

---

# Practical Bayesian Optimization of Objectives with Conditioning Variables

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Bayesian optimization is a class of data efficient model based algorithms typically  
2 focused on global optimization. We consider the more general case where a user is  
3 faced with multiple tasks that each need to be optimized, find the global optima  
4 within each task, all the task conditional optima. For example given a range of  
5 cities with different patient distributions, we optimize the ambulance locations for  
6 each and every city; given subclass partitions of CIFAR-10, we optimize CNN  
7 hyperparameters for each partition. Similarity across tasks boosts optimization of  
8 each task in two ways: in modelling by data sharing across objectives, and also in  
9 acquisition by quantifying how a single point on one task can help learn the optima  
10 of similar tasks. For this we propose a framework for conditional optimization:  
11 ConBO. This can be built on top of a range of acquisition functions and we propose  
12 a new Hybrid Knowledge Gradient acquisition function. The resulting method is  
13 intuitive and theoretically grounded, either matches or significantly outperforms  
14 recently published works on a range of problems, and thanks to the unique nature  
15 of conditional optimization, is easily parallelized to collect a batch of points.

## 16 1 Introduction

17 Expensive stochastic black box functions arise in many fields such as fluid simulations [1], engineering  
18 wing design [2], and machine learning parameter tuning [3]. Bayesian optimization is a powerful set  
19 of tools to optimize such functions, finding the input with highest long term average performance  
20  $x^* = \arg \max_x \mathbb{E}[f(x)]$ , (the expectation represents averaging over the performance noise). In this  
21 work we consider an under-explored generalization of the standard setting previously referred to as  
22 “conditional optimization” [4] where a user has a collection of functions, or *tasks*, and simply seeks  
23 the peak of each function/task. Formally, we have an expensive black box  $f$  that takes as arguments  
24 both a *task*  $s \in S$  and an *input*  $x \in X$ , typically box-constrained continuous variables, and returns a  
25 noisy scalar performance

$$f(s, x) : S \times X \rightarrow \mathbb{R}. \quad (1)$$

26 At each iteration an algorithm determines both task and input  $(s, x)$  then observes performance  
27  $y = f(s, x)$  and the goal is to learn the input with highest average performance for each task

$$x^*(s) = \arg \max_x \mathbb{E}[f(s, x)]. \quad (2)$$

28  $x^*(s)$  is referred to as the optima conditioned on  $s$ , see Figure 1 black line. In certain applications,  
29 one may want to give higher priority to particular tasks hence a task weighting function  $W(s)$  may  
30 also be specified. In this work we consider the following applications.

31 **CNN hyperparameters:** the CIFAR-10 dataset contains 10 classes, this is split into five mutually  
32 exclusive binary classification datasets and for each we train a CNN, each CNN is a task in  $S =$

33  $\{1, \dots, 5\}$ . For each CNN, we optimize dropout rates, batch size and Adam parameters, so  $X \subset \mathbb{R}^7$   
 34 and  $f(s, x)$  is the validation accuracy. We assume all five CNNs have equal priority,  $W(s) = 1/5$ .

35 **Ambulances in a square:** [5] given a range of  
 36  $30\text{km} \times 30\text{km}$  cities, each city is a task with a  
 37 different population centre  $s \in [0, 30]^2$ . For  
 38 a given city, we optimize the 2D location of  
 39 three ambulance bases,  $x \in [0, 30]^{3 \times 2} \subset \mathbb{R}^6$ .  
 40 Given a city  $s$  and ambulance bases  $x$ , a virtual  
 41 environment randomly generates patients  
 42 for a simulated day and the average ambulance  
 43 journey time,  $f(s, x)$ , is returned. Inland cities,  
 44 like Paris or London, with population centres  
 45 in the middle are more common than coastal  
 46 cities, like Singapore or Dubai, with a popula-  
 47 tion centre on a boundary. Thus,  $W(s)$  mirrors  
 48 this distribution of city centres, a Gaussian centred  
 49 at  $s = (15, 15)$  with 7km standard deviation  
 50 and truncated to  $[0, 30]^2$ .

51 **Assemble to order:** [6, 7] a company owns  
 52 many stock warehouses and each one faces a  
 53 different level of demand  $s \in [0.5, 1.5]$ . Equal  
 54 priority is given to all warehouses and so  $W(s)$   
 55 is uniform. At each warehouse, stock levels are  
 56 controlled by setting targets of a control policy  
 57  $x = [0, 20]^8$ . Given a demand level  $s$  and  
 58 control parameters  $x$ , a simulator generates cus-  
 59 tomer orders according to demand, the ware-  
 60 house stock is sold and replenished according to policy  $x$ , and profit over a simulated month is  
 61 returned as  $f(s, x)$ .

62 The closely related multi-task Bayesian optimisation or Bayesian Quadrature optimisation methods  
 63 [8, 9, 10, 11] aim to find a single peak that maximises the weighted sum/integral over tasks,  $x^* =$   
 64  $\arg \max_x \int \mathbb{E}[f(s, x)]W(s)ds$ . However, as they are fundamentally designed for a different purpose,  
 65 we observe that they struggle to find the peak of each and every task, see experiments in Supplementary  
 66 Material 6 (SM 6).

67 Also closely related are contextual optimization algorithms [12, 13, 14], conditional optimization  
 68 has also been referred to as ‘‘offline contextual’’ [15]. In contextual applications, at iteration  $n$  the  
 69 next task (called context)  $s^n$  is passed from the black box to a contextual algorithm that intelligently  
 70 determines  $x^n$ . Then the black box returns both performance  $y^n$  and the next task/context  $s^{n+1}$ . One  
 71 could ‘‘hotfix’’ a contextual algorithm for conditional optimization problems by randomly choosing  
 72  $s^{n+1} \sim \mathbb{P}[s] \propto W(s)$  at each iteration. However, we empirically show that, unsurprisingly, randomly  
 73 choosing the task in each iteration is significantly worse than intelligently choosing the task.

74 Multi-fidelity and multi-information source methods [16, 17, 18, 19] assume there is a single known  
 75 constant target task  $s^* \in S$  corresponding to the highest fidelity level or most expensive information  
 76 source that must be optimized,  $x^* = \arg \max_x \mathbb{E}[f(s^*, x)]$ . If the cost of evaluation varies across the  
 77 function domain,  $c(s, x)$ , cheaper regions of the domain can be evaluated improving sample efficiency.  
 78 Again, such methods may be hotfixed for the setting we consider by artificially designating a target  
 79 task, e.g. the highest priority task  $s^* = \arg \max_s W(s)$ . However, we observe that if cost is constant  
 80 across the domain (as is commonly assumed in the BO literature), these methods greedily optimise  
 81 the artificial target task  $s^*$  and blindly neglect all other tasks, see SM4. This desirable behaviour in  
 82 multi-fidelity optimization problems leads to failure in conditional optimization problems.

