
Bayesian Optimization with High-Dimensional Outputs

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Bayesian Optimization is a sample-efficient black-box optimization procedure
2 that is typically applied to problems with a small number of independent objec-
3 tives. However, in practice we often wish to optimize objectives defined over
4 many correlated outcomes (or “tasks”). For example, scientists may want to
5 optimize the coverage of a cell tower network across a dense grid of locations.
6 Similarly, engineers may seek to balance the performance of a robot across dozens
7 of different environments via constrained or robust optimization. However, the
8 Gaussian Process (GP) models typically used as probabilistic surrogates for multi-
9 task Bayesian Optimization scale poorly with the number of outcomes, greatly
10 limiting applicability. We devise an efficient technique for exact multi-task GP
11 sampling that combines exploiting Kronecker structure in the covariance matrices
12 with Matheron’s identity, allowing us to perform Bayesian Optimization using
13 exact multi-task GP models with tens of thousands of correlated outputs. In doing
14 so, we achieve substantial improvements in sample efficiency compared to existing
15 approaches that only model aggregate functions of the outcomes. We demonstrate
16 how this unlocks a new class of applications for Bayesian Optimization across
17 a range of tasks in science and engineering, including optimizing interference
18 patterns of an optical interferometer with more than 65,000 outputs.

19 1 Introduction

20 Many problems in science and engineering involve reasoning about multiple, correlated outputs. For
21 example, cell towers broadcast signal across an area, and thus signal strength is spatially correlated.
22 In randomized experiments, treatment effects on multiple outcomes are naturally correlated due
23 to shared causal mechanisms. Without further knowledge of the internal mechanisms (i.e., in a
24 “black-box” setting), Multi-task Gaussian processes (MTGPs) are a natural model for these types
25 of problems as they model the relationship between each output (or “task”), while maintaining the
26 gold standard predictive capability and uncertainty quantification of Gaussian processes (GPs). Many
27 downstream analyses require more of the model than just prediction; they also involve sampling from
28 the posterior distribution to estimate quantities of interest. For instance, we may be interested in the
29 performance of a complex stock trading strategy that requires modeling different stock prices jointly,
30 and want to characterize its conditional value at risk (CVaR) [43], which generally requires Monte
31 Carlo (MC) estimation strategies [10]. Or, we want to use MTGPs in Bayesian Optimization (BO), a
32 method for sample-efficient optimization of black-box functions. Many state of the art BO approaches
33 use MC acquisition functions [53, 3, 4], which require sampling from the posterior distribution over
34 new candidate data points.

35 Drawing posterior samples from MTGPs means sampling over all tasks and all new data points,
36 which typically scales *multiplicatively* in the number of tasks (t) and test data points (n), e.g. like

37 $\mathcal{O}(n^3t^3)$ [46, 6]. For problems with more than a few tasks, posterior sampling thus quickly becomes
 38 intractable due to the size of the posterior covariance matrix. This is especially problematic in the
 39 case of many real-world problems that can have hundreds or thousands of correlated outputs that
 40 should be modelled jointly in order to achieve the best performance.

41 For instance, the cell tower signal maps in Figure 1 each contain 2,500 outputs (pixels). In
 42 this problem, we aim to jointly tune the down-tilt angle and transmission power of the antennas
 43 on each cell tower (locations shown in red) to optimize a global coverage quality metric,
 44 which is a known function of power and interference at each location [19]. Since simulating the
 45 power and interference maps given a parameterization is computationally costly, traditionally
 46 one might apply BO to optimize the aggregate metric. At its core, this problem is a composite
 47 BO problem [3, 4], so we expect an approach that models the constituent outcomes at each
 48 pixel individually to achieve higher sample efficiency. However, modelling each pixel using
 49 existing approaches used for BO is completely intractable in this setting, as we would have to
 50 train and sample from a MTGP with over 5,000 tasks.

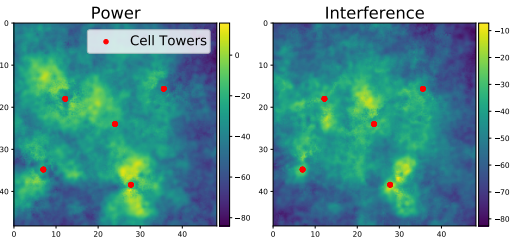


Figure 1: Map of radio signal power and interference for fixed locations of cell towers (red dots). Maps vary smoothly with respect to the towers’ down-tilt angle and transmission power. Our goal is to optimize statistics of these maps as to maximize the overall signal coverage across an area while minimizing interference.

51 To remedy the poor computational scaling with the number of tasks, we exploit Matheron’s rule for
 52 sampling from GP posterior distributions [13, 52]. We derive an efficient method for MTGP sampling
 53 that exploits Kronecker structure inherent to the posterior covariance matrices, thereby reducing the
 54 complexity of sampling from the posterior to become effectively *additive* in the combination of tasks
 55 of data points, i.e. $\mathcal{O}(n^3 + t^3)$, as compared to $\mathcal{O}(n^3t^3)$. Our implementation of Matheron’s rule
 56 draws from the *exact* posterior distribution and does not require random features or inducing points,
 57 unlike decoupled sampling [52]. More specifically, our contributions are as follows:

- 58 • We propose an exact sampling method for multi-task Gaussian processes that has additive time
 59 costs in the combination of tasks and data points, rather than multiplicative (Section 3).
- 60 • We demonstrate empirically how large-scale sampling from MTGPs can aid in challenging
 61 multi-objective, constrained, and contextual Bayesian Optimization problems (Section 4).
- 62 • We introduce a method for efficient posterior sampling for the High-Order Gaussian Process
 63 (HOGP) model [57], allowing it to be used for Bayesian Optimization (Section 3.2). This
 64 advance allows us to more efficiently perform BO on high-dimensional outputs such as images —
 65 including optimizing PDEs, optimizing placements of cell towers for cell coverage, and tuning
 66 the mirrors of an optical interferometer which optimizes over 65,000 tasks jointly (Section 4.3).

67 The rest of the paper is organized as follows: First, in Section 2 we review GPs, MTGPs, and sampling
 68 procedures from the posterior in both GPs and MTGPs. In Section 3, we review Matheron’s rule for
 69 sampling from GP posteriors and explain how to employ it for efficient sampling from MTGP models
 70 including the HOGP model. In Section 4, we illustrate the utility of our method on a wide suite of
 71 problems ranging from constrained BO to the first demonstration of large scale composite BO with
 72 the HOGP. Please see Appendix A for discussion of the limitations and broader impacts of our work.

81 2 Background

82 2.1 Bayesian Optimization

83 In Bayesian Optimization (BO), the goal is to minimize an expensive to evaluate black-box function,
 84 i.e., finding $\min_{x \in \mathcal{X}} f(x)$, by constructing a *surrogate model* to emulate that function. Gaussian
 85 processes (GPs) are often used as surrogates due to their flexibility and well-calibrated uncertainty
 86 estimates. BO optimizes an *acquisition function* defined on the predictive distribution of the surrogate
 87 model to select the next point(s) to evaluate on the true function. These acquisition functions are
 88 often written as intractable integrals and are typically evaluated using Monte Carlo (MC) integration
 89 [53, 3]. MC acquisition functions rely on posterior samples from the surrogate model, which should

90 support fast sampling capabilities for efficient optimization [4]. BO has been applied throughout
 91 machine learning, engineering, and the sciences, and many extensions to the setting described above
 92 exist. We focus on multi-task BO (MTBO), where $f(x)$ is composed of several correlated tasks [49].

