \mathbb{E}[\|Z - \mu_t^{\text{NLL}}(h)\|^2 | H_t = h].$$

1478 *Assume that there exists  $\rho > 0$  such that, posterior-a.s.,  $B(x_\theta^*, \rho) \cap \mathcal{X} \subseteq \mathcal{X}_\epsilon(\theta)$ , where  $B(x, r)$  is  
 1479 an Euclidean ball of radius  $r$  around  $x$ . Let  $\hat{\mu}_t(h) \in \mathcal{X}$  be the deployed mean recommendation and  
 1480 assume*

$$\|\hat{\mu}_t(h) - \mu_t^{\text{NLL}}(h)\| \leq e_t(h).$$

1481 *Then*

$$0 \leq r_t(h) - q_t(h, \hat{\mu}_t(h)) \leq \frac{V_t(h) + e_t(h)^2}{\rho^2}.$$

1482 *Moreover, if  $I_t^*(h) \in \arg \max_{x \in \mathcal{X}} q_t(h, x)$ ,  $q_t(h, \cdot)$  is  $L_t(h)$ -Lipschitz on  $B(I_t^*(h), R_t(h)) \cap \mathcal{X}$ , and  
 1483  $\|\hat{\mu}_t(h) - I_t^*(h)\| \leq R_t(h)$ , then*

$$0 \leq r_t(h) - q_t(h, \hat{\mu}_t(h)) \leq L_t(h) \|\hat{\mu}_t(h) - I_t^*(h)\|.$$

1484 *Proof.* The margin assumption implies

$$\{\|\hat{\mu}_t(h) - x_\theta^*\| \leq \rho\} \subseteq \{\hat{\mu}_t(h) \in \mathcal{X}_\epsilon(\theta)\}.$$

1485 Therefore

$$q_t(h, \hat{\mu}_t(h)) \geq 1 - \mathbb{P}(\|\hat{\mu}_t(h) - Z\| > \rho | H_t = h).$$

1486 By Markov's inequality,

$$q_t(h, \hat{\mu}_t(h)) \geq 1 - \frac{\mathbb{E}[\|\hat{\mu}_t(h) - Z\|^2 | H_t = h]}{\rho^2}.$$

1487 Since  $\hat{\mu}_t(h)$  is deterministic conditional on  $H_t = h$ ,

$$\mathbb{E}[\|\hat{\mu}_t(h) - Z\|^2 | H_t = h] = V_t(h) + \|\hat{\mu}_t(h) - \mu_t^{\text{NLL}}(h)\|^2 \leq V_t(h) + e_t(h)^2.$$

1488 The first claim follows because  $r_t(h) \leq 1$ . For the second claim, use  $r_t(h) = q_t(h, I_t^*(h))$  and  
 1489 Lipschitzness:

$$r_t(h) - q_t(h, \hat{\mu}_t(h)) = q_t(h, I_t^*(h)) - q_t(h, \hat{\mu}_t(h)) \leq L_t(h) \|\hat{\mu}_t(h) - I_t^*(h)\|.$$

1490  $\square$

1491 The proposition clarifies what is, and is not, implied by the Gaussian NLL. In general,  $\mu_t^{\text{NLL}}(h) =$   
1492  $\mathbb{E}[x_\theta^* | H_t = h]$  need not maximize  $q_t(h, \cdot)$ : if the posterior law of  $x_\theta^*$  is multimodal, the mean can lie  
1493 between modes. The first bound gives the guarantee tied to NLL training. If the posterior uncertainty  
1494 on  $x_\theta^*$  is small relative to the margin of  $\mathcal{X}_\epsilon(\theta)$ , then the deployed mean is near-optimal in posterior  
1495 success probability. The Lipschitz bound is different: it says that any deterministic recommendation  
1496 close to an ideal maximizer  $I_t^*(h)$  is near-optimal when  $q_t(h, \cdot)$  is locally regular. Thus the Lipschitz  
1497 argument also applies to the NLL mean, but only if one separately controls its distance to  $I_t^*(h)$ .

1498 For localization losses  $L_\theta(x) = \|x - x_\theta^*\|$ , the margin condition holds with  $\rho(\theta) = \epsilon$ . For smooth  
1499 value-gap losses, it follows from a local upper curvature bound near  $x_\theta^*$ . For example, if

$$f_\theta(x) - f_\theta(x_\theta^*) \leq \frac{M}{2} \|x - x_\theta^*\|^2$$

1500 near  $x_\theta^*$ , then

$$B\left(x_\theta^*, \sqrt{\frac{2\epsilon}{M}}\right) \cap \mathcal{X} \subseteq \mathcal{X}_\epsilon(\theta).$$

1501 **Remark 9** (Exact alignment under symmetric localization). *In special cases the NLL mean is exactly*  
1502 *Bayes-optimal. Suppose the success sets are translates of a fixed centrally symmetric convex set, i.e.,*

$$\mathcal{X}_\epsilon(\theta) = \{x \in \mathcal{X} : x - x_\theta^* \in S_\epsilon\},$$

1503 *and suppose  $x_\theta^* | H_t = h$  is Gaussian with mean  $\mu_t^{\text{NLL}}(h)$ , with  $\mu_t^{\text{NLL}}(h) \in \mathcal{X}$ . Then*

$$\mu_t^{\text{NLL}}(h) \in \arg \max_{x \in \mathcal{X}} q_t(h, x).$$

1504 *Indeed,  $q_t(h, x)$  is the posterior probability that a Gaussian random variable falls in a translate of*  
1505  *$S_\epsilon$ , and this probability is maximized when the translate is centered at the Gaussian mean. Without*  
1506 *this type of symmetry, the posterior mean, posterior mode, and maximizer of  $q_t(h, \cdot)$  can differ.*

## 1507 B.8 Robustness to prior misspecification

1508 We now study deployment under a misspecified task prior. The controller C-ICPE is trained under a  
1509 prior  $\nu$  on  $\Theta$  and then frozen. At deployment, the same learned inference network, critic, action rule,  
1510 cost parameter, and stopping rule are used, but the environment parameter is drawn from a different  
1511 prior  $\nu'$ . Thus, conditional on a fixed environment  $\theta$ , the trajectory law  $\mathbb{P}_\theta^{\pi_\phi}$  is unchanged; only the  
1512 outer averaging measure over  $\theta$  changes.

1513 Let

$$s(\theta) := \mathbb{P}_\theta^{\pi_\phi}(\mu_\phi(H_{\hat{\tau}}) \in \mathcal{X}_\epsilon(\theta)) \in [0, 1], \quad t(\theta) := \mathbb{E}_\theta^{\pi_\phi}[\hat{\tau}] \in [0, \infty].$$

1514 For a prior  $\eta$  on  $\Theta$ , define

$$p^{(\eta)} := \int_{\Theta} s(\theta) \eta(d\theta), \quad T^{(\eta)} := \int_{\Theta} t(\theta) \eta(d\theta).$$

1515 The training-prior guarantee is

$$p^{(\nu)} \geq 1 - \delta.$$

1516 The goal is to understand how  $p^{(\nu')}$  and  $T^{(\nu')}$  change when  $\nu$  is replaced by  $\nu'$ .

1517 **Assumption 10** (Measurable performance profiles). *The maps  $s : \Theta \rightarrow [0, 1]$  and  $t : \Theta \rightarrow [0, \infty]$*   
1518 *are Borel measurable. Moreover,  $t \in L^1(\eta)$  for every prior  $\eta$  considered below.*

1519 **Lemma 15** (Prior shift as reweighting). *Under Assumption 10, correctness and stopping time under*  
1520 *any prior  $\eta$  are given by*

$$p^{(\eta)} = \mathbb{E}_{\theta \sim \eta}[s(\theta)], \quad T^{(\eta)} = \mathbb{E}_{\theta \sim \eta}[t(\theta)].$$

1521 *If  $\eta \ll \nu$  with density ratio  $w_\eta = d\eta/d\nu$ , then*

$$p^{(\eta)} = \mathbb{E}_\nu[w_\eta s], \quad T^{(\eta)} = \mathbb{E}_\nu[w_\eta t].$$

1522 *Proof.* Under prior  $\eta$ , the joint law factors as

$$\eta(d\theta) \mathbb{P}_\theta^{\pi_\phi}(dh).$$

1523 Applying Tonelli's theorem to the success indicator gives

$$p^{(\eta)} = \int_{\Theta} \mathbb{P}_\theta^{\pi_\phi}(\mu_\phi(H_{\hat{\tau}}) \in \mathcal{X}_\varepsilon(\theta)) \eta(d\theta) = \int_{\Theta} s(\theta) \eta(d\theta).$$

1524 The identity for  $T^{(\eta)}$  follows similarly from Tonelli applied to the nonnegative random variable  $\hat{\tau}$ .  
 1525 The reweighting identities follow by the change of measure  $d\eta = w_\eta d\nu$ .  $\square$

1526 The next proposition gives several complementary ways of transferring the training-prior guarantee  
 1527  $p^{(\nu)} \geq 1 - \delta$  to a deployment prior  $\nu'$ . The density-ratio and  $\chi^2$  bounds are useful when  $\nu' \ll \nu$ ; the  
 1528 total-variation and good-set bounds do not require absolute continuity.

1529 **Proposition 9** (Success probability under prior shift). *Assume Assumption 10, and write  $p = p^{(\nu)}$ ,*  
 1530  *$p' = p^{(\nu')}$ . Then:*

1531 (a) *If  $\nu' \ll \nu$  and  $C = \|d\nu'/d\nu\|_\infty$ , then  $1 - p' \leq C(1 - p)$ . In particular,  $p \geq 1 - \delta$  implies*  
 1532  *$p' \geq 1 - C\delta$ .*

1533 (b) *If  $\nu' \ll \nu$  and  $\chi^2(\nu'|\nu) < \infty$ , then  $|p' - p| \leq \sqrt{\chi^2(\nu'|\nu) \text{Var}_\nu(s)}$ . Hence  $p \geq 1 - \delta$*   
 1534 *implies*

$$p' \geq 1 - \delta - \sqrt{\chi^2(\nu'|\nu)\delta}.$$

1535 (c) *For arbitrary priors,  $|p' - p| \leq \text{TV}(\nu', \nu)$ . Consequently,*

$$p' \geq 1 - \delta - \text{TV}(\nu', \nu).$$

1536 *By Pinsker's inequality, this also gives*

$$p' \geq 1 - \delta - \sqrt{\frac{1}{2} D_{\text{KL}}(\nu'|\nu)}$$

1537 *whenever  $D_{\text{KL}}(\nu'|\nu) < \infty$ , and the analogous bound with the KL arguments reversed.*

1538 (d) *Let  $G \subseteq \Theta$  be measurable. If  $s(\theta) \geq 1 - \eta$  on  $G$ , then*

$$p' \geq (1 - \eta)\nu'(G).$$

1539 *Moreover, if  $p \geq 1 - \delta$  and*

$$G_\eta := \{\theta : s(\theta) \geq 1 - \eta\},$$

1540 *then*

$$\nu(G_\eta) \geq 1 - \frac{\delta}{\eta}, \quad p' \geq (1 - \eta)\nu'(G_\eta).$$

1541 *Proof.* Let  $e(\theta) = 1 - s(\theta) \in [0, 1]$ .

1542 For (a), if  $w = d\nu'/d\nu$ , then

$$1 - p' = \mathbb{E}_{\nu'}[e] = \mathbb{E}_\nu[we] \leq \|w\|_\infty \mathbb{E}_\nu[e] = C(1 - p).$$

1543 For (b), since  $\mathbb{E}_\nu[w] = 1$ ,

$$p' - p = \mathbb{E}_\nu[(w - 1)s] = \mathbb{E}_\nu[(w - 1)(s - p)].$$

1544 Cauchy–Schwarz yields

$$|p' - p| \leq \sqrt{\mathbb{E}_\nu[(w - 1)^2]} \sqrt{\text{Var}_\nu(s)} = \sqrt{\chi^2(\nu'|\nu) \text{Var}_\nu(s)}.$$

1545 Since  $0 \leq s \leq 1$ ,  $s^2 \leq s$ , hence

$$\text{Var}_\nu(s) \leq p(1 - p) \leq 1 - p.$$

1546 If  $p \geq 1 - \delta$ , then  $\text{Var}_\nu(s) \leq \delta$ , giving the displayed bound.

1547 For (c), since  $s \in [0, 1]$ ,

$$|p' - p| = \left| \int s d(\nu' - \nu) \right| \leq \text{TV}(\nu', \nu).$$

1548 Pinsker's inequality gives the KL consequences.

1549 For (d),

$$p' = \int_G s d\nu' + \int_{G^c} s d\nu' \geq (1 - \eta)\nu'(G).$$

1550 If  $G = G_\eta$ , then  $G_\eta^c = \{e > \eta\}$ . Since  $\mathbb{E}_\nu[e] \leq \delta$ , Markov's inequality gives

$$\nu(G_\eta^c) \leq \frac{\delta}{\eta}.$$

1551

□

1552 The good-set formulation is often the most faithful explanation of empirical robustness. Average  
1553 correctness under  $\nu$  implies that the controller is accurate on a large  $\nu$ -measure set of environments.  
1554 Deployment remains accurate whenever  $\nu'$  continues to place most of its mass on that same set.

1555 We now state the analogous bounds for the stopping time. Since  $t$  is not bounded by one, the bounds  
1556 require either density-ratio control, a second-moment assumption, or a bounded horizon.

1557 **Proposition 10** (Stopping time under prior shift). *Assume Assumption 10, and write  $T = T^{(\nu)}$ ,*  
1558  *$T' = T^{(\nu')}$ . Then:*

1559 (a) *If  $\nu' \ll \nu$  and  $C = \|d\nu'/d\nu\|_\infty$ , then  $T' \leq CT$ .*

1560 (b) *If  $\nu' \ll \nu$ ,  $\chi^2(\nu'|\nu) < \infty$ , and  $t \in L^2(\nu)$ , then*

$$|T' - T| \leq \sqrt{\chi^2(\nu'|\nu) \text{Var}_\nu(t)}.$$

1561 (c) *If  $\hat{\tau} \leq T_{\max}$  almost surely under every  $\theta$ , then*

$$|T' - T| \leq T_{\max} \text{TV}(\nu', \nu).$$

1562 (d) *If  $\hat{\tau} \leq T_{\max}$  almost surely and  $t(\theta) \leq \tau_0$  on a measurable set  $F \subseteq \Theta$ , then*

$$T' \leq \tau_0 + (T_{\max} - \tau_0)\nu'(F^c).$$

1563 *Proof.* The first claim follows from

$$T' = \mathbb{E}_\nu[wt] \leq \|w\|_\infty \mathbb{E}_\nu[t] = CT.$$

1564 For the second, use

$$T' - T = \mathbb{E}_\nu[(w - 1)t] = \mathbb{E}_\nu[(w - 1)(t - T)]$$

1565 and apply Cauchy–Schwarz. For the third,  $t/T_{\max} \in [0, 1]$ , so the total-variation bound applies. For  
1566 the fourth, split

$$T' = \int_F t d\nu' + \int_{F^c} t d\nu'$$

1567 and use  $t \leq \tau_0$  on  $F$  and  $t \leq T_{\max}$  everywhere. □

1568 **Remark 10** (Beta–Uniform shifts). *For one-dimensional shifts with  $\nu = \text{Unif}(0, 1)$  and  $\nu' =$*   
1569 *Beta( $\alpha, \beta$ ), the constants in Proposition 9 can be evaluated in closed form. For example,*

$$\chi^2(\nu'|\nu) = \frac{B(2\alpha - 1, 2\beta - 1)}{B(\alpha, \beta)^2} - 1$$

1570 *when  $\alpha, \beta > 1/2$ , and is infinite otherwise. The essential supremum  $\|d\nu'/d\nu\|_\infty$  is finite exactly*  
1571 *when  $\alpha, \beta \geq 1$ , with the usual interior-mode formula when  $\alpha, \beta > 1$ . The reverse KL used to index*  
1572 *the robustness tables is*

$$D_{\text{KL}}(\text{Unif}||\text{Beta}(\alpha, \beta)) = \log B(\alpha, \beta) + \alpha + \beta - 2.$$

1573 *These constants are useful for interpreting the experimental tables, but the robustness mechanism*  
1574 *itself is entirely captured by the profile bounds above.*

1575 The results above are deliberately environment-agnostic and help identify generic failure modes:  
1576 robustness can fail only when the deployment prior emphasizes regions where the frozen controller  
1577 is inaccurate or slow, or when  $\nu'$  leaves the support region on which the controller was effectively  
1578 trained.

1579 **B.9 Sample Complexity: Value Estimation vs Argmax Localization**

1580 In this subsection we study the sample complexity of estimating the maximum value of a gaussian  
 1581 process  $F \sim \text{GP}(0, k_\ell)$  on  $D = [0, 1]^d$ , and consider an RBF kernel  $k_\ell(x, x') := \exp\left(-\frac{\|x-x'\|^2}{2\ell^2}\right)$   
 1582 with  $x, x' \in D$ .

1583 We let the unknown lengthscale be random:

$$\Lambda \sim \nu, \quad \text{supp}(\nu) \subseteq [\ell_-, \ell_+] \subset (0, \infty).$$

1584 The learner observes

$$Y_t = F(a_t) + \xi_t, \quad \xi_t \sim \mathcal{N}(0, \sigma^2),$$

1585 where  $a_t \in D$  is chosen adaptively from the past history and  $\sigma > 0$  is known.

1586 Define

$$F^* = \max_{x \in D} F(x), \quad D^* = \arg \max_{x \in D} F(x),$$

1587 and write  $X^*$  whenever the maximizer is unique.

1588 **Argmax localization complexity.** For  $r > 0$  and  $\delta \in (0, 1)$ , define the Bayesian argmax-  
 1589 localization complexity

$$T_{\text{arg}, \nu}(r, \delta) = \inf_{\mathcal{A}} \mathbb{E}_{\Lambda, F, \xi}[\tau],$$

1590 where the infimum is over all sequential algorithms  $\mathcal{A}$  that output  $\hat{X}$  and satisfy

$$\mathbb{P}_{\Lambda, F, \xi} \left( \inf_{X^* \in D^*} \|\hat{X} - X^*\| \leq r \right) \geq 1 - \delta.$$

1591 **Max-value estimation complexity.** Similarly, define the Bayesian max-value-estimation complex-  
 1592 ity

$$T_{\text{val}, \nu}(\epsilon, \delta) = \inf_{\mathcal{A}} \mathbb{E}_{\Lambda, F, \xi}[\tau],$$

1593 where the infimum is over all sequential algorithms that output  $\hat{v}$  and satisfy

$$\mathbb{P}_{\Lambda, F, \xi} (|\hat{v} - F^*| \leq \epsilon) \geq 1 - \delta.$$

1594 All probabilities and expectations are under the joint hierarchical law of  $(\Lambda, F)$ , the observation noise,  
 1595 and any internal randomness of the algorithm.

1596 **Main result.** To compare the sample complexity of argmax localization vs max-value estimation,  
 1597 we use the fact that under regularity assumptions we approximately have  $r \sim \sqrt{\epsilon}$ . This follows from  
 1598 a Taylor's expansion

$$F(X) = F(X^*) + \frac{1}{2}(X - X^*)^\top \nabla^2 F(X^*)(X - X^*) + o(\|X - X^*\|^2).$$

1599 from which we find  $F(X) - F(X^*) \approx c\|X - X^*\|_2^2$  for a suitable constant. Therefore, localizing to  
 1600 radius  $r$  gives an error of roughly  $cr^2$  on the value. Therefore, for  $\epsilon$  accuracy on the value, we can  
 1601 take the radius to be  $r \sim \sqrt{\epsilon}$ .

1602 Then, we have the following main result.

1603 **Theorem B.5** (Bayesian value-argmax separation). *Consider the hierarchical RBF-GP model:  
 1604  $\Lambda \sim \nu$  with  $\text{supp}(\nu) \subseteq [\ell_-, \ell_+] \subset (0, \infty)$ , and  $F \mid \Lambda = \ell \sim \text{GP}(0, k_\ell)$  on  $D = [0, 1]^d$ , with  
 1605 observations  $Y_t = F(a_t) + \xi_t$ ,  $\xi_t \sim \mathcal{N}(0, \sigma^2)$ .*

1606 *Fix  $\delta \in (0, 1)$ . Suppose there exists  $\eta > 0$  such that  $\beta_\eta = \mathbb{P}_{\Lambda, F}(\mathcal{I}_\eta^c) < \delta$ , where  $\beta_\eta$  is defined in  
 1607 Lemma 17. Then:*

1608 *1. (Argmax upper bound.) There exist constants  $B < \infty$  and  $C < \infty$ , depending on  $\nu, d, \delta, \sigma$   
 1609 but not on  $\epsilon$ , such that for all sufficiently small  $\epsilon > 0$ ,*

$$T_{\text{arg}, \nu}^{\text{Bayes}}(\sqrt{\epsilon}, \delta) \leq B + C \sigma^2 \epsilon^{-3/2} \log \log \frac{1}{\epsilon}.$$

1610 2. **(Value lower bound.)** There exist constants  $c > 0$  and  $c' < \infty$ , depending on  $\nu, d, \delta, \sigma$  but  
 1611 not on  $\epsilon$ , such that for all  $\epsilon > 0$ ,

$$T_{\text{val},\nu}^{\text{Bayes}}(\epsilon, \delta) \geq c \frac{\sigma^2}{\epsilon^2} \log \frac{1}{\delta} - c'.$$

1612 Consequently,

$$\lim_{\epsilon \rightarrow 0} \frac{T_{\text{val},\nu}^{\text{Bayes}}(\epsilon, \delta)}{T_{\text{arg},\nu}^{\text{Bayes}}(\sqrt{\epsilon}, \delta)} = \infty.$$

1613 Thus, under the high-probability interior regularity condition, max-value estimation is asymptotically  
 1614 harder than argmax localization in the fully Bayesian hierarchical RBF model.

1615 *Proof.* Set  $\alpha \in (\beta_\eta, \delta)$ . Since  $\beta_\eta < \alpha$ , Lemma 17 provides deterministic constants  $\rho, \mu, L, M, \Gamma > 0$   
 1616 and an event  $\mathcal{E}_\delta$  with  $\mathbb{P}(\mathcal{E}_\delta) \geq 1 - \delta$ .

1617 *Part 1.* We apply Theorem B.7 applies with  $r = \sqrt{\epsilon}$ , giving

$$T_{\text{arg},\nu}^{\text{Bayes}}(\sqrt{\epsilon}, \delta) \leq B + O\left(\sigma^2 \epsilon^{-3/2} \log \log \frac{1}{\epsilon}\right).$$

1618 *Part 2.* By Theorem B.6, with  $C_- = \exp(-d/(2\ell_-^2))$  and  $C_+ = \exp(d/(4\ell_-^2))$ ,

$$T_{\text{val},\nu}^{\text{Bayes}}(\epsilon, \delta) \geq \frac{\sigma^2}{C_+^2} \left[ \frac{C_-^2 z_{1-\delta/2}^2}{\epsilon^2} - 1 \right]_+.$$

1619 For fixed  $\delta$  and small  $\epsilon$ , using  $z_{1-\delta/2}^2 \geq c_0 \log(1/\delta)$ , this gives

$$T_{\text{val},\nu}^{\text{Bayes}}(\epsilon, \delta) \geq c \frac{\sigma^2}{\epsilon^2} \log \frac{1}{\delta} - c'.$$

1620 *Separation.*

$$\frac{T_{\text{val},\nu}^{\text{Bayes}}(\epsilon, \delta)}{T_{\text{arg},\nu}^{\text{Bayes}}(\sqrt{\epsilon}, \delta)} \geq \frac{c \sigma^2 \epsilon^{-2} \log(1/\delta) - c'}{B + C \sigma^2 \epsilon^{-3/2} \log \log(1/\epsilon)} \rightarrow \infty \quad \text{as } \epsilon \downarrow 0,$$

1621 since  $\epsilon^{-2} \gg \epsilon^{-3/2} \log \log(1/\epsilon)$ . □

1622 So value estimation can be asymptotically harder than argmax localization. Even at the radius where  
 1623 localizing the argmax would "in principle" tell you the value to accuracy  $\epsilon$ , the localization itself is  
 1624 cheaper than directly estimating the value. The reason is simple: the algorithm exploits the geometry  
 1625 (gradient information, smoothness) to localize at a fast rate, but it does not help with the problem of  
 1626 estimating the height of a function.

### 1627 B.9.1 Max-value estimation complexity

1628 To derive a lower bound on  $T_{\text{val},\nu}(\epsilon, \delta)$ , we convert the problem of estimating the maximum into the  
 1629 problem of estimating a parameter  $\Theta$ . We note that the lower bound is not tight, but for our purpose  
 1630 (of showing that the argmax localization is easier) this is not important, as our goal is to show that  
 1631 even this approximate lower bound still yields an harder problem. In particular, we assume the learner  
 1632 has access to a particular quantity  $W_\ell$  that appears in the proof. We obtain the following result (note  
 1633 that the constants are not optimized, and the only goal is to show the dependency on  $\epsilon$ ).

1634 **Theorem B.6.** Consider the problem of estimating  $F^* = \max_{x \in D} F(x)$  where  $D = [0, 1]^d$ ,  $F |$   
 1635  $\Lambda \sim \text{GP}(0, k_\Lambda)$  with kernel  $k_\ell(x, y) = \exp(-\|x - y\|_2^2 / (2\ell^2))$  and  $\Lambda \sim \nu$  with continuous support  
 1636 in  $[\ell_-, \ell_+] \subset (0, \infty)$ . Consider any sequential algorithm  $\mathcal{A}$  that in each round selects  $a_t$  and  
 1637 observes  $Y_t = F(a_t) + \xi_t$ , where  $\xi_t \sim \mathcal{N}(0, 1)$ . If the algorithm outputs  $\hat{v}$  at some stopping time  $\tau$   
 1638 satisfying  $\mathbb{P}_{\Lambda, F}(|\hat{v} - F^*| \leq \epsilon) \geq 1 - \delta$ , then we say that the algorithm is  $(\epsilon, \delta)$ -correct. Then, for  
 1639 any  $(\epsilon, \delta)$ -correct algorithm we have that

$$\inf_{\mathcal{A}: (\epsilon, \delta)\text{-correct}} \mathbb{E}[\tau] \geq \frac{\sigma^2}{C_+^2} \left[ \frac{C_-^2 z_{1-\delta/2}^2}{\epsilon^2} - 1 \right]_+,$$

1640 where  $z_{1-\delta/2} = \Phi^{-1}(1 - \delta/2)^2$ ,  $C_- = \exp(-d/(2\ell_-^2))$  and  $C_+ = \exp(d/(4\ell_-^2))$ .