83 To the best of our knowledge, there exist only a small number of works that propose algorithms  
 84 specifically designed for the conditional setting.

85 The Surrogate Collaborative Tuning (SCoT) algorithm [20] optimizes a finite set of tasks. It iteratively  
 86 visits each task,  $s$ , in a round-robin fashion and determines the input  $x$  by expected improvement (EI).  
 87 The authors apply this to optimize the hyperparameters of an ML model for multiple datasets. The

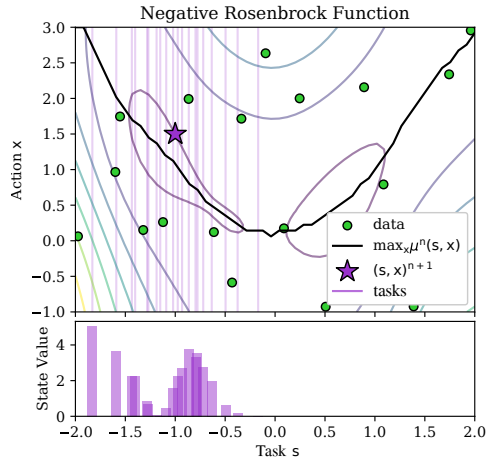


Figure 1: Top: GP model. Each task (vertical slice) is an objective function to maximise. A new sample  $(s, x)^{n+1}$  provides information about the optima of similar tasks. Bottom: the acquisition function value of each task using hybrid Knowledge Gradient.

88 Profile Expected Improvement (PEI) and Profile Expected Quantile Improvement (PEQI) algorithms  
 89 [21, 4] consider continuous tasks and significantly improve upon SCoT by dynamically determining  
 90 the task in each iteration. The acquisition value of a given point  $(s, x)$  is the expected improvement  
 91 of a new output over the best predicted output within the same task.

92 The REVI algorithm [22] at each iteration discretizes both task  $S$  and input  $X$  spaces. Given a new  
 93 hypothetical point,  $(s, x)^{n+1}$ , for each task in the discretization,  $s_i$ , the discrete Knowledge Gradient  
 94 over the  $X$  discretization quantifies how  $(s, x)^{n+1}$  will benefit task  $s_i$ . Summing the benefits over the  
 95 tasks yields the acquisition value of  $(s, x)^{n+1}$ . REVI was designed to account for how all tasks can be  
 96 optimized by each sample. However, this comes with exponential cost, increasing problem dimension  
 97 leads to exponentially increasing discretization size with corresponding exponential space and time  
 98 requirements. On the other hand, sparse discretizations can lead to poor and arbitrary measures  
 99 of acquisition value. Consequently, REVI has only successfully been applied to low-dimensional  
 100 synthetic problems.

101 Most recently, the Multi-task Thompson sampling (MTS) method [15] uses a novel kernel with a  
 102 length scale that varies across tasks. This is combined with a Thompson sampling method for col-  
 103 lecting new data and showed performance improvements over REVI. However, Thompson sampling  
 104 also does not account for how one sample provides benefit for optimizing all tasks, a fundamental  
 105 structural property of conditional optimization.

106 We propose a method that exploits the basic structure of conditional optimization while also being  
 107 highly scalable, and therefore applicable to more challenging real-world problems. We make the  
 108 following contributions:

- 109 1. A framework for conditional Bayesian optimization with theoretical guarantees: ConBO.
- 110 2. A new, fast global optimization method: Hybrid Knowledge Gradient.
- 111 3. State-of-the-art performance on open source problems including CNNs and simulators.

## 112 2 The Conditional Bayesian Optimization Algorithm

113 We first discuss the fitting of the Gaussian process model and the predicted conditional optima. We  
 114 then motivate the acquisition value for a single task and how this is integrated over tasks yielding the  
 115 acquisition function. Because integrating over tasks multiplies the computational burden, we propose  
 116 Hybrid Knowledge Gradient as a solution. See Algorithm 1 in the SM 2 for a high level summary.

117 At a stage after having observed  $n$  data points,  $\{(s^i, x^i, y^i)\}_{i=1}^n$  where  $y^i = f(s^i, x^i)$ , a Gaussian  
 118 process is fit over the joint space  $S \times X$  to scalar outputs  $y$ . Let  $\tilde{X}^n = ((s, x)^1, \dots, (s, x)^n)$  and  
 119  $Y^n = (y^1, \dots, y^n)$ . A Gaussian process is defined by a prior mean and prior covariance function,  
 120  $\mu^0(s, x)$ ,  $k^0((s, x), (s', x'))$  which are chosen for each application, for more information see [23].  
 121 Let the data covariance matrix be  $K = k^0(\tilde{X}^n, \tilde{X}^n) \in \mathbb{R}^{n \times n}$ , the posterior mean is given by

$$\mu^n(s, x) = \mu^0(s, x) + k^0((s, x), \tilde{X}^n)(K + \sigma_0^2 I)^{-1}(Y^n - \mu^0(\tilde{X}^n))$$

122 and the posterior covariance is given by

$$k^n((s, x), (s', x')) = k^0((s, x), (s', x')) + k^0((s, x), \tilde{X}^n)(K + \sigma_0^2 I)^{-1}k^0(\tilde{X}^n, (s', x')).$$

123 At the end of data collection, in standard BO, often the peak posterior mean is returned as the  
 124 predicted best input. Generalizing to the conditional setting we simply condition on  $s$ ,

$$x^{*N}(s) = \arg \max_x \mu^N(s, x). \quad (3)$$

125 During data collection, a new data point  $y^{n+1}$  at  $(s, x)^{n+1}$  will update the Gaussian process over the  
 126 whole domain  $S \times X$ . To construct an acquisition function for conditional optimization, we start  
 127 by looking for standard acquisition functions that account for how the model changes at *unsampled*  
 128 locations  $(s', x') \neq (s, x)^{n+1}$ . Specifically, the popular Expected improvement (EI) [24] and upper  
 129 confidence bound (UCB) [25] methods are both functions of the mean and kernel *at the sampled point*  
 130 *only*, i.e. of  $(\mu^n((s, x)^{n+1}), k^n((s, x)^{n+1}, (s, x)^{n+1}))$ , hence would require non-trivial modification  
 131 to be able to account for how a sample affects similar tasks. Methods that *do* utilise the mean and  
 132 kernel at unsampled points include Entropy search (ES [26, 27] and PES [28]) that measures the

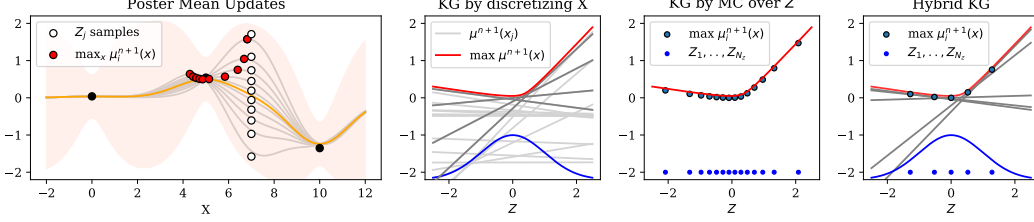


Figure 2: Methods for computing  $\text{KG}(x^{n+1})$  at  $x^{n+1} = 7$ . Left:  $\mu^n(x)$  and samples of  $\mu^{n+1}(x)$  determined by a scalar  $Z \sim N(0, 1)$ . Centre-left:  $\text{KG}_d$  replaces  $X$  with up to 3000 points  $x_i \in X_d$  and  $\mu^{n+1}(x_i)$  is linear in  $Z$ . Centre-right:  $\text{KG}_{MC}$  samples up to 1000 functions  $\mu^{n+1}(x)$  functions and maximises each of them numerically. Right:  $\text{KG}_h$  samples up to 5 functions  $\mu^{n+1}(x)$  and maximizes them numerically, the arg max points  $x_1^*, \dots, x_5^*$  are used as  $X_d$  in  $\text{KG}_d$ .

