

---

# Scalable Thompson Sampling using Sparse Gaussian Process Models

---

Sattar Vakili<sup>\*1</sup>, Henry Moss<sup>\*2</sup>, Artem Artemev<sup>2</sup>, Vincent Dutordoir<sup>2,3</sup>, Victor Picheny<sup>2</sup>

## Abstract

Thompson Sampling (TS) from Gaussian Process (GP) models is a powerful tool for the optimization of black-box functions. Although TS enjoys strong theoretical guarantees and convincing empirical performance, it incurs a large computational overhead that scales polynomially with the optimization budget. Recently, scalable TS methods based on sparse GP models have been proposed to increase the scope of TS, enabling its application to problems that are sufficiently multi-modal, noisy or combinatorial to require more than a few hundred evaluations to be solved. However, the approximation error introduced by sparse GPs invalidates all existing regret bounds. In this work, we perform a theoretical and empirical analysis of scalable TS. We provide theoretical guarantees and show that the drastic reduction in computational complexity of scalable TS can be enjoyed without loss in the regret performance over the standard TS. These conceptual claims are validated for practical implementations of scalable TS on synthetic benchmarks and as part of a real-world high-throughput molecular design task.

## 1 Introduction

Thompson sampling [TS, 1] is a popular algorithm for Bayesian optimization [BO, 2] — a sequential model-based approach for the optimization of expensive-to-evaluate black-box functions, typically characterised by limited prior knowledge and access to only a limited number of (possibly noisy) evaluations. By sequentially evaluating the maxima of random samples from a model of the objective function, TS provides a conceptually simple method for balancing exploration and exploitation.

TS is often paired with Gaussian Processes (GPs), which offers a spectrum of powerful and flexible modeling tools that provide probabilistic predictions of the objective function. The resulting GP-TS algorithms [3] have been found to provide highly efficient optimization under heavily restricted optimization budgets, with numerous successful applications including aerodynamic design [4], route planning [5] and web-streaming [6]. While most popular BO algorithms cannot query more than a handful of points at a time [7–10] without employing replicating designs [see 11, 12], TS has a natural ability to query large batches of points. Therefore, TS is a popular solution for optimization pipelines enjoying a large degree of parallelisation, for example in high-throughput chemical space exploration [13] and for the distributed tuning of machine learning models across cloud compute resources [14].

As BO incurs a substantial computational overhead between successive iterations, while updating models and choosing the next set of query points, standard BO methods are limited to optimization problems with small evaluation budgets [2]. However, with large batches, the computational overhead incurred by BO per individual function evaluation is considerably reduced. Therefore, considering large batches is a promising tactic to expand BO to larger optimization budgets, which are required to optimize highly noisy problems with rougher optimization landscapes [11, 12] or high dimensional

---

<sup>\*</sup>Equal contribution, <sup>1</sup>MediaTek Research, <sup>2</sup>Secondmind Labs, Cambridge, UK. Correspondence to Sattar Vakili <sattar.vakili@mtkresearch.com>, Henry Moss <henry.moss@secondmind.ai>.

and combinatorial search spaces [15, 13, 16]. Consequently, the highly-parallelizable TS is a promising candidate for BO under large optimization budgets.

Unfortunately, practical implementations of GP-TS suffer from two key computational bottlenecks that prevent the method from scaling in terms of total optimization budget. Not only does each update of the GP posterior distribution require a matrix inversion that incurs a cubic cost w.r.t. the number of observations  $t$  [17], but even sampling from this posterior can be a daunting task — the standard approach of drawing a joint sample across a  $N$  point discretization of the search space has an  $O(N^3)$  complexity [due to a Cholesky decomposition step, 18]. Alternative existing approaches for BO under large optimization budgets include using Neural Networks in lieu of GPs [15, 13] or to use local models [19] and ensembles [16].

A natural answer to the scalability issues of GP-TS is to rely on the recent advances in Sparse Variational GP models [SVGP, 20]. SVGPs provide a low rank  $O(m^2t)$  approximation of the GP posterior, where  $m$  is the number of the so-called *inducing variables* that grows at a rate much slower than  $t$ . Successful applications of SVGPs for BO under large optimization budgets include optimizing a free-electron laser [21], molecules under synthesis-ability constraints [22], and the composition of alloys [23]. Furthermore, [24] introduced an efficient sampling rule (referred to as *decoupled* sampling) which can be used to efficiently perform TS with SVGPs. In particular, [24] decomposes samples from the SVGP posterior into the sum of an approximate prior based on  $M$  features (see Sec. 3.3) and an SVGP model update, thus reducing the computational cost of drawing a Thompson sample to  $O((m + M)N)$ . Leveraging this sampling rule results in a scalable GP-TS algorithm (henceforth S-GP-TS) that can handle orders of magnitude greater optimization budgets.

While [3] proposed a comprehensive theoretical analysis of exact GP-TS, it does not apply to S-GP-TS. Indeed, using sparse models and decoupled sampling introduce two layers of approximation, that must be handled with care, as even a small constant error in the posterior can lead to poor performance by encouraging under-exploration in the vicinity of the optimum point [25]. Our primary contributions can be summarised as follows. First, we provide a theoretical analysis showing that batch TS from any approximate GP can achieve the same regret order as an exact GP-TS algorithm as long the quality of the posterior approximations satisfies certain conditions (Assumptions 3 and 4). Second, for the specific case of S-GP-TS (batch decoupled TS using a SVGP), we leverage the results of [26] to provide bounds in terms of GP’s kernel spectrum for the number of prior features and inducing variables required to guarantee low regret. Finally, we investigate empirically the performance of multiple practical implementations of S-GP-TS, considering synthetic benchmarks and a high-throughput molecular design task.

## 2 Problem Formulation

We consider the sequential optimization of an unknown function  $f$  over a compact set  $\mathcal{X} \subset \mathbb{R}^d$ . A sequential learning policy selects a batch of  $B$  observation points  $\{x_{t,b}\}_{b \in [B]}$  at each time step  $t = 1, 2, \dots, T$  and receives the corresponding real-valued and noisy rewards  $\{y_{t,b} = f(x_{t,b}) + \epsilon_{t,b}\}_{b \in [B]}$ , where  $\epsilon_{t,b}$  denotes the observation noise. Throughout the paper, we use the notation  $[n] = \{1, 2, \dots, n\}$ , for  $n \in \mathbb{N}$ . As is common in both the bandits and GP literature, our analysis uses the following sub-Gaussianity assumption, a direct consequence of which is that  $\mathbb{E}[\epsilon_{t,b}] = 0$ , for all  $t, b \in \mathbb{N}$ .

**Assumption 1.**  $\epsilon_{t,b}$  are i.i.d., over both  $t$  and  $b$ ,  $R$ -sub-Gaussian random variables, where  $R > 0$  is a fixed constant. Specifically,  $\mathbb{E}[e^{h\epsilon_{t,b}}] \leq \exp(\frac{h^2 R^2}{2})$ ,  $\forall h \in \mathbb{R}, \forall t, b \in \mathbb{N}$ .

Let  $x^* \in \operatorname{argmax}_{x \in \mathcal{X}} f(x)$  be an optimal point. We can then measure the performance of a sequential optimizer by its *strict regret*, defined as the cumulative loss compared to  $f(x^*)$  over a time horizon  $T$

$$R(T, B; f) = \mathbb{E} \left[ \sum_{t=1}^T \sum_{b=1}^B f(x^*) - f(x_{t,b}) \right], \quad (1)$$

where the expectation is with respect to the randomness in noise and the possible stochasticity in the sequence of the selected batch observation points  $\{x_{t,b}\}_{t \in [T], b \in [B]}$ . Note that our regret measure (1) is defined for the true unknown  $f$ . In contrast, the alternative Bayesian regret [see e.g. 27, 14] averages over a prior distribution for  $f$ . As upper bounds on strict regret directly apply to the Bayesian regret (but not necessarily the reverse), our results are stronger than those that can be achieved when

analysing just Bayesian regret, for example when applying the technique of [28] that equates TS’s Bayesian regret with that of the well-studied upper confidence bound policies.

Following [3, 29, 30], our analysis assumes a regularity condition on the objective function motivated by kernelized learning models and their associated reproducing kernel Hilbert spaces [RKHS, 31]:

**Assumption 2.** *Given an RKHS  $H_k$ , the norm of the objective function is bounded:  $\|f\|_{H_k} \leq \mathcal{B}$ , for some  $\mathcal{B} > 0$ , and  $k(x, x') \leq 1$ , for all  $x, x' \in \mathcal{X}$ .*

In the case of practically relevant kernels, Assumption 2 implies certain smoothness properties for the objective functions.

### 3 Gaussian Processes and Sparse Models

GPs are powerful non-parametric Bayesian models over the space of functions [17] with a distribution specified by a mean function  $\mu(x)$  (henceforth assumed to be zero for simplicity) and a positive definite kernel (or covariance function)  $k(x, x')$ . We provide here a brief description of the classical GP model and two sparse variational formulations.

#### 3.1 Exact Gaussian Process models

Suppose that we have collected a set of location-observation tuples  $\mathcal{H}_t = \{\mathbf{X}_t, \mathbf{y}_t\}$ , where  $\mathbf{X}_t$  is the  $tB \times d$  matrix of locations with rows  $[\mathbf{X}_t]_{(s-1)B+b} = x_{s,b}$ , and  $\mathbf{y}_t$  is the  $tB$ -dimensional column vector of observations with elements  $[\mathbf{y}_t]_{(s-1)B+b} = y_{s,b}$ , for all  $s \in [t]$ ,  $b \in [B]$ . Then, assuming a Gaussian observation noise, the posterior of the GP model  $\hat{f}$  given the set of past observations  $\mathcal{H}_t$ , is also a GP with mean  $\mu_t(\cdot)$ , variance  $\sigma_t^2(\cdot)$  and kernel function  $k_t(\cdot, \cdot)$  specified as

$$\mu_t(x) = k_{\mathbf{X}_t, x}^T (K_{\mathbf{X}_t, \mathbf{X}_t} + \tau \mathbf{I})^{-1} \mathbf{y}_t, \quad k_t(x, x') = k(x, x') - k_{\mathbf{X}_t, x}^T (K_{\mathbf{X}_t, \mathbf{X}_t} + \tau \mathbf{I})^{-1} k_{\mathbf{X}_t, x'}, \quad (2)$$

and  $\sigma_t^2(x) = k_t(x, x)$ , with  $k_{\mathbf{X}_t, x}$  the  $tB$  dimensional column vector with entries  $[k_{\mathbf{X}_t, x}]_{(s-1)B+b} = k(x_{s,b}, x)$ , and  $K_{\mathbf{X}_t, \mathbf{X}_t}$  the  $tB \times tB$  positive definite covariance matrix with entries  $[K_{\mathbf{X}_t, \mathbf{X}_t}]_{(s-1)B+b, (s'-1)B+b'} = k(x_{s,b}, x_{s',b'})$ . We directly see from (2) that accessing the posterior expressions require an  $O((tB)^3)$  matrix inversion, which is a computational bottleneck for large values of  $tB$ .

Note that in our problem formulation  $f$  is fixed and observation noise has an unknown sub-Gaussian distribution. Using a GP prior and assuming a Gaussian noise is merely for ease of modelling and does not affect our assumptions on  $f$  and  $\epsilon_{t,b}$ . The notation  $\hat{f}$  is thus used to distinguish the GP model from the fixed  $f$ .

#### 3.2 Sparse Variational Gaussian Process Models with Inducing Points

To overcome the cubic cost of exact GPs, SVGPs [20, 32] instead approximate the GP posterior through a set of *inducing points*  $\mathbf{Z}_t = \{z_1, \dots, z_{m_t}\}$  ( $z_i \in \mathcal{X}$ , with  $m_t \ll tB$ ). Conditioning on the *inducing variables*  $\mathbf{u}_t = \hat{f}(\mathbf{Z}_t)$  (rather than the  $tB$  observations in  $\mathbf{y}_t$ ) and specifying a prior Gaussian density  $q_t(\mathbf{u}_t) = \mathcal{N}(\mathbf{m}_t, \mathbf{S}_t)$ , yields an approximate posterior distribution that, crucially, is still a GP but with the significantly reduced computational complexity of  $O(m_t^2 t)$ . The posterior mean and covariance of the SVGP are given in closed form as

$$\mu_t^{(s)}(x) = k_{\mathbf{Z}_t, x}^T K_{\mathbf{Z}_t, \mathbf{Z}_t}^{-1} \mathbf{m}_t \quad k_t^{(s)}(x, x') = k(x, x') + k_{\mathbf{Z}_t, x}^T K_{\mathbf{Z}_t, \mathbf{Z}_t}^{-1} (\mathbf{S}_t - K_{\mathbf{Z}_t, \mathbf{Z}_t}) K_{\mathbf{Z}_t, \mathbf{Z}_t}^{-1} k_{\mathbf{Z}_t, x'}.$$

The variational parameters  $\mathbf{m}_t$  and  $\mathbf{S}_t$  are set as the maximizers of the evidence lower bound (ELBO, see Appendix A for the details) and can be optimized numerically with mini-batching [32]. There are various standard ways in practice to select the locations of the inducing points  $\mathbf{Z}_t$ , e.g. by using an experimental design, sampling from a k-DPP (that stands for determinantal point process), or by optimizing them along with the inducing variables.

