BOTORCH: A Framework for Efficient Monte-Carlo Bayesian Optimization

Maximilian Balandat

Brian Karrer Facebook Daniel R. Jiang

Samuel Daulton

Facebook balandat@fb.com

briankarrer@fb.com

Facebook drjiang@fb.com

Facebook sdaulton@fb.com

Benjamin Letham Facebook bletham@fb.com Andrew Gordon Wilson New York University andrewgw@cims.nyu.edu Eytan Bakshy Facebook ebakshy@fb.com

Abstract

Bayesian optimization provides sample-efficient global optimization for a broad range of applications, including automatic machine learning, engineering, physics, and experimental design. We introduce BoTorch, a modern programming framework for Bayesian optimization that combines Monte-Carlo (MC) acquisition functions, a novel sample average approximation optimization approach, autodifferentiation, and variance reduction techniques. BoTorch's modular design facilitates flexible specification and optimization of probabilistic models written in PyTorch, simplifying implementation of new acquisition functions. Our approach is backed by novel theoretical convergence results and made practical by a distinctive algorithmic foundation that leverages fast predictive distributions, hardware acceleration, and deterministic optimization. We also propose a novel "one-shot" formulation of the Knowledge Gradient, enabled by a combination of our theoretical and software contributions. In experiments, we demonstrate the improved sample efficiency of BoTorch relative to other popular libraries.

1 Introduction

Computational modeling and machine learning (ML) have led to an acceleration of scientific innovation in diverse areas, ranging from drug design to robotics to material science. These tasks often involve solving time- and resource-intensive global optimization problems to achieve optimal performance. Bayesian optimization (BO) [75, 46, 76], an established methodology for sample-efficient sequential optimization, has been proposed as an effective solution to such problems, and has been applied successfully to tasks ranging from hyperparameter optimization [24, 92, 110], robotic control [15, 5], chemical design [36, 60, 111], and tuning and policy search for internet-scale software systems [4, 58, 57, 23]. Meanwhile, ML research has been undergoing a revolution driven largely by new programming frameworks and hardware that reduce the time from ideation to execution [43, 16, 1, 81]. While BO has become rich with new methodologies, today there is no coherent framework that leverages these computational advances to simplify and accelerate BO research in the same way that modern frameworks have for deep learning. In this paper, we address this gap by introducing BoTorch, a modular and scalable Monte Carlo (MC) framework for BO that is built around modern paradigms of computation, and theoretically grounded in novel convergence results. Our contributions include:

A novel approach to optimizing MC acquisition functions that effectively combines with deterministic higher-order optimization algorithms and variance reduction techniques.

34th Conference on Neural Information Processing Systems (NeurIPS 2020), Vancouver, Canada.

- The first convergence results for sample average approximation (SAA) of MC acquisition functions, including novel general convergence results for SAA via randomized quasi-MC.
- A new, SAA-based "one-shot" formulation of the Knowledge Gradient, a look-ahead acquisition function, with improved performance over the state-of-the-art.
- Composable model-agnostic abstractions for MC BO that leverage modern computational technologies, including auto-differentiation and scalable parallel computation on CPUs and GPUs.

We discuss related work in Section 2 and then present the methodology underlying BOTORCH in Sections 3 and 4. Details of the BOTORCH framework, including its modular abstractions and implementation examples, are given in Section 5. Numerical results are provided in Section 6.

2 Background and Related Work

In BO, we aim to solve $\max_{x \in \mathbb{X}} f_{\text{true}}(x)$, where f_{true} is an expensive-to-evaluate function and $\mathbb{X} \subset \mathbb{R}^d$ is a feasible set. BO consists of two main components: a *probabilistic surrogate model* of the observed function—most commonly, a Gaussian process (GP)—and an *acquisition function* that encodes a strategy for navigating the exploration vs. exploitation trade-off [92]. Taking a model-agnostic view, our focus in this paper is on MC acquisition functions.

Popular libraries for BO include Spearmint [94], GPyOpt [98], Cornell-MOE [106], RoBO [52], Emukit [97], and Dragonfly [49]. We provide further discussion of these packages in Appendix A. Two other libraries, ProBO [72] and GPFlowOpt [55], are of particular relevance. ProBO is a recently suggested framework¹ for using general probabilistic programming in BO. While its model-agnostic approach is similar to ours, ProBO, unlike BOTORCH, does not benefit from gradient-based optimization provided by differentiable programming, or algebraic methods designed to exploit GPU acceleration. GPFlowOpt inherits support for auto-differentiation and hardware acceleration from TensorFlow [via GPFlow, 64], but unlike BOTORCH, it does not use algorithms designed to specifically exploit this potential. Neither ProBO nor GPFlowOpt naturally support MC acquisition functions. In contrast to all existing libraries, BOTORCH is a modular programming framework and employs novel algorithmic approaches that achieve a high degree of flexibility and performance.

The MC approach to optimizing acquisition functions has been considered in the BO literature to an extent, typically using stochastic methods for optimization [100, 106, 109, 104]. Our work takes the distinctive view of sample average approximation (SAA), an approach that combines sampling with deterministic optimization and variance reduction techniques. To our knowledge, we provide the first theoretical analysis and systematic implementation of this approach in the BO setting.

3 Monte-Carlo Acquisition Functions

We begin by describing a general formulation of BO in the context of MC acquisition functions. Suppose we have collected data $\mathcal{D}=\{(x_i,y_i)\}_{i=1}^n$, where $x_i\in\mathbb{X}$ and $y_i=f_{\mathrm{true}}(x_i)+v_i(x_i)$ with v_i some noise corrupting the true function value $f_{\mathrm{true}}(x_i)$. We allow f_{true} to be multi-output, in which case $y_i,v_i\in\mathbb{R}^m$. In some applications we may also have access to distributional information of the noise v_i , such as its (possibly heteroskedastic) variance. Suppose further that we have a probabilistic surrogate model f that for any $\mathbf{x}:=\{x_1,\ldots,x_q\}$ provides a distribution over $f(\mathbf{x}):=(f(x_1),\ldots,f(x_q))$ and $g(\mathbf{x}):=(g(x_1),\ldots,g(x_q))$. We denote by $f_{\mathcal{D}}(\mathbf{x})$ and $g(\mathbf{x})$ the respective posterior distributions conditioned on data $g(\mathbf{x})$. In BO, the model $g(\mathbf{x})$ are multivariate normal. The MC framework we consider here makes no particular assumptions about the form of these posteriors.

The next step in BO is to optimize an acquisition function evaluated on $f_{\mathcal{D}}(\mathbf{x})$ over the *candidate* set \mathbf{x} . Following [105, 7], many acquisition functions can be written as

$$\alpha(\mathbf{x}; \Phi, \mathcal{D}) = \mathbb{E}[a(g(f(\mathbf{x})), \Phi) | \mathcal{D}], \tag{1}$$

where $g: \mathbb{R}^{q \times m} \to \mathbb{R}^q$ is a (composite) objective function, $\Phi \in \Phi$ are parameters independent of \mathbf{x} in some set Φ , and $a: \mathbb{R}^q \times \Phi \to \mathbb{R}$ is a *utility function* that defines the acquisition function.

In some situations, the expectation over $f_{\mathcal{D}}(\mathbf{x})$ in (1) and its gradient $\nabla_{\mathbf{x}}\alpha(\mathbf{x}; \Phi, \mathcal{D})$ can be computed analytically, e.g. if one considers a single-output (m=1) model, a single candidate (q=1) point x, a

¹No implementation of ProBO is available at the time of this writing.

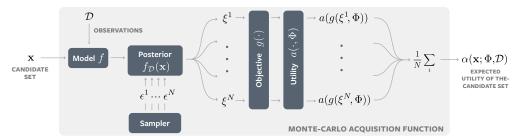


Figure 1: MC acquisition functions. Samples ξ_D^i from the posterior $f_D(\mathbf{x})$ provided by the model f at \mathbf{x} are evaluated in parallel and averaged as in (2). All operations are fully differentiable.

Gaussian posterior $f_{\mathcal{D}}(x) = \mathcal{N}(\mu_x, \sigma_x^2)$, and the identity objective $g(f) \equiv f$. Expected Improvement (EI) is a popular acquisition function that maximizes the expected difference between the currently observed best value f^* (assuming noiseless observations) and f at the next query point, through the utility $a(f, f^*) = \max(f - f^*, 0)$. EI and its gradient have a well-known analytic form [46].

In general, analytic expressions are not available for arbitrary objective functions $g(\cdot)$, utility functions $a(\cdot,\cdot)$, non-Gaussian model posteriors, or collections of points \mathbf{x} which are to be evaluated in a parallel or asynchronous fashion [32, 94, 106, 100, 104]. Instead, MC integration can be used to approximate the expectation (1) using samples from the posterior. An MC approximation $\hat{\alpha}_N(\mathbf{x}; \Phi, \mathcal{D})$ of (1) using N samples $\xi^i_{\mathcal{D}}(\mathbf{x}) \sim f_{\mathcal{D}}(\mathbf{x})$ is straightforward:

$$\hat{\alpha}_N(\mathbf{x}; \Phi, \mathcal{D}) = \frac{1}{N} \sum_{i=1}^N a(g(\xi_{\mathcal{D}}^i(\mathbf{x})), \Phi). \tag{2}$$

The obvious way to evaluate (2) is to draw i.i.d. samples $\xi_{\mathcal{D}}^i(\mathbf{x})$. Alternatively, randomized quasi-Monte Carlo (RQMC) techniques [14] can be used to significantly reduce the variance of the estimate and its gradient (see Appendix E for additional details).

4 MC Bayesian Optimization via Sample Average Approximation

To generate a new candidate set \mathbf{x} , one must optimize the acquisition function α . Doing this effectively, especially in higher dimensions, typically requires using gradient information. For differentiable analytic acquisition functions (e.g. EI, UCB), one can either manually implement gradients, or use auto-differentiation to compute $\nabla_x \alpha(x; \Phi, \mathcal{D})$, provided one can differentiate through the posterior parameters (as is the case for Gaussian posteriors).

4.1 Optimizing General MC Acquisition Functions

An unbiased estimate of the MC acquisition function gradient $\nabla_{\mathbf{x}}\alpha(\mathbf{x};\Phi,\mathcal{D})$ can often be obtained from (2) via the reparameterization trick [50, 85]. The basic idea is that $\xi \sim f_{\mathcal{D}}(\mathbf{x})$ can be expressed as a suitable (differentiable) deterministic transformation $\xi = h_{\mathcal{D}}(\mathbf{x},\epsilon)$ of an auxiliary random variable ϵ independent of \mathbf{x} . For instance, if $f_{\mathcal{D}}(\mathbf{x}) \sim \mathcal{N}(\mu_{\mathbf{x}}, \Sigma_{\mathbf{x}})$, then $h_{\mathcal{D}}(\mathbf{x},\epsilon) = \mu_{\mathbf{x}} + L_{\mathbf{x}}\epsilon$, with $\epsilon \sim \mathcal{N}(0,I)$ and $L_{\mathbf{x}}L_{\mathbf{x}}^T = \Sigma_{\mathbf{x}}$. If $a(\cdot,\Phi)$ and $g(\cdot)$ are differentiable, then $\nabla_{\mathbf{x}}a(g(\xi),\Phi) = \nabla_{g}a\nabla_{\xi}g\nabla_{\mathbf{x}}h_{\mathcal{D}}(\mathbf{x},\epsilon)$.

Our primary methodological contribution is to take a sample average approximation [53] approach to BO. The conventional way of optimizing MC acquisition functions of the form (2) is to re-draw samples from ϵ for each evaluation and apply stochastic first-order methods such as Stochastic Gradient Descent (SGD) [105]. In our SAA approach, rather than re-drawing samples from ϵ for each evaluation of the acquisition function, we draw a set of base samples $E := \{\epsilon^i\}_{i=1}^N$ once, and hold it fixed between evaluations throughout the course of optimization (this can be seen as a specific incarnation of the method of common random numbers). Conditioned on E, the resulting MC estimate $\hat{\alpha}_N(\mathbf{x}; \Phi, \mathcal{D})$ is deterministic. We then obtain the candidate set $\hat{\mathbf{x}}_N^*$ as

$$\hat{\mathbf{x}}_{N}^{*} \in \arg\max_{\mathbf{x} \in \mathbb{X}^{q}} \hat{\alpha}_{N}(\mathbf{x}; \Phi, \mathcal{D}). \tag{3}$$

The gradient $\nabla_{\mathbf{x}}\hat{\alpha}_N(\mathbf{x}; \Phi, \mathcal{D})$ can be computed as the average of the sample-level gradients, exploiting auto-differentiation. We emphasize that whether this average is a "proper" (i.e., unbiased, consistent) estimator of $\nabla_{\mathbf{x}}\alpha(\mathbf{x}; \Phi, \mathcal{D})$ is irrelevant for the convergence results we will derive below.

While the convergence properties of MC integration are well-studied [14], the respective literature on SAA, i.e., convergence of the *optimizer* (3), is far less comprehensive. Here, we derive what, to the best of our knowledge, are the first SAA convergence results for (RQ)MC acquisition functions in the context of BO. To simplify our exposition, we limit ourselves to GP surrogates and i.i.d. base samples; more general results and proofs are presented in Appendix D. For notational simplicity, we will drop the dependence of α and $\hat{\alpha}_N$ on Φ and \mathcal{D} for the remainder of this section. Let $\alpha^* := \max_{\mathbf{x} \in \mathbb{X}^q} \hat{\alpha}_N(\mathbf{x})$, and denote by \mathcal{X}^* the associated set of maximizers. Similarly, let $\hat{\alpha}_N^* := \max_{\mathbf{x} \in \mathbb{X}^q} \hat{\alpha}_N(\mathbf{x})$. With this we have the following key result:

Theorem 1. Suppose (i) $\mathbb X$ is compact, (ii) f has a GP prior with continuously differentiable mean and covariance functions, and (iii) $g(\cdot)$ and $a(\cdot,\Phi)$ are Lipschitz continuous. If the base samples $\{\epsilon^i\}_{i=1}^N$ are i.i.d. $\mathcal N(0,1)$, then (1) $\hat{\alpha}_N^* \to \alpha^*$ a.s., and (2) $\operatorname{dist}(\hat{\mathbf x}_N^*,\mathcal X^*) \to 0$ a.s.. Under additional regularity conditions, (3) $\forall \, \delta > 0, \, \exists \, K < \infty, \, \beta > 0$ s.t. $\mathbb P(\operatorname{dist}(\hat{\mathbf x}_N^*,\mathcal X^*) > \delta) \leq Ke^{-\beta N}, \, \forall \, N \geq 1$.

Under relatively weak conditions,² Theorem 1 ensures not only that the optimizer $\hat{\mathbf{x}}_N^*$ of $\hat{\alpha}_N$ converges to an optimizer of the true α with probability one, but also that the convergence (in probability) happens at an exponential rate. We stated Theorem 1 informally and for i.i.d. base samples for simplicity. In Appendix D.3 we give a formal statement, and extend it to base samples generated by a family of RQMC methods, leveraging recent theoretical advances [79]. While at this point we do not characterize improvements in theoretical convergence rates of RQMC over MC for SAA, we observe empirically that RQMC methods work remarkably well in practice (see Figures 2 and 3).

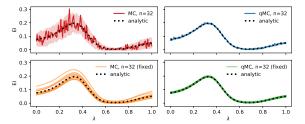


Figure 2: MC and RQMC acquisition functions, with and without ("fixed") re-drawing base samples between evaluations. The model is a GP fit on 15 points randomly sampled from $\mathbb{X} = [0,1]^6$ and evaluated on the Hartmann6 function along the slice $x(\lambda) = \lambda \mathbf{1}$.

Figure 3: Empirical convergence rates of the optimizer for EI using MC / RQMC sampling under SAA / stochastic optimization ("re-sample"). Appendix E provides additional detail and discussion.

The primary benefit from SAA comes from the fact that in order to optimize $\hat{\alpha}_N(\mathbf{x}; \Phi, \mathcal{D})$ for fixed base samples E, one can now employ the full toolbox of deterministic optimization, including quasi-Newton methods that provide faster convergence speeds and are generally less sensitive to optimization hyperparameters than stochastic first-order methods. By default, we use multi-start optimization via L-BFGS-B in conjunction with an initialization heuristic that exploits fast batch evaluation of acquisition functions (see Appendix F.1). We find that in practice the bias from using SAA only has a minor effect on the performance relative to using the analytic ground truth, and often improves performance relative to stochastic approaches (see Appendix E), while avoiding tedious tuning of optimization hyperparameters such as learning rates.

4.2 One-Shot Formulation of the Knowledge Gradient using SAA

The acquisition functions mentioned above, such as EI and UCB, are *myopic*, that is, they do not take into account the effect of an observation on the model in future iterations. In contrast, *look-ahead* methods do. Our SAA approach enables a novel formulation of a class of look-ahead acquisition functions. For the purpose of this paper we focus on the Knowledge Gradient (KG) [27], but our methods extend to other look-ahead acquisition functions such as two-step EI [107].