133 mutual information between the new output  $\mathbb{P}[y^{n+1}|x^{n+1}]$  and the (uncertain) location of the peak  
 134  $\mathbb{P}[x^*|\tilde{X}^n, Y^n]$ , Max-value entropy search (MES) [29] that measures mutual information between  
 135 the new output  $\mathbb{P}[y^{n+1}|x^{n+1}]$  and the (uncertain) largest output  $\mathbb{P}[\max y|\tilde{X}^n, Y^n]$ , and Knowledge  
 136 Gradient (KG) [30] that measures the expected peak of the new posterior mean  $\mathbb{E}[\max \mu^{n+1}(x)]$   
 137 caused by a new  $y^{n+1}$  at  $x^{n+1}$ .

138 Each single task  $s$  defines a single global optimization problem over  $x \in X$ . Given a proposed  
 139 sample  $(s, x)^{n+1} = (s^{n+1}, x^{n+1})$ , the acquisition benefit for task  $s$  may be computed using ES, PES,  
 140 MES or KG. In this work we adopt KG for its Bayesian decision theoretic derivation that extends  
 141 seamlessly to the conditional setting. For KG, the benefit for a given task is the expected increase in  
 142 peak predicted performance within the task. We denote the task-conditioned KG as

$$\text{KG}_c(s; (s, x)^{n+1}) = \mathbb{E}_{y^{n+1}}[\max_{x'} \mu^{n+1}(s, x')|(s, x)^{n+1}] - \max_{x''} \mu^n(s, x'').$$

143 We discuss numerical evaluation of  $\text{KG}_c(\cdot)$  in Section 2.1. Similar expressions for conditioned  
 144 entropy methods,  $\text{ES}_c(\cdot)$ ,  $\text{PES}_c(\cdot)$ ,  $\text{MES}_c(\cdot)$ , are derived in SM 7. Integrating over all tasks  $s$  yields  
 145 the total acquisition value

$$\int_S \text{KG}_c(s; (s, x)^{n+1}) W(s) ds. \quad (4)$$

146 For discrete  $S$  the integral is replaced by summation. For continuous  $S$ , the integral over tasks  $s$   
 147 cannot be computed analytically so we use Monte Carlo with importance sampling. When using  
 148 a kernel that factorises  $k(s, x, s', x') = \sigma_0^2 k_S(s, s') k_X(x, x')$ , like squared exponential or Matérn,  
 149 similarity across tasks is encoded in  $k_S(s, s')$ . This naturally leads to the proposal distribution  
 150  $q(s|s^{n+1}) \propto k_S(s, s^{n+1})$ . In our continuous task experiments, we use the Matérn kernel and a  
 151 Gaussian proposal distribution with mean  $s^{n+1}$  and the task kernel length scales,  $l_s$ , as standard  
 152 deviations,

$$q(s|s^{n+1}) \sim \mathcal{N}(s|s^{n+1}, \text{diag}(l_s^2)). \quad (5)$$

153 We generate  $n_s = 20$  samples  $S_{MC} = \{s_1, \dots, s_{n_s}\}$ , finally the acquisition function is

$$\text{ConBO}(s^{n+1}, x^{n+1}) = \sum_{s_i \in S_{MC}} \frac{W(s_i)}{q(s_i|s^{n+1})} \text{KG}_c(s_i; (s, x)^{n+1}).$$

154 Figure 1 shows a set of sampled tasks and the  $\text{KG}_c(\cdot)$  for each one. Each KG term directly measures  
 155 increase in predicted performance for one task, e.g. if  $y$  values are dollar amounts, ConBO with  
 156  $\text{KG}_c(\cdot)$  is the sum of dollar increases over all tasks. However, for entropy based methods, ConBO  
 157 becomes a sum of Shannon information units thereby indirectly optimizing the dollar amounts.  
 158 The randomly sampled tasks  $S_{MC}$  may be resampled with each call to  $\text{ConBO}(s, x)$  and gradients  
 159 estimated enabling the optimal  $(s, x)^{n+1}$  to be found with a stochastic gradient ascent optimizer such  
 160 as Adam [31]. Selecting each point according to maximising ConBO is also myopically optimal in a  
 161 value of information framework:

162 **Theorem 1** Let  $(s^*, x^*) \in \arg \max \text{ConBO}(s, x)$  be a point chosen for sampling.  $(s^*, x^*)$  is also  
 163 the point that maximises the myopic Value of Information, the increase in predicted performance.

164 Further, in finite search space, with an infinite sampling budget all points will be sampled infinitely:

---

**Algorithm 1** Computing ConBO( $s, x$ ). For each call with a candidate (task, input) point, similar tasks are sampled. For each sampled task a cheap acquisition function, hybrid KG, is evaluated. The output is the importance weighted average of hybrid KG values. Details and full numerical expressions are given in SM 3.

---

**Require:** candidate point  $(s, x)^{n+1}$ , parameters  $n_s, n_z$   
 Sample  $n_s$  tasks similar to  $s^{n+1}$ ,  $S_{MC} \sim q(s|s^{n+1})$   
 Initialize output value  $Q \leftarrow 0$   
**parfor**  $s_i$  **in**  $S_{MC}$  **do**

Initialize  $\tilde{X}_i^* \leftarrow \{\}$   
**parfor**  $j$  **in**  $1, \dots, n_z$  **do**  
 $z_j \leftarrow \Phi^{-1}((2j-1)/2n_z)$   
 $y_j^{n+1}$  computed with  $(s, x)^{n+1}$  and  $z_j$   
 $\mu_j^{n+1}(s, x)$  constructed by rank-1 update with  $(s, x, y_j)^{n+1}$   
 $x_j^* \leftarrow \arg \max_x \mu_j^{n+1}(s_i, x)$  with `Optimizer()`  
 $\tilde{X}_i^* \leftarrow \tilde{X}_i^* \cup (s_i, x_j^*)$   
**end parfor**  
 $\alpha_i \leftarrow \text{KG}_d(\tilde{X}_i^*)$

Hybrid Knowledge Gradient

$Q \leftarrow Q + \alpha_i W(s_i) / q(s_i|s^{n+1})$   
**end parfor**  
**return** average over tasks  $Q/n_s$

---

165 **Theorem 2** Let  $S$  and  $X$  be finite sets and  $N$  the budget to be sequentially allocated by ConBO.  
 166 Let  $n(s, x, N)$  be the number of samples allocated to a point  $(s, x)$  within budget  $N$ . Then for all  
 167  $(s, x) \in S \times X$  we have that  $\lim_{N \rightarrow \infty} n(s, x, N) = \infty$ .