1641 *Proof.* The proof relies on converting the problem into that of estimating a parameter  $\Theta$ . We use the  
 1642 property of independent Gaussian r.v.: for jointly Gaussian  $X, Y$  we have that  $X = \frac{\text{Cov}(X, Y)}{\text{Var}(Y)}Y + W$ ,  
 1643 where  $W$  is an independent zero-mean Gaussian.

1644 In this case, we take  $Y = \Theta_\ell$  and  $X = F$ , where we define  $\Theta_\ell$  as

$$\Theta_\ell := \frac{1}{s_\ell} \int_D F(x) dx, \quad s_\ell^2 := \int_{D \times D} k_\ell(x, y) dx dy.$$

1645 Computing the covariance, we obtain

$$\begin{aligned} \text{Cov}(F(x), \Theta_\ell) &= \mathbb{E}[F(x)\Theta_\ell], \\ &= \frac{1}{s_\ell} \mathbb{E} \left[ F(x) \int_D F(y) dy \right], \\ &= \frac{1}{s_\ell} \int_D \mathbb{E} [F(x)F(y)] dy, \\ &= \frac{1}{s_\ell} \int_D k_\ell(x, y) dy =: c_\ell(x). \end{aligned}$$

1646 Therefore, we can rewrite the observation using that  $F(x) = c_\ell(x)\Theta_\ell + W_\ell(x)$  for some GP  $W_\ell$   
 1647 with 0-mean, and thus

$$Y_t = [c_\ell(a_t)\Theta_\ell + W_\ell(a_t)] + \xi_t.$$

1648 From the expression of  $F(x)$  we observe that

$$V(\theta) = \max_x [c_\ell(x)\theta + W_\ell(x)],$$

1649 is increasing in  $\theta$ . Denote a maximizer by  $x_\theta$ , then, for  $\theta' \neq \theta$  we have

$$V(\theta') \geq c_\ell(x_\theta)\theta' + W_\ell(x_\theta) = V(\theta) + c_\ell(x_\theta)(\theta' - \theta).$$

1650 Similarly,

$$V(\theta) \geq c_\ell(x_{\theta'})\theta + W_\ell(x_{\theta'}) = V(\theta') + c_\ell(x_{\theta'})(\theta - \theta').$$

1651 We now derive bounds on  $c_\ell$ . From the definition

$$c_\ell(x) = \frac{1}{s_\ell} \int_D k_\ell(x, y) dy,$$

1652 since  $D$  is compact, and  $k_\ell$  is the RBF kernel, we have that  $k_\ell(x, y) \leq 1 \Rightarrow c_\ell(x) \leq 1/s_\ell$  and  
 1653  $k_\ell(x, y) \geq \exp(-d/(2\ell^2)) \Rightarrow c_\ell(x) \geq \exp(-d/(2\ell^2))/s_\ell$ . Since  $s_\ell^2 \geq \exp(-d/(2\ell^2))$  and  
 1654  $s_\ell^2 \leq 1$ , we find

$$\underbrace{\exp(-d/(2\ell_-^2))}_{=: C_-} \leq c_\ell(x) \leq \frac{1}{\underbrace{\exp(-d/(4\ell_-^2))}_{=: C_+}}.$$

1655 Then, from the bounds above on  $V$  we find

$$V(\theta') \geq V(\theta) + C_-(\theta' - \theta), \quad V(\theta) \geq V(\theta') - C_+(\theta' - \theta),$$

1656 leading to

$$C_-(\theta' - \theta) \leq V(\theta') - V(\theta) \leq C_+(\theta' - \theta).$$

1657 Therefore  $V$  is bi-Lipschitz, and thus invertible. Choosing  $u = V(\theta)$  and  $v = V(\theta')$ , we obtain

$$|V^{-1}(u) - V^{-1}(v)| \leq \frac{1}{C_-} |u - v|.$$

1658 Therefore, with  $\hat{v} = V(\hat{\Theta})$  and  $F^* = V(\Theta)$  we get

$$|\hat{\Theta} - \Theta| \leq \frac{1}{C_-} |\hat{v} - F^*|.$$

---

<sup>2</sup>Inverse of the CDF of a standard normal distribution.

1659 Hence, assuming the algorithm has access to  $W_\ell$ , we can obtain a lower bound on the sample  
 1660 complexity by using Lemma 16. For any  $(\epsilon', \delta)$ -algorithm that estimates  $\Theta$ , we can choose  $\epsilon' = \epsilon/C_-$   
 1661 to obtain  $|V^{-1}(\hat{\Theta}) - V^{-1}(\Theta)| \leq \epsilon'$  at the stopping time. Therefore,

$$\inf_{\mathcal{A}: (\epsilon, \delta)\text{-correct}} \mathbb{E}[\tau] \geq \frac{\sigma^2}{C_+^2} \left[ \frac{C_-^2 z_{1-\delta/2}^2}{\epsilon^2} - 1 \right]_+,$$

1662

□

1663 **Lemma 16** (Scalar estimation lemma). *Let  $\Theta \sim \mathcal{N}(0, 1)$ . Suppose to observe  $Z_t = C_t \Theta + \xi_t$  with  
 1664  $|C_t| \leq C$  chosen adaptively and  $\xi_t \sim \mathcal{N}(0, \sigma^2)$ . Let  $\tau$  be a stopping time such that at  $\tau$  the algorithm  
 1665 outputs  $\hat{\Theta}$  satisfying  $\mathbb{P}(|\hat{\Theta} - \Theta| \leq \epsilon) \geq 1 - \delta$  for  $\epsilon > 0, \delta \in (0, 1)$ . Then*

$$\inf_{\mathcal{A}: (\epsilon, \delta)\text{-correct}} \mathbb{E}[\tau] \geq \frac{\sigma^2}{C^2} \left[ \frac{z_{1-\delta/2}^2}{\epsilon^2} - 1 \right]_+,$$

1666 where  $z_{1-\delta/2} = \Phi^{-1}(1 - \delta/2)$ .

1667 *Proof.* Let  $\mathcal{H}_t$  denote the filtration history after  $t$  observations. Since everything is Gaussian, also the  
 1668 posterior law  $\Theta | \mathcal{H}_t$  is Gaussian, of parameter  $(m_t, v_t)$ . In particular, the precision is

$$P_t = v_t^{-1} = 1 + \frac{1}{\sigma^2} \sum_{s=1}^t C_s^2.$$

1669 After stopping, for any estimate  $u$  we have

$$\mathbb{P}(|u - \Theta| \leq \epsilon | \mathcal{H}_\tau) = \mathbb{P}(\Theta \in [u - \epsilon, u + \epsilon] | \mathcal{H}_\tau).$$

1670 This quantity is maximized when the interval is centered in  $m_\tau$ , therefore

$$\mathbb{P}(|u - \Theta| \leq \epsilon | \mathcal{H}_\tau) \leq \Phi\left(\frac{\epsilon}{\sqrt{v_\tau}}\right) - \Phi\left(-\frac{\epsilon}{\sqrt{v_\tau}}\right) = 2\Phi\left(\frac{\epsilon}{\sqrt{v_\tau}}\right) - 1.$$

1671 Then, for any  $(\epsilon, \delta)$ -PAC algorithm we have

$$1 - \delta \leq \mathbb{E}[\mathbb{P}(|u - \Theta| \leq \epsilon | \mathcal{H}_\tau)] \leq 2\mathbb{E}[\Phi(\epsilon\sqrt{P_\tau})] - 1.$$

1672 Since  $\Phi$  is concave, the argument is increasing and concave in  $P_\tau$ , then  $s \mapsto \Phi(\epsilon\sqrt{s})$  is concave, we  
 1673 also obtain

$$1 - \delta \leq 2\Phi\left(\epsilon\sqrt{\mathbb{E}[P_\tau]}\right) - 1.$$

1674 Taking the inverse, and defining  $z_{1-\delta/2} = \Phi^{-1}(1 - \delta/2)$ , we find

$$z_{1-\delta/2} \leq \epsilon\sqrt{\mathbb{E}[P_\tau]}.$$

1675 Hence, we are just left with bounding  $P_\tau$  :

$$P_\tau \leq 1 + \frac{1}{\sigma^2} C^2 \tau,$$

1676 from which we get

$$z_{1-\delta/2} \leq \epsilon\sqrt{1 + \frac{C^2}{\sigma^2} \mathbb{E}[\tau]}.$$

1677 Therefore

$$\mathbb{E}[\tau] \geq \frac{\sigma^2}{C^2} \left[ \frac{z_{1-\delta/2}^2}{\epsilon^2} - 1 \right]_+$$

1678

□

1679 **B.9.2 Argmax localization complexity**

1680 We now study the problem of locating the argmax of a GP. Recall that for  $r > 0$  and  $\delta \in (0, 1)$ , define  
 1681 the Bayesian argmax-localization complexity

$$T_{\text{arg},\nu}(r, \delta) = \inf_{\mathcal{A}} \mathbb{E}_{\Lambda, F, \xi}[\tau],$$

1682 where the infimum is over all sequential algorithms  $\mathcal{A}$  that output  $\hat{X}$  and satisfy

$$\mathbb{P}_{\Lambda, F, \xi} \left( \inf_{X^* \in D^*} \|\hat{X} - X^*\| \leq r \right) \geq 1 - \delta.$$

1683 Our goal is to provide a meaningful upperbound on  $T_{\text{arg},\nu}(r, \delta)$ . We provide an algorithm, T-BAL  
 1684 (Two-Stage Bayesian Argmax Localization; see Algorithm 2), that locates the argmax of a GP with a  
 1685 finite number of samples. The algorithm works in two phases: we first locate the nice region where  
 1686 the argmax lies, and then perform gradient ascent on that region. The analysis relies on the argmax  
 1687 being away from the boundary of  $D$ . Therefore, we introduce the following regularity event.

1688 **Definition 1** (Interior regularity event). For  $\eta > 0$ , let  $\text{dist}(D^*, \partial D) = \inf_{X \in D^*} \text{dist}(X, \partial D)$  and  
 1689 define

$$\mathcal{I}_\eta = \{\text{dist}(D^*, \partial D) \geq \eta\}, \quad \beta_\eta = \mathbb{P}_{\Lambda, F}(\mathcal{I}_\eta^c).$$

1690 One can show that on  $\mathcal{I}_\eta$ , the set of maximizer is a singleton almost surely (i.e.,  $D^* = \{X^*\}$ ), and  
 1691 similarly one can show that the Hessian is non-degenerate in  $X^*$ . In fact, we have the following.

1692 **Remark 11.** *An RBF-GP on  $D$  is a.s.  $C^\infty$ , and its restriction to the interior  $\text{int}(D)$  is a.s. a Morse*  
 1693 *function (every critical point has invertible Hessian, all critical values are distinct, and there are*  
 1694 *finitely many critical points). If the maximizer is not attained at the boundary, then it is unique almost*  
 1695 *surely.*

1696 We work under  $\mathcal{I}_\eta$ : this not only allows the maximizer to be unique, but also allows us to be at-least  
 1697 at a distance  $\eta$  from the boundary, where it is more degenerate. Then, we can show the existence of  
 1698 the following constants.

1699 **Lemma 17** (Regularity under the RBF prior). *Define  $B(x, \rho)$  to be the ball centered around  $x$  of*  
 1700 *radius  $\rho$  with some norm  $\|\cdot\|$ . Fix  $\eta > 0$  and suppose  $\beta_\eta < 1$ . For every  $\alpha \in (\beta_\eta, 1)$  there exist*  
 1701 *deterministic constants  $\rho_\alpha, \mu_\alpha, L_\alpha, M_\alpha, \Gamma_\alpha > 0$  and an event  $\mathcal{E}_\alpha \subseteq \mathcal{I}_\eta$  such that  $\mathbb{P}_{\Lambda, F}(\mathcal{E}_\alpha) \geq 1 - \alpha$ ,*  
 1702 *and on  $\mathcal{E}_\alpha$  the following properties hold:*

$$B(X^*, \rho_\alpha) \subset D,$$

1703

$$\mu_\alpha I \preceq -\nabla^2 F(x) \preceq L_\alpha I \quad \forall x \in B(X^*, \rho_\alpha),$$

1704

$$\sup_{x \in B(X^*, \rho_\alpha)} \|\nabla^3 F(x)\|_{\text{op}} \leq M_\alpha,$$

1705 and

$$F^* - \sup_{x \notin B(X^*, \rho_\alpha/2)} F(x) \geq \Gamma_\alpha.$$

1706 *Proof.* Let  $\gamma = \alpha - \beta_\eta$ .

1707 Under  $\mathcal{I}_\eta$  the maximizer  $X^*$  is in the interior, and unique almost surely. Define  $\lambda^* =$   
 1708  $\lambda_{\min}(-\nabla^2 F(X^*))$ : under  $\mathcal{I}_\eta$  we have that  $\lambda^* > 0$  almost surely, hence, there exists  $\mu_\alpha > 0$   
 1709 s.t.  $\mathbb{P}(\mathcal{I}_\eta \cap \{\lambda^*/2 < \mu_\alpha\}) \leq \gamma/5$ .

1710 Next, we use the fact that  $F \in C^\infty$  and  $D$  is compact to obtain an upper bound on the Hessian:

$$\sup_{x \in D} \|\nabla^2 F(x)\|_{\text{op}} \leq L_0 < \infty.$$

1711 Therefore, there exists  $L_\alpha < \infty$  such that  $\mathbb{P}(L_\alpha < L_0) \leq \gamma/5$ . With a similar reasoning, we also  
 1712 obtain  $\sup_{x \in D} \|\nabla^3 F(x)\|_{\text{op}} \leq M_0$ , and thus there exists  $M_\alpha < \infty$  such that  $\mathbb{P}(M_\alpha < M_0) \leq \gamma/5$ .

---

**Algorithm 2** Two-stage Bayesian argmax localization (T-BAL)
 

---

**Require:** Target radius  $r$ , confidence  $\delta$ , noise level  $\sigma^2$ , constants  $\rho, \mu, L, M, \Gamma$  from Lemma 17

▷ **Stage 1: Coarse grid search**

- 1: Set  $h \leftarrow \min\{\rho/8, \sqrt{3\Gamma/(2L)}\}$
  - 2: Construct  $h$ -net  $\mathcal{G}$  of  $D$  with  $|\mathcal{G}| \leq C_d h^{-d}$
  - 3: Set  $n_0 \leftarrow \lceil 128\sigma^2\Gamma^{-2} \log(2|\mathcal{G}|/\delta_0) \rceil$
  - 4: **for** each  $g \in \mathcal{G}$  **do**
  - 5:   Query  $F(g)$  exactly  $n_0$  times; compute sample mean  $\hat{F}(g)$
  - 6: **end for**
  - 7:  $x_0 \leftarrow \arg \max_{g \in \mathcal{G}} \hat{F}(g)$
- ▷ **Stage 2: Local finite-difference gradient ascent**
- 8: Set  $K \leftarrow \lceil \frac{4L}{3\mu} \log \frac{\rho}{2r} \rceil$ ,  $q \leftarrow 1 - \frac{3\mu}{4L}$ ,  $e_0 \leftarrow \rho/2$
  - 9: **for**  $k = 0, \dots, K - 1$  **do**
  - 10:    $e_k \leftarrow e_0 q^k$
  - 11:    $s_k \leftarrow \min \left\{ \rho/8, \sqrt{3\mu e_k / (4\sqrt{d} M)} \right\}$
  - 12:    $n_k \leftarrow \lceil 64d \cdot \sigma^2 \mu^{-2} e_k^{-2} s_k^{-2} \log(2dK/\delta_1) \rceil$
  - 13:   **for**  $j = 1, \dots, d$  **do**
  - 14:     Query  $F(x_k + s_k e_j)$  and  $F(x_k - s_k e_j)$  each  $n_k$  times
  - 15:      $\hat{g}_{k,j} \leftarrow \frac{\bar{Y}(x_k + s_k e_j) - \bar{Y}(x_k - s_k e_j)}{2s_k}$
  - 16:   **end for**
  - 17:    $x_{k+1} \leftarrow x_k + \frac{1}{L} \hat{g}_k$
  - 18: **end for**
  - 19: **return**  $\hat{X} = x_K$
- 

1713 Next, by continuity and the margin condition of  $\mathcal{I}_\eta$ , the following quantity exists and is strictly  
 1714 positive almost surely

$$\rho_0(\omega) = \sup \left\{ q < \eta : \sup_{X \in B(X^*, q)} \|\nabla^2 F(X) - \nabla^2 F(X^*)\|_{\text{op}} \leq \frac{\lambda^*}{2} \right\},$$

1715 where  $\omega$  denotes a realization of  $F$ . Therefore, there exists  $\rho_\alpha > 0$  such that  $\mathbb{P}(\mathcal{I}_\eta \cap \{\rho_0 < \rho_\alpha\}) \leq$   
 1716  $\gamma/5$ .

1717 Define then  $\Gamma(\rho_\alpha) = F^* - \sup_{X \notin B(X^*, \rho_\alpha/2)} F(X)$ : under  $\mathcal{I}_\eta$  this is strictly positive since the  
 1718 supremum does not attain  $F^*$ . Hence, there exists  $\Gamma_\alpha > 0$  such that  $\mathbb{P}(\mathcal{I}_\eta \cap \{\Gamma(\rho_\alpha) < \Gamma_\alpha\}) \leq \gamma/5$ .

1719 Define

$$\mathcal{E} = \mathcal{I}_\eta \cap \{\lambda^*/2 \geq \mu_\alpha\} \cap \{L_\alpha \geq L_0\} \cap \{M_\alpha \geq M_0\} \cap \{\rho_0 \geq \rho_\alpha\} \cap \{\Gamma(\rho_\alpha) \geq \Gamma_\alpha\}.$$

1720 Then, since

$$\mathcal{E}^c = \mathcal{I}_\eta^c \cup \{\lambda^*/2 < \mu_\alpha\} \cup \dots = \mathcal{I}_\eta^c \cup (\mathcal{I}_\eta \cap \{\lambda^*/2 < \mu_\alpha\}) \cup \dots$$

1721 we have  $\mathbb{P}(\mathcal{E}^c) \leq \beta_\eta + 5\frac{\gamma}{5} = \alpha$ . Hence, under  $\mathcal{E}$  we have that  $B(X^*, \rho_\alpha) \subset D$ , and all the properties  
 1722 follow quite immediately. We only show the lower bound on the Hessian: for any  $X \in B(X^*, \rho_\alpha)$ :

$$-\frac{\lambda^*}{2} I \preceq \nabla^2 F(x) - \nabla^2 F(X^*) \preceq \frac{\lambda^*}{2} I.$$

1723 Hence

$$-\nabla^2 F(x) \succeq -\nabla^2 F(X^*) - \frac{\lambda^*}{2} I \succeq \frac{\lambda^*}{2} I \succeq \mu_\alpha I.$$

1724

□

1725 **Argmax upper bound.** Under the event  $\mathcal{I}_\eta$ , we want to construct an algorithm that localizes  $X^*$  to  
 1726 within radius  $r$  with probability  $1 - \delta$ . The idea is to show an upper bound on the minimal lower  
 1727 bound of the type

$$T_{\text{arg},\nu}(r, \delta) := \inf_{\mathcal{A}} \mathbb{E}[\tau] \leq B_\alpha + O(\sigma^2 r^{-\gamma} \log(1/r)),$$

1728 for some  $\gamma > 0$ .

1729 Under  $\mathcal{I}_\eta$ , we have the guarantees from Lemma 17, and thus

$$F(X) = F(X^*) + \frac{1}{2}(X - X^*)^\top \nabla^2 F(X^*)(X - X^*) + o(\|X - X^*\|^2).$$

1730 from which we find  $F(X) - F(X^*) \approx \frac{\mu_\alpha}{2} \|X - X^*\|_2^2$ . Therefore, localizing to radius  $r$  gives an  
 1731 error of roughly  $\frac{\mu_\alpha r^2}{2}$  on the value. Therefore, for  $\epsilon$  accuracy on the value, we can take the radius to  
 1732 be  $r \sim \sqrt{\epsilon}$

1733 Therefore, under  $\mathcal{I}_\eta$  if  $\beta_\eta < 1$ , and

$$\liminf_{\epsilon \rightarrow 0} \frac{T_{\text{val},\nu}(\epsilon, \delta)}{T_{\text{arg},\nu}(\sqrt{\epsilon}, \delta)} \rightarrow \infty$$

1734 one can argue that the problem of estimating the max-value is intrinsically harder than the problem of  
 1735 estimating the argmax as the accuracy radius decreases. In particular, considering the lower bound on  
 1736 the max-value estimation problem Theorem B.6, and the proposed upper bound on  $T_{\text{arg},\nu}(r, \delta)$ , we  
 1737 obtain that

$$\frac{T_{\text{val},\nu}(\epsilon, \delta)}{T_{\text{arg},\nu}(\sqrt{\epsilon}, \delta)} \geq \frac{\Omega(\epsilon^{-2})}{O(\epsilon^{-\gamma/2} \log \log(1/\sqrt{\epsilon}))},$$

1738 which diverges for  $\gamma \in (0, 4)$  as  $\epsilon \rightarrow 0$ .

1739 So value estimation is asymptotically harder than argmax localization. Even at the radius where  
 1740 localizing the argmax would "in principle" tell you the value to accuracy  $\epsilon$ , the localization itself is  
 1741 cheaper than directly estimating the value. The reason is simple: the algorithm exploits the geometry  
 1742 (gradient information, smoothness) to localize at a fast rate, but it does not help with the problem of  
 1743 estimating the height of a function.

1744 **Analysis of T-BAL (Algorithm 2).** We provide now an algorithm for argmax localization in  
 1745 Gaussian processes. The algorithm, Two-Stage Bayesian Argmax Localization (T-BAL), outlined in  
 1746 Algorithm 2, works in two phases. In the first phase we try to find a point inside  $B(X^*, \rho/2)$ , where  
 1747  $\rho$  is described in Lemma 17. In the second phase, assuming we are inside the above ball, we perform  
 1748 gradient ascent to find the maximum using finite differences to approximate the gradients.

1749 The second phase analysed gradient descent when  $x_0 \in B(X^*, \rho)$ .

1750 **Theorem B.7** (Sample Complexity of T-BAL). *Consider the event  $\mathcal{E}_\alpha$  in Lemma 17. Set  $r > 0$ ,  $\delta \in$   
 1751  $(0, 1)$ ,  $\alpha \in (\beta_\eta, \delta)$  and  $\delta_0 = \delta_1 = (\delta - \alpha)/2$ . Then, T-BAL (Algorithm 2) satisfies  $\mathbb{P}_{F,\Lambda}(\|\hat{X} - X^*\| \leq$   
 1752  $r) \geq 1 - \delta$ , using at-most*

$$B_\alpha + O\left(\frac{\sigma^2}{\mu^2} \max(2, L/\mu) \left[\frac{d^2}{\rho^2 r^2} + \frac{d^{5/2} M}{\mu r^3}\right] \log \frac{Ld \log(\rho/r)}{\mu \delta_1}\right)$$

1753 samples, where  $B_\alpha$  is an appropriate finite constant for each  $\alpha$  that does not depend on  $r$ .

1754 *Proof.* In this proof we consider the analysis of the second stage, while the first stage and the constant  
 1755  $B_\alpha$  are provided in Proposition 11. The idea is to prove a bound on gradient ascent. We first show  
 1756 that one step of the ideal gradient ascent brings us closer to  $X^*$ . We then find the value of  $s_k$  and  $n_k$   
 1757 to compute the approximate gradient up to the desired accuracy. After that, we estimate how many  
 1758 iterations we need, and compute the total number of required samples.

1759 One step gradient ascent. Let  $x^+ = x + \frac{1}{L} \nabla F(x)$ .

1760 Define  $\phi(t) = \nabla F(X^* + t(x - X^*))$ . Since  $\nabla F(X^*) = 0$ , using the fundamental theorem of  
 1761 calculus we have

$$\nabla F(x) = \phi(1) - \phi(0) = \int_0^1 \phi'(t) dt = (x - X^*) \underbrace{\int_0^1 \nabla^2 F(x^* + t(x - X^*)) dt}_{=: -A(x)}.$$

1762 Then

$$x^+ - X^* = x - \frac{1}{L}(x - X^*)A(x) - X^* = (I - \frac{A(x)}{L})(x - X^*).$$

1763 The matrix  $(I - \frac{A(x)}{L})$  has eigenvalues in  $[0, 1 - \mu/L]$  for  $x \in B(X^*, \rho)$ . Therefore,

$$\|x^+ - X^*\| \leq \left(1 - \frac{\mu}{L}\right) \|x - X^*\|$$

1764 so the contraction factor is  $q_0 = 1 - \mu/L$ .

1765 *Noisy gradients.* Suppose we do not have access to the exact gradient, but only to a noisy gradient  $\hat{g}_k$   
1766 in round  $k$ . Assume the noisy gradient satisfy

$$\|\hat{g}_k - \nabla F(x_k)\| \leq \frac{\mu e_k}{4},$$

1767 where  $e_k = e_0 q^k$ . Then

$$\begin{aligned} \|x_{k+1} - X^*\| &= \|x_k + \frac{1}{L}\hat{g}_k - X^*\|, \\ &= \|x_k + \frac{1}{L}\nabla F(x_k) + \frac{1}{L}\hat{g}_k - \frac{1}{L}\nabla F(x_k) - X^*\|, \\ &\leq \|x_k + \frac{1}{L}\nabla F(x_k) - X^*\| + \frac{1}{L}\|\hat{g}_k - \nabla F(x_k)\|, \\ &\leq \left(1 - \frac{\mu}{L}\right) \|x_k - X^*\| + \frac{\mu e_k}{4L}, \\ &\leq \left(1 - \frac{3\mu}{4L}\right) e_k. \end{aligned}$$

1768 So noisy gradients still guarantee convergence as long as we can show  $\|\hat{g}_k - \nabla F(x_k)\| \leq \frac{\mu e_k}{4}$ .

1769 *Finite-difference gradient bound: bias term.* We now bound  $\|\hat{g}_k - \nabla F(x_k)\|$  through a bias-variance  
1770 decomposition:

$$\|\hat{g}_k - \nabla F(x_k)\| \leq \|\hat{g}_k - \mathbb{E}[\hat{g}_k]\| + \|\mathbb{E}[\hat{g}_k] - \nabla F(x_k)\|.$$

1771 We begin with the bias term  $\|\mathbb{E}[\hat{g}_k] - \nabla F(x_k)\|$ . Note that the noisy gradients are computed as follow  
1772 for each direction  $e_j$ :

$$\hat{g}_{k,j} \leftarrow \frac{\bar{Y}(x_k + s_k e_j) - \bar{Y}(x_k - s_k e_j)}{2s_k},$$

1773 where  $\bar{Y}$  is an average over  $n_k$  samples of the values observed. The idea is to bound  $|\mathbb{E}[\hat{g}_{k,j}] -$   
1774  $\partial_j F(x_k)|$  using Taylor series and the fact that  $\nabla^3 F$  is bounded in  $B(X^*, \rho)$  by Lemma 17.