KG quantifies the expected increase in the maximum of f from obtaining the additional (random) observation data $\{x, y_D(x)\}$. KG often shows improved BO performance relative to simpler, myopic acquisition functions such as EI [90], but in its traditional form it is computationally expensive and

 $^{^{2}}$ Many utility functions a are Lipschitz, including those representing (parallel) EI and UCB [104]. Lipschitzness is a sufficient condition, and convergence can also be shown in less restrictive settings (see Appendix D).

hard to implement, two challenges that we address in this work. Writing $\mathcal{D}_{\mathbf{x}} := \mathcal{D} \cup \{\mathbf{x}, \mathbf{y}_{\mathcal{D}}(\mathbf{x})\}$, we introduce a generalized variant of parallel KG (qKG) [106]:

$$\alpha_{\mathrm{KG}}(\mathbf{x}; \mathcal{D}) = \mathbb{E}\left[\max_{x' \in \mathbb{X}} \mathbb{E}\left[g(f(x')) \mid \mathcal{D}_{\mathbf{x}}\right] \mid \mathcal{D}\right] - \mu_{\mathcal{D}}^{*},\tag{4}$$

with $\mu_{\mathcal{D}}^* := \max_{x \in \mathbb{X}} \mathbb{E}[g(f(x)) | \mathcal{D}]$. Equation (4) quantifies the expected increase in the maximum posterior mean of $g \circ f$ after gathering samples at \mathbf{x} . For simplicity, we only consider standard BO here, but extensions for multi-fidelity optimization [83, 110] are also available in BOTORCH.

Maximizing KG requires solving a nested optimization problem. The standard approach is to optimize the inner and outer problems separately, in an iterative fashion. The outer problem is handled using stochastic gradient ascent, with each gradient observation potentially being an average over multiple samples [106, 109]. For each sample $y_{\mathcal{D}}^i(\mathbf{x}) \sim y_{\mathcal{D}}(\mathbf{x})$, the inner problem $\max_{x_i \in \mathbb{X}} \mathbb{E}\left[f(x_i) \mid \mathcal{D}_{\mathbf{x}}^i\right]$ is solved numerically, either via another stochastic gradient ascent [109] or multi-start L-BFGS-B [26]. An unbiased stochastic gradient can then be computed by leveraging the envelope theorem. Alternatively, the inner problem can be discretized [106]. The computational expense of this nested optimization can be quite large; our main insight is that it may also be unnecessary.

We treat optimizing $\alpha_{\rm KG}(\mathbf{x},\mathcal{D})$ in (4) as an entirely deterministic problem using SAA. Using the reparameterization trick, we express $y_{\mathcal{D}}(\mathbf{x}) = h_{\mathcal{D}}^y(\mathbf{x},\epsilon)$ for some deterministic $h_{\mathcal{D}}$, and draw N fixed base samples $\{\epsilon^i\}_{i=1}^N$ for the outer expectation. The resulting MC approximation of KG is:

$$\hat{\alpha}_{\mathrm{KG},N}(\mathbf{x};\mathcal{D}) = \frac{1}{N} \sum_{i=1}^{N} \max_{x_i \in \mathbb{X}} \mathbb{E}\left[g(f(x_i)) \mid \mathcal{D}_{\mathbf{x}}^i\right] - \mu^*. \tag{5}$$

Theorem 2. Suppose conditions (i) and (ii) of Theorem 1 hold, and that (iii) $g(\cdot)$ is affine. If the base samples $\{\epsilon^i\}_{i\geq 1}$ are drawn i.i.d from $\mathcal{N}(0,1)$, then (1) $\hat{\alpha}^*_{\mathrm{KG},N} \to \alpha^*_{\mathrm{KG}}$ a.s., (2) $\mathrm{dist}(\hat{\mathbf{x}}^*_{\mathrm{KG},N},\mathcal{X}^*_{\mathrm{KG}}) \to 0$ a.s., and (3) \forall δ > 0, \exists K < ∞ , β > 0 s.t. $\mathbb{P}(\mathrm{dist}(\hat{\mathbf{x}}^*_{\mathrm{KG},N},\mathcal{X}^*_{\mathrm{KG}}) > \delta) \leq Ke^{-\beta N}$ for all $N \geq 1$.

Theorem 2 also applies when using RQMC (Appendix D.3), in which case we again observe improved empirical convergence rates. In Appendix D.4, we prove that if f_{true} is drawn from the same GP prior as f and $g(f) \equiv f$, then the MC-approximated KG policy (i.e., when (5) is maximized in each period to select measurements) is $asymptotically\ optimal\ [27, 25, 83, 7]$, meaning that as the number of measurements tends to infinity, an optimal point $x^* \in \mathcal{X}_f^* := \arg\max_{x \in \mathbb{X}} f(x)$ is identified.

Conditional on the fixed base samples, (5) does not exhibit the nested structure used in the conventional formulation (which requires solving an optimization problem to get a noisy gradient estimate). Moving the maximization outside of the sample average yields the *equivalent* problem

$$\max_{\mathbf{x} \in \mathbb{X}} \hat{\alpha}_{\mathrm{KG},N}(\mathbf{x}, \mathcal{D}) \equiv \max_{\mathbf{x}, \mathbf{x}'} \frac{1}{N} \sum_{i=1}^{N} \mathbb{E} \left[g(f(x_i)) \, | \, \mathcal{D}_{\mathbf{x}}^i \right], \tag{6}$$

where $\mathbf{x}' := \{x^i\}_{i=1}^N \in \mathbb{X}^N$ represent "next stage" solutions, or "fantasy points." If g is affine, the expectation in (6) admits an analytical expression. If not, we use another MC approximation of the form (2) with N_I fixed inner based samples E_I . The key difference from the envelope theorem approach [109] is that we do not solve the inner problem to completion for every fantasy point for every gradient step w.r.t. \mathbf{x} . Instead, we solve (6) jointly over \mathbf{x} and the fantasy points \mathbf{x}' . The resulting optimization problem is of higher dimension, namely (q+N)d instead of qd, but unlike the envelope theorem formulation it can be solved as a single problem, using methods for deterministic optimization. Consequently, we dub this KG variant the "One-Shot Knowledge Gradient" (OKG). The ability to auto-differentiate the involved quantities (including the samples $y_{\mathcal{D}}^i(\mathbf{x})$ and $\xi_{\mathcal{D}_{\mathbf{x}}}^i(\mathbf{x})$ through the posterior updates) w.r.t. \mathbf{x} and \mathbf{x}' allows BOTORCH to solve this problem effectively. The main limitation of OKG is the linear growth of the dimension of the optimization problem in N, which can be challenging to solve - however, in practical settings, we observe good performance for moderate N. We provide a simplified implementation of OKG in the following section.

5 Programmable Bayesian Optimization with BoTorch

SAA provides an efficient and robust approach to optimizing MC acquisition functions through the use of deterministic gradient-based optimization. In this section, we introduce BOTORCH, a

³For a GP, $h_{\mathcal{D}}^y(\mathbf{x}, \epsilon) = \mu_{\mathcal{D}}(\mathbf{x}) + L_{\mathcal{D}}^{\sigma}(\mathbf{x})\epsilon$, with $L_{\mathcal{D}}^{\sigma}(\mathbf{x})$ a root decomposition of $\Sigma_{\mathcal{D}}^{\sigma}(\mathbf{x}) := \Sigma_{\mathcal{D}}(\mathbf{x}, \mathbf{x}) + \Sigma^{v}(\mathbf{x})$.

⁴Convergence results can be established in the same way, and will require that min $\{N, N_I\} \to \infty$.

complementary differentiable programming framework for Bayesian optimization research. Following the conceptual framework outlined in Figure 1, BOTORCH provides modular abstractions for representing and implementing sophisticated BO procedures. Operations are implemented as PyTorch modules that are highly parallelizable on modern hardware and end-to-end differentiable, which allows for efficient optimization of acquisition functions. Since the chain of evaluations on the sample level does not make any assumptions about the form of the posterior, BOTORCH's primitives can be directly used with any model from which re-parameterized posterior samples can be drawn, including probabilistic programs [99, 8], Bayesian neural networks [71, 87, 61, 41], and more general types of GPs [19]. In this paper, we focus on an efficient and scalable implementation of GPs, GPyTorch [29].

To illustrate the core components of BOTORCH, we demonstrate how both known and novel acquisition functions can readily be implemented. For the purposes of exposition, we show a set of simplified implementations here; details and additional examples are given in Appendices G and H.

5.1 Composing BOTORCH Modules for Multi-Objective Optimization

In our first example, we consider qParEGO [20], a variant of ParEGO [54], a method for multi-objective optimization.

```
1 weights = torch.distributions.Dirichlet(torch.ones(num_objectives)).sample()
2 scalarized_objective = GenericMCObjective(
3 lambda Y: 0.05 * (weights * Y).sum(dim=-1) + (weights * Y).min(dim=-1).values
4 )
5 qParEGO = qExpectedImprovement(model=model, objective=scalarized_objective)
6 candidates, values = optimize_acqf(qParEGO, bounds=box_bounds, q=1)
```

Code Example 1: Multi-objective optimization via augmented Chebyshev scalarizations.

Code Example 1 implements the inner loop of qParEGO. We begin by instantiating a GenericMCObjective module that defines an augmented Chebyshev scalarization. This is an instance of BoTorch's abstract MCObjective, which applies a transformation $g(\cdot)$ to samples ξ from a posterior in its forward(ξ) pass. In line 5, we instantiate an MCAcquisitionFunction module, in this case, qExpectedImprovement, parallel EI. Acquisition functions combine a model and the objective into a single module that assigns a utility $\alpha(\mathbf{x})$ to a candidate set \mathbf{x} in its forward pass. Models can be any PyTorch module implementing a probabilistic model conforming to BoTorch's basic Model API. Finally, candidate points are selected by optimizing the acquisition function, through the use of the optimize_acqf() utility function, which finds the candidates $\mathbf{x}^* \in \arg\max_{\mathbf{x}} \alpha(\mathbf{x})$. Auto-differentiation makes it straightforward to use gradient-based optimization even for complex acquisition functions and objectives. Our SAA approach permits the use of deterministic higher-order optimization to efficiently and reliably find \mathbf{x}^* .

In [6] it is shown how performing operations on independently modeled objectives yields better optimization performance when compared to modeling combined outcomes directly (e.g., for the case of calibrating the outputs of a simulator). MCObjective is a powerful abstraction that makes this straightforward. It can also be used to implement unknown (i.e. modeled) outcome constraints: BoTorch implements a ConstrainedMCObjective to compute a feasibility-weighted objective using a sample-level differentiable relaxation of the feasibility [89, 28, 31, 58].

5.2 Implementing Parallel, Asynchronous Noisy Expected Improvement

Noisy EI (NEI) [58] is an extension of EI that is well-suited to highly noisy settings, such as A/B tests. Here, we describe a novel *full MC* formulation of NEI that extends the original one from [58] to joint parallel optimization and generic objectives. Letting $(\xi, \xi_{\text{obs}}) \sim f_{\mathcal{D}}((\mathbf{x}, \mathbf{x}_{\text{obs}}))$, our implementation avoids the need to characterize the (uncertain) best observed function value explicitly by averaging improvements on samples from the joint posterior over new and previously evaluated points:

$$qNEI(\mathbf{x}; \mathcal{D}) = \mathbb{E}\left[\left(\max g(\xi) - \max g(\xi_{obs})\right)_{+} | \mathcal{D}\right]. \tag{7}$$

Code Example 2 provides an implementation of qNEI as formulated in (7). New MC acquisition functions are defined by extending an MCAcquisitionFunction base class and defining a forward

Code Example 2: Parallel Noisy EI

pass that compute the utility of a candidate x. In the constructor (not shown), the programmer sets X_baseline to an appropriate subset of the points at which the function was observed.

Like all MC acquisition functions, qNEI can be extended to support asynchronous candidate generation, in which a set $\tilde{\mathbf{x}}$ of pending points have been submitted for evaluation, but have not yet completed. This is done by concatenating pending points into \mathbf{x} with the @concatenate_pending_points decorator. This allows us to compute the joint utility $\alpha(\mathbf{x} \cup \tilde{\mathbf{x}}; \Phi, \mathcal{D})$ of all points, pending and new, but optimize only with respect to the new \mathbf{x} . This strategy also provides a natural way of generating parallel BO candidates using sequential greedy optimization [94]: We generate a single candidate, add it to the set of pending points, and proceed to the next. Due to submodularity of many common classes of acquisition functions (e.g., EI, UCB) [105], this approach can often yield better optimization performance compared to optimizing all candidate locations simultaneously (see Appendix F.2).

With the observed, pending, and candidate points (X_full) in hand, we use the Model's posterior() method to generate an object that represents the joint posterior across all points. The Posterior returned by posterior(x) represents $f_{\mathcal{D}}(\mathbf{x})$ (or $y_{\mathcal{D}}(\mathbf{x})$, if the observation_noise keyword argument is set to True), and may be be explicit (e.g. a multivariate normal in the case of GPs), or implicit (e.g. a container for a warmed-up MCMC chain). Next, samples are drawn from the posterior distribution p via a MCSampler, which employs the reparameterization trick [50, 85]. Given base samples $E \in \mathbb{R}^{N_s \times qm}$, a Posterior object produces N_s samples $\xi_{\mathcal{D}} \in \mathbb{R}^{N_s \times q \times m}$ from the joint posterior. Its forward(p) pass draws samples $\xi_{\mathcal{D}}^i$ from p by automatically constructing base samples E. By default, BoTorch uses RQMC via scrambled Sobol sequences [78]. Finally, these samples are mapped through the objective, and the expected improvement between the candidate point x and observed/pending points is computed by marginalizing the improvements on the sample level.

5.3 Look-ahead Bayesian Optimization with One-Shot KG

Code Example 3 shows a simplified OKG implementation, as discussed in Section 4.2.

Code Example 3: Implementation of One-Shot KG

Here, the input X to forward is a concatenation of x and N fantasy points \mathbf{x}' (this setup ensures that OKG can be optimized using the same APIs as all other acquisition functions). After X is split into its components, we utilize the Model's fantasize(x, sampler) method that, given x and a MCSampler, constructs a batched set of N fantasy models $\{f^i\}_{i=1}^N$ such that $f_{\mathcal{D}}^i(\mathbf{x}) \stackrel{d}{=} f_{\mathcal{D}_{\mathbf{x}}^i}(\mathbf{x}), \forall \mathbf{x} \in \mathbb{X}^q$, where $\mathcal{D}_{\mathbf{x}}^i := \mathcal{D} \cup \{\mathbf{x}, y_{\mathcal{D}}^i(\mathbf{x})\}$ is the original dataset augmented by a fantasy observation at x. The fantasy models provide a distribution over functions conditioned on future observations at x, which is used here to implement one-step look-ahead. SimpleRegret computes $\mathbb{E}\left[g(f(x_i)) \mid \mathcal{D}_{\mathbf{x}}^i\right]$ from (6) for each i in batch mode. The propagate_grads context enables auto-differentiation through both the generation of the fantasy models and the evaluation of their respective posteriors at the points \mathbf{x}' .

6 Experiments

6.1 Exploiting Parallelism and Hardware Acceleration

Botorch utilizes inference and optimization methods designed to exploit parallelization via batched computation, and integrates closely with GPyTorch [29]. These model have fast test-time (predictive) distributions and sampling. This is crucial for BO, where the same models are evaluated many times in order to optimize the acquisition function. GPyTorch makes use of structure-exploiting algebra and local interpolation for $\mathcal{O}(1)$ computations in querying the predictive distribution, and $\mathcal{O}(T)$ for drawing a posterior sample at T points, compared to the standard $\mathcal{O}(n^2)$ and $\mathcal{O}(T^3n^3)$ computations [82].

Figure 4 reports wall times for batch evaluation of qExpectedImprovement at multiple candidate sets $\{\mathbf{x}^i\}_{i=1}^b$ for different MC samples sizes N, on both CPU and GPU for a GPyTorch GP. We observe significant speedups from running on the GPU, with scaling essentially linear in the batch size b, except for very large b and N. Figure 5 shows between 10–40X speedups when using fast predictive covariance estimates over standard posterior inference in the same setting. The speedups grow slower on the GPU, whose cores do not saturate as quickly as on the CPU when doing standard posterior inference (for additional details see Appendix B). Together, batch evaluation and fast predictive distributions enable efficient, parallelized acquisition function evaluation for a very large number (tens of thousands) of points. This scalability allows us to implement and exploit novel highly parallelized initialization and optimization techniques.

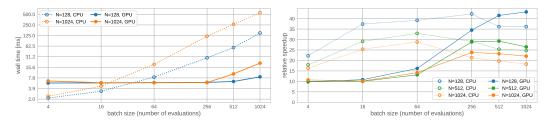


Figure 4: Wall times for batched evaluation of qEI

Figure 5: Fast predictive distributions speedups

6.2 Bayesian Optimization Performance Comparisons

We compare (i) the empirical performance of standard algorithms implemented in BoTorch with those from other popular BO libraries, and (ii) our novel acquisition function, OKG, against other acquisition functions, both within BoTorch and in other packages. We isolate three key frameworks—GPyOpt, Cornell MOE (*MOE EI, MOE KG*), and Dragonfly—because they are the most popular libraries with ongoing support⁵ and are most closely related to BoTorch in terms of state-of-the-art acquisition functions. GPyOpt uses an extension of EI with a local penalization heuristic (henceforth *GPyOpt LP-EI*) for parallel optimization [34]. For Dragonfly, we consider its default ensemble heuristic (henceforth *Dragonfly GP Bandit*) [49].

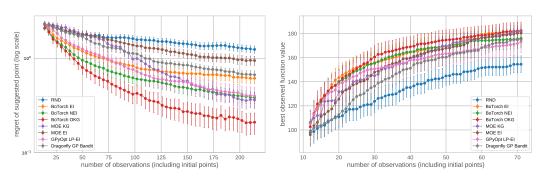


Figure 6: Hartmann (d = 6), noisy, best suggested

Figure 7: DQN tuning benchmark (Cartpole)

⁵We were unable to install GPFlowOpt due to its incompatibility with current versions of GPFlow/TensorFlow.

Our results provide three main takeaways. First, we find that BoTorch's algorithms tend to achieve greater sample efficiency compared to those of other packages (all packages use their default models and settings). Second, we find that OKG often outperforms all other acquisition functions. Finally, OKG is more computationally scalable than MOE KG (the gold-standard implementation of KG), showing significant reductions in wall time (up to 6X, see Appendix C.2) while simultaneously achieving improved optimization performance (Figure 6).