168 The law of large numbers ensures that the algorithm learns the true expected performance for all  
 169 points. Proofs are given in the SM 1.

## 170 2.1 Hybrid Knowledge Gradient

171 By definition, KG is more expensive than EI and UCB. Further, the function  $\text{KG}_c(s_i, (s, x)^{n+1})$   
 172 must be computed once for each sampled task  $s_i$ , the computational cost is therefore  $n_s$  times the  
 173 global acquisition function equivalent. To alleviate this cost, we propose a novel, efficient algorithm  
 174 for computing KG. In the following section we assume constant  $s$  for brevity, reducing to the  
 175 global optimization setting. Given a hypothetical location  $x^{n+1}$ , KG quantifies the value of a new  
 176 hypothetical observation  $y^{n+1}$  by the expected increase in the peak of the posterior mean

$$\text{KG}(x^{n+1}) = \mathbb{E}_{y^{n+1}} \left[ \max_{x'} \mu^{n+1}(x') | x^{n+1} \right] - \max_{x''} \mu^n(x''). \quad (6)$$

177 However,  $\max_{x'} \mu^{n+1}(x')$  has no explicit formula and approximations are required which we describe  
 178 next. At time  $n$ , the next posterior mean is unknown, however, it may be written as  $\mu^{n+1}(x) =$   
 179  $\mu^n(x) + \tilde{\sigma}(x; x^{n+1})Z$  where  $\tilde{\sigma}(x; x^{n+1})$  is a deterministic function and the scalar  $Z \sim \mathcal{N}(0, 1)$   
 180 captures the randomness of  $y^{n+1}$ , see SM 2. Previously,  $\text{KG}(x)$  has been computed in two ways.

**KG by discretization** [30, 32]: in Equation 6, the maximizations may be performed over a discrete set of  $d$  points  $X_d \subset X$ . Denoting  $\underline{\mu} = \mu^n(X_d) \in \mathbb{R}^d$  and  $\tilde{\sigma}(x^{n+1}) = \tilde{\sigma}(X_d; x^{n+1}) \in \mathbb{R}^d$ , then

$$\text{KG}_d(x^{n+1}) = \mathbb{E}_Z \left[ \max \{ \underline{\mu} + \tilde{\sigma}(x^{n+1})Z \} \right] - \max \underline{\mu}.$$

181 The  $\max \{ \underline{\mu} + \tilde{\sigma}(x^{n+1})Z \}$  is a piece-wise linear function of  $Z$  and the expectation is analytically  
 182 tractable. The output is a *lower bound* of the true  $\text{KG}(x)$ . REVI [22] and the MiSo algorithm  
 183 [17] used  $\text{KG}_d$  with 3000 uniformly random distributed points. This method suffers the curse of  
 184 dimensionality as  $X_d$  must grow exponentially with dimension. Further, the discretization may likely  
 185 contain many useless points in uneventful regions of  $X$ , see Figure 2 centre-left plot.

186

KG	Type of Estimate	$n_z = 3$	$n_z = 5$	$n_z = 7$	$n_z = 50$
Discrete	lower bound	0.24 (1.37)	0.41 (1.49)	0.55 (1.91)	2.03 (2.2)
MC	unbiased	3.54 (4.79)	3.73 (4.51)	3.50 (3.39)	3.36 (1.18)
MC++	unbiased	2.97 (3.78)	3.50 (2.68)	3.22 (1.69)	3.32 (0.2)
Hybrid	lower bound	3.15 (0.00)	3.28 (0.00)	3.31 (0.00)	3.34 (0.00)

Table 1: We collect 20 points on the Rosenbrock function, a point was randomly selected and KG computed by the different methods. Each calculation is repeated 50 times and we report mean and two standard deviations. The Monte Carlo methods suffer from high variance, the discrete method is volatile and returns a very loose lower bound. Hybrid KG is very stable with errors too small to show.

187 **KG by Monte Carlo** [33, 34]: given  $x^{n+1}$ , the method samples up to  $n_z = 1000$  Gaussian values  
188 of  $Z$ . For each  $Z_j$ , construct the posterior mean,  $\mu_j^{n+1}(x)$ , find the maximum with a continuous  
189 numerical `Optimizer()` like L-BFGS or CG. Averaging the maxima from all  $Z_j$  yields

$$\text{KG}_{MC}(x^{n+1}) = \frac{1}{n_s} \sum_j \text{Optimizer}_{x'}(\mu_j^{n+1}(x)) - \text{Optimizer}_{x''}(\mu^n(x'')).$$

190 The result is an *unbiased* estimate of true  $\text{KG}(x)$  and scales better to higher dimensional  $X$ , the  
191 univariate  $Z$  is discretized by Monte Carlo samples instead of  $X$ . However, for a good estimate,  $n_z$   
192 must be large, many `Optimizer()` calls are required. See Figure 2 centre right.

193 We instead propose a simple mixture of the two approaches above that both scales to higher dimen-  
194 sional  $X$  and drastically reduces the number of `Optimizer()` calls.

**Hybrid KG:** given  $x^{n+1}$ , following  $\text{KG}_{MC}$  we use  $n_z = 5$  values of  $Z_j$ , construct  $\mu_j^{n+1}(x)$  and use  
`Optimizer()` to find the peak location  $x_j^*$ . The set of peak locations  $X_d^* = \{x_1^*, \dots, x_{n_z}^*\}$  is used  
as a *dynamic optimized* discretization in  $\text{KG}_d$  thus analytically computing an extremely tight lower  
bound of the true  $\text{KG}(x)$ . Let  $\underline{\mu}^* = \mu^n(X_d^*)$  and similarly for  $\tilde{\sigma}^*(x^{n+1})$ , Hybrid KG is given by

$$\text{KG}_h(x^{n+1}) = \mathbb{E}_Z[\max \underline{\mu}^* + \tilde{\sigma}^*(x^{n+1})Z] - \max \underline{\mu}^*.$$

195 Compared to  $\text{KG}_d$ , the hybrid method removes redundant points in the discretization  $X_d$ , all  $X_d^*$   
196 points contribute to  $\max \mu^{n+1}(X_d^*)$  and there are far fewer points. Compared to  $\text{KG}_{MC}$  that samples  
197 many  $Z_j$  and optimizes many  $x_j^*$ , Hybrid KG uses far fewer  $Z_j$  and optimizes far fewer  $x_j^*$ . See  
198 Algorithm 1 for a summary of evaluating ConBO with Hybrid KG. ConBO can be reduced to the  
199 REVI algorithm by replacing the dynamic importance sampled tasks with a pre-frozen discretization  
200 of tasks and the (dynamic optimized) hybrid KG with discrete KG over a pre-frozen discretization of  
201 inputs. These changes drastically improve scalability enabling ConBO to be applied to a far broader  
202 range of applications.