1775 Then, fix a direction  $e_j$  and let  $\phi(t) = F(x_k + t e_j)$ . Then

$$\begin{aligned} \phi(s_k) &= \phi(0) + \phi'(0)s_k + \frac{1}{2}\phi''(0)s_k^2 + \frac{s_k^3}{6}\phi'''(\xi_+), \\ \phi(-s_k) &= \phi(0) - \phi'(0)s_k + \frac{1}{2}\phi''(0)s_k^2 - \frac{s_k^3}{6}\phi'''(\xi_-), \end{aligned}$$

1776 where  $\xi_+ \in (0, s_k)$ ,  $\xi_- \in (-s_k, 0)$ . Then

$$\frac{\phi(s_k) - \phi(-s_k)}{2s_k} = \phi'(0) + \frac{s_k^2}{12}[\phi'''(\xi_+) + \phi'''(\xi_-)].$$

1777 Hence, for  $s_k$  sufficiently small,  $s_k < \rho/2$  (we choose  $s_k < \rho/8$ ) we have  $|\phi'''| \leq M$ , and therefore

$$\left| \frac{\phi(s_k) - \phi(-s_k)}{2s_k} - \phi'(0) \right| \leq \frac{s_k^2 M}{6},$$

1778 leading to

$$\|\mathbb{E}[\hat{g}_k] - \nabla F(x_k)\| \leq \frac{\sqrt{d} s_k^2 M}{6}.$$

1779 Setting  $\frac{\sqrt{d}s_k^2 M}{6} \leq \mu e_k/8$  yields

$$s_k \leq \sqrt{\frac{3\mu e_k}{4\sqrt{d}M}}$$

1780 and we set  $s_k = \min \left\{ \rho/8, \sqrt{\frac{3\mu e_k}{4\sqrt{d}M}} \right\}$ . Furthermore, note that for this choice of  $s_k$  we can guarantee  
 1781 that we stay inside the nice region  $B(X^*, \rho)$ : since  $\|x_k - X^*\| \leq e_k \leq \rho/2$  and  $s_k \leq \rho/8$ , then

$$\|x_k \pm s_k e_j - X^*\| \leq \frac{5\rho}{8} < \rho.$$

1782 Finite-difference gradient bound: variance term. We now bound the variance term  $\|\hat{g}_k - \mathbb{E}[\hat{g}_k]\|$ . We  
 1783 have

$$\hat{g}_{k,j} - \mathbb{E}[\hat{g}_{k,j}] = \frac{\bar{\xi}(x_k + s_k e_j) - \bar{\xi}(x_k - s_k e_j)}{2s_k}, \quad \bar{\xi}(x) = \frac{1}{n_k} \sum_{i=1}^{n_k} \xi_i.$$

1784 Therefore,  $\text{Var}(\hat{g}_{k,j}) = \frac{2\sigma^2}{4s_k^2 n_k}$ . Then

$$\mathbb{P}(|\hat{g}_{k,j} - \mathbb{E}[\hat{g}_{k,j}]| > \mu e_k / (8\sqrt{d})) \leq 2 \exp\left(-\frac{\mu^2 e_k^2 s_k^2 n_k}{64d \cdot \sigma^2}\right).$$

1785 Set the right hand-side smaller than  $\delta_1/(dK)$ , where  $K$  is the total number of iterations  $k =$   
 1786  $0, \dots, K-1$ , to obtain

$$n_k \geq \frac{64d \cdot \sigma^2}{\mu^2 e_k^2 s_k^2} \log \frac{2dK}{\delta_1}.$$

1787 Then, a union bound over  $j = 1, \dots, d$  and  $k = 0, \dots, K-1$ , yields

$$\mathbb{P}(\exists j \in \{1, \dots, d\}, k \in \{0, \dots, K-1\} : |\hat{g}_{k,j} - \mathbb{E}[\hat{g}_{k,j}]| > \mu e_k / (8\sqrt{d})) \leq \delta_1.$$

1788 Hence, with probability  $1 - \delta_1$  we have that

$$\|\hat{g}_k - \nabla F(x_k)\| \leq \frac{\mu e_k}{8} + \frac{\mu e_k}{8} = \frac{\mu e_k}{4},$$

1789 which is what we wanted to show.

1790 Sample complexity. Since we need  $e_K \leq r$ , and  $e_K = e_0 q^K$ , we have

$$K = \left\lceil \frac{\log(e_0/r)}{\log(1/q)} \right\rceil.$$

1791 Since  $\log(1/q) - \log q = -\log(1 - 3\mu/4L) \geq 3\mu/4L$ , we have

$$K \leq 4L \frac{\log(e_0/r)}{3\mu} \leq 4L \frac{\log(\rho/(2r))}{3\mu},$$

1792 where we used  $e_0 \leq \rho/2$ . Then, summing the number of samples from  $k = 0, \dots, K-1$ , we obtain

$$\sum_{k=0}^{K-1} 2dn_k = \sum_{k=0}^{K-1} \frac{128d^2 \cdot \sigma^2}{\mu^2 e_k^2 s_k^2} \log \frac{2dK}{\delta_1}.$$

1793 Use that  $\frac{1}{s_k^2} \leq \frac{64}{\rho^2} + \frac{4\sqrt{d}M}{3\mu e_k}$  and that  $e_k = e_0 q^k$  with  $q \in (0, 1)$ :

$$\begin{aligned} \sum_{k=0}^{K-1} 2dn_k &\leq \sum_{k=0}^{K-1} \left[ \frac{128 \cdot 64 \cdot d^2 \cdot \sigma^2}{\mu^2 e_k^2 \rho^2} + \frac{512d^{5/2} M \cdot \sigma^2}{3\mu^3 e_k^3} \right] \log \frac{2dK}{\delta_1}, \\ &\lesssim \log \frac{2dK}{\delta_1} \sum_{k=0}^{K-1} \left[ \frac{1}{e_k^2} + \frac{1}{e_k^3} \right]. \end{aligned}$$

1794 Regarding the series,

$$\sum_{k=0}^{K-1} \frac{1}{e_k^p} = \frac{1}{e_0^p} \sum_{k=0}^{K-1} \frac{1}{q^{pk}} \leq e_0^{-p} \frac{1 - q^{-pK}}{1 - q^{-p}} \leq \frac{e_0^{-p} q^{-pK}}{q^{-p} - 1} = \frac{e_K^{-p}}{q^{-p} - 1}.$$

1795 Now, note that since  $K$  is the smallest integer achieving  $e_K = e_0 q^K \leq r$ , we have  $e_0 q^{K-1} \geq r$ , and  
 1796 thus  $e_K = q e_{K-1} \geq q r \Rightarrow (q r)^{-p} \geq e_K^{-p}$ . We obtain

$$\sum_{k=0}^{K-1} \frac{1}{e_k^p} \leq \frac{(q r)^{-p}}{q^{-p} - 1} = \frac{r^{-p}}{1 - q^p}.$$

1797 To lower bound the denominator, recall that  $q$  is of the form  $(1 - x)$ . Using that  $(1 - x)^p \leq e^{-xp}$  for  
 1798  $x \in (0, 1)$ , we have  $1 - q^p \geq 1 - e^{-p3\mu/(4L)}$ .

1799 Since  $1 - e^{-t} \geq t/2$  for  $t \in (0, 1)$ , for  $p$  sufficiently small we obtain  $1 - q^p \geq \frac{3p\mu}{8L}$ . If  $p$  is large,  
 1800 such that  $t \geq 1$ , then  $1 - e^{-t} \geq 1 - e^{-1} \geq 1/2$ . Therefore,  $\frac{1}{1 - q^p} \leq \max(2, \frac{8L}{3p\mu})$ , and

$$\sum_{k=0}^{K-1} \frac{1}{e_k^p} \leq r^{-p} \max\left(2, \frac{8L}{3p\mu}\right).$$

1801 Since  $p \in \{2, 3\}$  we also have  $\sum_{k=0}^{K-1} \frac{1}{e_k^p} \leq r^{-p} \max\left(2, \frac{8L}{6\mu}\right)$ .

1802 Then, we conclude that the second phase sample complexity is upper bounded by

$$\begin{aligned} \sum_{k=0}^{K-1} 2dn_k &\leq \frac{128\sigma^2}{\mu^2} \cdot \max\left(2, \frac{8L}{6\mu}\right) \cdot \left[\frac{64d^2}{\rho^2 r^2} + \frac{4d^{5/2}M}{3\mu r^3}\right] \log \frac{2dK}{\delta_1}, \\ &\leq \frac{128\sigma^2}{\mu^2} \cdot \max\left(2, \frac{8L}{6\mu}\right) \cdot \left[\frac{64d^2}{\rho^2 r^2} + \frac{4d^{5/2}M}{3\mu r^3}\right] \log \frac{8Ld \log(\rho/(2r))}{3\mu\delta_1}. \end{aligned}$$

1803 *Connecting everything together.* In phase 1 we have probability  $\delta_0$  of failure, while in phase 2 we  
 1804 have probability  $\delta_1$  of failure. Since we work under the event  $\mathcal{E}_\alpha$  with failure probability  $\alpha$ , we have  
 1805  $\mathbb{P}(\text{failure}) \leq \delta_0 + \delta_1 + \alpha = \delta$ .  $\square$

1806 The first phase analysis is provided in the following proposition. We construct an  $h$ -net  $\mathcal{G}$  of  $D$  that  
 1807 depends on the geometric of the problem (see Lemma 17). For each point in  $\mathcal{G}$ , we sample  $F(g)$   
 1808 exactly  $n_0$  times, such that we have good concentration, and we return the point that achieves the  
 1809 maximum.

1810 **Proposition 11.** *Consider Algorithm 2, and let  $\delta_0 \in (0, 1)$ . Under  $\mathcal{E}_\alpha$  (see Lemma 17), the first  
 1811 phase samples  $B_\alpha \leq 2^d h^{-d} \lceil 128\sigma^2 \Gamma^{-2} \log(2^{d+1} h^{-d} / \delta_0) \rceil$  queries. Furthermore, we have that  
 1812  $x_0 \in B(X^*, \rho/2)$  with probability  $1 - \delta_0$ .*

1813 *Proof.* Consider the event  $\mathcal{E}_\alpha$  in Lemma 17, and omit the subscript  $\alpha$  for simplicity. We construct an  
 1814  $h$ -net  $\mathcal{G}$  of  $D$  such that for every  $x \in D$  there exists  $g \in \mathcal{G}$  satisfying  $\|x - g\| \leq h$ . Hence, there  
 1815 exists a point  $g^* \in \mathcal{G}$  satisfying  $\|X^* - g^*\| \leq h$ .

1816 At each grid point  $g$  take  $n_0$  samples and compute the sample mean  $\hat{F}(g)$ : by the Gaussian tail bound,  
 1817 we have

$$\mathbb{P}(|\hat{F}(g) - F(g)| > \Gamma/8) \leq 2 \exp\left(-\frac{n_0 \Gamma^2}{128\sigma^2}\right).$$

1818 Choosing  $n_0 = \frac{128\sigma^2}{\Gamma^2} \log\left(\frac{2|\mathcal{G}|}{\delta_0}\right)$ , for some  $\delta_0 \in (0, 1)$ , and taking a union bound over  $g$ , we obtain

$$\mathbb{P}(\exists g \in \mathcal{G} : |\hat{F}(g) - F(g)| > \Gamma/8) \leq \delta_0.$$

1819 We now choose  $h$  small enough such that a lower bound on  $\hat{F}(g^*)$  upper bounds a valid upper bound  
 1820 on  $\hat{F}(g)$  for  $g \notin B(X^*, \rho/2)$ . Under the event  $\mathcal{E} = \{\forall g \in \mathcal{G} : |\hat{F}(g) - F(g)| \leq \Gamma/8\}$ , we have that  
 1821 for  $g \notin B(X^*, \rho/2)$

$$\hat{F}(g^*) \geq F(g^*) - \Gamma/8, \quad \text{and} \quad \hat{F}(g) \leq F(g) + \Gamma/8 \leq F^* - \Gamma + \Gamma/8 = F^* - \frac{7}{8}\Gamma,$$

1822 where we used the fact from Lemma 17 that  $F^* - F(g) \geq \Gamma$  for  $g \notin B(X^*, \rho/2)$ . To construct a  
 1823 lower bound on  $F(g^*)$ , we use the gradient properties of  $F$  in  $B(X^*, \rho)$ . To that aim, we need to  
 1824 ensure  $g^*$  is sufficiently inside the ball, that is, choose  $h$  small enough.

1825 We require  $h \leq \rho/2$ , and for simplicity, we just set  $h \leq \rho/8$ . Then,  $g^* \in B(X^*, \rho/2)$ , and since in  
 1826 the ball the function is smooth, with  $\nabla F(X^*) = 0$ , we have

$$F(g^*) \geq F(X^*) - \frac{L}{2} \|X^* - g^*\|^2 \geq F^* - \frac{L}{2} h^2.$$

1827 Hence, we require

$$F^* - \frac{L}{2} h^2 - \frac{\Gamma}{8} \geq F^* - \frac{7}{8} \Gamma,$$

1828 which is satisfied if  $h^2 \leq \frac{3}{2L} \Gamma$ . Hence, any point outside  $B(X^*, \rho/2)$  cannot upper bound  $\hat{F}(g^*)$ .

1829 Therefore,  $x_0 = \arg \max_{g \in \mathcal{G}} \hat{F}(g)$  satisfies  $x_0 \in B(X^*, \rho/2)$  and

$$\hat{F}(x_0) \geq \hat{F}(g^*) \geq F^* - \frac{L}{2} \left( \min\{\rho/8, \sqrt{3\Gamma/(2L)}\} \right)^2 - \frac{\Gamma}{8},$$

1830 Lastly, the number of points sampled depends on the number of points in  $\mathcal{G}$  is bounded by  $(1/h + 1)^d$ .

1831 Since  $h \leq 1$ , then  $|\mathcal{G}| \leq 2^d h^{-d}$ . □

1832 **C Appendix: Algorithms**

1833 This appendix describes the implementation of C-ICPE used in the experiments. We keep the notation  
 1834 of the main text:  $H_t$  is the current history,  $I_\phi(\cdot|H_t)$  is the inference distribution over the target  $x_\theta^*$ ,  
 1835  $Q_\psi(H_t, a)$  is the critic for a continuation action  $a \in \mathcal{A}$ , and  $Q_\psi(H_t, a_{\text{stop}})$  is the value of stopping.  
 1836 The implementation follows the Lagrangian view in Section B.3: the inference model learns a  
 1837 stochastic selector, the critic learns the stop/continue Bellman comparison, and the actor rule proposes  
 1838 the next continuation action. The main practical point is that all three objects are trained from replay.  
 1839 For this reason we use target networks, conservative critic targets, and simple regularizers that keep  
 1840 early noisy estimates from determining the stopping boundary.

1841 **C.1 History Encoder and Time Pooling Layer**

1842 The inference network, the critic, and the learned TD3 actors use the same sequential template,  
 1843 although their parameters are separate. Each interaction step is embedded as a token by concatenating  
 1844 the query and the next observation,  $u_s = [A_s; Y_{s+1}]$ , and passing it through a small embedding  
 1845 network. The resulting sequence  $(e_1, \dots, e_t)$  is processed by a causal sequential backbone. In  
 1846 the code this backbone can be an LSTM, an attention stack, or a recurrent/linear-attention variant.  
 1847 Padding is masked throughout, so replay batches can contain histories of different lengths without  
 1848 leaking future observations.

1849 The readout is not simply the last hidden state. After obtaining hidden states  $h_1, \dots, h_t \in \mathbb{R}^d$ , we  
 1850 use a query-conditioned pooling layer. Given a query  $q$ , it computes

$$\alpha_i(q, H_t) = \frac{\exp(\langle W_q q, W_k h_i \rangle / \sqrt{d})}{\sum_{j=1}^t \exp(\langle W_q q, W_k h_j \rangle / \sqrt{d})}, \quad v(q, H_t) = \sum_{i=1}^t \alpha_i(q, H_t) W_v h_i,$$

1851 and then applies a feature-wise gate,  $z(q, H_t) = v(q, H_t) \odot \text{silu}(W_m q)$ . This is useful because the  
 1852 relevant part of the history depends on what the model is asked to do. For inference and stopping, the  
 1853 query is a time embedding, since we need a representation of the current prefix. For the continuation  
 1854 critic, the query is the candidate action  $a$ , so the same history can be read differently when evaluating  
 1855 different future measurements. The critic therefore has two readouts from the same encoded history:  
 1856 a time-conditioned readout for  $Q_\psi(H_t, a_{\text{stop}})$  and an action-conditioned readout for  $Q_\psi(H_t, a)$ .

1857 **C.2 Replay Buffer and Prefix Sampling**

1858 Training is off-policy. A rollout samples a task  $\theta \sim \nu$ , interacts with the corresponding environment  
 1859 until the critic stops or the maximum horizon is reached, and stores the whole trajectory together  
 1860 with  $x_\theta^*$  and a flag indicating whether the trajectory ended by executing the stop action. Minibatches  
 1861 are built by first sampling trajectories and then sampling prefixes inside them. The prefix sampler  
 1862 uses a mixture of uniformly sampled prefixes, prefixes around the current average stopping time, and  
 1863 terminal prefixes. This gives the inference network examples from all time scales, but also gives the  
 1864 critic many prefixes close to the point where the decision changes from continue to stop.

1865 There is one convention that matters. If a trajectory terminates because the agent stops at time  $t$ , the  
 1866 training prefix for that terminal state is  $H_t$ , because no new observation is collected after the stop  
 1867 action. If the trajectory terminates only because the maximum horizon is reached, the last collected

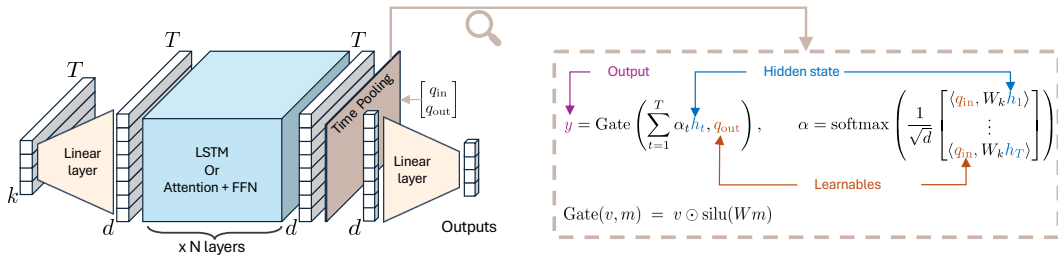


Figure 6: Sequential encoder and query-pooling readout used by the inference network  $I_\phi$ , the critic  $Q_\psi$ , and the learned TD3 actors.

1868 observation is included. This makes replay consistent with deployment: stopping is a decision made  
 1869 before paying for the next sample.

### 1870 C.3 Inference Update and Its Regularization

1871 For a prefix  $h$  and target  $x_\theta^*$ , the inference model outputs a diagonal Gaussian

$$I_\phi(\cdot|h) = \mathcal{N}(\mu_\phi(h), \text{diag}(\sigma_\phi^2(h))).$$

1872 The ideal objective in the main text is the negative log-likelihood in Eq. (5). In the implementation  
 1873 we use a robust version. Define the per-coordinate averaged NLL

$$\ell_\phi(h, x_\theta^*) = -\frac{1}{d_\mathcal{X}} \sum_{j=1}^{d_\mathcal{X}} \log \mathcal{N}((x_\theta^*)_j; \mu_{\phi,j}(h), \sigma_{\phi,j}^2(h)),$$

1874 and

$$\rho_\tau(u) = \min\{u, 0\} + \tau \log \left( 1 + \frac{[u]_+}{\tau} \right), \quad [u]_+ = \max\{u, 0\}.$$

1875 For Euclidean recommendation problems we train with

$$\mathcal{L}_{\text{inf}}(\phi) = \mathbb{E}_{(h, x_\theta^*) \sim \mathcal{B}} [\rho_\tau(\ell_\phi(h, x_\theta^*)) + \alpha_{\text{anc}} \text{SmoothL1}(\mu_\phi(h), x_\theta^*)]. \quad (22)$$

1876 The robust transform is a simple way of saying that early bad prefixes should not dominate the  
 1877 variance head. A standard Gaussian NLL can become very large when the model underestimates  
 1878 uncertainty for one minibatch, and this can push the log-variance to extremes. The logarithmic tail  
 1879 keeps the ranking of ordinary examples, but reduces the influence of rare outliers. The SmoothL1  
 1880 anchor is also important: the final decision is the mean  $\mu_\phi(H_\tau)$ , so we want the mean to remain a good  
 1881 deterministic recommendation even when the Gaussian still has large uncertainty. We additionally  
 1882 clamp the predicted log-standard deviations to fixed lower and upper bounds for numerical stability.

1883 For the  $\epsilon$ -best-arm problem, the target is directional. The magnitude of a vector is not meaningful  
 1884 once arms are represented on the sphere, and correctness is measured by cosine distance. Therefore  
 1885 recommendations and samples are normalized before they are evaluated. In this case we keep the  
 1886 robust NLL and replace the Euclidean anchor by a spherical alignment term. With  $Z \sim I_\phi(\cdot|h)$ ,

$$\mathcal{L}_{\text{sph}}(\phi) = -\mathbb{E} \left[ \left\langle \frac{Z}{\|Z\|_2}, \frac{x_\theta^*}{\|x_\theta^*\|_2} \right\rangle \right]. \quad (23)$$

1887 The reason is that an Euclidean anchor would penalize harmless radial errors, while the bandit loss  
 1888 only cares about the direction.

1889 A target copy  $I_{\bar{\phi}}$  is maintained by Polyak averaging. The critic never uses the online inference model  
 1890 inside its TD targets; it uses  $I_{\bar{\phi}}$ . This separation is important because the reward itself is learned  
 1891 through  $I_\phi$ , and bootstrapping from a rapidly moving reward makes the stop/continue comparison  
 1892 unstable.

### 1893 C.4 Reward, Critic Update, and Critic Regularization

1894 The critic learns two quantities from the same replay prefixes. The stopping head learns the value of  
 1895 recommending now, and the continuation head learns the value of paying for one more observation  
 1896 and then acting optimally. Since the ideal reward  $r_t(h) = \max_x q_t(h, x)$  is not available, we use the  
 1897 target inference model to estimate how likely the implemented stochastic selector is to be already  
 1898  $\epsilon$ -correct:

$$\hat{r}_{\bar{\phi}, m}(h, \theta) = \frac{1}{m} \sum_{k=1}^m \mathbf{1} \left\{ L_\theta \left( X^{(k)} \right) \leq \epsilon \right\}, \quad X^{(k)} \sim I_{\bar{\phi}}(\cdot|h). \quad (24)$$

1899 The reason for sampling is that the stopping decision should depend on posterior concentration, not  
 1900 only on the posterior mean. If  $\mu_\phi(h)$  is close to  $x_\theta^*$  but  $\sigma_\phi(h)$  is still large, stopping is risky. The  
 1901 sampled reward makes this visible to the critic, and this is exactly the gap controlled by Proposition 6.

1902 Let  $d$  be the terminal flag for a replay transition, and let  $a^+(h')$  be the target continuation action at the  
 1903 next prefix. This target action depends on the actor rule. For TS and TTPS it is the target inference

1904 mean plus small smoothing noise, projected to the feasible action set. For TD3 it is the target actor  
 1905 action, again with small target-policy smoothing noise. The smoothing noise is much smaller than  
 1906 the posterior sampling noise used during rollouts.

1907 When using twin critics, we scalarize target values by the conservative minimum

$$\bar{Q}_{\bar{\psi}}(h, a) = \min\{Q_{\bar{\psi},1}(h, a), Q_{\bar{\psi},2}(h, a)\},$$

1908 and define

$$V_{\bar{\psi}}(h') = \max\{\bar{Q}_{\bar{\psi}}(h', a_{\text{stop}}), \bar{Q}_{\bar{\psi}}(h', a^+(h'))\}.$$

1909 The TD targets are

$$y_{\text{stop}}(h, \theta) = \hat{r}_{\bar{\phi},m}(h, \theta), \quad (25)$$

$$y_{\text{cont}}(h, a, h', d, \theta) = -c(1-d) + d\hat{r}_{\bar{\phi},m}(h', \theta) + \gamma(1-d)V_{\bar{\psi}}(h'). \quad (26)$$

1910 Thus the stopping head is directly supervised by the current confidence, while the continuation head  
 1911 is supervised by the gain from collecting the next observation. The action-head loss is applied only to  
 1912 prefixes that did not already stop, because a stopped transition has no genuine continuation action.  
 1913 The stopping head is trained on every prefix, since stopping is a valid action at every prefix. With  
 1914  $M = 1$  for non-stopped replay transitions and  $M = 0$  for stopped transitions, the critic loss is

$$\mathcal{L}_Q(\psi) = \frac{1}{2}\mathbb{E}_{\mathcal{B}} \left[ M \sum_{j=1}^2 (Q_{\psi,j}(h, a) - y_{\text{cont}})^2 \right] + \frac{w_{\text{stop}}}{2}\mathbb{E}_{\mathcal{B}} \left[ \sum_{j=1}^2 (Q_{\psi,j}(h, a_{\text{stop}}) - y_{\text{stop}})^2 \right], \quad (27)$$

1915 with the obvious single-critic version when twin critics are disabled. We also optionally clip the  
 1916 bootstrap values in a minibatch to moderate quantiles before taking the maximum in  $V_{\bar{\psi}}$ . It prevents a  
 1917 few very optimistic target values from propagating through replay and moving the stopping boundary  
 1918 too early.

1919 At rollout time the default stopping test is

$$Q_{\psi}(H_t, a_{\text{stop}}) \geq Q_{\psi}(H_t, A_t),$$

1920 where both quantities are scalarized by the conservative twin rule. This is exactly the learned version  
 1921 of the Bellman comparison in Theorem 3.1.

## 1922 C.5 Actor Rules: TS, TTPS, and TD3

1923 The actor rule only chooses continuation actions. The stop action is handled by the critic comparison  
 1924 above. We use three actor rules depending on the relation between the query space  $\mathcal{A}$  and the  
 1925 recommendation space  $\mathcal{X}$ .