Synthetic Test Functions: We consider BO for parallel optimization of q=4 design points, on four noisy synthetic functions used in Wang et al. [100]: Branin, Rosenbrock, Ackley, and Hartmann. Figure 6 reports means and 95% confidence intervals over 100 trials for Hartmann; results for the other functions are qualitatively similar and are provided in Appendix C.1, together with details on the evaluation. Results for constrained BO using a differentiable relaxation of the feasibility indicator on the sample level are provided in Appendix C.3.

Hyperparameter Optimization: We illustrate the performance of BoTorch on real-world applications, represented by three hyperparameter optimization (HPO) experiments: (1) Tuning 5 parameters of a deep Q-network (DQN) learning algorithm [66, 67] on the *Cartpole* task from OpenAI gym [12] and the default DQN agent implemented in Horizon [30], Figure 7; (2) Tuning 6 parameters of a neural network surrogate model for the UCI Adult data set [56] introduced by Falkner et al. [22], available as part of HPOlib2 [21], Figure 17 in Appendix C.4; (3) Tuning 3 parameters of the recently proposed *Stochastic Weight Averaging* (SWA) procedure of Izmailov et al. [40] on the VGG-16 [93] architecture for CIFAR-10, which achieves superior accuracy compared to previously reported results. A more detailed description of these experiments is given in Appendix C.4.

7 Discussion and Outlook

We presented a novel strategy for effectively optimizing MC acquisition functions using SAA, and established strong theoretical convergence guarantees (in fact, our RQMC convergence results are novel more generally, and of independent interest). Our proposed OKG method, an extension of this approach to "one-shot" optimization of look-ahead acquisition functions, constitutes a significant development of KG, improving scalability and allowing for generic composite objectives and outcome constraints. This approach can naturally be extended to multi-step and other look-ahead approaches [44].

We make these methodological and theoretical contributions available in our open-source library BoTorch (https://botorch.org), a modern programming framework for BO that features a modular design and flexible API, our distinct SAA approach, and algorithms specifically designed to exploit modern computing paradigms such as parallelization and auto-differentiation. BoTorch is particularly valuable in helping researchers to rapidly assemble novel BO techniques. Specifically, the basic MC acquisition function abstraction provides generic support for batch optimization, asynchronous evaluation, RQMC integration, and composite objectives (including outcome constraints).

Our empirical results show that besides increased flexibility, our advancements in both methodology and computational efficiency translate into significantly faster and more accurate closed-loop optimization performance on a range of standard problems. While other settings such as high-dimensional [47, 102, 59], multi-fidelity [83, 110], or multi-objective [54, 80, 20] BO, and non-MC acquisition functions such as Max-Value Entropy Search [101], are outside the scope of this paper, these approaches can readily be realized in BOTORCH and are included in the open-source software package. One can also naturally generalize BO procedures to incorporate neural architectures in BOTORCH using standard PyTorch models. In particular, deep kernel architectures [103], deep Gaussian processes [19, 88], and variational auto-encoders [33, 68] can easily be incorporated into BoTorch's primitives, and can be used for more expressive kernels in high-dimensions.

In summary, BOTORCH provides the research community with a robust and extensible basis for implementing new ideas and algorithms in a modern computational paradigm, theoretically backed by our novel SAA convergence results.

Broader Impact

Bayesian optimization is a generic methodology for optimizing black-box functions, and therefore, by its very nature, not tied to any particular application domain. As mentioned earlier in the paper, Bayesian optimization has been used for various arguably good causes, including drug discovery or reducing the energy footprint of ML applications by reducing the computational cost of tuning hyperparameters. In the Appendix, we give an specific example for how our work can be applied in a public health context, namely to efficiently distribute survey locations for estimating malaria prevalence. BoTorch as a tool specifically has been used in various applications, including transfer learning for neural networks [62], high-dimensional Bayesian optimization [59], drug discovery [10], sim-to-real transfer [69], trajectory optimization [42], and nano-material design [73]. However, there is nothing inherent to this work and Bayesian optimization as a field more broadly that would preclude it from being abused in some way, as is the case with any general methodology.

Acknowledgments and Disclosure of Funding

We wish to thank Art Owen for insightful conversations on quasi-Monte-Carlo methods. We also express our appreciation to Peter Frazier and Javier Gonzalez for their helpful feedback on earlier versions of this paper.

Andrew Gordon Wilson is supported by NSF I-DISRE 193471, NIH R01 DA048764-01A1, NSF IIS-1910266, and NSF 1922658 NRT-HDR: FUTURE Foundations, Translation, and Responsibility for Data Science.

References

- [1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: A system for large-scale machine learning. In *OSDI*, volume 16, pages 265–283, 2016.
- [2] Robert J Adler. An introduction to continuity, extrema, and related topics for general Gaussian processes. IMS, 1990.
- [3] Robert J. Adler. The Geometry of Random Fields. Society for Industrial and Applied Mathematics, 2010.
- [4] Deepak Agarwal, Kinjal Basu, Souvik Ghosh, Ying Xuan, Yang Yang, and Liang Zhang. Online parameter selection for web-based ranking problems. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 23–32, 2018.
- [5] Rika Antonova, Akshara Rai, and Christopher G. Atkeson. Deep kernels for optimizing locomotion controllers. In *Proceedings of the 1st Conference on Robot Learning*, CoRL, 2017.
- [6] Raul Astudillo and Peter Frazier. Bayesian optimization of composite functions. In *Proceedings of the* 36th International Conference on Machine Learning, volume 97, pages 354–363. PMLR, 2019.
- [7] Julien Bect, François Bachoc, and David Ginsbourger. A supermartingale approach to gaussian process based sequential design of experiments. *Bernoulli*, 25(4A):2883–2919, 11 2019.
- [8] Eli Bingham, Jonathan P. Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofanis Karaletsos, Rohit Singh, Paul Szerlip, Paul Horsfall, and Noah D. Goodman. Pyro: Deep Universal Probabilistic Programming. *Journal of Machine Learning Research*, 2018.
- [9] Nikolay Bliznyuk, David Ruppert, Christine Shoemaker, Rommel Regis, Stefan Wild, and Pradeep Mugunthan. Bayesian calibration and uncertainty analysis for computationally expensive models using optimization and radial basis function approximation. *Journal of Computational and Graphical Statistics*, 17(2):270–294, 2008.
- [10] Jacques Boitreaud, Vincent Mallet, Carlos Oliver, and Jerome Waldispühl. Optimol: Optimization of binding affinities in chemical space for drug discovery. bioRxiv, 2020.
- [11] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013.
- [12] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.

- [13] Alexander Buchholz, Florian Wenzel, and Stephan Mandt. Quasi-Monte Carlo variational inference. In *Proceedings of the 35th International Conference on Machine Learning*. PMLR, 2018.
- [14] Russel E Caffisch. Monte carlo and quasi-monte carlo methods. Acta numerica, 7:1–49, 1998.
- [15] Roberto Calandra, André Seyfarth, Jan Peters, and Marc Peter Deisenroth. Bayesian optimization for learning gaits under uncertainty. *Annals of Mathematics and Artificial Intelligence*, 2016.
- [16] Tianqi Chen, Mu Li, Yutian Li, Min Lin, Naiyan Wang, Minjie Wang, Tianjun Xiao, Bing Xu, Chiyuan Zhang, and Zheng Zhang. MXNet: A flexible and efficient machine learning library for heterogeneous distributed systems. arXiv preprint arXiv:1512.01274, 2015.
- [17] Xi Chen and Qiang Zhou. Sequential experimental designs for stochastic kriging. In *Proceedings of the 2014 Winter Simulation Conference*, WSC '14, pages 3821–3832. IEEE Press, 2014.
- [18] Kurt Cutajar, Mark Pullin, Andreas Damianou, Neil Lawrence, and Javier González. Deep Gaussian Processes for Multi-fidelity Modeling. *arXiv preprint arXiv:1903.07320*, 2019.
- [19] Andreas Damianou and Neil Lawrence. Deep Gaussian Processes. In *Artificial Intelligence and Statistics*, pages 207–215, 2013.
- [20] Samuel Daulton, Maximilian Balandat, and Eytan Bakshy. Differentiable Expected Hypervolume Improvement for Parallel Multi-Objective Bayesian Optimization. In Advances in Neural Information Processing Systems 33, 2020.
- [21] Katharina Eggensperger, Matthias Feurer, Aaron Klein, and Stefan Falkner. Hpolib2 (development branch), 2019. URL https://github.com/automl/HPOlib2.
- [22] Stefan Falkner, Aaron Klein, and Frank Hutter. BOHB: robust and efficient hyperparameter optimization at scale. CoRR, abs/1807.01774, 2018.
- [23] Qing Feng, Benjamin Letham, Hongzi Mao, and Eytan Bakshy. High-dimensional contextual policy search with unknown context rewards using Bayesian optimization. In Advances in Neural Information Processing Systems 33, NeurIPS, 2020.
- [24] Matthias Feurer, Aaron Klein, Katharina Eggensperger, Jost Springenberg, Manuel Blum, and Frank Hutter. Efficient and robust automated machine learning. In Advances in Neural Information Processing Systems 28, pages 2962–2970. 2015.
- [25] Peter Frazier, Warren Powell, and Savas Dayanik. The knowledge-gradient policy for correlated normal beliefs. *INFORMS journal on Computing*, 21(4):599–613, 2009.
- [26] Peter I Frazier. A tutorial on Bayesian optimization. arXiv preprint arXiv:1807.02811, 2018.
- [27] Peter I Frazier, Warren B Powell, and Savas Dayanik. A knowledge-gradient policy for sequential information collection. SIAM Journal on Control and Optimization, 47(5):2410–2439, 2008.
- [28] Jacob Gardner, Matt Kusner, Zhixiang, Kilian Weinberger, and John Cunningham. Bayesian optimization with inequality constraints. In *Proceedings of the 31st International Conference on Machine Learning*, volume 32, pages 937–945, 2014.
- [29] Jacob Gardner, Geoff Pleiss, Kilian Q Weinberger, David Bindel, and Andrew G Wilson. GPytorch: Blackbox matrix-matrix Gaussian process inference with GPU acceleration. In *Advances in Neural Information Processing Systems*, pages 7576–7586, 2018.
- [30] Jason Gauci, Edoardo Conti, Yitao Liang, Kittipat Virochsiri, Yuchen He, Zachary Kaden, Vivek Narayanan, and Xiaohui Ye. Horizon: Facebook's open source applied reinforcement learning platform. arXiv preprint arXiv:1811.00260, 2018.
- [31] Michael A. Gelbart, Jasper Snoek, and Ryan P. Adams. Bayesian optimization with unknown constraints. In *Proceedings of the 30th Conference on Uncertainty in Artificial Intelligence*, UAI, 2014.
- [32] David Ginsbourger, Janis Janusevskis, and Rodolphe Le Riche. Dealing with asynchronicity in parallel Gaussian process based global optimization. Technical report, 2011. URL https://hal.archives-ouvertes.fr/hal-00507632.
- [33] Rafael Gómez-Bombarelli, Jennifer N. Wei, David Duvenaud, JoséMiguel Hernández-Lobato, Benjamín Sánchez-Lengeling, Dennis Sheberla, Jorge Aguilera-Iparraguirre, Timothy D. Hirzel, Ryan P. Adams, and Alán Aspuru-Guzik. Automatic chemical design using a data-driven continuous representation of molecules. ACS Central Science, 4(2):268–276, 02 2018.

- [34] Javier González, Zhenwen Dai, Philipp Hennig, and Neil D. Lawrence. Batch Bayesian optimization via local penalization. In *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, AISTATS, pages 648–657, 2016.
- [35] GPy. GPy: A gaussian process framework in python. http://github.com/SheffieldML/GPy, since 2012.
- [36] Ryan-Rhys Griffiths and José Miguel Hernández-Lobato. Constrained Bayesian optimization for automatic chemical design. arXiv preprint arXiv:1709.05501, 2017.
- [37] Nikolaus Hansen and Andreas Ostermeier. Completely derandomized self-adaptation in evolution strategies. *Evol. Comput.*, 9(2):159–195, June 2001.
- [38] José Miguel Hernández-Lobato, Michael A. Gelbart, Matthew W. Hoffman, Ryan P. Adams, and Zoubin Ghahramani. Predictive entropy search for Bayesian optimization with unknown constraints. In *Proceedings of the 32nd International Conference on Machine Learning*, ICML, 2015.
- [39] Tito Homem-de-Mello. On rates of convergence for stochastic optimization problems under non-independent and identically distributed sampling. SIAM Journal on Optimization, 19(2):524–551, 2008.
- [40] Pavel Izmailov, Dmitrii Podoprikhin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. Averaging weights leads to wider optima and better generalization. arXiv preprint arXiv:1803.05407, 2018.
- [41] Pavel Izmailov, Wesley Maddox, Timur Garipov, Polina Kirichenko, Dmitry Vetrov, and Andrew Gordon Wilson. Subspace inference for Bayesian deep learning. In *Uncertainty in Artificial Intelligence*, 2019.
- [42] Achin Jain and Manfred Morari. Computing the racing line using Bayesian optimization. *arXiv e-prints*, page arXiv:2002.04794, February 2020.
- [43] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. arXiv preprint arXiv:1408.5093, 2014.
- [44] Shali Jiang, Daniel R. Jiang, Maximilian Balandat, Brian Karrer, Jacob Gardner, and Roman Garnett. Efficient nonmyopic Bayesian optimization via one-shot multi-step trees. In *Advances in Neural Information Processing Systems*, 2020.
- [45] D. R. Jones, C. D. Perttunen, and B. E. Stuckman. Lipschitzian optimization without the Lipschitz constant. *Journal of Optimization Theory and Applications*, 79(1):157–181, Oct 1993.
- [46] Donald R. Jones, Matthias Schonlau, and William J. Welch. Efficient global optimization of expensive black-box functions. *Journal of Global Optimization*, 13:455–492, 1998.
- [47] Kirthevasan Kandasamy, Jeff Schneider, and Barnabás Póczos. High dimensional Bayesian optimisation and bandits via additive models. In *Proceedings of the 32nd International Conference on Machine Learning*, ICML, 2015.
- [48] Kirthevasan Kandasamy, Gautam Dasarathy, Junier Oliva, Jeff Schneider, and Barnab'as P'oczos. Gaussian process bandit optimisation with multi-fidelity evaluations. In Advances in Neural Information Processing Systems, NIPS, 2016.
- [49] Kirthevasan Kandasamy, Karun Raju Vysyaraju, Willie Neiswanger, Biswajit Paria, Christopher R. Collins, Jeff Schneider, Barnabas Póczos, and Eric P. Xing. Tuning hyperparameters without grad students: Scalable and robust Bayesian optimisation with Dragonfly. arXiv e-prints, art. arXiv:1903.06694, Mar 2019.
- [50] Diederik P Kingma and Max Welling. Auto-Encoding Variational Bayes. arXiv e-prints, page arXiv:1312.6114, Dec 2013.
- [51] A. Klein, S. Falkner, S. Bartels, P. Hennig, and F. Hutter. Fast Bayesian optimization of machine learning hyperparameters on large datasets. *CoRR*, 2016.
- [52] A. Klein, S. Falkner, N. Mansur, and F. Hutter. Robo: A flexible and robust Bayesian optimization framework in Python. In NIPS 2017 Bayesian Optimization Workshop, December 2017.
- [53] Anton J Kleywegt, Alexander Shapiro, and Tito Homem-de Mello. The sample average approximation method for stochastic discrete optimization. SIAM Journal on Optimization, 12(2):479–502, 2002.

- [54] J. Knowles. ParEGO: A hybrid algorithm with on-line landscape approximation for expensive multiobjective optimization problems. *IEEE Transactions on Evolutionary Computation*, 10(1):50–66, 2006.
- [55] Nicolas Knudde, Joachim van der Herten, Tom Dhaene, and Ivo Couckuyt. GPflowOpt: A Bayesian Optimization Library using TensorFlow. *arXiv* preprint *arXiv*:1711.03845, 2017.
- [56] R. Kohavi and B. Becker. UCI machine learning repository, 1996. URL http://archive.ics.uci.edu/ml.
- [57] Benjamin Letham and Eytan Bakshy. Bayesian optimization for policy search via online-offline experimentation. *Journal of Machine Learning Research*, 20(145):1–30, 2019.
- [58] Benjamin Letham, Brian Karrer, Guilherme Ottoni, and Eytan Bakshy. Constrained Bayesian optimization with noisy experiments. *Bayesian Analysis*, 14(2):495–519, 2019.
- [59] Benjamin Letham, Roberto Calandra, Akshara Rai, and Eytan Bakshy. Re-examining linear embeddings for high-dimensional Bayesian optimization. In Advances in Neural Information Processing Systems 33, NeurIPS, 2020.
- [60] Cheng Li, David Rubín de Celis Leal, Santu Rana, Sunil Gupta, Alessandra Sutti, Stewart Greenhill, Teo Slezak, Murray Height, and Svetha Venkatesh. Rapid Bayesian optimisation for synthesis of short polymer fiber materials. *Scientific reports*, 7(1):5683, 2017.
- [61] Wesley Maddox, Timur Garipov, Pavel Izmailov, Dmitry Vetrov, and Andrew Gordon Wilson. A simple baseline for Bayesian uncertainty in deep learning. In Advances in Neural Information Processing Systems, 2019.
- [62] Wesley J Maddox, Shuai Tang, Pablo Garcia Moreno, Andrew Gordon Wilson, and Andreas Damianou. On Transfer Learning via Linearized Neural Networks. In NeurIPS Workshop on Meta-Learning (MetaLearn 2019), 2019.
- [63] Malaria Atlas Project. Malaria atlas project, 2019. URL https://map.ox.ac.uk/malaria-burden-data-download.
- [64] Alexander G. de G. Matthews, Mark van der Wilk, Tom Nickson, Keisuke. Fujii, Alexis Boukouvalas, Pablo León-Villagrá, Zoubin Ghahramani, and James Hensman. GPflow: A Gaussian process library using TensorFlow. *Journal of Machine Learning Research*, 18(40):1–6, apr 2017.
- [65] Richard M Meyer. Essential mathematics for applied fields. Springer Science & Business Media, 2012.
- [66] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. arXiv preprint arXiv:1312.5602, 2013.
- [67] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015.
- [68] Riccardo Moriconi, K. S. Sesh Kumar, and Marc Peter Deisenroth. High-Dimensional Bayesian Optimization with Manifold Gaussian Processes. arXiv e-prints, page arXiv:1902.10675, Feb 2019.
- [69] Fabio Muratore, Christian Eilers, Michael Gienger, and Jan Peters. Bayesian Domain Randomization for Sim-to-Real Transfer. *arXiv e-prints*, page arXiv:2003.02471, March 2020.
- [70] Iain Murray. Differentiation of the Cholesky decomposition. arXiv e-prints, page arXiv:1602.07527, Feb 2016.
- [71] Radford M Neal. *Bayesian learning for neural networks*, volume 118. Springer Science & Business Media, 1996.
- [72] Willie Neiswanger, Kirthevasan Kandasamy, Barnabas Póczos, Jeff Schneider, and Eric Xing. ProBO: a Framework for Using Probabilistic Programming in Bayesian Optimization. *arXiv e-prints*, page arXiv:1901.11515, January 2019.
- [73] Thanh V. Nguyen, Youssef Mroueh, Samuel Hoffman, Payel Das, Pierre Dognin, Giuseppe Romano, and Chinmay Hegde. Nano-material configuration design with deep surrogate langevin dynamics. In ICLR 2020 Workshop on Integration of Deep Neural Models and Differential Equations, 2019.