203 To ensure asymptotic convergence, in a discrete domain, we require that the acquisition function  
204 is non-negative,  $\text{KG}_h(x) \geq 0$ , and the acquisition function is zero where GP variance is zero,  
205  $\text{KG}_h(x) = 0 \iff k^n(x, x) = 0$ . Therefore, always choosing  $x^{n+1} = \arg \max \text{KG}_h(x)$  ensures  
206 only points with GP variance will be revisited until all points have no variance i.e. the true function is  
207 known for all points. We can ensure these properties by setting  $n_z \geq 2$  and at least one  $Z_j$  is equal to  
208 zero.

209 **Theorem 3** Let  $n_z \geq 2$  and let  $\underline{Z} = \{Z_j | j = 1, \dots, n_z\}$ . If  $0 \in \underline{Z}$  then  $\text{KG}_h(x) \geq 0$  for all  $x \in X$   
210 and if  $x$  is sampled infinitely often  $\text{KG}_h(x) = 0$ .

211 Proof is in the SM 1. The  $Z_j$  values can be fixed, for  $n_z = 5$  we use equal Gaussian quantiles  
212  $\underline{Z} = \{\Phi^{-1}(0.1), \Phi^{-1}(0.3), \dots, \Phi^{-1}(0.9)\}$  where  $\Phi(\cdot)$  is the Gaussian CDF. Using quantile spacing  
213 and odd  $n_z$  ensures  $Z_j = \Phi^{-1}(0.5) = 0$  is included which satisfies the assumptions of asymptotic  
214 convergence. See Fig.2 (right).

215 The computational complexity of a single call to ConBO requires the posterior variance ( $O(n^2)$ ) and  
216  $n_s n_z$  runs of `Optimizer()`. Let  $n_{calls}$  be the number of times that `Optimizer()` calls the posterior  
217 mean costing  $O(n)$ . Thus, ConBO total complexity is  $O(n^2 + n n_{calls} n_s n_z)$ . Note this is linear in  
218  $n_s n_z$ , the size of the (small) dynamic optimal discretization over  $S \times X$ . Thompson sampling with

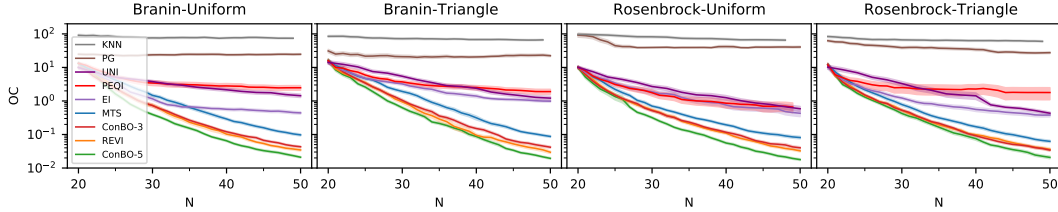


Figure 3: Opportunity Cost across a range of synthetic test problems. The dummy baseline, KNN, is worst in all cases, policy gradient is better however the Gaussian process based methods all perform better. UNI and EI are not conditional algorithms yet outperform PEQI. Amongst other conditional algorithms, MTS, REVI and ConBO methods all perform significantly better. ConBO-3 is outperformed by ConBO-5 demonstrating the improvement with more accurate  $KG_h$ .

219 discretization uses one operation that scales cubically as  $O(n^2(n_s n_z) + n(n_s n_z)^2 + (n_s n_z)^3)$  in  
 220 the worst case and to reduce cost special techniques are required e.g. Fourier features, CG matrix  
 221 inversion.

### 222 3 Experiments

223 We consider synthetic benchmarks and the three applications described in Section 1. In SM 4 we  
 224 also present parallelization batch sampling results, and experiments with contextual, multi-task and  
 225 multi-fidelity methods in SM 6. We also briefly test Hybrid KG for global optimization, SM 5, and  
 226 observe that  $KG_{MC}$  performs worse in the same computation time hence we exclude  $KG_{MC}$  methods  
 227 from the conditional experiments. For each benchmark, for evaluation, held out test tasks are sampled  
 228 from  $\mathbb{P}[s] \propto W(s)$ , for each test task,  $s_i^{test}$ , the predicted optimal input is computed,  $x^*(s_i^{test})$  and  
 229 we report the true black box output averaged over all test tasks, see SM 3 for details.

#### 230 3.1 How Accurate is Hybrid Knowledge Gradient?

231 The theoretical Knowledge Gradient cannot be analytically computed and must be approximated  
 232 (similarly, entropy based methods are not fully tractable). To illustrate the quality of the approxima-  
 233 tions, we collected 20 points from the Rosenbrock test function and fit a Gaussian process. Another  
 234 point was then randomly selected and true KG was estimated by the different methods, keeping the  
 235  $20 + 1$  points fixed and only varying the  $n_z$  parameter. Each estimate of KG was recalculated 50  
 236 times and Table 1 shows the mean and two standard deviations. Discrete KG results in an extremely  
 237 loose, volatile bound. The (state of the art) MC method with  $n_z = 50$  has a run-to-run variance of  
 238  $\pm 35\%$ , variance reduction techniques used in MC++ (latin hypercube inverse sampling and control  
 239 variates) can reduce this to  $\pm 6\%$ . Hybrid KG with  $n_z = 3$  is extremely stable, run-to-run variance  
 240 is too small to show, and is within the error margins of MC++ with  $n_z = 50$  which is  $\sim 17\times$  more  
 241 expensive to compute. Increasing  $n_z$  tightens the bound, hybrid KG with  $n_z = 5$  has 98.2% of the  
 242 value of  $n_z = 50$ . As a result, Hybrid KG is much easier to optimize, the Adam optimizer may be  
 243 used with a much higher learning rate and lower momentum and thus converges much faster.

244 As an aside, the recently proposed one-shot KG [34] enhances the optimization of  $KG_{MC}$ . By  
 245 freezing the  $Z$  values between calls to  $KG_{MC}$  thus  $\tilde{X}^*$  may be reused, this enables *joint* gradient  
 246 ascent of  $(x^{n+1}, \tilde{X}^*)$ . In the global optimization use case (a constant single task) one-shot KG and  
 247 Hybrid KG can be combined. In the conditional setting, the  $\tilde{X}_i^*$  of past sampled tasks may be saved,  
 248 however each call to ConBO must sample new tasks that are not in the saved history, one-shot KG  
 249 style joint optimization is not possible. In our implementation we utilise caching of  $\tilde{X}_i^*$  from old  
 250 tasks to heuristically warm start finding  $\tilde{X}_i^*$  for new tasks, see SM 3.

