
Scalable Thompson Sampling using Sparse Gaussian Process Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

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

15 1 Introduction

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

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

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

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

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

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

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

70 2 Problem Formulation

71 We consider the sequential optimization of an unknown function f over a compact set $\mathcal{X} \subset \mathbb{R}^d$.
 72 A sequential learning policy selects a batch of B observation points $\{x_{t,b}\}_{b \in [B]}$ at each time step
 73 $t = 1, 2, \dots, T$ and receives the corresponding real-valued and noisy rewards $\{y_{t,b} = f(x_{t,b}) +$
 74 $\epsilon_{t,b}\}_{b \in [B]}$, where $\epsilon_{t,b}$ denotes the observation noise. Throughout the paper, we use the notation
 75 $[n] = \{1, 2, \dots, n\}$, for $n \in \mathbb{N}$. As is common in both the bandits and GP literature, our analysis
 76 uses the following sub-Gaussianity assumption, a direct consequence of which is that $\mathbb{E}[\epsilon_{t,b}] = 0$, for
 77 all $t, b \in \mathbb{N}$.

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

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

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

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

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

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

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

93 In the case of practically relevant kernels, Assumption 2 implies certain smoothness properties for
 94 the objective functions. For the details on the RKHS and its norm, see the supplementary material.

95 3 Gaussian Processes and Sparse Models

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

100 3.1 Exact Gaussian Process models

101 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
 102 $tB \times d$ matrix of locations with rows $[\mathbf{X}_t]_{(s-1)B+b} = x_{s,b}$, and \mathbf{y}_t is the tB -dimensional column
 103 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
 104 Gaussian observation noise, the posterior of the GP model \hat{f} given the set of past observations \mathcal{H}_t , is
 105 also a GP with mean $\mu_t(\cdot)$, variance $\sigma_t^2(\cdot)$ and kernel function $k_t(\cdot, \cdot)$ specified as

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

106 and $\sigma_t^2(x) = k_t(x, x)$, with $k_{\mathbf{X}_t, x}$ the tB dimensional column vector with entries $[k_{\mathbf{X}_t, x}]_{(s-1)B+b} =$
 107 $k(x_{s,b}, x)$, and $K_{\mathbf{X}_t, \mathbf{X}_t}$ the $tB \times tB$ positive definite covariance matrix with entries
 108 $[K_{\mathbf{X}_t, \mathbf{X}_t}]_{(s-1)B+b, (s'-1)B+b'} = k(x_{s,b}, x_{s',b'})$. We directly see from (2) that accessing the pos-
 109 terior expressions require an $O((tB)^3)$ matrix inversion, which is a computational bottleneck for
 110 large values of tB .

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

114 3.2 Sparse Variational Gaussian Process Models with Inducing Points

115 To overcome the cubic cost of exact GPs, SVGPs [20, 32] instead approximate the GP posterior
 116 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
 117 the *inducing variables* $\mathbf{u}_t = \hat{f}(\mathbf{Z}_t)$ (rather than the tB observations in \mathbf{y}_t) and specifying a prior
 118 Gaussian density $q_t(\mathbf{u}_t) = \mathcal{N}(\mathbf{m}_t, \mathbf{S}_t)$, yields an approximate posterior distribution that, crucially,
 119 is still a GP but with the significantly reduced computational complexity of $O(m_t^2 t)$. The posterior
 120 mean and covariance of the SVGP is given in closed form as

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

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

126 3.3 Sparse Variational Gaussian Process Models with Inducing Features

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

129 features [35] $\psi_1(x), \dots, \psi_{m_t}(x)$, i.e we set our i^{th} inducing variable as $u_{t,i} = \int_{\mathcal{X}} \hat{f}(x) \psi_i(x) dx$.
 130 Courtesy of Mercer’s theorem, we can decompose our chosen kernel k as the inner product of
 131 possibly infinite dimensional feature maps (see Theorem 4.1 in [36]) to provide the expansion
 132 $k(x, x') = \sum_{j=1}^{\infty} \lambda_j \phi_j(x) \cdot \phi_j(x')$ for eigenvalues $\{\lambda_j \in \mathbb{R}^+\}_{j=1}^{\infty}$ and eigenfunctions $\{\phi_j \in H_k\}_{j=1}^{\infty}$.
 133 If we set our inducing features to be the m_t eigenfunctions with largest eigenvalues, it can be shown
 134 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
 135 Process model with posterior mean and covariance given by

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

136 Here, \mathbf{m}_t and \mathbf{S}_t are inducing parameters (as above), $\boldsymbol{\phi}_m(x) \triangleq [\phi_1(x), \dots, \phi_m(x)]^{\text{T}}$ is the truncated
 137 feature vector and Λ_m is the $m \times m$ diagonal matrix of eigenvalues, $[\Lambda_m]_{i,j} = \lambda_i \delta_{i,j}$.