93 There are many sub-classes of MTBO problems: constrained BO uses surrogate models to optimize
 94 an objective subject to black-box constraints [25, 21, 26], contextual BO models a single function
 95 that varies across different contexts or environments [33, 11, 22], multi-objective BO learns a
 96 Pareto frontier across several objectives [31, 32, 15], and composite BO considers the setting of a
 97 differentiable objective function defined on the outputs of a vector-valued black-box function [3, 4].
 98 In all of these problems, the setting is similar: several outputs are modelled by the surrogate, whether
 99 the output is a constraint, another objective, or a separate context. As the outputs are often correlated,
 100 multi-task Gaussian processes, which model the relationships between the outputs in a data-efficient
 101 manner, are a natural and common modeling choice.

102 2.2 Gaussian Processes

103 **Single Output Gaussian Processes:** We briefly review single output GPs, see Rasmussen and
 104 Williams [42] for a more detailed introduction. We assume that $y = f(x) + \varepsilon$, $f \sim \mathcal{GP}(0, k_\theta(x, x'))$,
 105 and $\varepsilon \sim \mathcal{N}(0, \sigma^2)$, where f is the noiseless latent function and y are noisy observations of f with
 106 standard deviation σ . $k_\theta(x, x')$ is the kernel with hyperparameters θ (we will drop the dependence
 107 on θ for simplicity); we use $K_{AB} := k_\theta(A, B)$ to refer to the evaluated kernel function on data
 108 points A and B , a matrix which has size $|A| \times |B|$. The predictive distributions over new data points,
 109 x_{test} , is given by the conditional distribution of the Gaussian distribution. That is, $p(f(x_{\text{test}})|\mathcal{D}, \theta) =$
 110 $\mathcal{N}(\mu_{f|\mathcal{D}}^*, \Sigma_{f|\mathcal{D}}^*)$, where $\mathcal{D} := \{X, \mathbf{y}\}$ is the training dataset of size $n = |X|$ and

$$\mu_{f|\mathcal{D}}^* = K_{x_{\text{test}}X}(K_{\text{train}} + \sigma^2 I)^{-1} \mathbf{y}, \quad (1)$$

$$\Sigma_{f|\mathcal{D}}^* = K_{x_{\text{test}}x_{\text{test}}} - K_{x_{\text{test}}X}(K_{\text{train}} + \sigma^2 I)^{-1} K_{Xx_{\text{test}}}, \quad (2)$$

111 with $K_{\text{train}} := K_{XX}$. For simplicity, we will drop the subscripts $f|\mathcal{D}$ in all future statements.
 112 Computing the predictive mean μ^* and variance Σ^* requires $\mathcal{O}(n^3)$ time and $\mathcal{O}(n^2)$ space when
 113 using Cholesky decompositions for the linear solves [42]. Sampling is usually performed by

$$f(\mathbf{x}_{\text{test}})|(Y = \mathbf{y}) = \mu^* + (\Sigma^*)^{1/2} z, \quad (3)$$

114 where $z \sim \mathcal{N}(0, I)$. Computing s samples at n_{test} test points from the predictive distribution costs
 115 $\mathcal{O}(n^3 + sn_{\text{test}}^2 + n^2 n_{\text{test}} + n_{\text{test}}^3)$, computed by adding up all of the matrix vector multiplications (MVMs)
 116 and matrix solves. We can incur the cubic terms only once by re-using Cholesky factorizations of Σ^*
 117 and $K_{XX} + \sigma^2 I$ for each sample.

118 To reduce the time complexity, we can replace all matrix solves with $r < n$ steps of conjugate
 119 gradients (CG) and the Cholesky decompositions with rank $r < n$ Lanczos decompositions (an
 120 approach called LOVE). These change the major scaling from n^3 down to rn^2 and the overall time
 121 complexity to $\mathcal{O}(rn^2 + srn_{\text{test}} + rnn_{\text{test}} + rn_{\text{test}}^2)$ [39, 24]. In general, $r \ll n$ is used and is usually
 122 accurate to nearly numerical precision [24]. We provide additional details in Appendix B.

123 **Multi-Output Gaussian Processes:** One straightforward way of modelling multiple outputs is
 124 to consider each output as an independent GP, modelled in batch together, either with shared or
 125 independent hyperparameters [42, 24, 21]. However, there are two major drawbacks to this approach:
 126 (i) the model is not able to model correlations between the outputs, and (ii) if there are many outputs
 127 then maintaining a separate model for each can result in high memory usage and inference times. To
 128 remedy these issues, Higdon et al. [28] propose the PCA-GP, using principal component analysis
 129 (PCA) to project the outputs to a low-dimensional subspace and then use batch GPs to model the
 130 lower-dimensional outputs.

131 We consider **multi-task Gaussian processes** (MTGPs) with the intrinsic co-regionalization model
 132 (ICM), which considers the relationship between responses as the product of the data and task features
 133 [27, 6]. We focus on this model due to its popularity and simplicity. Given inputs x and x' belonging
 134 to tasks i and j respectively, the covariance under the ICM is $k([x, i], [x', j]) = k_{\mathcal{D}}(x, x')k_t(i, j)$.
 135 Given n data points X with the n associated task indices \mathcal{I} , the covariance is a Hadamard product
 136 of the data kernel and the task kernel, $K_{\text{train}} = K_{XX} \odot K_{\mathcal{I}\mathcal{I}}$, and the response \mathbf{y} is still of size n .
 137 We term this implementation of multi-task GPs ‘‘Hadamard MTGPs’’ in our experiments. In general,
 138 there is no easily exploitable structure in this model.

Table 1: Time complexities for posterior sampling in single-output, multi-task, and high-order Gaussian Process (HOGP) models. **Time complexities shown in blue** are our contributions that have not yet been considered by the literature. Standard sampling from MTGPs scales multiplicatively in the combination of the number of tasks, t , and the number of data points, n , while using Matheron’s rule reduces the combination to effectively become additive in these components.

Model	Distributional (Standard) (Eq. 3)	With Matheron’s rule (Eq. 5)
Single-Output	$\mathcal{O}(n^3 + n_{\text{test}}^3)$	$\mathcal{O}((n^3 + n_{\text{test}}^3))$
Multi-Task	$\mathcal{O}((n^3 + n_{\text{test}}^3)t^3)$	$\mathcal{O}((n^3 + n_{\text{test}}^3) + t^3)$
HOGP	$\mathcal{O}((n^3 + n_{\text{test}}^3) \prod_{i=2}^d d_i^3)$	$\mathcal{O}((n^3 + n_{\text{test}}^3) + \sum_{i=2}^k d_i^3)$

139 If we observe each task at each data point (the so-called *block design* case), the covariance matrix
140 becomes Kronecker structured, e.g. $K_{\text{train}} = K_{XX} \otimes K_T$, where K_{XX} is the data covariance matrix
141 and K_T is the $t \times t$ task covariance matrix between tasks [6], and we now have nt scalar responses.
142 To simplify our exposition, we assume that K_T is full-rank (this is not required as we can use
143 pseudo-inverses in place of inverses). Thus, the GP prior is $\text{vec}(\mathbf{y}) \sim \mathcal{N}(0, K_{XX} \otimes K_T)$, where \mathbf{y} is
144 a matrix of size $n \times t$, and $\text{vec}(\mathbf{y})$ is a vector of shape nt . The GP predictive distribution is given by
145 $p(f^*|x_{\text{test}}, \mathcal{D}) = \mathcal{N}(\mu^*, \Sigma^*)$, where

$$\begin{aligned} \mu^* &= (K_{x_{\text{test}}, X} \otimes K_T)(K_{XX} \otimes K_T + \sigma^2 I_{nT})^{-1} \text{vec}(\mathbf{y}) \\ \Sigma^* &= (K_{x_{\text{test}}, x_{\text{test}}} \otimes K_T) - (K_{x_{\text{test}}, X} \otimes K_T)(K_{XX} \otimes K_T + \sigma^2 I_{nT})^{-1} (K_{x_{\text{test}}, X}^\top \otimes K_T). \end{aligned} \quad (4)$$

146 The kernel matrix on the training data, $K_{XX} \otimes K_T$, is of size $nt \times nt$, which under standard approaches
147 yields inference cubic and multiplicative in n and t , that is $\mathcal{O}(n^3 t^3)$; however, the Kronecker structure
148 can be exploited to compute the posterior mean and variance in $\mathcal{O}(nt(n+t) + n^3 + t^3)$, which is
149 dominated by the individual cubic terms [44, 48].