#### 3.3 Sparse Variational Gaussian Process Models with Inducing Features

An alternative approximation strategy is using inducing feature approximations [33, 26, 34]. Here, we define inducing variables as the linear integral transform of  $\hat{f}$  with respect to some *inducing*

features [35]  $\psi_1(x), \dots, \psi_{m_t}(x)$ , i.e we set our  $i^{\text{th}}$  inducing variable as  $u_{t,i} = \int_{\mathcal{X}} \hat{f}(x) \psi_i(x) dx$ . Courtesy of Mercer’s theorem, we can usually decompose our chosen kernel  $k$  as the inner product of possibly infinite dimensional feature maps (see Theorem 4.1 in [36]) to provide the expansion  $k(x, x') = \sum_{j=1}^{\infty} \lambda_j \phi_j(x) \cdot \phi_j(x')$  for eigenvalues  $\{\lambda_j \in \mathbb{R}^+\}_{j=1}^{\infty}$  and eigenfunctions  $\{\phi_j \in H_k\}_{j=1}^{\infty}$ . If we set our inducing features to be the  $m_t$  eigenfunctions with largest eigenvalues, it can be shown that  $\text{cov}(u_{t,i}, u_{t,j}) = \lambda_j \delta_{i,j}$  and  $\text{cov}(u_{t,j}, \hat{f}(x)) = \lambda_j \phi_j(x)$ , yielding an approximate Gaussian Process model with posterior mean and covariance given by

$$\mu_t^{(s)}(x) = \phi_{m_t}^T(x) \mathbf{m}_t \quad k_t^{(s)}(x, x') = k(x, x') + \phi_{m_t}^T(x) (\mathbf{S}_t - \Lambda_{m_t}) \phi_{m_t}(x').$$

Here,  $\mathbf{m}_t$  and  $\mathbf{S}_t$  are inducing parameters (as above),  $\phi_{m_t}(x) \triangleq [\phi_1(x), \dots, \phi_{m_t}(x)]^T$  is the truncated feature vector and  $\Lambda_{m_t}$  is the  $m_t \times m_t$  diagonal matrix of eigenvalues,  $[\Lambda_{m_t}]_{i,j} = \lambda_i \delta_{i,j}$ .

Inducing feature approximations have strong advantages, in particular a reduced computational cost and the fact that no inducing points need to be specified. However, accessing these eigenfeatures require the Mercer decomposition of the used kernel, which is available for certain kernels on manifolds [37, 34], but limited to low dimensions for others [38, 39].

#### 4 Scalable Thompson Sampling using Gaussian Process Models (S-GP-TS)

At each BO step  $t$ , GP-TS proceeds by drawing  $B$  i.i.d. samples  $\{\hat{f}_{t,b}\}_{b \in [B]}$  from the posterior distribution of  $\hat{f}$  and finding their maximizers, i.e. we select samples  $x_{t,b}$  satisfying

$$\{x_{t,b} = \arg\max_{x \in \mathcal{X}} \hat{f}_{t,b}(x)\}_{b \in [B]}. \quad (3)$$

However, since  $\hat{f}_{t,b}$  is an infinite dimensional object, generating such samples is computationally challenging. Consequently, it is common to resort to approximate strategies, the most simple of which is to sample across an  $N_t$  point discretization  $D_t$  of  $\mathcal{X}$  [14] which can be obtained with an  $O(N_t^3)$  cost (due to a required Cholesky decomposition).

To improve the computational efficiency of TS, a classical strategy [40, 41] is to rely on kernel decompositions. For instance, a sample  $\hat{f}$  from a GP can be expressed as a randomly weighted sum of the kernel’s eigenfunctions  $\hat{f}(x) = \sum_{j=1}^{\infty} \sqrt{\lambda_j} w_j \phi_j(x)$ , or, in the case of shift-invariant kernels, the kernel’s Fourier features  $\psi_j(x)$  (see [42]) as  $\hat{f}(x) = \sum_{j=1}^{\infty} w_j \psi_j(x)$ . By truncating these infinite expansions to contain only the  $M$  eigenfunctions with largest eigenvalues or  $M$  random Fourier features, we have access to approximate but analytically tractable samples. For both expansions, the weights  $w_j$  are sampled independently from a standard normal distribution. Conditioned on current  $tB$  observations, the posterior distribution of  $w_j$  are Gaussian with mean and covariance functions that can be calculated with an  $O(M^3)$  computations, resulting in an  $O(M^3 + BNM)$  cost to draw  $B$  Thompson samples.

Fast approximation strategies described above avoid costly matrix operations and work best only when sampling from GP priors. Posterior GP distributions are often too complex to be well-approximated by a finite feature representation [16, 43, 30]. The recent work of [24] tackled this issue by using truncated feature representations only to approximate the prior GP and a separate model update term to approximate posterior samples. For SVGP models, this has been shown to yield more accurate Thompson samples whilst incurring only an  $O((m_t + M)BN)$ , on top of the  $O(tBm_t^2)$  SVGP model fit, per optimization step  $t$ .

For our theoretical analysis, we consider two distinct decoupled sampling rules inspired by [24], one for each of the two SVGP formulations presented above [see 24, for derivations and similar expressions for Fourier decompositions]. The first rule is referred to as *Decoupled Sampling with Inducing Points* and is defined as

$$\tilde{f}_t(x) = \sum_{j=1}^M \alpha_t \sqrt{\lambda_j} w_j \phi_j(x) + \sum_{j=1}^{m_t} v_{t,j} k(x, z_j), \quad (4)$$

where we have coefficients  $v_{t,j} = [K_{\mathbf{z}_t, \mathbf{z}_t}^{-1} (\alpha_t (\mathbf{u}_t - \mathbf{m}_t) + \mathbf{m}_t - \alpha_t \Phi_{m_t, M} \Lambda_M^{\frac{1}{2}} \mathbf{w}_M)]_j$  for  $\Phi_{m_t, M} = [\phi_M(z_1), \dots, \phi_M(z_{m_t})]^T$  and  $\mathbf{w}_M = [w_1, \dots, w_M]^T$ . The weights  $w_i$  are drawn i.i.d from  $\mathcal{N}(0, 1)$ .

(4) is a modification of the sampling rule of [24] where we have added a scaling parameter  $\alpha_t \in \mathbb{R}$  (with  $\alpha_t = 1$ , the sampling rule of [24] is recovered). When set to be greater than one,  $\alpha_t$  serves to increase the variability of the approximate function samples (without changing their mean) and is used in our analysis to ensure sufficient exploration.

To efficiently sample from our second class of SVGP models, we also consider *Decoupled Sampling with Inducing Features*:

$$\tilde{f}_t(x) = \sum_{j=1}^M \alpha_t \sqrt{\lambda_j} w_j \phi_j(x) + \sum_{j=1}^{m_t} v_{t,j} \lambda_j \phi_j(x), \quad (5)$$

where  $v_{t,j} = [\Lambda_{m_t}^{-1}(\alpha_t(\mathbf{u}_t - \mathbf{m}_t) + \mathbf{m}_t - \alpha_t \Lambda_{m_t}^{\frac{1}{2}} \mathbf{w}_{m_t})]_j$  for  $\Lambda_{m_t}$  defined in Section 3.3.

## 5 Regret Analysis of S-GP-TS

Here, we first establish an upper bound on the regret of any approximate GP model (Theorem 1) based on the quality of their approximate posterior, as parameterized in Assumptions 3 and 4. We then discuss the consequences of Theorem 1 for the regret bounds and the computational complexity of S-GP-TS methods based on SVGPs and the decoupled sampling rules (4) and (5).

### 5.1 Regret Bounds Based on the Quality of Approximations

Consider a TS algorithm using an approximate GP model. In particular, assume an approximate model is provided where  $\tilde{k}_t$ ,  $\tilde{\sigma}_t$  and  $\tilde{\mu}_t$  are approximations of  $k_t$ ,  $\sigma_t$  and  $\mu_t$ , respectively. At each time  $t$ , a batch of  $B$  samples  $\{\tilde{f}_{t,b}\}_{b=1}^B$  is drawn from a GP with mean  $\tilde{\mu}_{t-1}$  and the scaled covariance  $\alpha_t^2 \tilde{k}_{t-1}$ . The batch of observation points  $\{x_{t,b}\}_{b=1}^B$  are selected as the maximizers of  $\{\tilde{f}_{t,b}\}_{b=1}^B$  over a discretization  $D_t$  of the search space.

We start our analysis by making two assumptions on the *quality* of approximations  $\tilde{\mu}_t$ ,  $\tilde{\sigma}_t$  of the posterior mean and the standard deviation. This parameterization is agnostic to the particular sampling rule (governing  $\tilde{\mu}_t$  and  $\tilde{\sigma}_t$ ) and provides valuable intuition that can be applied to any approximate method. When it comes to S-GP-TS (as the model governing  $\tilde{\mu}_t$ ,  $\tilde{\sigma}_t$ ), we show, in Sec. 5.2, that these assumptions are satisfied under some conditions on the value of the parameters of the sampling rules.

**Assumption 3** (quality of the approximate standard deviation). *For the approximate  $\tilde{\sigma}_t$ , the exact  $\sigma_t$ , and for all  $x \in \mathcal{X}$ ,*

$$\frac{1}{\underline{a}_t} \sigma_t(x) - \epsilon_t \leq \tilde{\sigma}_t(x) \leq \bar{a}_t \sigma_t(x) + \epsilon_t,$$

where  $1 \leq \underline{a}_t \leq \underline{a}$ ,  $1 \leq \bar{a}_t \leq \bar{a}$  for all  $t \geq 1$  and some constants  $\underline{a}, \bar{a} \in \mathbb{R}$ , and  $0 \leq \epsilon_t \leq \epsilon$  for all  $t \geq 1$  and some small constant  $\epsilon \in \mathbb{R}$ .

**Assumption 4** (quality of the approximate prediction). *For the approximate  $\tilde{\mu}_t$ , the exact  $\mu_t$  and  $\sigma_t$ , and for all  $x \in \mathcal{X}$ ,*

$$|\tilde{\mu}_t(x) - \mu_t(x)| \leq c_t \sigma_t(x),$$

where  $0 \leq c_t \leq c$  for all  $t \geq 1$  and some constant  $c \in \mathbb{R}$ .

The following Lemma establishes a concentration inequality for the approximate statistics using the one for exact statistics [3, Theorem 2].

**Lemma 1.** *Under Assumptions 1, 2, 3 and 4, with probability at least  $1 - \delta$ ,  $|f(x) - \tilde{\mu}_t(x)| \leq \tilde{u}_t(\tilde{\sigma}_t(x) + \epsilon_t)$ , where  $\tilde{u}_t(\delta) = \underline{a}_t \left( \mathcal{B} + R\sqrt{2(\gamma_{tB} + 1 + \log(1/\delta))} + c_t \right)$ .*

Proof is provided in Appendix B. Here,  $\gamma_s$  is the *maximal information gain*:  $\gamma_s = \max_{A \subset \mathcal{X}, |A|=s} \mathcal{I}([y(x)]_{x \in A}; [\hat{f}(x)]_{x \in A})$ , where  $\mathcal{I}([y(x)]_{x \in A}; [\hat{f}(x)]_{x \in A})$  denotes the mutual information [44, Chapter 2] between observations and the underlying GP model. The maximal information gain can itself be bounded for a specific kernel (see Sec. 5.3).

Following [29] and [3], we consider a discretization  $D_t$  of the search space satisfying the following assumption.

**Assumption 5.** The discretization  $D_t$  is designed in a way that  $|f(x) - f(\mathbf{x}^{(t)})| \leq 1/t^2$  for all  $x \in \mathcal{X}$ , where  $\mathbf{x}^{(t)} = \operatorname{argmin}_{x' \in D_t} \|x - x'\|$  is the closest point (in Euclidean norm) to  $x$  in  $D_t$ . The size of this discretization satisfies  $|D_t| = N_t \leq C(d, B)t^{2d}$  where  $C(d, B)$  is independent of  $t$  ([3, 29]).

We are now in a position to present regret bounds based on the quality of GP approximations:

**Theorem 1.** Consider S-GP-TS with  $\alpha_t = 2\bar{u}_t(1/(t^2))$ . Under Assumptions 1, 2, 3, 4 and 5, the regret defined in (1), satisfies

$$\begin{aligned} R(T, B; f) &\leq 30\bar{a}\beta_TB\sqrt{\frac{2T\gamma_T}{\log(1 + \frac{1}{\tau})}} + (31\beta_T + \alpha_T)\epsilon TB + 15BB + 2B \\ &= O\left(\underline{a}\bar{a}BR\sqrt{d\gamma_T(\gamma_TB + \log(T))T\log(T)} + \underline{a}\epsilon TBR\sqrt{d(\gamma_TB + \log(T))\log(T)}\right), \end{aligned} \quad (6)$$

where  $\beta_t = \alpha_t(b_t + \frac{1}{2})$  with  $b_t = \sqrt{2\log(N_t t^2)}$ .

See the proof in Appendix B. This regret bound scales with the product of the ratios  $\underline{a}$  and  $\bar{a}$ , with an additive term depending on the additive approximation error in the standard deviation.

## 5.2 Approximation Quality of the Decoupled Sampling Rule

For S-GP-TS with inducing points, we assume, as in [26], that the inducing points are sampled according to a discrete k-DPP. While this might be costly in practice, [26] showed that  $\mathbf{Z}_t$  can be efficiently sampled from  $\epsilon_0$  close sampling methods without compromising the predictive quality of SVGP. For both sampling rules, we also assume in our analysis that the Mercer decomposition of the kernel is used.