#### 251 3.2 Synthetic Functions

252 We perform low-dimensional toy experiments in an ideal setting as a sanity check where we expect  
 253 all conditional methods to perform similarly. We use the popular Branin-Hoo and Rosenbrock test  
 254 functions in 2D defining the (task, input) domain as displayed in Figure 1.

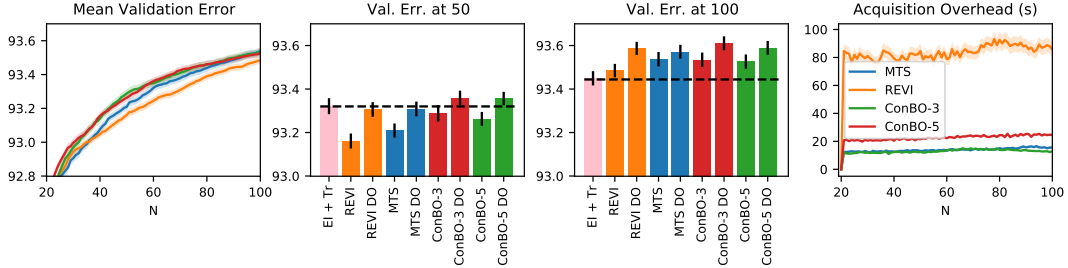


Figure 4: Left: validation accuracy. Centre-left: validation accuracy after 50 samples. Centre-right: validation accuracy after 100 samples. Right: algorithm overhead in seconds. After 50 samples, none of the multi-task models outperform the baseline, EI + Tr (dashed line) suggesting all datasets can use similar hyperparameters. For the larger budget 100, all models outperform the baseline by 0.1% suggesting that for more fine-tuning, each dataset requires different hyperparameters. In all cases, performing data optimization significantly increases performance.

255 **Synthetic Functions** Using the Rosenbrock and Branin-Hoo functions, we consider a uniform and a  
 256 triangular task weighting  $W(s)$ . For baselines, we adopt two policy based methods. **KNN**: (dummy  
 257 baseline) randomly collect data,  $x_{KNN}^*(s)$  takes a task,  $s$ , and returns the best input from 10 nearest  
 258 neighbor tasks. **PG**: policy gradient, a parametric quadratic policy  $x_{PG}^*(s) = \pi_\theta(s)$  is learnt by  
 259 maximising observed performance values, each iteration samples a task from  $\mathbb{P}[s] \propto W(s)$  then  $x$  is  
 260 sampled with an  $\epsilon$ -greedy strategy. For a controlled ablation study, all the BO methods fit the same  
 261 GP and for evaluation use the same definition of  $x^*(s)$  (Equation 3), methods only differ by their  
 262 data acquisition strategy. **UNI**: random data collection, the most recent conditional methods **PEQI**,  
 263 **REVI**, **MTS** with all parameters given in SM 3. **ConBO- $n_z$** : given  $(s, x)$ , 20 tasks are importance  
 264 sampled,  $KG_h$  with  $n_z = 3, 5$  points is used. **EI**: expected improvement that treats  $(s, x)$  as inputs to  
 265 be optimized.

266 Results are shown in Figure 3. Policy based methods KNN and PG consistently perform worse than  
 267 the Gaussian process methods. Surprisingly, the conditional BO algorithm PEQI performs similarly to  
 268 UNI and much worse than EI. All other conditional methods outperform all non-conditional methods.

### 269 3.3 CNN Training Hyperparameters

270 We apply MTS, REVI, and ConBO variants, and we adopt the recently proposed kernel used for BO  
 271 with Common Random Numbers [35],

$$k((s, x), (s', x')) = \sigma_0^2 M(x, x'; \underline{l}) + \delta_{s's} (\sigma_1^2 M(x, x'; \underline{l}) + \sigma_3^2).$$

272 where  $M(x, x'; \underline{l})$  is a Matérn  $\frac{5}{2}$  kernel with length scales  $\underline{l}$ . The first term models a common trend  
 273 function across all tasks and the second term models how each task independently differs from the  
 274 trend. The differences are composed of another Matérn and the constant kernel to model a global  
 275 offset e.g. one dataset may have universally higher validation accuracy. This kernel has far fewer  
 276 parameters than a full multi-task product kernel, it is easy to fit and scales to an arbitrary number of  
 277 tasks (or datasets) without adding extra parameters.

278 In this problem setting, learning hyperparameters over similar datasets, one may expect that the  
 279 optimal hyperparameters would be the same for all datasets. Therefore, as a baseline we apply EI  
 280 to learn the hyperparameters of the first dataset (task 1). We then evaluate the objective function  
 281 (validation accuracy) on the rest of the datasets using the best observed hyperparameters from dataset  
 282 1, we refer to this as **EI + Transfer**.

283 **Argument Optimization versus Data Optimization** In continuous task settings, it is not possible to  
 284 evaluate every task. In discrete task settings with large sampling budgets  $N \gg |S|$ , a user may desire  
 285 a single high (although stochastic) output value,  $\max y$ , for every task. For example, in network  
 286 hyperparameter optimization, the network with the best validation error,  $\max y$ , will be deployed  
 287 (and the hyperparameters  $x$  may or may not be reused). We refer to this is *data optimization* (DO).  
 288 For simulated environments, the input (or “action”) that generalizes providing the best long-term  
 289 average performance,  $\max_x \mathbb{E}[f(s, x)]$ , is deployed and we refer to this is *argument optimization*  
 290 (AO). Past work [22] has shown that argument optimization finds inputs that generalize better but



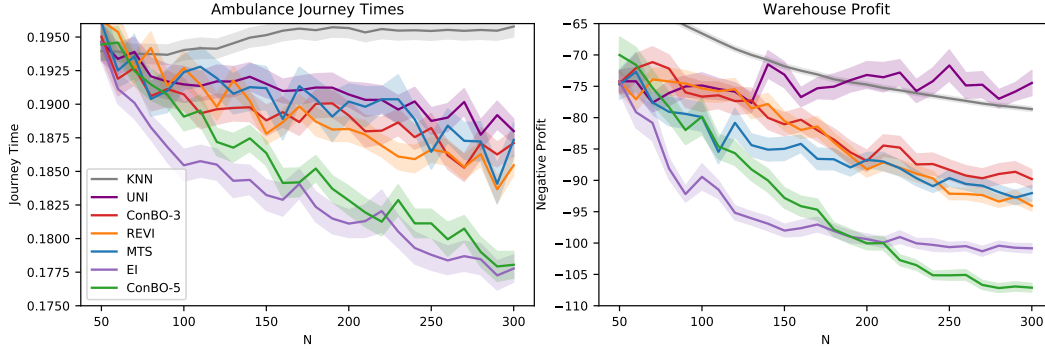


Figure 5: Left: average journey times across a range of cities. Right : average profit across a range of warehouses. ConBO-5 and EI perform best on these benchmarks.

