
Efficient Prior Selection in Gaussian Process Bandits with Thompson Sampling

Jack Sandberg, Morteza Haghiri Chehreghani

Department of Computer Science and Engineering
Chalmers University of Technology and Gothenburg University
{jack.sandberg, morteza.chehreghani}@chalmers.se

Abstract

Gaussian process (GP) bandits provide a powerful framework for solving blackbox optimization of unknown functions. The characteristics of the unknown function depend heavily on the assumed GP prior. Most work in the literature assume that this prior is known but in practice this seldom holds. Instead, practitioners often rely on maximum likelihood estimation to select the hyperparameters of the prior - which lacks theoretical guarantees. In this work, we propose two algorithms for joint prior selection and regret minimization in GP bandits based on GP Thompson sampling (GP-TS): Prior-Elimination GP-TS (PE-GP-TS) and HyperPrior GP-TS (HP-GP-TS). We theoretically analyze the algorithms and establish upper bounds for the regret of HP-GP-TS. In addition, we demonstrate the effectiveness of our algorithms compared to the alternatives through experiments with synthetic and real-world data.

1 Introduction

In Gaussian process bandits, we consider a variant of the multi-armed bandit problem where the arms are correlated and their expected reward is sampled from a Gaussian process (GP). The flexibility of GPs have made GP bandits applicable in a wide range of areas that need to optimize blackbox functions with noisy estimates, including machine learning hyperparameter tuning (Turner et al., 2021), drug discovery (Hernández-Lobato et al., 2017; Pyzer-Knapp, 2018), chemical design (Griffiths & Hernández-Lobato, 2020), battery charging protocols (Jiang et al., 2022), online advertising (Nuara et al., 2018), portfolio optimization (Gonzalez et al., 2019) and energy-efficient navigation (Sandberg et al., 2025). However, most of the theoretical results in the literature assume that the GP prior is known but this is seldom the case in practical applications. Even with expert domain knowledge, selecting the exact prior to use can be a difficult task. Most practitioners tend to utilize maximum likelihood estimation (MLE) to identify suitable prior parameters. However, in a sequential decision making problem MLE is not guaranteed to recover the correct parameters.

In the literature, Wang & de Freitas (2014); Berkenkamp et al. (2019); Ziomek et al. (2024) provided algorithms with theoretical guarantees when the kernel lengthscale is unknown. More recently, Ziomek et al. (2025) introduced an elimination-based algorithm with theoretical guarantees for an arbitrary set of discrete priors. Their algorithm, Prior-Elimination GP-UCB (PE-GP-UCB), selects the arm and prior which provide the most optimistic upper confidence bound (UCB). If a prior generates too many incorrect predictions, then it may be eliminated. The previous work has focused on optimistic UCB methods which are known to over-explore.

In this work, we investigate the use of Thompson sampling for solving GP-bandit problems with unknown priors and we propose two algorithms. The first algorithm, Prior-Elimination GP-TS (PE-GP-TS), is an extension of PE-GP-UCB that replaces the doubly optimistic selection rule with posterior sampling and one less layer of optimism. We analyze the regret of PE-GP-TS. For the terms

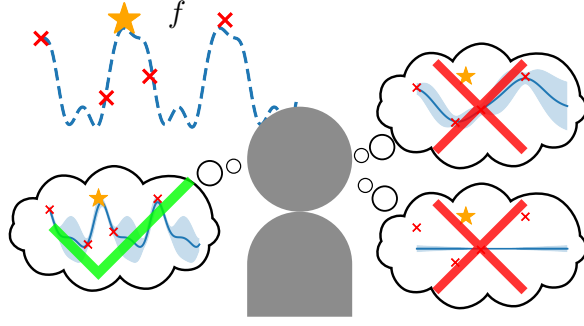


Figure 1: Agent deciding between priors based on observed data. The unknown function f is sampled from a periodic prior and the observed data is in the top left. The agent discards the priors on the right since they do not match the new observation (star).

we can bound, we obtain a regret bound for PE-GP-TS of order $\mathcal{O}(\sqrt{T\beta_T|P|\hat{\gamma}_T})$ where T is the horizon, $|P|$ is the number of priors and $\hat{\gamma}_T$ is the worst-case maximum information gain, which matches that of PE-GP-UCB. The second algorithm, HyperPrior GP-TS (HP-GP-TS), uses bi-level posterior sampling to efficiently explore the priors and arms. We analyze the regret using both a UCB-based and an information-theoretic framework and obtain an information-theoretic regret bound of order $\mathcal{O}(\sqrt{T|\mathcal{X}|\log|\mathcal{X}|})$ where T is the horizon and $|\mathcal{X}|$ is the number of arms.

We evaluate our methods on three sets of synthetic experiments and three experiments with real-world data. Across the experiments, our Thompson sampling based methods outperform PE-GP-UCB. Additionally, we find that the regret of HP-GP-TS does not increase with $|P|$ in our experiments. Finally, we analyze the priors selected by the algorithms and observe that HP-GP-TS selects the correct prior more often than the other algorithms.

The contributions of this work can be summarized as:

- We propose two novel algorithms for GP-bandits with unknown prior: PE- and HP-GP-TS.
- We analyze the regret of HP-GP-TS using a UCB framework and an information-theoretic framework and provide a regret bound of order $\mathcal{O}(\sqrt{T|\mathcal{X}|\log|\mathcal{X}|})$ for HP-GP-TS. Additionally, we analyze the regret of PE-GP-TS.
- We experimentally evaluate our algorithms on both synthetic and real-world data, demonstrating that they achieve superior performance and that the regret of HP-GP-TS does not increase with $|P|$.

2 Background and problem statement

Problem statement We consider a sequential decision making problem where an agent repeatedly selects among a set of arms and receives a random reward whose mean depends on the selected arm and is unknown to the agent. The goal of the agent is to maximize the cumulative sum of rewards over a finite time horizon. We assume that the distribution of the means, the *prior*, is sampled from a set of priors, the *hyperprior*. An effective agent must distinguish which prior the means are sampled from to ensure it explores efficiently.

Now, let us formally state the problem. Let $\mathcal{X} \subseteq [0, r]^d \subset \mathbb{R}^d$ denote the finite set of arms and P a finite set of priors with associated prior mean and kernel functions $\mu_{1,p} : \mathcal{X} \mapsto \mathbb{R}$ and $k_{1,p} : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$, $\forall p \in P$. Let $p^* \in P$ denote the true prior and assume the expected reward function $f : \mathcal{X} \mapsto \mathbb{R} \sim \mathcal{GP}(\mu_{1,p^*}, k_{1,p^*})$ is a sample from a Gaussian process with prior p^* . Both the function f and the true prior p^* are considered unknown. We will consider two settings: In the frequentist selection setting, the prior $p^* \in P$ is picked arbitrarily. In the Bayesian selection setting, the prior is sampled from a hyperprior $p^* \sim \mathcal{P}(P)$. To simplify notation, let P_1 denote the hyperprior.

Let T denote the horizon. For time step $t = 1, 2, \dots, T$, the agent selects an arm $x_t \in \mathcal{X}$ and observes the reward $y_t = f(x_t) + \epsilon_t$ where $\{\epsilon_t\}_{t=1}^T$ are i.i.d. zero-mean Gaussian noise with variance σ^2 . The goal of the agent is to select a sequence of arms $\{x_t\}_{t=1}^T$ that minimizes the

regret $R(T) = \sum_{t \in [T]} f(x^*) - f(x_t)$ where $[T] = \{1, \dots, T\}$ and $x^* = \arg \max_{x \in \mathcal{X}} f(x)$. In the Bayesian selection setting, we evaluate the agent based on the Bayesian regret $\text{BR}(T) = \mathbb{E}[R(T)]$ where the expectation is taken over the prior p^* , the expected reward function f , the noise $\{\epsilon_t\}_{t=1}^T$ and the (potentially) stochastic selection of arms.

Gaussian processes A Gaussian process $f(x) \sim \mathcal{GP}(\mu, k)$ is a collection of random variables such that for any subset $\{x_1, \dots, x_n\} \subset \mathcal{X}$, the vector $[f(x_1), \dots, f(x_n)] \in \mathbb{R}^n$ has a multivariate Gaussian distribution. The probabilistic nature of GPs makes them very useful for defining and solving bandit problems where the arms are correlated. Given the history $H_t = \{(x_i, y_i)\}_{i=1}^{t-1}$, the posterior mean and kernel functions of a Gaussian process $\mathcal{GP}(\mu, k)$ are given by

$$\mu_t(x) = \mu(x) + \mathbf{k}^\top (\mathbf{K} + \sigma^2 I)^{-1} (\mathbf{y} - \boldsymbol{\mu}), \quad (1)$$

$$k_t(x, \tilde{x}) = k(x, \tilde{x}) - \mathbf{k}^\top (\mathbf{K} + \sigma^2 I)^{-1} \tilde{\mathbf{k}}. \quad (2)$$

Above, $\mathbf{k}, \tilde{\mathbf{k}} \in \mathbb{R}^{t-1}$ are vectors such that $(\mathbf{k})_i = k(x_i, x)$ and $(\tilde{\mathbf{k}})_i = k(x_i, \tilde{x})$. Additionally, $\mathbf{y}, \boldsymbol{\mu} \in \mathbb{R}^{t-1}$ are also vectors such that $(\mathbf{y})_i = y_i$ and $(\boldsymbol{\mu})_i = \mu(x_i)$. The gram matrix is denoted by $\mathbf{K} \in \mathbb{R}^{(t-1) \times (t-1)}$ where $(\mathbf{K})_{i,j} = k(x_i, x_j)$. Let $\mu_{t,p}$ and $k_{t,p}$ denote the posterior mean and kernel for a Gaussian process with prior $p \in P$ at time t and let $\sigma_{t,p}^2(x) = k_{t,p}(x, x)$ denote the posterior variance at time t . The kernel k determines important characteristics of the functions f , see Appendix B for more details and examples.

Information gain The maximal information gain (MIG) is a measure of reduction in uncertainty of f after observing the most informative data points up to a specified size. The MIG commonly occurs in regret bounds for GP bandit algorithms (Srinivas et al., 2012; Vakili et al., 2021) and its growth rate is strongly determined by the prior kernel of the GP. Hence, we will define the MIG for any fixed GP prior $p \in P$. Let \mathbf{y}_A denote noisy observations of f at the locations $A \subset \mathcal{X}$. Then, the MIG given prior $p \in P$, $\gamma_{T,p}$, is defined as

$$\gamma_{T,p} := \sup_{A \subset \mathcal{X}, |A| \leq T} I_p(\mathbf{y}_A; f), \quad (3)$$

where $I_p(\mathbf{y}_A; f) = H(\mathbf{y}_A | p) - H(\mathbf{y}_A | f, p)$ is the mutual information between \mathbf{y}_A and f given p , and $H(\cdot)$ denotes the entropy. To aid our analysis later, we also define the worst-case MIG as $\hat{\gamma}_T := \max_{p \in P} \gamma_{T,p}$ and the average MIG as $\bar{\gamma}_T := \mathbb{E}_{p \sim P_1}[\gamma_{T,p}]$. For the RBF and Matérn kernels, $\gamma_{T,p} = \mathcal{O}(\log^{d+1}(T))$ and $\gamma_{T,p} = \mathcal{O}(T^{\frac{d}{2\nu+d}} \log^{\frac{2\nu}{2\nu+d}}(T))$ (Srinivas et al., 2012; Vakili et al., 2021).

Previous work Plenty of previous work have proposed fully Bayesian approaches that integrate the acquisition function over the hyperposterior (Osborne et al., 2009; Benassi et al., 2011; Snoek et al., 2012; Hernández-Lobato et al., 2014; Wang & Jegelka, 2017; De Ath et al., 2021; Hvarfner et al., 2023). In contrast, HP-GP-TS optimizes a single hyperposterior sample instead of expected values over the hyperposterior.

Wang & de Freitas (2014) first derived regret bounds for GP bandits with unknown lengthscale for the Expected Improvement algorithm (Moćkus, 1975). However, the proposed algorithm requires a lower bound on the lengthscale and the regret bound depends on the worst-case MIG. Later work by Berkenkamp et al. (2019) introduced Adaptive GP-UCB (A-GP-UCB) that continually lowers the lengthscale parameter. Given a sufficiently small lengthscale, the function f lies within the reproducing kernel Hilbert space (RKHS) and the regular GP-UCB theory can be applied. However, A-GP-UCB lacks a stopping mechanism and will overexplore as the lengthscale continues to shrink. Recent work by Ziomek et al. (2025) introduced Prior-Elimination GP-UCB (PE-GP-UCB) for time-varying GP-bandits with unknown prior. Unlike the work before, the regret bound of PE-GP-UCB holds for arbitrary types of hyperparameters in the GP prior. PE-GP-UCB is doubly optimistic and selects the prior *and* arm with the highest upper confidence bound. PE-GP-UCB tracks the cumulative prediction error made by the selected priors and eliminates priors that exceed a threshold level.

Other works have introduced regret balancing algorithms that maintain a set of base learning algorithms and balance their selection frequency to achieve close to optimal regret (Abbasi-Yadkori et al., 2020; Pacchiano et al., 2020). Ziomek et al. (2024) built on this idea and introduced length-scale balancing GP-UCB which can adaptively explore smaller lengthscales but can return to longer ones, unlike A-GP-UCB.