1926 **Thompson sampling.** When  $\mathcal{A} = \mathcal{X}$ , we can use the inference distribution itself as the actor. The TS  
 1927 rule samples

$$A_t = \mu_{\phi}(H_t) + \sigma_{\phi}(H_t) \odot \xi_t, \quad \xi_t \sim \mathcal{N}(0, I). \quad (28)$$

1928 At evaluation in greedy mode we set  $\xi_t = 0$ , so the action is the projected posterior mean. This choice  
 1929 is deliberately simple. Early in a task, the learned posterior is broad and TS explores. Later, when the  
 1930 inference distribution contracts, the same rule automatically becomes exploitative. No separate actor  
 1931 is trained, which removes one source of approximation error and is appropriate when informative  
 1932 queries are themselves plausible recommendations. In the critic target, the next action for TS is the  
 1933 target inference mean, with only the small target-smoothing noise described above.

1934 **Top-two posterior sampling.** TTPS is also inference-based and therefore also assumes  $\mathcal{A} = \mathcal{X}$ .  
 1935 It first draws a Thompson sample  $Z_1$ . With probability  $1/2$  this sample is used directly. With the  
 1936 remaining probability, the rule draws a challenger  $Z_2$  and keeps the sample that is farther from  
 1937 the posterior mean, provided the current sample is too close to the mean. In Euclidean tasks the  
 1938 distance is  $\|z - \mu_{\phi}(H_t)\|_2$ ; in  $\epsilon$ -best-arm it is  $1 - \langle z, \mu_{\phi}(H_t) \rangle$ . The intuition is that the posterior  
 1939 mean is the current recommendation, while farther posterior samples represent plausible alternatives.  
 1940 TTPS spends part of the sampling budget checking those alternatives before the critic decides that  
 1941 stopping is safe. As for TS, the target action used by the critic is the target inference mean plus small  
 1942 smoothing noise.

---

**Algorithm 3** Implementation of C-ICPE
 

---

```

1: Initialize replay buffer  $\mathcal{B}$ , inference network  $I_\phi$ , critic  $Q_\psi$ , actor rule  $\text{Act} \in \{\text{TS}, \text{TTPS}, \text{TD3}\}$ , target networks
    $I_{\bar{\phi}}, Q_{\bar{\psi}}$ , and cost  $c$ .
// Training phase
2: while training is not over do
3:   Sample a batch of tasks  $\theta \sim \nu$  and initialize their histories  $H_1$ .
4:   for  $t = 1, \dots, T_{\max}$  do
5:     Propose continuation actions  $A_t$  using TS, TTPS, or the TD3 actor with its current exploration schedule.
6:     Stop tasks satisfying  $Q_\psi(H_t, a_{\text{stop}}) \geq Q_\psi(H_t, A_t)$  after the warmup and minimum-time gates.
7:     Execute  $A_t$  on the remaining tasks, observe  $Y_{t+1}$ , and append  $(A_t, Y_{t+1})$  to the histories.
8:   end for
9:   Store the completed trajectories, stop flags, terminal times, and targets  $x_\theta^*$  in  $\mathcal{B}$ .
10:  Sample replay prefixes and update  $I_\phi$  with Eq. (22); for  $\epsilon$ -best-arm also use Eq. (23).
11:  Estimate sampled rewards with  $I_{\bar{\phi}}$  using Eq. (24) and update  $Q_\psi$  using Eqs. (25) to (27).
12:  If using TD3, update  $\pi_\vartheta$  on the delayed actor schedule using Eqs. (29) and (30).
13:  Polyak-update  $I_{\bar{\phi}}, Q_{\bar{\psi}}$ , and, when present,  $\pi_{\bar{\vartheta}}$ ; update  $c$  with Eq. (31).
14: end while

// Deployment phase
15: Freeze the learned networks and initialize a fresh test task.
16: for  $t = 1, \dots, T_{\max}$  do
17:   Propose  $A_t$  using the selected actor rule.
18:   if  $Q_\psi(H_t, a_{\text{stop}}) \geq Q_\psi(H_t, A_t)$  then
19:     return  $\hat{x} = \mu_\phi(H_t)$ .
20:   end if
21:   Execute  $A_t$ , observe  $Y_{t+1}$ , and update  $H_{t+1}$ .
22: end for
23: return  $\hat{x} = \mu_\phi(H_{T_{\max}+1})$ .

```

---

1943 **TD3 actor.** TD3 is the actor rule used when we want to learn the continuation action from the critic,  
 1944 and especially when  $\mathcal{A}$  and  $\mathcal{X}$  are different objects. A target actor  $\pi_{\bar{\vartheta}}$  is used in the critic target, while  
 1945 the online actor  $\pi_\vartheta$  is updated on a delayed schedule. The target action is

$$a^+(h') = \tanh(\pi_{\bar{\vartheta}}(h')) + \zeta, \quad \zeta \sim \mathcal{N}(0, \sigma_{\text{tgt}}^2 I),$$

1946 with clipping or normalization depending on the domain. The target-policy smoothing noise  $\zeta$   
 1947 prevents the critic from learning sharp artificial peaks in action space, and the delayed actor update  
 1948 prevents the actor from chasing a critic that is still changing after every minibatch.

1949 For a deterministic actor, the main update is

$$\mathcal{L}_{\text{act}}^{\text{det}}(\vartheta) = -\mathbb{E}_{h \sim \mathcal{B}} [Q_{\psi,1}(h, \tanh(\pi_\vartheta(h)))]. \quad (29)$$

1950 The actor uses the first critic for the policy gradient, while the target uses the first critic.

1951 When  $\mathcal{A} = \mathcal{X}$  and the actor is Gaussian, we regularize the TD3 actor toward the inference distribution,

$$\mathcal{L}_{\text{act}}^{\text{KL}}(\vartheta) = -\mathbb{E}_{h \sim \mathcal{B}} [Q_{\psi,1}(h, A_\vartheta(h))] + \beta_{\text{KL}} \mathbb{E}_{h \sim \mathcal{B}} [\text{KL}(I_\phi(\cdot|h) \parallel \pi_\vartheta(\cdot|h))], \quad A_\vartheta(h) \sim \pi_\vartheta(\cdot|h). \quad (30)$$

1952

1953 During training we also perform random exploration to encourage parametric exploration (with small  
 1954 probability we sample a random action), and void collapsing of the actor.

### 1955 C.6 Cost Update for Fixed Confidence

1956 The scalar cost  $c$  is the implemented Lagrange tradeoff between confidence and sample complexity.  
 1957 We update it from the empirical stopped success rate. For a batch of completed rollouts, let

$$D_i = \begin{cases} \|\mu_\phi(H_{\tau_i}^{(i)}) - x_{\theta_i}^*\|_2, & \text{Euclidean tasks,} \\ 1 - \langle \mu_\phi(H_{\tau_i}^{(i)}), x_{\theta_i}^* \rangle, & \epsilon\text{-best-arm tasks,} \end{cases}$$

1958 and use the smooth accuracy proxy

$$\hat{p} = \frac{1}{B} \sum_{i=1}^B \sigma\left(\frac{\epsilon - D_i}{\kappa\epsilon}\right).$$

1959 The update is

$$c \leftarrow \text{Proj}_{[0,1]} (c - \eta_c ((1 - \delta) - \hat{p})). \quad (31)$$

1960 Thus, if the observed correctness is below  $1 - \delta$ , the cost decreases and continuing becomes cheaper,  
1961 so trajectories become longer. If correctness is above the target, the cost increases and stopping  
1962 becomes more aggressive. The sign of the update is still exactly the dual intuition in Eq. (9).

## 1963 D Appendix: Numerical Results

1964 In this section we present more details on the numerical results. We refer the reader to the code for  
1965 more details (see the README.md file), especially regarding the hyperparameters. We now present the  
1966 synthetic benchmarks with additional numerical results. These additional results include sweeps over  
1967 various values of  $\epsilon, \sigma$ , as well as checking robustness to prior misspecification. We conclude with  
1968 details and additional results regarding the geochemical exploration task.

1969 **Computational resources.** All experiments were run on NVIDIA V100 GPU or NVIDIA L40S  
1970 GPU. For the synthetic benchmarks, each C-ICPE training run takes approximately 12 hours. With 3  
1971 random seeds, at least 4  $(\epsilon, \sigma)$  configurations, and 3 dimensionalities per benchmark, the total training  
1972 budget per synthetic task is  $12 \times 3 \times 4 \times 3 = 432$  GPU-hours. For the geochemical exploration  
1973 task, each training run takes approximately 70 hours; with 3 seeds, 2  $\epsilon$  configurations, and a single  
1974 dimensionality ( $d = 2$ ), the total is  $70 \times 3 \times 2 = 420$  GPU-hours.

1975 **Confidence intervals via hierarchical bootstrap.** To account for variability across both task  
1976 instances and trajectory randomness, we report 95% confidence intervals computed via hierarchical  
1977 bootstrap [15]. For each trained model (seed), we sample 300 test environments  $\theta \sim \nu$  and collect  
1978 15 independent trajectories per environment. The total variance of a statistic  $\hat{\mu}$  (e.g., accuracy)  
1979 decomposes as

$$\text{Var}(\hat{\mu}) = \underbrace{\text{Var}_s(\mathbb{E}[\hat{\mu} | s])}_{\text{between-seed}} + \underbrace{\mathbb{E}_s[\text{Var}_\theta(\mathbb{E}[\hat{\mu} | s, \theta])]}_{\text{between-environment}} + \underbrace{\mathbb{E}_{s,\theta}[\text{Var}(\hat{\mu} | s, \theta)]}_{\text{within-environment}},$$

1980 where the first term captures variability due to training, the second due to which tasks are drawn from  
1981  $\nu$ , and the third due to observation noise and policy stochasticity within a fixed task. A single bootstrap  
1982 replicate is constructed by (i) resampling seeds with replacement, then (ii) for each resampled seed,  
1983 resampling environments with replacement, then (iii) for each resampled environment, resampling  
1984 trajectories with replacement, and computing the statistic on the resampled dataset. This three-level  
1985 resampling preserves all components of variance. We draw 10,000 bootstrap replicates and report the  
1986 2.5% and 97.5% percentiles as the confidence interval.

1987 **Remark 12** (On testing the inference model). *Several of our synthetic localization benchmarks fall*  
1988 *close to the symmetric case described in Section B.7 (after Proposition 8): the loss is a distance*  
1989 *to a selected target  $x_\theta^*$ , so  $\mathcal{X}_\epsilon(\theta)$  is a ball, interval, or cap around this target. In the optimization*  
1990 *benchmarks, such as the value estimation task and the geochemical task, the loss is instead induced*  
1991 *by value gaps and the success sets need not be symmetric or convex. These experiments therefore*  
1992 *test the method beyond the setting where the Gaussian NLL mean has an exact Bayes-optimality*  
1993 *interpretation.*

### 1994 D.1 Synthetic Benchmarks: description

1995 The synthetic benchmarks in Section 5 are designed to isolate different aspects of the continuous  
1996 fixed-confidence problem. Binary search is the cleanest localization problem: every query returns  
1997 a noisy comparison with the unknown target. The  $\epsilon$ -best-arm problem keeps the same idea of  
1998 identifying  $x_\theta^*$ , but changes the geometry to a sphere and makes the loss directional rather than  
1999 Euclidean. Ackley minimization adds nuisance parameters and a multimodal response surface, so the  
2000 agent has to learn an exploration rule that is not purely local. Finally, GP max-value estimation is  
2001 included because it is the case where the query space and the recommendation space are genuinely  
2002 different: the agent queries a location, but it recommends a scalar value. This is the setting where  
2003 a TD3 actor is necessary, since posterior samples from the inference network are no longer valid  
2004 actions.

2005 **Common protocol.** Each training episode starts by sampling a fresh task parameter  $\theta$  from the task  
2006 prior  $\nu$ . The agent then observes a sequential history  $H_t = (A_1, Y_2, \dots, A_{t-1}, Y_t)$  and either stops  
2007 or selects a new action. We use a maximum horizon  $t_{\max}$ ; if the learned stopping rule does not stop  
2008 before this horizon, the episode is truncated and the final recommendation is still evaluated. For the  
2009 synthetic experiments reported in the survival plots and correctness tables, we use  $t_{\max} = 100$  unless  
2010 otherwise stated.

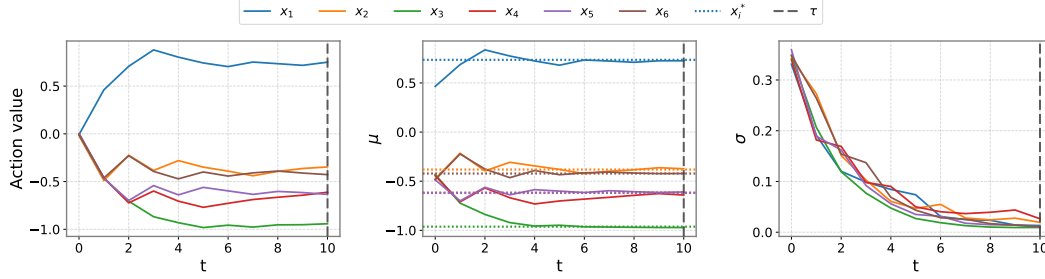


Figure 7: Binary search: visualization of how C-ICPE explores in 6 dimensions. From left to right: the query action (left), posterior mean (middle), and posterior standard deviation (right) along an exploration trajectory in the noisy binary search problem.

2011 All fixed-confidence runs use  $\delta = 0.1$ . The reported accuracy is

$$\widehat{\text{Acc}} = \frac{1}{n} \sum_{i=1}^n \mathbf{1} \{L_{\theta_i}(\hat{x}_i) \leq \epsilon\},$$

2012 and the goal is to achieve accuracy at least  $1 - \delta$  while minimizing the expected stopping time  $\mathbb{E}[\tau]$ .  
 2013 For Euclidean localization tasks we use  $L_{\theta}(x) = \|x - x_{\theta}^*\|_2$ . For  $\epsilon$ -best-arm we use the directional  
 2014 loss  $L_{\theta}(x) = 1 - \langle x, x_{\theta}^* \rangle$  after normalizing both vectors to the unit sphere. For GP max-value  
 2015 estimation the loss is the scalar absolute error  $L_{\theta}(x) = |x - v_{\theta}^*|$ . Confidence intervals are computed  
 2016 with hierarchical bootstrap over test episodes and random seeds.

### 2017 D.1.1 Noisy binary search.

2018 In binary search the unknown target is  $x_{\theta}^* = \theta \in [-1, 1]^d$ , sampled uniformly. A query  $a \in [-1, 1]^d$   
 2019 returns one noisy comparison per coordinate,

$$Y_{t,i} = \xi_{t,i} \text{sign}(\theta_i - A_{t,i}), \quad \mathbb{P}(\xi_{t,i} = 1) = 1 - p, \quad \mathbb{P}(\xi_{t,i} = -1) = p, \quad (32)$$

2020 independently over  $i$  and  $t$ . We use this problem because the statistically useful action is interpretable:  
 2021 a good policy should place queries near the current posterior median in each coordinate and shrink  
 2022 the feasible region. This makes binary search a sanity check for the inference network and critic. If  
 2023 the inference model does not contract its posterior, or if the critic cannot recognize when the posterior  
 2024 radius is below  $\epsilon$ , the method will fail even in this simple setting. Since  $\mathcal{A} = \mathcal{X} = [-1, 1]^d$ , TS and  
 2025 TTPS can act directly by sampling from the learned posterior over the target.

### 2026 D.1.2 $\epsilon$ -best-arm identification on the sphere.

2027 For the continuous  $\epsilon$ -best-arm problem, the task parameter is a direction  $x_{\theta}^* = \theta \in \mathbb{S}^{d-1}$  sampled by  
 2028 normalizing a standard Gaussian vector. The agent queries a vector  $a$  and observes

$$Y_t = \theta^{\top} A_t + \xi_t, \quad \xi_t \sim \mathcal{N}(0, \sigma^2). \quad (33)$$

2029 The recommendation is correct if  $\theta^{\top} \hat{x} \geq 1 - \epsilon$ . Although the implementation stores the enclosing  
 2030 action bounds as  $[-1, 1]^d$ , the inference mean, posterior samples, TTPS candidates, and uniform  
 2031 baseline actions are projected to  $\mathbb{S}^{d-1}$  for this benchmark. This projection is important: Euclidean  
 2032 uncertainty is not the right object near the sphere, and two vectors with the same direction but different  
 2033 norms should not be treated as different hypotheses. The reason for this benchmark is that each  
 2034 observation is a scalar projection. The agent must choose directions that disambiguate the posterior  
 2035 over  $\theta$ , while the stopping rule must reason in terms of cosine error rather than Euclidean error.

2036 This experiment also lets us compare to a specialized frequentist fixed-confidence baseline. Lazy  
 2037 Track-and-Stop uses the known linear observation structure and an analytic generalized-likelihood-  
 2038 ratio stopping rule. In our implementation it queries canonical directions and keeps a least-squares  
 2039 estimate of  $\theta$ , stopping only after both the likelihood-ratio condition and the spectral coverage  
 2040 condition are satisfied. This is not a general baseline for all our tasks, but it is a useful reference point  
 2041 on the one benchmark where a specialized fixed-confidence method is available. However, note that  
 2042 in this problem the optimal exploration strategy is uniform [30]. Therefore, it shows to what degree  
 2043 C-ICPE is able to learn a good inference model.

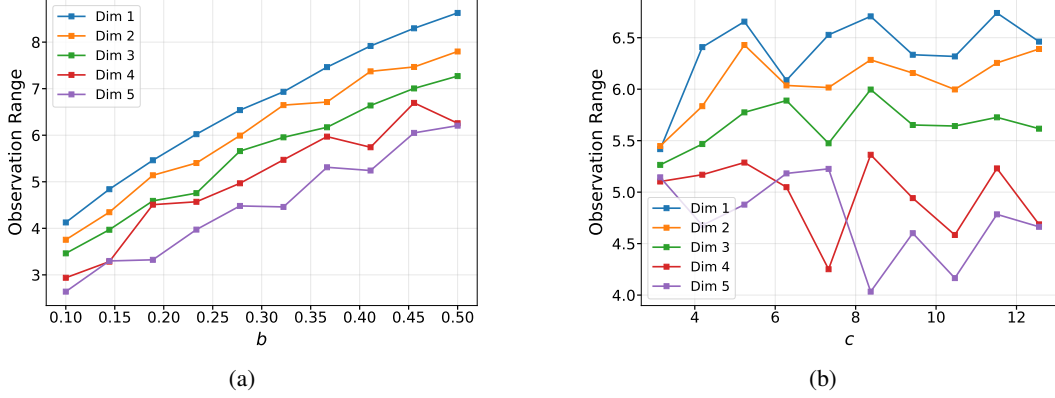


Figure 8: Effect of Ackley function’s parameters on output range across multiple dimensions: (a) range vs  $b$ ; (b) range vs  $c$ .

2044 **D.1.3 Ackley minimizer identification.**

2045 The Ackley task is a shifted and randomly parametrized global-optimization problem. The task  
 2046 parameter is

$$\theta = (a, b, c, \theta^*), \quad \theta^* \sim \text{Unif}([-1, 1]^d),$$

2047 where  $\theta^*$  is the global minimizer and  $(a, b, c)$  control the shape of the response surface. In the  
 2048 reported runs we fix  $a = 10$  and sample  $b \sim \text{Unif}[0.1, 0.5]$  and  $c \sim \text{Unif}[\pi, 4\pi]$ . Given  $u = A_t - \theta^*$ ,  
 2049 the unnormalized Ackley value is

$$F_{a,b,c}(u) = a + e - a \exp\left(-b \sqrt{\frac{1}{d} \sum_{j=1}^d u_j^2}\right) - \exp\left(\frac{1}{d} \sum_{j=1}^d \cos(cu_j)\right). \quad (34)$$

2050 We observe the sign-inverted and normalized value (more on this in the next page)

$$Y_t = 1 - 2 \frac{F_{a,b,c}(A_t - \theta^*)}{Z_{\text{norm}}(b, c, d)} + \xi_t, \quad \xi_t \sim \mathcal{N}(0, \sigma^2), \quad (35)$$

2051 so that larger observations are better and the target remains the minimizer  $\theta^*$ . The nuisance parameters  
 2052 are not provided to the agent. Thus, across episodes, C-ICPE must infer not only where the optimum  
 2053 is but also how observations should be interpreted for that episode.

2054 Ackley is included because it is deliberately hostile to naive local search. The function has many  
 2055 oscillations near the optimum and a broad outer region where observations can be weakly informative.  
 2056 The active policy therefore has to balance broad exploration with local refinement, and the critic has  
 2057 to stop based on whether the inferred minimizer is accurate, not based on whether the last observed  
 2058 function value was high.

2059 **Ackley function output normalization.** The Ackley function’s global minimum is always at the  
 2060 origin with a value of 0, but the maximum value within our defined recommendation space  $\mathcal{X}$  depends  
 2061 on the function parameters and dimensionality. Figs. 8a and 8b show how the values of  $b$  and  $c$  affect  
 2062 the function output ranges. Larger  $b$  values consistently increase the range, while  $c$  has less significant  
 2063 effect on the output ranges. From the figures, we also see that the ranges depend on dimensionality,  
 2064 where lower dimensions tend to have larger ranges. The issue with varying output ranges is that  
 2065 the influence of noise can vary across different priors sampled, and we want the noise effect to be  
 2066 on the same scale. Additionally, without normalization, C-ICPE must not only learn the relative  
 2067 patterns from  $H_t$  but also account for the scale differences across different  $H_t$ . For these reasons, we  
 2068 derive a normalization constant empirically from multiple samples across different  $b, c$  values and  
 2069 dimensionalities:  $Z_{\text{norm}}(b, c, d) = \pi - 0.21 \cdot D + 9.68 \cdot b + 0.04 \cdot c$ .

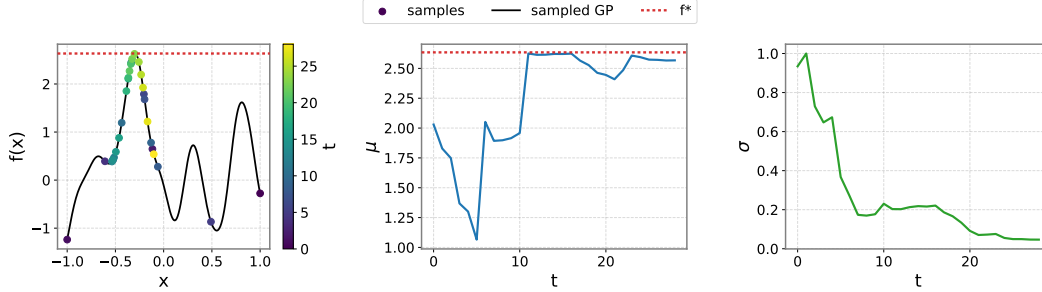


Figure 9: Visualization of how C-ICPE explores in the max-value estimation problem. From left to right: the query action (left), posterior mean (middle), and posterior standard deviation (right) along an exploration trajectory in the GP value estimation problem. Darker samples are queried earlier in the trajectory, while lighter samples are queried later.

#### 2070 D.1.4 GP max-value estimation

2071 The GP benchmark separates the action and recommendation spaces. At the beginning of an episode  
 2072 we sample a latent function

$$f \sim \text{GP}(0, k_{\text{RBF}}(\ell, \sigma_f)), \quad \ell \sim \text{Unif}[0.05, 0.2], \quad \sigma_f = 1,$$

2073 on  $[0, 1]^d$ , with  $d = 1$  in the implementation. The sample path is generated on a dense grid using  
 2074 circulant embedding. Queries are continuous points  $A_t \in [0, 1]^d$ ; the observed value is obtained  
 2075 by linear interpolation in one dimension or bilinear interpolation in two dimensions, followed by  
 2076 Gaussian noise:

$$Y_t = f(A_t) + \xi_t, \quad \xi_t \sim \mathcal{N}(0, \sigma^2). \quad (36)$$

2077 The target is the scalar maximum value

$$v_\theta^* = \max_{u \in [0, 1]^d} f(u), \quad \mathcal{X} \subseteq \mathbb{R}, \quad \mathcal{A} = [0, 1]^d,$$

2078 where the maximum is computed on the same grid used to generate the episode. The final recommen-  
 2079 dation is the inference mean for this scalar value, and success is  $|\hat{x} - v_\theta^*| \leq \epsilon$ .

2080 The reason for using max-value estimation rather than another argmax-localization task is that it tests  
 2081 the decoupling of inference and exploration. In binary search,  $\epsilon$ -best-arm, and Ackley, a posterior  
 2082 sample of  $x_\theta^*$  is itself a reasonable query, so TS and TTPS can be implemented directly from the  
 2083 inference distribution. For GP value estimation this would be meaningless: a sample from the  
 2084 inference model is a scalar value, not a point in  $[0, 1]^d$ . Therefore the action must be learned through  
 2085 the critic. We use the TD3 actor for this benchmark because the critic can assign value to a query  
 2086 according to how much it is expected to improve the future estimate of  $v_\theta^*$ , even though the query is  
 2087 not itself a recommendation. This experiment is consequently the main empirical check that C-ICPE  
 2088 handles the general  $\mathcal{A} \neq \mathcal{X}$ .

#### 2089 D.2 Synthetic Benchmarks: baselines

2090 The most important baseline is C-ICPE-Uniform, which keeps the same inference network, critic,  
 2091 stopping rule, replay buffer, and fixed-confidence cost update as C-ICPE, but replaces the learned  
 2092 active query rule by uniform exploration. This isolates the value of active experimentation: if C-ICPE  
 2093 improves over C-ICPE-Uniform, the gain cannot be explained by the inference model alone, because  
 2094 both methods use the same form of inference and the same stopping mechanism.

2095 For tasks with  $\mathcal{A} = \mathcal{X}$ , we evaluate TS and TTPS because they use the learned posterior in the most  
 2096 direct way. TS samples a plausible target and queries it. TTPS keeps the posterior mean as the current  
 2097 recommendation and intentionally samples a plausible challenger that is sufficiently different. This is  
 2098 useful when many posterior samples are small perturbations around the current mean: such samples  
 2099 do not test the remaining uncertainty, whereas a challenger query can reveal whether another region  
 2100 is still plausible. For GP value estimation, TS and TTPS are not the right action rules for the reason  
 2101 described above, so we use TD3.