291 may not provide optimal  $\max y$  for all tasks during the optimization run. The authors propose a “DO  
 292 trick”, use an AO method for  $N - |S|$  iterations, then finally allocate one sample per task with input,  
 293  $x^{n+1}$ , determined by EI within the task. We apply this trick to all algorithms in this experiment.

294 Results are shown in Figure 4. For the medium budget of 50 samples, ConBO performs best of  
 295 the standard algorithms yet it is still worse than the EI+Tr baseline. Applying the final round of  
 296 DO improves all results to match the baseline. For the large budget of 100 samples, all methods  
 297 outperform the baseline suggesting dataset specific fine-tuning of hyperparameters is required to  
 298 achieve best results. Again, DO provides a significant boost to performance for all methods.

299 Gaussian process kernel parameter learning required approximately 2–5 seconds using Tensorflow. In  
 300 Figure 4 (right) we show the runtime of each algorithm excluding model fitting and network training,  
 301 purely acquisition function optimization time. MTS and ConBO-3 are quickest while Conbo-5  
 302 increases linearly over ConBO-3 and REVI takes much longer.

### 303 3.4 Ambulances and Warehouses

304 We apply all methods from Section 3.2 to two benchmarks from the `www.SimOpt.org` library for  
 305 simulation optimization problems. The ambulance problem (AMB) is 8-dimensional and consists of  
 306 a range of cities and one must optimize ambulance locations for each city. The Assemble-to-order  
 307 problem (ATO) is 9-dimensional consisting of a range of warehouses and one must optimize target  
 308 stock level for each warehouse. Results are shown in Figure 5. Of the policy based methods, PG  
 309 performs poorly and does not show on the plots whilst KNN performs poorly on AMB and performs  
 310 well on ATO suggesting that AMB is a more difficult problem. Of the GP based methods, EI performs  
 311 well for smaller budgets. Although it is not a conditional algorithm we include it to highlight that  
 312 sometimes the simplest idea can also work. Of the conditional methods, MTS, REVI, and ConBO-3  
 313 all perform similarly, either slightly (AMB) or largely (ATO) outperforming UNI. These methods  
 314 struggle in higher dimensions while ConBO-5 uses a more accurate acquisition function and is the  
 315 only method that consistently performs well *across all problems*. We hypothesize that these problems  
 316 are more difficult than the synthetics and CNN and truly stress test conditional algorithms.

## 317 4 Conclusion

318 **Potential Limitations and Broader Impact** there are multiple ways in which ConBO may fail, in  
 319 this work we have not investigated how ConBO or Hybrid KG suffers with poorly learnt Gaussian  
 320 process hyperparameters. In many applications, a poorly chosen kernel or unoptimized hyperparam-  
 321 eters can lead to poor performance and our proposed methods may be more sensitive or more robust to  
 322 these failures than alternative approaches. We propose a general purpose optimization algorithm and  
 323 analysis, we do not currently see any immediate societal impact.

324 We investigate Conditional Bayesian optimization and propose ConBO. ConBO is designed from  
 325 the ground up to fully exploit the structure of conditional problems, namely that optimizing one task  
 326 helps optimize similar tasks. Hence every point should be collected to maximise the benefit of all  
 327 tasks. However, this can lead to excessive computational cost, particularly in higher dimensions.  
 328 Thus we also propose Hybrid KG that mixes past methods to be both fast and scalable.

329 **References**

- 330 [1] Victor Picheny and David Ginsbourger. A nonstationary space-time gaussian process model for  
331 partially converged simulations. *SIAM/ASA Journal on Uncertainty Quantification*, 1(1):57–78,  
332 2013.
- 333 [2] Shinkyu Jeong, Mitsuhiro Murayama, and Kazuomi Yamamoto. Efficient optimization design  
334 method using kriging model. *Journal of aircraft*, 42(2):413–420, 2005.
- 335 [3] Jasper Snoek, Hugo Larochelle, and Ryan P Adams. Practical bayesian optimization of machine  
336 learning algorithms. In *Advances in neural information processing systems*, pages 2951–2959,  
337 2012.
- 338 [4] Diariétou Sambakhé, Lauriane Rouan, Jean-Noël Bacro, and Eric Gozé. Conditional opti-  
339 mization of a noisy function using a kriging metamodel. *Journal of Global Optimization*,  
340 73(3):615–636, 2019.
- 341 [5] Naijia Anna Dong, David J Eckman, Xueqi Zhao, Shane G Henderson, and Matthias Poloczek.  
342 Empirically comparing the finite-time performance of simulation-optimization algorithms. In  
343 *2017 Winter Simulation Conference (WSC)*, pages 2206–2217. IEEE, 2017.
- 344 [6] Jing Xie, Peter I Frazier, and Stephen E Chick. Bayesian optimization via simulation with  
345 pairwise sampling and correlated prior beliefs. *Operations Research*, 64(2):542–559, 2016.
- 346 [7] Matthias Poloczek, Jialei Wang, and Peter I Frazier. Warm starting bayesian optimization. In  
347 *2016 Winter Simulation Conference (WSC)*, pages 770–781. IEEE, 2016.
- 348 [8] Kevin Swersky, Jasper Snoek, and Ryan Prescott Adams. Multi-task bayesian optimization. In  
349 *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2013.
- 350 [9] Supratik Paul, Konstantinos Chatzilygeroudis, Kamil Ciosek, Jean-Baptiste Mouret, Michael  
351 Osborne, and Shimon Whiteson. Alternating optimisation and quadrature for robust control. In  
352 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- 353 [10] Saul Toscano-Palmerin and Peter I Frazier. Bayesian optimization with expensive integrands.  
354 *arXiv preprint arXiv:1803.08661*, 2018.
- 355 [11] Michael Pearce and Juergen Branke. Bayesian simulation optimization with input uncertainty.  
356 In *2017 Winter Simulation Conference (WSC)*, pages 2268–2278. IEEE, 2017.
- 357 [12] Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff  
358 functions. In Geoffrey Gordon, David Dunson, and Miroslav Dudík, editors, *Proceedings of  
359 the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of  
360 *Proceedings of Machine Learning Research*, pages 208–214, Fort Lauderdale, FL, USA, 11–13  
361 Apr 2011. JMLR Workshop and Conference Proceedings.
- 362 [13] Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff  
363 functions. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence  
364 and Statistics*, pages 208–214. JMLR Workshop and Conference Proceedings, 2011.
- 365 [14] Supratik Paul, Michael A Osborne, and Shimon Whiteson. Contextual policy optimisation.  
366 *arXiv preprint arXiv:1805.10662*, 2018.
- 367 [15] Ian Char, Youngseog Chung, Willie Neiswanger, Kirthevasan Kandasamy, Andrew Oakleigh  
368 Nelson, Mark Boyer, Egemen Kolemen, and Jeff Schneider. Offline contextual bayesian  
369 optimization. In *Advances in Neural Information Processing Systems*, pages 4629–4640, 2019.
- 370 [16] Deng Huang, Theodore T Allen, William I Notz, and R Allen Miller. Sequential kriging  
371 optimization using multiple-fidelity evaluations. *Structural and Multidisciplinary Optimization*,  
372 32(5):369–382, 2006.
- 373 [17] Matthias Poloczek, Jialei Wang, and Peter Frazier. Multi-information source optimization. In  
374 *Advances in Neural Information Processing Systems*, pages 4289–4299, 2017.