The quality of the approximation can be characterized using the spectral properties of the GP kernel. Let us define the tail mass of eigenvalues  $\delta_M = \sum_{i=M+1}^{\infty} \lambda_i \bar{\phi}_i^2$  where  $\bar{\phi}_i = \max_{x \in \mathcal{X}} \phi_i(x)$ . With decaying eigenvalues, including sufficient eigenfunctions in the feature representation results in a small  $\delta_M$ . In addition, [26] showed that, for an SVGP, a sufficient number of inducing variables ensures that the Kullback–Leibler (KL) divergence between the approximate and the true posterior distributions diminishes. Consequently, the approximate posterior mean and the approximate posterior variance converge to the true ones. Building on this result, we are able to prove Proposition 1 on the quality of approximations.

**Proposition 1.** For S-GP-TS based on sampling rule (4) with  $\alpha_t = 1$  and an SVGP using an  $\epsilon_0$  close k-DPP for selecting  $\mathbf{Z}_t$ , with probability at least  $1 - \delta$ , Assumptions 3 and 4 hold with parameters  $c_t = \sqrt{\kappa_t}$ ,  $\underline{a}_t = \frac{1}{\sqrt{1 - \sqrt{3\kappa_t}}}$ ,  $\bar{a}_t = \sqrt{1 + \sqrt{3\kappa_t}}$ , and  $\epsilon_t = \sqrt{C_1 m_t \delta_M}$ , where  $C_1$  is a constant specified in the appendix and  $\kappa_t = \frac{2tB(m_t+1)\delta_{m_t}}{\tau\delta} + \frac{4tB\epsilon_0}{\tau\delta}$ .

For S-GP-TS based on sampling rule (5) with  $\alpha_t = 1$ , Assumptions 3 and 4 hold with parameters  $c_t = \sqrt{\kappa_t}$ ,  $\underline{a}_t = \frac{1}{\sqrt{1 - \sqrt{3\kappa_t}}}$ ,  $\bar{a}_t = \sqrt{1 + \sqrt{3\kappa_t}}$ , and  $\epsilon_t = \sqrt{C_1 m_t \delta_M}$ , where  $C_1$  is the same constant as above and  $\kappa_t = \frac{2tB\delta_{m_t}}{\tau}$ .

Note that our proposition requires extending the results of [26] in two non-trivial ways. First, the decoupled sampling rules introduce an additional error. Secondly, [26] built their convergence results on the assumption that the observation points  $x_{t,b}$  are drawn from a prefixed distribution, which is not the case in S-GP-TS, where  $x_{t,b}$  are selected according to an experimental design method. A detailed proof of Proposition 1 is provided in Appendix B.

## 5.3 Application of Regret Bounds to Matérn and SE Kernels

We now investigate the application of Theorem 1 to the Squared Exponential (SE) and Matérn kernels, widely used in practice [see, e.g., 17, 45]. In the case of a Matérn kernel with smoothness parameter  $\nu > \frac{d}{2}$  it is known that  $\lambda_j = O(j^{-\frac{2\nu+d}{d}})$  [46]. For the SE kernel, we have  $\lambda_j = O(\exp(-j^{\frac{1}{d}}))$  [47, 48]. With these bounds on the spectrum of the kernels and the specific bounds on the maximal information gain [e.g.,  $\gamma_s \leq O(\log(s)^{d+1})$  for SE and  $\gamma_s \leq O(s^{d/(2\nu+d)} \log(s))$  for Matérn, 49], Theorem 1 and Proposition 1 result in the following theorem.



**Theorem 2.** Under Assumptions 1 and 2, with the algorithmic parameters, kernels and sampling rules specified in Table 1, S-GP-TS offers  $R(T, B; f) = O(B\sqrt{\gamma_T\gamma_{TB}T\log(T)})$ .

With a batch size  $B = 1$  Theorem 2 recovers the same regret bounds as the exact GP-TS [3]. We also note that for a fair comparison in terms of both the batch size and number of samples we should consider  $T' = TB$  as the number of samples. In that case, our regret bound becomes  $O(\sqrt{B\gamma_{T'}/B\gamma_{T'}T'\log(T'/B)})$ , which scales at most with  $\sqrt{B}$ . That is  $\sqrt{B}$  tighter than the trivial scaling with  $B$ .

In order to prove Theorem 2, the algorithmic parameters  $M$  and  $m_t$  must be selected large enough such that approximation parameters  $\underline{a}, \bar{a}, c, \epsilon$  in Assumptions 3 and 4 are sufficiently small. Using the relation between the algorithmic parameters, the approximation parameters and  $m_t$  provided by Proposition 1, the regret bound follows from Theorem 1. See Appendix B for a detailed proof.

The values of  $M$  and  $m_t$  required for Theorem 2 are summarized in Table 1. We also show the resulting computational cost of each sampling rule (as given by  $O(B(M + m_T)N_T T + Bm_T^2 T^2)$ ), explicitly demonstrating the improvement of S-GP-TS over the  $O(BN_T^3 T + B^3 T^4)$  computational cost of the vanilla GP-TS. Note that, for the Matérn kernel under sampling rule (4),  $\nu$  is required to be sufficiently larger than  $\frac{d}{2}$  in order for  $m_t$  to grow slower than  $t$ .

		Inducing points (4)	Inducing features (5)
Matérn	Condition	$m_t \sim T^{\frac{2d}{2\nu-d}}, M \sim T^{\frac{(2\nu+d)d}{2(2\nu-d)\nu}}$	$m_t \sim T^{\frac{d}{2\nu}}, M \sim T^{\frac{(2\nu+d)d}{4\nu^2}}$
	Cost	$O\left(BN_T T^{\frac{4\nu^2+d^2}{2(2\nu-d)\nu}} + BT^2 \min\{T^{\frac{4d}{2\nu-d}}, T^2\}\right)$	$O\left(BN_T T^{\frac{(2\nu+d)^2-2\nu d}{4\nu^2}} + BT^{\frac{2\nu+d}{\nu}}\right)$
SE	Condition	$m_t, M \sim (\log(T))^d$	$m_t, M \sim (\log(T))^d$
	Cost	$O(BN_T T \log^d(T) + BT^2 \log^{2d}(T))$	$O(BN_T T \log^d(T) + BT^2 \log^{2d}(T))$

Table 1: Conditions on the number of features  $m_t$  and inducing variables  $M_T$  required for Theorem 2, alongside the resulting cost of each decoupled sampling method.

## 6 Experiments

We now provide an empirical evaluation of S-GP-TS. As [24] have already comprehensively demonstrated the practical advantage of decoupled sampling for problems with small optimization budgets, we focus here on scalability of S-GP-TS, and in particular a) its efficiency with large batch size, b) its ability to handle large data volumes. We first investigate a collection of classical synthetic problems for BO, before demonstrating S-GP-TS in a challenging real-world high-throughput molecular design considered by [13]. Our synthetic experiments focus on multi-modal problems with substantial observation noise, as these cannot be solved accurately with a small budget yet are still unsuitable for local, exhaustive, or deterministic optimization routines. Our implementation is provided as part of the open-source toolbox `trieste` [50]<sup>2</sup> and relies also on `gpflow` [51] and `gpflux` [52].

As is often the case, our regret-based analysis applies to a version of the algorithm that is slightly different to a practically viable BO method. Rather than focusing on recreating our algorithm exactly in the very limited settings (e.g. a 1-d RBF kernel for which we can calculate eigen-features exactly) that are of little interest to the BO community, we instead choose to demonstrate the practical strength and unprecedented scalability of S-GP-TS by investigating an implementation that could be used by practitioners. The resulting algorithms demonstrated in this section are still well-aligned with our work through their use of sparse GP surrogate models and decoupled Thompson sampling.

### 6.1 Synthetic Benchmarks

We first consider two toy problems: Hartmann (6 dim, moderately multi-modal) with a large additive noise and Shekel (4 dim, highly multi-modal) with moderate noise, see Appendix C for the full description. Our SVGP models use inducing points and a Matérn kernel with smoothness parameter

<sup>2</sup><https://github.com/secondmind-labs/trieste>

$\nu = 2.5$ . As eigenfunctions for this kernel are limited to small dimensions [39], we implement decoupled TS using the easily accessible random Fourier Features (RFF). Note that [24] have shown decoupled sampling to significantly alleviate the *variance starvation* phenomenon (underestimating the variance of points far from the observations [16, 43]) that typically hampers the efficacy of RFFs. We use  $M = 1000$  features and maximise each sample as in (3) using L-BFGS-B, starting from the best point among a large sample.

As sampling inducing points from a k-DPP is prohibitively costly for the repeated model fitting required by BO loops, we use the greedy variance selection method of [26] which is  $\epsilon_0$  close to k-DPP and has been shown to outperform optimisation of inducing points in practice. We also consider the practical alternative of choosing inducing points chosen by a k-means clustering of the observations. As the optimisation progresses, observations are likely to be concentrated in the optimal regions, so clustering would result in somehow “targeted” inducing points for BO. In order to control the computational overhead of S-GP-TS and to allow an efficient computational implementation (i.e. avoiding Tensorflow recompilation issues), we use a fixed number  $m_t$  of points, set to either 250 or 500. Similarly, we set the covariance scaling parameter  $\alpha_t = 1$  to avoid having dynamic tunable parameters, like those that plague UCB-based approaches.

For each experiment, we run  $t = 50$  steps of S-GP-TS with  $B = 100$  (i.e. 5,000 total observations). For baselines, we compare against  $t = 750$  steps of standard sequential non-batch BO routines with an exact GP model: Expected Improvement [EI, 53], Augmented Expected Improvement [AEI, 54], and an extension of Max-value Entropy search suitable for noisy observations [GIBBON, 10]. Due to the large number of steps, we only consider low-cost but high-performance acquisitions, following the cost-benefit analysis of [10], and exclude the popular knowledge gradient [9] or classical entropy search [55, 40]. Popular existing batch acquisition functions do not scale to batches as large as  $B = 100$ , however, we present their performance on smaller batches across additional experiments in Appendix C. We report simple regret of the current believed best solution (maximizer of the current model mean) across the previously queried data points. All results are averaged over 30 runs and reported as a function of either the number of function evaluations ( $tB$  for S-GP-TS and  $t$  for the baselines), or the number of BO iterations, in Figure 1.

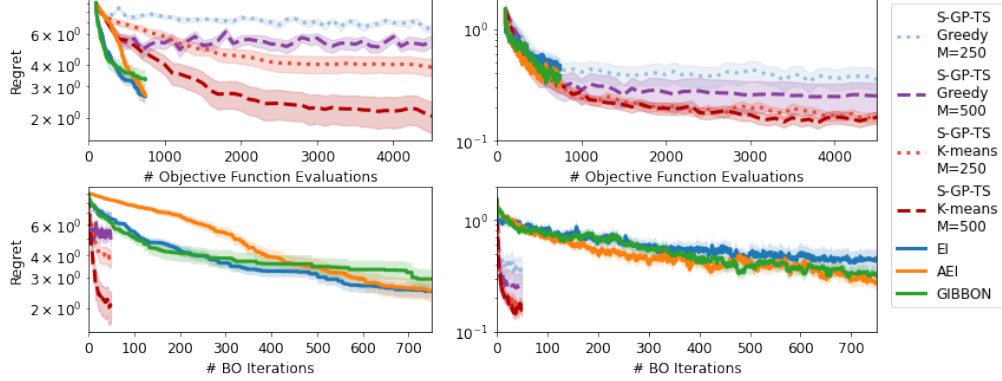


Figure 1: Simple regret on Shekel (4D, left) and Hartmann (6D, right). When considering regret with respect to the total number of objective function evaluations  $tB$  (top panels, purely sequential setting), all S-GP-TS methods are initially less efficient (Shekel) or match the performance (Hartmann) of the best baselines, however the best S-GP-TS approach is able to efficiently allocate its additional budget to achieve lower final regret. When considering regret with respect to the BO iteration (bottom panels, idealised parallel setting), S-GP-TS achieves low regret in a fraction of the iterations required by the standard BO routines.

The fact that S-GP-TS is able to find solutions on both benchmarks with substantially improved regret than found by standard BO, provides strong evidence that S-GP-TS is effectively leveraging parallel resources. Moreover, as these higher-quality solutions were only found after large number of total evaluations, Figure 1 also highlights the necessity for BO routines, like S-GP-TS, that can handle these larger (heavily parallelized) optimization budgets. We reiterate that the existing BO baselines cannot handle as many evaluations as S-GP-TS, becoming prohibitively slow once we surpass 750 data-points). When considering the regret achieved per individual function evaluation, we typically expect batch routines to be less efficient than purely sequential BO routines. However, in the case of the Hartmann function (the benchmark with the largest observation noise), we see that our best



S-GP-TS exactly matches (before going on to exceed) the performance of the sequential routines, suggesting that S-GP-TS is a particularly effective optimizer for functions with significant levels of observation noise.

Note that the performance of S-GP-TS is sensitive to its chosen inducing points, with k-means providing the most effective routines. On Hartmann, 250 inducing points is sufficient to deliver good performances, while on Shekel, which is much more multimodal, using a larger number is critical.

## 6.2 High-throughput Molecular Search

Finally, we investigate the performance of S-GP-TS with respect to an established baseline for high-throughput molecular screening. Although molecular search has been tackled many times with BO [56, 22, 57], only the approach of [13] - standard (non-decoupled) TS over a Bayesian neural network (BNN-TS) - is truly scalable. We now recreate the largest experiment considered by [13], where the objective is to uncover the top 10% of molecules in terms of power conversion efficiency among a library of 2.3 million candidate from the Harvard Clean Energy Project [58]. Molecules are encoded as Morgan circular fingerprints of Bond radius 3 (i.e. 512-dimensional bit vectors, see [59]).