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

142 4 Scalable Thompson Sampling using Gaussian Process Models (S-GP-TS)

143 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
 144 distribution of \hat{f} and finding their maximizers, i.e. we select samples $x_{t,b}$ satisfying

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

145 However, since $\hat{f}_{t,b}$ is an infinite dimensional object, generating such samples is computationally
 146 challenging. Consequently, it is common to resort to approximate strategies, the most simple of
 147 which is to sample across an N_t point discretization D_t of \mathcal{X} [14] which can be obtained with an
 148 $O(N_t^3)$ cost (due to a required Cholesky decomposition).

149 To improve the computational efficiency of TS, a classical strategy [40, 41] is to rely on kernel
 150 decompositions. For instance, a sample \hat{f} from a GP can be expressed as a randomly weighted sum
 151 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,
 152 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
 153 expansions to contain only the M eigenfunctions with largest eigenvalues or M random Fourier
 154 features, we have access to approximate but analytically tractable samples. For both expansions, the
 155 weights w_j are sampled independently from a standard normal distribution. Conditioned on current
 156 tB observations, the posterior distribution of w_j are Gaussian with mean and covariance functions
 157 that can be calculated with an $O(M^3)$ computations, resulting in an $O(M^3 + BNM)$ cost to draw
 158 B Thompson samples.

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

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

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

170 where we have coefficients $v_{t,j} = [K_{\mathbf{z}_t, \mathbf{z}_t}^{-1} (\alpha_t (\mathbf{u}_t - \mathbf{m}_t) + \mathbf{m}_t - \alpha_t \boldsymbol{\Phi}_{m_t, M} \Lambda_M^{\frac{1}{2}} \mathbf{w}_M)]_j$ for $\boldsymbol{\Phi}_{m_t, M} =$
 171 $[\boldsymbol{\phi}_M(z_1), \dots, \boldsymbol{\phi}_M(z_{m_t})]^{\text{T}}$ and $\mathbf{w}_M = [w_1, \dots, w_M]^{\text{T}}$. The weights w_i are drawn i.i.d from $\mathcal{N}(0, 1)$.

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

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

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

178 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 Λ_{m_t} defined in Section 3.3.

179 5 Regret Analysis of S-GP-TS

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

184 5.1 Regret Bounds Based on the Quality of Approximations

185 Consider a TS algorithm using an approximate GP model. In particular, assume an approximate
 186 model is provided where \tilde{k}_t , $\tilde{\sigma}_t$ and $\tilde{\mu}_t$ are approximations of k_t , σ_t and μ_t , respectively. At each
 187 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
 188 $\alpha_t^2 \tilde{k}_{t-1}$. 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
 189 a discretization D_t of the search space.

190 We start our analysis by making two assumptions on the *quality* of approximations $\tilde{\mu}_t$, $\tilde{\sigma}_t$ of the
 191 posterior mean and the standard deviation. This parameterization is agnostic to the particular sampling
 192 rule (governing $\tilde{\mu}_t$ and $\tilde{\sigma}_t$) and provides valuable intuition that can be applied to any approximate
 193 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
 194 assumption are satisfied under some conditions on the value of the parameters of the sampling rules.

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

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

197 *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*
 198 *$t \geq 1$ and some small constant $\epsilon \in \mathbb{R}$.*

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

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

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

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

204 **Lemma 1.** *Under Assumptions 1, 2, 3 and 4, with probability at least $1 - \delta$, $|f(x) - \tilde{\mu}_t(x)| \leq$
 205 $\tilde{u}_t(\tilde{\sigma}_t(x) + \epsilon_t)$, where $\tilde{u}_t(\delta) = \underline{a}_t \left(\mathcal{B} + R\sqrt{2(\gamma_{tB} + 1 + \log(1/\delta))} + c_t \right)$.*

206 Proof is provided in the appendix. Here, γ_s is the *maximal information gain*: $\gamma_s =$
 207 $\max_{A \subset \mathcal{X}, |A|=s} \mathcal{I}([y(x)]_{x \in A}; [\hat{f}(x)]_{x \in A})$, where $\mathcal{I}([y(x)]_{x \in A}; [\hat{f}(x)]_{x \in A})$ denotes the mutual infor-
 208 mation [44, Chapter 2] between observations and the underlying GP model. The maximal information
 209 gain can itself be bounded for a specific kernel (see Sec. 5.3).

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

212 **Assumption 5.** The discretization D_t is designed in a way that $|f(x) - f(\mathbf{x}^{(t)})| \leq 1/t^2$ for all
 213 $x \in \mathcal{X}$, where $\mathbf{x}^{(t)} = \operatorname{argmin}_{x' \in D_t} \|x - x'\|$ is the closest point (in Euclidean norm) to x in D_t .
 214 The size of this discretization satisfies $|D_t| = N_t \leq C(d, B)t^{2d}$ where $C(d, B)$ is independent of t
 215 ([3, 29]).