- 375 [18] Kirthevasan Kandasamy, Gautam Dasarathy, Jeff Schneider, and Barnabás Póczos. Multi-fidelity  
376 bayesian optimisation with continuous approximations. In *Proceedings of the 34th International*  
377 *Conference on Machine Learning-Volume 70*, pages 1799–1808. JMLR. org, 2017.
- 378 [19] Aaron Klein, Stefan Falkner, Simon Bartels, Philipp Hennig, and Frank Hutter. Fast Bayesian  
379 Optimization of Machine Learning Hyperparameters on Large Datasets. In Aarti Singh and  
380 Jerry Zhu, editors, *Proceedings of the 20th International Conference on Artificial Intelligence*  
381 *and Statistics*, volume 54 of *Proceedings of Machine Learning Research*, pages 528–536, Fort  
382 Lauderdale, FL, USA, 20–22 Apr 2017. PMLR.
- 383 [20] R. Bardenet, Mátyás Brendel, Balázs Kégl, and Michele Sebag. Collaborative hyperparameter  
384 tuning. In *International Conference on Machine Learning*, pages 199–207, 2013.
- 385 [21] David Ginsbourger, Jean Baccou, Clément Chevalier, Frédéric Perales, Nicolas Garland, and  
386 Yann Monerie. Bayesian adaptive reconstruction of profile optima and optimizers. *SIAM/ASA*  
387 *Journal on Uncertainty Quantification*, 2(1):490–510, 2014.
- 388 [22] Michael Pearce and Juergen Branke. Continuous multi-task bayesian optimisation with correla-  
389 tion. *European Journal of Operational Research*, 270(3):1074–1085, 2018.
- 390 [23] C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. MIT Press,  
391 2004.
- 392 [24] D. R. Jones, M. Schonlau, and W. J. Welch. Efficient global optimization of expensive black-box  
393 functions. *Journal of Global optimization*, 13(4):455–492, 1998.
- 394 [25] Andreas Krause and Cheng S Ong. Contextual gaussian process bandit optimization. In  
395 *Advances in neural information processing systems*, pages 2447–2455, 2011.
- 396 [26] Philipp Hennig and Christian J Schuler. Entropy search for information-efficient global opti-  
397 mization. *Journal of Machine Learning Research*, 13(Jun):1809–1837, 2012.
- 398 [27] J. Villemonteix, E. Vazquez, and E. Walter. An informational approach to the global optimization  
399 of expensive-to-evaluate functions. *Journal of Global Optimization*, 44(4):509–534, 2009.
- 400 [28] J. M. Hernández-Lobato, M. W. Hoffman, and Z. Ghahramani. Predictive entropy search  
401 for efficient global optimization of black-box functions. In *Advances in Neural Information*  
402 *Processing Systems*, pages 918–926. Curran Associates, Inc., 2014.
- 403 [29] Zi Wang and Stefanie Jegelka. Max-value entropy search for efficient bayesian optimization.  
404 In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages  
405 3627–3635. JMLR. org, 2017.
- 406 [30] P. Frazier, W. Powell, and S. Dayanik. The knowledge-gradient policy for correlated normal  
407 beliefs. *INFORMS Journal on Computing*, 21(4):599–613, 2009.
- 408 [31] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*  
409 *arXiv:1412.6980*, 2014.
- 410 [32] Warren Scott, Peter Frazier, and Warren Powell. The correlated knowledge gradient for  
411 simulation optimization of continuous parameters using gaussian process regression. *SIAM*  
412 *Journal on Optimization*, 21(3):996–1026, 2011.
- 413 [33] Jian Wu and Peter I Frazier. Discretization-free knowledge gradient methods for bayesian  
414 optimization. *arXiv preprint arXiv:1707.06541*, 2017.
- 415 [34] Maximilian Balandat, Brian Karrer, Daniel Jiang, Samuel Daulton, Ben Letham, Andrew G Wil-  
416 son, and Eytan Bakshy. Botorch: A framework for efficient monte-carlo bayesian optimization.  
417 *Advances in Neural Information Processing Systems*, 33, 2020.
- 418 [35] Michael Pearce, Matthias Poloczek, and Juergen Branke. Bayesian optimization allowing for  
419 common random numbers. *arXiv preprint arXiv:1910.09259*, 2019.

420 **Checklist**

- 421 1. For all authors...
- 422 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
423 contributions and scope? [Yes] ConBO is theoretically grounded, see Theorems 1-3.  
424 ConBO outperforms recently published works on a range of problems, see Section 3,  
425 ConBO is easily parallelizable, see SM Section 4.
- 426 (b) Did you describe the limitations of your work? [Yes] see Section 4
- 427 (c) Did you discuss any potential negative societal impacts of your work? [Yes] see Section  
428 4
- 429 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
430 them? [Yes]
- 431 2. If you are including theoretical results...
- 432 (a) Did you state the full set of assumptions of all theoretical results? [Yes] see SM Section  
433 1
- 434 (b) Did you include complete proofs of all theoretical results? [Yes] see SM Section 1
- 435 3. If you ran experiments...
- 436 (a) Did you include the code, data, and instructions needed to reproduce the main exper-  
437 imental results (either in the supplemental material or as a URL)? [Yes] All code is  
438 submitted with the SM
- 439 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
440 were chosen)? [Yes] all parameters for all baselines are given in SM section 3.
- 441 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
442 ments multiple times)? [Yes]
- 443 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
444 of GPUs, internal cluster, or cloud provider)? [Yes] SM section 3
- 445 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 446 (a) If your work uses existing assets, did you cite the creators? [Yes] code from Simopt.org  
447 is cited.
- 448 (b) Did you mention the license of the assets? [N/A]
- 449 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
450 the benchmark test problems.
- 451 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
452 using/curating? [N/A]
- 453 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
454 information or offensive content? [N/A]
- 455 5. If you used crowdsourcing or conducted research with human subjects...
- 456 (a) Did you include the full text of instructions given to participants and screenshots, if  
457 applicable? [N/A]
- 458 (b) Did you describe any potential participant risks, with links to Institutional Review  
459 Board (IRB) approvals, if applicable? [N/A]
- 460 (c) Did you include the estimated hourly wage paid to participants and the total amount  
461 spent on participant compensation? [N/A]