As the standard GP kernels considered above are not suitable for sparse and high-dimensional molecule inputs [60], we instead build our SVGP with a zeroth order ArcCosine kernel [61], chosen due to its strong empirical performance under sparsity and as it permits a random decomposition that can be exploited to perform decoupled TS. In particular, we use the  $M$ -feature decomposition investigated by [62] of

$$k_{arc}(\mathbf{x}, \mathbf{x}') = 2 \int d\mathbf{w} \frac{e^{-\frac{\|\mathbf{w}\|^2}{2}}}{(2\pi)^{d/2}} \Theta(\mathbf{w}^T \mathbf{x}) \Theta(\mathbf{w}^T \mathbf{x}') \approx \frac{2}{M} \sum_{j=1}^M \Theta(\mathbf{w}_j^T \mathbf{x}) \Theta(\mathbf{w}_j^T \mathbf{x}'),$$

where  $\Theta(\cdot)$  is the Heaviside step function and  $\mathbf{w}_j \sim \mathcal{N}(0, I)$ .

In our experiments, we use  $M = 1\,000$  random features and, to avoid memory issues, we compute our GP samples over a random subset of 100 000 molecules (renewed at each sample). We run S-GP-TS twice, once with  $m_t = 500$  and once with 2000 inducing points. We chose inducing points as uniform samples from the already evaluated molecules (for each model step), as preliminary experiments showed that neither the k-means nor greedy selection routines discussed above were effective when applied to sparse and high-dimensional molecular fingerprint inputs.

Following [13], we report the recall (fraction of the top 10% of molecules so far chosen by the BO loop) for S-GP-TS, along with the performance of BNN-TS, a greedy BNN (that queries the  $B$  maximizers of the BNN’s posterior mean), and a random search baseline (all taken from [13]). All routines (including our S-GP-TS) are ran for  $t = 250$  successive batches of  $B = 500$  molecules. Figure 2 shows that S-GP-TS is able to perform effective batch optimization over very large optimization budgets (120,000 total evaluations) and, when using  $m = 2000$  or even just  $m = 500$  inducing points, S-GP-TS matches the performance of [13]’s BNN-based TS and greedy sampling approaches, respectively. Note that due to the high computational demands of this experiment, we report just a single replication of S-GP-TS (a limitation also of [13]’s results). However, we stress that an additional realization of the  $m = 500$  experiment returned indistinguishable results.

## 7 Discussion

We have shown that S-GP-TS enjoys the same regret order as exact GP-TS but with a greatly reduced  $O(N_t M)$  computation per step  $t$ , compared to the  $O(N_t^3)$  cost of the standard sampling. However, the discretization size  $N_t$  is exponential in the dimension  $d$  of the search space and so remains a limiting computational factor when optimizing over high dimensional search spaces. Hence, while S-GP-TS with decoupled sampling rule allows orders of magnitude larger optimization budgets compared to vanilla GP-TS, it still suffers from the *curse of dimensionality*. Intuitively, this seems inevitable due to NP-Hardness of non-convex optimization problems [see, e.g., 63] as required to find the maximizer of the GP-UCB acquisition function [see, e.g., 30], or even in the application of UCB to linear bandits [64]. In particular, the computational cost of the *state-of-the-art* adaptive sketching method for implementing GP-UCB [30] was reported as  $O(N_T d_{\text{eff}}^2)$  where  $d_{\text{eff}}$ , referred to as the effective dimension of the problem, is upper bounded by  $\gamma_T$ .

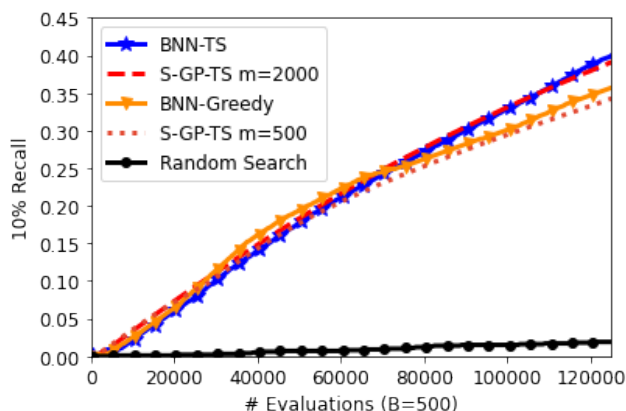


Figure 2: Proportion of the top 10% of molecules found by each of the search routines. S-GP-TS is able to process substantial data volumes and effectively allocates large batches, matching the performance of the well-established BNN baselines.

An important practical consideration when using S-GP-TS in practice is how to choose its inducing points. The performance improvement provided by choosing inducing points by k-means rather than greedy variance selection, as demonstrated in our experiments, raises the possibility that BO-specific routines for choosing inducing points could allow even better performance. This is an important avenue for future work.

## References

- [1] William Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 1933.
- [2] Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P. Adams, and Nando de Freitas. Taking the human out of the loop: A review of Bayesian optimization. *Proceedings of the IEEE*, 2016.
- [3] Sayak Ray Chowdhury and Aditya Gopalan. On kernelized multi-armed bandits. In *International Conference on Machine Learning*, 2017.
- [4] Ricardo Baptista and Matthias Poloczek. Bayesian optimization of combinatorial structures. In *International Conference on Machine Learning*, 2018.
- [5] David Eriksson and Matthias Poloczek. Scalable constrained bayesian optimization. In *International Conference on Artificial Intelligence and Statistics*, 2021.
- [6] Samuel Daulton, Shaun Singh, Vashist Avadhanula, Drew Dimmery, and Eytan Bakshy. Thompson sampling for contextual bandit problems with auxiliary safety constraints. *arXiv preprint arXiv:1911.00638*, 2019.
- [7] Clément Chevalier and David Ginsbourger. Fast computation of the multi-points expected improvement with applications in batch selection. In *International Conference on Learning and Intelligent Optimization*, 2013.
- [8] Javier González, Zhenwen Dai, Philipp Hennig, and Neil Lawrence. Batch Bayesian optimization via local penalization. In *Artificial intelligence and statistics*, 2016.
- [9] Jian Wu and Peter I Frazier. The parallel knowledge gradient method for batch Bayesian optimization. In *Advances in Neural Information Processing Systems*, 2016.
- [10] Henry B Moss, David S Leslie, Javier Gonzalez, and Paul Rayson. Gibbon: General-purpose information-based bayesian optimisation. *arXiv preprint arXiv:2102.03324*, 2021.
- [11] Hamed Jalali, Inneke Van Nieuwenhuijse, and Victor Picheny. Comparison of kriging-based algorithms for simulation optimization with heterogeneous noise. *European Journal of Operational Research*, 2017.

- [12] Mickaël Binois, Jiangeng Huang, Robert B Gramacy, and Mike Ludkovski. Replication or exploration? sequential design for stochastic simulation experiments. *Technometrics*, 2019.
- [13] José Miguel Hernández-Lobato, James Requeima, Edward O Pyzer-Knapp, and Alán Aspuru-Guzik. Parallel and distributed Thompson sampling for large-scale accelerated exploration of chemical space. In *International Conference on Machine Learning*, 2017.
- [14] Kirthivasan Kandasamy, Akshay Krishnamurthy, Jeff Schneider, and Barnabás Póczos. Parallelised Bayesian optimisation via Thompson sampling. In *International Conference on Artificial Intelligence and Statistics*, 2018.
- [15] Jasper Snoek, Oren Rippel, Kevin Swersky, Ryan Kiros, Nadathur Satish, Narayanan Sundaram, Mostofa Patwary, Mr Prabhat, and Ryan Adams. Scalable Bayesian optimization using deep neural networks. In *International conference on machine learning*, pages 2171–2180, 2015.
- [16] Zi Wang, Clement Gehring, Pushmeet Kohli, and Stefanie Jegelka. Batched large-scale Bayesian optimization in high-dimensional spaces. In *International Conference on Artificial Intelligence and Statistics*, 2018.
- [17] Carl E Rasmussen and Christopher KI Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- [18] Peter J Diggle, Jonathan A Tawn, and Rana A Moyeed. Model-based geostatistics. *Journal of the Royal Statistical Society: Series C*, 1998.
- [19] David Eriksson, Michael Pearce, Jacob Gardner, Ryan D Turner, and Matthias Poloczek. Scalable global optimization via local Bayesian optimization. In *Advances in Neural Information Processing Systems*, 2019.
- [20] M. K. Titsias. Variational learning of inducing variables in sparse Gaussian processes. In *International Conference on Artificial Intelligence and Statistics*, 2009.
- [21] Mitchell McIntire, Daniel Ratner, and Stefano Ermon. Sparse gaussian processes for bayesian optimization. In *Association for Uncertainty in Artificial Intelligence*, 2016.
- [22] Ryan-Rhys Griffiths and José Miguel Hernández-Lobato. Constrained Bayesian optimization for automatic chemical design. *arXiv preprint arXiv:1709.05501*, 2017.
- [23] Ang Yang, Cheng Li, Santu Rana, Sunil Gupta, and Svetha Venkatesh. Sparse spectrum Gaussian process for Bayesian optimization. *arXiv preprint arXiv:1906.08898*, 2019.
- [24] James T. Wilson, Viacheslav Borovitskiy, Alexander Terenin, Peter Mostowsky, and Marc Peter Deisenroth. Efficiently sampling functions from Gaussian process posteriors. *International Conference on Machine Learning*, 2020.
- [25] My Phan, Yasin Abbasi Yadkori, and Justin Domke. Thompson sampling and approximate inference. In *Advances in Neural Information Processing Systems*, 2019.
- [26] David Burt, Carl Edward Rasmussen, and Mark Van Der Wilk. Rates of convergence for sparse variational Gaussian process regression. In *International Conference on Machine Learning*, 2019.
- [27] Daniel J. Russo, Benjamin Van Roy, Abbas Kazerouni, Ian Osband, and Zheng Wen. A tutorial on thompson sampling. *Foundational Trends in Machine Learning*, 2018.
- [28] D. Russo and B. Van Roy. Learning to optimize via posterior sampling. *Mathematics of Operations Research*, 2014.
- [29] Niranjan Srinivas, Andreas Krause, Sham Kakade, and Matthias Seeger. Gaussian process optimization in the bandit setting: no regret and experimental design. In *International Conference on Machine Learning*, 2010.
- [30] Daniele Calandriello, Luigi Carratino, Alessandro Lazaric, Michal Valko, and Lorenzo Rosasco. Gaussian process optimization with adaptive sketching: scalable and no regret. In *Conference on Learning Theory*, 2019.

- [31] Alain Berlinet and Christine Thomas-Agnan. *Reproducing kernel Hilbert spaces in probability and statistics*. Springer Science & Business Media, 2011.
- [32] James Hensman, Nicolás Fusi, and Neil D. Lawrence. Gaussian processes for big data. In *Uncertainty in Artificial Intelligence (UAI 2013)*, 2013.
- [33] James Hensman, Nicolas Durrande, Arno Solin, et al. Variational Fourier features for Gaussian processes. *Journal of Machine Learning Research*, 2017.
- [34] Vincent Dutordoir, Nicolas Durrande, and James Hensman. Sparse Gaussian Processes with Spherical Harmonic Features. In *Proceedings of the 37th International Conference on Machine Learning*, 2020.
- [35] Miguel Lázaro-Gredilla and Anibal Figueiras-Vidal. Inter-domain Gaussian processes for sparse inference using inducing features. In *Advances in Neural Information Processing Systems*, 2009.
- [36] Motonobu Kanagawa, Philipp Hennig, Dino Sejdinovic, and Bharath K Sriperumbudur. Gaussian processes and kernel methods: A review on connections and equivalences. *arXiv preprint arXiv:1807.02582*, 2018.
- [37] Viacheslav Borovitskiy, Alexander Terenin, Peter Mostowsky, and Marc P. Deisenroth. Matern Gaussian processes on Riemannian manifolds. In *Advances in Neural Information Processing Systems*, 2020.
- [38] Huaiyu Zhu, Christopher KI Williams, Richard Rohwer, and Michal Morciniec. Gaussian regression and optimal finite dimensional linear models, 1997.
- [39] Arno Solin and Simo Särkkä. Hilbert space methods for reduced-rank Gaussian process regression. *Statistics and Computing*, 2020.
- [40] José Miguel Hernández-Lobato, Matthew W Hoffman, and Zoubin Ghahramani. Predictive entropy search for efficient global optimization of black-box functions. *Advances in Neural Information Processing Systems*, 2014.
- [41] José Miguel Hernández-Lobato, Matthew W Hoffman, and Zoubin Ghahramani. Predictive entropy search for efficient global optimization of black-box functions. In *Advances in Neural Information Processing Systems*, 2014.
- [42] Salomon Bochner et al. *Lectures on Fourier integrals*. Princeton University Press, 1959.
- [43] Mojmir Mutny and Andreas Krause. Efficient high dimensional Bayesian optimization with additivity and quadrature fourier features. In *Advances in Neural Information Processing Systems 31*, 2018.
- [44] Thomas M Cover. *Elements of information theory*. John Wiley & Sons, 1999.
- [45] Jasper Snoek, Hugo Larochelle, and Ryan P Adams. Practical Bayesian optimization of machine learning algorithms. In *Advances in Neural Information Processing Systems 25*, 2012.
- [46] Gabriele Santin and Robert Schaback. Approximation of eigenfunctions in kernel-based spaces. *Advances in Computational Mathematics*, 2016.
- [47] Mikhail Belkin. Approximation beats concentration? an approximation view on inference with smooth radial kernels. In *Conference On Learning Theory*, 2018.
- [48] Gabriel Riutort-Mayol, Paul-Christian Burkner, Michael R. Andersen, Arno Solin, and Aki Vehtari. Practical Hilbert space approximate Bayesian Gaussian processes for probabilistic programming. *arXiv preprint arXiv:2004.11408*, 2020.
- [49] Sattar Vakili, Kia Khezeli, and Victor Picheny. On information gain and regret bounds in Gaussian process bandits. In *International Conference on Artificial Intelligence and Statistics*, 2021.