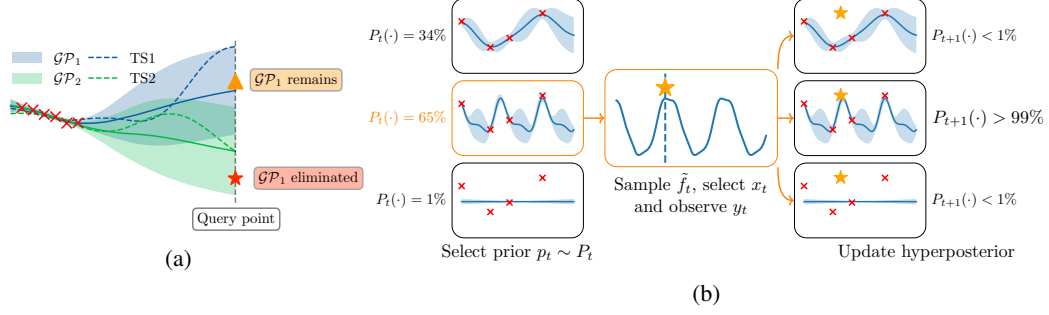


Figure 2: a) Elimination procedure of PE-GP-TS. The solid lines correspond to posterior means and the shaded regions are confidence intervals. The figure has been adapted from Ziomek et al. (2025). The dashed lines are samples from the posteriors. b) Overview of HP-GP-TS. The orange star corresponds to y_t .

The aforementioned works are based on UCB, EI, PI or regret balancing. However, another line of work has studied Thompson sampling in standard and linear bandits with unknown prior distribution (Kveton et al., 2021; Basu et al., 2021; Hong et al., 2022; Li et al., 2024). In their setting (meta or hierarchical bandits), the agent plays multiple bandit instances, either simultaneously or sequentially. The unknown means are sampled from the same (unknown) prior and by gathering knowledge across instances, the agent can solve later instances more efficiently once it has identified the prior. We emphasize that these methods have been studied only for standard stochastic and linear bandits, not for GP bandits.

3 Algorithms

As discussed by Russo & Roy (2014), TS can offer advantages over UCB algorithms for problems where constructing tight confidence bounds is difficult. In addition, Thompson sampling is often observed to perform better than UCB in practice (Chapelle & Li, 2011; Wen et al., 2015; Kandasamy et al., 2018; Åkerblom et al., 2023b,a). Motivated by this, we present two algorithms for efficient prior selection based on TS.

3.1 Prior-Elimination with Thompson sampling

Our first algorithm is an extension of PE-GP-UCB (Ziomek et al., 2025) to be employed with Thompson sampling - instead of UCB. The key difference is that instead of maximizing the upper confidence bound $U_t(x, p) = \mu_{t,p}(x) + \sqrt{\beta_t} \sigma_{t,p}(x)$ over $\mathcal{X} \times P_{t-1}$, we instead sample $\tilde{f}_{t,p}$ from the posterior $\mathcal{GP}(\mu_{t,p}, k_{t,p})$ for all priors $p \in P_{t-1}$ where P_{t-1} is the set of active priors. Then, we select the arm and prior x_t, p_t such that $x_t, p_t = \arg \max_{x,p \in \mathcal{X} \times P_{t-1}} \tilde{f}_{t,p}(x)$. Whilst PE-GP-UCB has two layers of optimism, the upper confidence bound and joint maximization of x and p , PE-GP-TS has only a single layer of optimism - which should alleviate potential overexploration issues.

The elimination procedure of PE-GP-TS is illustrated in Fig. 2. Samples $\tilde{f}_{t,p}$ are drawn from the active prior $p \in P_{t-1}$. Then, the unknown function f is queried at the selected arm x_t . If the observed value differs too much from the

Algorithm 1 Prior Elimination GP-TS (PE-GP-TS)

input Horizon T , prior functions $\{\mu_{1,p}, k_{1,p}\}_{p \in P}$, confidence parameters $\{\beta_t\}_{t=1}^T$ and $\{\xi_t\}_{t=1}^T$.

- 1: $P_1 = P, S_{0,p} = \emptyset \forall p \in P$
- 2: **for** $t = 1, 2, \dots, T$ **do**
- 3: Sample $\tilde{f}_{t,p} \sim \mathcal{GP}(\mu_{t,p}, k_{t,p}) \forall p \in P_t$
- 4: Set $x_t, p_t = \arg \max_{x,p \in \mathcal{X} \times P_{t-1}} \tilde{f}_{t,p}(x)$
- 5: $S_{t,p_t} = S_{t-1,p_t} \cup \{t\}$ and $S_{t,p} = S_{t-1,p}$ for $p \in P \setminus \{p_t\}$
- 6: Observe $y_t = f(x_t) + \epsilon_t$
- 7: Set $\eta_t = y_t - \mu_{t,p_t}(x_t)$
- 8: Set $V_t = \sqrt{\xi_t |S_{t,p_t}|} + \sum_{i \in S_{t,p_t}} \sqrt{\beta_i} \sigma_{i,p_t}(x_i)$
- 9: **if** $\left| \sum_{i \in S_{t,p_t}} \eta_i \right| > V_t$ **then**
- 10: $P_{t+1} = P_t \setminus \{p_t\}$
- 11: **else**
- 12: $P_{t+1} = P_t$

prediction made by the selected prior, then the selected prior is eliminated. Otherwise, it remains active.

The PE-GP-TS algorithm is presented in Algorithm 1. Similar to PE-GP-UCB, the set $S_{t,p}$ is used to store the time steps where prior p was selected up to and including time t . When prior p_t is selected, the prediction error $\eta_t = y_t - \mu_{t,p_t}(x_t)$ between the observed and predicted value made by the prior p_t is computed. If the sum of prediction errors made by the prior p_t exceeds the threshold value V_t , then p_t is eliminated from the active priors P_t , see line 9. Note that at time step t , only the selected prior p_t can be eliminated. As such, if a prior is very pessimistic it may never be selected and therefore will never be eliminated. Thus, the final set of active priors P_T should be viewed as non-eliminated priors rather than necessarily being reasonable priors.

3.2 HyperPrior Thompson sampling

In our first algorithm, we removed one layer of optimism. In our second algorithm, we adopt a fully Bayesian algorithm by using a hyperposterior sampling scheme where both the prior and the mean function are sampled from their respective posteriors. By shedding the optimism over the selected prior p_t , HP-GP-TS should be able avoid costly exploration by selecting likely priors instead of optimistic ones.

The algorithm is visualized in Fig. 2 and presented with more details in Algorithm 2. In the first step, the current prior p_t is sampled from the hyperposterior P_{t-1} . Then, a single sample \tilde{f}_t is taken from the selected posterior $\mathcal{GP}(\mu_{t,p_t}, k_{t,p_t})$ and is used to select the current arm: $x_t = \arg \max_{x \in \mathcal{X}} \tilde{f}_t(x)$. After observing y_t , the hyperposterior is updated by computing the likelihood of y_t under the different priors. Note that since the set of priors P is finite, computing the posterior is tractable albeit computationally costly with a complexity of $\mathcal{O}(t^3|P|)$. The likelihood $\mathbb{P}(y_t|x_t, \{x_i, y_i\}_{i=1}^{t-1}, p) = \mathcal{N}(y_t; \mu_{t,p}(x_t), \sigma_{t,p}^2(x_t) + \sigma^2)$ is simply the Gaussian likelihood of the posterior at x_t plus Gaussian noise with variance σ^2 .

Algorithm 2 HyperPrior GP-TS (HP-GP-TS)

input Horizon T , prior functions $\{\mu_{1,p}, k_{1,p}\}_{p \in P}$, hyperprior P_1 .
1: **for** $t = 1, 2, \dots, T$ **do**
2: Sample $p_t \sim P_t$
3: Sample $\tilde{f}_t \sim \mathcal{GP}(\mu_{t,p_t}, k_{t,p_t})$
4: Set $x_t = \arg \max_x \tilde{f}_t$
5: Observe $y_t = f(x_t) + \epsilon_t$
6: Set $P_{t+1}(p) \propto \mathbb{P}(y_t|x_t, \{x_i, y_i\}_{i=1}^{t-1}, p) \cdot P_t(p)$
 ▷ Update hyperposterior

4 Regret analysis

In this section, we analyze the regret for the proposed algorithms. Recall from the problem statement that we consider two slightly different settings for the two algorithms. Specifically, for PE-GP-TS we assume the unknown prior p^* is selected arbitrarily from P whilst for HP-GP-TS we assume that the unknown prior p^* is selected from a known hyperprior distribution P_1 .

Ziomek et al. (2025) structured the proof of the regret bound of PE-GP-UCB into 4 larger steps; First, showing that p^* is never eliminated with high probability. Second, establishing a bound on the simple regret. Third, bounding the cumulative regret. Finally, the cumulative bound is re-expressed in terms of the worst-case MIG. For PE-GP-TS, we establish a new bound on the simple regret and then adapt the steps of Ziomek et al. to accommodate the new simple regret bound.

To bound the simple regret, we require two concentration inequalities to hold for both the posteriors and the posterior samples which we present in the following lemma.

Lemma 4.1. *If $f(x) \sim \mathcal{GP}(\mu_{1,p^*}, k_{1,p^*})$ and $\beta_t = 2 \log \left(\frac{|\mathcal{X}||P|\pi^2 t^2}{3\delta} \right)$. Then, with probability at least $1 - \delta$, the following holds for all $t, x, p \in [T] \times \mathcal{X} \times P$:*

$$|f(x) - \mu_{t,p^*}(x)| \leq \sqrt{\beta_t \sigma_{t,p^*}(x)}, \quad (4)$$

$$|\tilde{f}_{t,p}(x) - \mu_{t,p}(x)| \leq \sqrt{\beta_t \sigma_{t,p}(x)}. \quad (5)$$

All proofs can be found in Appendix A. Lemma 4.1 is based on Lemma 5.1 of Srinivas et al. (2012) but adapted to TS by specifying that it holds for any sequence of x_1, \dots, x_T , as discussed by Russo

& Roy (2014). Additionally, we add Eq. (5) which can be shown through the same steps and an additional union bound over P . Next, we state our bound for the simple regret of PE-GP-TS.

Lemma 4.2. *If the event of Lemma 4.1 holds, then the following holds for the simple regret of PE-GP-TS for all $t \in [T]$:*

$$f(x^*) - f(x_t) \leq 2\sqrt{\beta_t \sigma_{t,p^*}(x^*)} + \sqrt{\beta_t \sigma_{t,p_t}(x_t)} - \eta_t + \epsilon_t. \quad (6)$$

Compared to the simple regret bound for PE-GP-UCB, we obtain the additional term $2\sqrt{\beta_t \sigma_{t,p^*}(x^*)}$. Since p^* is fixed, the sum over t of the new term is $\mathcal{O}(\sqrt{T\beta_T\gamma_T})$ and we obtain the following regret bound:

Theorem 4.3. *If $p^* \in P$ and $f \sim \mathcal{GP}(\mu_{1,p^*}, k_{1,p^*})$, then PE-GP-TS with confidence parameters $\beta_t = 2\log(2|\mathcal{X}||P|\pi^2 t^2/3\delta)$ and $\xi_t = 2\sigma^2 \log(|P|\pi^2 t^2/3\delta)$, satisfies the following regret bound with probability at least $1 - \delta$:*

$$R(T) \leq 2|P|B_{p^*} + 2\sqrt{\xi_T|P|T} + 2\sqrt{CT\beta_T\hat{\gamma}_T|P|} + 2\sqrt{CT\beta_T \sum_{t \in [T]} \sigma_{t,p^*}^2(x^*)} \quad (7)$$

where $B_{p^*} = \beta_1 + \sup_{x \in \mathcal{X}} |\mu_{1,p^*}(x)|$ and $C = 2/\log(1 + \sigma^{-2})$.

The bound of the first three terms is of order $\mathcal{O}(\sqrt{T\beta_T\hat{\gamma}_T})$ w.r.t. T which matches that of PE-GP-UCB. To our knowledge, the best lower bound for standard GP bandits in the Bayesian setting, where f is sampled from a GP, is $\Omega(\sqrt{T})$ for $d = 1$ (Scarlett, 2018). This would suggest that our bound is tight up to a factor $\mathcal{O}(\sqrt{\beta_T\hat{\gamma}_T})$.