- [74] Vu Nguyen, Sunil Gupta, Santu Rana, Cheng Li, and Svetha Venkatesh. Predictive variance reduction search. In NIPS 2017 Workshop on Bayesian Optimization, Dec 2017.
- [75] A O'Hagan. On curve fitting and optimal design for regression. J. Royal Stat. Soc. B, 40:1–32, 1978.
- [76] Michael A Osborne. Bayesian Gaussian processes for sequential prediction, optimisation and quadrature. PhD thesis, Oxford University, UK, 2010.
- [77] Art B. Owen. Randomly permuted (t,m,s)-nets and (t, s)-sequences. In Harald Niederreiter and Peter Jau-Shyong Shiue, editors, *Monte Carlo and Quasi-Monte Carlo Methods in Scientific Computing*, pages 299–317, New York, NY, 1995. Springer New York.
- [78] Art B Owen. Quasi-monte carlo sampling. Monte Carlo Ray Tracing: Siggraph, 1:69–88, 2003.
- [79] Art B. Owen and Daniel Rudolf. A strong law of large numbers for scrambled net integration. *arXiv e-prints*, page arXiv:2002.07859, February 2020.
- [80] B. Paria, K. Kandasamy, and B. Póczos. A Flexible Multi-Objective Bayesian Optimization Approach using Random Scalarizations. ArXiv e-prints, May 2018.
- [81] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch. 2017.
- [82] Geoff Pleiss, Jacob R Gardner, Kilian Q Weinberger, and Andrew Gordon Wilson. Constant-time predictive distributions for gaussian processes. In *International Conference on Machine Learning*, 2018.
- [83] Matthias Poloczek, Jialei Wang, and Peter Frazier. Multi-information source optimization. In *Advances in Neural Information Processing Systems*, pages 4288–4298, 2017.
- [84] Carl Edward Rasmussen and Christopher KI Williams. Gaussian processes for machine learning. 2006. The MIT Press, Cambridge, MA, USA, 38:715–719, 2006.
- [85] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *Proceedings of the 31st International Conference on International Conference on Machine Learning Volume 32*, ICML'14, pages II–1278–II–1286. JMLR.org, 2014.
- [86] Mark Rowland, Krzysztof M Choromanski, François Chalus, Aldo Pacchiano, Tamas Sarlos, Richard E Turner, and Adrian Weller. Geometrically coupled monte carlo sampling. In Advances in Neural Information Processing Systems 31, pages 195–206. 2018.
- [87] Yunus Saatci and Andrew G Wilson. Bayesian GAN. In Advances in neural information processing systems, pages 3622–3631, 2017.
- [88] Hugh Salimbeni and Marc Peter Deisenroth. Doubly stochastic variational inference for deep gaussian processes. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 4588–4599. Curran Associates, Inc., 2017.
- [89] Matthias Schonlau, William J. Welch, and Donald R. Jones. Global versus local search in constrained optimization of computer models. *Lecture Notes-Monograph Series*, 34:11–25, 1998.
- [90] Warren Scott, Peter Frazier, and Warren Powell. The correlated knowledge gradient for simulation optimization of continuous parameters using Gaussian process regression. *SIAM Journal of Optimization*, 21:996–1026, 2011.
- [91] S. Seo, M. Wallat, T. Graepel, and K. Obermayer. Gaussian process regression: active data selection and test point rejection. In *Proceedings of the IEEE-INNS-ENNS International Joint Conference on Neural Networks. IJCNN 2000. Neural Computing: New Challenges and Perspectives for the New Millennium*, volume 3, pages 241–246 vol.3, July 2000.
- [92] Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P. Adams, and Nando de Freitas. Taking the human out of the loop: A review of Bayesian optimization. *Proceedings of the IEEE*, 104:1–28, 2016.
- [93] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556, 2014.
- [94] Jasper Snoek, Hugo Larochelle, and Ryan P Adams. Practical Bayesian optimization of machine learning algorithms. In *Advances in neural information processing systems*, pages 2951–2959, 2012.

- [95] Jasper Snoek, Kevin Swersky, Richard Zemel, and Ryan P Adams. Input warping for Bayesian optimization of non-stationary functions. In *Proceedings of the 31st International Conference on Machine Learning*, ICML'14, 2014.
- [96] J. T. Springenberg, A. Klein, S.Falkner, and F. Hutter. Bayesian optimization with robust Bayesian neural networks. In *Advances in Neural Information Processing Systems* 29, December 2016.
- [97] The Emukit authors. Emukit: Emulation and uncertainty quantification for decision making. https://github.com/amzn/emukit, 2018.
- [98] The GPyOpt authors. GPyOpt: A Bayesian optimization framework in Python. http://github.com/ SheffieldML/GPyOpt, 2016.
- [99] Dustin Tran, Matthew D Hoffman, Rif A Saurous, Eugene Brevdo, Kevin Murphy, and David M Blei. Deep probabilistic programming. *Proceedings of the International Conference on Learning Representations* (*ICLR*), 2017.
- [100] Jialei Wang, Scott C Clark, Eric Liu, and Peter I Frazier. Parallel Bayesian global optimization of expensive functions. *arXiv preprint arXiv:1602.05149*, 2016.
- [101] Zi Wang and Stefanie Jegelka. Max-value entropy search for efficient Bayesian optimization. volume 70 of *Proceedings of Machine Learning Research*, pages 3627–3635, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR.
- [102] Ziyu Wang, Frank Hutter, Masrour Zoghi, David Matheson, and Nando De Freitas. Bayesian optimization in a billion dimensions via random embeddings. J. Artif. Int. Res., 55(1):361–387, January 2016.
- [103] Andrew Gordon Wilson, Zhiting Hu, Ruslan Salakhutdinov, and Eric P Xing. Deep kernel learning. In Artificial Intelligence and Statistics, pages 370–378, 2016.
- [104] J. T. Wilson, R. Moriconi, F. Hutter, and Marc Peter Deisenroth. The reparameterization trick for acquisition functions. *ArXiv e-prints*, December 2017.
- [105] James Wilson, Frank Hutter, and Marc Peter Deisenroth. Maximizing acquisition functions for Bayesian optimization. In *Advances in Neural Information Processing Systems 31*, pages 9905–9916. 2018.
- [106] Jian Wu and Peter Frazier. The parallel knowledge gradient method for batch Bayesian optimization. In *Advances in Neural Information Processing Systems*, pages 3126–3134, 2016.
- [107] Jian Wu and Peter Frazier. Practical two-step lookahead Bayesian optimization. In *Advances in Neural Information Processing Systems* 32, 2019.
- [108] Jian Wu and Peter I. Frazier. Discretization-free Knowledge Gradient Methods for Bayesian Optimization. arXiv e-prints, page arXiv:1707.06541, Jul 2017.
- [109] Jian Wu, Matthias Poloczek, Andrew Gordon Wilson, and Peter I Frazier. Bayesian optimization with gradients. In *Advances in Neural Information Processing Systems*, pages 5267–5278, 2017.
- [110] Jian Wu, Saul Toscano-Palmerin, Peter I. Frazier, and Andrew Gordon Wilson. Practical multi-fidelity Bayesian optimization for hyperparameter tuning. *CoRR*, abs/1903.04703, 2019.
- [111] Yichi Zhang, Daniel W Apley, and Wei Chen. Bayesian optimization for materials design with mixed quantitative and qualitative variables. *Scientific Reports*, 10(1):1–13, 2020.

Appendix to:

BOTORCH: A Framework for Efficient Monte-Carlo Bayesian Optimization

A Brief Overview of Other Software Packages for BO

One of the earliest commonly-used packages is **Spearmint** [94], which implements a variety of modeling techniques such as MCMC hyperparameter sampling and input warping [95]. Spearmint also supports parallel optimization via fantasies, and constrained optimization with the expected improvement and predictive entropy search acquisition functions [31, 38]. Spearmint was among the first libraries to make BO easily accessible to the end user.

GPyOpt [98] builds on the popular GP regression framework GPy [35]. It supports a similar set of features as Spearmint, along with a local penalization-based approach for parallel optimization [34]. It also provides the ability to customize different components through an alternative, more modular API.

Cornell-MOE [106] implements the Knowledge Gradient (KG) acquisition function, which allows for parallel optimization, and includes recent advances such as large-scale models incorporating gradient evaluations [109] and multi-fidelity optimization [110]. Its core is implemented in C++, which provides performance benefits but renders it hard to modify and extend.

RoBO [52] implements a collection of models and acquisition functions, including Bayesian neural nets [96] and multi-fidelity optimization [51].

Emukit [97] is a Bayesian optimization and active learning toolkit with a collection of acquisition functions, including for parallel and multi-fidelity optimization. It does not provide specific abstractions for implementing new algorithms, but rather specifies a model API that allows it to be used with the other toolkit components.

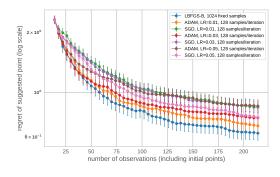
The recent **Dragonfly** [49] library supports parallel optimization, multi-fidelity optimization [48], and high-dimensional optimization with additive kernels [47]. It takes an ensemble approach and aims to work out-of-the-box across a wide range of problems, a design choice that makes it relatively hard to extend.

B Parallelism and Hardware Acceleration

B.1 Batch Evaluation

Batch evaluation, an important element of modern computing, enables automatic dispatch of independent operations across multiple computational resources (e.g. CPU and GPU cores) for parallelization and memory sharing. All BoTorch components support batch evaluation, which makes it easy to write concise and highly efficient code in a platform-agnostic fashion. Batch evaluation enables fast queries of acquisition functions at a large number of candidate sets in parallel, facilitating novel initialization heuristics and optimization techniques.

Specifically, instead of sequentially evaluating an acquisition function at a number of candidate sets $\mathbf{x}_1,\ldots,\mathbf{x}_b$, where $\mathbf{x}_k\in\mathbb{R}^{q\times d}$ for each k, BOTORCH evaluates a batched tensor $\mathbf{x}\in\mathbb{R}^{b\times q\times d}$. Computation is automatically distributed so that, depending on the hardware used, speedups can be close to linear in the batch size b. Batch evaluation is also heavily used in computing MC acquisition functions, with the effect that significantly increasing the number of MC samples often has little impact on wall time. In Figure 4 we observe significant speedups from running on the GPU, with



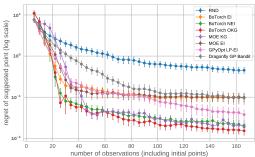
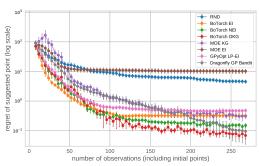


Figure 8: Stochastic/deterministic opt. of EI on Hartmann6

Figure 9: Branin (d=2)



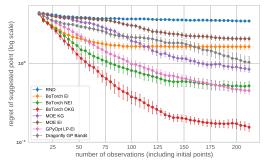


Figure 10: Rosenbrock (d = 3)

Figure 11: Ackley (d = 5)

scaling essentially linear in the batch size, except for very large b and N. The fixed cost due to communication overhead renders CPU evaluation faster for small batch and sample sizes.

C Additional Empirical Results

This section describes a number of empirical results that were omitted from the main paper due to space constraints.

C.1 Synthetic Functions

Algorithms start from the same set of 2d+2 QMC sampled initial points for each trial, with d the dimension of the design space. We evaluate based on the true noiseless function value at the "suggested point" (i.e., the point to be chosen if BO were to end at this batch). OKG, MOE KG, and NEI use "out-of-sample" suggestions (introduced as χ_n in Section D.4), while the others use "in-sample" suggestions [26].

All functions are evaluated with noise generated from a $\mathcal{N}(0, .25)$ distribution. Figures 9-11 give the results for all synthetic functions from Section 6. The results show that BOTORCH's NEI and OKG acquisition functions provide highly competitive performance in all cases.

C.2 One-Shot KG Computational Scaling

Figure 12 shows the wall time for generating a set of q=8 candidates as a function of the number of total data points n for both standard (Cholesky-based) as well as scalable (Linear CG) posterior inference methods, on both CPU and GPU. While the GPU variants have a significant overhead for small models, they are significantly faster for larger models. Notably, our SAA based OKG is significantly faster than $MOE\ KG$, while at the same time achieving much better optimization performance (Figure 13).

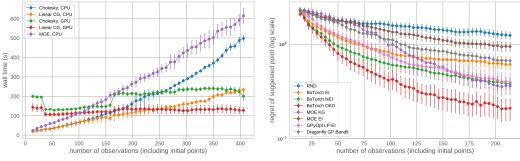


Figure 12: KG wall times

Figure 13: Hartmann (d = 6), noisy, best suggested

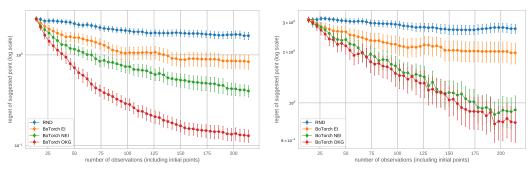


Figure 14: Constrained Hartmann6, $f_2(x) = ||x||_1 - 3$ Figure 15: Constrained Hartmann6, $f_1(x) = ||x||_2 - 1$

C.3 Constrained Bayesian Optimization

We present results for constrained BO on a synthetic function. We consider a multi-output function $f = (f_1, f_2)$ and the optimization problem:

$$\max_{x \in \mathbb{X}} f_1(x) \quad \text{s.t.} \quad f_2(x) \le 0. \tag{8}$$

Both f_1 and f_2 are observed with $\mathcal{N}(0,0.5^2)$ noise and we model the two components using independent GP models. A constraint-weighted composite objective is used in each of the BoTorch acquisition functions EI, NEI, and OKG.

Results for the case of a Hartmann6 objective and two types of constraints are given in Figures 14-15 (we only show results for BoTorch's algorithms, since the other packages do not natively support optimization subject to unknown constraints).

The regret values are computed using a feasibility-weighted objective, where "infeasible" is assigned an objective value of zero. For random search and EI, the suggested point is taken to be the best feasible noisily observed point, and for NEI and OKG, we use out-of-sample suggestions by optimizing the feasibility-weighted version of the posterior mean. The results displayed in Figure 15 are for the constrained Hartmann6 benchmark from [58]. Note, however, that the results here are not directly comparable to the figures in [58] because (1) we use feasibility-weighted objectives to compute regret and (2) they follow a different convention for suggested points. We emphasize that our contribution of outcome constraints for the case of KG has not been shown before in the literature.

C.4 Hyperparameter Optimization Details

This section gives further detail on the experimental settings used in each of the hyperparameter optimization problems. As HPO typically involves long and resource intensive training jobs, it is standard to select the configuration with the best observed performance, rather than to evaluate a "suggested" configuration (we cannot perform noiseless function evaluations).

DQN and Cartpole: We consider the case of tuning a deep Q-network (DQN) learning algorithm [66, 67] on the *Cartpole* task from OpenAI gym [12] and the default DQN agent implemented in

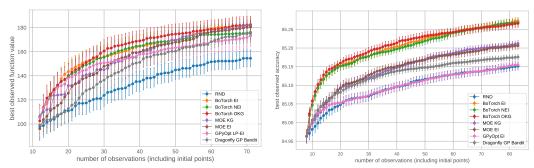


Figure 16: DQN tuning benchmark (Cartpole)

Figure 17: NN surrogate model, best observed accuracy

Horizon [30]. Figure 16 shows the results of tuning five hyperparameters, exploration parameter ("epsilon"), the target update rate, the discount factor, the learning rate, and the learning rate decay. We allow for a maximum of 60 training episodes or 2000 training steps, whichever occurs first. To reduce noise, each "function evaluation" is taken to be an average of 10 independent training runs of DQN. Figure 16 presents the optimization performance of various acquisition functions from the different packages, using 15 rounds of parallel evaluations of size q=4, over 100 trials. While in later iterations all algorithms achieve reasonable performance, BoTorch OKG, EI, NEI, and GPyOpt LP-EI show faster learning early on.