- [50] Joel Berkeley, Henry B. Moss, Artem Artemev, Sergio Pascual-Diaz, Uri Granta, Hrvoje Stojic, Ivo Couckuyt, Jixiang Quing, Loka Satrio, and Victor Picheny. Trieste, 2021.
- [51] Alexander G de G Matthews, Mark van der Wilk, Tom Nickson, Keisuke Fujii, Alexis Boukouvalas, Pablo León-Villagr , Zoubin Ghahramani, and James Hensman. Gpflow: A gaussian process library using tensorflow. *Journal of Machine Learning Research*, 2017.
- [52] Vincent Dutordoir, Hugh Salimbeni, Eric Hambro, John McLeod, Felix Leibfried, Artem Artemev, Mark van der Wilk, James Hensman, Marc P Deisenroth, and ST John. Gpflux: A library for deep Gaussian processes. *arXiv preprint arXiv:2104.05674*, 2021.
- [53] Donald R Jones, Matthias Schonlau, and William J Welch. Efficient global optimization of expensive black-box functions. *Journal of Global optimization*, 1998.
- [54] Deng Huang, Theodore T Allen, William I Notz, and Ning Zeng. Global optimization of stochastic black-box systems via sequential kriging meta-models. *Journal of global optimization*, 2006.
- [55] Philipp Hennig and Christian J Schuler. Entropy search for information-efficient global optimization. *Journal of Machine Learning Research*, 2012.
- [56] Rafael G mez-Bombarelli, Jorge Aguilera-Iparraguirre, Timothy D Hirzel, David Duvenaud, Dougal Maclaurin, Martin A Blood-Forsythe, Hyun Sik Chae, Markus Einzinger, Dong-Gwang Ha, Tony Wu, et al. Design of efficient molecular organic light-emitting diodes by a high-throughput virtual screening and experimental approach. *Nature Materials*, 2016.
- [57] Henry B Moss, Daniel Beck, Javier Gonz lez, David S Leslie, and Paul Rayson. Boss: Bayesian optimization over string spaces. *Advances in Neural Information Processing Systems*, 2020.
- [58] Johannes Hachmann, Roberto Olivares-Amaya, Sule Atahan-Evrenk, Carlos Amador-Bedolla, Roel S S nchez-Carrera, Aryeh Gold-Parker, Leslie Vogt, Anna M Brockway, and Al n Aspuru-Guzik. The Harvard clean energy project: large-scale computational screening and design of organic photovoltaics on the world community grid. *The Journal of Physical Chemistry Letters*, 2011.
- [59] David Rogers and Mathew Hahn. Extended-connectivity fingerprints. *Journal of chemical information and modeling*, 2010.
- [60] Henry B Moss and Ryan-Rhys Griffiths. Gaussian process molecule property prediction with flowmo. *Advances in Neural Information Processing Systems: Workshop on Machine Learning for Molecules.*, 2020.
- [61] Youngmin Cho. *Kernel methods for deep learning*. PhD thesis, UC San Diego, 2012.
- [62] Kurt Cutajar, Edwin V Bonilla, Pietro Michiardi, and Maurizio Filippone. Random feature expansions for deep Gaussian processes. In *International Conference on Machine Learning*, 2017.
- [63] Prateek Jain and Purushottam Kar. Non-convex optimization for machine learning. *Foundational Trends in Machine Learning*, 2017.
- [64] V. Dani, T. P. Hayes, and S. M. Kakade. Stochastic linear optimization under bandit feedback. In *Conference on Learning Theory*, 2008.
- [65] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [66] Rodolphe Le Riche and Victor Picheny. Revisiting Bayesian optimization in the light of the coco benchmark. *arXiv preprint arXiv:2103.16649*, 2021.
- [67] Dong C Liu and Jorge Nocedal. On the limited memory bfgs method for large scale optimization. *Mathematical programming*, 1989.



## A Complements on SVGPs

As discussed in Section 3, SVGPs approximate the posterior of exact GPs through either a set of inducing points  $\mathbf{Z}_t \triangleq \{z_1, \dots, z_m\}$  or through a set of inducing features  $\phi_m(x) \triangleq \{\phi_1(x), \dots, \phi_m(x)\}$ . The resulting inducing features, defined as  $u_{t,i} = \hat{f}(z_{t,i})$  (for inducing points) or  $u_{t,i} = \int_{\mathcal{X}} \hat{f}(x) \phi_i(x) dx$  (for inducing features), are assumed to follow a prior Gaussian density  $q_t(\mathbf{u}_t) = \mathcal{N}(\mathbf{m}_t, \mathbf{S}_t)$ . We now discuss how to set these variational parameters  $\mathbf{m}_t$  and  $\mathbf{S}_t$  for a given dataset.

For SVGPS, the posterior mean and covariance is given in closed form as

$$\mu_t^{(s)}(x) = k_{\mathbf{Z}_t, x}^T K_{\mathbf{Z}_t, \mathbf{Z}_t}^{-1} \mathbf{m}_t \quad k_t^{(s)}(x, x') = k(x, x') + k_{\mathbf{Z}_t, x}^T K_{\mathbf{Z}_t, \mathbf{Z}_t}^{-1} (\mathbf{S}_t - K_{\mathbf{Z}_t, \mathbf{Z}_t}) K_{\mathbf{Z}_t, \mathbf{Z}_t}^{-1} k_{\mathbf{Z}_t, x'}$$

and

$$\mu_t^{(s)}(x) = \phi_{m_t}^T(x) \mathbf{m}_t \quad k_t^{(s)}(x, x') = k(x, x') + \phi_{m_t}^T(x) (\mathbf{S}_t - \Lambda_{m_t}) \phi_{m_t}(x'),$$

for the inducing point and inducing features representations, respectively. See Section 3 or [26] for more details. However, the marginal likelihood for these models is intractable, and so, as is common practice in variational inference methods, we set of our variational parameters (as-well as the SVGP's kernel parameters) to maximize instead the tractable Evidence-based Lower Bound (ELBO).

For inducing point SVGPS, the ELBO can be written as

$$\text{ELBO}(t) = -\frac{1}{2} \mathbf{y}_t^T (Q_t + \tau \mathbf{I}_t)^{-1} \mathbf{y}_t - \frac{1}{2} \log |Q_t + \tau \mathbf{I}_t| - \frac{t}{2} \log(2\pi) - \frac{\theta_t}{2\tau},$$

where  $Q_t = K_{\mathbf{Z}_t, \mathbf{X}_t}^T K_{\mathbf{Z}_t, \mathbf{Z}_t}^{-1} K_{\mathbf{Z}_t, \mathbf{X}_t}$ ,  $K_{\mathbf{Z}_t, \mathbf{X}_t} = [k_{z_i, x_j}]_{i,j}$ ,  $i = 1, \dots, m_t$ ,  $j = 1, \dots, t$ ,  $\mathbf{I}_t$  is the  $t \times t$  identity matrix and  $\theta_t = \text{Tr}(K_{\mathbf{X}_t, \mathbf{X}_t} - Q_t)$ . See [32] for a full derivation.

For inducing feature SVGP, the expression of ELBO is the same but with  $Q_t = K_{\phi_{m_t}, \mathbf{X}_t}^T \Lambda_{m_t}^{-1} K_{\phi_{m_t}, \mathbf{X}_t}$ ,  $K_{\phi_{m_t}, \mathbf{X}_t} = [\lambda_i \phi_i(x_j)]_{i,j}$ ,  $i = 1, \dots, m_t$ ,  $j = 1, \dots, t$ .

To optimize the ELBO in practice, [32] proposed a numerical solution allowing for mini-batching [see also 26] and the use of stochastic gradient descent algorithms such as Adam [65]. In addition, [20] provides an explicit solution for the convex optimization problem of finding  $(\mathbf{m}_t, \mathbf{S}_t)$ , allowing more involved alternate optimization schemes.

## B Detailed Proofs

In this section, we provide detailed proofs for Theorem 1, Lemma 1, Proposition 1 and Theorem 2, in order.

### B.1 Proof of Theorem 1

Before presenting the proof of Theorem 1, we first overview the regret bound for vanilla GP-TS [3, Theorem 4].

*The Existing Regret Bound for Vanilla GP-TS.* [3] proved that, with probability at least  $1 - \delta$ ,  $|f(x) - \mu_t(x)| \leq u_t \sigma_t(x)$ , where  $u_t = \left( B + R \sqrt{2(\gamma_t + 1 + \log(1/\delta))} \right)$  and  $\gamma_t$  is the maximal information gain. Based on this concentration inequality, [3] showed that the regret of GP-TS scales with the cumulative uncertainty at the observation points measured by the standard deviation:  $O(\sum_{t=1}^T u_t \sigma_{t-1}(x_t))$ . Furthermore, [29] showed that  $\sum_{i=1}^t \sigma_{i-1}^2(x_i) \leq \gamma_t$ . Using this result and applying Cauchy-Schwarz inequality to  $O(\sum_{t=1}^T u_t \sigma_{t-1}(x_t))$ , [3] proved that  $R(T, 1; f) = O\left(\gamma_T \sqrt{T \log(T)}\right)$ , for vanilla GP-TS.

□

We build on the analysis of GP-TS in [3] to prove the regret bounds for S-GP-TS. We stress that despite some similarities in the proof, the analysis of standard GP-TS does not extend to S-GP-TS. This proof characterizes the behavior of the upper bound on regret in terms of the approximation

constants, namely  $\underline{a}, \bar{a}, c$  and  $\epsilon$ . A notable difference is that the additive approximation error in the posterior standard deviation ( $\epsilon_t$ ) can cause under-exploration which is an issue the analysis of exact GP-TS cannot address. In addition, we account for the effect of batch sampling on the regret bounds.

We first focus on the instantaneous regret at each time  $t$  within the discrete set,  $f(\mathbf{x}^{*(t)}) - f(x_{t,b})$ . Recall  $\mathbf{x}^{*(t)} \triangleq \operatorname{argmin}_{x' \in D_t} \|x^* - x'\|$  from Assumption 5. It is then easy to upper bound the cumulative regret by the cumulative value of  $f(\mathbf{x}^{*(t)}) - f(x_{t,b}) + \frac{1}{t^2}$  as our discretization ensures that  $f(x^*) - f(\mathbf{x}^{*(t)}) \leq \frac{1}{t^2}$ . For upper bounds on instantaneous regret, we start with concentration of GP samples  $\tilde{f}_{t,b}$  around their predicted values and the concentration of the prediction around the true objective function. We then consider the anti-concentration around the optimum point. The necessary anti-concentration may fail due to approximation error in the standard deviation around the optimum point. We thus consider two cases of low and sufficiently high standard deviation at  $\mathbf{x}^{*(t)}$  separately. While a low standard deviation implies good prediction at  $\mathbf{x}^{*(t)}$ , a sufficiently high standard deviation guarantees sufficient exploration. We use these results to upper bound the instantaneous regret at each time  $t$  with uncertainties measured by the standard deviation.

#### Concentration events $\mathcal{E}_t$ and $\tilde{\mathcal{E}}_t$ :

**Define**  $\mathcal{E}_t$  as the event that at time  $t$ , for all  $x \in D_t$ ,  $|f(x) - \tilde{\mu}_{t-1}(x)| \leq \frac{1}{2}\alpha_t(\tilde{\sigma}_{t-1}(x) + \epsilon_t)$ . Recall  $\alpha_t = 2\tilde{u}_t(1/(t^2))$ . Applying lemma 1, we have  $\Pr[\mathcal{E}_t] \geq 1 - \frac{1}{t^2}$ .

**Define**  $\tilde{\mathcal{E}}_t$  as the event that for all  $x \in D_t$ , and for all  $b \in [B]$ ,  $|\tilde{f}_{t,b}(x) - \tilde{\mu}_{t-1}(x)| \leq \alpha_t b_t \tilde{\sigma}_{t-1}(x)$  where  $b_t = \sqrt{2\ln(BN_t t^2)}$ . We have  $\Pr[\tilde{\mathcal{E}}_t] \geq 1 - \frac{1}{t^2}$ .

*Proof.* For a fixed  $x \in D_t$ , and a fixed  $b \in [B]$ ,

$$\Pr \left[ |\tilde{f}_{t,b}(x) - \tilde{\mu}_{t-1}(x)| > \alpha_t b_t \tilde{\sigma}_{t-1}(x) \right] < \exp\left(-\frac{b_t^2}{2}\right) = \frac{1}{BN_t t^2}.$$

The inequality holds because of the following bound on the CDF of a normal random variable  $1 - \text{CDF}_{\mathcal{N}(0,1)}(c) \leq \frac{1}{2} \exp(-\frac{c^2}{2})$  and the observation that  $\frac{\tilde{f}_{t,b}(x) - \tilde{\mu}_{t-1}(x)}{\alpha_t \tilde{\sigma}_{t-1}(x)}$  has a normal distribution.