150 Sampling from the posterior distribution in (4) produces additional computational challenges as we
151 must compute a root (e.g. Cholesky) decomposition of Σ^* , which naively costs $\mathcal{O}((n_{\text{test}}t)^3)$ plus
152 an additional Cholesky decomposition of $(K_{XX} \otimes K_T + \sigma^2 I)^{-1}$, which similarly costs $\mathcal{O}((nt)^3)$
153 time [6]. Thus, the time complexity of drawing s samples is *multiplicative* in n and t , $\mathcal{O}((nt)^3 +$
154 $(n_{\text{test}}t)^3 + s((nt)^2 + (n_{\text{test}}t)^2))$. Using CG and LOVE reduces the complexity; see Appendix B.4.

155 **High-Order Gaussian Processes:** Recently, Zhe et al. [57] proposed the high-order Gaussian process
156 (HOGP) model, a MTGP designed for matrix- and tensor-valued outputs. Given outputs $\mathbf{y} \in$
157 $\mathbb{R}^{d_1 \times \dots \times d_k}$ (e.g. a matrix or tensor), the covariance is the product of the data dimension and each out-
158 put index (i_l and j_l respectively) $k([x, i_1, \dots, i_k], [x', j_1, \dots, j_k]) = k(x, x')k(v_1, v'_1) \dots k(v_k, v'_k)$,
159 where i_1, \dots, i_k are the indices for the output tensor and v_1, \dots, v_k are latent parameters that are
160 optimized along with the kernel hyper-parameters. Thus, K_T in the MTGP framework is replaced
161 by a chain of Kronecker products, so that the GP prior is $\text{vec}(\mathbf{y}) \sim \mathcal{N}(0, K_{XX} \otimes K_2 \otimes \dots \otimes K_k)$.
162 Exploiting the Kronecker structure, computation of posterior mean, posterior variance and hyper-
163 parameter learning takes $\mathcal{O}(n^3 + \sum_{i=2}^k d_i^3 + n \prod_{i=1}^k d_i)$ time as $d_1 = n$ [57]. The experiments
164 of Zhe et al. [57] measure only the error of the predictive mean, rather than quantities that use the
165 predictive variance (such as the negative log likelihood or calibration), and they do not provide a way
166 to sample from the posterior distribution. They demonstrate that the HOGP outperforms PCA-GPs
167 [28], among other high-dimensional output GP methods, with respect to prediction accuracy.

168 2.3 Matheron’s Rule for Single-Task Gaussian Processes

169 Matheron’s rule provides an alternative way of sampling from GP posterior distributions: rather
170 than decomposing the GP predictive covariance for sampling, one can jointly sample from the prior
171 distribution over both train and test points and then update the training samples using the observed
172 data. Matheron’s rule is well known in the geostatistics literature where it is used for “prediction by
173 conditional simulation” [14, 13]. Wilson et al. [52] revitalized interest in Matheron’s rule within the
174 machine learning community by developing a decoupled sampling approach that exploits Matheron’s
175 rule to use both random Fourier features (RFFs) and inducing point approaches in the context of
176 sampling from the posterior in single task GPs. Wilson et al. [54] extended decoupled sampling to
177 use combinations of RFFs, inducing points, and iterative methods. They applied these approaches to
178 approximate Thompson sampling, simulations of dynamical systems, and deep GPs.

179 Matheron’s rule states that if two random variables, X and Y , are jointly Gaussian, then

$$X|Y=y \stackrel{d}{=} X + \text{Cov}(X, Y)\text{Cov}(Y, Y)^{-1}(y - Y),$$

180 where $\text{Cov}(a, b)$ is the covariance of a and b and $\stackrel{d}{=}$ denotes equality in distribution [18, 52]. The
 181 identity can easily be shown by computing the mean and variance of both sides. Following Wilson
 182 et al. [52], we can use this identity to draw posterior samples

$$f^*|Y+\epsilon=y \stackrel{d}{=} f^* + K_{x_{\text{test}}X}(K_{XX} + \sigma^2 I)^{-1}(y - Y - \epsilon), \quad (5)$$

183 where $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$, f^* is the random variable and $f^*|Y+\epsilon=y$ is the conditional random
 184 variable we are interested in drawing samples from. To implement this procedure, we first draw from
 185 the joint prior (of size $n + n_{\text{test}}$):

$$(f, Y) \sim \mathcal{N}\left(0, \begin{pmatrix} K_{XX} & K_{x_{\text{test}}X} \\ K_{Xx_{\text{test}}} & K_{x_{\text{test}}x_{\text{test}}} \end{pmatrix}\right). \quad (6)$$

186 For shorthand, we denote the joint covariance matrix in (6) by $\mathcal{K}_{\text{joint}}$. We next sample $\bar{\epsilon} \sim \mathcal{N}(0, \sigma^2 I)$
 187 and compute: $\bar{f} = f + K_{x_{\text{test}}X}(K_{XX} + \sigma^2 I)^{-1}(y - Y - \bar{\epsilon})$. The solve is against a matrix of size
 188 $n \times n$ so that \bar{f} is our realization of the random variable $f^*|Y+\epsilon=y$. The total time requirements
 189 are then $\mathcal{O}((n + n_{\text{test}})^3 + n^3)$, which is slightly slower than $\mathcal{O}(n_{\text{test}}^3 + n^3)$. Thus, sampling from the
 190 single-task GP posterior using (5) is slower than the standard method based on (1) and (2).

191 3 Matheron’s Rule for Multi-Task Sampling

192 While the time complexity of Matheron-based posterior sampling is inferior for single-task GPs,
 193 the opposite holds true for multi-task GPs, provided that one exploits the special structure in the
 194 covariance matrices. In the following, we demonstrate the advantages in using Matheron-based MTGP
 195 sampling in Section 3.1, extend it to the HOGP [57] in Section 3.2, and identify further advantages of
 196 it for Bayesian Optimization in Section 3.3. This approach allows for further pre-computations than
 197 distributional sampling and that it maintains the same convergence results.

198 3.1 Extending Matheron’s Rule for Multi-task GPs

199 The primary bottleneck that we wish to avoid is the multiplicative scaling in n (or n_{test}) and t . Ideally,
 200 we hope to achieve posterior sampling that is additive in n and t , similar to how Kronecker matrix
 201 vector multiplication is additive in its components. Unlike in the single task case, Matheron’s rule
 202 can substantially reduce the time and memory complexity sampling for MTGPs. For brevity we
 203 limit our presentation here to the core ideas, please see Appendix C for further discussion of the
 204 implementation, as well as Appendix B.1 for a description of the Kronecker identities we exploit.