2102 For all benchmarks, we set the sample budget to  $t_{\max} = 100$  by default. If the trained C-ICPE policy  
 2103 failed to reach the target  $(1 - \delta)$ -accuracy, we extended the sample budget by 50. For easier setting,  
 2104 such as low dimension, we instead used a smaller budget. The settings in which  $t_{\max} \neq 100$  are  
 2105 listed in Table 1, all other configurations use  $t_{\max} = 100$ .

Environment	$d$	$\varepsilon$	$t_{\max}$
$\varepsilon$ -Best-Arm	15	0.005	150
GP-value estimation	1	0.2	60
Geochemical	2	0.15	150

Table 1: Sample budget  $t_{\max}$  for setting where  $t_{\max} \neq 100$

2106 We also compare against standard fixed-budget optimization methods implemented through Optuna  
 2107 [2]: TPE, GP-based Bayesian optimization, and CMA-ES. We also compare with GP-UCB via  
 2108 BoTorch [7]. These baselines do not have a learned stopping rule and are not optimized for  $(\varepsilon, \delta)$ -  
 2109 correctness. To make the comparison conservative, we give them budgets tied to the empirical  
 2110 stopping time of the corresponding C-ICPE variant.

2111 They are then evaluated under the same success criterion  $L_{\theta}(\hat{x}) \leq \varepsilon$ . This comparison asks whether  
 2112 a generic optimizer, with a fixed budget equal to C-ICPE’s expected sample complexity, already  
 2113 reaches the fixed-confidence target.

- 2114 • **TPE** [8]: splits observations into good and bad groups based on a quantile threshold and  
 2115 models the objective by building a density estimator for each group, then selects candidates  
 2116 that maximize the ratio of good-to-bad.
- 2117 • **CMA-ES** [24]: an evolutionary algorithm that samples a population of candidates and  
 2118 iteratively updates a multivariate Gaussian distribution by adapting its covariance matrix  
 2119 based on the successful candidates.
- 2120 • **GP-logEI** [5]: updates the Matérn kernel’s hyperparameters by maximizing the marginal log-  
 2121 likelihood on the past observations, and uses log expected improvement as the acquisition  
 2122 function.
- 2123 • **GP-UCB** [57]: a variant of GP-based Bayesian Optimization with a Matérn kernel that uses  
 2124 upper confidence bound as the acquisition function.
- 2125 • **Uniform bin**: partitions the query space into  $\lceil \sqrt{\mathbb{E}[\tau]} \rceil$  bins, where  $\mathbb{E}[\tau]$  is the corresponding  
 2126 C-ICPE variant’s expected sample complexity, and uniformly sample within each bin. We  
 2127 compute the average value of each bin and report the maximum.
- 2128 • **Uniform top 5%**: queries uniformly and return the mean of the top 5% values. The number  
 2129 of query matches the corresponding C-ICPE variant’s expected sample complexity.
- 2130 • **Lazy Track-and-Stop round robin** [30]: queries the canonical basis in a round-robin  
 2131 fashion. The method maintains a least square estimate  $\hat{\theta}_t$ , and stops whenever

$$Z_t \geq \beta(\delta, t) \quad \text{and} \quad \min_j N_t(j) \geq \max\left(c, \frac{\rho(\delta, t)}{\|\hat{\theta}_t\|^2}\right),$$

2132 where  $Z_t$  is the generalized likelihood ratio (GLR) for the  $\varepsilon_t$  best-arm hypothesis evaluated  
 2133 against the worst-case competitor on the  $\varepsilon_t$ -boundary. This measures how much the current  
 2134 guess is better than the closest alternative.  $N_t(j)$  is the number of pulls of the  $j$ -th element  
 2135 of the canonical basis, and  $\beta(\delta, t)$ ,  $\rho(\delta, t)$ ,  $\varepsilon_t$  are the threshold rule, spectral threshold, and  
 2136 gap-relaxation threshold. Lastly  $c$  is a constant defined in [30]. See also [30] for more  
 2137 details. Upon stopping, the agent recommends  $\hat{a}_t = \hat{\theta}_t / \|\hat{\theta}_t\|$ .

- 2138 • **Lazy Track-and-Stop uniform** [30]: Same as Lazy Track-and-Stop round robin except it  
 2139 samples the basis uniformly.

2140 **D.3 Synthetic Benchmarks: numerical results**

2141 We now present detailed accuracy and sample complexity results across all benchmarks, sweeping over  
 2142  $(\varepsilon, \sigma, d)$  configurations. Tables report mean accuracy and sample complexity with 95% confidence  
 2143 intervals for every method and parameter combination. To complement these aggregate statistics,  
 2144 we examine the stopping behavior of each actor through two diagnostics: (i) the survival function  
 2145  $P(\tau > t)$ , which reveals how quickly the learned critic commits to stopping, and (ii) the standard  
 2146 deviation of the inference model over the horizon, which tracks how rapidly the posterior uncertainty  
 2147 around  $x_\theta^*$  contracts. Across all benchmarks, C-ICPE with learned exploration (TS, TTPS, or TD3)  
 2148 consistently meets the  $1 - \delta$  accuracy target while C-ICPE-Uniform degrades as dimension increases,  
 2149 particularly on Ackley and binary search where directed exploration is essential.

2150 **D.3.1 Noisy Binary Search**

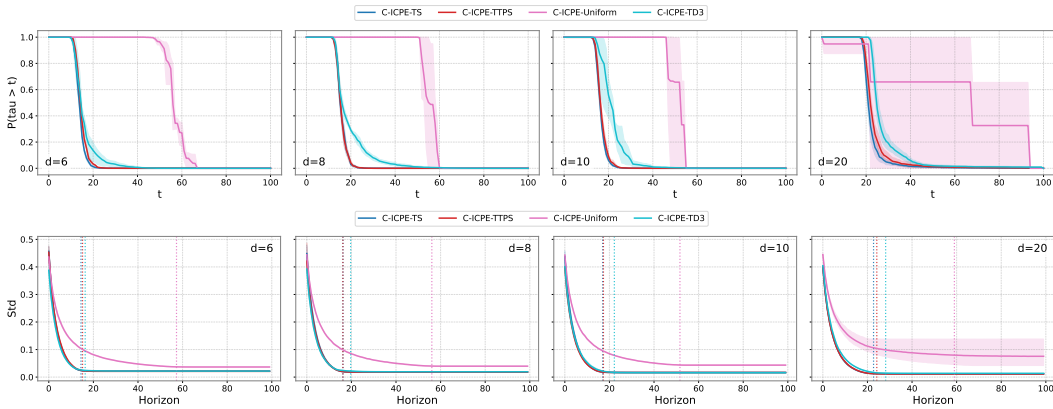


Figure 10: Results for Binary Search Problem with fixed confidence  $\delta = 0.1$  and  $N = 100$  across different dimensions at the most hardest  $(\varepsilon, \sigma)$  setting: (top) survival function of  $\tau$ ; (bottom) inference uncertainty convergence.

2151 Tables 2 and 3 report accuracy and sample complexity across all  $(d, \varepsilon, \sigma)$  configurations; Figure 10  
 2152 shows the survival function of  $\tau$  and the convergence of the inference model’s standard deviation. In  
 2153 Fig. 7 we also show how C-ICPE explores in 5 dimensions, depicting the queries chosen by the actor,  
 2154 the posterior mean, and the posterior standard deviation over timesteps.

2155 *Accuracy.* All three active actors (TS, TTPS, TD3) meet the  $1 - \delta = 0.90$  target across every  
 2156 configuration tested, with mean accuracy between 0.895 and 0.916. The confidence intervals confirm  
 2157 that the target is met reliably: even the lower bounds remain at or above 0.886. C-ICPE-Uniform  
 2158 matches the active methods at  $d = 6$  with  $\varepsilon = 0.2$  (accuracy  $\geq 0.901$ ) but degrades sharply as either  
 2159  $d$  increases or  $\varepsilon$  decreases. At  $(\varepsilon, \sigma) = (0.1, 0.05)$ , accuracy drops from 0.688 at  $d = 6$  to 0.440  
 2160 at  $d = 8$ , 0.171 at  $d = 10$ , and 0.006 at  $d = 20$ . This confirms that passive exploration cannot  
 2161 accumulate sufficient directional information per coordinate to localize the target within the allowed  
 2162 horizon in high dimensions.

2163 *Sample complexity.* Among active methods, C-ICPE-TS and C-ICPE-TTPS achieve comparable  
 2164 sample complexity across all settings: at  $d = 20$ ,  $(\varepsilon, \sigma) = (0.1, 0.05)$ , C-ICPE-TS stops in 22.7  
 2165 queries on average and C-ICPE-TTPS in 24.2, both with tight confidence intervals. C-ICPE-TD3  
 2166 is competitive at moderate dimensions ( $d \leq 10$ ) but exhibits higher mean stopping times and  
 2167 substantially wider confidence intervals at  $d = 20$  (e.g., 35.8 [19.2, 67.6] at  $\varepsilon = 0.2, \sigma = 0.05$ ),  
 2168 suggesting that the learned actor is less stable in high dimensions. Sample complexity scales  
 2169 sublinearly in  $d$  for the active actors: C-ICPE-TS increases from 11.1 ( $d = 6$ ) to 16.4 ( $d = 20$ ) at  
 2170  $(\varepsilon, \sigma) = (0.2, 0.05)$ .

2171 *Stopping behavior.* The survival functions (Figure 10, top) corroborate the sample complexity  
 2172 results. At  $d \leq 10$ , C-ICPE-TS and C-ICPE-TTPS exhibit sharp transitions:  $P(\tau > t)$  drops  
 2173 from 1 to 0 within a narrow window, indicating that the critic identifies a consistent stopping point.  
 2174 C-ICPE-Uniform has a heavy-tailed survival function that extends to the horizon, and at  $d = 20$  it  
 2175 rarely stops before  $t_{\max}$ . The inference standard deviation (bottom row) confirms that the active actors’

2176 posteriors contract rapidly, reaching near-zero uncertainty before the median stopping time (dashed  
2177 vertical lines), while C-ICPE-Uniform at  $d = 20$  retains high residual uncertainty throughout the  
2178 episode.

$d$	Method	$\varepsilon = 0.2$		$\varepsilon = 0.1$	
		$\sigma = 0.05$	$\sigma = 0.1$	$\sigma = 0.025$	$\sigma = 0.05$
6	C-ICPE-TD3	0.895 [886,902]	0.898 [886,908]	0.903 [893,914]	<b>0.901</b> [889,912]
	C-ICPE-TS	<b>0.904</b> [895,912]	<b>0.907</b> [899,913]	<b>0.909</b> [899,920]	0.900 [893,909]
	C-ICPE-TTPS	0.897 [890,905]	0.903 [893,913]	0.900 [893,907]	0.898 [889,905]
	C-ICPE-uniform	0.901 [894,908]	0.904 [893,914]	0.854 [835,874]	0.688 [666,709]
8	C-ICPE-TD3	0.903 [897,911]	0.898 [890,907]	0.901 [891,912]	0.899 [890,906]
	C-ICPE-TS	0.905 [892,917]	0.898 [889,905]	<b>0.916</b> [902,930]	<b>0.909</b> [899,918]
	C-ICPE-TTPS	0.903 [896,911]	<b>0.905</b> [893,915]	0.907 [895,917]	0.901 [891,909]
	C-ICPE-uniform	<b>0.907</b> [897,918]	0.894 [886,901]	0.724 [693,756]	0.440 [406,477]
10	C-ICPE-TD3	0.901 [892,908]	0.898 [889,904]	0.904 [893,914]	<b>0.915</b> [891,948]
	C-ICPE-TS	0.896 [886,903]	<b>0.903</b> [890,916]	<b>0.916</b> [905,925]	0.904 [892,915]
	C-ICPE-TTPS	<b>0.902</b> [895,910]	0.900 [891,906]	0.901 [891,911]	0.904 [897,912]
	C-ICPE-uniform	0.873 [853,891]	0.860 [847,873]	0.493 [471,514]	0.171 [136,201]
20	C-ICPE-TD3	0.897 [888,905]	0.905 [896,913]	<b>0.925</b> [908,940]	0.903 [895,912]
	C-ICPE-TS	<b>0.911</b> [893,931]	0.903 [892,914]	0.901 [891,910]	0.903 [889,919]
	C-ICPE-TTPS	0.896 [887,903]	<b>0.907</b> [897,917]	0.900 [890,909]	<b>0.913</b> [906,921]
	C-ICPE-uniform	0.662 [630,698]	0.079 [001,136]	0.036 [022,049]	0.006 [000,013]

Table 2: Binary search: accuracy (mean and 95% CI) for every  $(d, \varepsilon, \sigma)$  configuration.

$d$	Method	$\varepsilon = 0.2$		$\varepsilon = 0.1$	
		$\sigma = 0.05$	$\sigma = 0.1$	$\sigma = 0.025$	$\sigma = 0.05$
6	C-ICPE-TD3	<b>10.4</b> [10.2,10.6]	15.5 [15.2,15.9]	13.3 [12.5,14.2]	16.3 [15.7,16.8]
	C-ICPE-TS	11.1 [11.0,11.3]	<b>15.3</b> [15.2,15.4]	<b>12.0</b> [11.9,12.0]	<b>14.4</b> [14.3,14.6]
	C-ICPE-TTPS	11.3 [11.0,11.6]	16.2 [16.1,16.3]	12.4 [12.1,12.6]	15.1 [14.9,15.4]
	C-ICPE-uniform	36.3 [36.0,36.7]	54.1 [52.4,56.9]	60.0 [58.2,62.3]	57.3 [56.1,58.1]
8	C-ICPE-TD3	<b>11.8</b> [11.7,11.9]	<b>17.3</b> [17.0,17.5]	15.8 [15.0,16.8]	19.7 [19.1,20.2]
	C-ICPE-TS	12.0 [11.5,12.4]	17.4 [16.9,17.7]	<b>13.0</b> [12.6,13.4]	16.2 [15.8,16.5]
	C-ICPE-TTPS	12.3 [12.2,12.4]	17.8 [17.6,18.0]	13.4 [13.2,13.6]	<b>16.2</b> [15.9,16.5]
	C-ICPE-uniform	44.0 [43.4,44.7]	62.1 [61.6,62.7]	62.1 [60.0,64.4]	56.0 [53.8,58.8]
10	C-ICPE-TD3	<b>12.7</b> [12.5,12.8]	19.4 [18.5,20.7]	18.2 [16.2,20.4]	22.2 [20.5,23.6]
	C-ICPE-TS	12.8 [12.7,12.8]	<b>18.5</b> [18.1,18.8]	<b>14.0</b> [13.8,14.1]	<b>17.0</b> [16.7,17.3]
	C-ICPE-TTPS	13.4 [12.9,13.8]	18.8 [18.7,18.9]	14.1 [14.0,14.2]	17.3 [17.0,17.7]
	C-ICPE-uniform	47.6 [45.2,49.3]	68.1 [67.2,68.9]	59.2 [57.9,60.4]	51.7 [47.1,55.0]
20	C-ICPE-TD3	35.8 [19.2,67.6]	33.0 [27.2,42.2]	26.7 [23.1,32.3]	28.2 [26.9,29.5]
	C-ICPE-TS	<b>16.4</b> [16.2,16.6]	<b>23.3</b> [23.0,23.5]	<b>17.6</b> [17.3,17.9]	<b>22.7</b> [21.8,23.5]
	C-ICPE-TTPS	16.4 [16.1,16.6]	24.4 [24.1,24.7]	18.7 [18.0,19.4]	24.2 [23.2,25.5]
	C-ICPE-uniform	64.0 [62.0,66.0]	47.5 [29.5,62.4]	81.1 [71.6,88.2]	59.6 [19.4,91.9]

Table 3: Binary search: sample complexity (mean and 95% CI) for every  $(d, \varepsilon, \sigma)$  configuration.

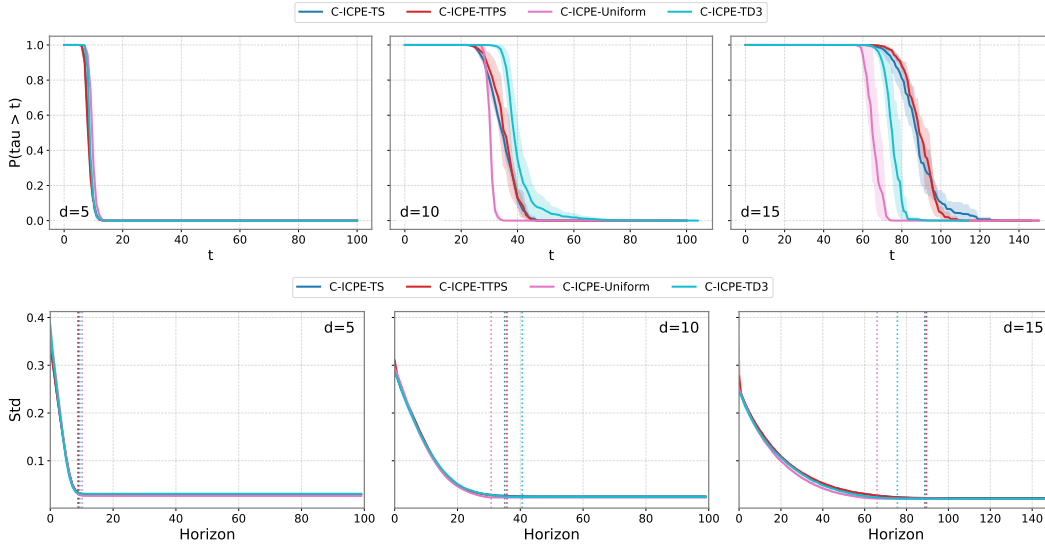


Figure 11: Results for  $\epsilon$ -Best-Arm Identification Problem with fixed confidence  $\delta = 0.005$  and  $N = 100/150$  across different dimensions at the most hardest  $(\epsilon, \sigma)$  setting: (top) survival function of  $\tau$ ; (bottom) inference uncertainty convergence.

2180 Tables 4 and 5 report accuracy and sample complexity; Figure 11 shows the survival function and  
 2181 inference uncertainty convergence.

2182 *Accuracy.* All C-ICPE variants meet the  $1 - \delta$  target across every  $(d, \epsilon, \sigma)$  configuration, with mean  
 2183 accuracy between 0.897 and 0.939. Notably, C-ICPE-Uniform also meets the target throughout  
 2184 (accuracy 0.897–0.936), consistent with the rotational symmetry of the problem: since the loss  
 2185  $L_\theta(x) = 1 - \theta^\top x$  is invariant to orthogonal transformations, no query direction is intrinsically  
 2186 more informative than another, and isotropic exploration is in general optimal. Lazy Track-and-Stop  
 2187 achieves perfect accuracy (1.000) in all configurations by exploiting the known linear observation  
 2188 structure and an analytic likelihood-ratio stopping rule.

2189 *Sample complexity.* The gap between C-ICPE and LT&S is substantial and grows with dimension.  
 2190 LT&S queries exactly  $d$  samples in every configuration, one per canonical direction, achieving the  
 2191 minimum for a full-rank linear estimator. C-ICPE uses 2–6 $\times$  more queries: at  $d = 15$ ,  $(\epsilon, \sigma) =$   
 2192  $(0.005, 0.0025)$ , C-ICPE-TS stops at 88.8 and C-ICPE-TTPS at 89.5, compared to 15.0 for LT&S.  
 2193 This gap reflects the cost of a model-agnostic stopping rule: C-ICPE does not know the observation  
 2194 model is linear and must learn when to stop from interaction data alone. Among C-ICPE variants,  
 2195 C-ICPE-Uniform is competitive with and sometimes more sample-efficient than the active actors.  
 2196 At  $d = 15$ ,  $(\epsilon, \sigma) = (0.005, 0.0025)$ , C-ICPE-Uniform stops at 65.9 while C-ICPE-TS requires  
 2197 88.8 and C-ICPE-TTPS 89.5. This inversion occurs because the problem’s symmetry makes directed  
 2198 exploration unnecessary, and the overhead of posterior-driven action selection, which occasionally  
 2199 concentrates queries in already well-estimated directions, slows convergence relative to uniform  
 2200 coverage. We also note that the current noise levels are low relative to  $\epsilon$ , placing the problem in a  
 2201 regime where LT&S resolves the target direction in a single pass of  $d$  orthogonal queries. At higher  
 2202 noise levels, where multiple measurement rounds are necessary, we expect the relative performance  
 2203 of C-ICPE to improve.

2204 *Stopping behavior.* The survival functions (Figure 11, top) show that all C-ICPE actors have similar  
 2205 stopping profiles at  $d = 5$ , where episodes terminate within  $t \approx 10$ . At  $d = 15$ , the curves spread:  
 2206 C-ICPE-Uniform stops earlier (median  $\approx 60$ ) than C-ICPE-TS and C-ICPE-TTPS (median  $\approx 75$ –  
 2207 85), again reflecting the advantage of isotropic coverage in this symmetric problem. The inference  
 2208 standard deviation (bottom row) converges at comparable rates across all actors, confirming that the  
 2209 posterior contracts uniformly regardless of the exploration strategy, the sample complexity differences  
 2210 are driven by stopping calibration, not by differences in information acquisition.

$d$	Method	$\varepsilon = 0.02$		$\varepsilon = 0.005$	
		$\sigma = 0.005$	$\sigma = 0.01$	$\sigma = 0.00125$	$\sigma = 0.0025$
5	C-ICPE-TD3	0.925 [0.907,0.943]	0.915 [0.905,0.926]	0.898 [0.890,0.907]	0.902 [0.893,0.912]
	C-ICPE-TS	0.939 [0.912,0.958]	0.939 [0.929,0.947]	0.908 [0.890,0.930]	0.904 [0.891,0.917]
	C-ICPE-TTPS	0.934 [0.922,0.944]	0.938 [0.930,0.944]	0.902 [0.891,0.911]	0.900 [0.893,0.909]
	C-ICPE-uniform	0.936 [0.925,0.946]	0.931 [0.919,0.943]	0.922 [0.914,0.928]	0.927 [0.918,0.935]
	LT&S round-robin	<b>1.000</b> [1.000,1.000]	<b>1.000</b> [1.000,1.000]	<b>1.000</b> [1.000,1.000]	<b>1.000</b> [1.000,1.000]
	LT&S uniform	1.000 [1.000,1.000]	1.000 [1.000,1.000]	1.000 [1.000,1.000]	1.000 [1.000,1.000]
10	C-ICPE-TD3	0.920 [0.908,0.931]	0.901 [0.891,0.909]	0.908 [0.896,0.919]	0.919 [0.906,0.932]
	C-ICPE-TS	0.902 [0.893,0.911]	0.909 [0.897,0.920]	0.908 [0.892,0.922]	0.915 [0.899,0.930]
	C-ICPE-TTPS	0.917 [0.907,0.926]	0.904 [0.897,0.911]	0.911 [0.900,0.921]	0.906 [0.892,0.920]
	C-ICPE-uniform	0.911 [0.901,0.921]	0.905 [0.895,0.913]	0.905 [0.894,0.916]	0.900 [0.890,0.910]
	LT&S round-robin	<b>1.000</b> [1.000,1.000]	<b>1.000</b> [1.000,1.000]	<b>1.000</b> [1.000,1.000]	<b>1.000</b> [1.000,1.000]
	LT&S uniform	1.000 [1.000,1.000]	1.000 [1.000,1.000]	1.000 [1.000,1.000]	1.000 [1.000,1.000]
15	C-ICPE-TD3	0.904 [0.894,0.913]	0.906 [0.896,0.915]	0.909 [0.898,0.920]	0.910 [0.902,0.917]
	C-ICPE-TS	0.918 [0.898,0.937]	0.918 [0.900,0.932]	0.908 [0.895,0.921]	0.898 [0.889,0.907]
	C-ICPE-TTPS	0.907 [0.894,0.919]	0.907 [0.894,0.919]	0.910 [0.887,0.934]	0.923 [0.893,0.959]
	C-ICPE-uniform	0.897 [0.887,0.905]	0.902 [0.889,0.913]	0.906 [0.893,0.919]	0.901 [0.892,0.910]
	LT&S round-robin	<b>1.000</b> [1.000,1.000]	<b>1.000</b> [1.000,1.000]	<b>1.000</b> [1.000,1.000]	<b>1.000</b> [1.000,1.000]
	LT&S uniform	1.000 [1.000,1.000]	1.000 [1.000,1.000]	1.000 [1.000,1.000]	1.000 [1.000,1.000]

Table 4:  $\varepsilon$ -best arm: accuracy (mean and 95% CI) for every  $(d, \varepsilon, \sigma)$  configuration.

$d$	Method	$\varepsilon = 0.02$		$\varepsilon = 0.005$	
		$\sigma = 0.005$	$\sigma = 0.01$	$\sigma = 0.00125$	$\sigma = 0.0025$
5	C-ICPE-TD3	7.4 [7.3,7.5]	7.5 [7.4,7.5]	9.4 [9.2,9.6]	9.4 [9.3,9.5]
	C-ICPE-TS	7.3 [7.2,7.5]	7.6 [7.5,7.6]	8.9 [8.6,9.1]	9.0 [8.8,9.2]
	C-ICPE-TTPS	7.3 [7.1,7.4]	7.5 [7.4,7.5]	8.7 [8.6,8.9]	8.8 [8.7,9.0]
	C-ICPE-uniform	8.3 [8.2,8.4]	8.2 [8.0,8.4]	10.0 [9.9,10.1]	10.1 [10.0,10.3]
	LT&S round-robin	<b>5.0</b> [5.0,5.0]	<b>5.0</b> [5.0,5.0]	<b>5.0</b> [5.0,5.0]	<b>5.0</b> [5.0,5.0]
	LT&S uniform	5.0 [5.0,5.0]	5.0 [5.0,5.0]	5.0 [5.0,5.0]	5.0 [5.0,5.0]
10	C-ICPE-TD3	21.7 [21.1,22.6]	21.9 [21.6,22.3]	37.3 [35.3,39.0]	40.6 [38.2,43.3]
	C-ICPE-TS	20.0 [19.5,20.7]	23.9 [22.8,25.3]	31.4 [30.1,33.7]	35.0 [34.7,35.5]
	C-ICPE-TTPS	20.7 [20.4,21.0]	22.8 [22.6,23.0]	32.7 [30.9,34.6]	35.8 [33.8,37.3]
	C-ICPE-uniform	22.1 [21.9,22.3]	22.4 [22.2,22.5]	31.1 [30.8,31.4]	30.7 [30.5,31.1]
	LT&S round-robin	<b>10.0</b> [10.0,10.0]	<b>10.0</b> [10.0,10.0]	<b>10.0</b> [10.0,10.0]	<b>10.0</b> [10.0,10.0]
	LT&S uniform	10.0 [10.0,10.0]	10.0 [10.0,10.0]	10.0 [10.0,10.0]	10.0 [10.0,10.0]
15	C-ICPE-TD3	46.9 [46.5,47.2]	48.6 [48.1,49.2]	71.4 [70.2,72.6]	75.7 [73.4,78.9]
	C-ICPE-TS	46.5 [46.0,47.1]	53.8 [53.1,54.7]	74.8 [73.3,76.3]	88.8 [85.4,92.7]
	C-ICPE-TTPS	45.2 [44.0,47.0]	52.7 [50.9,54.1]	78.9 [74.6,87.0]	89.5 [87.7,91.7]
	C-ICPE-uniform	44.5 [43.3,45.8]	45.1 [43.6,46.9]	68.1 [65.6,70.8]	65.9 [64.1,69.4]
	LT&S round-robin	<b>15.0</b> [15.0,15.0]	<b>15.0</b> [15.0,15.0]	<b>15.0</b> [15.0,15.0]	<b>15.0</b> [15.0,15.0]
	LT&S uniform	15.0 [15.0,15.0]	15.0 [15.0,15.0]	15.0 [15.0,15.0]	15.0 [15.0,15.0]

Table 5:  $\varepsilon$ -best arm: sample complexity (mean and 95% CI) for every  $(d, \varepsilon, \sigma)$  configuration.