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

217 **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
 218 regret defined in (1), satisfies

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

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

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

222 5.2 Approximation Quality of the Decomposed Sampling Rule

223 For S-GP-TS with inducing points, we assume, as in [26], that the inducing points are sampled
 224 according to a discrete k-DPP. While this might be costly in practice, [26] showed that \mathbf{Z}_t can be
 225 efficiently sampled from ϵ_0 close sampling methods without compromising the predictive quality of
 226 SVGP. For both sampling rules, we also assume in our analysis that the Mercer decomposition of the
 227 kernel is used.

228 The quality of the approximation can be characterized using the spectral properties of the GP kernel.
 229 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
 230 decaying eigenvalues, including sufficient eigenfunctions in the feature representation results in a
 231 small δ_M . In addition, [26] showed that, for an SVGP, a sufficient number of inducing variables
 232 ensures that the Kullback–Leibler (KL) divergence between the approximate and the true posterior
 233 distributions diminishes. Consequently, the approximate posterior mean and the approximate posterior
 234 variance converge to the true ones. Building on this result, we are able to prove Proposition 1 on the
 235 quality of approximations.

236 **Proposition 1.** For S-GP-TS based on sampling rule (4) with $\alpha_t = 1$ and an SVGP using an ϵ_0 close
 237 k-DPP for selecting \mathbf{Z}_t , with probability at least $1 - \delta$, Assumptions 3 and 4 hold with parameters
 238 $c_t = \sqrt{\kappa_t}$, $\underline{a}_t = \frac{1}{\sqrt{1 - \sqrt{3\kappa_t}}}$, $\bar{a}_t = \sqrt{1 + \sqrt{3\kappa_t}}$, and $\epsilon_t = \sqrt{C_1 m_t \delta_M}$, where C_1 is a constant
 239 specified in the appendix and $\kappa_t = \frac{2tB(m_t+1)\delta_{m_t}}{\tau\delta} + \frac{4tB\epsilon_0}{\tau\delta}$.

240 For S-GP-TS based on sampling rule (5) with $\alpha_t = 1$, Assumptions 3 and 4 hold with parameters
 241 $c_t = \sqrt{\kappa_t}$, $\underline{a}_t = \frac{1}{\sqrt{1 - \sqrt{3\kappa_t}}}$, $\bar{a}_t = \sqrt{1 + \sqrt{3\kappa_t}}$, and $\epsilon_t = \sqrt{C_1 m_t \delta_M}$, where C_1 is the same constant
 242 as above and $\kappa_t = \frac{2tB\delta_{m_t}}{\tau}$.

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

248 5.3 Application of Regret Bounds to Matérn and SE Kernels

249 We now investigate the application of Theorem 1 to the Squared Exponential (SE) and Matérn
 250 kernels, widely used in practice [see, e.g., 17, 45]. In the case of a Matérn kernel with smoothness
 251 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 $\lambda_j =$
 252 $O(\exp(-j^{\frac{1}{d}}))$ [47, 48]. With these bounds on the spectrum of the kernels and the specific bounds
 253 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
 254 Matérn, 49], Theorem 1 and Proposition 1 result in the following theorem.

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

257 With a batch size $B = 1$ Theorem 2 recovers the same regret bounds as the exact GP-TS [3].

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

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

		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	$O(BN_T T \log^d(T) + BT^2 \log^{2d}(T))$	$O(BN_T T \log^d(T) + BT^2 \log^{2d}(T))$

266

267 6 Experiments

268 We now provide an empirical evaluation of S-GP-TS. As [24] have already comprehensively demon-
 269 strated the practical advantage of decoupled sampling for problems with small optimization budgets,
 270 we focus here on scalability of S-GP-TS, and in particular a) its efficiency with large batch size, b) its
 271 ability to handle large data volumes. We first investigate a collection of classical synthetic problems
 272 for BO, before demonstrating S-GP-TS in a challenging real-world high-throughput molecular design
 273 considered by [13]. Our synthetic experiments focus on multi-modal problems with substantial
 274 observation noise, as these cannot be solved accurately with a small budget yet are still unsuitable
 275 for local, exhaustive, or deterministic optimization routines. Our implementation is based on the
 276 open-source toolboxes `gpflow` [50] and `gpflux` [51] for modelling and `trieste` [52] for BO. We
 277 provide an implementation at https://anonymous.4open.science/r/S_GP_TS.

278 6.1 Synthetic Benchmarks

279 We first consider two toy problems: Hartmann (6 dim, moderately multi-modal) with a large additive
 280 noise and Shekel (4 dim, highly multi-modal) with moderate noise, see appendix for the full descrip-
 281 tion. Our SVGP models use inducing points and a Matérn kernel with smoothness parameter $\nu = 2.5$.
 282 As eigenfunctions for this kernel are limited to small dimensions [39], we implement decoupled TS
 283 using the easily accessible random Fourier Features (RFF). Note that [24] have shown decoupled
 284 sampling to significantly alleviate the *variance starvation* phenomenon (underestimating the variance
 285 of points far from the observations [16, 43]) that typically hampers the efficacy of RFFs. We use
 286 $M = 1000$ features and maximise each sample as in (3) using L-BFGS-B, starting from the best
 287 point among a large sample.