Applying a union bound we get  $\Pr[\tilde{\mathcal{E}}_t] \leq \frac{1}{t^2}$  which gives us the bound on probability of  $\tilde{\mathcal{E}}_t$ .  $\square$

We thus proved  $\mathcal{E}_t$  and  $\tilde{\mathcal{E}}_t$  are high probability events. This will facilitate the proof by conditioning on  $\mathcal{E}_t$  and  $\tilde{\mathcal{E}}_t$ . Also notice that when both  $\mathcal{E}_t$  and  $\tilde{\mathcal{E}}_t$  hold true, we have, for all  $x \in D_t$ , and for all  $b \in [B]$ ,

$$|\tilde{f}_{t,b}(x) - f(x)| \leq \beta_t \tilde{\sigma}_{t-1}(x) + \frac{1}{2}\alpha_t \epsilon_t \quad (7)$$

where  $\beta_t = \alpha_t(b_t + \frac{1}{2})$ .

**Anti Concentration Bounds.** It is standard in the analysis of TS methods to prove sufficient exploration using an anti-concentration bound. That establishes a lower bound on the probability of a sample being sufficiently large (so that the corresponding point is likely to be selected by TS rule). For this purpose, we use the following bound on the CDF of a normal distribution:  $1 - \text{CDF}_{\mathcal{N}(0,1)}(c) \geq \frac{\exp(-\frac{c^2}{2})}{4c\sqrt{\pi}}$ . The underestimation of the posterior standard deviation at the optimum point however might result in an under exploration. On the other hand, a low standard deviation at the optimum point implies a low prediction error. We use this observation in our regret analysis by considering the two cases separately. Specifically, the regret  $f(\mathbf{x}^{*(t)}) - f(x_{t,b})$  at each time  $t$  for each sample  $b$  is bounded differently under the conditions: I.  $\tilde{\sigma}_{t-1}(\mathbf{x}^{*(t)}) \leq \epsilon_t$  and II.  $\tilde{\sigma}_{t-1}(\mathbf{x}^{*(t)}) > \epsilon_t$ .

**Under Condition I** ( $\tilde{\sigma}_{t-1}(\mathbf{x}^{*(t)}) \leq \epsilon_t$ ), when both  $\mathcal{E}_t$  and  $\tilde{\mathcal{E}}_t$  hold true, we have

$$\begin{aligned}
& f(\mathbf{x}^{*(t)}) - f(x_{t,b}) \\
& \leq \tilde{f}_{t,b}(\mathbf{x}^{*(t)}) + \beta_t \tilde{\sigma}_{t-1}(\mathbf{x}^{*(t)}) + \frac{1}{2} \alpha_t \epsilon_t \\
& \quad - \tilde{f}_t(x_{t,b}) + \beta_t \tilde{\sigma}_{t-1}(x_{t,b}) + \frac{1}{2} \alpha_t \epsilon_t \quad \text{by (7),} \\
& \leq \beta_t \tilde{\sigma}_{t-1}(\mathbf{x}^{*(t)}) + \beta_t \tilde{\sigma}_{t-1}(x_{t,b}) + \alpha_t \epsilon_t \quad \text{by the selection rule of TS,} \\
& \leq \beta_t \tilde{\sigma}_{t-1}(x_{t,b}) + (\beta_t + \alpha_t) \epsilon_t \quad \text{by Condition I.}
\end{aligned} \tag{8}$$

that upper bounds the instantaneous regret at time  $t$  by a factor of approximate standard deviation up to an additive term caused by approximation error. Since  $f(\mathbf{x}^{*(t)}) - f(x_{t,b}) \leq 2B$ , under Condition I,

$$\mathbb{E}[f(\mathbf{x}^{*(t)}) - f(x_{t,b})] \leq \beta_t \tilde{\sigma}_{t-1}(x_{t,b}) + (\beta_t + \alpha_t) \epsilon_t + \frac{4B}{t^2}. \tag{9}$$

where the inequality holds by  $\Pr[\bar{\mathcal{E}}_t \text{ or } \tilde{\mathcal{E}}_t] \leq \frac{2}{t^2}$ .

**Under Condition II** ( $\tilde{\sigma}_{t-1}(\mathbf{x}^{*(t)}) > \epsilon_t$ ), we can show sufficient exploration by anti-concentration at the optimum point. In particular under Condition II, if  $\mathcal{E}_t$  holds true, we have

$$\Pr[\tilde{f}_{t,b}(\mathbf{x}^{*(t)}) > f(\mathbf{x}^{*(t)})] \geq p, \tag{10}$$

where  $p = \frac{1}{4\sqrt{\pi}}$ .

*Proof.* Applying the anti-concentration of a normal distribution

$$\begin{aligned}
\Pr[\tilde{f}_{t,b}(\mathbf{x}^{*(t)}) > f(\mathbf{x}^{*(t)})] &= \Pr\left[\frac{\tilde{f}_{t,b}(\mathbf{x}^{*(t)}) - \tilde{\mu}_{t-1}(\mathbf{x}^{*(t)})}{\alpha_t \tilde{\sigma}_{t-1}(\mathbf{x}^{*(t)})} > \frac{f(\mathbf{x}^{*(t)}) - \tilde{\mu}_{t-1}(\mathbf{x}^{*(t)})}{\alpha_t \tilde{\sigma}_{t-1}(\mathbf{x}^{*(t)})}\right] \\
&\geq p.
\end{aligned}$$

As a result of the observation that the right hand side of the inequality inside the probability argument is upper bounded by 1:

$$\begin{aligned}
\frac{f(\mathbf{x}^{*(t)}) - \tilde{\mu}_{t-1}(\mathbf{x}^{*(t)})}{\alpha_t \tilde{\sigma}_{t-1}(\mathbf{x}^{*(t)})} &\leq \frac{\frac{1}{2} \alpha_t \tilde{\sigma}_{t-1}(\mathbf{x}^{*(t)}) + \frac{1}{2} \alpha_t \epsilon_t}{\alpha_t \tilde{\sigma}_{t-1}(\mathbf{x}^{*(t)})} \quad \text{By } \mathcal{E}_t \\
&\leq 1. \quad \text{By Condition II} \quad \square
\end{aligned}$$

**Sufficiently Explored Points.** Let  $\mathcal{S}_t$  denote the set of sufficiently explored points which are unlikely to be selected by S-GP-TS if  $\tilde{f}_{t,b}(\mathbf{x}^{*(t)})$  is higher than  $f(\mathbf{x}^{*(t)})$ . Specifically, we use the notation

$$\mathcal{S}_t = \{x \in D_t : f(x) + \beta_t \tilde{\sigma}_{t-1}(x) + \frac{1}{2} \alpha_t \epsilon_t \leq f(\mathbf{x}^{*(t)})\}. \tag{11}$$

Recall  $\beta_t = \alpha_t(b_t + \frac{1}{2})$ . In addition, we define

$$\bar{x}_t = \operatorname{argmin}_{x \in D_t \setminus \mathcal{S}_t} \tilde{\sigma}_{t-1}(x). \tag{12}$$

We showed in equation (8) that the instantaneous regret can be upper bounded by the sum of standard deviations at  $x_{t,b}$  and  $\mathbf{x}^{*(t)}$ . The standard method based on information gain can be used to bound the cumulative standard deviations at  $x_{t,b}$ . This is not sufficient however because the cumulative standard deviations at  $\mathbf{x}^{*(t)}$  do not converge unless there is sufficient exploration around  $x^*$ . To address this, we use  $\bar{x}_t$  as an intermediary to be able to upper bound the instantaneous regret by a factor of  $\tilde{\sigma}_{t-1}(x_{t,b})$  through the following lemma.

**Lemma 2.** Under Condition II, for  $t \geq \sqrt{\frac{2}{p}}$ , if  $\mathcal{E}_t$  holds true

$$\tilde{\sigma}_{t-1}(\bar{x}_t) \leq \frac{2}{p} \mathbb{E}[\tilde{\sigma}_{t-1}(x_{t,b})], \tag{13}$$

where the expectation is taken with respect to the randomness in the sample  $\tilde{f}_{t,b}$ .

*Proof of Lemma 2.* First notice that when both  $\mathcal{E}_t$  and  $\tilde{\mathcal{E}}_t$  hold true, for all  $x \in S_t$

$$\begin{aligned}\tilde{f}_{t,b}(x) &\leq f(x) + \beta_t \tilde{\sigma}_{t-1}(x) + (\alpha_t - 1)\epsilon_t && \text{by (7)} \\ &\leq f(\mathbf{x}^{*(t)}), && \text{by definition of } S_t.\end{aligned}\quad (14)$$

Also, if  $\tilde{f}_{t,b}(\mathbf{x}^{*(t)}) > \tilde{f}_{t,b}(x), \forall x \in S_t$ , the rule of selection in TS ( $x_{t,b} = \operatorname{argmax}_{x \in \mathcal{X}} \tilde{f}_{t,b}(x)$ ) ensures  $x_{t,b} \in D_t \setminus S_t$ . So we have

$$\begin{aligned}\Pr[x_{t,b} \in D_t \setminus S_t] &\geq \Pr[\tilde{f}_{t,b}(\mathbf{x}^{*(t)}) > \tilde{f}_{t,b}(x), \forall x \in S_t] \\ &\geq \Pr[\tilde{f}_{t,b}(\mathbf{x}^{*(t)}) > \tilde{f}_{t,b}(x), \forall x \in S_t, \tilde{\mathcal{E}}_t] - \Pr[\tilde{\mathcal{E}}_t] \\ &\geq \Pr[\tilde{f}_{t,b}(\mathbf{x}^{*(t)}) > f(\mathbf{x}^{*(t)})] - \Pr[\tilde{\mathcal{E}}_t] && \text{by (14)} \\ &\geq p - \frac{1}{t^2} && \text{by (10)} \\ &\geq \frac{p}{2}, && \text{for } t \geq \sqrt{2/p}.\end{aligned}$$

Finally, we have

$$\begin{aligned}\mathbb{E}[\tilde{\sigma}_{t-1}(x_{t,b})] &\geq \mathbb{E}\left[\tilde{\sigma}_{t-1}(x_{t,b}) \middle| x_{t,b} \in D_t \setminus S_t\right] \Pr[x_{t,b} \in D_t \setminus S_t] \\ &\geq \frac{p\tilde{\sigma}_{t-1}(\bar{x}_t)}{2},\end{aligned}\quad (15)$$

where the expectation is taken with respect to the randomness in the sample  $\tilde{f}_{t,b}$  at time  $t$ .  $\square$

Now we are ready to bound the simple regret under Condition II using  $\bar{x}_t$  as an intermediary. Under Condition II, when both  $\mathcal{E}_t$  and  $\tilde{\mathcal{E}}_t$  hold true,

$$\begin{aligned}f(\mathbf{x}^{*(t)}) - f(x_{t,b}) &= f(\mathbf{x}^{*(t)}) - f(\bar{x}_t) + f(\bar{x}_t) - f(x_{t,b}) \\ &\leq \beta_t \tilde{\sigma}_{t-1}(\bar{x}_t) + \frac{1}{2}\alpha_t \epsilon_t + f(\bar{x}_t) - f(x_{t,b}) && \text{by definition of } S_t \\ &\leq \beta_t \tilde{\sigma}_{t-1}(\bar{x}_t) + \frac{1}{2}\alpha_t \epsilon_t \\ &\quad + \tilde{f}_{t,b}(\bar{x}_t) + \beta_t \tilde{\sigma}_{t-1}(\bar{x}_t) - \tilde{f}_{t,b}(x_{t,b}) + \beta_t \tilde{\sigma}_{t-1}(x_{t,b}) + \alpha_t \epsilon_t && \text{by (7)} \\ &\leq \beta_t (2\tilde{\sigma}_{t-1}(\bar{x}_t) + \tilde{\sigma}_{t-1}(x_{t,b})) + \frac{3}{2}\alpha_t \epsilon_t, && \text{by the rule of selection in TS.}\end{aligned}$$

Thus, since  $f(\mathbf{x}^*) - f(x_{t,b}) \leq 2B$ , under Condition II, for  $t \geq \sqrt{\frac{2}{p}}$

$$\mathbb{E}[f(\mathbf{x}^{*(t)}) - f(x_{t,b})] \leq \frac{(4+p)\beta_t}{p} \mathbb{E}[\tilde{\sigma}_{t-1}(x_{t,b})] + \frac{3}{2}\alpha_t \epsilon_t + \frac{4B}{t^2} \quad (16)$$

where we used Lemma 2 and  $\Pr[\bar{\mathcal{E}}_t \text{ or } \tilde{\mathcal{E}}_t] \leq \frac{2}{t^2}$ .