We analyze the regret of HP-GP-TS using both the UCB-based framework of Russo & Roy (2014) and the information-theoretic framework of Russo & Roy (2016). First, note that HP-GP-TS inherits the probability matching property of GP-TS that $x_t|H_t \stackrel{d}{=} x^*|H_t$ where $\stackrel{d}{=}$ denotes equal in distribution. In addition, $p_t|H_t \stackrel{d}{=} p^*|H_t$ since p_t is sampled from the posterior distribution of p^* . In the UCB-based framework, we decompose the regret into three terms:

$$\mathbb{E}[f(x^*) - f(x_t)] = \mathbb{E}\left[\underbrace{f(x^*) - U_{t,p^*}(x^*)}_{(1)} + \underbrace{U_{t,p^*}(x^*) - U_{t,p^*}(x_t)}_{(2)} + \underbrace{U_{t,p^*}(x_t) - f(x_t)}_{(3)}\right] \quad (8)$$

where $U_{t,p}(x) = \mu_{t,p}(x) + \sqrt{\beta_t \sigma_{t,p}(x)}$. Summing over t , the first term can be bounded by a constant whilst the third term is bounded by $\sqrt{CT\beta_T\hat{\gamma}_T}$. Together, these two terms match the Bayesian regret bound for GP-TS with known prior. For the second term, one can utilize that $p^*, x^*|H_t \stackrel{d}{=} p_t, x_t|H_t$ and $p^*, x_t|H_t \stackrel{d}{=} p_t, x^*|H_t$ to re-express it as $\mathbb{E}[U_{t,p^*}(x^*) - U_{t,p_t}(x^*)]$. Hence, the second term can be seen as the cost of learning the true prior. If the priors are sufficiently different, one would intuitively expect $P_t(p^*) \rightarrow 1$ quickly. However, if the priors are similar then the upper-confidence bounds of the true and selected prior will not differ significantly.

Theorem 4.4. *If $p^* \sim P_1$, $f \sim \mathcal{GP}(\mu_{1,p^*}, k_{1,p^*})$ and $\beta_t = 2\log(|\mathcal{X}|t^2/\sqrt{2\pi})$, then the Bayesian regret of HP-GP-TS is bounded by*

$$BR(T) \leq \pi^2/6 + \sum_{t \in [T]} \mathbb{E}[U_{t,p^*}(x^*) - U_{t,p_t}(x^*)] + \sqrt{CT\beta_T\hat{\gamma}_T}. \quad (9)$$

Unlike PE-GP-TS and PE-GP-UCB, the regret bound of the first and third term for HP-GP-TS depends on the average MIG $\sqrt{\hat{\gamma}_T}$ rather than the worst case $\sqrt{|P|\hat{\gamma}_T}$ which can impact the theoretical regret significantly if the complexity of the priors differ and the prior is weighted towards simple priors. This is reasonable since the elimination methods assume arbitrary selection of p^* as opposed to sampling from a hyperprior. If the hyperprior is deterministic then the regret bound for HP-GP-TS matches that of GP-TS up to a factor $\mathcal{O}(\sqrt{\log T})$ (Takeno et al., 2024) and $\hat{\gamma}_T$ would be equal to the worst case $\hat{\gamma}_T$. Again, using the lower bound of Scarlett (2018), our upper bound would be tight up to a factor of $\mathcal{O}(\sqrt{\beta_T\hat{\gamma}_T})$.

The information-theoretic framework of Russo & Roy (2016) can be applied generally if the probability matching property is satisfied and the rewards are subgaussian (Vashishtha & Maillard, 2025).

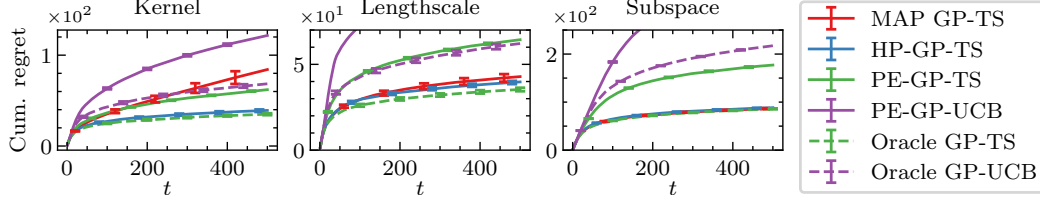


Figure 3: Cumulative regret for synthetic experiments with varying kernel (left), lengthscale (center) and mean function (right). The average final regret for PE-GP-UCB is 116.5 and 389.0 in the lengthscale and subspace experiments. Errorbars correspond to ± 1 standard error.

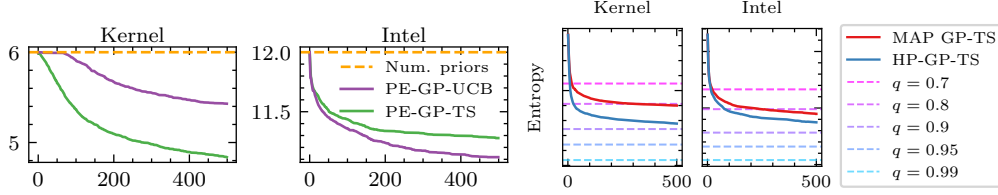


Figure 4: Mean number of priors remaining in P_t over time for PE-GP-UCB and -TS (left). Mean entropy in the hyperposterior P_t over time for HP- and MAP GP-TS (right). The dashed reference values correspond to entropies of discrete distributions with prob. q on one choice and prob. $\frac{1-q}{|P|-1}$ on the other $|P| - 1$ choices.

Lemma 4.5. Fix $x \in \mathcal{X}$ and let $Z = f(x) + \epsilon_t - \mathbb{E}_t[f(x)]$. If $p^* \sim P_1$, $f \sim \mathcal{GP}(\mu_{1,p^*}, k_{1,p^*})$ with $k_{1,p^*}(\cdot, \cdot) \leq \sigma_0^2$, then $Z|H_t$ is $\sqrt{\sigma_0^2 + \sigma^2}$ -subgaussian.

Theorem 4.6. If $p^* \sim P_1$, $f \sim \mathcal{GP}(\mu_{1,p^*}, k_{1,p^*})$, then the Bayesian regret of HP-GP-TS is bounded by

$$BR(T) \leq \sqrt{2|\mathcal{X}| \log(|\mathcal{X}|)(\sigma_0^2 + \sigma^2)T}. \quad (10)$$

The proof of Theorem 4.6 follows the proof in section D.2 of Russo & Roy (2016) with subgaussian noise. For completeness, we provide a proof using the GP-bandit notation in Algorithm 2. The information-theoretic regret bound is $\mathcal{O}(\sqrt{T})$ which improves upon the $\mathcal{O}(\sqrt{T\beta_T\gamma_T})$ obtained previously and would match the lower bound of Scarlett (2018) in terms of T . However, we also obtain a $\mathcal{O}(\sqrt{|\mathcal{X}| \log |\mathcal{X}|})$ dependency.

5 Experiments

Synthetic experiments We consider three synthetic setups with different choices of priors in P . For the first setup, the priors have varying kernels selected according to one of the following kernels: i) RBF kernel, Matérn kernel with $\nu = 5/2$, ii) Matérn kernel with $\nu = 3/2$, iii) periodic kernel with period $\rho = 5$, iv) linear kernel with $v = 0.05^2$, and, v) the rational quadratic kernel with $\alpha = 0.5$. All kernels use a lengthscale of 1.0 and are scaled s.t. $k(x, \tilde{x}) \leq 1$. In addition, the mean function for all priors are zero everywhere. For the second setup, the priors use the RBF kernel with lengthscales 4, 2, 1 or 1/2. For the third setup, the total dimensions $d = 16$ but each prior p_i assumes $f(x)$ depends on $d_s = 4$ subdimensions: $[i, i+1, i+2, i+3]$ for $i \in [5]$. Dimensions larger than 5 are wrapped around 1, i.e. $(j \bmod 5) + 1$, such that the priors are equally difficult to distinguish and optimize. All priors use the RBF kernel with lengthscale $\ell = 8$. For all three setups, the true prior p^* is sampled uniformly from P , the noise variance $\sigma^2 = 0.25^2$, and the horizon $T = 500$. For the first two setups, 500 arms are equidistantly spaced in $[0, 20]$ and for the third 500 arms are sampled uniformly from $[0, 20]^{16}$. The prior elimination methods use $\delta = 0.05$. All models are evaluated on 500 seeds on each setup. As baselines, we use PE-GP-UCB and Maximum A Posteriori (MAP) GP-TS where MAP GP-TS is identical to HP-GP-TS except for greedily selecting p_t from the posterior: $p_t = \arg \max_p P_{t-1}(p)$ ¹.

¹Note that since the hyperprior is uniform, MAP is equivalent to discrete maximum likelihood estimation.

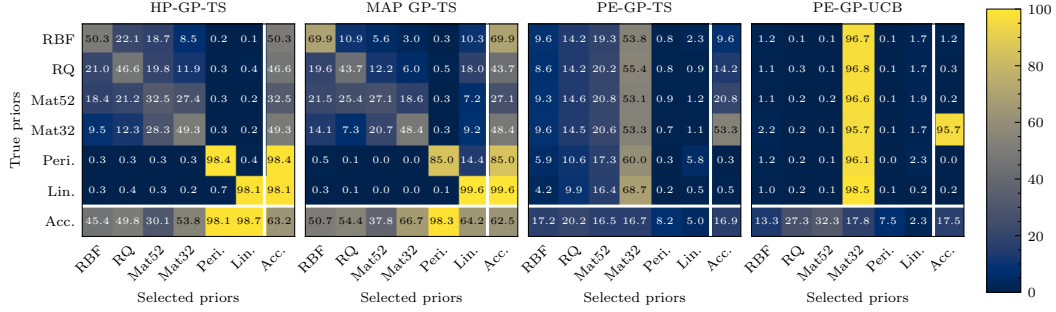


Figure 5: Confusion matrices for the true prior p^* and the selected priors p_t for the kernel experiment. Row-wise normalized to 100%.

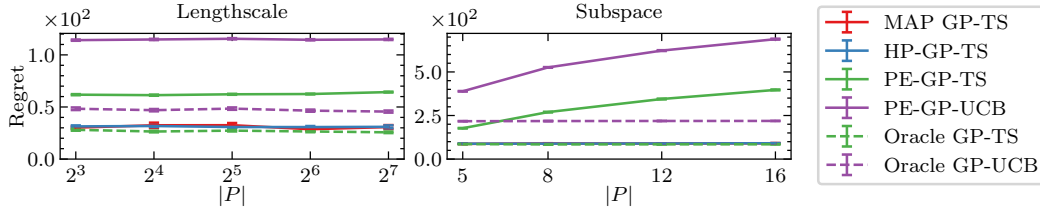


Figure 6: Total regret for the lengthscale and subspace experiments as $|P|$ increases.

Regardless of the selected prior, HP-GP-TS observes an observation y_t to update the hyperposterior P_t . Hence, greedily selecting the prior could reduce unnecessary exploration. In addition, we investigate the oracle variants of PE-GP-TS and PE-GP-UCB with $\delta = 0.05$ that are only given the true prior: $P_1 = \{p^*\}$.

The cumulative regret for the three synthetic experiments is shown in Fig. 3. Across all three experiments, we observe that HP-GP-TS has lower regret than the other methods and performs close to the oracle GP-TS. For the kernel and subspace experiments, PE-GP-TS has lower regret than the oracle GP-UCB. Hence, even if PE-GP-UCB was optimized to perform as well as the oracle, it would still not achieve the regret of our proposed methods. MAP GP-TS has slightly higher regret than HP-GP-TS for the lengthscale and subspace experiments but has significantly higher regret and variance for the kernel experiment. The greedy selection of MAP (MLE) leads to under-exploration for MAP GP-TS in certain instances.

The number of priors remaining $|P_t|$ and the hyperposterior entropy for the kernel experiment is shown in Fig. 4. The PE-methods eliminate at most one prior on average. In contrast, the hyperposterior entropy of HP-GP-TS is equivalent to 80-90% of the probability mass being assigned to one prior. HP- and MAP-GP-TS thus effectively discards priors at a much faster rate. The same pattern holds for the lengthscale and subspace experiments, see Figs. 9 and 10 in Appendix D.

In Fig. 5, we visualize how often the methods select the true prior p^* (or kernel) in the kernel experiment as confusion matrices. PE-GP-UCB selects the Matérn-3/2 kernel more than 96% of the rounds. The Matérn-3/2 kernel induces a distribution over functions that are less smooth compared to the other kernels and produces much higher confidence intervals outside the observed data leading to excessive optimistic exploration. PE-GP-TS also shows a bias towards the Matérn-3/2 kernel but does not select it as frequently as PE-GP-UCB - demonstrating that one layer of optimism has been removed. The overall “accuracy” of the selected priors, i.e. $\sum_{t \in [T]} \mathbb{1}\{p_t = p^*\}/T$, for the elimination-based methods is around 17% in the kernel experiment compared to 62.5% and 63.2% for MAP and HP-GP-TS respectively. For HP-GP-TS, we observe that it can easily identify the periodic and linear kernels. However, the RBF, Matérn and RQ kernels are often confused with each other. These kernels do not have as easily distinguishable characteristics and are likely to produce similar posteriors even with a small amount of data. See Fig. 11 in Appendix D for confusion matrices in the lengthscale and subspace experiments.