288 As generating inducing points using a k-DPP is prohibitively costly for the repeated model fitting
 289 required by BO loops, we use the greedy variance selection method of [26] which is ϵ_0 close to k-DPP
 290 and has been shown to outperform optimisation of inducing points in practice. We also consider the
 291 practical alternative of choosing inducing points chosen by a k-means clustering of the observations.
 292 As the optimisation progresses, observations are likely to be concentrated in the optimal regions,
 293 so clustering would result in somehow “targeted” inducing points for BO. In order to control the
 294 computational overhead of S-GP-TS, we use a fixed number m_t of points, set to either 250 or 500.
 295 We also set the covariance scaling parameter $\alpha_t = 1$.

296 For each experiment, we run $t = 50$ steps of S-GP-TS with $B = 100$ (i.e. 5,000 total observations).
 297 For baselines, we compare against $t = 750$ steps of standard sequential non-batch BO routines with
 298 an exact GP model: Expected Improvement [EI, 53], Augmented Expected Improvement [AEI, 54],
 299 and an extension of Max-value Entropy search suitable for noisy observations [GIBBON, 10]. Due to
 300 the large number of steps, we only consider low-cost but high-performance acquisitions, following
 301 the cost-benefit analysis of [10], and exclude the popular knowledge gradient [9] or classical entropy
 302 search [55, 40]. Popular existing batch acquisition functions do not scale to batches as large as
 303 $B = 100$, however, we present their performance on smaller batches across additional experiments
 304 in our Supplement. We report simple regret of the current believed best solution (maximizer of the
 305 current model mean) across the previously queried data points. All results are averaged over 30 runs
 306 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.

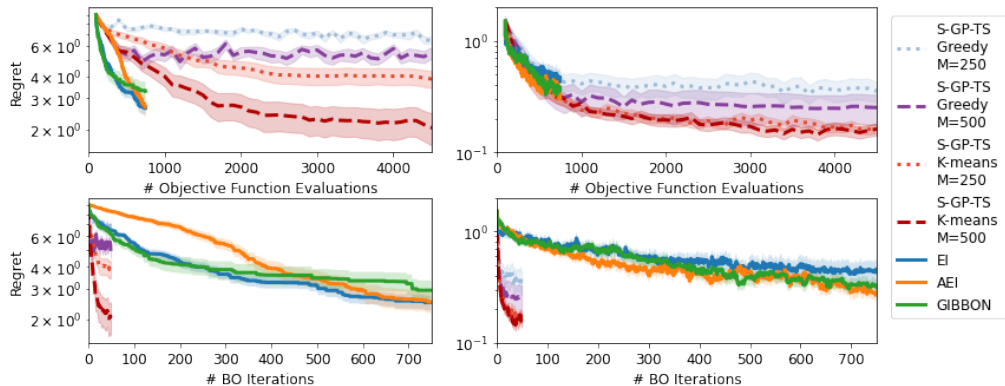


Figure 1: Simple regret on Shekel (4D, left) and Hartmann (6D, right). When considering regret with respect to the total number of objective function evaluations tB (top panels, purely sequential setting), all S-GP-TS methods are initially less efficient (Shekel) or match the performance (Hartmann) of 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.

307 The fact that S-GP-TS is able to find solutions on both benchmarks with substantially improved
 308 regret than found by standard BO, provides strong evidence that S-GP-TS is effectively leveraging
 309 parallel resources. Moreover, as these higher-quality solutions were only found after large number
 310 of total evaluations, Figure 1 also highlights the necessity for BO routines, like S-GP-TS, that
 311 can handle these larger (heavily parallelized) optimization budgets. When considering the regret
 312 achieved per individual function evaluation, we typically expect batch routines to be less efficient
 313 than purely sequential BO routines. However, in the case of the Hartmann function (the benchmark
 314 with the largest observation noise), we see that our best S-GP-TS exactly matches (before going
 315 on to exceed) the performance of the sequential routines, suggesting that S-GP-TS is a particularly
 316 effective optimizer for functions with significant levels of observation noise.
 317

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

321 6.2 High-throughput Molecular Search

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

329 As the standard GP kernels considered above are not suitable for sparse and high-dimensional
 330 molecule inputs [60], we instead build our SVGP with a zeroth order ArcCosine kernel [61], chosen

331 due to its strong empirical performance under sparsity and as it permits a random decomposition
 332 that can be exploited to perform decoupled TS. In particular, we use the M -feature decomposition
 333 investigated by [62] of