Neural Network Surrogate: We consider the neural network surrogate model for the UCI Adult data set introduced by Falkner et al. [22], which is available as part of HPOlib2 [21]. We use a surrogate model to achieve a high level of precision in comparing the performance of the algorithms without incurring excessive computational training costs. This is a six-dimensional problem over network parameters (*number of layers*, *units per layer*) and training parameters (*initial learning rate*, *batch size*, *dropout*, *exponential decay factor for learning rate*). Figure 17 shows optimization performance in terms of best observed classification accuracy. Results are means and 95% confidence intervals computed from 200 trials with 75 iterations of size q=1. All BoTorch algorithms perform quite similarly here, with OKG doing slightly better in earlier iterations. Notably, they all achieve significantly better accuracy than all other algorithms.

Stochastic Weight Averaging on CIFAR-10: Our final example is for the recently proposed Stochastic Weight Averaging (SWA) procedure of Izmailov et al. [40], for which good hyperparameter settings are not fully understood. The setting is 300 epochs of training on the VGG-16 [93] architecture for CIFAR-10. We tune three SWA hyperparameters: learning rate, update frequency, and starting iteration using OKG. Izmailov et al. [40] report the mean and standard deviation of the test accuracy over three runs to be 93.64 and 0.18, respectively, which corresponds to a 95% confidence interval of 93.64 ± 0.20 . We tune the problem to an average accuracy of 93.84 ± 0.03 .

D Additional Theoretical Results and Omitted Proofs

D.1 General SAA Results

Recall that we assume that $f(\mathbf{x}) \sim h(\mathbf{x}, \epsilon)$ for some $h : \mathbb{X} \times \mathbb{R}^s \to \mathbb{R}^{q \times m}$ and base random variable $\epsilon \in \mathbb{R}^s$ (c.f. Section 4 for an explicit expression for h in case of a GP model). We write

$$A(\mathbf{x}, \epsilon) := a(g(h(\mathbf{x}, \epsilon))). \tag{9}$$

Theorem 3 (Homem-de-Mello [39]). Suppose that (i) \mathbb{X} is a compact metric space, (ii) $\hat{\alpha}_N(\mathbf{x}) \xrightarrow{a.s.} \alpha(\mathbf{x})$ for all $\mathbf{x} \in \mathbb{X}^q$, and (iii) there exists an integrable function $\ell : \mathbb{R}^s \mapsto \mathbb{R}$ such that for almost every ϵ and all $\mathbf{x}, \mathbf{y} \in \mathbb{X}$,

$$|A(\mathbf{x}, \epsilon) - A(\mathbf{y}, \epsilon)| \le \ell(\epsilon) \|\mathbf{x} - \mathbf{y}\|. \tag{10}$$

Then $\hat{\alpha}_N^* \xrightarrow{a.s.} \alpha^*$ and $\operatorname{dist}(\hat{\mathbf{x}}_N^*, \mathcal{X}_f^*) \xrightarrow{a.s.} 0$.

Proposition 1. Suppose that (i) \mathbb{X} is a compact metric space, (ii) f is a GP with continuously differentiable prior mean and covariance functions, and (iii) $g(\cdot)$ and $a(\cdot, \Phi)$ are Lipschitz continuous. Then, condition (10) in Theorem 3 holds.

The following proposition follows directly from Proposition 2.1, Theorem 2.3, and remarks on page 528 of [39].

Proposition 2 (Homem-de-Mello [39]). Suppose that, in addition to the conditions in Theorem 3, (i) the base samples $E = \{\epsilon^i\}_{i=1}^N$ are i.i.d., (ii) for all $\mathbf{x} \in \mathbb{X}^q$ the moment generating function $M_{\mathbf{x}}^A(t) := \mathbb{E}[e^{tA(\mathbf{x},\epsilon)}]$ of $A(\mathbf{x},\epsilon)$ is finite in an open neighborhood of t=0 and (iii) the moment generating function $M^\ell(t) := \mathbb{E}[e^{t\ell(\epsilon)}]$ is finite in an open neighborhood of t=0. Then, there exist $K < \infty$ and $\beta > 0$ such that $\mathbb{P}(\operatorname{dist}(\hat{\mathbf{x}}_N, \mathcal{X}_f^*)) \le Ke^{-\beta N}$ for all $N \ge 1$.

D.2 Formal Statement of Theorem 1

Theorem 1 (Formal Version). Suppose (i) \mathbb{X} is compact, (ii) f has a GP prior with continuously differentiable mean and covariance functions, and (iii) $g(\cdot)$ and $a(\cdot, \Phi)$ are Lipschitz continuous. If the base samples $\{\epsilon^i\}_{i=1}^N$ are drawn i.i.d. from $\mathcal{N}(0,1)$, then

- (1) $\hat{\alpha}_N^* \rightarrow \alpha^*$ a.s., and
- (2) $\operatorname{dist}(\hat{\mathbf{x}}_N^*, \mathcal{X}^*) \to 0$ a.s.

If, in addition, (iii) for all $\mathbf{x} \in \mathbb{X}^q$ the moment generating function $M_{\mathbf{x}}^A(t) := \mathbb{E}[e^{tA(\mathbf{x},\epsilon)}]$ of $A(\mathbf{x},\epsilon)$ is finite in an open neighborhood of t=0 and (iv) the moment generating function $M^{\ell}(t) := \mathbb{E}[e^{t\ell(\epsilon)}]$ is finite in an open neighborhood of t=0, then

(3)
$$\forall \delta > 0, \exists K < \infty, \beta > 0 \text{ s.t. } \mathbb{P}(\operatorname{dist}(\hat{\mathbf{x}}_N^*, \mathcal{X}^*) > \delta) < Ke^{-\beta N} \text{ for all } N > 1.$$

D.3 Randomized Quasi-Monte Carlo Sampling for Sample Average Approximation

In order to use randomized QMC methods with SAA for MC acquisition function, the base samples $E = \{\epsilon^i\}$ will need to be generated via RQMC. For the case of Normal base samples, this can be achieved in various ways, e.g. by using inverse CDF methods or a suitable Box-Muller transform of samples $\epsilon^i \in [0,1]^s$ (both approaches are implemented in BOTORCH). In the language of Section 4, such a transform will become part of the base sample transform $\epsilon \mapsto h(\mathbf{x},\epsilon)$ for any fixed \mathbf{x} .

For the purpose of this paper, we consider scrambled (t,d)-sequences as discussed by Owen [77], which are a particular class of RQMC method (Botorch uses PyTorch's implementation of scrambled Sobol sequences, which are (t,d)-nets in base 2). Using recent theoretical advances from Owen and Rudolf [79], it is possible to generalize the convergence results from Theorems 1 and 2 to the RQMC setting (to our knowledge, this is the first practical application of these theoretical results).

Let $(N_i)_{i\geq 1}$ be a sequence with $N_i\in\mathbb{N}$ s.t. $N_i\to\infty$ as $i\to\infty$. Then we have the following (see Appendix D.5 for the proofs):

Theorem 1(q). In the setting of Theorem 1, let $\{\epsilon^i\}$ be samples from a (t,d)-sequence in base b with gain coefficients no larger than $\Gamma < \infty$, randomized using a nested uniform scramble as in [77]. Then, the conclusions of Theorem 1 still hold. In particular,

- (1) $\hat{\alpha}_{N_i}^* \to \alpha^* \text{ a.s. as } i \to \infty$,
- (2) $\operatorname{dist}(\hat{\mathbf{x}}_{N_i}^*, \mathcal{X}^*) \to 0 \text{ a.s. as } i \to \infty,$

(3)
$$\forall \delta > 0$$
, $\exists K < \infty$, $\beta > 0$ s.t. $\mathbb{P}(\operatorname{dist}(\hat{\mathbf{x}}_{N_i}^*, \mathcal{X}^*) > \delta) \leq Ke^{-\beta N_i}$ for all $i \geq 1$.

Theorem 2(q). In the setting of Theorem 2, let $\{\epsilon^i\}$ be samples from (t,d)-sequence in base b, with gain coefficients no larger than $\Gamma < \infty$, randomized using a nested uniform scramble as in [77]. Then.

(1)
$$\hat{\alpha}_{\mathrm{KG},N_i}^* \xrightarrow{a.s.} \alpha_{\mathrm{KG}}^* \text{ as } i \to \infty$$
,

(2)
$$\operatorname{dist}(\hat{\mathbf{x}}_{\mathrm{KG},N_i}^*, \mathcal{X}_{\mathrm{KG}}^*) \xrightarrow{a.s.} 0 \text{ as } i \to \infty.$$

Theorem 2(q) as stated does not provide a rate on the convergence of the optimizer. We believe that such result is achievable, but leave it to future work.

Note that while the above results hold for any sequence $(N_i)_i$ with $N_i \to \infty$, in practice the RQMC integration error can be minimized by using sample sizes that exploit intrinsic symmetry of the (t,d)-sequences. Specifically, for integers $b \ge 2$ and $M \ge 1$, let

$$\mathcal{N} := \{ mb^k \mid m \in \{1, \dots, M\}, k \in \mathbb{N}_+ \}. \tag{11}$$

In practice, we chose the MC sample size N from the unique elements of \mathcal{N} .

D.4 Asymptotic Optimality of OKG

Consider the case where f_{true} is drawn from a GP prior with $f \stackrel{d}{=} f_{\text{true}}$, and that $g(f) \equiv f$. The KG policy (i.e., when used to select sequential measurements in a dynamic setting) is known to be $asymptotically\ optimal\ [27,25,83,7]$, meaning that as the number of measurements tends to infinity, an optimal point $x^* \in \mathcal{X}_f^* := \arg\max_{x \in \mathbb{X}} f(x)$ is identified. Although it does not necessarily signify good finite sample performance, this is considered a useful property for acquisition functions [25]. In this section, we state two results showing that OKG also possesses this property, providing further theoretical justification for the MC approach taken by BoTorch.

Let \mathcal{D}_0 be the initial data and \mathcal{D}_n for $n \geq 1$ be the data generated by taking measurements according to OKG using N_n MC samples in iteration n, i.e., $\mathbf{x}_{n+1} \in \arg\max_{\mathbf{x} \in \mathbb{X}^q} \hat{\alpha}_{\mathrm{KG},N_n}(\mathbf{x};\mathcal{D}_n)$ for all n, and let $\chi_n \in \arg\max_{\mathbf{x} \in \mathbb{X}} \mathbb{E}[f(x) \mid \mathcal{D}_n]$. Then we can show the following:

Theorem 4. Suppose conditions (i) and (ii) of Theorem 1 and (iii) of Theorem 2 are satisfied. In addition, suppose that $\limsup_n N_n = \infty$. Then, $f(\chi_n) \to f(x^*)$ a.s. and in L^1 .

Theorem 4 shows that OKG is asymptotically optimal if the number of fantasies N_n grows asymptotically with n (this assumes we have an analytic expression for the inner expectation. If not, a similar condition must be imposed on the number of inner MC samples). In the special case of finite \mathbb{X} , we can quantify the sample sizes $\{N_n\}$ that ensure asymptotic optimality of OKG:

Theorem 5. Along with conditions (i) and (ii) of Theorem 1, suppose that $|\mathbb{X}| < \infty$ and q = 1. Then, if for some $\delta > 0$, $N_n \ge A_n^{-1} \log(K_n/\delta)$ a.s., where A_n and K_n are a.s. finite and depend on \mathcal{D}_n (these quantities can be computed), we have $f(\chi_n) \to \max_{x \in \mathbb{X}} f(x)$ a.s..

D.5 Proofs

In the following, we will denote by $\mu_{\mathcal{D}}(x) := \mathbb{E}[f(x) \mid \mathcal{D}]$ and $K_{\mathcal{D}}(x,y) := \mathbb{E}[(f(x) - \mathbb{E}[f(x)])(f(y) - \mathbb{E}[f(y)])^T \mid \mathcal{D}]$ the posterior mean and covariance functions of f conditioned on data \mathcal{D} , respectively. Under some abuse of notation, we will use $\mu_{\mathcal{D}}(\mathbf{x})$ and $K_{\mathcal{D}}(\mathbf{x},\mathbf{y})$ to denote multi point (vector / matrix)-valued variants of $\mu_{\mathcal{D}}$ and $K_{\mathcal{D}}$, respectively. If f has a GP prior, then the posterior mean and covariance $\mu_{\mathcal{D}}(\mathbf{x})$ and $K_{\mathcal{D}}(\mathbf{x},\mathbf{y})$ have well-known explicit expressions [84].

For notational simplicity and without loss of generality, we will focus on single-output GP case (m=1) in this section. Indeed, in the multi-output case (m>1), we have a GP f on $\mathbb{X} \times \mathbb{M}$ with $\mathbb{M} = \{1,\ldots,m\}$, and covariance function $(x_1,i_1),(x_2,i_2) \mapsto \tilde{K}((x_1,i_1),(x_2,i_2))$. For q=1 we then define $x\mapsto \tilde{f}(x):=[f(x,0),...,f(x,m)]$, and then stack these for q>1: $\mathbf{x}\mapsto [\tilde{f}(\mathbf{x}_1)^T,...,\tilde{f}(\mathbf{x}_q)^T]^T$. Then the analysis in the proofs below can be done on mq-dimensional and $mq\times mq$ -dimensional posterior mean and covariance matrices (instead of q and $q\times q$ dimensional ones for m=1). Differentiability assumptions are needed only to establish certain boundedness results (e.g. in the proof of Proposition 1), but \mathbb{M} is finite, so we will require differentiability of $K((\cdot,i_1),(\cdot,i_2))$ for each i_1 and i_2 . Assumptions on other quantities can be naturally extended (e.g. for Theorem 2 g will need to be Lipschitz on $\mathbb{R}^{q\times m}$ rather than on \mathbb{R}^q , etc.).

Proof of Proposition 1. Without loss of generality, we may assume m=1 (the multi-output GP case follows immediately from applying the result below to q'=qm and re-arranging the output). For a GP, we have $h_{\mathcal{D}}(\mathbf{x},\epsilon)=\mu_{\mathcal{D}}(\mathbf{x})+L_{\mathcal{D}}(\mathbf{x})\epsilon$ with $\epsilon\sim\mathcal{N}(0,I_q)$, where $\mu_{\mathcal{D}}(\mathbf{x})$ is the posterior mean and $L_{\mathcal{D}}(\mathbf{x})$ is the Cholesky decomposition of the posterior covariance $M_{\mathcal{D}}(\mathbf{x})$. It is easy to verify from the classic GP inference equations [84] that if prior mean and covariance function are continuously differentiable, then so are posterior mean $\mu_{\mathcal{D}}(\cdot)$ and covariance $K_{\mathcal{D}}(\cdot)$. Since the Cholesky decomposition is also continuously differentiable [70], so is $L_{\mathcal{D}}(\cdot)$. As $\mathbb X$ is compact and $\mu_{\mathcal{D}}(\cdot)$ are continuously differentiable, their derivatives are bounded. It follows from

the mean value theorem that there exist $C_{\mu}, C_L < \infty$ s.t. $\|\mu_{\mathcal{D}}(\mathbf{x}) - \mu_{\mathcal{D}}(\mathbf{y})\| \le C_{\mu} \|\mathbf{x} - \mathbf{y}\|$ and $\|(L_{\mathcal{D}}(\mathbf{x}) - L_{\mathcal{D}}(\mathbf{y}))\epsilon\| \le C_L \|\epsilon\| \|\mathbf{x} - \mathbf{y}\|$. Thus,

$$||h_{\mathcal{D}}(\mathbf{x}, \epsilon) - h_{\mathcal{D}}(\mathbf{y}, \epsilon)|| = ||\mu_{\mathcal{D}}(\mathbf{x}) - \mu(\mathbf{y}) + (L_{\mathcal{D}}(\mathbf{x}) - L_{\mathcal{D}}(\mathbf{y}))\epsilon||$$

$$\leq ||\mu_{\mathcal{D}}(\mathbf{x}) - \mu_{\mathcal{D}}(\mathbf{y})|| + ||(L_{\mathcal{D}}(\mathbf{x}) - L_{\mathcal{D}}(\mathbf{y}))\epsilon||$$

$$\leq \ell_h(\epsilon)||\mathbf{x} - \mathbf{y}||$$

where $\ell_h(\epsilon) := C_\mu + C_L \|\epsilon\|$. Since, by assumption, $g(\cdot)$ and $a(\cdot; \Phi)$ are Lipschitz (say with constants L_a and L_g , respectively), it follows that $\|A(\mathbf{x}, \epsilon) - A(\mathbf{y}, \epsilon)\| \le L_a L_g \ell_h(\epsilon) \|\mathbf{x} - \mathbf{y}\|$. It thus suffices to show that $\ell_h(\epsilon)$ is integrable. To see this, note that $|\ell_h(\epsilon)| \le C_\mu + C_L C \sum_i |\epsilon_i|$ for some $C < \infty$ (equivalence of norms), and that $\epsilon_i \sim \mathcal{N}(0, 1)$ is integrable.

Lemma 1. Suppose that (i) f is a GP with continuously differentiable prior mean and covariance function, and (ii) that $a(\cdot, \Phi)$ and $g(\cdot)$ are Lipschitz. Then, for all $\mathbf{x} \in \mathbb{X}^q$ the moment generating functions $M_{\mathbf{x}}^A(t) := \mathbb{E}[e^{tA(\mathbf{x},\epsilon)}]$ of $A(\mathbf{x},\epsilon)$ and $M^\ell(t) := \mathbb{E}[e^{t\ell(\epsilon)}]$ are finite for all $t \in \mathbb{R}$.