Scaling $|P|$ We perform two experiments to understand how the regret of our algorithms scale with the number of priors. In both experiments, the average difficulty of the problem is kept constant

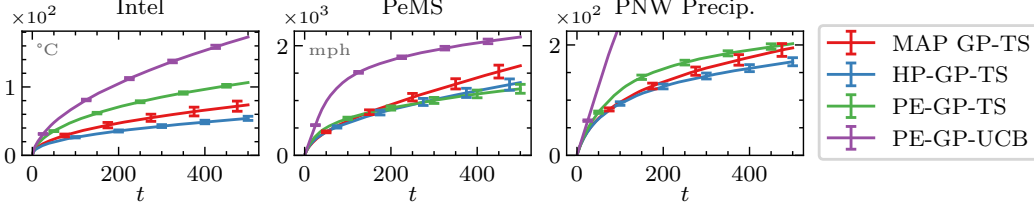


Figure 7: Cumulative regret on Intel temperature data (top) and PeMS data (bottom). Errorbars correspond to ± 1 standard error. The average final regret for PE-GP-UCB is 511.9 for PNW.

such that the regret of the oracle models is constant. In the first experiment, we increase the discretization of the lengthscale values. The lengthscales are equidistantly spaced in $[0.5, 4]$ with $|P| \in \{8, 16, 32, 64, 128\}$. As $|P|$ increases, the difference between similar priors is reduced. In the second experiment, we increase the number of priors in the subspace experiment from 5 up to 16. Each prior can share at most $3/4$ dimensions with other priors which ensures the priors remain meaningfully different. The total regret as the number of priors increases is shown in Fig. 6. For the lengthscale experiment, increasing the number of priors above 8 does not affect the regret for any algorithm. This is likely due to the redundancy in the priors. However in the subspace experiment, the regret of the prior elimination algorithms scales approximately as $\sqrt{|P|}$ whilst MAP- and HP-GP-TS are consistently close to the constant regret of the oracle.

Real-world data We perform three experiments with real-world data from the Intel Berkeley dataset (Madden et al., 2004), California Performance Measurement System (PeMS) (Chen et al., 2001; California Department of Transportation, 2024) and Pacific Northwest (PNW) daily precipitation dataset (Widmann & Bretherton, 1999, 2000). Each dataset contains measurements from a set of sensors over time. We split each dataset into a training and test set where the test set contains the last third of the data. Hence, the distribution of the test data may have shifted slightly from the training data. The training sets are split further into separate buckets to define our priors. For each bucket p , we compute the empirical mean $\hat{\mu}_p$ and covariance $\hat{\Sigma}_p$ which defines the prior $\mathcal{GP}(\hat{\mu}_p, \hat{\Sigma}_p)$. The buckets in the Intel data corresponds to the 12 days in the training dataset. For the PeMS data, each hour between 06:00 and 13:00 defines one prior, giving 7 priors. For the daily precipitation data, each month in the year constitutes a prior, yielding 12 priors. When running the experiments, we select a measurement of all sensors from the test data uniformly at random. The selected measurements correspond to the unknown function $f(x)$ where x is the sensor index and the goal is then to identify sensors measuring large temperatures, small speeds or high precipitation respectively for the three datasets. When the algorithms select an arm to evaluate, we add Gaussian noise to y_t with variance σ^2 around 5% of the signal variance, similar to Srinivas et al. (2012); Bogunovic et al. (2016). See Appendix C for more details about the experimental setup.

The cumulative regret for the experiments with real-world data is presented in Fig. 7. For the Intel and PNW data, HP- and MAP GP-TS obtain the lowest and second lowest cumulative regret respectively. HP- and MAP GP-TS have lower regret than PE-GP-TS initially but PE-GP-TS catches up and has the lowest total regret. To understand this better, we visualize quantiles of the total regret in Fig. 13. MAP- and HP-GP-TS have the lowest median regret for all three experiments and hence perform best in a majority of instances. However, the 90th and 95th quantiles are considerably larger for the PeMS data which impacts the average regret significantly. Hence, for the PeMS data, the prior elimination methods seem to yield more stable results.

In Fig. 4, the number of priors remaining in $|P_t|$ and the hyperposterior entropy is shown for the Intel experiment. Similar to the synthetic experiment, on average, the prior elimination methods eliminate less than 1 prior whereas the hyperposteriors of HP- and MAP GP-TS concentrate the equivalent of 80-90% of the probability mass to one prior. The results for PeMS and PNW experiment are shown in Figs. 9 and 10 in Appendix D. Here, effectively no priors are eliminated. For the PNW experiment, the hyperposteriors do not concentrate as much compared to the other experiments. This could indicate that knowing the exact prior is not as important for the PNW data.

6 Discussion

Limitations The main theoretical limitation is that the regret bound of PE-GP-TS and the UCB-based regret bound of HP-GP-TS each contain a term that we are unable to bound sublinearly. The main computational limitation of all the methods we have considered is that they require the set of priors P to be discrete. In addition, their computational cost scales linearly with the number of priors considered. In theory, the elimination-based methods could have lower computational cost compared to HP-GP-TS. However, in practice, priors are rarely eliminated as shown by Fig. 9. At best, one sixth of the priors is eliminated in the kernel experiment for PE-GP-TS. Since only priors that are selected can be eliminated, the confidence parameters β_t , ξ_t increase over time and ξ_t increases with $|P|$, including more priors could lead to less priors being eliminated overall.

7 Conclusion

In this paper, we have proposed two algorithms for joint prior selection and regret minimization in GP bandits based on GP-TS. We have analyzed the algorithms theoretically and have experimentally evaluated both algorithms on both synthetic and real-world data. We find that they both select the true prior more often and obtain lower regret than previous work due to lowering the amount of optimistic exploration.

Acknowledgments

The work of Jack Sandberg and Morteza Haghiri Chehreghani was partially supported by the Wallenberg AI, Autonomous Systems and Software Program (WASP) funded by the Knut and Alice Wallenberg Foundation. The computations were enabled by resources provided by the National Academic Infrastructure for Supercomputing in Sweden (NAISS), partially funded by the Swedish Research Council through grant agreement no. 2022-06725.

References

- Abbasi-Yadkori, Y., Pacchiano, A., and Phan, M. Regret Balancing for Bandit and RL Model Selection, June 2020. URL <https://arxiv.org/abs/2006.05491>.
- Åkerblom, N., Chen, Y., and Haghiri Chehreghani, M. Online learning of energy consumption for navigation of electric vehicles. *Artificial Intelligence*, 317:103879, April 2023a. doi: 10.1016/j.artint.2023.103879.
- Åkerblom, N., Hoseini, F. S., and Haghiri Chehreghani, M. Online learning of network bottlenecks via minimax paths. *Machine Learning*, 112(1):131–150, January 2023b. doi: 10.1007/s10994-022-06270-0.
- Basu, S., Kveton, B., Zaheer, M., and Szepesvari, C. No Regrets for Learning the Prior in Bandits. In *Advances in Neural Information Processing Systems*, volume 34, pp. 28029–28041. Curran Associates, Inc., 2021.
- Benassi, R., Bect, J., and Vazquez, E. Robust Gaussian Process-Based Global Optimization Using a Fully Bayesian Expected Improvement Criterion. In Coello, C. A. C. (ed.), *Learning and Intelligent Optimization*, volume 6683, pp. 176–190. Springer Berlin Heidelberg, Berlin, Heidelberg, 2011. ISBN 978-3-642-25565-6 978-3-642-25566-3. doi: 10.1007/978-3-642-25566-3_13. URL http://link.springer.com/10.1007/978-3-642-25566-3_13. Series Title: Lecture Notes in Computer Science.
- Berkenkamp, F., Schoellig, A. P., and Krause, A. No-Regret Bayesian Optimization with Unknown Hyperparameters. *Journal of Machine Learning Research*, 20(50):1–24, 2019. ISSN 1533-7928.
- Bogunovic, I., Scarlett, J., and Cevher, V. Time-Varying Gaussian Process Bandit Optimization. In *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, pp. 314–323. PMLR, May 2016. URL <https://proceedings.mlr.press/v51/bogunovic16.html>. ISSN: 1938-7228.

- California Department of Transportation. Caltrans performance measurement system, 2024. URL <https://pems.dot.ca.gov/>.
- Chapelle, O. and Li, L. An Empirical Evaluation of Thompson Sampling. In *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc., 2011.
- Chen, C., Petty, K., Skabardonis, A., Varaiya, P., and Jia, Z. Freeway Performance Measurement System: Mining Loop Detector Data. *Transportation Research Record*, 1748(1):96–102, January 2001. doi: 10.3141/1748-12.
- De Ath, G., Everson, R. M., and Fieldsend, J. E. How Bayesian should Bayesian optimisation be? In *Proceedings of the Genetic and Evolutionary Computation Conference Companion, GECCO '21*, pp. 1860–1869, New York, NY, USA, July 2021. Association for Computing Machinery. ISBN 978-1-4503-8351-6. doi: 10.1145/3449726.3463164. URL <https://dl.acm.org/doi/10.1145/3449726.3463164>.
- Gonzalez, J., Lezmi, E., Roncalli, T., and Xu, J. Financial Applications of Gaussian Processes and Bayesian Optimization, 2019. URL <https://arxiv.org/abs/1903.04841>.
- Griffiths, R.-R. and Hernández-Lobato, J. M. Constrained Bayesian optimization for automatic chemical design using variational autoencoders. *Chemical Science*, 11(2):577–586, 2020. doi: 10.1039/C9SC04026A.
- Hernández-Lobato, J. M., Requeima, J., Pyzer-Knapp, E. O., and Aspuru-Guzik, A. Parallel and Distributed Thompson Sampling for Large-scale Accelerated Exploration of Chemical Space. In *Proceedings of the 34th International Conference on Machine Learning*, pp. 1470–1479. PMLR, July 2017.
- Hernández-Lobato, J. M., Hoffman, M. W., and Ghahramani, Z. Predictive Entropy Search for Efficient Global Optimization of Black-box Functions. In *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014. URL https://proceedings.neurips.cc/paper_files/paper/2014/hash/6488484c982e9af5c35689523ba1abfe-Abstract.html.
- Hong, J., Kveton, B., Zaheer, M., and Ghavamzadeh, M. Hierarchical Bayesian Bandits. In *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, pp. 7724–7741. PMLR, May 2022.
- Hvarfner, C., Hellsten, E., Hutter, F., and Nardi, L. Self-Correcting Bayesian Optimization through Bayesian Active Learning. *Advances in Neural Information Processing Systems*, 36:79173–79199, December 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/hash/fa55bf1947530fc9567059ff42a806c2-Abstract-Conference.html.
- Jiang, B., Berliner, M. D., Lai, K., Asinger, P. A., Zhao, H., Herring, P. K., Bazant, M. Z., and Braatz, R. D. Fast charging design for Lithium-ion batteries via Bayesian optimization. *Applied Energy*, 307:118244, February 2022. doi: 10.1016/j.apenergy.2021.118244.
- Kandasamy, K., Krishnamurthy, A., Schneider, J., and Poczos, B. Parallelised Bayesian Optimisation via Thompson Sampling. In *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, pp. 133–142. PMLR, March 2018.
- Krause, A., Singh, A., and Guestrin, C. Near-Optimal Sensor Placements in Gaussian Processes: Theory, Efficient Algorithms and Empirical Studies. *Journal of Machine Learning Research*, 9(8): 235–284, 2008. ISSN 1533-7928. URL <http://jmlr.org/papers/v9/krause08a.html>.
- Kveton, B., Konobeev, M., Zaheer, M., Hsu, C.-W., Mladenov, M., Boutilier, C., and Szepesvari, C. Meta-Thompson Sampling. In *Proceedings of the 38th International Conference on Machine Learning*, pp. 5884–5893. PMLR, July 2021.
- Li, H., Liang, D., and Xie, Z. Modified Meta-Thompson Sampling for Linear Bandits and Its Bayes Regret Analysis, September 2024. URL <https://arxiv.org/abs/2409.06329>.
- Mackay, D. J. C. Introduction to Gaussian processes. In *NATO ASI Series. Series F : Computer and System Sciences*, pp. 133–165, 1998. ISBN 978-3-540-64928-1.