205 The covariance $\mathcal{K}_{\text{joint}}$ in (6) is structured as the Kronecker product of the joint test-train covariance
 206 matrix and the inter-task covariance; that is, $\mathcal{K}_{\text{joint}} = K_{(X, x_{\text{test}}), (X, x_{\text{test}})} \otimes K_T$, where $K_{(X, x_{\text{test}}), (X, x_{\text{test}})}$
 207 appends n_{test} rows and columns to the joint training data covariance matrix K_{XX} . To sample from
 208 the prior distribution, we need to compute a root decomposition of $\mathcal{K}_{\text{joint}}$.

209 We assume that we have pre-computed $RR^\top = K_{XX}$ with either a Cholesky decomposition
 210 ($\mathcal{O}(n^3)$ time) or a Lanczos decomposition ($\mathcal{O}(rn^2)$ time). Then, we update R to compute
 211 $\tilde{R}\tilde{R}^\top \approx K_{(x_{\text{test}}, X), (x_{\text{test}}, X)}$. Chevalier et al. [12] used a similar strategy to update samples in the
 212 context of single task kriging. Following Jiang et al. [29, Prop. 2], computing \tilde{R} from R costs
 213 $\mathcal{O}(rn_{\text{test}}n + rn_{\text{test}}^2)$ time, dominated by the $rn_{\text{test}}n$ terms for small n_{test} . Using a Cholesky decompo-
 214 sition, this is $\mathcal{O}(n_{\text{test}}n^2)$ time following the same procedure. We then have

$$\mathcal{K}_{\text{joint}} = K_{(x_{\text{test}}X), (x_{\text{test}}X)} \otimes K_T = (\tilde{R} \otimes L)(\tilde{R} \otimes L)^\top, \quad (7)$$

215 where $LL^\top = K_T$. To use the root to sample from the joint random variables, (f, Y) , we need
 216 to compute only $(f, Y) = (\tilde{R}^\top \otimes L^\top)z$, where $z \sim \mathcal{N}(0, I_{nt})$; this computation is a Kronecker
 217 matrix vector multiplication (MVM), which costs just $\mathcal{O}(nt(n+t))$ time. Thus, the overall cost of
 218 sampling to compute the joint random variables is $\mathcal{O}(n^3 + t^3 + nt(n+t) + n_{\text{test}}n^2)$ time, reduced
 219 to $\mathcal{O}(n^2 + rt^2 + rnt + r^2t + rn_{\text{test}}n)$ if using Lanczos decompositions throughout. L only needs
 220 to be computed once, so samples at new test points only require re-computing \tilde{R} and further MVMs.

221 Computing the solve $w = (K_{XX} \otimes K_T + \sigma^2 I_{nT})^{-1}(y - Y - \epsilon)$ takes $\mathcal{O}(n^3 + t^3 + nt(n+t))$ time
 222 using eigen-decompositions and Kronecker MVMs. The cost of the eigen-decomposition dominates
 223 the Kronecker MVM cost so the major cost in this is $\mathcal{O}(n^3 + t^3)$. Finally, there is one more Kronecker
 224 MVM $(K_{x_{\text{test}}, X} \otimes K_T)w$ which takes $\mathcal{O}(nt^2 + n_{\text{test}}nt)$ time. We then only need to reshape \tilde{f} to be
 225 of the correct size.

226 Therefore, the total time complexity of using Matheron’s rule to sample from the MTGP posterior is
 227 $\mathcal{O}(n^3 + t^3)$ for small n_{test} , as the cubic terms dominate due to the Cholesky and eigen-decompositions.
 228 We show the improvements from using Matheron’s rule in Table 2, where the dominating terms are
 229 now *additive in the combination of the tasks and the number of data points, rather than multiplicative*.
 230 Memory complexities, which are also much reduced, are provided in Table 4 in Appendix C.

231 3.2 Extension to HOGP

232 The HOGP model can be seen as a special case of the general procedure for sampling MTGPs. We
 233 replace K_T with kernels across each tensor dimension of the response, so that $K_T = \otimes_{i=2}^k K_i$.
 234 Therefore, the time complexities for sampling go from cubic dependence, $n^3 \prod_{i=2}^k d_i^3$, to a
 235 $(n \prod_{i=1}^k d_i) + (n^3 + \sum_{i=2}^k d_i^3)$ dependence, which again will generally be dominated by the additive
 236 terms for $k \lesssim 5$ as generally n is much larger than the tensor sizes so their product will generally be
 237 less than n^2 . Prior to this work, sampling from the posterior was infeasible for all but the smallest
 238 HOGP models due to the time complexity. By using Matheron’s rule, we can sample from HOGP
 239 posterior distributions over large tensors, unlocking the model for a much wider range of applications
 240 such as BO with composite objectives computed on high-dimensional simulator outputs.

241 3.3 Sample Average Approximation (SAA) Convergence Results

242 The convergence guarantees for Sample Average Approximation (SAA) of MC acquisition functions
 243 from Balandat et al. [4] still apply if posterior samples are drawn via Matheron’s rule. Consider
 244 acquisition functions of the form $\alpha(x; y) = \mathbb{E}[h(f(x)) \mid Y = y]$ with $h : \mathbb{R}^{n_{\text{test}} \times t} \rightarrow \mathbb{R}$, and their
 245 MC approximation $\hat{\alpha}_N(x; y) := \frac{1}{N} \sum_{i=1}^N a(g(f_i^*))$, where f_i^* are i.i.d. samples drawn from the
 246 model posterior at $x \in \mathbb{R}^{n_{\text{test}}}$ using Matheron’s rule. Then, under sufficient regularity conditions,
 247 the optimizer $\arg \max_x \hat{\alpha}_N(x; y)$ converges to the true optimizer $\arg \max_x \alpha(x; y)$ almost surely as
 248 $N \rightarrow \infty$. For a more formal statement and proof of this result see Appendix C.3.

249 4 Experiments

250 We first demonstrate the computational efficiencies gained by posterior sampling using Matheron’s
 251 rule. While this contribution is much more broadly useful and applicable, we focus on the benefits it
 252 brings for BO. Namely, we show that accounting for correlations between outcomes in the model
 253 improves performance in multi-objective and large-scale constrained optimization problems. Finally,
 254 we perform composite BO with tens of thousands of tasks using HOGPs with Matheron-based
 255 sampling. Additional experiments on contextual policy optimization are given in Appendix D.2.

256 4.1 Drawing Samples from Multi-Task Models

257 To demonstrate the performance gains of sampling using Matheron’s rule, we first vary the number of
 258 test points and tasks for a fixed number of samples and training points. Following Feng et al. [22],
 259 we consider a multi-task version of the Hartmann-6 function, generated by splitting the response
 260 surface into tasks using slices of the last dimension. In Figures 2a and 2c, we vary the number of test
 261 points for 5 and 50 tasks, demonstrating that Matheron’s rule sampling is faster and more memory
 262 efficient than distributional sampling with either Kronecker or Hadamard MTGPs. In Figures 2b and
 263 2d, we vary the number of tasks for fixed test points, again finding that distributional sampling is only
 264 competitive for 5 tasks on the GPU. See Appendix D for sampling with LOVE predictive covariances.