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

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

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

341 Following [13], we report the recall (fraction of the top 10% of molecules so far chosen by the
 342 BO loop) for S-GP-TS, along with the performance of BNN-TS, a greedy BNN (that queries
 343 the B maximizers of the BNN’s posterior mean), and a random search baseline (all taken from
 344 [13]). All routines (including our S-GP-TS) are ran for $t = 250$ successive batches of $B = 500$
 345 molecules. Figure 2 shows that S-GP-TS is able to perform effective batch optimization over very large
 346 optimization budgets (120,000 total evaluations) and, when using $m = 2000$ or even just $m = 500$
 347 inducing points, S-GP-TS matches the performance of [13]’s BNN-based TS and greedy sampling
 348 approaches, respectively. Note that due to the high computational demands of this experiment, we
 349 report just a single replication of S-GP-TS (a limitation also of [13]’s results). However, we stress
 350 that an additional realization of the $m_t = 500$ experiment returned indistinguishable results.

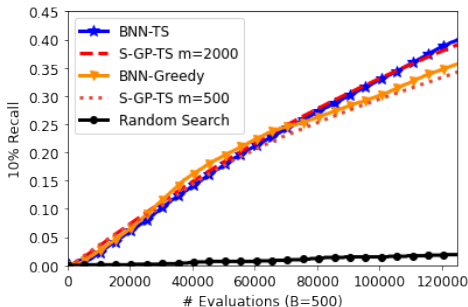


Figure 2: Proportion of the top 10% of molecules found by each of the search routines. S-GP-TS is able to process substantial data volumes and effectively allocates large batches, matching the performance of the well-established BNN baselines.

351 7 Discussion

352 We have shown that S-GP-TS enjoys the same regret order as exact GP-TS but with a greatly reduced
 353 $O(N_t M)$ computation per step t , compared to the $O(N_t^3)$ cost of the standard sampling. However,
 354 the discretization size N_t is exponential in the dimension d of the search space and so remains a
 355 limiting computational factor when optimizing over high dimensional search spaces. Hence, while
 356 S-GP-TS with decoupled sampling rule allows orders of magnitude larger optimization budgets
 357 compared to vanilla GP-TS, it still suffers from the *curse of dimensionality*. Intuitively, this seems
 358 inevitable due to NP-Hardness of non-convex optimization problems [see, e.g., 63] as required to
 359 find the maximizer of the GP-UCB acquisition function [see, e.g., 30], or even in the application
 360 of UCB to linear bandits [64]. In particular, the computational cost of the *state-of-the-art* adaptive
 361 sketching method for implementing GP-UCB [30] was reported as $O(N_T d_{\text{eff}}^2)$ where d_{eff} , referred to
 362 as the effective dimension of the problem, is upper bounded by γ_T .

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

368 **References**

- 369 [1] William Thompson. On the likelihood that one unknown probability exceeds another in view of
370 the evidence of two samples. *Biometrika*, 1933.
- 371 [2] Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P. Adams, and Nando de Freitas. Taking
372 the human out of the loop: A review of Bayesian optimization. *Proceedings of the IEEE*, 2016.
- 373 [3] Sayak Ray Chowdhury and Aditya Gopalan. On kernelized multi-armed bandits. In *International
374 Conference on Machine Learning*, 2017.
- 375 [4] Ricardo Baptista and Matthias Poloczek. Bayesian optimization of combinatorial structures. In
376 *International Conference on Machine Learning*, 2018.
- 377 [5] David Eriksson and Matthias Poloczek. Scalable constrained bayesian optimization. In *International
378 Conference on Artificial Intelligence and Statistics*, 2021.
- 379 [6] Samuel Daulton, Shaun Singh, Vashist Avadhanula, Drew Dimmery, and Eytan Bakshy. Thompson
380 sampling for contextual bandit problems with auxiliary safety constraints. *arXiv preprint
381 arXiv:1911.00638*, 2019.
- 382 [7] Clément Chevalier and David Ginsbourger. Fast computation of the multi-points expected
383 improvement with applications in batch selection. In *International Conference on Learning and
384 Intelligent Optimization*, 2013.
- 385 [8] Javier González, Zhenwen Dai, Philipp Hennig, and Neil Lawrence. Batch Bayesian optimiza-
386 tion via local penalization. In *Artificial intelligence and statistics*, 2016.
- 387 [9] Jian Wu and Peter I Frazier. The parallel knowledge gradient method for batch Bayesian
388 optimization. In *Advances in Neural Information Processing Systems*, 2016.
- 389 [10] Henry B Moss, David S Leslie, Javier Gonzalez, and Paul Rayson. Gibbon: General-purpose
390 information-based bayesian optimisation. *arXiv preprint arXiv:2102.03324*, 2021.
- 391 [11] Hamed Jalali, Inneke Van Nieuwenhuysse, and Victor Picheny. Comparison of kriging-based
392 algorithms for simulation optimization with heterogeneous noise. *European Journal of Opera-
393 tional Research*, 2017.
- 394 [12] Mickaël Binois, Jiangeng Huang, Robert B Gramacy, and Mike Ludkovski. Replication or
395 exploration? sequential design for stochastic simulation experiments. *Technometrics*, 2019.
- 396 [13] José Miguel Hernández-Lobato, James Requeima, Edward O Pyzer-Knapp, and Alán Aspuru-
397 Guzik. Parallel and distributed Thompson sampling for large-scale accelerated exploration of
398 chemical space. In *International Conference on Machine Learning*, 2017.
- 399 [14] Kirthevasan Kandasamy, Akshay Krishnamurthy, Jeff Schneider, and Barnabás Póczos. Paral-
400 lised Bayesian optimisation via Thompson sampling. In *International Conference on Artificial
401 Intelligence and Statistics*, 2018.
- 402 [15] Jasper Snoek, Oren Rippel, Kevin Swersky, Ryan Kiros, Nadathur Satish, Narayanan Sundaram,
403 Mostofa Patwary, Mr Prabhat, and Ryan Adams. Scalable Bayesian optimization using deep
404 neural networks. In *International conference on machine learning*, pages 2171–2180, 2015.
- 405 [16] Zi Wang, Clement Gehring, Pushmeet Kohli, and Stefanie Jegelka. Batched large-scale Bayesian
406 optimization in high-dimensional spaces. In *International Conference on Artificial Intelligence
407 and Statistics*, 2018.
- 408 [17] Carl E Rasmussen and Christopher KI Williams. *Gaussian Processes for Machine Learning*.
409 MIT Press, 2006.
- 410 [18] Peter J Diggle, Jonathan A Tawn, and Rana A Moyeed. Model-based geostatistics. *Journal of
411 the Royal Statistical Society: Series C*, 1998.
- 412 [19] David Eriksson, Michael Pearce, Jacob Gardner, Ryan D Turner, and Matthias Poloczek.
413 Scalable global optimization via local Bayesian optimization. In *Advances in Neural Information
414 Processing Systems*, 2019.