- Madden, S. et al. Intel lab data, 2004. URL <https://db.csail.mit.edu/labdata/labdata.html>.
- Matérn, B. Spatial Variation. In Brillinger, D., Fienberg, S., Gani, J., Hartigan, J., and Krickeberg, K. (eds.), *Spatial Variation*, volume 36 of *Lecture Notes in Statistics*. Springer, New York, NY, 1986. ISBN 978-0-387-96365-5 978-1-4615-7892-5. doi: 10.1007/978-1-4615-7892-5.
- Moćkus, J. On bayesian methods for seeking the extremum. In Marchuk, G. I. (ed.), *Optimization Techniques IFIP Technical Conference Novosibirsk, July 1–7, 1974*, pp. 400–404, Berlin, Heidelberg, 1975. Springer. ISBN 978-3-540-37497-8. doi: 10.1007/3-540-07165-2_55.
- Nuara, A., Trovò, F., Gatti, N., and Restelli, M. A Combinatorial-Bandit Algorithm for the Online Joint Bid/Budget Optimization of Pay-per-Click Advertising Campaigns. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1), April 2018. doi: 10.1609/aaai.v32i1.11888.
- Osborne, M. A., Garnett, R., and Roberts, S. J. Gaussian processes for global optimization. 2009.
- Pacchiano, A., Dann, C., Gentile, C., and Bartlett, P. Regret Bound Balancing and Elimination for Model Selection in Bandits and RL, December 2020. URL <https://arxiv.org/abs/2012.13045>.
- Pyzer-Knapp, E. O. Bayesian optimization for accelerated drug discovery. *IBM Journal of Research and Development*, 62(6):2:1–2:7, November 2018. doi: 10.1147/JRD.2018.2881731.
- Russo, D. and Roy, B. V. Learning to Optimize via Posterior Sampling. *Mathematics of Operations Research*, April 2014. doi: 10.1287/moor.2014.0650.
- Russo, D. and Roy, B. V. An Information-Theoretic Analysis of Thompson Sampling. *Journal of Machine Learning Research*, 17(68):1–30, 2016. ISSN 1533-7928. URL <http://jmlr.org/papers/v17/14-087.html>.
- Sandberg, J., Åkerblom, N., and Chehreghani, M. H. Bayesian Analysis of Combinatorial Gaussian Process Bandits. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Scarlett, J. Tight Regret Bounds for Bayesian Optimization in One Dimension. In *Proceedings of the 35th International Conference on Machine Learning*, pp. 4500–4508. PMLR, July 2018.
- Snoek, J., Larochelle, H., and Adams, R. P. Practical Bayesian Optimization of Machine Learning Algorithms. In *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc., 2012. URL https://papers.nips.cc/paper_files/paper/2012/hash/05311655a15b75fab86956663e1819cd-Abstract.html.
- Srinivas, N., Krause, A., Kakade, S. M., and Seeger, M. W. Information-Theoretic Regret Bounds for Gaussian Process Optimization in the Bandit Setting. *IEEE Transactions on Information Theory*, 58(5):3250–3265, May 2012. doi: 10.1109/TIT.2011.2182033.
- Takeno, S., Inatsu, Y., Karasuyama, M., and Takeuchi, I. Posterior Sampling-Based Bayesian Optimization with Tighter Bayesian Regret Bounds. In *Proceedings of the 41st International Conference on Machine Learning*. PMLR, June 2024.
- Turner, R., Eriksson, D., McCourt, M., Kiili, J., Laaksonen, E., Xu, Z., and Guyon, I. Bayesian Optimization is Superior to Random Search for Machine Learning Hyperparameter Tuning: Analysis of the Black-Box Optimization Challenge 2020. In *Proceedings of the NeurIPS 2020 Competition and Demonstration Track*, pp. 3–26. PMLR, August 2021.
- Vakili, S., Khezeli, K., and Picheny, V. On Information Gain and Regret Bounds in Gaussian Process Bandits. In *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, pp. 82–90. PMLR, March 2021.
- Vashishtha, S. and Maillard, O.-A. Leveraging priors on distribution functions for multi-arm bandits. *Reinforcement Learning Journal*, 2025.
- Wang, Z. and de Freitas, N. Theoretical Analysis of Bayesian Optimisation with Unknown Gaussian Process Hyper-Parameters, June 2014. URL <https://arxiv.org/abs/1406.7758>.

- Wang, Z. and Jegelka, S. Max-value Entropy Search for Efficient Bayesian Optimization. In *Proceedings of the 34th International Conference on Machine Learning*, pp. 3627–3635. PMLR, July 2017. URL <https://proceedings.mlr.press/v70/wang17e.html>. ISSN: 2640-3498.
- Wen, Z., Kveton, B., and Ashkan, A. Efficient Learning in Large-Scale Combinatorial Semi-Bandits. In *Proceedings of the 32nd International Conference on Machine Learning*, pp. 1113–1122. PMLR, June 2015.
- Widmann, M. and Bretherton, C. S. "50" km resolution daily precipitation for the Pacific Northwest, 1949-94, May 1999. URL <http://research.jisao.washington.edu/data/widmann/>.
- Widmann, M. and Bretherton, C. S. Validation of Mesoscale Precipitation in the NCEP Reanalysis Using a New Gridcell Dataset for the Northwestern United States. *Journal of Climate*, 13(11): 1936–1950, June 2000. ISSN 0894-8755, 1520-0442. doi: 10.1175/1520-0442(2000)013<1936:VOMPIT>2.0.CO;2. URL https://journals.ametsoc.org/view/journals/clim/13/11/1520-0442_2000_013_1936_vompit_2.0.co_2.xml. Publisher: American Meteorological Society Section: Journal of Climate.
- Williams, C. K. and Rasmussen, C. E. *Gaussian Processes for Machine Learning*, volume 2. MIT press Cambridge, MA, 2006.
- Ziomek, J., Adachi, M., and Osborne, M. A. Bayesian Optimisation with Unknown Hyperparameters: Regret Bounds Logarithmically Closer to Optimal. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, November 2024.
- Ziomek, J., Adachi, M., and Osborne, M. A. Time-varying Gaussian Process Bandits with Unknown Prior. In *The 28th International Conference on Artificial Intelligence and Statistics*, February 2025.

A Proofs

In the following section, we state and prove the results shown in the main text.

A.1 PE-GP-TS

First, we state and prove concentration inequalities for $f(x)$ and $\tilde{f}_{t,p}(x)$.

Lemma 4.1. *If $f(x) \sim \mathcal{GP}(\mu_{1,p^*}, k_{1,p^*})$ and $\beta_t = 2 \log \left(\frac{|\mathcal{X}||P|\pi^2 t^2}{3\delta} \right)$. Then, with probability at least $1 - \delta$, the following holds for all $t, x, p \in [T] \times \mathcal{X} \times P$:*

$$|f(x) - \mu_{t,p^*}(x)| \leq \sqrt{\beta_t} \sigma_{t,p^*}(x), \quad (4)$$

$$|\tilde{f}_{t,p}(x) - \mu_{t,p}(x)| \leq \sqrt{\beta_t} \sigma_{t,p}(x). \quad (5)$$

Proof. Follows by the same steps as Lemma 5.1 of Srinivas except we condition on the complete history H_t instead of only $y_{1:t-1}$. Additionally, for Eq. (5) we must take an additional union bound over $p \in P$.

Fix $t, x, p \in [T] \times \mathcal{X} \times P$. Given the history H_t , $\tilde{f}_{t,p}(x) \sim \mathcal{N}(\mu_{t,p}(x), \sigma_{t,p}^2(x))$. Using that $\mathbb{P}(Z > c) \leq 1/2e^{-c^2/2}$ for $Z \sim \mathcal{N}(0, 1)$, we get that

$$\mathbb{P} \left(\left| \frac{\tilde{f}_{t,p}(x) - \mu_{t,p^*}(x)}{\sigma_{t,p^*}(x)} \right| > \sqrt{\beta_t} \right) \leq \exp(-\beta_t/2) \quad (11)$$

$$= \frac{3\delta}{|\mathcal{X}||P|\pi^2 t^2} \quad (12)$$

Note that $\sum_{t \geq 1} \frac{1}{t^2} = \frac{\pi^2}{6}$. By taking the union bound over \mathcal{X} , P and $t \geq 1$, Eq. (5) holds w.p. at least $1 - \delta/2$. By the same reasoning and skipping the union bound over P , Eq. (4) holds w.p. at least $1 - \delta/2$. Thus, both events hold w.p. at least $1 - \delta$. \square

Next, we state three lemmas from Ziomek et al. (2025) that are used in the proof of our regret bound.

Lemma A.1. (Lemma 5.1 of Ziomek et al. (2025)) *If $\xi_t = 2\sigma^2 \log \left(\frac{|P|\pi^2 t^2}{6\delta} \right)$, then the following holds with probability at least $1 - \delta$:*

$$\left| \sum_{i \in S_{t,p}} \epsilon_i \right| \leq \sqrt{\xi_t |S_{t,p}|} \quad \forall t, p \in [T] \times P. \quad (13)$$

Lemma A.2. (Lemma 5.2 of Ziomek et al. (2025)) *Let $B_{p^*} = \beta_1 + \sup_{x \in \mathcal{X}} |\mu_{1,p^*}(x)|$, then if μ_{1,p^*} and k_{1,p^*} satisfy $|\mu_{1,p^*}(\cdot)| < \infty$ and $k_{1,p^*}(\cdot, \cdot) \leq 1$ and Lemma 4.1 holds, then*

$$\sup_{x \in \mathcal{X}} |f(x)| \leq B_{p^*}. \quad (14)$$

Lemma A.3. (Lemma 5.3 of Ziomek et al. (2025)) *For $C = 2/\log(1 + \sigma^{-2})$, $\sum_{t \notin \mathcal{C}} \sqrt{\beta_t} \sigma_{t,p_t}(x_t) \leq \sqrt{CT\beta_T\hat{\gamma}_T|P|}$ where $\beta_T = \max_{p \in P} \beta_T$ and $\hat{\gamma}_T = \max_{p \in P} \gamma_{T,p}$.*

Then, we state and prove the new simple regret bound for PE-GP-TS.

Lemma 4.2. *If the event of Lemma 4.1 holds, then the following holds for the simple regret of PE-GP-TS for all $t \in [T]$:*

$$f(x^*) - f(x_t) \leq 2\sqrt{\beta_t} \sigma_{t,p^*}(x^*) + \sqrt{\beta_t} \sigma_{t,p_t}(x_t) - \eta_t + \epsilon_t. \quad (6)$$

Proof. First, we upper bound $f(x^*)$ as follows

$$f(x^*) \leq \mu_{t,p^*}(x^*) + \sqrt{\beta_t} \sigma_{t,p^*}(x^*) \quad (\text{Eq. (4)}) \quad (15)$$

$$\leq \tilde{f}_{t,p^*}(x^*) + 2\sqrt{\beta_t} \sigma_{t,p^*}(x^*) \quad (\text{Eq. (5)}) \quad (16)$$

$$\leq \tilde{f}_{t,p_t}(x_t) + 2\sqrt{\beta_t} \sigma_{t,p^*}(x^*). \quad (\text{TS selection rule}) \quad (17)$$

Then, we lower bound $f(x_t)$

$$f(x_t) = \mu_{t,p_t}(x_t) + \eta_t - \epsilon_t \quad (\text{Def. of } \eta_t) \quad (18)$$

$$\geq \tilde{f}_{t,p_t}(x_t) - \sqrt{\beta_t} \sigma_{t,p_t}(x_t) + \eta_t - \epsilon_t. \quad (\text{Eq. (5)}) \quad (19)$$

Combining, Eqs. (17) and (19) we obtain

$$f(x^*) - f(x_t) \leq 2\sqrt{\beta_t} \sigma_{t,p^*}(x^*) + \sqrt{\beta_t} \sigma_{t,p_t}(x_t) - \eta_t + \epsilon_t. \quad (20)$$

□

Finally, we state and prove the cumulative regret bound for PE-GP-TS.

Theorem 4.3. *If $p^* \in P$ and $f \sim \mathcal{GP}(\mu_{1,p^*}, k_{1,p^*})$, then PE-GP-TS with confidence parameters $\beta_t = 2 \log(2|\mathcal{X}||P|\pi^2 t^2 / 3\delta)$ and $\xi_t = 2\sigma^2 \log(|P|\pi^2 t^2 / 3\delta)$, satisfies the following regret bound with probability at least $1 - \delta$:*

$$R(T) \leq 2|P|B_{p^*} + 2\sqrt{\xi_T|P|T} + 2\sqrt{CT\beta_T\hat{\gamma}_T|P|} + 2\sqrt{CT\beta_T \sum_{t \in [T]} \sigma_{t,p^*}^2(x^*)} \quad (7)$$

where $B_{p^*} = \beta_1 + \sup_{x \in \mathcal{X}} |\mu_{1,p^*}(x)|$ and $C = 2/\log(1 + \sigma^{-2})$.