2211 **D.3.3 Ackley minimization**

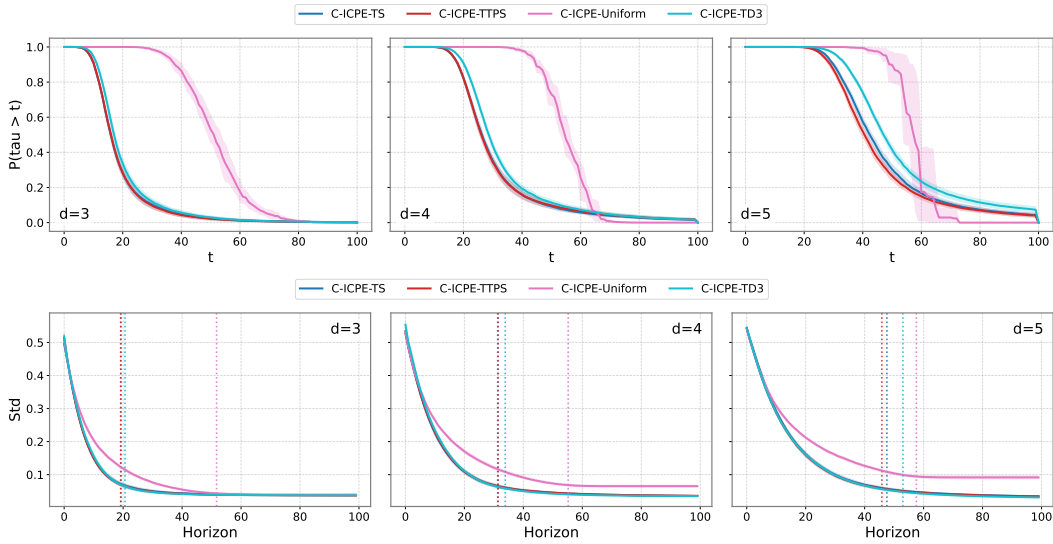


Figure 12: Results for Ackley function with fixed confidence  $\delta = 0.1$  and  $N = 100$  across different dimensions at the most hardest  $(\epsilon, \sigma)$  setting: (top) survival function of  $\tau$ ; (bottom) inference uncertainty convergence.

2212 Tables 7 and 6 report accuracy and sample complexity; Figure 12 shows the survival function and  
 2213 inference uncertainty convergence.

2214 *Accuracy.* All active C-ICPE variants meet the  $1 - \delta = 0.90$  target across every  $(d, \epsilon, \sigma)$  configuration,  
 2215 with mean accuracy between 0.896 and 0.928. C-ICPE-Uniform meets the target at  $d = 3$  (accuracy  
 2216  $\geq 0.890$ ) but degrades as dimension increases: at  $(\epsilon, \sigma) = (0.1, 0.05)$ , accuracy drops from 0.890  
 2217 at  $d = 3$  to 0.589 at  $d = 4$  and 0.237 at  $d = 5$ . Unlike the  $\epsilon$ -best arm problem, the Ackley function  
 2218 has no symmetry that makes uniform exploration competitive: the multimodal landscape and flat  
 2219 outer region require directed queries to distinguish the global minimizer from local optima. The  
 2220 Bayesian optimization baselines fail across the board. GP-UCB achieves the highest accuracy among  
 2221 them (0.344–0.426 at  $d = 3$ ) but remains far below the 0.90 target even at the easiest configuration.  
 2222 TPE and CMA-ES are near zero at  $d \geq 4$ . These methods optimize a fixed-budget objective without  
 2223 a stopping rule and are not designed for  $(\epsilon, \delta)$ -correctness; the comparison confirms that standard  
 2224 continuous optimization does not yield fixed-confidence guarantees at comparable sample budgets.

2225 *Sample complexity.* Active actors use roughly half the samples of C-ICPE-Uniform across all  
 2226 dimensions. At  $d = 5$ ,  $(\epsilon, \sigma) = (0.1, 0.05)$ , C-ICPE-TS stops at 47.5 and C-ICPE-TTPS at 45.8,  
 2227 compared to 57.5 for C-ICPE-Uniform. C-ICPE-TS and C-ICPE-TTPS achieve similar sample  
 2228 complexity throughout, while C-ICPE-TD3 is slightly higher (e.g., 53.0 at the same setting). Sample  
 2229 complexity grows with dimension for all methods: C-ICPE-TS increases from 14.9 ( $d = 3$ ) to 34.4  
 2230 ( $d = 5$ ) at  $(\epsilon, \sigma) = (0.2, 0.05)$ . Increasing  $\sigma$  at fixed  $\epsilon$  consistently raises sample complexity, as  
 2231 expected from the noisier observations.

2232 *Stopping behavior.* The survival functions (Figure 12, top) show a clear separation between active  
 2233 and passive exploration. At  $d = 3$ , C-ICPE-TS and C-ICPE-TTPS exhibit sharp transitions around  
 2234  $t \approx 15$ –25, while C-ICPE-Uniform has a heavy tail extending past  $t = 80$ . At  $d = 5$ , the active  
 2235 actors stop around  $t \approx 35$ –50 while C-ICPE-Uniform rarely stops before  $t = 60$  and retains  
 2236 substantial probability mass near the horizon. The inference standard deviation (bottom row) reveals  
 2237 a qualitative difference from binary search: the posterior uncertainty plateaus around 0.05–0.1 rather  
 2238 than converging to zero. This reflects the inherent difficulty of the Ackley landscape, the multimodal  
 2239 structure and observation noise prevent the inference model from achieving the same posterior  
 2240 concentration as in the unimodal binary search setting. Nevertheless, the critic learns to stop at an  
 2241 appropriate uncertainty level that is sufficient for  $\epsilon$ -correctness.

$d$ Method	$\varepsilon = 0.2$		$\varepsilon = 0.1$	
	$\sigma = 0.05$	$\sigma = 0.1$	$\sigma = 0.025$	$\sigma = 0.05$
3 C-ICPE-TD3	15.3 [14.6,15.9]	18.9 [18.2,19.6]	18.9 [18.2,19.8]	20.6 [19.9,21.3]
C-ICPE-TS	<b>14.9</b> [14.5,15.5]	18.8 [18.1,19.7]	<b>17.5</b> [16.9,18.0]	19.3 [18.5,20.2]
C-ICPE-TTPS	15.4 [14.8,16.1]	<b>18.7</b> [18.1,19.4]	18.0 [17.3,18.7]	<b>19.2</b> [18.6,19.9]
C-ICPE-uniform	37.6 [34.4,40.2]	47.9 [46.3,49.4]	50.5 [48.6,52.2]	51.7 [49.8,53.7]
4 C-ICPE-TD3	25.9 [25.0,26.7]	32.8 [31.8,33.9]	31.5 [30.5,32.6]	33.9 [33.1,34.7]
C-ICPE-TS	24.4 [23.2,25.8]	<b>31.1</b> [29.7,32.6]	<b>27.5</b> [26.3,28.9]	<b>31.3</b> [30.4,32.1]
C-ICPE-TTPS	<b>24.2</b> [22.7,25.8]	32.2 [31.0,33.6]	27.9 [27.0,29.0]	31.4 [30.4,32.4]
C-ICPE-uniform	62.4 [60.3,64.6]	59.4 [56.8,61.4]	57.8 [57.3,58.3]	55.2 [54.0,56.5]
5 C-ICPE-TD3	38.1 [37.0,39.3]	50.7 [49.3,52.3]	44.1 [42.2,46.2]	53.0 [51.9,54.3]
C-ICPE-TS	<b>34.4</b> [33.3,35.9]	<b>45.8</b> [44.2,47.1]	38.5 [37.5,39.5]	47.5 [46.3,48.6]
C-ICPE-TTPS	35.6 [33.4,38.0]	47.0 [45.7,48.1]	<b>37.7</b> [36.9,38.9]	<b>45.8</b> [44.8,46.9]
C-ICPE-uniform	70.0 [63.9,79.2]	58.2 [57.1,59.7]	61.0 [55.3,68.6]	57.5 [55.1,60.1]

Table 6: Ackley: sample complexity (mean and 95% CI) for every  $(d, \varepsilon, \sigma)$  configuration.

$d$ Method	$\varepsilon = 0.2$		$\varepsilon = 0.1$	
	$\sigma = 0.05$	$\sigma = 0.1$	$\sigma = 0.025$	$\sigma = 0.05$
3 C-ICPE-TD3	<b>0.915</b> [.902,.927]	0.904 [.894,.914]	0.898 [.888,.906]	0.908 [.899,.918]
C-ICPE-TS	0.896 [.886,.904]	0.899 [.890,.909]	0.899 [.887,.906]	0.907 [.895,.917]
C-ICPE-TTPS	0.902 [.894,.913]	<b>0.907</b> [.894,.919]	<b>0.911</b> [.900,.920]	<b>0.912</b> [.899,.924]
C-ICPE-uniform	0.902 [.886,.915]	0.905 [.890,.918]	0.906 [.894,.914]	0.890 [.868,.910]
TPE	0.098 [.072,.124]	0.136 [.106,.166]	0.016 [.005,.027]	0.028 [.014,.042]
CMA-ES	0.086 [.061,.111]	0.124 [.095,.153]	0.020 [.008,.032]	0.034 [.018,.050]
GP-logEI	0.322 [.281,.363]	0.356 [.314,.398]	0.218 [.182,.254]	0.194 [.159,.229]
GP-UCB	0.400 [.357,.443]	0.426 [.383,.469]	0.344 [.302,.386]	0.350 [.308,.392]
4 C-ICPE-TD3	0.907 [.895,.919]	<b>0.906</b> [.895,.914]	0.899 [.888,.908]	0.899 [.891,.907]
C-ICPE-TS	0.907 [.899,.918]	0.904 [.892,.916]	<b>0.914</b> [.905,.924]	<b>0.913</b> [.900,.923]
C-ICPE-TTPS	0.898 [.886,.908]	0.898 [.888,.908]	0.908 [.897,.919]	0.906 [.894,.915]
C-ICPE-uniform	<b>0.908</b> [.897,.919]	0.739 [.713,.765]	0.735 [.712,.762]	0.589 [.560,.616]
TPE	0.048 [.029,.067]	0.084 [.060,.108]	0.006 [.000,.013]	0.006 [.000,.013]
CMA-ES	0.020 [.008,.032]	0.056 [.036,.076]	0.002 [.000,.006]	0.002 [.000,.006]
GP-logEI	0.270 [.231,.309]	0.206 [.171,.241]	0.132 [.102,.162]	0.104 [.077,.131]
GP-UCB	0.344 [.302,.386]	0.308 [.267,.349]	0.336 [.295,.377]	0.276 [.237,.315]
5 C-ICPE-TD3	0.904 [.888,.918]	0.904 [.891,.916]	0.907 [.892,.920]	0.904 [.891,.916]
C-ICPE-TS	<b>0.909</b> [.900,.920]	<b>0.912</b> [.899,.921]	0.913 [.905,.925]	0.911 [.897,.921]
C-ICPE-TTPS	0.905 [.893,.915]	0.904 [.889,.915]	<b>0.928</b> [.916,.939]	<b>0.912</b> [.895,.927]
C-ICPE-uniform	0.690 [.618,.756]	0.416 [.401,.437]	0.428 [.332,.538]	0.237 [.203,.279]
TPE	0.030 [.015,.045]	0.042 [.024,.060]	0.000 [.000,.000]	0.002 [.000,.006]
CMA-ES	0.014 [.004,.024]	0.040 [.023,.057]	0.002 [.000,.006]	0.004 [.000,.010]
GP-logEI	0.154 [.122,.186]	0.074 [.051,.097]	0.068 [.046,.090]	0.044 [.026,.062]
GP-UCB	0.344 [.302,.386]	0.214 [.178,.250]	0.284 [.244,.324]	0.220 [.184,.256]

Table 7: Ackley: accuracy (mean and 95% CI) for every  $(d, \varepsilon, \sigma)$  configuration.

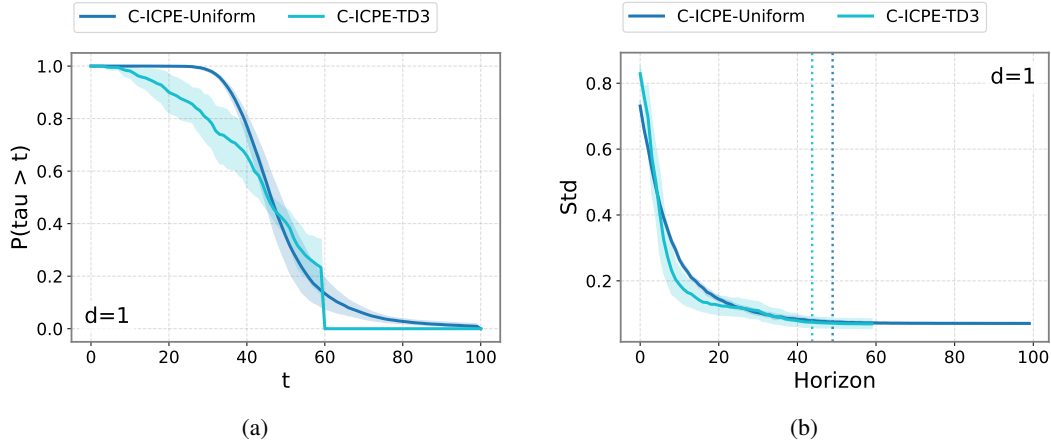


Figure 13: Results for GP value estimation with fixed confidence  $\delta = 0.1$  and  $N = 100$  across different dimensions at the most hardest  $(\epsilon, \sigma)$  setting: (a) survival function of  $\tau$ ; (b) inference uncertainty convergence.

2243 Tables 8 and 9 report accuracy and sample complexity; Figure 13 shows the survival function and  
 2244 inference uncertainty convergence. This is the  $\mathcal{X} \neq \mathcal{A}$  setting: the agent queries locations  $a \in [0, 1]^d$   
 2245 but recommends a scalar value estimate  $\hat{x} \in \mathbb{R}$ . Only C-ICPE-TD3 and C-ICPE-Uniform are  
 2246 evaluated, since TS and TTPS require  $\mathcal{X} = \mathcal{A}$ . In Fig. 9 we also show how C-ICPE explores in this  
 2247 problem, depicting the queries chosen by the actor, the posterior mean, and the posterior standard  
 2248 deviation over timesteps.

2249 *Accuracy.* Both C-ICPE variants meet the  $1 - \delta$  target across all  $(\epsilon, \sigma)$  configurations. C-ICPE-TD3  
 2250 achieves 0.906–0.930 and C-ICPE-Uniform 0.901–0.931; the two methods are comparable in  
 2251 accuracy, with C-ICPE-Uniform slightly higher at the noisier settings (e.g., 0.927 vs. 0.910 at  
 2252  $\epsilon = 0.2, \sigma = 0.1$ ). The non-parametric baselines fail to meet the target: Uniform bin achieves  
 2253 0.130–0.230 and Uniform top 5% reaches 0.648–0.762. Bayesian optimization baselines are not  
 2254 applicable to this task, as they return locations rather than value estimates.

2255 *Sample complexity.* C-ICPE-TD3 uses substantially fewer samples than C-ICPE-Uniform, particu-  
 2256 larly at the larger tolerance. At  $(\epsilon, \sigma) = (0.2, 0.05)$ , C-ICPE-TD3 stops at 13.7 queries on average  
 2257 compared to 31.7 for C-ICPE-Uniform — a  $2.3\times$  reduction. At  $(\epsilon, \sigma) = (0.2, 0.1)$  the ratio is  $1.8\times$   
 2258 (22.7 vs. 41.0). The gap narrows at  $\epsilon = 0.15$ : 39.9 vs. 41.6 at  $\sigma = 0.0375$  and 43.7 vs. 49.0 at  
 2259  $\sigma = 0.075$ . This pattern indicates that the learned exploration policy provides the largest benefit when  
 2260 the tolerance is generous enough that a well-chosen sequence of queries can resolve the max-value  
 2261 quickly, whereas at tighter tolerances both methods require extensive coverage and the advantage of  
 2262 directed exploration diminishes.

2263 *Stopping behavior.* The survival function (left plot in Fig. 13) shows that C-ICPE-TD3 stops earlier  
 2264 than C-ICPE-Uniform, with the bulk of episodes terminating between  $t = 20$  and  $t = 60$ . Both  
 2265 methods exhibit gradual transitions rather than the sharp drops observed in binary search, reflecting  
 2266 the greater variability in task difficulty under the GP prior (functions with short lengthscales require  
 2267 more queries to resolve the peak value). The inference standard deviation (right plot in Fig. 13) starts  
 2268 high ( $\approx 0.8$ ) and contracts rapidly in the first 20–40 queries, then plateaus. This residual uncertainty  
 2269 is consistent with the difficulty of estimating a function’s global maximum from noisy pointwise  
 2270 observations: even after localizing the region of high values, the precise peak height remains uncertain  
 2271 until sufficient samples accumulate near the optimum.

$d$ Method	$\varepsilon = 0.2$		$\varepsilon = 0.15$	
	$\sigma = 0.05$	$\sigma = 0.1$	$\sigma = 0.0375$	$\sigma = 0.075$
1 C-ICPE-TD3	<b>0.906</b> [.878,.928]	0.910 [.883,.932]	<b>0.920</b> [.890,.942]	0.930 [.903,.948]
C-ICPE-uniform	0.902 [.889,.914]	<b>0.927</b> [.915,.938]	0.901 [.887,.916]	<b>0.931</b> [.915,.946]
Uniform bin	0.130 [.100,.160]	0.176 [.143,.209]	0.230 [.193,.267]	0.226 [.189,.263]
Uniform top 5%	0.648 [.606,.690]	0.672 [.631,.713]	0.762 [.725,.799]	0.696 [.656,.736]

Table 8: GP value estimation: accuracy (mean and 95% CI) for every  $(d, \varepsilon, \sigma)$  configuration.

$d$ Method	$\varepsilon = 0.2$		$\varepsilon = 0.15$	
	$\sigma = 0.05$	$\sigma = 0.1$	$\sigma = 0.0375$	$\sigma = 0.075$
1 C-ICPE-TD3	<b>13.7</b> [12.3,15.4]	<b>22.7</b> [21.8,23.8]	<b>39.9</b> [36.6,43.4]	<b>43.7</b> [41.3,45.8]
C-ICPE-uniform	31.7 [30.6,32.8]	41.0 [39.8,42.1]	41.6 [39.5,43.9]	49.0 [46.7,51.2]

Table 9: GP value estimation: sample complexity (mean and 95% CI) for every  $(d, \varepsilon, \sigma)$  configuration.

2272 **D.4 Synthetic Benchmarks: robustness**

2273 When training C-ICPE, the task prior  $\nu$  is uniform over the parameter space. We investigate robustness  
 2274 to prior misspecification by evaluating frozen C-ICPE models under Beta( $\alpha, \beta$ ) deployment priors  
 2275 with varying concentration parameters. Figs. 14 to 17 report accuracy and confidence intervals for all  
 2276 actor variants at the most challenging ( $\epsilon, \sigma$ ) configuration and highest dimensionality per benchmark.

2277 **D.4.1 Noisy Binary Search**

2278 The training prior is Uniform $[-1, 1]^{20}$ , corresponding to  $\alpha = \beta = 1$  (white box). In Fig. 14 we report  
 2279 the results. C-ICPE-TS and C-ICPE-TTPS are robust across the full grid of Beta priors: accuracy  
 2280 remains above 0.83 in all configurations, even under substantial distributional shift ( $\alpha = 0.5, \beta = 7$   
 2281 or vice versa). Performance improves mildly when both  $\alpha, \beta \geq 3$ , since concentrated priors place  
 2282 more mass in the interior of  $[-1, 1]^d$ , where localization is easier. C-ICPE-TD3 matches or exceeds  
 2283 the other actors when the deployment prior is concentrated ( $\alpha, \beta \geq 3$ , reaching 0.943), but degrades  
 2284 sharply when either parameter is small: at  $\alpha = 0.5, \beta = 7$  accuracy drops to 0.612. Small  $\alpha$  or  $\beta$   
 2285 produces a U-shaped Beta distribution that concentrates mass near the boundary of the domain, where  
 2286 targets are harder to disambiguate and the learned actor generalizes poorly. C-ICPE-Uniform fails  
 2287 entirely in  $d = 20$  (accuracy  $< 0.01$ ), confirming that passive exploration cannot localize a target in  
 2288 high-dimensional binary search within the allowed horizon.

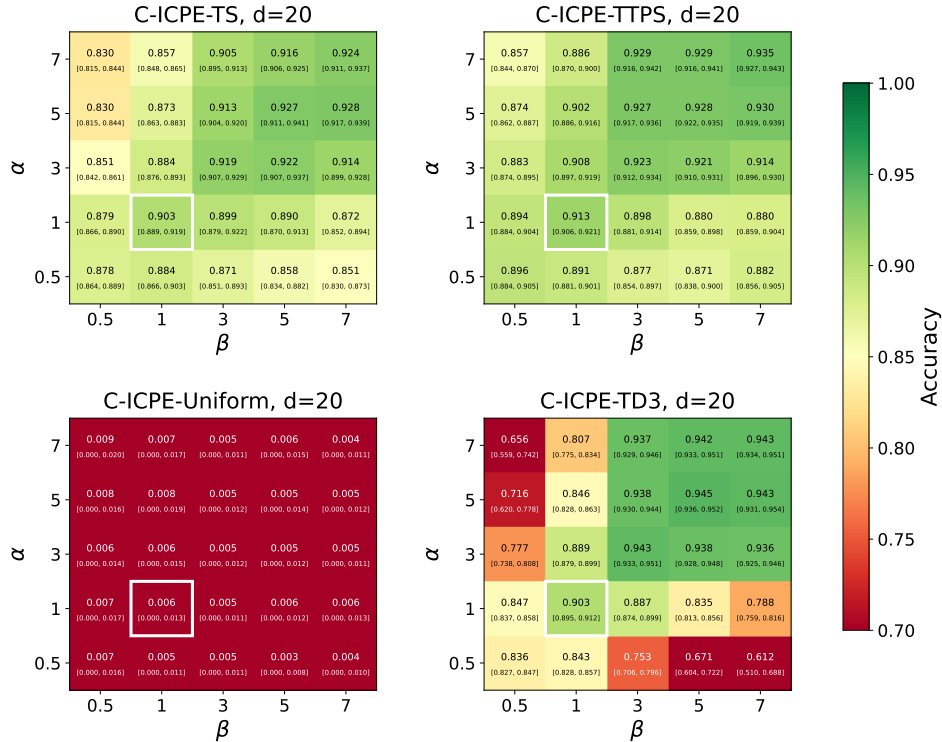


Figure 14: Robustness to prior misspecification on the 20D noisy binary search problem ( $\epsilon = 0.1$  and  $\sigma = 0.05$ ). Each heatmap reports the mean accuracy and the confidence intervals under varying Beta prior hyperparameters  $(\alpha, \beta) \in \{0.5, 1, 3, 5, 7\}$ . The white box indicates the matched prior ( $\alpha = \beta = 1$ ) during training.

2289 **D.4.2  $\epsilon$ -best arm problem**

2290 We report results in Fig. 15. All four C-ICPE variants are remarkably stable across the full Beta  
 2291 prior grid at  $d = 15$ : accuracy varies by less than 0.04 across all  $(\alpha, \beta)$  configurations for each  
 2292 actor. This robustness is a direct consequence of the rotational symmetry of the problem. Since  $\theta$  is  
 2293 drawn on  $\mathbb{S}^{d-1}$  and the loss  $L_\theta(x) = 1 - \theta^\top x$  is invariant to orthogonal transformations, the intrinsic

2294 difficulty of each task instance does not depend on the location of  $\theta$  on the sphere. Reweighting  
 2295 the prior therefore has little effect on the distribution of problem difficulty, unlike binary search or  
 2296 Ackley where boundary effects create heterogeneous difficulty across the parameter space. Notably,  
 2297 C-ICPE-Uniform performs well here ( $\approx 0.89$ – $0.90$  uniformly), consistent with the observation that  
 2298 isotropic exploration is optimal [30]. C-ICPE-TTPS achieves the highest accuracy overall (0.92–  
 2299 0.94), suggesting that the challenger mechanism provides a modest benefit even in a setting where  
 2300 uniform exploration is already near-optimal.