- 415 [20] M. K. Titsias. Variational learning of inducing variables in sparse Gaussian processes. In
416 *International Conference on Artificial Intelligence and Statistics*, 2009.
- 417 [21] Mitchell McIntire, Daniel Ratner, and Stefano Ermon. Sparse gaussian processes for bayesian
418 optimization. In *Association for Uncertainty in Artificial Intelligence*, 2016.
- 419 [22] Ryan-Rhys Griffiths and José Miguel Hernández-Lobato. Constrained Bayesian optimization
420 for automatic chemical design. *arXiv preprint arXiv:1709.05501*, 2017.
- 421 [23] Ang Yang, Cheng Li, Santu Rana, Sunil Gupta, and Svetha Venkatesh. Sparse spectrum
422 Gaussian process for Bayesian optimization. *arXiv preprint arXiv:1906.08898*, 2019.
- 423 [24] James T. Wilson, Viacheslav Borovitskiy, Alexander Terenin, Peter Mostowsky, and Marc Peter
424 Deisenroth. Efficiently sampling functions from Gaussian process posteriors. *International
425 Conference on Machine Learning*, 2020.
- 426 [25] My Phan, Yasin Abbasi Yadkori, and Justin Domke. Thompson sampling and approximate
427 inference. In *Advances in Neural Information Processing Systems*, 2019.
- 428 [26] David Burt, Carl Edward Rasmussen, and Mark Van Der Wilk. Rates of convergence for sparse
429 variational Gaussian process regression. In *International Conference on Machine Learning*,
430 2019.
- 431 [27] Daniel J. Russo, Benjamin Van Roy, Abbas Kazerouni, Ian Osband, and Zheng Wen. A tutorial
432 on thompson sampling. *Foundational Trends in Machine Learning*, 2018.
- 433 [28] D. Russo and B. Van Roy. Learning to optimize via posterior sampling. *Mathematics of
434 Operations Research*, 2014.
- 435 [29] Niranjan Srinivas, Andreas Krause, Sham Kakade, and Matthias Seeger. Gaussian process opti-
436 mization in the bandit setting: no regret and experimental design. In *International Conference
437 on Machine Learning*, 2010.
- 438 [30] Daniele Calandriello, Luigi Carratino, Alessandro Lazaric, Michal Valko, and Lorenzo Rosasco.
439 Gaussian process optimization with adaptive sketching: scalable and no regret. In *Conference
440 on Learning Theory*, 2019.
- 441 [31] Alain Berlinet and Christine Thomas-Agnan. *Reproducing kernel Hilbert spaces in probability
442 and statistics*. Springer Science & Business Media, 2011.
- 443 [32] James Hensman, Nicolás Fusi, and Neil D. Lawrence. Gaussian processes for big data. In
444 *Uncertainty in Artificial Intelligence (UAI 2013)*, 2013.
- 445 [33] James Hensman, Nicolas Durrande, Arno Solin, et al. Variational Fourier features for Gaussian
446 processes. *Journal of Machine Learning Research*, 2017.
- 447 [34] Vincent Dutordoir, Nicolas Durrande, and James Hensman. Sparse Gaussian Processes with
448 Spherical Harmonic Features. In *Proceedings of the 37th International Conference on Machine
449 Learning*, 2020.
- 450 [35] Miguel Lázaro-Gredilla and Anibal Figueiras-Vidal. Inter-domain Gaussian processes for sparse
451 inference using inducing features. In *Advances in Neural Information Processing Systems*,
452 2009.
- 453 [36] Motonobu Kanagawa, Philipp Hennig, Dino Sejdinovic, and Bharath K Sriperumbudur. Gaus-
454 sian processes and kernel methods: A review on connections and equivalences. *arXiv preprint
455 arXiv:1807.02582*, 2018.
- 456 [37] Viacheslav Borovitskiy, Alexander Terenin, Peter Mostowsky, and Marc P. Deisenroth. Matern
457 Gaussian processes on Riemannian manifolds. In *Advances in Neural Information Processing
458 Systems*, 2020.
- 459 [38] Huaiyu Zhu, Christopher KI Williams, Richard Rohwer, and Michal Morciniec. Gaussian
460 regression and optimal finite dimensional linear models, 1997.