Proof. First, we show that the true prior p^* is never rejected if Lemmas 4.1 and A.1 hold.

$$\left| \sum_{i \in S_{t,p^*}} \eta_i \right| = \left| \sum_{i \in S_{t,p^*}} (y_i - f(x_i) + f(x_i) - \mu_{i,p^*}(x_i)) \right| \quad (21)$$

$$\leq \left| \sum_{i \in S_{t,p^*}} \epsilon_i \right| + \sum_{i \in S_{t,p^*}} |f(x_i) - \mu_{i,p^*}(x_i)| \quad (\text{Triangle ineq.}) \quad (22)$$

$$\leq \sqrt{\xi_t |S_{t,p^*}|} + \sum_{i \in S_{t,p^*}} \sqrt{\beta_i} \sigma_{i,p^*}(x_i). \quad (\text{Lemmas 4.1 and A.1}) \quad (23)$$

Next, we bound the cumulative regret. To establish a bound on the cumulative regret, we must separate out the rounds where priors are eliminated. Hence, define the set of critical iterations as

$$\mathcal{C} = \left\{ t \in [T] : \left| \sum_{i \in S_{t,p_t}} \eta_i \right| > \sqrt{\xi_t |S_{t,p_t}|} + \sum_{i \in S_{t,p_t}} \sqrt{\beta_i} \sigma_{i,p_t}(x_i) \right\}. \quad (24)$$

Using Lemma A.2 and Eq. (20), we can bound the cumulative regret as follows:

$$\text{BR}(T) = \sum_{t \in \mathcal{C}} \text{BR}_t + \sum_{t \notin \mathcal{C}} \text{BR}_t \quad (25)$$

$$\leq 2|P|B_{p^*} + \sum_{t \notin \mathcal{C}} 2\sqrt{\beta_t} \sigma_{t,p^*}(x^*) + \sum_{t \notin \mathcal{C}} \sqrt{\beta_t} \sigma_{t,p_t}(x_t) + \sum_{p \in P} \sum_{t \in S_{T,p} \setminus \mathcal{C}} (\epsilon_t - \eta_t). \quad (26)$$

where $B_{p^*} := \beta_1 + \sup_{x \in \mathcal{X}} |\mu_{1,p^*}(x)|$. If $t \notin \mathcal{C}$, line 9 in Algorithm 1 evaluates to `false` and hence

$$\sum_{p \in P} \sum_{t \in S_{T,p} \setminus \mathcal{C}} -\eta_t \leq \sum_{p \in P} \sqrt{\xi_T |S_{T,p}|} + \sum_{p \in P} \sum_{t \in S_{T,p} \setminus \mathcal{C}} \sqrt{\beta_t} \sigma_{t,p}(x_t). \quad (27)$$

Additionally, using Lemma A.1, we can bound the Gaussian noise:

$$\sum_{p \in P} \sum_{t \in S_{T,p} \setminus \mathcal{C}} \epsilon_t \leq \sum_{p \in P} \left| \sum_{t \in S_{T,p} \setminus \mathcal{C}} \epsilon_t \right| \leq \sum_{p \in P} \left| \sum_{t \in S_{T,p}} \epsilon_t \right| \quad (28)$$

$$\leq \sum_{p \in P} \sqrt{\xi_T |S_{T,p}|} \quad (\text{Lemma A.1}) \quad (29)$$

$$\leq \sqrt{\xi_T |P|T} \quad (\text{Cauchy-Schwarz}) \quad (30)$$

Combining the above, the cumulative regret is bounded by

$$\text{BR}(T) \leq 2|P|B_{p^*} + 2\sqrt{\xi_T|P|T} + 2 \sum_{t \notin \mathcal{C}} \sqrt{\beta_t \sigma_{t,p^*}(x^*)} + 2 \sum_{t \notin \mathcal{C}} \sqrt{\beta_t \sigma_{t,p_t}(x_t)}. \quad (31)$$

Finally, applying Lemma A.3, we obtain the result

$$\text{BR}(T) \leq 2|P|B_{p^*} + 2\sqrt{\xi_T|P|T} + 2\sqrt{CT\beta_T \sum_{t \in [T]} \sigma_{t,p^*}^2(x^*)} + 2\sqrt{CT\beta_T \hat{\gamma}_T|P|}. \quad (32)$$

□

A.2 UCB-analysis of HP-GP-TS

Next, we state and prove our UCB-based regret bound for HP-GP-TS.

Theorem 4.4. *If $p^* \sim P_1$, $f \sim \mathcal{GP}(\mu_{1,p^*}, k_{1,p^*})$ and $\beta_t = 2\log(|\mathcal{X}|t^2/\sqrt{2\pi})$, then the Bayesian regret of HP-GP-TS is bounded by*

$$\text{BR}(T) \leq \pi^2/6 + \sum_{t \in [T]} \mathbb{E}[U_{t,p^*}(x^*) - U_{t,p_t}(x^*)] + \sqrt{CT\beta_T \bar{\gamma}_T}. \quad (9)$$

Proof. To begin, we note that $x_t|H_t \stackrel{d}{=} x^*|H_t$ and $p_t|H_t \stackrel{d}{=} p^*|H_t$ since both x_t and p_t are sampled from their respective posteriors. Let $U_{t,p}(x) := \mu_{t,p}(x) + \sqrt{\beta_t \sigma_{t,p}(x)}$. Then, we start decomposing the instant regret into two terms:

$$\text{BR}(T) = \sum_{t \in [T]} \mathbb{E}[f(x^*) - f(x_t)] \quad (33)$$

$$= \sum_{t \in [T]} \mathbb{E}[f(x^*) - U_{t,p^*}(x^*) + U_{t,p^*}(x^*)] \quad (34)$$

$$\begin{aligned} & - U_{t,p_t}(x^*) + U_{t,p^*}(x_t) - f(x_t) \quad (x^*, p_t|H_t \stackrel{d}{=} x_t, p^*|H_t) \quad (35) \\ & = \underbrace{\sum_{t \in [T]} \mathbb{E}[f(x^*) - U_{t,p^*}(x^*)]}_{(1)} + \underbrace{\sum_{t \in [T]} \mathbb{E}[U_{t,p_t}(x_t) - U_{t,p^*}(x_t)]}_{(2)} + \underbrace{\sum_{t \in [T]} \mathbb{E}[U_{t,p^*}(x_t) - f(x_t)]}_{(3)} \end{aligned} \quad (36)$$

We begin by bounding term (1),

$$(1) = \sum_{t \in [T]} \mathbb{E}\left[f(x^*) - \mu_{t,p^*}(x^*) - \sqrt{\beta_t \sigma_{t,p^*}(x^*)}\right] \quad (37)$$

$$\leq \sum_{t \in [T]} \mathbb{E}\left[\left[f(x^*) - \mu_{t,p^*}(x^*) - \sqrt{\beta_t \sigma_{t,p^*}(x^*)}\right]_+\right] \quad ([\cdot]_+ := \max(\cdot, 0)) \quad (38)$$

$$\leq \sum_{t \in [T]} \sum_{x \in \mathcal{X}} \mathbb{E}\left[\left[f(x) - \mu_{t,p^*}(x) - \sqrt{\beta_t \sigma_{t,p^*}(x)}\right]_+\right] \quad (x^* \in \mathcal{X}, [\cdot]_+ \geq 0) \quad (39)$$

$$\leq \sum_{t \in [T]} \sum_{x \in \mathcal{X}} \mathbb{E}_{p^*, H_t} \left[\mathbb{E}_t \left[\left[f(x) - \mu_{t,p^*}(x) - \sqrt{\beta_t \sigma_{t,p^*}(x)} \right]_+ \mid p^*, H_t \right] \right] \quad (\text{Tower rule}) \quad (40)$$

Recall that for $Z \sim \mathcal{N}(\mu, \sigma)$ with $\mu \leq 0$, $\mathbb{E}[[Z]_+] = \frac{\sigma}{\sqrt{2\pi}} \exp\left(\frac{-\mu^2}{2\sigma^2}\right)$. In our case, note that $f(x)|p^*, H_t \sim \mathcal{N}(\mu_{t,p^*}(x), \sigma_{t,p^*}^2(x))$ and $-\mu_{t,p^*}(x) - \sqrt{\beta_t \sigma_{t,p^*}(x)}$ is deterministic given p^*, H_t .

Hence,

$$(1) \leq \sum_{t \in [T]} \sum_{x \in \mathcal{X}} \mathbb{E}_{p^*, H_t} \left[\frac{\sigma_{t,p^*}(x)}{\sqrt{2\pi}} \exp\left(-\frac{\beta_t}{2}\right) \right] \quad (41)$$

$$\leq \sum_{t \in [T]} \sum_{x \in \mathcal{X}} \mathbb{E}_{p^*, H_t} \left[\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{\beta_t}{2}\right) \right] \quad (\sigma_{t,p^*}(x) \leq \sigma_{0,p^*}(x) \leq 1) \quad (42)$$

$$= \sum_{t \in [T]} \sum_{x \in \mathcal{X}} \frac{1}{\sqrt{2\pi}} \exp(-\beta_t/2) \quad (43)$$

$$\leq \sum_{t \in [T]} \frac{1}{t^2} \leq \frac{\pi^2}{6}. \quad (\beta_t = 2 \log(|\mathcal{X}|t^2/\sqrt{2\pi})) \quad (44)$$

Next, we bound (3) as follows:

$$(3) = \sum_{t \in [T]} \mathbb{E} [U_{t,p^*}(x_t) - f(x_t)] \quad (45)$$

$$= \sum_{t \in [T]} \mathbb{E}_{H_t} [\mathbb{E}_t [U_{t,p^*}(x_t) - f(x_t) | H_t]] \quad (\text{Tower rule}) \quad (46)$$

$$= \sum_{t \in [T]} \mathbb{E}_{H_t} [\mathbb{E}_t [U_{t,p^*}(x_t) - \mu_{t,p^*}(x_t) | H_t]] \quad (\mathbb{E}[f(x_t) | H_t] = \mathbb{E}[\mu_{t,p^*}(x_t) | H_t]) \quad (47)$$

$$= \sum_{t \in [T]} \mathbb{E}_{H_t} [\mathbb{E}_t [\sqrt{\beta_t} \sigma_{t,p^*}(x_t) | H_t]] \quad (U_{t,p^*}(\cdot) = \mu_{t,p^*}(\cdot) + \sqrt{\beta_t} \sigma_{t,p^*}(\cdot)) \quad (48)$$

$$= \sum_{t \in [T]} \mathbb{E} [\sqrt{\beta_t} \sigma_{t,p^*}(x_t)] \quad (49)$$

Continuing,

$$(3) = \mathbb{E} \left[\sum_{t \in [T]} \sqrt{\beta_t} \sigma_{t,p^*}(x_t) \right] \quad (50)$$

$$\leq \mathbb{E} \left[\sqrt{\sum_{t \in [T]} \beta_t \sum_{t \in [T]} \sigma_{t,p^*}^2(x_t)} \right] \quad (\text{Cauchy-Schwarz}) \quad (51)$$

$$= \sqrt{\sum_{t \in [T]} \beta_t} \mathbb{E} \left[\sqrt{\sum_{t \in [T]} \sigma_{t,p^*}^2(x_t)} \right] \quad (\beta_t \text{ deterministic}) \quad (52)$$

$$\leq \sqrt{\sum_{t \in [T]} \beta_t} \sqrt{\mathbb{E} \left[\sum_{t \in [T]} \sigma_{t,p^*}^2(x_t) \right]} \quad (\text{Jensen's inequality}) \quad (53)$$

$$\leq \sqrt{\beta_T T} \sqrt{\mathbb{E}_{p^*} \left[\mathbb{E} \left[\sum_{t \in [T]} \sigma_{t,p^*}^2(x_t) \mid p^* \right] \right]} \quad (\beta_t \text{ increasing}) \quad (54)$$

$$\leq \sqrt{\beta_T T} \sqrt{C \mathbb{E}_{p^*} [\gamma_{T,p^*}]} \quad (\text{Lemma 5.4 of Srinivas et al. (2012)}) \quad (55)$$

$$\leq \sqrt{\beta_T T} \sqrt{C \bar{\gamma}_T} \quad (56)$$

Combining the bounds for (1) and (3), we obtain the desired result

$$\text{BR}(T) \leq \frac{\pi^2}{6} + \sum_{t \in [T]} \mathbb{E} [U_{t,p_t}(x_t) - U_{t,p^*}(x_t)] + \sqrt{CT\beta_T\bar{\gamma}_T}. \quad (57)$$

□

A.3 Information-theoretic regret bound for GP-TS

We begin by showing that the rewards are subgaussian.

Lemma 4.5. *Fix $x \in \mathcal{X}$ and let $Z = f(x) + \epsilon_t - \mathbb{E}_t[f(x)]$. If $p^* \sim P_1$, $f \sim \mathcal{GP}(\mu_{1,p^*}, k_{1,p^*})$ with $k_{1,p^*}(\cdot, \cdot) \leq \sigma_0^2$, then $Z|H_t$ is $\sqrt{\sigma_0^2 + \sigma^2}$ -subgaussian.*

Proof. Start by considering $Z = \sum_{p \in P} \mathbb{1}\{p^* = p\}(f_p(x) - \mu_{t,p}(x))$ where $f_p \sim \mathcal{GP}_p$ such that $f_p(x) - \mu_{t,p}(x)$ is σ_0 -subgaussian. Then,

$$\mathbb{E}[\lambda Z | H_t] = \mathbb{E}[\mathbb{E}[\exp(\lambda Z) | p^* = p, H_t]] \quad (58)$$

$$= \sum_{p \in P} \mathbb{P}_t(p^* = p) \mathbb{E}[\exp(\lambda(f_p(x) - \mu_{t,p}(x))) | p^* = p, H_t] \quad (59)$$

$$\leq \sum_{p \in P} \mathbb{P}_t(p^* = p) \exp(\sigma_{t,p}^2(x) \lambda^2 / 2) \quad (60)$$

$$\leq \sum_{p \in P} \mathbb{P}_t(p^* = p) \exp(\sigma_0^2 \lambda^2 / 2) \quad (61)$$

$$= \exp(\sigma_0^2 \lambda^2 / 2) \quad (62)$$

Hence, Z is σ -subgaussian. Then, $Z + \epsilon_t | H_t$ is $\sqrt{\sigma_0^2 + \sigma^2}$ -subgaussian since $Z | H_t$ and $\epsilon_t | H_t$ are independent. \square

Then, we state and prove the information-theoretic regret bound for HP-GP-TS.