**Upper bound on regret.** From the upper bounds on instantaneous regret under Condition I and Condition II we conclude that, for  $t \geq \sqrt{\frac{2}{p}}$

$$\begin{aligned}\mathbb{E}[f(\mathbf{x}^{*(t)}) - f(x_{t,b})] &\leq \max \left\{ \beta_t \tilde{\sigma}_{t-1}(x_{t,b}) + (\beta_t + \alpha_t)\epsilon_t + \frac{4B}{t^2}, \right. \\ &\quad \left. \frac{(4+p)\beta_t}{p} \mathbb{E}[\tilde{\sigma}_{t-1}(x_{t,b})] + \frac{3}{2}\alpha_t \epsilon_t + \frac{4B}{t^2} \right\} \\ &\leq \frac{(4+p)\beta_t}{p} \mathbb{E}[\tilde{\sigma}_{t-1}(x_{t,b})] + (\beta_t + \alpha_t)\epsilon_t + \frac{4B}{t^2}.\end{aligned}\quad (17)$$

We can now upper bound the cumulative regret. Noticing  $\lceil \sqrt{\frac{2}{p}} \rceil = 4$ .

$$\begin{aligned}
R(T, B; f) &= \sum_{t=1}^T \sum_{b=1}^B \mathbb{E}[f(x^*) - f(x_{t,b})] \\
&= \sum_{t=1}^4 \sum_{b=1}^B \mathbb{E}[f(x^*) - f(x_{t,b})] + \sum_{t=5}^T \sum_{b=1}^B \mathbb{E}[f(x^*) - f(x_{t,b})] \\
&\leq 8BB + \sum_{t=5}^T (\mathbb{E}[f(x^{*(t)}) - f(x_{t,b})] + \frac{1}{t^2}) \\
&\leq 8BB + \sum_{t=5}^T \sum_{b=1}^B \left( \frac{(4+p)\beta_t}{p} \mathbb{E}[\tilde{\sigma}_{t-1}(x_{t,b})] + (\beta_t + \alpha_t)\epsilon_t + \frac{4B+1}{t^2} \right) \\
&\leq 8BB + \frac{\pi^2 B(4B+1)}{6} + \frac{(4+p)\beta_T}{p} \sum_{t=1}^T \sum_{b=1}^B \mathbb{E}[\tilde{\sigma}_{t-1}(x_{t,b})] + (\beta_T + \alpha_T) \sum_{t=1}^T \sum_{b=1}^B \epsilon_t \\
&\leq 15BB + 2B + 30\beta_T \sum_{t=1}^T \sum_{b=1}^B (\bar{a}\mathbb{E}[\sigma_{t-1}(x_{t,b})] + \epsilon_t) + (\beta_T + \alpha_T)\epsilon TB \\
&\leq 15BB + 2B + 30\bar{a}\beta_T \sum_{t=1}^T \sum_{b=1}^B \mathbb{E}[\sigma_{t-1}(x_{t,b})] + 30\beta_T\epsilon TB + (\beta_T + \alpha_T)\epsilon TB \\
&\leq 15BB + 2B + 30\bar{a}\beta_T \sum_{t=1}^T \sum_{b=1}^B \mathbb{E}[\sigma_{t-1}(x_{t,b})] + (31\beta_T + \alpha_T)\epsilon TB.
\end{aligned}$$

We simplified the expressions by  $\frac{4+p}{p} \leq 30$ ,  $\frac{4\pi^2}{6} \leq 7$  and  $\frac{\pi^2}{6} \leq 2$ .

We now use a technique based on information gain to upper bound  $\sum_{t=1}^T \sum_{b=1}^B \mathbb{E}[\sigma_{t-1}(x_{t,b})]$  as formalized in the following lemma.

**Lemma 3.** *For all batch observation sequences  $\{x_{t,b}\}_{t \in [T], b \in [B]}$ , we have*

$$\sum_{t=1}^T \sum_{b=1}^B \sigma_{t-1}(x_{t,b}) \leq B \sqrt{\frac{2T\gamma_T}{\log(1 + \frac{1}{\tau})}} \quad (18)$$

*Proof of Lemma 3.* Without loss of generality assume that at each time instance  $t = 1, 2, \dots, T$ , the batch observations are ordered such that  $\sigma_{t-1}(x_{t,1}) \geq \sigma_{t-1}(x_{t,b})$ , for all  $b \in [B]$ . We thus have

$$\sum_{t=1}^T \sum_{b=1}^B \sigma_{t-1}(x_{t,b}) \leq B \sum_{t=1}^T \sigma_{t-1}(x_{t,1}). \quad (19)$$

For the sequence of observations  $\{x_{t,1}\}_{t=1}^T$ , define the conditional posterior mean and variance

$$\begin{aligned}
\bar{\mu}_t(x) &= \mathbb{E}[\hat{f}(x) | \{x_{s,1}\}_{s=1}^t] \\
\bar{\sigma}_t^2(x) &= \mathbb{E}[(\hat{f}(x) - \bar{\mu}_t(x))^2 | \{x_{s,1}\}_{s=1}^t].
\end{aligned}$$

By the expression of posterior variance of multivariate Gaussian random variables and by positive definiteness of the covariance matrix, we know that conditioning on a larger set reduces the posterior variance. Thus  $\bar{\sigma}_t(x) \geq \sigma_t(x)$ . Notice that  $\sigma_t(x)$  is the posterior variance conditioned on full batches of the observations while  $\bar{\sigma}_t(x)$  is the posterior variance conditioned on only the first observation at each batch. We thus have

$$\sum_{t=1}^T \sigma_{t-1}(x_{t,1}) \leq \sum_{t=1}^T \bar{\sigma}_{t-1}(x_{t,1}) \quad (20)$$



We can now follow the standard steps in bounding the cumulative standard deviation in the non-batch setting. In particular using Cauchy-Schwarz inequality, we have

$$\sum_{t=1}^T \bar{\sigma}_{t-1}(x_{t,1}) \leq \sqrt{T \sum_{t=1}^T \bar{\sigma}_{t-1}^2(x_{t,1})}. \quad (21)$$

In addition, [29] showed that

$$\sum_{t=1}^T \bar{\sigma}_{t-1}^2(x_{t,1}) \leq \frac{2\gamma_T}{\log(1 + \frac{1}{\tau})}. \quad (22)$$

Combining (19), (20), (21) and (22), we arrive at the lemma.  $\square$

We thus have

$$R(T; \text{S-GP-TS}) \leq 30\bar{a}\beta_TB \sqrt{\frac{2T\gamma_T}{\log(1 + \frac{1}{\tau})}} + (31\beta_T + \alpha_T)\epsilon TB + 15BB + 2B \quad (23)$$

which can be simplified to

$$R(T; \text{S-GP-TS}) = \tilde{O}\left(\underline{a}\bar{a}(1+c)B\sqrt{T\gamma_T} + \underline{a}^2(1+c^2)\epsilon TB\right). \quad (24)$$

$\square$

## B.2 Proof of Lemma 1

It remains to prove the concentration inequality for the approximate statistics given in Lemma 1.

*Proof of Lemma 1.* By triangle inequality we have

$$\begin{aligned} |f(x) - \tilde{\mu}_t(x)| &\leq |f(x) - \mu_t(x)| + |\tilde{\mu}_t(x) - \mu_t(x)| \\ &\leq |f(x) - \mu_t(x)| + c_t\sigma_t(x) \quad \text{by Assumptions 4.} \end{aligned}$$

From Theorem 2 of [3], with probability at least  $1 - \delta$ ,

$$f(x) - \mu_t(x) \leq \left(B + R\sqrt{2(\gamma_t + 1 + \log(1/\delta))}\right)\sigma_t(x).$$

Thus,

$$\begin{aligned} |f(x) - \tilde{\mu}_t(x)| &\leq \left(B + R\sqrt{2(\gamma_t + 1 + \log(1/\delta))}\right)\sigma_t(x) + c_t\sigma_t(x) \\ &\leq \underline{a}_t(B + R\sqrt{\frac{2\ln(1/\delta)}{\tau}} + c_t)(\tilde{\sigma}_t(x) + \epsilon_t), \end{aligned}$$

where the last inequality holds by Assumption 3.  $\square$

## B.3 Proof of Proposition 1

Here, we use  $\tilde{\mu}_t$  and  $\tilde{\sigma}_t$  to specifically denote the approximate posterior mean and the approximate posterior standard deviations of the decomposed sampling rules (4) and (5) in contrast to Sec. 5.1 where we used the notation more generally for any approximate model. We also use  $\mu_t^{(s)}$  and  $\sigma_t^{(s)}$  to refer to the posterior mean and the posterior standard deviation of SVGP models, and  $\mu^{(w)}$  and  $\sigma^{(w)}$  to refer to the priors generated from an  $M$ -truncated feature vector. For the approximate posterior mean, we have  $\tilde{\mu}_t = \mu_t^{(s)}$ . However, the approximate posterior standard deviations  $\sigma^{(s)}$  and  $\tilde{\sigma}$  are not the same.

By the triangle inequality we have

$$|\tilde{\sigma}_t(x) - \sigma_t(x)| \leq |\tilde{\sigma}_t(x) - \sigma_t^{(s)}(x)| + |\sigma_t^{(s)}(x) - \sigma_t(x)|. \quad (25)$$

For the first term, following the exact same lines as in the proof of Proposition 7 in [24], we have

$$|\tilde{\sigma}_t^2(x) - \sigma_t^{(s)2}(x)| \leq C_1 m_t |\sigma^2(x) - \sigma^{(w)2}(x)| \quad (26)$$

where  $C_1 = \max_{1 \leq t \leq T} (1 + \|K_{\mathbf{Z}_{m_t}, \mathbf{Z}_{m_t}}^{-1}\|_{C(\mathcal{X}^2)})$ . [24] proceed to upper bound  $|\sigma^2(x) - \sigma^{(w)2}(x)|$  by a constant divided by  $\sqrt{M}$ . We use a tighter bound based on feature representation of the kernel. Specifically from definition of  $\delta_M$  we have that

$$\begin{aligned} |\sigma^2(x) - \sigma^{(w)2}(x)| &\leq \sum_{i=M+1}^{\infty} \lambda_i \bar{\phi}_i^2 \\ &= \delta_M, \end{aligned} \quad (27)$$

which results in the following upper bound

$$|\tilde{\sigma}_t^2(x) - \sigma_t^{(s)2}(x)| \leq C_1 m_t \delta_M. \quad (28)$$

For the standard deviations we have

$$\begin{aligned} |\tilde{\sigma}_t(x) - \sigma_t^{(s)}(x)| &= \sqrt{|\tilde{\sigma}_t(x) - \sigma_t^{(s)}(x)|^2} \\ &\leq \sqrt{|\tilde{\sigma}_t(x) - \sigma_t^{(s)}(x)| |\tilde{\sigma}_t(x) + \sigma_t^{(s)}(x)|} \\ &= \sqrt{|\tilde{\sigma}_t^2(x) - \sigma_t^{(s)2}(x)|} \\ &\leq \sqrt{C_1 m_t \delta_M}, \end{aligned} \quad (29)$$

where the first inequality holds because  $|\tilde{\sigma}_t(x) - \sigma_t^{(s)}(x)| \leq |\tilde{\sigma}_t(x) + \sigma_t^{(s)}(x)|$  for positive  $\tilde{\sigma}_t(x)$  and  $\sigma_t^{(s)}(x)$ .

We can efficiently bound the error in the SVGP approximation based on the convergence of SVGP methods. Let us first focus on the inducing features. It was shown that (Lemma 2 in [26]), for the SVGP with inducing features

$$\text{KL} \left( \text{GP}(\mu_t, \sigma_t), \text{GP}(\mu_t^{(s)}, k_t^{(s)}) \right) \leq \frac{\theta_t}{\tau}. \quad (30)$$

where  $\text{GP}(\mu_t, \sigma_t)$  and  $\text{GP}(\mu_t^{(s)}, k_t^{(s)})$  are the true and the SVGP approximate posterior distributions at time  $t$ , and KL denotes the Kullback-Leibler divergence between them. On the right hand side,  $\theta_t$  is the trace of the error in the covariance matrix. Specifically,  $\theta_t = \text{Tr}(E_t)$  where  $E_t = K_{\mathbf{X}_t, \mathbf{X}_t} - K_{\mathbf{Z}_t, \mathbf{X}_t}^T K_{\mathbf{Z}_t, \mathbf{Z}_t} K_{\mathbf{Z}_t, \mathbf{X}_t}$ . Using the Mercer expansion of the kernel matrix, [26] showed that  $[E_t]_{i,i} = \sum_{j=m_t+1}^{\infty} \lambda_j \phi_j^2(x_i)$ . Thus

$$\begin{aligned} \theta_t &= \sum_{i=1}^t \sum_{j=m+1}^{\infty} \lambda_j \phi_j^2(x_i) \\ &\leq t \sum_{j=m_t+1}^{\infty} \lambda_j \bar{\phi}_j^2 \\ &= t \delta_{m_t} \end{aligned} \quad (31)$$

Thus,

$$\text{KL} \left( \text{GP}(\mu_t, \sigma_t), \text{GP}(\mu_t^{(s)}, k_t^{(s)}) \right) \leq \kappa_t / 2. \quad (32)$$

where  $\kappa_t = 2tB\delta_m/\tau$  that is determined by the number of current observations. In comparison, [26] proceed by introducing a prior distribution on  $x_i$  and bounding  $[E_t]_{i,i}$ , differently.

For the case of inducing points drawn from an  $\epsilon_0$  close k-DPP distribution, similarly following the exact lines as [26] except for the upper bound on  $[E_t]_{i,i}$ , with probability at least  $1 - \delta$ , (32) holds with  $\kappa_t = \frac{2tB(m_t+1)\delta_{m_t}}{\delta\tau} + \frac{4tB\epsilon_0}{\delta\tau}$  where  $\epsilon_0$  that is determined by the number of current observations.