Proof of Lemma 1. Recall that $h_{\mathcal{D}}(\mathbf{x}, \epsilon) = \mu_{\mathcal{D}}(\mathbf{x}) + L_{\mathcal{D}}(\mathbf{x})\epsilon$ for the case of f being a GP, where $\mu_{\mathcal{D}}(\mathbf{x})$ is the posterior mean and $L_{\mathcal{D}}(\mathbf{x})$ is the Cholesky decomposition of the posterior covariance $K_{\mathcal{D}}(\mathbf{x})$. Mirroring the argument from the proof of Proposition 1, it is clear that $A(\mathbf{x}, \epsilon)$ is Lipschitz in ϵ for each $\mathbf{x} \in \mathbb{X}^q$, say with constant \tilde{C}_L . Note that this implies that $\mathbb{E}[|A(\mathbf{x}, \epsilon)|] < \infty$ for all \mathbf{x} . We can now appeal to results pertaining to the concentration of Lipschitz functions of Gaussian random variables: the Tsirelson-Ibragimov-Sudakov inequality [11, Theorem 5.5] implies that

$$\log M_{\mathbf{x}}^{A}(t) \leq \frac{t^{2} \tilde{C}_{L}^{2}}{2} + t \, \mathbb{E}[A(\mathbf{x}, \epsilon)]$$

for any $t \in \mathbb{R}$, which is clearly finite for all t since $\mathbb{E}[A(\mathbf{x}, \epsilon)] \leq \mathbb{E}[|A(\mathbf{x}, \epsilon)|]$. From the proof of Proposition 1, we know that $A(\mathbf{x}, \epsilon)$ is $\ell(\epsilon)$ -Lipschitz in \mathbf{x} , where $\ell(\epsilon)$ is itself Lipschitz in ϵ . Hence, the concentration result in Theorem 5.5 of [11] applies again, and we are done.

Proof of Theorem 1. Under the stated assumptions, Lemma 1 ensures that condition (10) in Theorem 3 holds. Further, note that the argument about Lipschitzness of $A(\mathbf{x}, \epsilon)$ in ϵ in the proof of Lemma 1 implies that $\mathbb{E}[|A(\mathbf{x}, \epsilon)|] < \infty$ for all $\mathbf{x} \in \mathbb{X}^q$. Since the $\{\epsilon^i\}_{i=1}^N$ are i.i.d, the strong law of large numbers implies that $\hat{\alpha}_N(\mathbf{x}) \to \alpha(\mathbf{x})$ a.s. for all $x \in \mathbb{X}$. Claims (1) and (2) then follow by applying Theorem 3, and claim (3) follows by applying Proposition 2.

Proof of Theorem 1(q). Mirroring the proof of Theorem 1, we need to show that $\hat{\alpha}_{N_i}(\mathbf{x}) \to \alpha(\mathbf{x})$ a.s. as $i \to \infty$ for all $\mathbf{x} \in \mathbb{X}^q$. For any $\mathbf{x} \in \mathbb{X}^q$ and any $\epsilon_0 \in \mathbb{R}^q$, we have (by convexity and monotonicity of $|x| \mapsto |x|^2$ and the Lipschitz assumption on a and g) that

$$|A(\mathbf{x}, \epsilon)|^2 = |A(\mathbf{x}, \epsilon_0) + A(\mathbf{x}, \epsilon) - A(\mathbf{x}, \epsilon_0)|^2$$

$$\leq |A(\mathbf{x}, \epsilon_0)|^2 + |A(\mathbf{x}, \epsilon) - A(\mathbf{x}, \epsilon_0)|^2$$

$$\leq |A(\mathbf{x}, \epsilon_0)|^2 + L_a^2 L_a^2 ||h_{\mathcal{D}}(\mathbf{x}, \epsilon) - h_{\mathcal{D}}(\mathbf{x}, \epsilon_0)||^2$$

where $h_{\mathcal{D}}(\mathbf{x}, \epsilon) = \mu_{\mathcal{D}}(\mathbf{x}) + L_{\mathcal{D}}(\mathbf{x})\Phi^{-1}(\epsilon)$ with Φ^{-1} the inverse CDF of $\mathcal{N}(0, 1)$, applied elementwise to the vector ϵ of qMC samples. Now choose $\epsilon_0 = (0.5, \dots, 0.5)$, then

$$|A(\mathbf{x}, \epsilon)|^2 \le |a(g(0))|^2 + L_a^2 L_g^2 ||h_{\mathcal{D}}(\mathbf{x}, \epsilon)||^2$$

Since the $\{\epsilon^i\}$ are generated by a nested uniform scramble, we know from Owen [77] that $\epsilon \sim U[0,1]^q$, and therefore $\Phi^{-1}(\epsilon) \sim \mathcal{N}(0,I_q)$. Since affine transformations of Gaussians remain Gaussian, we have that $\mathbb{E}\left[\|h_{\mathcal{D}}(\mathbf{x},\epsilon)\|^2\right] < \infty$. This shows that $A(\mathbf{x},\epsilon) \in L^2([0,1]^q)$. That $\hat{\alpha}_{N_i}(\mathbf{x}) \to 0$ a.s. as $i \to \infty$ for all $\mathbf{x} \in \mathbb{X}^q$ now follows from Owen and Rudolf [79, Theorem 3]. \square

Lemma 2. If f is a GP, then $f_{\mathcal{D}_{\mathbf{x}}}(x') = h(x', \mathbf{x}, \epsilon, \epsilon_I)$, where $\epsilon \sim \mathcal{N}(0, I_q)$ and $\epsilon_I \sim \mathcal{N}(0, 1)$ are independent and h is linear in both ϵ and ϵ_I .

Proof of Lemma 2. This essentially follows from the property of a GP that the covariance conditioned on a new observation (x,y) is independent of y. We can write $f_{\mathcal{D}_{\mathbf{x}}}(x') = \mu_{\mathcal{D}_{\mathbf{x}}}(x') + L_{\mathcal{D}_{\mathbf{x}}}^{\sigma}(x')\epsilon_{I}$, where

$$\mu_{\mathcal{D}_{\mathbf{x}}}(x') := \mu_{\mathcal{D}}(x') + K_{\mathcal{D}}(x', \mathbf{x}) K_{\mathcal{D}}^{\sigma}(\mathbf{x})^{-1} L_{\mathcal{D}}^{\sigma}(\mathbf{x}) \epsilon,$$

 $L_{\mathcal{D}}^{\sigma}(\mathbf{x})$ is the Cholesky decomposition of $K_{\mathcal{D}}^{\sigma}(\mathbf{x}) := K_{\mathcal{D}}(\mathbf{x}, \mathbf{x}) + \operatorname{diag}(\sigma^{2}(\mathbf{x}_{1}), \dots, \sigma^{2}(\mathbf{x}_{q}))$, and $L_{\mathcal{D}_{\mathbf{x}}}^{\sigma}(x')$ is the Cholesky decomposition of

$$K_{\mathcal{D}_{\mathbf{x}}}(x',x') := K(x',x') - K_{\mathcal{D}}(x',\mathbf{x})K_{\mathcal{D}}^{\sigma}(\mathbf{x})^{-1}K_{\mathcal{D}}(\mathbf{x},x').$$

Hence, we see that $f_{\mathcal{D}_{\mathbf{x}}}(x') = h(x', \mathbf{x}, \epsilon, \epsilon_I)$, with

$$h(x', \mathbf{x}, \epsilon, \epsilon_I) = \mu_{\mathcal{D}}(x') + K_{\mathcal{D}}(x', \mathbf{x}) K_{\mathcal{D}}^{\sigma}(\mathbf{x})^{-1} L_{\mathcal{D}}^{\sigma}(\mathbf{x}) \epsilon + L_{\mathcal{D}_{\mathbf{x}}}^{\sigma}(x') \epsilon_I, \tag{12}$$

which completes the argument.

Theorem 6. Let $(a_n)_{n\geq 1}$ be a sequence of non-negative real numbers such that $a_n \to 0$. Suppose that $(i) \mathbb{X}$ is a compact metric space, (ii) f is a GP with continuous sample paths and continuous variance function $x \mapsto \sigma^2(x)$, and (iii) $(\mathbf{x}_n)_{n\geq 1}$ is such that $\alpha_{\mathrm{KG}}^n(\mathbf{x}_n) > \sup_{\mathbf{x} \in \mathbb{X}^q} \alpha_{\mathrm{KG}}^n(\mathbf{x}) - a_n$ infinitely often almost surely. Then $\alpha_{\mathrm{KG}}^n(\mathbf{x}) \to 0$ a.s. for all $\mathbf{x} \in \mathbb{X}^q$.

Proof of Theorem 6. Bect et al. [7] provide a proof for the case q=1. Following their exposition, one finds that the only thing that needs to be verified in order to generalize their results to q>1 is that condition (c) in their Definition 3.18 holds also for the case q>1. What follows is the multi-point analogue of step (f) in the proof of their Theorem 4.8, which establishes this.

Let $\mu: \mathbb{X} \to \mathbb{R}$ and $K: \mathbb{X} \times \mathbb{X} \to \mathbb{R}_+$ denote mean and covariance function of f. Let $Z_{\mathbf{x}} := f(\mathbf{x}) + \operatorname{diag}(\sigma(\mathbf{x}))$, where $\sigma(\mathbf{x}) := (\sigma(\mathbf{x}_1), \dots, \sigma(\mathbf{x}_q))$, with $\epsilon \sim \mathcal{N}(0, I_d)$ independent of f. Moreover, let $x^* \in \arg\max\mu(x)$. Following the same argument as Bect et al. [7], we arrive at the intermediate conclusion that $\mathbb{E}[\max\{0, W_{\mathbf{x},y}\}] = 0$, where $W_{\mathbf{x},y} := \mathbb{E}[f(y) \mid Z_{\mathbf{x}}] - \mathbb{E}[f(x^*) \mid Z_{\mathbf{x}}]$. We need to show that this implies that $\max_{x \in \mathbb{X}} f(x) = m(x^*)$.

Under some abuse of notation we will use μ and K also as the vector / matrix-valued mean / kernel function. Let $K^{\sigma}(\mathbf{x}) := K(\mathbf{x}, \mathbf{x}) + \text{diag}(\sigma(\mathbf{x}))$ and observe that

$$W_{\mathbf{x},y} = \mu(y) - \mu(x^*) + \mathbb{1}_{\{C(\mathbf{x}) \succ 0\}} (K(y, \mathbf{x}) - K(x^*, \mathbf{x})) K^{\sigma}(\mathbf{x})^{-1} (Z_{\mathbf{x}} - \mu(\mathbf{x})),$$

i.e., $W_{\mathbf{x},y}$ is Gaussian with $\mathrm{Var}(W_{\mathbf{x},y}) = V(\mathbf{x},y,x^*)V(\mathbf{x},y,x^*)^T$, where $V(\mathbf{x},y,x^*) := (K(y,\mathbf{x}) - K(x^*,\mathbf{x}))K^\sigma(\mathbf{x})^{-1}$. Since $\mathbb{E}[\max\{0,W_{\mathbf{x},y}\}] = 0$, we must have that $\mathrm{Var}(W_{\mathbf{x},y}) = 0$. If $K^\sigma(\mathbf{x}) \succ 0$, this means that $(K(y,\mathbf{x}) - K(x^*,\mathbf{x})) = 0_q$. But if $K^\sigma(\mathbf{x}) \not\succ 0$, then $K(\mathbf{x},\mathbf{x}) \not\succ 0$, which in turn implies that $K(y,\mathbf{x}) = K(x^*,\mathbf{x}) = 0_q$. This shows that $K(y,\mathbf{x}) = K(x^*,\mathbf{x})$ for all $y \in \mathbb{X}$ and all $\mathbf{x} \in \mathbb{X}$. In particular, $K(x,y) = K(x,x^*)$ for all $y \in \mathbb{X}$. Thus, $K(x,x) - K(x,y) = K(x,x^*) - K(x,x^*)$ for all $y,x \in \mathbb{X}$, and therefore $\mathrm{Var}(f(x) - f(y)) = K(x,x) - K(x,y) - K(y,x) + K(y,y) = 0$. As in [7] we can conclude that this means that the sample paths of $f - \mu$ are constant over \mathbb{X} , and therefore $\max_{x \in \mathbb{X}} f(x) = m(x^*)$.

Proof of Theorem 2. From Lemma 2 we have that $f_{\mathcal{D}_{\mathbf{x}}}(x') = h(x', \mathbf{x}, \epsilon, \epsilon_I)$ with h as in (12). Without loss of generality, we can absorb ϵ_I into ϵ for the purposes of showing that condition (10) holds for the mapping $A_{\mathrm{KG}}(\mathbf{x}, \epsilon) := \max_{x' \in \mathbb{X}} \mathbb{E}\left[g(f(x')) \mid \mathcal{D}_{\mathbf{x}}\right]$. Since the affine (and thus, continuously differentiable) transformation g preserves the necessary continuity and differentiability properties, we can follow the same argument as in the proof of Theorem 1 of [108]. In particular, using continuous differentiability of GP mean and covariance function, compactness of \mathbb{X} , and continuous differentiability of g, we can apply the envelope theorem in the same fashion. From this, it follows that for any $\epsilon \in \mathbb{R}^q$ and for each $1 \leq l \leq q, 1 \leq k \leq d$, the restriction of $\mathbf{x} \mapsto A_{\mathrm{KG}}(\mathbf{x}, \epsilon)$ to the k, l-th coordinate is absolutely continuous for all \mathbf{x} , thus the partial derivative $\partial_{\mathbf{x}_{lk}} A_{\mathrm{KG}}(\mathbf{x}, \epsilon)$ exists a.e. Further, for each l there exist $\Lambda_l \in \mathbb{R}^q$ with $\|\Lambda_l\| < \infty$ s.t. $|\partial_{\mathbf{x}_{kl}} A_{\mathrm{KG}}(\mathbf{x}, \epsilon)| \leq \Lambda_l^T |\varepsilon|$ a.e. on \mathbb{X}^q (here $|\cdot|$ denotes the element-wise absolute value of a vector). This uniform bound on the partial derivatives can be used to show that A_{KG} is $\ell(\epsilon)$ -Lipschitz. Indeed, writing the difference $A_{\mathrm{KG}}(\mathbf{y}, \epsilon) - A_{\mathrm{KG}}(\mathbf{x}, \epsilon)$ as a

⁶In some cases we may consider constructing a heteroskedastic noise model that results in the function $\sigma^2(\mathbf{x})$ changing depending on observations y, in which case this argument does not hold true anymore. We will not consider this case further here.

sum of differences in each of the qd components of x and y, respectively, using the triangle inequality, absolute continuity of the element-wise restrictions, and uniform bound on the partial derivatives, we have that

$$|A_{\mathrm{KG}}(\mathbf{y}, \epsilon) - A_{\mathrm{KG}}(\mathbf{x}, \epsilon)| \leq \sum_{k=1}^{q} \sum_{l=1}^{d} \Lambda_{l}^{T} |\epsilon| |\mathbf{y}_{kl} - \mathbf{x}_{kl}| \leq \max_{1 \leq l \leq d} \left\{ \Lambda_{l}^{T} |\epsilon| \right\} ||\mathbf{y} - \mathbf{x}||_{1}$$

and so $\ell(\epsilon) = \max_{l} \{\Lambda_{l}^{T} | \epsilon | \}$. Going back to viewing ϵ as a random variable, it is straightforward to verify that $\ell(\epsilon)$ is integrable. Indeed,

$$\textstyle \mathbb{E}[|\ell(\epsilon)|] \leq \max_l \left\{ \sum_{k=1}^q \Lambda_{lk} \mathbb{E}[|\epsilon_k|] \right\} = \sqrt{2/\pi} \max_l \left\{ \|\Lambda_l\|_1 \right\}.$$

Since g is assumed to be affine in (iii), we can apply Lemma 2 to see that $\mathbb{E}[g(f(x')) \mid \mathcal{D}_{\mathbf{x}}]$ is a GP. Therefore, $A_{\mathrm{KG}}(\mathbf{x}, \epsilon)$ represents the maximum of a GP and its moment generating function $\mathbb{E}[e^{tA_{\mathrm{KG}}(\mathbf{x}, \epsilon)}]$ is finite for all t by Lemma 4. This implies finiteness of its absolute moments [65, Exercise 9.15] and we have that $\mathbb{E}[|A_{\mathrm{KG}}(\mathbf{x}, \epsilon)|] < \infty$ for all $\mathbf{x} \in \mathbb{X}$. Since the $\{\epsilon^i\}$ are i.i.d, the strong law of large numbers ensures that $\hat{\alpha}_{\mathrm{KG},N}(\mathbf{x}) \to \alpha_{\mathrm{KG}}(\mathbf{x})$ a.s. Theorem 3 now applies to obtain (1) and (2).

Moreover, by the analysis above, it holds that

$$\ell(\epsilon) = \max_{l} \{\Lambda_{l}^{T} | \epsilon| \} \le q \max_{l} \|\Lambda_{l}^{T}\|_{\infty} \|\epsilon\|_{\infty} =: \ell'(\epsilon),$$

so $\ell'(\epsilon)$ is also a Lipschitz constant for $A_{\mathrm{KG}}(\cdot,\epsilon)$. Here, the absolute value version (the second result) of Lemma 4 applies, so we have that $\mathbb{E}[e^{t\ell'(\epsilon)}]$ is finite for all t. The conditions of Proposition 2 are now satisfied and we have the desired conclusion.