265 4.2 Multi-Objective Bayesian Optimization

266 We next consider constrained multi-objective BO, where the goal is to find the *Pareto Frontier*, i.e.,
 267 the set of objective values for which no objective can be improved without deteriorating another

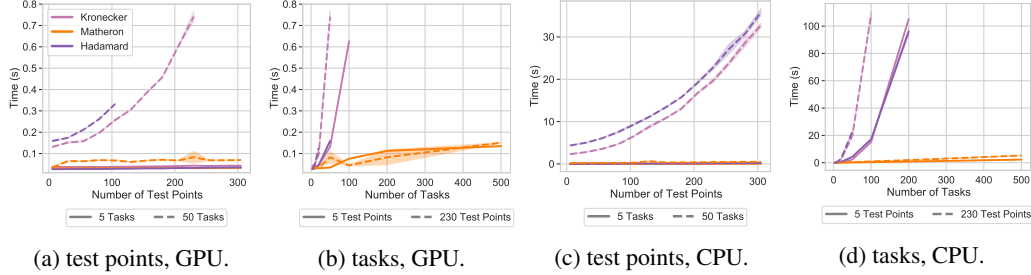


Figure 2: Timings for distributional sampling with Hadamard and Kronecker MTGPs as well as Matheron’s rule sampling for a MTGPs as the number of test points vary for fixed tasks (a,c) and as the number of tasks vary for fixed test points (b,d) on a single Tesla V100 GPU (a,b) and on a single CPU (c,d). The multiplicative scaling of the number of tasks and data points creates significant timing and memory overhead, causing the Kronecker and Hadamard implementations to run out of memory very quickly for all but the smallest numbers of tasks and data points, whereas sampling using Matheron’s rule is efficient even in the many-task large-data regime. Mean and two standard errors are shown, over 10 trials on the GPU, 6 on the CPU.

268 while satisfying all constraints. To measure the quality of the Pareto Frontier, we compute the
 269 hypervolume (HV) of the non-dominated objectives [55]. The optimization problem is made more
 270 difficult by the presence of black-box constraints (and hence additional outcomes) that must be
 271 modelled. We use MC batch versions of the ParEGO and EHVI acquisition functions (qParEGO and
 272 qEHVI) [15]. To generate new candidate points, qParEGO maximizes expected improvement using a
 273 random Chebyshev scalarization of the objectives, while qEHVI maximizes expected hypervolume
 274 improvement. As far as we are aware, Shah and Ghahramani [46] are the only authors to investigate
 275 the use of MTGPs in combination with multi-objective optimization, but only consider 2-4 tasks in
 276 the sequential setting (generating one point at a time). Here, we use full rank inter-task covariances
 277 which work well even in the low data regime. Our Matheron-based sampling scales to large batches
 278 and tasks, and is more sample-efficient on both tasks.

279 We compare Matheron sampled MTGPs to batch
 280 independent MTGPs on the C2DTLZ2 [17] (two
 281 objectives, one constraint for a total of three
 282 modelled tasks) and OSY [37] (two objectives,
 283 six constraints for a total of eight modelled tasks)
 284 test problems. On OSY, we also compare to
 285 PCA GPs [28] due to the larger number of out-
 286 puts. Following Daulton et al. [15] we use both
 287 qParEGO and qEHVI with $q = 2$, for C2DTLZ2
 288 and optimize for 200 iterations. In Figure 3a,
 289 we see that the MTGPs outperform their batch
 290 counterparts by coming closer to the known opti-
 291 mal HV. On OSY, in Figure 3b, we plot the
 292 maximum HV achieved for each method, using
 293 a batch size of $q = 10$, optimizing for 30 it-
 294 erations, where again we see that the MTGPs
 295 significantly outperform their batch counterparts
 296 as well as the PCA GPs, which stagnate quickly.

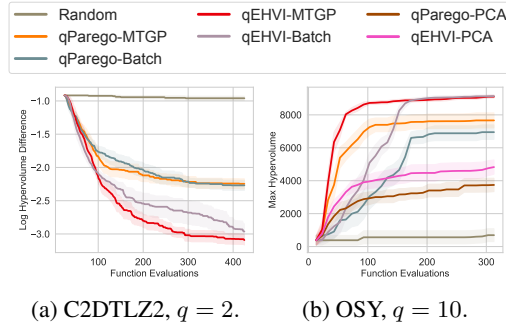


Figure 3: Constrained multi-objective Bayesian Optimization tasks. MTGPs outperform batch models in both the (a) small batch ($q = 2$) and the (b) large batch ($q = 10$) setting. The latter was previously computationally infeasible for MTGPs.

297 4.3 Scalable Constrained Bayesian Optimization

298 We next extend scalable constrained Bayesian Optimization [SCBO, 21], a state of the art algorithm
 299 for constrained BO in high-dimensional problems, to use MTGPs instead of independent GPs. In
 300 constrained BO, the goal is to minimize the objective, f , subject to black box constraints, c_i ; e.g.,

$$\arg \min_x f(x) \text{ s.t. } c_i(x) \leq 0 \quad \forall i \in \{1, \dots, m\}. \quad (8)$$

301 SCBO is a trust region based method and uses batched independent GPs to model the outcome
 302 and transformed constraints. We compare to their results on their two largest problems — the 12-

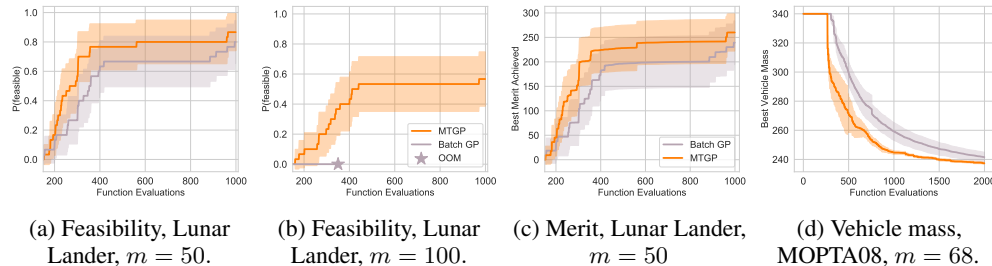


Figure 4: Scalable constrained Bayesian Optimization on Lunar Lander $m = 50, 100$ (a-c) and on MOPTA08 (d). On all three problems, a multi-task GP provides better solutions with better data efficiency. The batch GP reaches feasibility on lunar lander with 50 constraints (a) and a competitive solution (c) but requires more trials, while on 100 constraints, it simply runs out of memory while the MTGP succeeds. On MOPTA08, the MTGP reaches a better solution faster (d).

303 dimensional lunar lander and the 124-dimensional MOPTA08 problem, using the same benchmark
 304 hyper-parameters. To extend their approach, we replace the independent GPs with a single MTGP
 305 model with a full rank ICM kernel over objectives and constraints.

306 **Lunar Lander:** We first consider the lunar lander problem with both 50 and 100 constraints from
 307 the OpenAI Gym [8]. Following Eriksson and Poloczek [21], we initialize with 150 data points
 308 and use TuRBO with Thompson sampling with batches of $q = 20$ for a total of 1000 function
 309 evaluations and repeat over 30 trials. Using a multi-task GP reduces the number of iterations
 310 to achieve at least a 50% chance of feasibility by about 75 steps for the 50 constraint problem (Figure
 311 4a). On the 100 constraint problem, the batch GP runs out of memory after 350 steps on a single
 312 GPU and never achieves feasibility, as indicated in Figure 4b. In Figure 4c, we show the best merit
 313 ($f(x) * \prod_{i=1}^m 1_{c_i(x) \leq 0}$) achieved where the MTGPs are able to achieve feasibility in fewer samples,
 314 but do not reach significantly higher reward. Wall clock times for the $m = 50$ constraint problem, a
 315 table of the steps to achieve feasibility, and a comparison to PCA-GPs [28] are given in Appendix D.