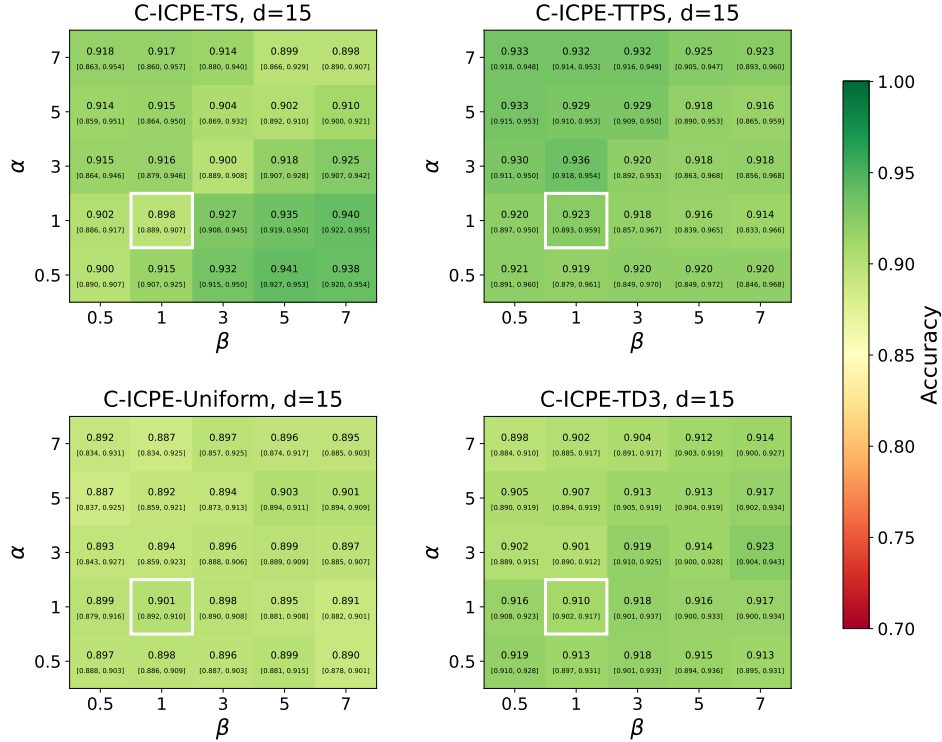


Figure 15: Robustness to prior misspecification on the 15D  $\varepsilon$ -best-arm identification problem ( $\varepsilon = 0.005$  and  $\sigma = 0.0025$ ). Each heatmap reports the mean accuracy and the confidence intervals under varying Beta prior hyperparameters  $(\alpha, \beta) \in \{0.5, 1, 3, 5, 7\}$ . The white box indicates the matched prior ( $\alpha = \beta = 1$ ) during training.

### 2301 D.4.3 Ackley minimization

2302 Results are reported in Fig. 16. The Ackley benchmark at  $d = 5$  exhibits the strongest sensitivity  
 2303 to prior misspecification among all tasks. At the training prior ( $\alpha = \beta = 1$ ), C-ICPE-TS achieves  
 2304 0.911 and C-ICPE-TTPS 0.912. When both  $\alpha, \beta \geq 3$ , i.e., the deployment prior concentrates mass  
 2305 toward the interior of  $[-1, 1]^d$ , all active methods improve substantially, with C-ICPE-TS reaching  
 2306 0.975 and C-ICPE-TD3 0.963. The gains reflect the structure of the Ackley function: targets near  
 2307 the center of the domain sit in a region of higher curvature where observations are more informative,  
 2308 making identification easier.

2309 Conversely, when either  $\alpha$  or  $\beta$  is small (0.5), the Beta prior becomes U-shaped or boundary-  
 2310 skewed, placing significant mass on targets near the edges of  $[-1, 1]^d$ . In the flat outer region of the  
 2311 Ackley function, observations carry little signal, and the learned policies, trained under a uniform  
 2312 prior that rarely produces such extreme configurations, degrade. The effect is most pronounced  
 2313 for  $\alpha = 0.5, \beta = 7$ . This asymmetry between interior and boundary targets is specific to the  
 2314 Ackley geometry and is absent in the rotationally symmetric  $\varepsilon$ -best arm problem. C-ICPE-Uniform  
 2315 fails across the board (accuracy  $\leq 0.45$ ), confirming that directed exploration is essential for this  
 2316 multimodal benchmark regardless of the prior.

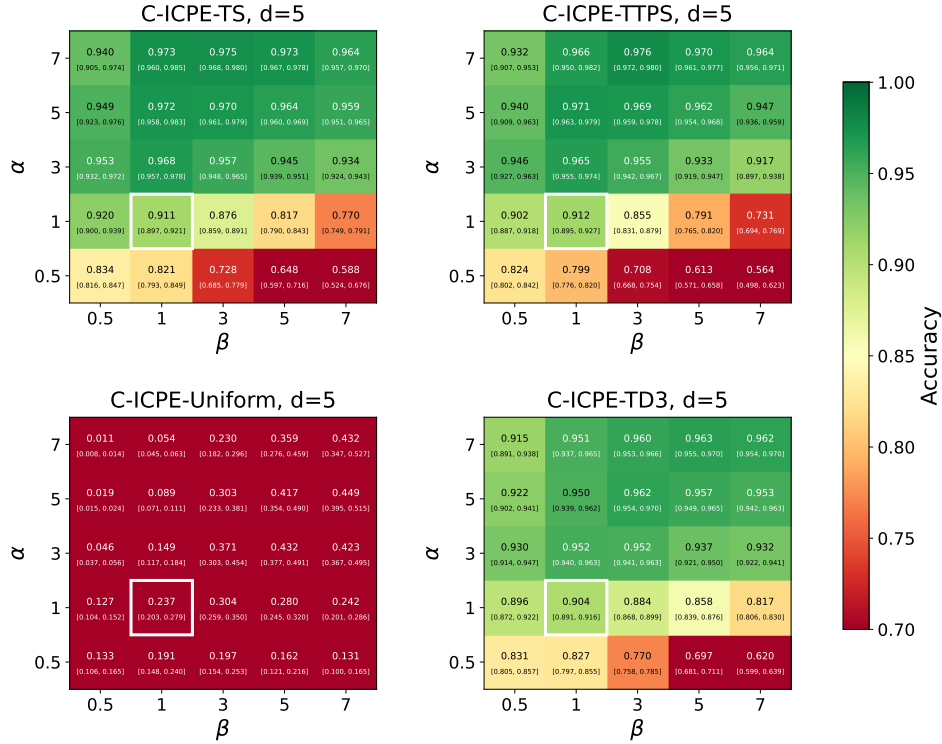


Figure 16: Robustness to prior misspecification on the 5D Ackley function ( $\varepsilon = 0.1$  and  $\sigma = 0.05$ ). Each heatmap reports the mean accuracy and the confidence intervals under varying Beta prior hyperparameters  $(\alpha, \beta) \in \{0.5, 1, 3, 5, 7\}$ . The white box indicates the matched prior ( $\alpha = \beta = 1$ ) during training.

#### 2317 D.4.4 GP max-value estimation

2318 We report results in Fig. 17. This is the  $\mathcal{X} \neq \mathcal{A}$  setting ( $d = 1$ ), so only C-ICPE-Uniform and  
 2319 C-ICPE-TD3 are applicable. C-ICPE-TD3 achieves the highest accuracy when the deployment prior  
 2320 is well-matched with small  $\beta$ : it reaches 0.992 at  $(\alpha, \beta) = (7, 0.5)$  and remains  
 2321 above 0.94 throughout the upper-left triangle of the grid. However, it degrades when  $\beta$  is large (0.772  
 2322 at  $\alpha = 0.5, \beta = 7$ ; 0.817 at  $\alpha = 1, \beta = 7$ ), indicating that the learned exploration policy is sensitive  
 2323 to prior shifts that alter the distribution of task difficulty. C-ICPE-Uniform, by contrast, is more  
 2324 robust: accuracy stays between 0.837 and 0.965 across the entire grid, with a milder gradient from  
 2325 upper-left to lower-right. Because the uniform actor does not depend on a learned exploration policy,  
 2326 its performance varies only through the stopping criterion and inference model, both of which appear  
 2327 stable under moderate prior shift. At the training prior ( $\alpha = \beta = 1$ ), the two methods are comparable  
 2328 (0.920 vs. 0.913), but C-ICPE-TD3 offers a clear advantage when the deployment prior concentrates  
 2329 mass on tasks where directed exploration helps most ( $\alpha \geq 3, \beta \leq 1$ ). The overall pattern suggests  
 2330 that the TD3 actor learns an exploration strategy well-adapted to the training distribution but with  
 2331 limited extrapolation to deployment priors that shift the typical task structure.

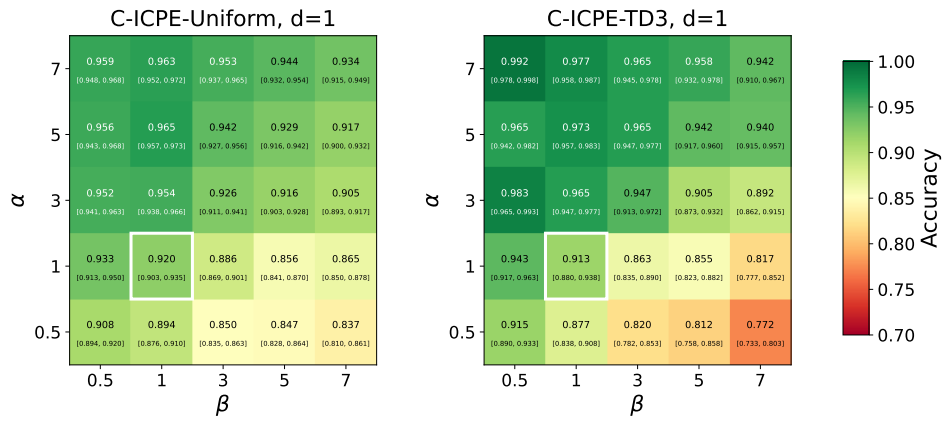


Figure 17: Robustness to prior misspecification on the GP value estimation problem ( $\varepsilon = 0.15$  and  $\sigma = 0.075$ ). Each heatmap reports the mean accuracy and the confidence intervals under varying Beta prior hyperparameters  $(\alpha, \beta) \in \{0.5, 1, 3, 5, 7\}$ . The white box indicates the matched prior ( $\alpha = \beta = 1$ ) during training.

## 2332 D.5 Geochemical Exploration: Experimental Details and Numerical Results

2333 The geochemical exploration task is a stylized version of a real problem in mineral exploration: given  
2334 a budget of field samples (each requiring physical collection, transport, and laboratory analysis),  
2335 identify the most promising location for further investigation. In practice, each sample costs hundreds  
2336 to thousands of dollars, and field campaigns are logistically constrained. A method that can identify  
2337 the target location with  $(\epsilon, \delta)$ -guarantees while minimizing the number of samples has direct economic  
2338 value.

### 2339 D.5.1 Dataset and Motivation

2340 We use data from the USGS National Geochemical Survey [1], which provides soil and sediment  
2341 measurements of element concentrations across the continental United States. We focus on copper  
2342 (Cu) concentrations, which are of direct interest in mineral exploration: copper deposits are spatially  
2343 heterogeneous, and identifying regions of peak concentration from sparse, noisy field measurements  
2344 is a costly sequential decision problem.

2345 The dataset contains point measurements at irregularly spaced locations, each reporting the con-  
2346 centration of multiple elements. We extract copper concentration values and apply a log-transform  
2347 followed by z-score normalization per region, yielding standardized log-concentrations that serve as  
2348 observations.

2349 As a concrete example, Fig. 19 shows the Kingman region in southeastern California and southern  
2350 Nevada. This region contains the Mountain Pass rare earth mine (35.5°N, 115.5°W), an open-pit mine  
2351 of rare earth elements in the Mojave Desert. A satellite image of the mine and surrounding terrain,  
2352 acquired by the Advanced Spaceborne Thermal Emission and Reflection Radiometer (ASTER)  
2353 instrument on NASA’s Terra spacecraft on March 28, 2010, is shown in Fig. 19 [39]. The region  
2354 exhibits spatially varying copper concentrations with a clear peak near the mining district, making it  
2355 a representative example of the kind of localization problem C-ICPE is designed to solve.

### 2356 D.5.2 Region Partitioning and GP Fitting

2357 We partition the geochemical survey data into 17 geographic regions, each covering approximately  
2358  $1^\circ \times 2^\circ$  in latitude and longitude. The regions span diverse geological settings across the western  
2359 and southeastern United States: Gadsden, Bozeman, Billings, Wells, Needles, Jenkins, Montgomery,  
2360 Millett, Prescott, Lovelock, Aurora, Holbrook, Atlanta, Ely, Kingman, Winnemucca, and Baker.  
2361 Fig. 20 shows all 17 regions with sample locations colored by normalized log-copper concentration.  
2362 The regions vary substantially in sample density (from  $\sim 50$  to  $\sim 500$  measurements), spatial structure,  
2363 and concentration range, providing a diverse task distribution for meta-training.

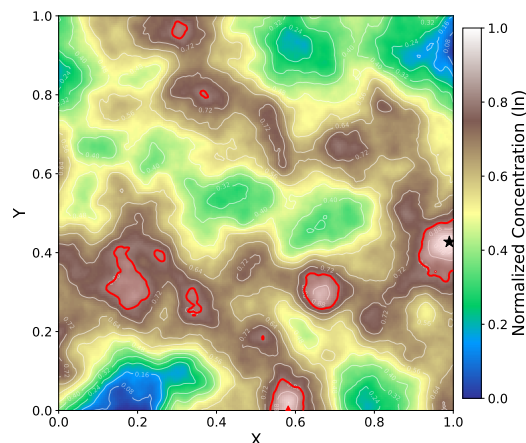


Figure 18: Example copper concentration in a 2D region in the geochemical exploration task. Red regions indicate concentration of copper within  $\epsilon$  of the maximum value.

2364 For each region, we fit a sparse variational Gaussian process (SVGP) with a Matérn-3/2 ARD kernel  
 2365 and Gaussian likelihood. The model is:

$$f \sim \text{GP}(0, \sigma_f^2 k_{\text{Matérn-3/2}}(\cdot, \cdot; \ell_1, \ell_2)), \quad y_i = f(\mathbf{s}_i) + \xi_i, \quad \xi_i \sim \mathcal{N}(0, \sigma_n^2), \quad (37)$$

2366 where  $\mathbf{s}_i \in \mathbb{R}^2$  are UTM coordinates (normalized to  $[0, 1]^2$  for numerical stability),  $y_i$  is the stan-  
 2367 dardized log-copper concentration, and  $(\ell_1, \ell_2, \sigma_f, \sigma_n)$  are per-region hyperparameters learned by  
 2368 maximizing the variational evidence lower bound (ELBO). We use  $M = 500$  inducing points initial-  
 2369 ized via  $k$ -means clustering of the observation locations, and optimize for 1000 Adam iterations at  
 2370 learning rate 0.01.

2371 Fig. 21 shows the ELBO training curves for all 17 regions. All regions converge smoothly, confirming  
 2372 that the SVGP fits are well-behaved. The fitted hyperparameters, in particular the lengthscales  
 2373  $(\ell_1, \ell_2)$ , vary across regions, reflecting different spatial correlation structures: some regions exhibit  
 2374 short-range variability (small lengthscales) while others have smoother concentration surfaces (large  
 2375 lengthscales).

### 2376 D.5.3 Ground Truth Construction

2377 For each fitted SVGP, we evaluate the posterior mean on a dense  $200 \times 200$  grid over the normalized  
 2378 domain  $[0, 1]^2$ . The ground truth target is defined as:

$$\theta^* = \arg \max_{\mathbf{s} \in \text{grid}} \mu_{\text{GP}}(\mathbf{s}), \quad (38)$$

2379 i.e., the grid location with the highest posterior mean copper concentration.

### 2380 D.5.4 Task Prior and Train/Test Split

2381 The 17 regions are split into training and evaluation sets. Training regions define the task prior  $\nu$ :  
 2382 during meta-training, each episode samples a region uniformly from the training set and presents  
 2383 C-ICPE with the corresponding fitted GP as the unknown function. The agent queries 2D locations  
 2384  $a \in [0, 1]^2$  and observes noisy evaluations  $y = \mu_{\text{GP}}(a) + \xi$ ,  $\xi \sim \mathcal{N}(0, \sigma_n^2)$ , where  $\sigma_n$  is the fitted  
 2385 noise standard deviation for that region. The goal is to identify the location  $\theta^*$  of peak concentration  
 2386 to within  $\epsilon$  with probability at least  $1 - \delta$ .

2387 Evaluation is performed on held-out regions whose spatial structure, lengthscales, and concentration  
 2388 patterns were not seen during training. This tests two properties simultaneously:

- 2389 1.  **$(\epsilon, \delta)$ -correct identification on realistic functions.** The GP posterior means are spatially  
 2390 structured, non-stationary (due to irregular sampling), and vary in smoothness across regions,  
 2391 a substantial departure from the synthetic benchmarks.
- 2392 2. **Robustness to distribution shift.** The evaluation regions have different hyperparame-  
 2393 ters  $(\ell_1, \ell_2, \sigma_f, \sigma_n)$  from the training regions, so the agent must generalize across spatial  
 2394 correlation structures it has not encountered during meta-training.

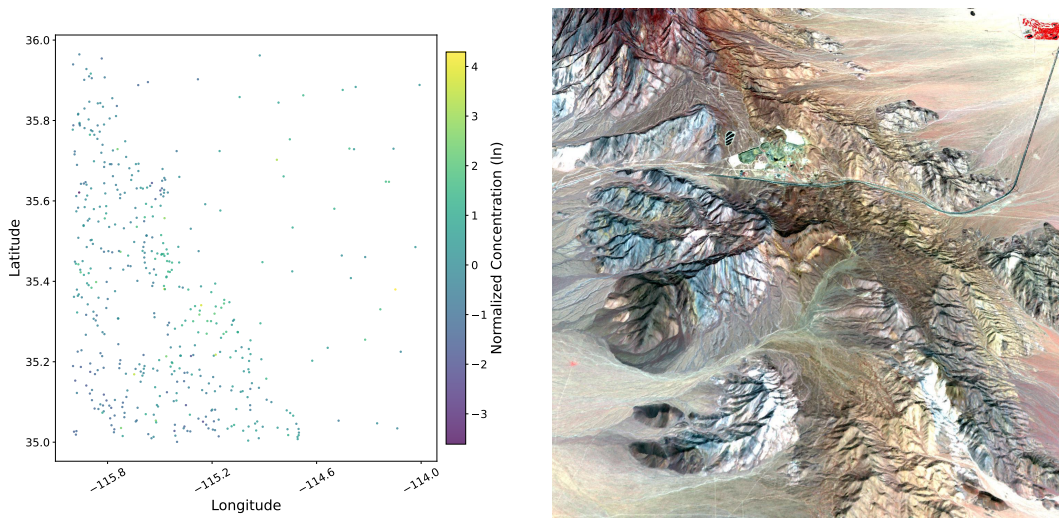


Figure 19: **Left:** Normalized log-copper concentration in the Kingman region (southeastern California / southern Nevada). Each point is a soil sample from the USGS National Geochemical Survey [1]; the red star marks the location of peak GP posterior mean. **Right:** ASTER satellite image of the Mountain Pass rare earth mine (35.5°N, 115.5°W) within this region, acquired March 28, 2010. The mine area is visible as the light-colored open pit in the upper center of the image. Credit: NASA/GSFC/METI/ERSDAC/JAROS, and U.S./Japan ASTER Science Team [39].

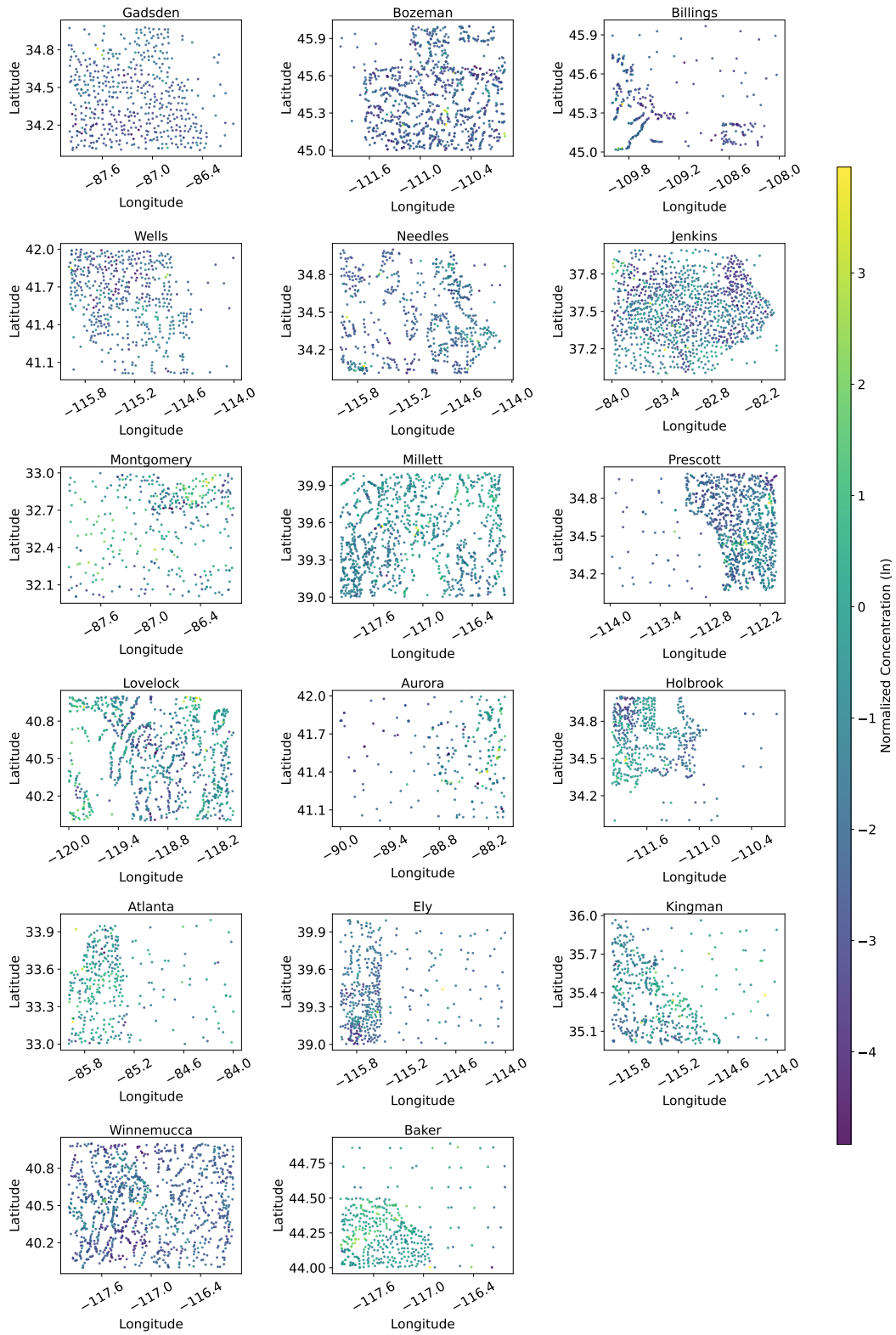


Figure 20: All 17 geographic regions used in the geochemical experiment. Each panel shows soil sample locations colored by normalized log-copper concentration. Regions are split into training and evaluation sets; evaluation regions have spatial structure not seen during meta-training.

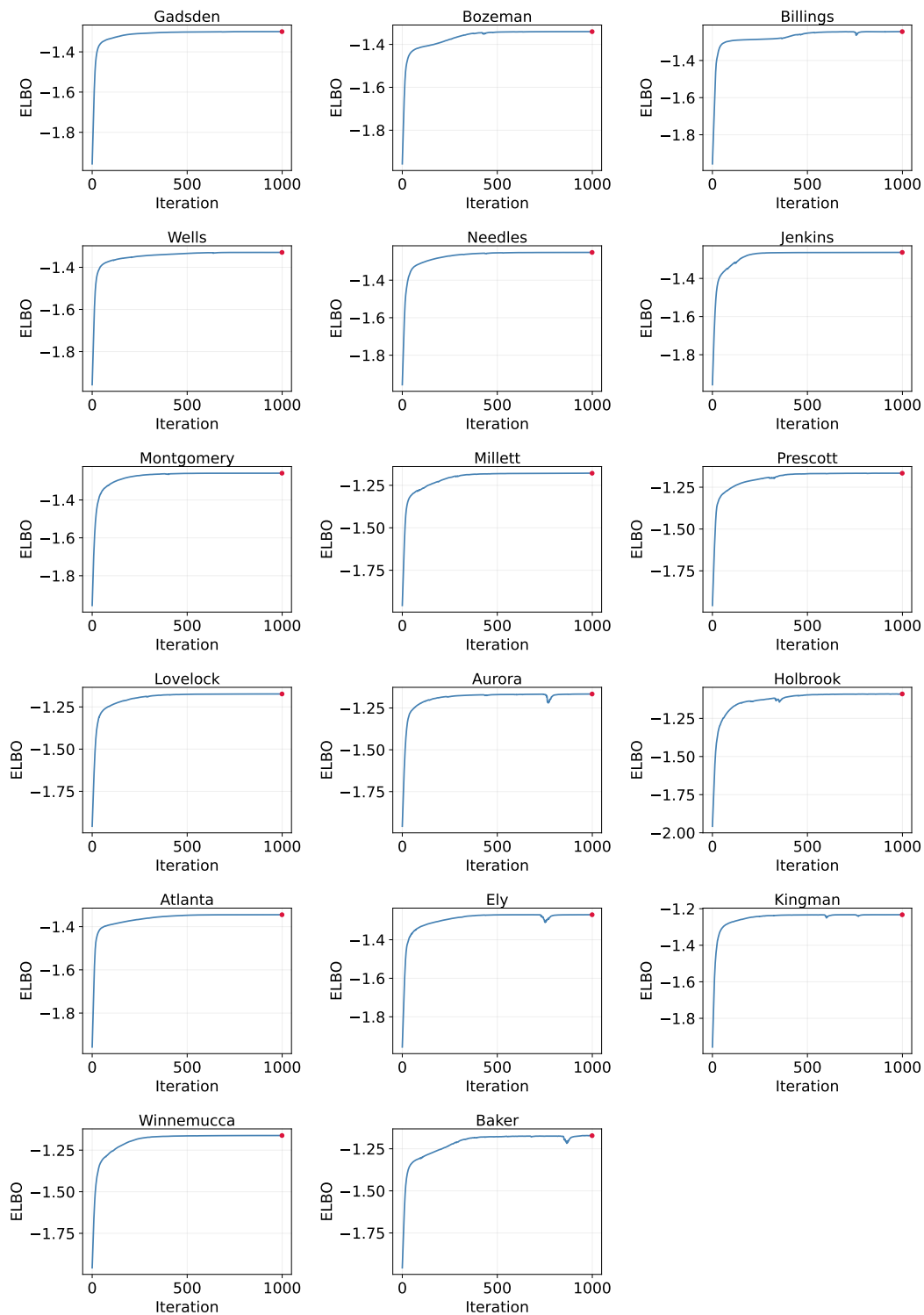


Figure 21: ELBO training curves for sparse variational GP fits across all 17 regions (1000 Adam iterations,  $M = 500$  inducing points). All regions converge smoothly, indicating well-behaved GP fits.