Proof of Theorem 2(q). In the RQMC setting, we have by Owen [77] that $\epsilon \sim U[0,1]^q$. Therefore, we are now interested in examining $\tilde{A}_{KG}(\mathbf{x},\epsilon) := A_{KG}(\mathbf{x},\Phi^{-1}(\epsilon))$, since $\Phi^{-1}(\epsilon) \sim \mathcal{N}(0,I_q)$. Following the same analysis as in the proof of Theorem 2, we have Lipschitzness of $\tilde{A}_{KG}(\cdot,\epsilon)$:

$$|\tilde{A}_{KG}(\mathbf{y}, \epsilon) - \tilde{A}_{KG}(\mathbf{x}, \epsilon)| \le \ell(\Phi^{-1}(\epsilon)) \|\mathbf{y} - \mathbf{x}\|_1,$$

where $\ell(\cdot)$ is as defined in the proof of Theorem 2. As before, $\ell(\Phi^{-1}(\epsilon))$ is integrable. Like in the proof of Theorem 2, $\tilde{A}_{\mathrm{KG}}(\mathbf{x},\epsilon)$ is the maximum of a GP and its moment generating function $\mathbb{E}[e^{t\tilde{A}_{\mathrm{KG}}(\mathbf{x},\epsilon)}]$ is finite for all t by Lemma 4, implying finiteness of its second moment: $\mathbb{E}[\tilde{A}_{\mathrm{KG}}(\mathbf{x},\epsilon)^2] < \infty$ for all $\mathbf{x} \in \mathbb{X}$. Thus, that $\tilde{A}_{\mathrm{KG}}(\mathbf{x},\epsilon) \in L^2([0,1]^q)$ and $\hat{\alpha}_{\mathrm{KG},N_i}(\mathbf{x}) \to \alpha_{\mathrm{KG}}(\mathbf{x})$ a.s. as $i \to \infty$ for all $\mathbf{x} \in \mathbb{X}^q$ follows from Owen and Rudolf [79, Theorem 3]. Theorem 3 now allows us to conclude (1) and (2).

The following Lemma will be used to prove Theorem 4:

Lemma 3. Consider a Gaussian Process f on $\mathbb{X} \subset \mathbb{R}^d$ with covariance function $K(\cdot, \cdot) : \mathbb{X} \times \mathbb{X} \to \mathbb{R}$. Suppose that (i) \mathbb{X} is compact, and (ii) K is continuously differentiable. Then f has continuous sample paths.

Proof of Lemma 3. Since K is continuously differentiable and $\mathbb X$ is compact, K is Lipschitz on $\mathbb X \times \mathbb X$, i.e., $\exists \ L < \infty$ such that $|K(x,y) - K(x',y')| \le L \big(\|x - x'\| + \|y - y'\| \big)$ for all $(x,y), (x',y') \in \mathbb X \times \mathbb X$. Thus

$$\mathbb{E}|f(x) - f(x)|^2 = K(x, x) - 2K(x, y) + K(y, y)$$

$$\leq |K(x, x) - K(x, y)| + |K(y, y) - K(x, y)|$$

$$< 2L||x - y||$$

Since $\mathbb X$ is compact, there exists $C:=\max_{x,y\in\mathbb X}\|x-y\|<\infty$. With this it is easy to verify that there exist $C'<\infty$ and $\eta>0$ such that $2L\|x-y\|< C'|\log\|x-y\||^{-(1+\eta)}$ for all $x,y\in\mathbb X$. Continuity of the sample paths then follows from Theorem 3.4.1 in [3].

Proof of Theorem 4. From Lemma 3 we know that the GP has continuous sample paths. If $\mathbf{x}_{n+1} \in \arg\max_{\mathbf{x} \in \mathbb{X}^q} \hat{\alpha}^n_{\mathrm{KG},N_n}(\mathbf{x})$ for all n, the almost sure convergence of $\hat{\mathbf{x}}^n_{\mathrm{KG},N_n}$ to the set of optimizers of α^n_{KG} from Theorem 2 together with continuity of α^n_{KG} (established in the proof of Theorem 2) implies that for all $\delta > 0$ and each $n \geq 1, \exists N_n < \infty$ such that $\alpha^n_{\mathrm{KG}}(\mathbf{x}_{n+1}) > \sup_{\mathbf{x} \in \mathbb{X}^q} \alpha^n_{\mathrm{KG}}(\mathbf{x}) - \delta$. As $\limsup_n N_n = \infty, \exists (a_n)_{n \geq 1}$ with $a_n \to 0$ such that $\alpha^n_{\mathrm{KG}}(\mathbf{x}_{n+1}) > \sup_{\mathbf{x} \in \mathbb{X}^q} \alpha^n_{\mathrm{KG}}(\mathbf{x}) - a_n$ infinitely often. That $\alpha^n_{\mathrm{KG}}(\mathbf{x}) \to 0$ a.s. for all $\mathbf{x} \in \mathbb{X}^q$ then follows from Theorem 6. The convergence result for $f(\chi_n)$ then follows directly from Proposition 4.9 in [7].

Lemma 4. Let f be a mean zero GP defined on \mathbb{X} such that $|f(x)| < \infty$ almost surely for each $x \in \mathbb{X}$. It holds that the moment generating functions of $\sup_{x \in \mathbb{X}} f(x)$ and $\sup_{x \in \mathbb{X}} |f(x)|$ are both finite, i.e.,

 $\mathbb{E}\big[e^{t\sup_{x\in\mathbb{X}}f(x)}\big]<\infty\quad \text{and}\quad \mathbb{E}\big[e^{t\sup_{x\in\mathbb{X}}|f(x)|}\big]<\infty$

for any $t \in \mathbb{R}$.

Proof of Lemma 4. Let $||f|| := \sup_{x \in \mathbb{X}} f$. Since the sample paths of f are almost surely finite, the Borell-TIS inequality [2, Theorem 2.1] states that $\mathbb{E} ||f|| < \infty$. We first consider t > 0 and begin by re-writing the expectation as

$$\mathbb{E}\left[e^{t\,\|f\|}\right] = \int_{0}^{\infty} \mathbb{P}\left(e^{t\,\|f\|} > u\right) du
\leq 1 + \int_{1}^{\infty} \mathbb{P}\left(e^{t\,\|f\|} > u\right) du
= 1 + \int_{1}^{\infty} \mathbb{P}\left(\|f\| - \mathbb{E}\|f\| > t^{-1}\log u - \mathbb{E}\|f\|\right) du
= 1 + te^{t\,\mathbb{E}\|f\|} \int_{-\mathbb{E}\|f\|}^{\infty} \mathbb{P}\left(\|f\| - \mathbb{E}\|f\| > u\right) e^{tu} du
\leq 1 + te^{t\,\mathbb{E}\|f\|} \left[\int_{\min\{-\mathbb{E}\|f\|, 0\}}^{0} + \int_{0}^{\infty}\right] \mathbb{P}\left(\|f\| - \mathbb{E}\|f\| > u\right) e^{tu} du
\leq 1 + \left|\mathbb{E}\|f\|\right| te^{t\,\mathbb{E}\|f\|} + te^{t\,\mathbb{E}\|f\|} \int_{0}^{\infty} \mathbb{P}\left(\|f\| - \mathbb{E}\|f\| > u\right) e^{tu} du, \tag{13}$$

where a change of variables is performed in the third equality. Let $\sigma_{\mathbb{X}}^2 := \sup_{x \in \mathbb{X}} \mathbb{E}[f(x)^2]$. We can now use the Borell-TIS inequality to bound the tail probability in (13) by $2e^{-u^2/(2\sigma_{\mathbb{X}}^2)}$, obtaining:

$$\mathbb{E}\left[e^{t\,\|f\|}\right] \leq 1 + \left|\mathbb{E}\,\|f\|\right| te^{t\,\mathbb{E}\|f\|} + te^{t\,\mathbb{E}\|f\|} \int_0^\infty 2e^{-u^2/(2\sigma_{\mathbb{X}}^2) + tu} \, du < \infty.$$

Similarly, for t < 0, we have:

$$\mathbb{E}\left[e^{t \, |||f|||}\right] = \int_{0}^{\infty} \mathbb{P}\left(e^{t \, |||f|||} > u\right) du
\leq 1 + \int_{1}^{\infty} \mathbb{P}\left(e^{t \, |||f|||} > u\right) du
= 1 + \int_{1}^{\infty} \mathbb{P}\left(|||f||| - \mathbb{E} \, ||f|| < t^{-1} \log u - \mathbb{E} \, ||f||\right) du
= 1 - te^{t \, \mathbb{E} \, ||f||} \int_{-\infty}^{-\mathbb{E} \, ||f||} \mathbb{P}\left(|||f||| - \mathbb{E} \, ||f|| < u\right) e^{tu} du
\leq 1 - te^{t \, \mathbb{E} \, ||f||} \left[\int_{0}^{\max\{-\mathbb{E} \, ||f||, 0\}} + \int_{-\infty}^{0} \mathbb{P}\left(|||f||| - \mathbb{E} \, ||f|| < u\right) e^{tu} du
\leq 1 - \left|\mathbb{E} \, ||f|| te^{t \, \mathbb{E} \, ||f||} - te^{t \, \mathbb{E} \, ||f||} \int_{-\infty}^{0} \mathbb{P}\left(|||f||| - \mathbb{E} \, ||f|| < u\right) e^{tu} du, \tag{14}$$

The same can be done for (14) to conclude that $\mathbb{E}\left[e^{t\,\|f\|}\right]<\infty$ for all t. For the case of $\mathbb{E}\left[e^{t\,\||f\|\|}\right]$ and t>0, we use a similar line of analysis as (13) along with the observation that

$$\mathbb{P}(\||f|\| - \mathbb{E}\|f\| > u) \le 2 \,\mathbb{P}(\|f\| - \mathbb{E}\|f\| > u).$$

Proof of Theorem 5. Since we are in the case of finite \mathbb{X} , let μ_n and Σ_n denote the posterior mean vector and covariance matrix of our GP after conditioning on \mathcal{D}_n . First, we give a brief outline of the argument. We know from previous work (Lemma A.6 of [25] or Lemma 3 of [83]) that given a posterior distribution parameterized by μ and Σ , if $\alpha_{KG}(x;\mu,\Sigma)=0$ for all $x\in\mathbb{X}$, then an optimal design is identified:

$$\mathop{\arg\max}_{x\in\mathbb{X}}\mu(x)=\mathop{\arg\max}_{x\in\mathbb{X}}f(x)$$

almost surely. Thus, we can use the true KG values as a "potential function" to quantify how the OKG policy performs asymptotically, even though we are never using the KG acquisition function for selecting points. We emphasize that the data that induce $\{\mu_n\}_{n\geq 0}$ and $\{\Sigma_n\}_{n\geq 0}$ are collected using the OKG policy.

By a martingale convergence argument, there exists a limiting posterior distribution described by random variables $(\mu_{\infty}, \Sigma_{\infty})$, i.e., $\mu_n \to \mu_{\infty}$ and $\Sigma_n \to \Sigma_{\infty}$ almost surely [25, Lemma A.5]. Let $A \subseteq \mathbb{X}$ be a subset of the feasible space. As was done in the proof of Theorem 4 of [25], we define the event:

$$H_A = \left\{ \alpha_{KG}(x; \mu_{\infty}, \Sigma_{\infty}) > 0, \ x \in A \right\} \cap \left\{ \alpha_{KG}(x; \mu_{\infty}, \Sigma_{\infty}) = 0, \ x \notin A \right\}. \tag{15}$$

Note that H_A , for all possible subsets A, partition the sample space. Consider some $A \neq \emptyset$. By Lemma A.7 of [25], if $\alpha_{\rm KG}(x;\mu_\infty,\Sigma_\infty)>0$, then x is measured a finite number of times, meaning that there exists an almost surely finite random variable M_0 such that on iterations after N_0 , OKG stops sampling from A. By the definition of H_A in (15), there must exist another random iteration index $M_1\geq M_0$ such that when $n\geq N_1$,

$$\min_{x \in \mathcal{A}} \alpha_{KG}(x; \mu_n, \Sigma_n) > \max_{x \notin \mathcal{A}} \alpha_{KG}(x; \mu_n, \Sigma_n),$$

implying that the exact KG policy must prefer points in A over all others after iteration M_1 . This implies that

$$H_A \subseteq \left\{ \operatorname*{arg\,max}_{x \in \mathbb{X}} \hat{\alpha}_{\mathrm{KG}, N_n}(x, \mu_n, \Sigma_n) \not\subseteq \operatorname*{arg\,max}_{x \in \mathbb{X}} \alpha_{\mathrm{KG}}(x, \mu_n, \Sigma_n), \, \forall \, n \geq M_1 - 1 \right\} =: E,$$

because if not, then there exists an iteration after M_0 where an element from A is selected, which is a contradiction. As shown in the proof of Lemma 2, the next period posterior mean $\mathbb{E}\big[g(f(x')) \mid \mathcal{D}_{\mathbf{x}}\big]$ is a GP. Therefore, by Lemma 4, the moment generating function of $\max_{x' \in \mathbb{X}} \mathbb{E}\big[g(f(x')) \mid \mathcal{D}_{\mathbf{x}}\big]$ is finite. Theorem 2.6 of [39] establishes that our choice of N_n guarantees

$$\mathbb{P}\Big[\arg\max_{x\in\mathbb{X}}\hat{\alpha}_{\mathrm{KG},N_n}(x,\mu_n,\Sigma_n)\not\subseteq\arg\max_{x\in\mathbb{X}}\alpha_{\mathrm{KG}}(x,\mu_n,\Sigma_n)\,|\,\mathcal{F}_n\Big]\leq\delta,$$

from which it follows that

$$\sum_{n=0}^{\infty} \log \mathbb{P} \Big[\argmax_{x \in \mathbb{X}} \hat{\alpha}_{\mathrm{KG}, N_n}(x, \mu_n, \Sigma_n) \not\subseteq \argmax_{x \in \mathbb{X}} \alpha_{\mathrm{KG}}(x, \mu_n, \Sigma_n) \, | \, \mathcal{F}_n \Big] = -\infty.$$

After writing the probability of E as an infinite product and performing some manipulation, we see that the above condition implies that the probability of event E is zero, and we conclude that $\mathbb{P}(H_A)=0$ for any nonempty A. Therefore, $\mathbb{P}(H_\emptyset)=1$ and $\alpha_{\mathrm{KG}}(x;\mu_\infty,\Sigma_\infty)=0$ for all x almost surely.

E Illustration of Sample Average Approximation

QMC methods have been used in other applications in machine learning, including variational inference [13] and evolutionary strategies [86], but rarely in BO. Letham et al. [58] use QMC in the context of a specific acquisition function. BOTORCH's abstractions make it straightforward (and mostly automatic) to use QMC integration with any acquisition function.

Using SAA, i.e., fixing the base samples $E = \{\epsilon^i\}$, introduces a consistent bias in the function approximation. While i.i.d. re-sampling in each evaluation ensures that $\hat{\alpha}_N(\mathbf{x}, \Phi, \mathcal{D})$ and $\hat{\alpha}_N(\mathbf{y}, \Phi, \mathcal{D})$ are conditionally independent given (\mathbf{x}, \mathbf{y}) , this no longer holds when fixing the base samples.

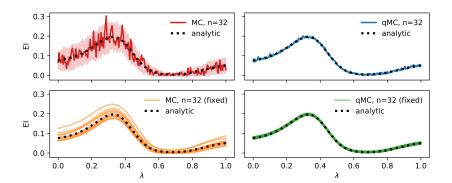


Figure 18: MC and QMC acquisition functions, with and without re-drawing the base samples between evaluations. The model is a GP fit on 15 points randomly sampled from $\mathbb{X} = [0,1]^6$ and evaluated on the (negative) Hartmann6 test function. The acquisition functions are evaluated along the slice $x(\lambda) = \lambda \mathbf{1}$.

Figure 18 illustrates this behavior for EI (we consider the simple case of q=1 for which we have an analytic ground truth available). The top row shows the MC and QMC version, respectively, when re-drawing base samples for every evaluation. The solid lines correspond to a single realization, and the shaded region covers four standard deviations around the mean, estimated across 50 evaluations. It is evident that QMC sampling significantly reduces the variance of the estimate. The bottom row shows the same functions for 10 different realizations of fixed base samples. Each of these realizations is differentiable w.r.t. x (and hence λ in the slice parameterization). In expectation (over the base samples), this function coincides with the true function (the dashed black line). Conditional on the base sample draw, however, the estimate displays a consistent bias. The variance of this bias (across re-drawing the base samples) is much smaller for the QMC versions.

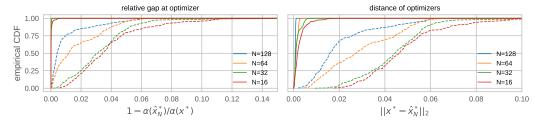


Figure 19: Performance for optimizing QMC-based EI. Solid lines: fixed base samples, optimized via L-BFGS-B. Dashed lines: re-sampling base samples, optimized via Adam (lr=0.025).

Even thought the function values may show noticeable bias, the bias of the maximizer (in $\mathbb X$) is typically very small. Figure 19 illustrates this behavior, showing empirical cdfs of the relative gap $1-\alpha(\hat x_N^*)/\alpha(x^*)$ and the distance $\|x^*-\hat x_N^*\|_2$ over 250 optimization runs for different numbers of samples, where x^* is the optimizer of the analytic function EI, and $\hat x_N^*$ is the optimizer of the QMC approximation. The quality of the solution of the deterministic problem is excellent even for relatively small sample sizes, and generally better than of the stochastic optimizer.

Figure 20 shows empirical mean and variance of the metrics from Figure 19 as a function of the number of MC samples N on a log-log scale. The stochastic optimizer used is Adam with a learning rate of 0.025. Both for the SAA and the stochastic version we use the same number of random restart initial conditions generated from the same initialization heuristic.