Theorem 4.6. *If $p^* \sim P_1$, $f \sim \mathcal{GP}(\mu_{1,p^*}, k_{1,p^*})$, then the Bayesian regret of HP-GP-TS is bounded by*

$$BR(T) \leq \sqrt{2|\mathcal{X}| \log(|\mathcal{X}|)(\sigma_0^2 + \sigma^2)T}. \quad (10)$$

Proof. Here, we analyze the regret of HP-GP-TS using the information-theoretic framework of [Russo & Roy \(2016\)](#). The proof presented here follows the proof in section D.2 of [Russo & Roy \(2016\)](#) with subgaussian noise. The idea of the information-theoretic framework is to express the regret as follows

$$\mathbb{E} \left[\sum_{t \in [T]} f(x^*) - f(x_t) \right] = \mathbb{E} \left[\sum_{t \in [T]} \underbrace{\mathbb{E}[f(x^*) - f(x_t) | H_t]}_{\geq 0} \right] \quad (\text{Tower rule}) \quad (63)$$

$$= \mathbb{E} \left[\sum_{t \in [T]} \sqrt{\mathbb{E}[f(x^*) - f(x_t) | H_t]^2} \cdot \frac{I(\cdot; \cdot | H_t)}{I(\cdot; \cdot | H_t)} \right] \quad (64)$$

where $I(\cdot; \cdot | H_t)$ is the mutual information between two carefully chosen variables such that $\mathbb{E}[f(x^*) - f(x_t) | H_t]^2 \leq CI(\cdot; \cdot | H_t)$ for some C .

Let $\mathbb{E}_t[\cdot] = \mathbb{E}[\cdot | H_t]$, $\mathbb{P}_t[\cdot] = \mathbb{P}[\cdot | H_t]$. Then,

$$\mathbb{E}_t[f(x^*) - f(x_t)] = \sum_{x \in \mathcal{X}} \mathbb{P}_t(x^* = x) \mathbb{E}_t[f(x) | x^* = x] - \sum_{x \in \mathcal{X}} \mathbb{P}_t(x_t = x) \mathbb{E}_t[f(x) | x_t = x] \quad (65)$$

$$= \sum_{x \in \mathcal{X}} \mathbb{P}_t(x^* = x) (\mathbb{E}_t[f(x) | x^* = x] - \mathbb{E}_t[f(x) | x_t = x]) \quad (\mathbb{P}_t(x_t = x) = \mathbb{P}_t(x^* = x)) \quad (66)$$

$$= \sum_{x \in \mathcal{X}} \mathbb{P}_t(x^* = x) (\mathbb{E}_t[f(x) | x^* = x] - \mathbb{E}_t[f(x)]) \cdot \left(\begin{array}{l} \mathbb{E}_t[f(x) | x_t = x] = \mathbb{E}_t[f(x)] \\ \text{since } f | H_t \perp\!\!\!\perp x_t | H_t \end{array} \right) \quad (67)$$

Let $Z = f(x) + \epsilon_t - \mathbb{E}_t[f(x)]$. Note that $Z|H_t$ is $\sqrt{\sigma_0^2 + \sigma^2}$ -subgaussian. Consequently, $\log \mathbb{E}_t[\exp(\lambda Z)] \leq \frac{\lambda^2(\sigma_0^2 + \sigma^2)}{2} \forall \lambda \in \mathbb{R}$.

Fix $x^*, x \in \mathcal{X}$. Then,

$$\lambda \left(\mathbb{E}_t[f(x) + \epsilon_t | x^* = x^*] - \mathbb{E}_t[f(x)] \right) - \frac{\lambda^2(\sigma_0^2 + \sigma^2)}{2} \quad (68)$$

$$\leq \lambda \left(\mathbb{E}_t[f(x) + \epsilon_t | x^* = x^*] - \mathbb{E}_t[f(x)] \right) \quad (69)$$

$$- \log \mathbb{E}_t[\exp(\lambda(f(x) + \epsilon_t - \mathbb{E}_t[f(x)]))] \quad (Z|H_t \text{ is } \sqrt{\sigma_0^2 + \sigma^2}\text{-subgaussian}) \quad (70)$$

$$= \lambda \mathbb{E}_t[Z | x^* = x^*] - \log \mathbb{E}_t[\exp(\lambda Z)] \quad (71)$$

$$\leq D(\mathbb{P}_t(f(x) + \epsilon_t | x^* = x^*) || \mathbb{P}_t(f(x) + \epsilon_t)). \quad (\text{Fact 12 of Russo \& Roy (2016)}) \quad (72)$$

Now, let $\lambda = \frac{1}{\sigma_0^2 + \sigma^2} (\mathbb{E}_t[f(x) + \epsilon_t | x^* = x^*] - \mathbb{E}_t[f(x)])$, then the following holds for all $x^*, x \in \mathcal{X}$

$$\mathbb{E}_t[f(x) | x^* = x^*] - \mathbb{E}_t[f(x)] \leq \sqrt{2(\sigma_0^2 + \sigma^2) D(\mathbb{P}_t(f(x) + \epsilon_t | x^* = x^*) || \mathbb{P}_t(f(x) + \epsilon_t))}. \quad (73)$$

Let $I_t(\cdot; \cdot) = I(\cdot; \cdot | H_t = h_t)$. Then,

$$I_t(x^*; (x_t, y_t)) = I_t(x^*; x_t) + I_t(x^*; y_t | x_t) \quad (\text{Chain rule}) \quad (74)$$

$$= I_t(x^*; y_t | x_t) \quad (x^* | H_t \perp\!\!\!\perp x_t | H_t) \quad (75)$$

$$= \sum_{x \in \mathcal{X}} \mathbb{P}_t(x_t = x) I_t(x^*; y_t | x_t = x) \quad (76)$$

$$= \sum_{x \in \mathcal{X}} \mathbb{P}_t(x_t = x) I_t(x^*; f(x_t) + \epsilon_t | x_t = x) \quad (77)$$

$$= \sum_{x \in \mathcal{X}} \mathbb{P}_t(x_t = x) I_t(x^*; f(x) + \epsilon_t) \quad (f, x^* | H_t \perp\!\!\!\perp x_t | H_t) \quad (78)$$

$$= \sum_{x \in \mathcal{X}} \mathbb{P}_t(x_t = x) \left(\sum_{x^* \in \mathcal{X}} \mathbb{P}_t(x^* = x^*) \cdot \right. \quad (79)$$

$$\left. D(\mathbb{P}_t(f(x) + \epsilon_t | x^* = x^*) || \mathbb{P}_t(f(x) + \epsilon_t)) \right) \quad (80)$$

$$= \sum_{x, x^* \in \mathcal{X}} \mathbb{P}_t(x^* = x) \mathbb{P}_t(x^* = x^*) \cdot \quad (x_t | H_t \stackrel{d}{=} x^* | H_t) \quad (81)$$

$$D(\mathbb{P}_t(f(x) + \epsilon_t | x^* = x^*) || \mathbb{P}_t(f(x) + \epsilon_t)) \quad (82)$$

Putting the above together, we now bound $\mathbb{E}_t[f(x^*) - f(x_t)]^2$:

$$\mathbb{E}_t[f(x^*) - f(x_t)]^2 = \left(\sum_{x \in \mathcal{X}} \mathbb{P}_t(x^* = x) (\mathbb{E}_t[f(x) | x^* = x] - \mathbb{E}_t[f(x)]) \right)^2 \quad (83)$$

$$\leq |\mathcal{X}| \sum_{x \in \mathcal{X}} \mathbb{P}_t(x^* = x)^2 (\mathbb{E}_t[f(x) | x^* = x] - \mathbb{E}_t[f(x)])^2 \quad (\text{Cauchy-Schwarz}) \quad (84)$$

$$\leq |\mathcal{X}| \sum_{x, x^* \in \mathcal{X}} \mathbb{P}_t(x^* = x) \mathbb{P}_t(x^* = x^*) (\mathbb{E}_t[f(x) | x^* = x^*] - \mathbb{E}_t[f(x)])^2 \quad (85)$$

$$\leq 2|\mathcal{X}|(\sigma_0^2 + \sigma^2) \sum_{x, x^* \in \mathcal{X}} \left(\mathbb{P}_t(x^* = x) \mathbb{P}_t(x^* = x^*) \cdot \right. \quad (86)$$

$$\left. D(\mathbb{P}_t(f(x) + \epsilon_t | x^* = x^*) || \mathbb{P}_t(f(x) + \epsilon_t)) \right) \quad (87)$$

$$\leq 2|\mathcal{X}|(\sigma_0^2 + \sigma^2) I_t(x^*; (x_t, y_t)). \quad (88)$$

Returning to the full regret,

$$\mathbb{E} \left[\sum_{t \in [T]} f(x^*) - f(x_t) \right] = \mathbb{E} \left[\sum_{t \in [T]} \sqrt{\mathbb{E}_t [f(x^*) - f(x_t)]^2} \cdot \frac{I_t(x^*; (x_t, y_t))}{I_t(x^*; (x_t, y_t))} \right] \quad (89)$$

$$\leq \sqrt{2|\mathcal{X}|(\sigma_0^2 + \sigma^2)} \mathbb{E} \left[\sum_{t \in [T]} \sqrt{I_t(x^*; (x_t, y_t))} \right] \quad (90)$$

$$(91)$$

Then, the expectation can then be bounded as

$$\mathbb{E} \left[\sum_{t \in [T]} \sqrt{I_t(x^*; (x_t, y_t))} \right] \leq \sum_{t \in [T]} \mathbb{E} \left[\sqrt{I_t(x^*; (x_t, y_t))} \right] \quad (92)$$

$$\leq \sum_{t \in [T]} 1 \cdot \sqrt{\mathbb{E} [I_t(x^*; (x_t, y_t))]} \quad (\text{Jensen's inequality}) \quad (93)$$

$$= \sum_{t \in [T]} \sqrt{\mathbb{E}_{H_t} [I(x^*; (x_t, y_t)) | H_t = h_t]} \quad (94)$$

$$= \sum_{t \in [T]} \sqrt{I(x^*; (x_t, y_t)) | H_t} \quad (95)$$

$$\leq \sqrt{\sum_{t \in [T]} 1^2 \sum_{t \in [T]} I(x^*; (x_t, y_t)) | H_t} \quad (\text{Cauchy-Schwarz}) \quad (96)$$

Finally, the

$$\mathbb{E} \left[\sum_{t \in [T]} f(x^*) - f(x_t) \right] = \sqrt{2|\mathcal{X}|(\sigma_0^2 + \sigma^2)} \sqrt{TI(x^*; H_T)} \quad (97)$$

$$\leq \sqrt{2|\mathcal{X}|(\sigma_0^2 + \sigma^2)} \sqrt{TH(x^*)} \quad (98)$$

$$\leq \sqrt{2|\mathcal{X}|(\sigma_0^2 + \sigma^2)} \sqrt{T \log |\mathcal{X}|} \quad (99)$$

$$(100)$$

□

B Description of kernels

The RBF kernel, $k(x, \tilde{x}) = \exp(-\|x - \tilde{x}\|^2/\ell^2)$ guarantees that f is smooth. The length-scale parameter $\ell > 0$ determines how quickly f changes, smaller values lead to more fluctuations. The rational quadratic (RQ) kernel $k(x, \tilde{x}) = \left(1 + \frac{\|x - \tilde{x}\|^2}{2\alpha\ell^2}\right)^{-\alpha}$ where $\alpha > 0$ is a mixture of RBF kernels with varying lengthscales. The Matérn kernel (Matérn, 1986) $k(x, \tilde{x}) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}\|x - \tilde{x}\|}{\ell}\right)^\nu K_\nu\left(\frac{\sqrt{2\nu}\|x - \tilde{x}\|}{\ell}\right)$ where $\nu > 0$ is the smoothness parameter that imposes that f is k -times differentiable if $\nu > k$ for integer k . The functions $\Gamma(\nu)$ and K_ν correspond to the gamma function and a modified Bessel function (Williams & Rasmussen, 2006). The periodic kernel $k(x, \tilde{x}) = \exp\left(-\frac{1}{2} \sum_{i=1}^d \sin^2\left(\frac{\pi}{\rho}(x - \tilde{x})\right)/\ell\right)$ generates smooth and periodic functions with period $\rho > 0$ (Mackay, 1998). The linear kernel $k(x, \tilde{x}) = vx^\top \tilde{x}/\nu$ generates linear functions where v is the variance parameter.

C Additional experimental details

In this section, we provide some additional details about the experiments.

For all the real-world datasets, sensors containing any null measurements have been filtered out.

The Intel Berkeley dataset consists of measurements from 46 temperature sensors across 19 days. The training set consists of the first 12 days of measurements and the remaining 7 days constitute the test set. The noise variance is set to $\sigma^2 = 0.7^2$.

The PeMS data consists of measurements from 211 sensors along the I-880 highway from all of 2023. The goal is to find the sensors with low speeds to identify congestions. We use the 5-min averages provided by PeMS. Data between 2023-01-01 and 2023-09-01 is put into the training set whilst the data until 2023-12-31 is put into the test set. The noise variance is set to $\sigma^2 = 2.25^2$.

The PNW precipitation data consists of daily precipitation data from 1949 to 1950 across $167 50 \times 50$ km regions in the Pacific Northwest. The goal is to find the region with the highest precipitation for any given day. The training data consists of the measurements made prior to 1980 and the test data consists of the measurements between 1980 and 1994. The original data is stated to be given in mm/day however the data seems to be off by a factor of 10. We rescale the data to a log-scale using $\log(\cdot/10 + 0.1)$, similar to Krause et al. (2008). The noise variance is set to $\sigma^2 = 0.41^2$.