In addition, if the KL divergence between two Gaussian distributions is bounded by  $\kappa_t/2$ , we have the following bound on the means and variances of the marginals [Proposition 1 in [26]]

$$\begin{aligned} |\mu_t^{(s)}(x) - \mu_t(x)| &\leq \sigma_t(x)\sqrt{\kappa_t}, \\ |1 - \frac{\sigma_t^{(s)2}(x)}{\sigma_t^2(x)}| &\leq \sqrt{3\kappa_t}, \end{aligned} \quad (33)$$

which by algebraic manipulation gives

$$\sqrt{1 - \sqrt{3\kappa_t}\sigma_t(x)} \leq \sigma_t^{(s)}(x) \leq \sqrt{1 + \sqrt{3\kappa_t}\sigma_t(x)} \quad (34)$$

Combining the bounds on  $\sigma_t^{(s)}$  with (29), we get

$$\sqrt{1 - \sqrt{3\kappa_t}\sigma_t(x)} - \sqrt{C_1 m_t \delta_M} \leq \tilde{\sigma}_t(x) \leq \sqrt{1 + \sqrt{3\kappa_t}\sigma_t(x)} + \sqrt{C_1 m_t \delta_M}$$

Comparing this bound with Assumption 3, we have  $\underline{a}_t = \frac{1}{\sqrt{1 - \sqrt{3\kappa_t}}}$ ,  $\bar{a}_t = \sqrt{1 + \sqrt{3\kappa_t}}$ , and  $\epsilon_t = \sqrt{C_1 m_t \delta_M}$ . Also, since  $\mu_t^{(s)} = \tilde{\mu}_t$ , comparing (33) with Assumption 4, we have  $c_t = \sqrt{\kappa_t}$ .

□

#### B.4 Proof of Theorem 2

In Theorem 1, we proved that

$$R(T, B; f) = O\left(\underline{a}\bar{a}BR\sqrt{d\gamma_T(\gamma_{TB} + \log(T))T\log(T)} + \underline{a}\epsilon TBR\sqrt{d(\gamma_{TB} + \log(T))\log(T)}\right)$$

We thus need to show that  $\underline{a}\bar{a}$  is a constant independent of  $T$  and  $\underline{a}\epsilon$  is small so that the second term is dominated by the first term.

In the case of Matérn kernel,  $\lambda_j = O(j^{-\frac{2\nu+d}{d}})$  implies that  $\delta_m = O(m^{-\frac{2\nu}{d}})$ . Under sampling rule (4), we select  $\delta = \frac{1}{T}$  and  $\epsilon_0 = \frac{1}{T^2 \log(T)}$  in Proposition 1. We thus need  $\kappa_T = O(T^2 m_T \delta_{m_T})$  and  $\epsilon_T \sqrt{T} = O(\sqrt{m_T \delta_M T})$  be sufficiently small constants. That is achieved by selecting  $m_T = T^{\frac{2d}{2\nu-d}}$  and  $M = T^{\frac{(2\nu+d)d}{2(2\nu-d)\nu}}$ .

Under sampling rule (5), we need  $\kappa_T = O(T\delta_{m_T})$  and  $\epsilon_T \sqrt{T} = O(\sqrt{m_T \delta_M T})$  be sufficiently small constants. That is achieved by selecting  $m_T = T^{\frac{d}{2\nu}}$  and  $M = \frac{(2\nu+d)d}{4\nu^2}$ .

In the case of SE kernel,  $\lambda_j = O(\exp(-j^{\frac{1}{d}}))$  implies that  $\delta_m = O(\exp(-m^{\frac{1}{d}}))$ . Under sampling rule (4), we select  $\delta = \frac{1}{T}$  and  $\epsilon_0 = \frac{1}{T^2 \log(T)}$  in Proposition 1. We thus need  $\kappa_T = O(T^2 m_T \delta_{m_T})$  and  $\epsilon_T \sqrt{T} = O(\sqrt{m_T \delta_M T})$  be sufficiently small constants. That is achieved by selecting  $m_T = (\log(T))^d$  and  $M = (\log(T))^d$ . We obtain the same results under sampling rule (5) where we need  $\kappa_T = O(T\delta_{m_T})$  and  $\epsilon_T \sqrt{T} = O(\sqrt{m_T \delta_M T})$  be sufficiently small constants.

□

## C Additional Experiments and Experimental Details

In Section 6, we tested S-GP-TS across popular synthetic benchmarks from the BO literature. We considered the Shekel, Hartmann and Ackley (see Figure 3) functions, each contaminated by Gaussian noise with variance 0.1, 0.5 and 0.5, respectively. Note that for Hartmann and Ackley, we chose our observation noise to be an order of magnitude larger than usually considered for these problems in order to demonstrate the suitability of S-GP-TS for controlling large optimization budgets (as required to optimize these highly noisy functions). We now provide explicit forms for these synthetic functions and list additional experimental details left out from the main paper.

**Shekel function.** A four-dimensional function with ten local and one global minima defined on  $\mathcal{X} \in [0, 10]^4$ :

$$f(\mathbf{x}) = - \sum_{i=1}^{10} \left( \sum_{j=1}^4 (x_j - A_{j,i})^2 + \beta_i \right)^{-1},$$

where

$$\beta = \begin{pmatrix} 1 \\ 2 \\ 2 \\ 4 \\ 4 \\ 6 \\ 3 \\ 7 \\ 5 \\ 5 \end{pmatrix} \quad \text{and} \quad A = \begin{pmatrix} 4 & 1 & 8 & 6 & 3 & 2 & 5 & 8 & 6 & 7 \\ 4 & 1 & 8 & 6 & 7 & 9 & 3 & 1 & 2 & 3.6 \\ 4 & 1 & 8 & 6 & 3 & 2 & 5 & 8 & 6 & 7 \\ 4 & 1 & 8 & 6 & 7 & 9 & 3 & 1 & 2 & 3.6 \end{pmatrix}.$$

**Ackley function.** A five-dimensional function with many local minima surrounding a single global minima defined on  $\mathcal{X} \in [-2, 1]^5$ :

$$f(\mathbf{x}) = -20 \exp \left( -0.2 * \sqrt{\frac{1}{4} \sum_{i=1}^d x_i^2} \right) - \exp \left( \frac{1}{4} \sum_{i=1}^4 \cos(2\pi x_i) \right) + 20 + \exp(1).$$

**Hartmann 6 function.** A six-dimensional function with six local minima and a single global minima defined on  $\mathcal{X} \in [0, 1]^6$ :

$$f(\mathbf{x}) = - \sum_{i=1}^4 \alpha_i \exp \left( - \sum_{j=1}^6 A_{i,j} (x_j - P_{i,j})^2 \right),$$

where

$$A = \begin{pmatrix} 10 & 3 & 17 & 3.5 & 1.7 & 8 \\ 0.05 & 10 & 17 & 0.1 & 8 & 14 \\ 3 & 3.5 & 1.7 & 10 & 17 & 8 \\ 17 & 8 & 0.05 & 10 & 0.1 & 14 \end{pmatrix}, \quad \alpha = \begin{pmatrix} 1 \\ 1.2 \\ 3 \\ 3.2 \end{pmatrix},$$

$$P = 10^{-4} \begin{pmatrix} 1312 & 1696 & 5569 & 124 & 8283 & 5886 \\ 2329 & 4135 & 8307 & 3736 & 1004 & 9991 \\ 2348 & 1451 & 3522 & 2883 & 3047 & 6650 \\ 4047 & 8828 & 8732 & 5743 & 1091 & 381 \end{pmatrix}.$$

For all our synthetic experiments (both for S-GP-TS and the baseline BO methods), we follow the implementation advice of [66] regarding constraining length-scales (to stabilize model fitting) and by maximizing acquisition functions (and Thompson samples) using L-BFGS [67] starting from the best location found across a random sample of  $500 * d$  locations (where  $d$  is the problem dimension). Our SVGP models are fit with an ADAM optimizer [65] with an initial learning rate of 0.1, ran for at most 10,000 iterations but with an early stopping criteria (if 100 successive steps lead to a loss less than 0.1). We also implemented a learning rate reduction factor of 0.5 with a patience of 10. Our implementation of the GIBBON acquisition function follows [10] and is built on 10 Gumbel samples built across a grid of  $10,000 * d$  query points. For BO's initialization step, our S-GP-TS models are given a single random sample of the same size as the considered batches and standard BO routines are given  $d + 4$  initial samples (again following the advice of [66]). The function evaluations required for these initialization are included in our Figures.

### C.1 S-GP-TS on the Ackley Function

To supplement the synthetic examples included in the main body of the paper, we now consider the performance of S-GP-TS when used to optimize the challenging Ackley function, defined over 5 dimensions and under very high levels of observation noise (Gaussian with variance 0.5). The Ackley function (in 5D) has thousands of local minima and a single global optima in the centre. As this global optima has a very small volume, achieving high precision optimization on this benchmark requires high levels of exploration (akin to an active learning task). Figure 3 demonstrates the performance of S-GP-TS on the Ackley benchmark, where we see that S-GP-TS is once again able to find solutions with lower regret than the sequential benchmarks and effectively allocate batch resources. In contrast to our other experiments, where the K-means inducing point selection routine significantly outperforms greedy variance reduction, our Ackley experiment shows little difference between the different inducing point selection routines. In fact, greedy variance selection slightly outperforms selection by k-means. We hypothesize that the strong repulsion properties of DPPs (as approximated by greedy variance selection) are advantageous for optimization problems requiring high levels of exploration.

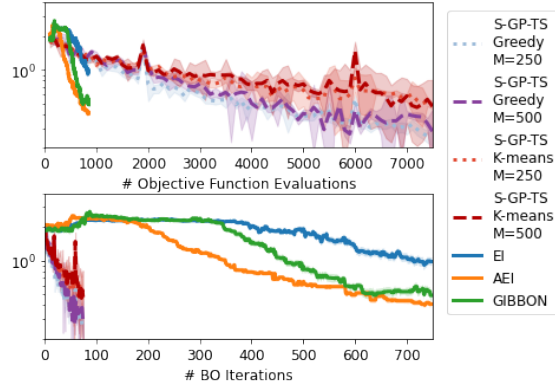


Figure 3: Simple regret on 5D Ackley function. The best S-GP-TS approaches are able to efficiently allocate additional optimization budgets to achieve lower final regret than the sequential baselines. When considering regret with respect to the BO iteration (bottom panels, idealised parallel setting), S-GP-TS achieves low regret in a fraction of the iterations required by standard BO routines. For this task, the choice of inducing point selection strategy (and number of inducing points) is not as crucial as for our other synthetic benchmarks, however, greedy variance selection provides a small improvement over selection by k-means.

### C.2 A Comparison of S-GP-TS with other batch BO routines

To accompany Figures 1 and 3 (our comparison of S-GP-TS with sequential BO routines), we also now compare S-GP-TS with popular batch BO routines. Once again, we stress that these existing BO routines do not scale to the large batch sizes that we consider for S-GP-TS, and so we plot their performance for  $B = 25$  (a batch size considered large in the context of these existing BO methods). We consider two well-known batch extensions of EI: Locally Penalized EI [LP, 8] and the multi-point EI (known as qEI) of [7]. We also consider with a recently proposed batch information-theoretic approach known as General-purpose Information-Based Bayesian OptimizationN [GIBBON, 10]. The large optimization budgets considered in these problems prevent our use of batch extensions of other popular but high-cost acquisition functions such as those based on knowledge gradients [9] or entropy search [13]. Figure 4 compares our S-GP-TS methods ( $B=100$ ) with the popular batch routines ( $B=25$ ), where we see that S-GP-TS achieves lower regret than existing batch BO methods for our most noisy synthetic function (Hartmann).



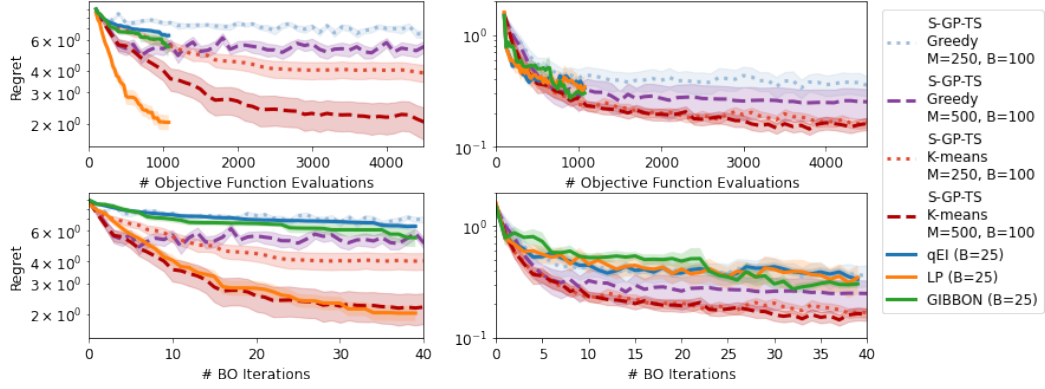


Figure 4: Simple regret on Shekel (4D, left) and Hartmann (6D, right) as a function of either the number of evaluations (top) or BO iterations (bottom). S-GP-TS methods are ran for batches of size  $B = 100$  and the batch BO methods for batches of size  $B = 25$ . We see that S-GP-TS is particularly effective when performing the batch optimization of particularly noisy functions (Hartmann), exceeding the regret of the batch baselines. In our synthetic benchmark with low observation noise (Shekel), S-GP-TS is less efficient in terms of individual function evaluations, however, S-GP-TS’s ability to control larger batches means that it can match the performance of the highly performant LP with respect to the number of BO iterations (the idealised parallel setting).