2395 **D.5.5 Numerical Results**

2396 **Baselines.** We compare C-ICPE-TD3 and C-ICPE-TS against four black-box optimization baselines  
 2397 that do not possess a stopping rule for  $(\varepsilon, \delta)$ -correct identification: TPE [8], CMA-ES [24], GP-  
 2398 logEI [5], and GP-UCB [57]. Since these methods optimize a fixed-budget objective and have no  
 2399 principled mechanism for adaptive stopping, we allocate each a fixed sample budget calibrated to  
 2400 the mean sample complexity of C-ICPE-TS:  $N = 22$  for  $\varepsilon = 0.2$  and  $N = 39$  for  $\varepsilon = 0.15$ . After  
 2401 exhausting this budget, each baseline returns the best-observed location as its recommendation. This  
 2402 protocol is deliberately generous: the baselines receive as many samples as C-ICPE-TS typically  
 2403 needs on average, yet bear no cost for deciding when to stop.

	$\varepsilon = 0.2$	$\varepsilon = 0.15$
	$\sigma = 0$	$\sigma = 0$
2 C-ICPE-TD3	<b>0.913</b> [.894,.931]	0.916 [.904,.927]
C-ICPE-TS	0.913 [.894,.930]	<b>0.925</b> [.905,.944]
TPE	0.438 [.394,.482]	0.526 [.482,.570]
CMA-ES	0.484 [.440,.528]	0.514 [.470,.558]
GP-logEI	0.564 [.520,.608]	0.724 [.685,.763]
GP-UCB	0.730 [.691,.769]	0.720 [.681,.759]

Table 10: Geochem: accuracy (mean and 95% CI) for every  $(d, \varepsilon, \sigma)$  configuration.

	$\varepsilon = 0.2$	$\varepsilon = 0.15$
	$\sigma = 0$	$\sigma = 0$
2 C-ICPE-TD3	23.9 [21.6,26.2]	44.4 [42.4,46.5]
C-ICPE-TS	<b>22.0</b> [20.0,23.9]	<b>38.7</b> [34.2,43.5]

Table 11: Geochem: sample complexity (mean and 95% CI) for every  $(d, \varepsilon, \sigma)$  configuration.

2404 **Accuracy and sample complexity.** Table 10 reports identification accuracy (mean and 95% CI  
 2405 over held-out regions) at the two tolerance levels, with  $\delta = 0.1$  and maximum horizon  $N = 150$ .  
 2406 Both C-ICPE variants exceed the  $1 - \delta = 0.90$  correctness target in every configuration. At  $\varepsilon = 0.2$ ,  
 2407 C-ICPE-TD3 and C-ICPE-TS are tied at 0.913 [.894, .931] and 0.913 [.894, .930], respectively. At  
 2408 the harder  $\varepsilon = 0.15$ , C-ICPE-TS pulls ahead with 0.925 [.905, .944] versus 0.916 [.904, .927] for  
 2409 C-ICPE-TD3. Among the baselines, GP-UCB is the strongest, reaching 0.730 [.691, .769] at  $\varepsilon = 0.2$   
 2410 and 0.720 [.681, .759] at  $\varepsilon = 0.15$ , still roughly 19–20 percentage points below C-ICPE despite  
 2411 receiving a comparable sample budget. GP-logEI performs comparably to GP-UCB at  $\varepsilon = 0.15$   
 2412 (0.724) but falls to 0.564 at  $\varepsilon = 0.2$ . TPE and CMA-ES remain below 0.53 in both settings, indicating  
 2413 that gradient-free search without a surrogate is ineffective on these spatially structured surfaces.

2414 Table 11 reports sample complexity for the two C-ICPE variants. C-ICPE-TS is more sample-efficient  
 2415 in both settings: 22.0 [20.0, 23.9] versus 23.9 [21.6, 26.2] at  $\varepsilon = 0.2$ , and 38.7 [34.2, 43.5] versus  
 2416 44.4 [42.4, 46.5] at  $\varepsilon = 0.15$ . The gap widens at the tighter tolerance, suggesting that Thompson  
 2417 sampling’s implicit exploration adapts more efficiently to the difficulty of each region.

2418 **Survival function and uncertainty convergence.** Figure 22a displays the survival function  $\mathbb{P}(\tau >$   
 2419  $t)$  at the hardest setting ( $\varepsilon = 0.15, d = 2$ ). Both variants exhibit a rapid initial decline. The shaded  
 2420 confidence bands for C-ICPE-TD3 are noticeably slightly wider, reflecting higher variance in stopping  
 2421 times, consistent with the wider confidence interval in Table 11.

2422 Figure 22b shows the posterior standard deviation of the recommendation as a function of the  
 2423 horizon. Both methods converge to approximately  $\text{Std} \approx 0.19$  by  $t = 150$ . The vertical dashed  
 2424 lines mark each method’s median stopping time; C-ICPE-TS stops earlier than C-ICPE-TD3, and  
 2425 at both stopping points the standard deviation has already dropped below 0.25. This confirms that  
 2426 C-ICPE-TS stops at a point where the inference model’s posterior is sufficiently concentrated, rather  
 2427 than stopping prematurely. However, it’s interesting to note that the posterior variance has an overall  
 2428 larger decrease with TD3, confirming that C-ICPE-TD3 is learning a good explorative policy. The

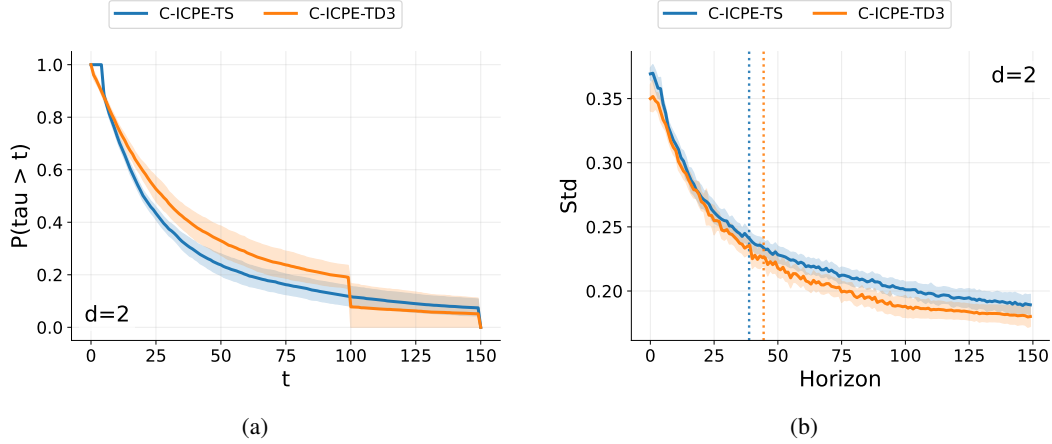


Figure 22: Results for geochemical problem with fixed confidence  $\delta = 0.1$  and  $N = 150$  across different dimensions at the most hardest  $\varepsilon$  setting: (a) survival function of  $\tau$ ; (b) inference uncertainty convergence.

2429 slightly larger stopping time may then be due to the fact that training was stopped early, and one  
 2430 could have trained for longer for better performance of C-ICPE-TD3.

2431 **Robustness to prior misspecification.** Figure 23 reports accuracy under misspecified Beta priors  
 2432  $(\alpha, \beta) \in \{0.5, 1, 3, 5, 7\}^2$  on the normalized  $[0, 1]^2$  domain, with  $\varepsilon = 0.15$ . The matched prior  
 2433 used during meta-training corresponds to  $\alpha = \beta = 1$  (uniform). C-ICPE-TS maintains accuracy  
 2434 between 0.90 and 0.93 across the entire  $5 \times 5$  grid, with no discernible degradation even at extreme  
 2435 configurations such as  $(\alpha, \beta) = (0.5, 7)$  or  $(7, 0.5)$ . C-ICPE-TD3 is similarly stable in the upper  
 2436 portion of the grid ( $\alpha \geq 3$ ), but shows a mild decline to 0.90 at  $(\alpha, \beta) = (0.5, 5)$  and  $(0.5, 7)$ . Across  
 2437 all 25 configurations, every cell remains at or above 0.90, satisfying the  $1 - \delta$  correctness target. This  
 2438 degree of robustness is notably stronger than what is observed on the Ackley benchmark (Figure 16),  
 2439 where boundary-skewed priors cause substantial degradation.

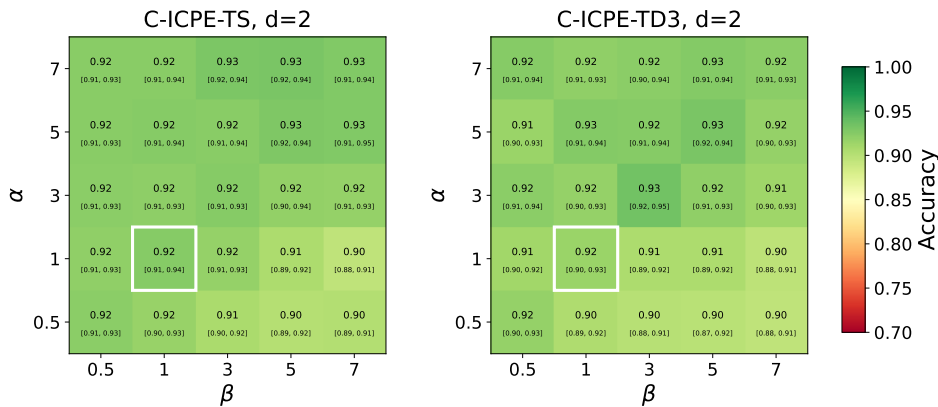


Figure 23: Robustness to prior misspecification on the Geochemical exploration ( $\varepsilon = 0.15$ ). Each heatmap reports the mean accuracy and the confidence intervals under varying Beta prior hyperparameters  $(\alpha, \beta) \in \{0.5, 1, 3, 5, 7\}$ . The white box indicates the matched prior ( $\alpha = \beta = 1$ ) during training.

2440 Taken together, these results demonstrate that C-ICPE transfers to a real-data task involving genuine  
 2441 distribution shift: the evaluation regions have spatial correlation structures, lengthscales, and noise  
 2442 levels not seen during meta-training, yet both variants maintain  $(\varepsilon, \delta)$ -correctness while using fewer  
 2443 samples than fixed-budget baselines that fail to meet the accuracy target. This validates C-ICPE as a  
 2444 practical tool for sequential experimental design.

2445 **NeurIPS Paper Checklist**

2446 **1. Claims**

2447 Question: Do the main claims made in the abstract and introduction accurately reflect the  
2448 paper’s contributions and scope?

2449 Answer: [Yes]

2450 Justification: The abstract and introduction claim three contributions: (i) formulating  
2451 Bayesian fixed-confidence pure exploration with continuous recommendations via the  
2452 posterior success probability, (ii) establishing Bellman optimality and  $(\epsilon, \delta)$ -correctness  
2453 under a local closedness condition, and (iii) a practical algorithm C-ICPE. All three are  
2454 delivered in sections 3-4 and proved in Section B of the appendix. Experimental evaluation  
2455 covers the scope described in the introduction.

2456 Guidelines:

- 2457 • The answer [N/A] means that the abstract and introduction do not include the claims  
2458 made in the paper.
- 2459 • The abstract and/or introduction should clearly state the claims made, including the  
2460 contributions made in the paper and important assumptions and limitations. A [No] or  
2461 [N/A] answer to this question will not be perceived well by the reviewers.
- 2462 • The claims made should match theoretical and experimental results, and reflect how  
2463 much the results can be expected to generalize to other settings.
- 2464 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
2465 are not attained by the paper.

2466 **2. Limitations**

2467 Question: Does the paper discuss the limitations of the work performed by the authors?

2468 Answer: [Yes]

2469 Justification: Limitations are discussed in Section A.

2470 Guidelines:

- 2471 • The answer [N/A] means that the paper has no limitation while the answer [No] means  
2472 that the paper has limitations, but those are not discussed in the paper.
- 2473 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 2474 • The paper should point out any strong assumptions and how robust the results are to  
2475 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
2476 model well-specification, asymptotic approximations only holding locally). The authors  
2477 should reflect on how these assumptions might be violated in practice and what the  
2478 implications would be.
- 2479 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
2480 only tested on a few datasets or with a few runs. In general, empirical results often  
2481 depend on implicit assumptions, which should be articulated.
- 2482 • The authors should reflect on the factors that influence the performance of the approach.  
2483 For example, a facial recognition algorithm may perform poorly when image resolution  
2484 is low or images are taken in low lighting. Or a speech-to-text system might not be  
2485 used reliably to provide closed captions for online lectures because it fails to handle  
2486 technical jargon.
- 2487 • The authors should discuss the computational efficiency of the proposed algorithms  
2488 and how they scale with dataset size.
- 2489 • If applicable, the authors should discuss possible limitations of their approach to  
2490 address problems of privacy and fairness.
- 2491 • While the authors might fear that complete honesty about limitations might be used by  
2492 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
2493 limitations that aren’t acknowledged in the paper. The authors should use their best  
2494 judgment and recognize that individual actions in favor of transparency play an impor-  
2495 tant role in developing norms that preserve the integrity of the community. Reviewers  
2496 will be specifically instructed to not penalize honesty concerning limitations.

2497 **3. Theory assumptions and proofs**

2498 Question: For each theoretical result, does the paper provide the full set of assumptions and  
2499 a complete (and correct) proof?

2500 Answer: [Yes]

2501 Justification: All assumptions are stated explicitly in the appendix. Complete proofs of  
2502 all theorems, propositions, and lemmas are provided in the appendix (Section B) with a  
2503 roadmap at the start distinguishing standard tools from new results.

2504 Guidelines:

- 2505 • The answer [N/A] means that the paper does not include theoretical results.
- 2506 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
2507 referenced.
- 2508 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 2509 • The proofs can either appear in the main paper or the supplemental material, but if  
2510 they appear in the supplemental material, the authors are encouraged to provide a short  
2511 proof sketch to provide intuition.
- 2512 • Inversely, any informal proof provided in the core of the paper should be complemented  
2513 by formal proofs provided in appendix or supplemental material.
- 2514 • Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 2515 4. Experimental result reproducibility

2516 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
2517 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
2518 of the paper (regardless of whether the code and data are provided or not)?

2519 Answer: [Yes]

2520 Justification: Full source code with a README is provided in the supplementary material.  
2521 Benchmark constructions and baseline configurations are detailed in the appendix. Hyper-  
2522 parameter values (learning rates, architecture sizes, training schedules) are specified in the  
2523 configuration files included with the code. The geochemical dataset is publicly available  
2524 from the USGS [1].

2525 Guidelines:

- 2526 • The answer [N/A] means that the paper does not include experiments.
- 2527 • If the paper includes experiments, a [No] answer to this question will not be perceived  
2528 well by the reviewers: Making the paper reproducible is important, regardless of  
2529 whether the code and data are provided or not.
- 2530 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
2531 to make their results reproducible or verifiable.
- 2532 • Depending on the contribution, reproducibility can be accomplished in various ways.  
2533 For example, if the contribution is a novel architecture, describing the architecture fully  
2534 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
2535 be necessary to either make it possible for others to replicate the model with the same  
2536 dataset, or provide access to the model. In general, releasing code and data is often  
2537 one good way to accomplish this, but reproducibility can also be provided via detailed  
2538 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
2539 of a large language model), releasing of a model checkpoint, or other means that are  
2540 appropriate to the research performed.
- 2541 • While NeurIPS does not require releasing code, the conference does require all submis-  
2542 sions to provide some reasonable avenue for reproducibility, which may depend on the  
2543 nature of the contribution. For example
  - 2544 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
2545 to reproduce that algorithm.
  - 2546 (b) If the contribution is primarily a new model architecture, the paper should describe  
2547 the architecture clearly and fully.
  - 2548 (c) If the contribution is a new model (e.g., a large language model), then there should  
2549 either be a way to access this model for reproducing the results or a way to reproduce  
2550 the model (e.g., with an open-source dataset or instructions for how to construct  
2551 the dataset).

2552 (d) We recognize that reproducibility may be tricky in some cases, in which case  
2553 authors are welcome to describe the particular way they provide for reproducibility.  
2554 In the case of closed-source models, it may be that access to the model is limited in  
2555 some way (e.g., to registered users), but it should be possible for other researchers  
2556 to have some path to reproducing or verifying the results.

## 2557 5. Open access to data and code

2558 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
2559 tions to faithfully reproduce the main experimental results, as described in supplemental  
2560 material?

2561 Answer: [Yes]

2562 Justification: Source code and instructions are included in the supplementary material (see  
2563 README.md). The geochemical dataset is publicly available from the USGS National  
2564 Geochemical Survey [1]. All synthetic benchmarks are fully specified in the code and in  
2565 Section D.

2566 Guidelines:

- 2567 • The answer [N/A] means that paper does not include experiments requiring code.
- 2568 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/  
2569 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 2570 • While we encourage the release of code and data, we understand that this might not  
2571 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not  
2572 including code, unless this is central to the contribution (e.g., for a new open-source  
2573 benchmark).
- 2574 • The instructions should contain the exact command and environment needed to run to  
2575 reproduce the results. See the NeurIPS code and data submission guidelines ([https:  
2576 //neurips.cc/public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 2577 • The authors should provide instructions on data access and preparation, including how  
2578 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 2579 • The authors should provide scripts to reproduce all experimental results for the new  
2580 proposed method and baselines. If only a subset of experiments are reproducible, they  
2581 should state which ones are omitted from the script and why.
- 2582 • At submission time, to preserve anonymity, the authors should release anonymized  
2583 versions (if applicable).
- 2584 • Providing as much information as possible in supplemental material (appended to the  
2585 paper) is recommended, but including URLs to data and code is permitted.

## 2586 6. Experimental setting/details

2587 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-  
2588 rameters, how they were chosen, type of optimizer) necessary to understand the results?

2589 Answer: [Yes]

2590 Justification: Benchmark constructions, baseline configurations, and reporting protocols  
2591 are described in Section D of the appendix. The geochemical train/test split is detailed in  
2592 Section D.5. Hyperparameter values are specified in the configuration files included with  
2593 the source code in the supplementary material.

2594 Guidelines:

- 2595 • The answer [N/A] means that the paper does not include experiments.
- 2596 • The experimental setting should be presented in the core of the paper to a level of detail  
2597 that is necessary to appreciate the results and make sense of them.
- 2598 • The full details can be provided either with the code, in appendix, or as supplemental  
2599 material.

## 2600 7. Experiment statistical significance

2601 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
2602 information about the statistical significance of the experiments?

2603 Answer: [Yes]

2604  
2605  
2606  
2607  
2608  
2609  
2610  
2611  
2612  
2613  
2614  
2615  
2616  
2617  
2618  
2619  
2620  
2621  
2622  
2623  
2624  
2625  
2626  
2627  
2628  
2629  
2630  
2631  
2632  
2633  
2634  
2635  
2636  
2637  
2638  
2639  
2640  
2641  
2642  
2643  
2644  
2645  
2646  
2647  
2648  
2649  
2650  
2651  
2652  
2653  
2654  
2655  
2656

Justification: All experiments report 95% confidence intervals computed via hierarchical bootstrap over three levels: random seeds (each corresponding to a trained model), sampled environments, and trajectories within each environment. Confidence intervals are shown in all figures and tables. Details are explained in Section D.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The authors should answer [Yes] if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g., negative error rates).
- If error bars are reported in tables or plots, the authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

**8. Experiments compute resources**

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Computational resources and training details are reported in the numerical results section of the appendix (Section D).

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

**9. Code of ethics**

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: [Yes]

Justification: The research involves no human subjects, no private data, and no dual-use risks. The geochemical data is publicly available. The work is methodological in nature and raises no ethical concerns.

Guidelines:

- The answer [N/A] means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer [No], they should explain the special circumstances that require a deviation from the Code of Ethics.

- 2657 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
2658 eration due to laws or regulations in their jurisdiction).

## 2659 10. Broader impacts

2660 Question: Does the paper discuss both potential positive societal impacts and negative  
2661 societal impacts of the work performed?

2662 Answer: [Yes]

2663 Justification: Broader impact is discussed in Section A.

2664 Guidelines:

- 2665 • The answer [N/A] means that there is no societal impact of the work performed.
- 2666 • If the authors answer [N/A] or [No], they should explain why their work has no societal  
2667 impact or why the paper does not address societal impact.
- 2668 • Examples of negative societal impacts include potential malicious or unintended uses  
2669 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
2670 (e.g., deployment of technologies that could make decisions that unfairly impact specific  
2671 groups), privacy considerations, and security considerations.
- 2672 • The conference expects that many papers will be foundational research and not tied  
2673 to particular applications, let alone deployments. However, if there is a direct path to  
2674 any negative applications, the authors should point it out. For example, it is legitimate  
2675 to point out that an improvement in the quality of generative models could be used to  
2676 generate Deepfakes for disinformation. On the other hand, it is not needed to point out  
2677 that a generic algorithm for optimizing neural networks could enable people to train  
2678 models that generate Deepfakes faster.
- 2679 • The authors should consider possible harms that could arise when the technology is  
2680 being used as intended and functioning correctly, harms that could arise when the  
2681 technology is being used as intended but gives incorrect results, and harms following  
2682 from (intentional or unintentional) misuse of the technology.
- 2683 • If there are negative societal impacts, the authors could also discuss possible mitigation  
2684 strategies (e.g., gated release of models, providing defenses in addition to attacks,  
2685 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
2686 feedback over time, improving the efficiency and accessibility of ML).

## 2687 11. Safeguards

2688 Question: Does the paper describe safeguards that have been put in place for responsible  
2689 release of data or models that have a high risk for misuse (e.g., pre-trained language models,  
2690 image generators, or scraped datasets)?

2691 Answer: [N/A]

2692 Justification: The method addresses pure exploration in continuous spaces and does not pose  
2693 risks for misuse. The geochemical data is publicly available from the USGS.

2694 Guidelines:

- 2695 • The answer [N/A] means that the paper poses no such risks.
- 2696 • Released models that have a high risk for misuse or dual-use should be released with  
2697 necessary safeguards to allow for controlled use of the model, for example by requiring  
2698 that users adhere to usage guidelines or restrictions to access the model or implementing  
2699 safety filters.
- 2700 • Datasets that have been scraped from the Internet could pose safety risks. The authors  
2701 should describe how they avoided releasing unsafe images.
- 2702 • We recognize that providing effective safeguards is challenging, and many papers do  
2703 not require this, but we encourage authors to take this into account and make a best  
2704 faith effort.

## 2705 12. Licenses for existing assets

2706 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
2707 the paper, properly credited and are the license and terms of use explicitly mentioned and  
2708 properly respected?

2709 Answer: [Yes]

2710 Justification: The USGS National Geochemical Survey dataset [1] is publicly available and  
2711 properly cited. All baseline methods and libraries used are cited. Our code will be released  
2712 under the MIT license upon publication.

2713 Guidelines:

- 2714 • The answer [N/A] means that the paper does not use existing assets.
- 2715 • The authors should cite the original paper that produced the code package or dataset.
- 2716 • The authors should state which version of the asset is used and, if possible, include a  
2717 URL.
- 2718 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 2719 • For scraped data from a particular source (e.g., website), the copyright and terms of  
2720 service of that source should be provided.
- 2721 • If assets are released, the license, copyright information, and terms of use in the  
2722 package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets)  
2723 has curated licenses for some datasets. Their licensing guide can help determine the  
2724 license of a dataset.
- 2725 • For existing datasets that are re-packaged, both the original license and the license of  
2726 the derived asset (if it has changed) should be provided.
- 2727 • If this information is not available online, the authors are encouraged to reach out to  
2728 the asset's creators.

### 2729 13. New assets

2730 Question: Are new assets introduced in the paper well documented and is the documentation  
2731 provided alongside the assets?

2732 Answer: [Yes]

2733 Justification: The supplementary material includes the source code with a README  
2734 documenting setup, dependencies, and instructions for reproducing all experiments. The  
2735 geochemical benchmark construction pipeline is included and documented. The code is  
2736 provided as an anonymized zip file.

2737 Guidelines:

- 2738 • The answer [N/A] means that the paper does not release new assets.
- 2739 • Researchers should communicate the details of the dataset/code/model as part of their  
2740 submissions via structured templates. This includes details about training, license,  
2741 limitations, etc.
- 2742 • The paper should discuss whether and how consent was obtained from people whose  
2743 asset is used.
- 2744 • At submission time, remember to anonymize your assets (if applicable). You can either  
2745 create an anonymized URL or include an anonymized zip file.

### 2746 14. Crowdsourcing and research with human subjects

2747 Question: For crowdsourcing experiments and research with human subjects, does the paper  
2748 include the full text of instructions given to participants and screenshots, if applicable, as  
2749 well as details about compensation (if any)?

2750 Answer: [N/A]

2751 Justification: This work does not involve crowdsourcing or human subjects.

2752 Guidelines:

- 2753 • The answer [N/A] means that the paper does not involve crowdsourcing nor research  
2754 with human subjects.
- 2755 • Including this information in the supplemental material is fine, but if the main contribu-  
2756 tion of the paper involves human subjects, then as much detail as possible should be  
2757 included in the main paper.
- 2758 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
2759 or other labor should be paid at least the minimum wage in the country of the data  
2760 collector.

2761 **15. Institutional review board (IRB) approvals or equivalent for research with human**  
2762 **subjects**

2763 Question: Does the paper describe potential risks incurred by study participants, whether  
2764 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
2765 approvals (or an equivalent approval/review based on the requirements of your country or  
2766 institution) were obtained?

2767 Answer: [N/A]

2768 Justification: This work does not involve human subjects or study participants.

2769 Guidelines:

- 2770 • The answer [N/A] means that the paper does not involve crowdsourcing nor research  
2771 with human subjects.
- 2772 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
2773 may be required for any human subjects research. If you obtained IRB approval, you  
2774 should clearly state this in the paper.
- 2775 • We recognize that the procedures for this may vary significantly between institutions  
2776 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
2777 guidelines for their institution.
- 2778 • For initial submissions, do not include any information that would break anonymity (if  
2779 applicable), such as the institution conducting the review.

2780 **16. Declaration of LLM usage**

2781 Question: Does the paper describe the usage of LLMs if it is an important, original, or  
2782 non-standard component of the core methods in this research? Note that if the LLM is used  
2783 only for writing, editing, or formatting purposes and does *not* impact the core methodology,  
2784 scientific rigor, or originality of the research, declaration is not required.

2785 Answer: [N/A]

2786 Justification: LLMs are not used as a component of the methodology. They were used only  
2787 for writing and editing assistance.

2788 Guidelines:

- 2789 • The answer [N/A] means that the core method development in this research does not  
2790 involve LLMs as any important, original, or non-standard components.
- 2791 • Please refer to our LLM policy in the NeurIPS handbook for what should or should not  
2792 be described.