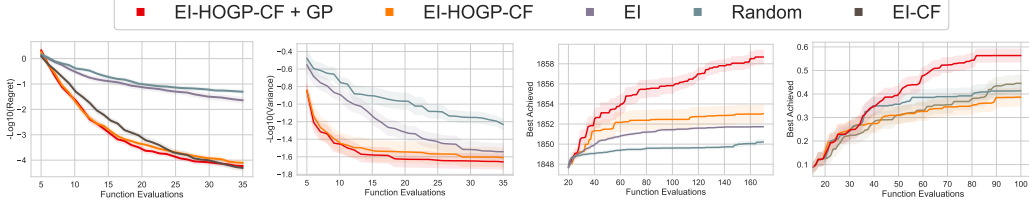
316 **MOPTA08:** We next compare to batch GPs on the MOPTA08 benchmark problem [30] which has 68
 317 constraints that measure the feasibility of a vehicle’s design, while each dimension involves gauges,
 318 materials, and shapes. We use Thompson sampling to acquire points with a batch size of $q = 10$,
 319 130 initial points, and optimize for 2000 iterations repeating over 9 trials. The results are shown in
 320 Figure 4d, where we again observe that SCBO with MTGPs significantly improves both the time
 321 to feasibility and the best overall objective found. Using MTGPs would have been computationally
 322 infeasible in this setting without our efficient posterior sampling approach.

323 4.4 Composite Bayesian Optimization with the HOGP

324 Finally, we push well beyond current scalability limits by extending BO to deal with *many thousands*
 325 of tasks, enabling us to perform sample-efficient optimization in the space of images. *Composite BO*
 326 is a form of BO where the objective is of the form $\max_x g(h(x))$, where h is a multi-output black-box
 327 function modelled by a surrogate and g is cheap to evaluate and differentiable. Decomposing a single
 328 objective into constituent objectives in this way can provide substantial improvements in sample
 329 complexity [3]. Balandat et al. [4] gave convergence guarantees for optimizing general sampled
 330 composite acquisition function objectives that extend to MTGP models under the same regularity
 331 conditions. However, both works only experimentally evaluate batch independent multi-output GPs.

332 We compare to three different baselines: random points (Random), expected improvement on the
 333 metric (EI), and batch GPs optimizing EI in the composite setting (EI-CF). We consider the HOGP
 334 [57] with Matheron’s rule sampling (HOGP-CF) as well as an extension of HOGPs with a prior over
 335 the latent parameters that encourages smoothly varying latent dimensions (HOGP-CF + GP); see
 336 Appendix C.2 for details. More detailed descriptions of the problems are provided in Appendix D.

337 **Chemical Pollutants:** Following Astudillo and Frazier [3], we start with a simple spatial problem
 338 in which environmental pollutant concentrations are observed on a 3×4 grid originally defined in
 339 Bliznyuk et al. [5]. The goal is to optimize a set of four parameters to achieve the true observed value
 340 by minimizing the mean squared error of the output grid to the output grid of the true parameters.



(a) Pollutant concentration $(t = 3 \times 4, d = 4)$. (b) PDE control $(t = 2 \times 64 \times 64, d = 4)$. (c) Cell-tower coverage $(t = 2 \times 50 \times 50, d = 30)$. (d) Optics $(t = 16 \times 64 \times 64, d = 4)$.

Figure 5: Performance of BO with and without HOGP-based composite objectives (EI-HOGP-CF) on four scientific problems. EI-HOGP-CF outperforms a standard BO model directly on the metric itself (EI) and a random baseline (Random). Composite BO with the HOGP also outperforms composite BO with independent GPs (EI-CF). Our smooth latents for the HOGP (EI-HOGP + GP) typically outperform the HOGP itself. EI-CF is only feasible on the smallest problem.

341 The results, over 50 trials, are shown in Figure 5a, where we find that the HOGP models with these
 342 few tasks outperform both independent batch GPs (but slightly) and BO on the metric itself.

343 **Optimizing PDEs:** As a larger experimental problem, we consider optimizing the two diffusivity and
 344 two rate parameters of a spatial Brusselator PDE (solved in `py-pde` [58]) to minimize the weighted
 345 variance of the PDE output as an example of control of a dynamical system. Here, we solve the PDE
 346 on 64×64 grid, producing output solutions of size $2 \times 64 \times 64$. Results over 20 trials are shown in
 347 Figure 5b, where the HOGP models outperform EI fit on the metric and the random baseline.

348 **Cell-Tower Coverage:** Following Dreifuerst et al. [19], we optimize the simulated “coverage map”
 349 resulting from the transmission power and down-tilt settings of 15 cell towers (for a total of 30
 350 parameters) based on a scalarized quality metric combining signal power and inference at each
 351 location so as to maximize total coverage, while minimizing total interference. To reduce model
 352 complexity, we down-sample the simulator output to 50×50 , initializing the optimization with
 353 20 points. Figure 5c presents the results over 20 trials, where the HOGP models with composite
 354 objective EI outperform EI, indicating that modeling the full high-dimensional output is valuable.

355 **Optical Interferometer:** Finally, we consider the tuning of an optical interferometer by the alignment
 356 of two mirrors as in Sorokin et al. [47]. Here, the problem is to optimize the mirror coordinates
 357 to align the interferometer so that it reflects light without interference. There is a sequence of 16
 358 different interference patterns and the simulation outputs are 64×64 images (a tensor of shape
 359 $16 \times 64 \times 64$). Thus, we jointly model 65,536 output dimensions. Scaling composite BO to a problem
 360 of this size would be impossible without the step change in scalability our method provides. Results
 361 are shown in Figure 5d over 20 trials, where we find that the HOGP-CF + GP models outperform EI,
 362 with the HOGP + GP under-performing (perhaps due to high frequency variation in the images).

363 Across all of our experiments, we consistently find that composite BO is considerably more sample
 364 efficient than BO on the metric itself, with significant computational improvements gained from
 365 using the HOGP as compared to batch GPs, which are infeasible much beyond batch sizes of 100.
 366 Furthermore, our structured prior approach for the latent parameters of the HOGP tends to outperform
 367 the random initialization strategy in the original work of Zhe et al. [57].

368 5 Conclusion

369 We demonstrated the utility of Matheron’s rule for sampling the posterior in multi-task Gaussian
 370 processes. Combining Matheron’s rule with scalable Kronecker algebra enables posterior sampling
 371 in $\mathcal{O}(n^3 + t^3)$ rather than the previous $\mathcal{O}(n^3 t^3)$ time. This renders posterior sampling from high-
 372 order Gaussian processes [57] practical, for the first time unlocking Bayesian Optimization with
 373 composite objectives defined on high-dimensional outputs. This increase in computational efficiency
 374 dramatically reduces the time required to do multi-task Bayesian Optimization, and thus enables
 375 practitioners to achieve better automated Bayesian decision making. While we focus on the application
 376 to Bayesian Optimization in this work, our contribution is much broader and provides a step change
 377 in scalability to all methods that involve sampling from MTGP posteriors. We hope in the future to
 378 explore stronger inter-task covariance priors to make MTGP model fits more sample efficient.