- 461 [39] Arno Solin and Simo Särkkä. Hilbert space methods for reduced-rank Gaussian process
462 regression. *Statistics and Computing*, 2020.
- 463 [40] José Miguel Hernández-Lobato, Matthew W Hoffman, and Zoubin Ghahramani. Predictive
464 entropy search for efficient global optimization of black-box functions. *Advances in Neural
465 Information Processing Systems*, 2014.
- 466 [41] José Miguel Hernández-Lobato, Matthew W Hoffman, and Zoubin Ghahramani. Predictive
467 entropy search for efficient global optimization of black-box functions. In *Advances in Neural
468 Information Processing Systems*, 2014.
- 469 [42] Salomon Bochner et al. *Lectures on Fourier integrals*. Princeton University Press, 1959.
- 470 [43] Mojmir Mutny and Andreas Krause. Efficient high dimensional Bayesian optimization with
471 additivity and quadrature fourier features. In *Advances in Neural Information Processing
472 Systems 31*, 2018.
- 473 [44] Thomas M Cover. *Elements of information theory*. John Wiley & Sons, 1999.
- 474 [45] Jasper Snoek, Hugo Larochelle, and Ryan P Adams. Practical Bayesian optimization of machine
475 learning algorithms. In *Advances in Neural Information Processing Systems 25*, 2012.
- 476 [46] Gabriele Santin and Robert Schaback. Approximation of eigenfunctions in kernel-based spaces.
477 *Advances in Computational Mathematics*, 2016.
- 478 [47] Mikhail Belkin. Approximation beats concentration? an approximation view on inference with
479 smooth radial kernels. In *Conference On Learning Theory*, 2018.
- 480 [48] Gabriel Riutort-Mayol, Paul-Christian Burkner, Michael R. Andersen, Arno Solin, and Aki
481 Vehtari. Practical Hilbert space approximate Bayesian Gaussian processes for probabilistic
482 programming. *arXiv preprint arXiv:2004.11408*, 2020.
- 483 [49] Sattar Vakili, Kia Khezeli, and Victor Picheny. On information gain and regret bounds in
484 Gaussian process bandits. In *International Conference on Artificial Intelligence and Statistics*,
485 2021.
- 486 [50] Alexander G de G Matthews, Mark van der Wilk, Tom Nickson, Keisuke Fujii, Alexis Bouk-
487 ouvalas, Pablo León-Villagrà, Zoubin Ghahramani, and James Hensman. Gpflow: A gaussian
488 process library using tensorflow. *Journal of Machine Learning Research*, 2017.
- 489 [51] Vincent Dutordoir, Hugh Salimbeni, Eric Hambro, John McLeod, Felix Leibfried, Artem
490 Artemev, Mark van der Wilk, James Hensman, Marc P Deisenroth, and ST John. Gpflux: A
491 library for deep Gaussian processes. *arXiv preprint arXiv:2104.05674*, 2021.
- 492 [52] trieste. <https://github.com/secondmind-labs/trieste>, 2021.
- 493 [53] Donald R Jones, Matthias Schonlau, and William J Welch. Efficient global optimization of
494 expensive black-box functions. *Journal of Global optimization*, 1998.
- 495 [54] Deng Huang, Theodore T Allen, William I Notz, and Ning Zeng. Global optimization of
496 stochastic black-box systems via sequential kriging meta-models. *Journal of global optimization*,
497 2006.
- 498 [55] Philipp Hennig and Christian J Schuler. Entropy search for information-efficient global opti-
499 mization. *Journal of Machine Learning Research*, 2012.
- 500 [56] Rafael Gómez-Bombarelli, Jorge Aguilera-Iparraguirre, Timothy D Hirzel, David Duvenaud,
501 Dougal Maclaurin, Martin A Blood-Forsythe, Hyun Sik Chae, Markus Einzinger, Dong-Gwang
502 Ha, Tony Wu, et al. Design of efficient molecular organic light-emitting diodes by a high-
503 throughput virtual screening and experimental approach. *Nature Materials*, 2016.
- 504 [57] Henry B Moss, Daniel Beck, Javier González, David S Leslie, and Paul Rayson. Boss: Bayesian
505 optimization over string spaces. *Advances in Neural Information Processing Systems*, 2020.