In the Intel experiment, we removed one outlier seed. All methods had a final cumulative regret around 6000°C, note that the average for the worst performing model across the other seeds was $\approx 150^\circ\text{C}$. The outlier is shown Fig. 8. We can see that one of the sensors display very high temperatures compared to all other sensors, which is why all methods performed poorly on this seed. It should be noted that many of the sensors in the Intel data logged degrees above 100°C after a certain time - likely due to sensor failure rather boiling temperatures in an office environment. Also note that these days were excluded from both our training and test data. The outlier could be an indication that this particular sensor was starting to fail earlier than others.

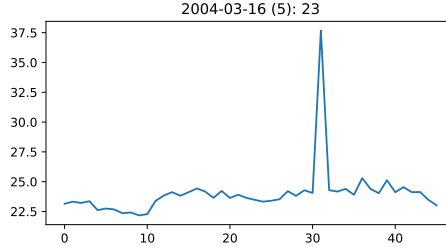


Figure 8: Removed sample from the test data in the Intel experiment. One of the sensors display very high temperatures.

D Additional experimental results

In this section, we provide some additional experimental results.

First, we include the mean number of priors in P_t for all experiments in Fig. 9. Similarly, we include the average entropy of the hyperposterior for all experiments in Fig. 10. For the lengthscale, subspace, PeMS and PNW precipitation experiments, hardly any priors are eliminated. In contrast, the hyperposterior entropy concentrates rapidly across all experiments with the subspace and PNW precipitation having the most and least concentrated hyperposterior.

We include the full set of confusion matrices for the lengthscale and subspace experiments in Fig. 11. In the lengthscale experiments, we observe that PE-GP-UCB and -TS oversample the shortest lengthscale. This is similar to the kernel experiment where the Matérn 3/2 kernel was also oversampled. However, we see that HP-GP-TS and MAP GP-TS do not suffer from this optimistic bias. In the subspace experiment, HP- and MAP GP-TS have an accuracy of around 96% where as PE-GP-TS and -UCB have accuracies 30% and 36% respectively. The priors are equivalent up to coordinate permutations and are therefore difficult. The PE-methods do not oversample any specific prior but commit to much time to exploring along the irrelevant dimensions.

In Tables 1 and 2, the total regret for the lengthscale and subspace scaling experiments are shown.

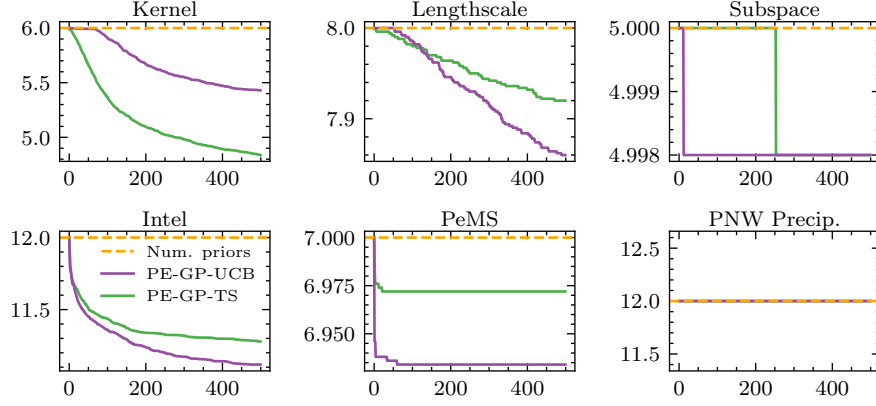


Figure 9: Mean number of priors remaining in P_t over time for PE-GP-UCB and -TS.

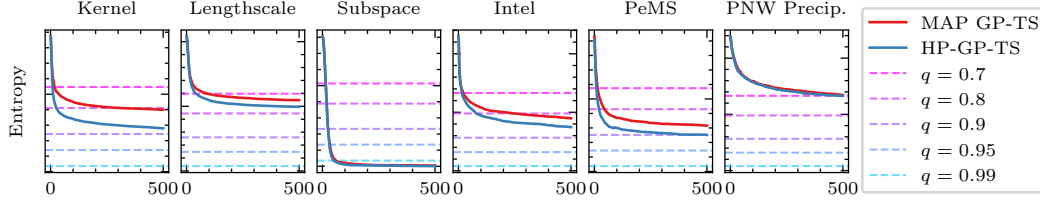


Figure 10: Average entropy in the hyperposterior P_t over time for HP- and MAP GP-TS. The dashed reference values correspond to entropies of discrete distributions with prob. q on one choice and prob. $\frac{1-q}{|P|-1}$ on the other $|P| - 1$ choices.

In Fig. 12, we visualize the median cumulative regret on the real-world data experiments. Similarly, in Fig. 13, we show further quantiles of the final cumulative regret. MAP and HP-GP-TS exhibit similar median regret across the experiments but MAP GP-TS has a larger variance and tail values.

Table 1: Average total regret and ± 1 standard error for the lengthscale experiment as $|P|$ increases.

Algorithm	Lengthscales, $ P $				
	8	16	32	64	128
MAP GP-TS	30.2 ± 1.2	32.4 ± 2.5	32.5 ± 2.1	28.7 ± 1.1	30.8 ± 1.9
HP-GP-TS	31.4 ± 1.0	31.7 ± 0.9	30.8 ± 0.8	30.7 ± 1.0	31.0 ± 1.4
PE-GP-TS	61.8 ± 0.5	61.3 ± 0.5	62.2 ± 0.5	62.4 ± 0.4	64.3 ± 0.4
PE-GP-UCB	114.2 ± 0.6	114.8 ± 0.6	115.5 ± 0.6	114.5 ± 0.6	114.8 ± 0.6
Oracle GP-TS	28.1 ± 0.8	26.4 ± 0.8	27.3 ± 0.8	26.5 ± 0.7	25.7 ± 0.7
Oracle GP-UCB	48.3 ± 1.2	46.9 ± 1.1	48.4 ± 1.1	46.5 ± 1.0	45.6 ± 1.0

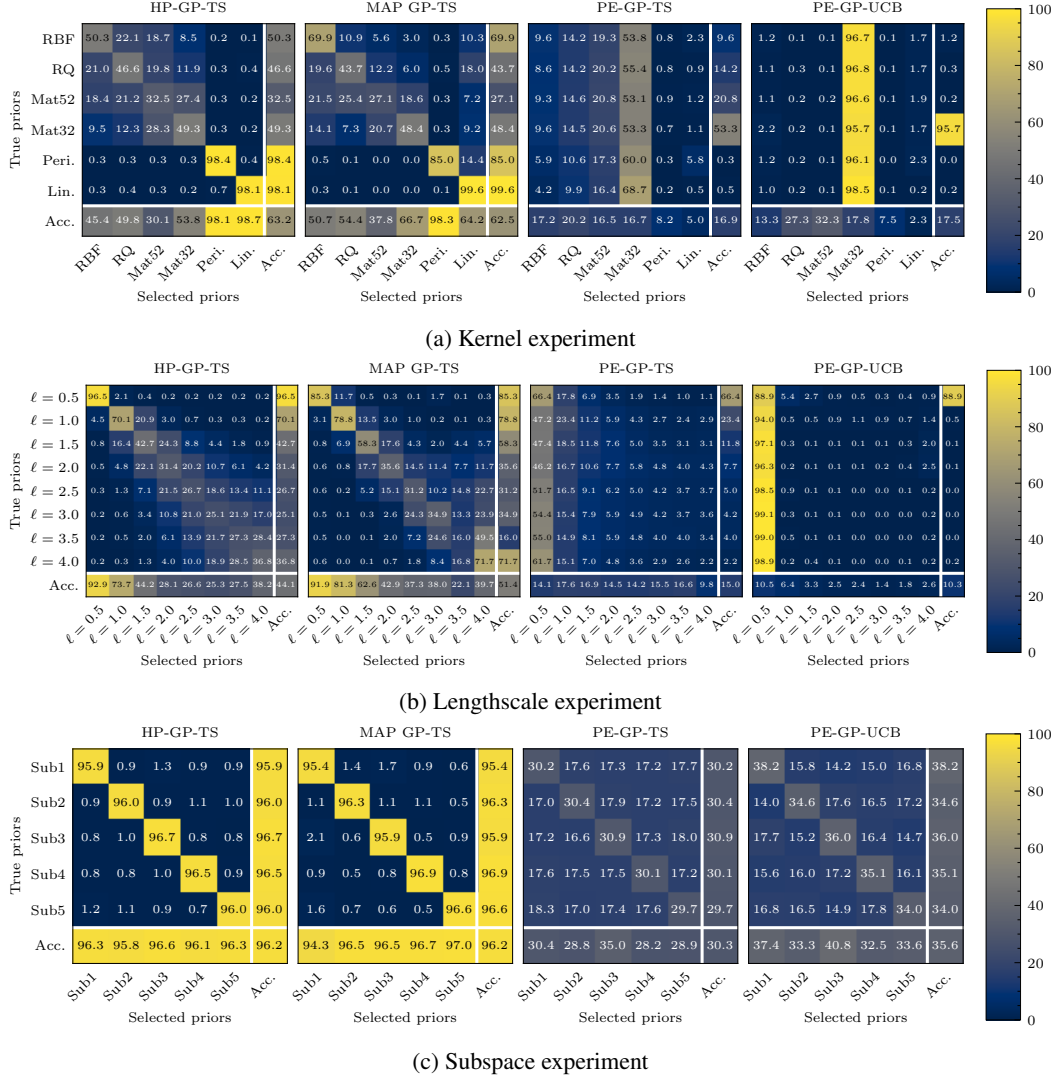


Figure 11: Confusion matrices for the true prior p^* and p_t across all time steps of the synthetic experiments.

Table 2: Average total regret and ± 1 standard error for the subspace experiment as $|P|$ increases.

Algorithm	Subspaces, $ P $			
	5	8	12	16
MAP GP-TS	87.2 ± 1.0	89.9 ± 1.1	89.1 ± 0.9	90.9 ± 1.2
HP-GP-TS	88.3 ± 0.9	88.8 ± 0.9	89.5 ± 0.9	90.8 ± 0.9
PE-GP-TS	177.1 ± 1.4	269.5 ± 1.9	344.7 ± 2.3	396.9 ± 2.5
PE-GP-UCB	389.0 ± 1.5	526.0 ± 1.8	622.4 ± 2.3	688.0 ± 2.7
Oracle GP-TS	86.0 ± 1.0	84.1 ± 0.9	84.6 ± 1.0	84.8 ± 1.0
Oracle GP-UCB	217.3 ± 1.0	218.2 ± 1.0	218.6 ± 1.0	218.9 ± 0.9

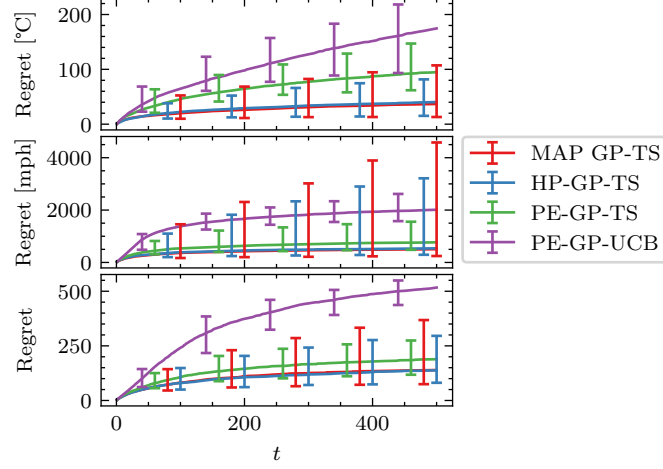


Figure 12: Median cumulative regret on Intel temperature data (top), PeMS speed data (middle) and PNW precipitation data (bottom). Errorbars correspond to first and last decile.

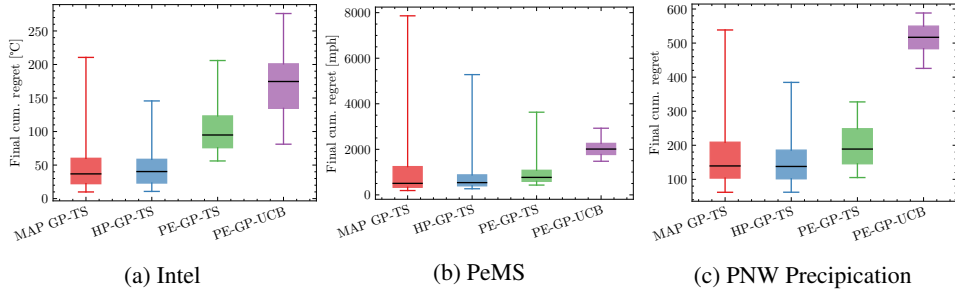


Figure 13: Quantiles of the final cumulative regret on the real-world data experiments. The median is shown with a black line. The whiskers correspond to the 5th and 95th percentile and the lower and upper edges of the boxes show the first and third quartile.