References

- 379
- 380 [1] Álvarez, M. A., Rosasco, L., Lawrence, N. D., et al. (2012). Kernels for vector-valued functions:
381 A review. *Foundations and Trends® in Machine Learning*, 4(3):195–266.
- 382 [2] Andrade-Pacheco, R., Rerolle, F., Lemoine, J., Hernandez, L., Meité, A., Juziwelo, L., Bibaut,
383 A. F., van der Laan, M. J., Arnold, B. F., and Sturrock, H. J. (2020). Finding hotspots: development
384 of an adaptive spatial sampling approach. *Scientific reports*, 10(1):1–12.
- 385 [3] Astudillo, R. and Frazier, P. (2019). Bayesian Optimization of Composite Functions. In
386 *International Conference on Machine Learning*, pages 354–363. PMLR. ISSN: 2640-3498.
- 387 [4] Balandat, M., Karrer, B., Jiang, D., Daulton, S., Letham, B., Wilson, A. G., and Bakshy, E.
388 (2020). BoTorch: A Framework for Efficient Monte-Carlo Bayesian Optimization. In *Advances in*
389 *Neural Information Processing Systems*, volume 33.
- 390 [5] Bliznyuk, N., Ruppert, D., Shoemaker, C., Regis, R., Wild, S., and Mugunthan, P. (2008).
391 Bayesian calibration and uncertainty analysis for computationally expensive models using op-
392 timization and radial basis function approximation. *Journal of Computational and Graphical*
393 *Statistics*, 17(2):270–294.
- 394 [6] Bonilla, E. V., Chai, K., and Williams, C. (2007). Multi-task Gaussian Process Prediction. In
395 *Advances in Neural Information Processing Systems*, volume 20, pages 153–160.
- 396 [7] Boustati, A., Damoulas, T., and Savage, R. S. (2019). Non-linear multitask learning with deep
397 gaussian processes. *arXiv preprint arXiv:1905.12407*.
- 398 [8] Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba, W.
399 (2016). Openai gym.
- 400 [9] Bruinsma, W., Perim, E., Tebbutt, W., Hosking, S., Solin, A., and Turner, R. (2020). Scalable
401 exact inference in multi-output gaussian processes. In *International Conference on Machine*
402 *Learning*, pages 1190–1201. PMLR.
- 403 [10] Cakmak, S., Astudillo Marban, R., Frazier, P., and Zhou, E. (2020). Bayesian optimization
404 of risk measures. In *Advances in Neural Information Processing Systems*, volume 33, pages
405 20130–20141. Curran Associates, Inc.
- 406 [11] Char, I., Chung, Y., Neiswanger, W., Kandasamy, K., Nelson, A. O., Boyer, M., Kolemen,
407 E., and Schneider, J. (2019). Offline contextual bayesian optimization. In *Advances in Neural*
408 *Information Processing Systems*, volume 32, pages 4627–4638.
- 409 [12] Chevalier, C., Emery, X., and Ginsbourger, D. (2015). Fast update of conditional simulation
410 ensembles. *Mathematical Geosciences*, 47(7):771–789.
- 411 [13] Chiles, J.-P. and Delfiner, P. (2009). *Geostatistics: modeling spatial uncertainty*, volume 497.
412 John Wiley & Sons.
- 413 [14] Chilès, J.-P. and Lantuéjoul, C. (2005). Prediction by conditional simulation: models and
414 algorithms. In *Space, Structure and Randomness*, pages 39–68. Springer.
- 415 [15] Daulton, S., Balandat, M., and Bakshy, E. (2020). Differentiable Expected Hypervolume
416 Improvement for Parallel Multi-Objective Bayesian Optimization. *Advances in Neural Information*
417 *Processing Systems*, 33.
- 418 [16] de Fouquet, C. (1994). Reminders on the conditioning kriging. In *Geostatistical simulations*,
419 pages 131–145. Springer.
- 420 [17] Deb, K. (2019). Constrained multi-objective evolutionary algorithm. In *Evolutionary and*
421 *Swarm Intelligence Algorithms*, pages 85–118. Springer.
- 422 [18] Doucet, A. (2010). A Note on Efficient Conditional Simulation of Gaussian Distributions.
423 https://www.cs.ubc.ca/~arnaud/doucet_simulationconditionalgaussian.pdf.

- 424 [19] Dreifuerst, R. M., Daulton, S., Qian, Y., Varkey, P., Balandat, M., Kasturia, S., Tomar, A.,
425 Yazdan, A., Ponnampalam, V., and Heath, R. W. (2020). Optimizing coverage and capacity in
426 cellular networks using machine learning. *arXiv preprint arXiv:2010.13710*.
- 427 [20] Emery, X. (2007). Conditioning simulations of gaussian random fields by ordinary kriging.
428 *Mathematical Geology*, 39(6):607–623.
- 429 [21] Eriksson, D. and Poloczek, M. (2021). Scalable constrained bayesian optimization. In *International
430 Conference on Artificial Intelligence and Statistics*. PMLR.
- 431 [22] Feng, Q., Letham, B., Mao, H., and Bakshy, E. (2020). High-Dimensional Contextual Policy
432 Search with Unknown Context Rewards using Bayesian Optimization. *Advances in Neural
433 Information Processing Systems*, 33.
- 434 [23] Feydy, J., Glaunès, J., Charlier, B., and Bronstein, M. (2020). Fast geometric learning with
435 symbolic matrices. *Advances in Neural Information Processing Systems*, 33(4):6.
- 436 [24] Gardner, J., Pleiss, G., Weinberger, K. Q., Bindel, D., and Wilson, A. G. (2018). GPyTorch:
437 Blackbox Matrix-Matrix Gaussian Process Inference with GPU Acceleration. In *Advances in
438 Neural Information Processing Systems*, volume 31, pages 7576–7586.
- 439 [25] Gardner, J. R., Kusner, M. J., Xu, Z. E., Weinberger, K. Q., and Cunningham, J. P. (2014).
440 Bayesian optimization with inequality constraints. In *International Conference on Machine
441 Learning*, volume 32, pages 937–945. PMLR.
- 442 [26] Gelbart, M. A., Snoek, J., and Adams, R. P. (2014). Bayesian optimization with unknown
443 constraints. In *Uncertainty in Artificial Intelligence*, volume 30, pages 250–259. AUAI Press.
- 444 [27] Goovaerts, P. et al. (1997). *Geostatistics for natural resources evaluation*. Oxford University
445 Press on Demand.
- 446 [28] Higdon, D., Gattiker, J., Williams, B., and Rightley, M. (2008). Computer model calibration
447 using high-dimensional output. *Journal of the American Statistical Association*, 103(482):570–
448 583.
- 449 [29] Jiang, S., Jiang, D., Balandat, M., Karrer, B., Gardner, J., and Garnett, R. (2020). Efficient
450 Nonmyopic Bayesian Optimization via One-Shot Multi-Step Trees. In *Advances in Neural
451 Information Processing Systems*, volume 33.
- 452 [30] Jones, D. R. (2008). Large-scale multi-disciplinary mass optimization in the auto industry. In
453 *MOPTA 2008 Conference (20 August 2008)*.
- 454 [31] Khan, N., Goldberg, D. E., and Pelikan, M. (2002). Multi-objective bayesian optimization
455 algorithm. In *Proceedings of the 4th Annual Conference on Genetic and Evolutionary Computation*,
456 pages 684–684.
- 457 [32] Knowles, J. (2006). Parego: A hybrid algorithm with on-line landscape approximation for
458 expensive multiobjective optimization problems. *IEEE Transactions on Evolutionary Computation*,
459 10(1):50–66.
- 460 [33] Krause, A. and Ong, C. S. (2011). Contextual gaussian process bandit optimization. In *Advances
461 in neural information processing systems*, pages 2447–2455.
- 462 [34] Larocque, G., Dutilleul, P., Pelletier, B., and Fyles, J. W. (2006). Conditional gaussian co-
463 simulation of regionalized components of soil variation. *Geoderma*, 134(1-2):1–16.
- 464 [35] Lorch, L., Kremer, H., Trouleau, W., Tsirtsis, S., Szanto, A., Schölkopf, B., and Gomez-
465 Rodriguez, M. (2020). Quantifying the Effects of Contact Tracing, Testing, and Containment
466 Measures in the Presence of Infection Hotspots. *arXiv e-prints*, page arXiv:2004.07641.
- 467 [36] Nguyen, T. V., Bonilla, E. V., et al. (2014). Collaborative multi-output gaussian processes. In
468 *Uncertainty in Artificial Intelligence*, volume 30, pages 643–652. AUAI Press.
- 469 [37] Osyczka, A. and Kundu, S. (1995). A new method to solve generalized multicriteria optimization
470 problems using the simple genetic algorithm. *Structural optimization*, 10(2):94–99.