- 506 [58] Johannes Hachmann, Roberto Olivares-Amaya, Sule Atahan-Evrenk, Carlos Amador-Bedolla,
507 Roel S Sánchez-Carrera, Aryeh Gold-Parker, Leslie Vogt, Anna M Brockway, and Alán Aspuru-
508 Guzik. The Harvard clean energy project: large-scale computational screening and design of
509 organic photovoltaics on the world community grid. *The Journal of Physical Chemistry Letters*,
510 2011.
- 511 [59] David Rogers and Mathew Hahn. Extended-connectivity fingerprints. *Journal of chemical*
512 *information and modeling*, 2010.
- 513 [60] Henry B Moss and Ryan-Rhys Griffiths. Gaussian process molecule property prediction with
514 flowmo. *Advances in Neural Information Processing Systems: Workshop on Machine Learning*
515 *for Molecules.*, 2020.
- 516 [61] Youngmin Cho. *Kernel methods for deep learning*. PhD thesis, UC San Diego, 2012.
- 517 [62] Kurt Cutajar, Edwin V Bonilla, Pietro Michiardi, and Maurizio Filippone. Random feature
518 expansions for deep Gaussian processes. In *International Conference on Machine Learning*,
519 2017.
- 520 [63] Prateek Jain and Purushottam Kar. Non-convex optimization for machine learning. *Foundational*
521 *Trends in Machine Learning*, 2017.
- 522 [64] V. Dani, T. P. Hayes, and S. M. Kakade. Stochastic linear optimization under bandit feedback.
523 In *Conference on Learning Theory*, 2008.

524 Checklist

- 525 1. For all authors...
- 526 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
527 contributions and scope? [Yes] The last paragraph of the introduction states our 3 main
528 contributions.
- 529 (b) Did you describe the limitations of your work? [Yes] See Section 7
- 530 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 531 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
532 them? [Yes]
- 533 2. If you are including theoretical results...
- 534 (a) Did you state the full set of assumptions of all theoretical results? [Yes] Section 2
535 describes our problem formulation and the regularity assumptions on the objective
536 function and noise (Assumptions 1 and 2). Additional assumptions required for Theorem
537 1 are stated in Assumptions 3, 4 and 5. The first paragraph of Section 5.2 specifies
538 the restrictions for which Proposition 1 is valid (Mercer decomposition of the kernel;
539 inducing points chosen by a k -DPP).
- 540 (b) Did you include complete proofs of all theoretical results? [Yes] All proofs are given
541 in the supplementary material.
- 542 3. If you ran experiments...
- 543 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
544 mental results (either in the supplemental material or as a URL)? [No] The code will
545 be made available upon publication to avoid breaking anonymity.
- 546 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
547 were chosen)? [Yes] See section 6
- 548 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
549 ments multiple times)? [Yes] the experiments of Section 6.1 report confidence intervals
550 over 30 replicated experiments. However, due to the high cost of experiments, a single
551 run has been used in the experiment of Section 6.2 and no error bar is reported.
- 552 (d) Did you include the total amount of compute and the type of resources used (e.g., type
553 of GPUs, internal cluster, or cloud provider)? [N/A] We have chosen to report our
554 results in terms of number of calls to the objective function rather than wall-clock time,
555 which is difficult to report fairly for very different routines, particularly when using
556 parallel computing resources.

- 557 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 558 (a) If your work uses existing assets, did you cite the creators? [Yes] the 3 python libraries
- 559 used in the experiments are cited at the beginning of Section 6.
- 560 (b) Did you mention the license of the assets? [Yes] We mentioned that all 3 libraries are
- 561 open-source, see Section 6. The full information is available following the references.
- 562 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 563 See introduction of Section 6
- 564 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 565 using/curating? [N/A]
- 566 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 567 information or offensive content? [No] Does not apply to the data used here.
- 568 5. If you used crowdsourcing or conducted research with human subjects...
- 569 (a) Did you include the full text of instructions given to participants and screenshots, if
- 570 applicable? [N/A]
- 571 (b) Did you describe any potential participant risks, with links to Institutional Review
- 572 Board (IRB) approvals, if applicable? [N/A]
- 573 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 574 spent on participant compensation? [N/A]