Empirical asymptotic convergence rates can be obtained as the slopes of the OLS fit (dashed lines), and are given in Table 1. It is quite remarkable that in order to achieve the same error as the MC approximation with 4096 samples, the QMC approximation only requires 64 samples. This holds true for the bias and variance of the (relativized) optimal value as well as for the distance from the true optimizer. That said, as we are in a BO setting, we are not necessarily interested in the estimation error $\hat{\alpha}_N^* - \alpha^*$ of the optimum, but primarily in how far x_N^* is from the true optimizer x^* .

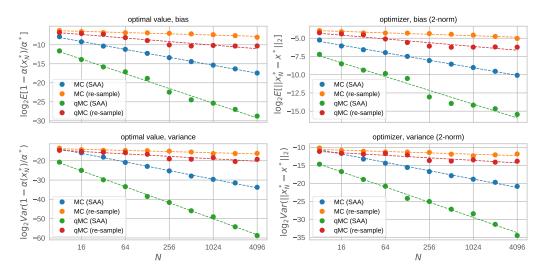


Figure 20: Bias and variance of optimizer x_N^* and true EI value $\mathrm{EI}(x_N^*)$ evaluated at the optimizer as a function of the number of (Q)MC samples for both SAA and stochastic optimization ("re-sample").

	MC	QMC	MC^{\dagger}	QMC^{\dagger}
$\mathbb{E}[1-\hat{\alpha}_N^*/\alpha^*]$	-0.52	-0.95	-0.10	-0.26
$Var(1-\hat{\alpha}_N^*/\alpha^*)$	-1.16	-2.11	-0.19	-0.35
$\mathbb{E}[\ x_N^* - x^*\ _2]$	-1.04	-1.94	-0.16	-0.47
$Var(x_N^* - x^* _2)$	-2.24	-4.14	-0.30	-0.63

Table 1: Empirical asymptotic convergence rates for the setting in Figure 20 († denotes re-sampling + optimization with Adam).

A somewhat subtle point is that whether better optimization of the acquisition function results in improved closed-loop BO performance depends on the acquisition function as well as the underlying problem. More exploitative acquisition functions, such as EI, tend to show worse performance for problems with high noise levels. In these settings, not solving the EI maximization exactly adds randomness and thus induces additional exploration, which can improve closed-loop performance. While a general discussion of this point is outside the scope of this paper, BOTORCH does provide a framework for optimizing acquisition functions well, so that these questions can be compartmentalized and acquisition function performance can be investigated independently from the quality of optimization.

Perhaps the most significant advantage of using deterministic optimization algorithms is that, unlike for algorithms such as SGD that require tuning the learning rate, the optimization procedure is essentially hyperparameter-free. Figure 8 shows the closed-loop optimization performance of qEI for both deterministic and stochastic optimization for different optimizers and learning rates. While some of the stochastic variants (e.g. ADAM with learning rate 0.01) achieve performance similar to the deterministic optimization, the type of optimizer and learning rate matters. In fact, the rank order of SGD and ADAM w.r.t. to the learning rate is reversed, illustrating that selecting the right hyperparameters for the optimizer is itself a non-trivial problem.

F Additional Implementation Details

F.1 Batch Initialization for Multi-Start Optimization

For most acquisition functions, the optimization surface is highly non-convex, multi-modal, and (especially for "improvement-based" ones such as EI or KG) often flat (i.e. has zero gradient) in much of the domain \mathbb{X} . Therefore, optimizing the acquisition function is itself a challenging problem.

The simplest approach is to use zeroth-order optimizers that do not require gradient information, such as DIRECT or CMA-ES [45, 37]. These approaches are feasible for lower-dimensional problems, but do not scale to higher dimensions. Note that performing parallel optimization over q candidates in a d-dimensional feature space means solving a qd-dimensional optimization problem.

A more scalable approach incorporates gradient information into the optimization. As described in Section 4, BOTORCH by default uses quasi-second order methods, such as L-BFGS-B. Because of the complex structure of the objective, the initial conditions for the algorithm are extremely important so as to avoid getting stuck in a potentially highly sub-optimal local optimum. To reduce this risk, one typically employs multi-start optimization (i.e. start the solver from multiple initial conditions and pick the best of the final solutions). To generate a good set of initial conditions, BOTORCH heavily exploits the fast batch evaluation discussed in the previous section. Specifically, BOTORCH by default uses $N_{\rm opt}$ initialization candidates generated using the following heuristic:

- 1. Sample \tilde{N}_0 quasi-random q-tuples of points $\tilde{\mathbf{x}}_0 \in \mathbb{R}^{\tilde{N}_0 \times q \times d}$ from \mathbb{X}^q using quasi-random Sobol sequences.
- 2. Batch-evaluate the acquisition function at these candidate sets: $\tilde{v} = \alpha(\tilde{\mathbf{x}}_0; \Phi, \mathcal{D})$.
- 3. Sample N_0 candidate sets $\mathbf{x} \in \mathbb{R}^{N_0 \times q \times d}$ according to the weight vector $p \propto \exp(\eta v)$, where $v = (\tilde{v} \hat{\mu}(\tilde{v}))/\hat{\sigma}(\tilde{v})$ with $\hat{\mu}$ and $\hat{\sigma}$ the empirical mean and standard deviation, respectively, and $\eta > 0$ is a temperature parameter. Acquisition functions that are known to be flat in large parts of \mathbb{X}^q are handled with additional care in order to avoid starting in locations with zero gradients.

Sampling initial conditions this way achieves an exploration/exploitation trade-off controlled by the magnitude of η . As $\eta \to 0$ we perform Sobol sampling, while $\eta \to \infty$ means the initialization is chosen in a purely greedy fashion. The latter is generally not advisable, since for large \tilde{N}_0 the highest-valued points are likely to all be clustered together, which would run counter to the goal of multi-start optimization. Fast batch evaluation allows evaluating a large number of samples (\tilde{N}_0 in the tens of thousands is feasible even for moderately sized models).

F.2 Sequential Greedy Batch Optimization

The pending points approach discussed in Section 5 provides a natural way of generating parallel BO candidates using *sequential greedy* optimization, where candidates are chosen sequentially, while in each step conditioning on selected points and integrating over the uncertainty in their outcome (using MC integration). By using a full MC formulation, in which we jointly sample at new and pending points, we avoid constructing an individual "fantasy" model for each sampled outcome, a common (and costly) approach in the literature [94]. In practice, the sequential greedy approach often performs well, and may even outperform the joint optimization approach, since it involves a sequence of small, simpler optimization problems, rather than a larger and complex one that is harder to solve.

[105] provide a theoretical justification for why the sequential greedy approach works well with a class of acquisition functions that are submodular.

G Active Learning Example

Recall from Section 5 the negative integrated posterior variance (NIPV) [91, 17] of the model:

$$NIPV(\mathbf{x}) = -\int_{\mathbb{X}} \mathbb{E}\left[Var(f(x) \mid \mathcal{D}_{\mathbf{x}}) \mid \mathcal{D}\right] dx.$$
 (16)

We can implement (16) using standard Botorch components, as shown in Code Example 4. Here mc_points is the set of points used for MC-approximating the integral. In the most basic case, one can use QMC samples drawn uniformly in X. By allowing for arbitrary mc_points, we permit weighting regions of X using non-uniform sampling. Using mc_points as samples of the maximizer of the posterior, we recover the recently proposed Posterior Variance Reduction Search [74] for BO.

This acquisition function supports both parallel selection of points and asynchronous evaluation. Since MC integration requires evaluating the posterior variance at a large number of points, this acquisition function benefits significantly from the fast predictive variance computations in GPyTorch [82, 29].

```
\verb|class| q Negative Integrated Posterior Variance (Analytic Acquisition Function): \\
3
      @concatenate_pending_points
4
      @t_batch_mode_transform()
      def forward(self, X: Tensor) -> Tensor:
6
           fant_model = self.model.fantasize(
               X=X, sampler=self._dummy_sampler,
               observation_noise=True
9
10
           sz = [1] * len(X.shape[:-2]) + [-1, X.size(-1)]
          mc_points = self.mc_points.view(*sz)
11
12
           with settings.propagate_grads(True):
               posterior = fant_model.posterior(mc_points)
               = posterior.variance.mean(dim=-2)
           return -ivar.view(X.shape[:-2])
```

Code Example 4: Active Learning (NIPV)

To illustrate how NIPV may be used in combination with scalable probabilistic modeling, we examine the problem of efficient allocation of surveys across a geographic region. Inspired by Cutajar et al. [18], we utilize publicly-available data from The Malaria Atlas Project (2019) dataset, which includes the yearly mean parasite rate (along with standard errors) of *Plasmodium falciparum* at a 4.5km² grid spatial resolution across Africa. In particular, we consider the following active learning problem: given a spatio-temporal probabilistic model fit to data from 2011-2016, which geographic locations in and around Nigeria should one sample in 2017 in order to minimize the model's error for 2017 across all of Nigeria?

We fit a heteroskedastic GP model to 2500 training points prior to 2017 (using a noise model that is itself a GP fit to the provided standard errors). We then select q=10 sample locations for 2017 using the NIPV acquisition function, and make predictions across the entirety of Nigeria using this new data. Compared to using no 2017 data, we find that our new dataset reduces MSE by 16.7% on average (SEM = 0.96%) across 60 subsampled datasets. By contrast, sampling the new 2017 points at a regularly spaced grid results only in a 12.4% reduction in MSE (SEM = 0.99%). The mean relative improvement in MSE reduction from NIPV optimization is 21.8% (SEM = 6.64%). Figure 21 shows the NIPV-selected locations on top of the base model's estimated parasite rate and standard deviation.

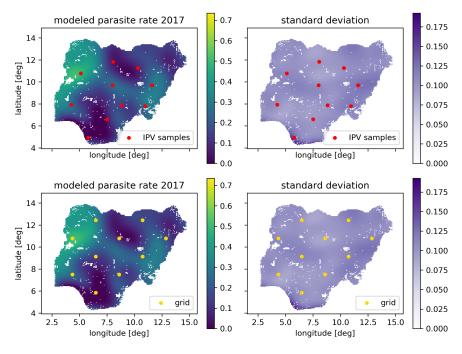


Figure 21: Locations for 2017 samples from IPV minimization and the base grid. Observe how the NIPV samples cluster in higher variance areas.

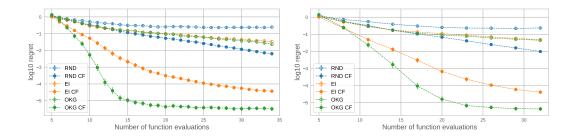


Figure 22: Composite function optimization for q=1 Figure 23: Composite function optimization for q=3

H Additional Implementation Examples

Comparing Implementation Complexity

Many of BoTorch's benefits are qualitative, including the simplification and acceleration of implementing new acquisition functions. Quantifying this in a meaningful way is very challenging. Comparisons are often made in terms of Lines of Code (LoC) - while this metric is problematic when comparing across different design philosophies, non-congruent feature sets, or even programming languages, it does provides a general idea of the effort required for developing and implementing new methods.

MOE's KG involves thousands of LoC in C++ and python spread across a large number of files, while our more efficient implementation is <30 LoC. Astudillo and Frazier [6] is a full paper in last year's installment of this conference, whose composite function method we implement and significantly extend (e.g to support KG) in 7 LoC using BoTorch's abstractions. The original NEI implementation is >250 LoC, while the one from Code Example 2 is 14 LoC.

H.1 Composite Objectives

We consider the Bayesian model calibration of a simulator with multiple outputs from Section 5.3 of Astudillo and Frazier [6]. In this case, the simulator from Bliznyuk et al. [9] models the concentrations of chemicals at 12 positions in a one-dimensional channel. Instead of modeling the overall loss function (which measures the deviation of the simulator outputs with a set of observations) directly, we follow Astudillo and Frazier [6] and model the underlying concentrations while utilizing a composite objective approach. A powerful aspect of BoTorch's modular design is the ability to easily combine different approaches into one. For the composite function problem in this section this means that we can easily extend the work by Astudillo and Frazier [6] not only to use the Knowledge Gradient, but also to the "parallel BO" setting of jointly selecting q>1 points. Figures 22 and 22 show results for this with q=1 and q=3, repspectively. The plots show log regret evaluated at the maximizer of the posterior mean averaged over 250 trials. While the performance of EI-CF is similar for q=1 and q=3, KG-CF reaches lower regret significantly faster for q=1 compared to q=3, suggesting that "looking ahead" is beneficial in this context.

H.2 Generalized UCB

Code Example 5 presents a generalized version of parallel UCB from Wilson et al. [104] supporting pending candidates, generic objectives, and QMC sampling. If no sampler is specified, a default QMC sampler is used. Similarly, if no objective is specified, the identity objective is assumed.

H.3 Full Code Examples

In this section we provide full implementations for the code examples. Specifically, we include parallel Noisy EI (Code Example 6), OKG (Code Example 7), and (negative) Integrated Posterior Variance (Code Example 8).

⁷https://github.com/wujian16/Cornell-MOE

 $^{^8\}mathrm{Code}$ available at https://github.com/RaulAstudillo06/BOCF

```
1 class qUpperConfidenceBound(MCAcquisitionFunction):
3
      def __init__(
4
          self,
5
          model: Model,
          beta: float,
6
7
          sampler: Optional[MCSampler] = None,
8
          objective: Optional[MCAcquisitionObjective] = None,
9
          X_pending: Optional[Tensor] = None,
10
      ) -> None:
11
          super().__init__(model, sampler, objective, X_pending)
12
          self.beta_prime = math.sqrt(beta * math.pi / 2)
13
14
      @concatenate_pending_points
15
      @t_batch_mode_transform()
16
      def forward(self, X: Tensor) -> Tensor:
          posterior = self.model.posterior(X)
17
18
          samples = self.sampler(posterior)
19
          obj = self.objective(samples)
20
          mean = obj.mean(dim=0)
21
          z = mean + self.beta_prime * (obj - mean).abs()
22
          return z.max(dim=-1).values.mean(dim=0)
```

Code Example 5: Generalized Parallel UCB

```
1 class qNoisyExpectedImprovement(MCAcquisitionFunction):
3
      def __init__(
4
          self,
5
          model: Model,
          X_baseline: Tensor,
6
7
          sampler: Optional[MCSampler] = None,
8
          objective: Optional[MCAcquisitionObjective] = None,
9
          X_pending: Optional[Tensor] = None,
10
      ) -> None:
11
          super().__init__(model, sampler, objective, X_pending)
12
          self.register_buffer("X_baseline", X_baseline)
13
14
      @concatenate_pending_points
15
      @t_batch_mode_transform()
16
      def forward(self, X: Tensor) -> Tensor:
17
          q = X.shape[-2]
18
          X_bl = match_shape(self.X_baseline, X)
19
          X_full = torch.cat([X, X_bl], dim=-2)
20
          posterior = self.model.posterior(X_full)
21
          samples = self.sampler(posterior)
22
          obj = self.objective(samples)
          obj_n = obj[...,:q].max(dim=-1).values
23
24
          obj_p = obj[...,q:].max(dim=-1).values
25
          return (obj_n - obj_p).clamp_min(0).mean(dim=0)
```

Code Example 6: Parallel Noisy EI (full)

```
1 class qKnowledgeGradient(MCAcquisitionFunction):
      def __init__(
3
4
           self,
5
           model: Model,
           sampler: MCSampler,
6
7
           objective: Optional[Objective] = None,
8
           inner_sampler: Optional[MCSampler] = None,
9
           X_pending: Optional[Tensor] = None,
10
      ) -> None:
11
           super().__init__(model, sampler, objective, X_pending)
           self.inner_sampler = inner_sampler
12
13
14
      def forward(self, X: Tensor) -> Tensor:
           splits = [X.size(-2) - self.Nf, self.N_f]
X, X_fantasies = torch.split(X, splits, dim=-2)
15
16
           # [...] some re-shaping for batch evaluation purposes
17
18
           if self.X_pending is not None:
19
               X_p = match_shape(self.X_pending, X)
20
               X = torch.cat([X, X_p], dim=-2)
21
           fmodel = self.model.fantasize(
22
               X = X,
23
               sampler=self.sampler,
24
               observation_noise=True,
25
           )
26
           obj = self.objective
27
           if isinstance(obj, MCAcquisitionObjective):
28
               inner_acqf = SimpleRegret(
29
                    fmodel, sample=self.inner_sampler, objective=obj,
30
               )
31
           else:
32
               inner_acqf = PosteriorMean(fmodel, objective=obj)
33
           with settings.propagate_grads(True):
34
               values = inner_acqf(X_fantasies)
35
           return values.mean(dim=0)
```

Code Example 7: One-Shot Knowledge Gradient (full)

```
1 \; {\tt class} \; \; {\tt qNegIntegratedPosteriorVariance(AnalyticAcquisitionFunction)}:
      def __init__(
3
           self,
4
          model: Model,
5
          mc_points: Tensor,
          X_pending: Optional[Tensor] = None,
      ) -> None:
7
8
          super().__init__(model=model)
9
           self._dummy_sampler = IIDNormalSampler(1)
10
           self.X_pending = X_pending
           self.register_buffer("mc_points", mc_points)
11
12
13
      @concatenate_pending_points
      @t_batch_mode_transform()
14
15
      def forward(self, X: Tensor) -> Tensor:
           fant_model = self.model.fantasize(
16
17
               X = X,
18
               sampler=self._dummy_sampler,
19
               observation_noise=True,
          )
20
21
          batch_ones = [1] * len(X.shape[:-2])
22
          mc_points = self.mc_points.view(*batch_ones, -1, X.size(-1))
23
           with settings.propagate_grads(True):
               posterior = fant_model.posterior(mc_points)
24
25
          ivar = posterior.variance.mean(dim=-2)
          return -ivar.view(X.shape[:-2])
26
```

Code Example 8: Active Learning (full)