- 471 [38] Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z.,
472 Gimelshein, N., Antiga, L., et al. (2019). Pytorch: An imperative style, high-performance deep
473 learning library. *Advances in Neural Information Processing Systems*, 32:8026–8037.
- 474 [39] Pleiss, G., Gardner, J. R., Weinberger, K. Q., and Wilson, A. G. (2018). Constant-Time
475 Predictive Distributions for Gaussian Processes. In *Artificial Intelligence and Statistics*, volume
476 arXiv:1803.06058 [cs, stat].
- 477 [40] Pleiss, G., Jankowiak, M., Eriksson, D., Damle, A., and Gardner, J. (2020). Fast Matrix Square
478 Roots with Applications to Gaussian Processes and Bayesian Optimization. In *Advances in Neural
479 Information Processing Systems*, volume 33.
- 480 [41] Rakitsch, B., Lippert, C., Borgwardt, K., and Stegle, O. (2013). It is all in the noise: Efficient
481 multi-task gaussian process inference with structured residuals. In *Advances in Neural Information
482 Processing Systems*, volume 26, pages 1466–1474. Curran Associates, Inc.
- 483 [42] Rasmussen, C. E. and Williams, C. K. I. (2008). *Gaussian processes for machine learning*.
484 Adaptive computation and machine learning. MIT Press, Cambridge, Mass., 3. print edition.
- 485 [43] Rockafellar, R. T., Uryasev, S., et al. (2000). Optimization of conditional value-at-risk. *Journal
486 of risk*, 2:21–42.
- 487 [44] Rougier, J. (2008). Efficient emulators for multivariate deterministic functions. *Journal of
488 Computational and Graphical Statistics*, 17(4):827–843.
- 489 [45] Saatchi, Y. (2011). *Scalable inference for structured Gaussian process models*. PhD thesis,
490 University of Cambridge.
- 491 [46] Shah, A. and Ghahramani, Z. (2016). Pareto frontier learning with expensive correlated
492 objectives. In *Proceedings of The 33rd International Conference on Machine Learning*, volume 48,
493 pages 1919–1927, New York, New York, USA. PMLR.
- 494 [47] Sorokin, D., Ulanov, A., Sazhina, E., and Lvovsky, A. (2020). Interferobot: aligning an optical
495 interferometer by a reinforcement learning agent. *arXiv preprint arXiv:2006.02252*.
- 496 [48] Stegle, O., Lippert, C., Mooij, J., Lawrence, N., and Borgwardt, K. (2011). Efficient inference
497 in matrix-variate gaussian models with iid observation noise. *Advances in Neural Information
498 Processing Systems*, 24:1–9.
- 499 [49] Swersky, K., Snoek, J., and Adams, R. P. (2013). Multi-task Bayesian optimization. In *Advances
500 in Neural Information Processing Systems*, volume 26, pages 2004–2012, Red Hook, NY, USA.
501 Curran Associates Inc.
- 502 [50] Wang, K. A., Pleiss, G., Gardner, J. R., Tyree, S., Weinberger, K. Q., and Wilson, A. G. (2019).
503 Exact Gaussian Processes on a Million Data Points. In *Advances in Neural Information Processing
504 Systems*. arXiv: 1903.08114.
- 505 [51] Wilson, A. G., Gilboa, E., Cunningham, J. P., and Nehorai, A. (2014). Fast kernel learning for
506 multidimensional pattern extrapolation. In *Advances in Neural Information Processing Systems*,
507 volume 27, pages 3626–3634.
- 508 [52] Wilson, J., Borovitskiy, V., Terenin, A., Mostowsky, P., and Deisenroth, M. (2020a). Efficiently
509 sampling functions from Gaussian process posteriors. In *International Conference on Machine
510 Learning*, pages 10292–10302. PMLR. ISSN: 2640-3498.
- 511 [53] Wilson, J., Hutter, F., and Deisenroth, M. P. (2018). Maximizing acquisition functions for
512 Bayesian optimization. In *Advances in Neural Information Processing Systems 31*, pages 9905–
513 9916.
- 514 [54] Wilson, J. T., Borovitskiy, V., Terenin, A., Mostowsky, P., and Deisenroth, M. P. (2020b).
515 Pathwise conditioning of gaussian processes. *arXiv preprint arXiv:2011.04026*.

- 516 [55] Yang, K., Emmerich, M., Deutz, A., and Fonseca, C. M. (2017). Computing 3-d expected
517 hypervolume improvement and related integrals in asymptotically optimal time. In *9th International
518 Conference on Evolutionary Multi-Criterion Optimization - Volume 10173*, EMO 2017,
519 page 685–700, Berlin, Heidelberg. Springer-Verlag.
- 520 [56] Zhang, H. (2007). Maximum-likelihood estimation for multivariate spatial linear coregional-
521 ization models. *Environmetrics: The official journal of the International Environmetrics Society*,
522 18(2):125–139.
- 523 [57] Zhe, S., Xing, W., and Kirby, R. M. (2019). Scalable High-Order Gaussian Process Regression.
524 In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2611–2620.
525 PMLR. ISSN: 2640-3498.
- 526 [58] Zwicker, D. (2020). py-pde: A python package for solving partial differential equations. *Journal
527 of Open Source Software*, 5(48):2158.

528 **Checklist**

- 529 1. For all authors...
- 530 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
531 contributions and scope? [Yes]
- 532 (b) Did you describe the limitations of your work? [Yes] Appendix A.
- 533 (c) Did you discuss any potential negative societal impacts of your work? [Yes] Appendix
534 A.
- 535 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
536 them? [Yes]
- 537 2. If you are including theoretical results...
- 538 (a) Did you state the full set of assumptions of all theoretical results? [Yes] Appendix C.3.
539 (b) Did you include complete proofs of all theoretical results? [Yes] Appendix C.3.
- 540 3. If you ran experiments...
- 541 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
542 mental results (either in the supplemental material or as a URL)? [Yes] See links in
543 supplement. Code attached. We are releasing as much as we can that is non-proprietary.
- 544 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
545 were chosen)? [Yes] Appendix D.
- 546 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
547 ments multiple times)? [Yes] Appendix D.
- 548 (d) Did you include the total amount of compute and the type of resources used (e.g., type
549 of GPUs, internal cluster, or cloud provider)? Appendix D [Yes]
- 550 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 551 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 552 (b) Did you mention the license of the assets? [Yes] Appendix D
- 553 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
554 New codebase attached. MIT License.
- 555 (d) Did you discuss whether and how consent was obtained from people whose data you're
556 using/curating? [N/A]
- 557 (e) Did you discuss whether the data you are using/curating contains personally identifiable
558 information or offensive content? [N/A]
- 559 5. If you used crowdsourcing or conducted research with human subjects...
- 560 (a) Did you include the full text of instructions given to participants and screenshots, if
561 applicable? [N/A]
- 562 (b) Did you describe any potential participant risks, with links to Institutional Review
563 Board (IRB) approvals, if applicable? [N/A]
- 564 (c) Did you include the estimated hourly wage paid to participants and the total amount
565 spent on participant compensation? [N/A]